

UNIVERSIDADE FEDERAL DO PARANÁ

IRAPURU HARUO FLÓRIDO

ANÁLISE TEXTUAL DE LETRAS DE MÚSICAS BRASILEIRAS
APLICANDO MÉTODOS DE BIOINFORMÁTICA E MINERAÇÃO DE TEXTO

CURITIBA

2025

IRAPURU HARUO FLÓRIDO

ANÁLISE TEXTUAL DE LETRAS DE MÚSICAS BRASILEIRAS APLICANDO
MÉTODOS DE BIOINFORMÁTICA E DE MINERAÇÃO DE TEXTOS

Tese apresentada ao curso de Pós-graduação em
Gestão da Informação, do Setor de Sociais
Aplicadas, da Universidade Federal do Paraná,
como requisito parcial para a obtenção do título de
Doutor em Gestão da Informação.

Orientador: Prof. Dr. José Simão de Paula Pinto
Coorientador: Prof. Dr. Roberto Tadeu Raittz

CURITIBA

2025

DADOS INTERNACIONAIS DE CATALOGAÇÃO NA PUBLICAÇÃO (CIP)
UNIVERSIDADE FEDERAL DO PARANÁ
SISTEMA DE BIBLIOTECAS – BIBLIOTECA CIÊNCIAS SOCIAIS APLICADAS

Flório, Irapuru Haruo

Análise textual de letras de músicas brasileiras aplicando métodos de bioinformática e de mineração de textos / Irapuru Haruo Flório. – Curitiba, 2025.

1 recurso on-line : PDF.

Tese (Doutorado) – Universidade Federal do Paraná, Setor de Ciências Sociais Aplicadas, Programa de Pós-Graduação em Gestão da Informação.

Orientador: Prof. Dr. José Simão de Paula Pinto.

Coorientador: Prof. Dr. Roberto Tadeu Raittz.

1. Gestão da Informação. 2. Música. 3. Bioinformática 4. Mineração de dados (Computação). I. Pinto, José Simão de Paula. II. Raittz, Roberto Tadeu. III. Universidade Federal do Paraná. Programa de Pós-Graduação em Gestão da Informação. IV. Título

Bibliotecário: Nilson Carlos Vieira Junior - CRB-9/1797



MINISTÉRIO DA EDUCAÇÃO
SETOR DE CIÊNCIAS SOCIAIS E APLICADAS
UNIVERSIDADE FEDERAL DO PARANÁ
PRÓ-REITORIA DE PÓS-GRADUAÇÃO
PROGRAMA DE PÓS-GRADUAÇÃO GESTÃO DA
INFORMAÇÃO - 40001016058P1

TERMO DE APROVAÇÃO

Os membros da Banca Examinadora designada pelo Colegiado do Programa de Pós-Graduação GESTÃO DA INFORMAÇÃO da Universidade Federal do Paraná foram convocados para realizar a arguição da tese de Doutorado de IRAPURU HARUO FLÓRIDO, intitulada: **ANÁLISE TEXTUAL DE LETRAS DE MÚSICAS BRASILEIRAS APLICANDO MÉTODOS DE BIOINFORMÁTICA E DE MINERAÇÃO DE TEXTOS**, sob orientação do Prof. Dr. JOSÉ SIMÃO DE PAULA PINTO, que após terem inquirido o aluno e realizada a avaliação do trabalho, são de parecer pela sua **APROVAÇÃO** no rito de defesa.

A outorga do título de doutor está sujeita à homologação pelo colegiado, ao atendimento de todas as indicações e correções solicitadas pela banca e ao pleno atendimento das demandas regimentais do Programa de Pós-Graduação.

CURITIBA, 12 de Dezembro de 2025.

Assinatura Eletrônica
16/12/2025 11:07:19.0
JOSÉ SIMÃO DE PAULA PINTO
Presidente da Banca Examinadora

Assinatura Eletrônica
16/12/2025 21:56:54.0
RONAN ASSUMPTÃO SILVA
Avaliador Interno (INSTITUTO FEDERAL DO PARANÁ)

Assinatura Eletrônica
19/12/2025 18:11:23.0
JERONIZA NUNES MARCAUKOSKI
Avaliador Externo (UNIVERSIDADE FEDERAL DO PARANÁ)

Assinatura Eletrônica
16/12/2025 10:46:42.0
WILSON LEMOS JUNIOR
Avaliador Externo (INSTITUTO FEDERAL DE EDUC., CIÊNCIA E
TECNOLOGIA DO PARANÁ)

DEDICATÓRIA

Aos meus pais (*in memoriam*), minha eterna gratidão. À minha esposa, Débora, dedico especialmente esta obra, pelo suporte incansável, pelo incentivo e por ser meu porto seguro diante das dificuldades desta jornada.

AGRADECIMENTOS

A Deus, pela vida e pela força para superar os desafios desta jornada.

À Universidade Federal do Paraná (UFPR) e ao Programa de Pós-Graduação em Gestão da Informação (PPGGI), pela oportunidade de realizar este doutorado e pela excelência do ambiente acadêmico proporcionado.

Ao meu orientador, Prof. Dr. **José Simão de Paula Pinto**, pela orientação segura, pela generosidade no compartilhamento de conhecimentos e pela confiança depositada em meu trabalho ao longo desta trajetória.

Ao meu coorientador, Prof. Dr. **Roberto Tadeu Raittz**, pelas valiosas contribuições, pelas discussões técnicas e pelo suporte fundamental ao desenvolvimento desta pesquisa.

Aos professores do PPGGI e aos membros da banca examinadora, pelas sugestões pertinentes que contribuíram para o aprimoramento e o amadurecimento deste estudo.

Aos colegas e amigos dos laboratórios do Setor de Educação Profissional e Tecnológica (BiolInfo) e do PPGGI, pela parceria, pelo convívio diário e pelas trocas de experiências que tornaram este percurso mais leve e colaborativo.

À minha esposa, **Débora**, às minhas filhas, **Nara e Ana**, e aos meus genros, **Jonathan e Thiago**, pelo apoio constante e pelo incentivo.

Aos meus netos, **Jonathan Filho e Valentina**, que trazem alegria e renovação aos meus dias e servem de inspiração para o futuro.

Aos meus pais, **Célia e Zito**, e aos meus sogros, **Delmar e Jandira** (*in memoriam*). Embora não estejam fisicamente aqui, fazem-se presentes em outra dimensão.

A todos que, direta ou indiretamente, contribuíram para a concretização deste estudo, meu muito obrigado!

RESUMO

A expansão dos serviços de *streaming* ampliou o acervo musical, tornando impraticável a classificação manual de atributos como gênero, instrumentação e sentimento devido ao volume de dados. No entanto, essa caracterização detalhada é indispensável para a eficácia dos sistemas de recomendação musical e da recuperação de informação musical, áreas fundamentais no cenário digital atual. Este estudo propôs o método de Análise Textual em Músicas Brasileiras (ATMBR), uma metodologia para a recuperação de informações musicais e a identificação automática de rótulos. O diferencial desta pesquisa foi a seleção e a aplicação de um algoritmo de bioinformática, o SWeeP, em conjunto com algoritmos de mineração de texto e de aprendizagem de máquina. A escolha do SWeeP se justifica por sua característica intrínseca de baixo custo computacional, permitindo o processamento massivo de dados textuais sequenciais com eficiência e velocidade de 10 a 100 vezes superiores às dos métodos tradicionais de alinhamento. A metodologia ATMBR envolve a realização de experimentos utilizando um corpus de letras de músicas brasileiras extraído de sítios de letras. A amostra original bruta do *corpora* continha 138 mil títulos musicais. Após um processo intenso de pré-processamento, cura e normalização dos dados, o método aplica o algoritmo SWeeP para gerar representações vetoriais (*embeddings*) detalhadas das letras. O objetivo foi desenvolver modelos robustos para a classificação textual de músicas. Como contribuições, a pesquisa apresenta a proposição do método ATMBR e a disponibilização de uma nova base de dados de músicas brasileiras, curada e rotulada, que servirá como recurso valioso para pesquisas futuras. O experimento resultou na criação da plataforma Ritmo Brasil, um ambiente de consulta e pesquisa voltado a músicos e apreciadores da música popular brasileira. O método proposto gerou modelos de classificação para a análise de sentimentos e para o mapeamento de gêneros musicais. A validação demonstrou que a estratégia híbrida de classificação alcançou desempenho superior, com 92% de acurácia na classificação de emoções básicas (Alegria, Tristeza, Raiva e Medo). Além disso, o método comprovou ser capaz de recuperar, de forma não supervisionada, a estrutura hierárquica e a genealogia dos gêneros musicais a partir do conteúdo semântico das letras. O estudo, de natureza interdisciplinar (Computação, Arte e Biologia), reforça a importância de explorar novas tecnologias para lidar com o volume exponencial de dados na ciência da informação, especialmente no contexto musical brasileiro.

Palavras-chave: Gestão da informação; Letras de músicas; Rotulação de músicas; Recuperação de informação musical; Aprendizagem de máquina.

ABSTRACT

The expansion of streaming services has broadened the musical repertoire, making manual classification of attributes such as genre, instrumentation, and sentiment impractical due to the volume of data. However, this detailed characterization is indispensable for the effectiveness of music recommendation systems and music information retrieval, crucial areas in the current digital landscape. This study proposed Análise Textual em Músicas Brasileiras (ATMBR), a methodology for retrieving musical information and automatically identifying labels. The distinguishing feature of this research was the selection and application of the bioinformatics algorithm SWeeP, in conjunction with text-mining and machine-learning algorithms. The choice of SWeeP is justified by its low computational cost, which enables efficient processing of sequential textual data at 10-100 times the speed of traditional alignment methods. The ATMBR methodology involves conducting experiments using a corpus of Brazilian song lyrics extracted from lyric websites. The original raw corpus contained 138,000 musical titles. After preprocessing, curation, and data normalization, the method applies the SWeeP algorithm to generate detailed vector representations (embeddings) of the lyrics. The objective was to develop robust models for classifying songs. As contributions, the research presents the proposed ATMBR method and a new, curated, labeled database of Brazilian songs, which will serve as a valuable resource for future research. The experiment led to the creation of the Ritmo Brasil platform, a consultation and research environment for musicians and enthusiasts of Brazilian popular music. The proposed method generated classification models for sentiment analysis and the mapping of musical genres. Validation demonstrated that the hybrid classification strategy achieved superior performance, with 90.5% accuracy in classifying basic emotions (Joy, Sadness, Anger, and Fear). Furthermore, the method proved capable of unsupervised recovery of the hierarchical structure and genealogy of musical genres from the lyrics' semantic content. This interdisciplinary study (Computer Science, Art, and Biology) underscores the importance of exploring recent technologies to manage the exponential volume of data in information science, particularly in the Brazilian musical context.

Keywords: Information management; Song lyrics; Song labeling; Machine learning; Music information retrieval (MIR).

LISTA DE FIGURAS

Figura 1 – Vendas de música por segmento em 2022, 2023 e 2024.	19
Figura 2 – Áreas do conhecimento envolvidas no método proposto.	25
Figura 3 - Ciclo da Interdisciplinaridade	27
Figura 4 – Faturamento mundial da indústria fonográfica (1999 a 2024)	29
Figura 5 – Variedade de rótulos em músicas	30
Figura 6 – Amostra de trecho de uma música executada em flauta.....	31
Figura 7 – Estrutura da Música Popular Brasileira.	32
Figura 8 – Tarefas e Métodos de AM.	39
Figura 9 – Inteligência Artificial aplicada à Bioinformática	41
Figura 10 - Diagrama do funcionamento do método ATM.....	58
Figura 11 – Árvore dos gêneros musicais brasileiros	60
Figura 12 – Pré-processamento e cura do corpora musical.....	61
Figura 13 - Diagrama da ATMBR recorte processo SWeeP.....	64
Figura 14 - Consulta do termo “saudades” na ferramenta HTML-TM	65
Figura 15 – Consulta das 20 músicas mais relacionadas com o termo “saudades”....	66
Figura 16 – Consulta do termo “Saudades” e as distâncias entre os termos.....	66
Figura 17 - Frequência do termo “saudades” na linha do tempo.....	67
Figura 18 – Plataforma RITMO BRASIL.....	73
Figura 19 – Tela de consulta opção “Músicas”	74
Figura 20 - Tela de consulta opção “Termos”	74
Figura 21- Incidência da termo Saudades durante as décadas 60 a 2020	75
Figura 22 – Correlação de termos versus títulos.....	76
Figura 23 – Consulta com o termo “ALEGRIA” em todo o corpora.	76
Figura 24 - Termos relacionados com o termo principal	78
Figura 25 - Listas de letras relacionadas ao termo principal “ALEGRIA.”	79
Figura 26 - Dendrograma dos termos relacionados com “ALEGRIA.”	79
Figura 27 – Termo “ALEGRIA” aplicado às músicas ao longo das décadas.....	79
Figura 28 – Dendrograma do termo GAIVOTA	80
Figura 29 – Pesquisa da memória histórica, termo indígena “TUPI”	82
Figura 30 – Memória histórica, lista de músicas relacionadas “TUPI”	83
Figura 31 – Pesquisa na dimensão regionalismo termo “MINEIRO”	83

Figura 32 – Regionalismo, lista de músicas relacionadas com “MINEIRO”	84
Figura 33 – Matriz de confusão, análise de emoções VADER	92
Figura 34 – Matriz de confusão da MLP (PCA=34, HD =377).....	93
Figura 35 – Dendrograma dos gêneros musicais.....	94

LISTA DE QUADROS

Quadro 1- Palavras-chave utilizadas para pesquisa e quantitativos	24
Quadro 2 - Vantagens do uso do algoritmo SWeeP.....	47
Quadro 3 - Recuperação de informação musical em letras de músicas	50
Quadro 4 – Informações do conjunto de dados extraído do streaming Vagalume	59
Quadro 5 – Exemplo de Conversão de termos para pesquisa	63
Quadro 6 - Produção criativa artificial dos termos relacionados com “GAIVOTA”.....	81
Quadro 7 – Amostra de respondentes	86
Quadro 8 – Corpora selecionado para os experimentos	91

LISTA DE TABELAS

Tabela 1 – Quantitativos de títulos musicais, termos versus gênero.....	72
Tabela 2 – Resultado aplicação do algoritmo VADER	92
Tabela 3 – Resultados comparativos algoritmos supervisionados	93

LISTA DE ABREVIATURAS OU SIGLAS

AM	Aprendizagem de Máquina
AMTA	American Music Therapy Association
ATMBr	Método de Análise Textual em Músicas Brasileiras
BERT	Bidirectional Encoder Representations from Transformers
BSF	Biological Sequence Format
CBOW	Continuous Bag of Words
CNN	Redes Neurais Convolucionais
DL	Deep Learning
DNA	DeoxyriboNucleic Acid
DNAbits	Conversão de caracteres de texto dividindo cada byte em pares de bits
FASTA	Formato de Sequência Biológica
GloVe	Global Vectors for Word Representation
IA	Inteligência Artificial
IFPI	International Federation of the Phonographic Industry
KNN	K-vizinhos mais próximos
LDA	Latent Dirichlet Allocation
LSTM	Long Short-Term Memory
ML	Machine Learning (Aprendizagem de Máquina)
MPB	Música Popular Brasileira
MT	Mineração de Texto
NB	Näive Bayes
NGS	Next Generation Sequencing
NMF	Non-Negative Matrix Factorization
PCA	Análise de Componentes Principais
PIBIC	Programa Institucional de Bolsas de Iniciação Científica
PLN	Processamento de Linguagem Natural
RIM	Recuperação de Informação Musical
RNN	Redes Neurais Recorrentes
SRM	Sistemas de Recomendação Musical
SVM	Máquinas de Vetores de Suporte
SWeeP	Spaced Words Projection
UNESCO	United Nations Educational, Scientific, and Cultural Organization

VADER Valence Aware Dictionary and sEntiment Reasoner

SUMÁRIO

1	INTRODUÇÃO	18
1.1	JUSTIFICATIVA.....	22
1.2	OBJETIVOS.....	22
1.3	METODOLOGIA DE PESQUISA.....	23
1.4	ORGANIZAÇÃO DO ESTUDO	24
2	FUNDAMENTAÇÃO TEÓRICA	25
2.1	MÚSICA.....	28
2.1.1	Músicas Digitais.....	29
2.1.2	Características da Música Digital.....	31
2.1.3	Música Popular Brasileira	32
2.1.4	Importâncias dos rótulos.....	33
2.2	EMOÇÕES E SENTIMENTOS HUMANOS	33
2.3	SISTEMAS DE RECOMENDAÇÃO MUSICAL (SRM).....	34
2.4	MINERAÇÃO DE TEXTO	36
2.4.1	Aplicação da mineração de texto.....	36
2.4.2	Métodos de Mineração de Texto	37
2.5	APRENDIZAGEM DE MÁQUINA	38
2.5.1	Tipos de aprendizagem de máquina.....	39
2.5.2	Pesquisas Recentes	40
2.6	ALGORITMOS DE BIOINFORMÁTICA	41
2.6.1	Tipos de Algoritmos de Bioinformática.....	42
2.6.2	Pesquisas Mais Recentes em Algoritmos de Bioinformática	43
2.6.3	Algoritmo de Bioinformática SWeeP	45
2.6.4	Método de mineração de texto com algoritmos de bioinformática	48
2.7	ESTADO DA ARTE DA RECUPERAÇÃO DE INFORMAÇÃO MUSICAL EM LETRAS DE MÚSICAS.....	49
2.7.1	Aplicações atuais na recuperação de informação textual em músicas	49
2.7.2	Algoritmos de aprendizagem de máquina para análise textual.....	51
2.8	MÉTODOS DE INCORPORAÇÃO DE PALAVRAS	53
2.8.1	Modelo Word2Vec	53
2.8.2	Modelo GloVe	54

2.8.3	Modelo BERT	55
3	METODOLOGIA PARA ANÁLISE TEXTUAL DE MÚSICAS	57
3.1	VISÃO GERAL DO MÉTODO ATMBR.....	57
3.1.1	Fonte do <i>corpora</i> musical.	58
3.1.2	Pré-processamento para a cura e normalização do <i>corpora</i> musical	60
3.1.1	Aplicação do algoritmo de bioinformática SWeeP	62
3.1.2	Geração de modelos do <i>corpora</i> processado	67
3.1.3	Métricas de avaliação	68
3.1.4	Validação do método ATMBR	68
3.2	CONSIDERAÇÕES DO MÉTODO ATMBR	69
4	EXPERIMENTO DO MÉTODO ATMBR	70
4.1	Vieses de análises do experimento	70
4.2	Ambiente de desenvolvimento computacional	70
4.3	Seleção do <i>corpora</i> musical para o experimento	71
4.4	Plataforma de letras de músicas Ritmo Brasil	73
4.5	Correlação dos termos com as músicas	75
4.6	Análise da dimensão emoção no experimento	76
4.7	Análise da dimensão criativa no experimento.....	79
5	APRESENTAÇÃO DOS RESULTADOS	86
5.1	Avaliação da plataforma Ritmo Brasil	86
5.1.1	Visão Geral da Amostra.....	86
5.1.2	Correlações entre Formação Escolar e Gostos Musicais	86
5.1.3	Preferência de gêneros e relação com a faixa etária.....	87
5.1.4	O quanto é útil a plataforma para os músicos.....	87
5.1.5	Possibilidades a serem exploradas pelos músicos	88
5.1.6	Recursos que deveriam ser adicionados	89
5.1.7	Finalidade da Plataforma e Uso para Pesquisas	89
5.2	Resultados dos experimentos da análise de sentimentos.	91
5.3	Definição do dendrograma dos gêneros musicais	94
5.4	Resultados da dimensão descritiva	95
6	CONSIDERAÇÕES FINAIS	96
	REFERÊNCIAS	97
	Apêndice A – Artigo aceito na revista Advances Knowledge Representation (AKR).....	97

Apêndice B – Rotulação manual de músicas do conjunto de dados da MSD.....	122
Apêndice C – Formulário de pesquisa de opinião da plataforma	134
Apêndice D – Tabela de conversão entre AMNOcode e DNAbits.....	140
Apêndice E – Métricas de avaliação.....	141
Apêndice F – Manual de instrução de consulta Ritmo Brasil	142

1 INTRODUÇÃO

A música, ao longo da história da civilização, tem sido valorizada como meio de preservar a memória cultural. A UNESCO considera-a patrimônio cultural imaterial da humanidade porque é uma expressão que grupos e indivíduos transmitem de geração em geração, preservando a identidade do povo em suas comunidades (Peralta, 2013). Como aconteceu com o gênero musical “samba de roda”, que passou a ser reconhecido como patrimônio imaterial nacional pelo Instituto do Patrimônio Histórico e Artístico Nacional (IPHAN) e como patrimônio oral da humanidade pela UNESCO Silveira (2016).

Além da função de entretenimento e de memória cultural das épocas da nossa sociedade, a música tem sido usada também no âmbito educacional. Especificamente em processos metacognitivos, facilitando a aprendizagem em muitas áreas do conhecimento (Florido et al., 2023). Cabe ressaltar que, por vezes, a música também está presente em terapias voltadas ao tratamento de algumas doenças, como o *Alzheimer*¹. A memória musical ocupa, no cérebro, um espaço diferente das demais memórias; por isso, mesmo em estágios avançados da doença de Alzheimer, pacientes conseguem recuperar lembranças de sua música favorita, inclusive relacionadas à infância. A música proporciona uma experiência emocional mais intensa, e esses tipos de memórias afetivas são preservadas de maneira mais duradoura (Baird; Samson, 2009)

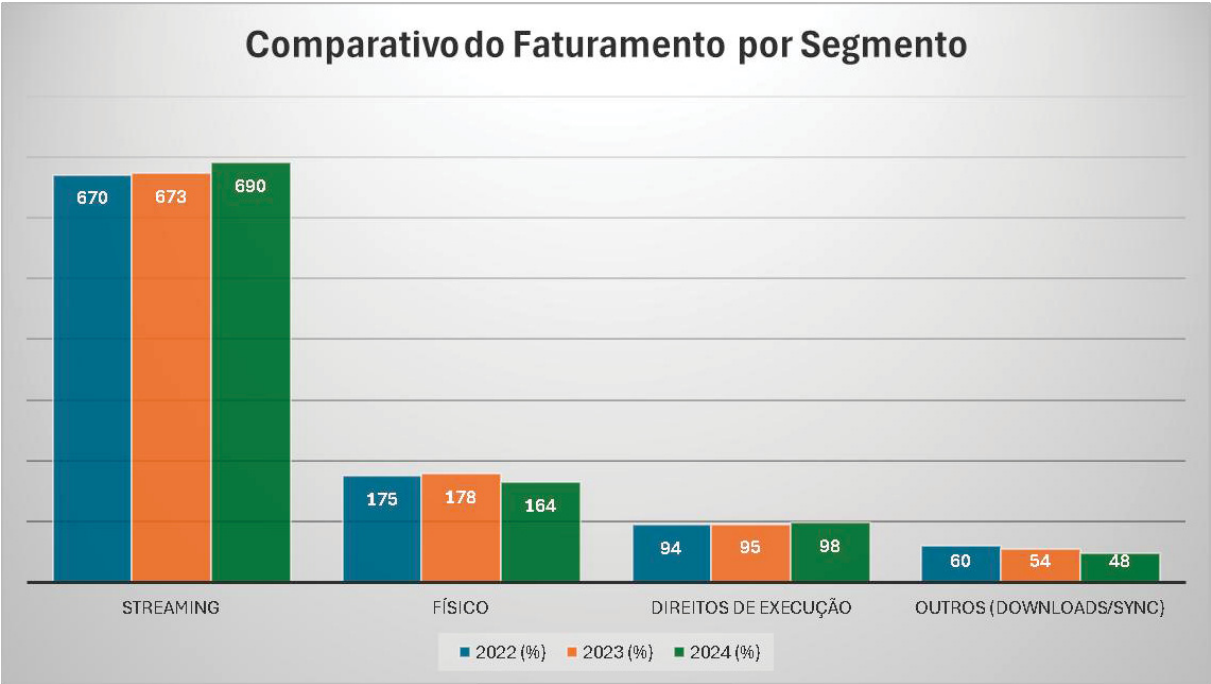
A música é uma forma de terapia que pode ajudar as pessoas a relaxar e a se sentirem melhor, sendo usada como recurso de conforto mental por indivíduos em todo o mundo. Embora existam músicas para aliviar o estresse, diferentes gêneros musicais têm efeitos diferentes nas pessoas (Thoma et al., 2013). Pesquisas clínicas específicas sobre o público que ouve músicas mostram que cerca de 30% das pessoas escolhem tipos inadequados de música para relaxar e, como resultado, os seus níveis de estresse aumentam, por isso, torna-se necessária a escolha adequada do tipo de música de acordo com o propósito de uso (Chang; Huang; Wu, 2017).

¹ *Alzheimer* - é um transtorno neurodegenerativo progressivo e fatal que se manifesta pela deterioração cognitiva e da memória, comprometimento progressivo das atividades de vida diária e uma variedade de sintomas neuropsiquiátricos e de alterações comportamentais

A transformação digital vivida nas últimas décadas e a recente tecnologia de *streaming*¹ incorporada às mídias de vídeo e música, facilitou e tornou o acesso à música mais rápido e imediato (Song, 2024).

Atualmente, as músicas estão disponíveis em vários reprodutores, como smartphones, *tablets*, *notebooks* e outros dispositivos portáteis. A música, como forma de entretenimento, alimenta uma indústria promissora que cresce exponencialmente a cada dia. Conforme levantamento apresentado nos relatórios anuais de 2022, 2023 e 2024 da IFPI², o crescimento e a participação das músicas no formato de *streaming* são expressivos a cada ano (Figura 1).

Figura 1 – Vendas de música por segmento em 2022, 2023 e 2024.



Fonte: Relatórios de 2022, 2023 e 2024 da IFPI (2025).

Os serviços de *streaming* de música dependem, de forma intrínseca, da Recuperação de Informação Musical (RIM) e dos Sistemas de Recomendação de Música (SRM), que os alimentam. Esse conjunto de tecnologias permite ao usuário

¹ *Streaming* - Tecnologia de transmissão de conteúdo online que permite consumir filmes, séries e músicas sem a necessidade de baixar no dispositivo de acesso.

² IFPI é uma das vozes da indústria fonográfica mundial, representando mais de 8.000 gravadoras em todo o mundo. Promovem o valor da música gravada, fazem campanha pelos direitos dos produtores de discos e expandem o comércio de música gravada em todo o mundo.

encontrar lançamentos e músicas antigas em meio a um repertório de milhões de obras musicais disponíveis nessas plataformas (Webster, 2020).

Para que os sistemas de recomendação de músicas funcionem, é necessária uma recuperação eficiente de informações musicais para fornecer dados básicos, como o nome do artista, o título da obra, o álbum, o ano de lançamento e o gênero musical. Além destes dados, é importante ampliar a disponibilidade de atributos dos *corpora*¹ musicais. Para que seja possível conhecer novas músicas e artistas até então desconhecidos, e que atuem em determinados grupos de gêneros musicais, vocais, instrumentos e sentimentos, é necessário classificá-los no contexto mencionado ao se levantar o “DNA”² dos *corpora* musicais (Sordo et al., 2013; Café; Barros, 2018).

Em 1999, Tim Westergren iniciou um projeto na área da música, denominado *Music Genome Project*. O projeto demandou um grande esforço para "capturar a essência da música em nível fundamental", identificando em torno de 400 atributos e em alguns casos até 450, dependendo do gênero musical, para descrever as músicas (Oramas, 2025). Participaram mais de quarenta musicólogos, que rotularam manualmente milhares de arquivos musicais a partir de 2000. Com base nesse projeto, um sítio de música muito conhecido nos EUA, intitulado Pandora, foi criado, explorando essas rotulações manuais, resultantes do projeto. (Joyce, 2006).

A aplicação de métodos computacionais de mineração de texto e de aprendizagem de máquina pode auxiliar os musicólogos a reduzir o tempo e o custo, tanto em recursos quanto em trabalho manual de rotulação. Esse estudo frequentemente demanda conhecimentos específicos e muito tempo para a sua realização. A fim de compreender a extensão temporal dedicada à rotulação manual de músicas, um experimento que selecionou 250 músicas para identificar rótulos de gêneros, emoções e instrumentos demandou 384 horas, envolvendo quatro pessoas — não músicos nem musicólogos — no projeto PIBIC (2014). Ver Apêndice B, parte do resultado do trabalho.

A recuperação de informação musical lida com a dificuldade de definir uma "unidade de significado" básica, com métodos de segmentação e com a necessidade

¹ *Corpora* - coletânea, reunião de textos ou documentos sobre um assunto ou tema, repertório ou aquilo que registra toda a obra de um autor (termo latino no plural; singular: *corpus*).

² DNA – Ácido Desoxirribonucleico: funciona como um “livro de receitas” para o desenvolvimento, o funcionamento e a reprodução da vida.

de focar na satisfação do usuário, e não apenas na precisão. A área enfrenta barreiras, como a escassez de fontes de informação semântica, especialmente para idiomas além do inglês. Há ainda a necessidade de anotação para a web semântica, bem como a lentidão dos algoritmos semânticos em comparação com os métodos tradicionais (Ren; Bracewell, 2009). Nesse contexto, pode-se dizer que a maior dificuldade em sistemas de recuperação de informação não é apenas extrair a informação em si, mas também decidir sobre a relevância das características dessas informações em relação à necessidade do usuário, o que torna esse desafio ainda mais complexo.

Portanto, a relevância é um fator importante também nas análises textuais dos sistemas de recuperação de informação musical. Além desse aspecto e de outros citados, há uma carência de *corpora* disponíveis para validar os métodos de recuperação de informação, especificamente nas letras de músicas de gêneros e subgêneros da música popular brasileira. Com o intuito de preencher essa lacuna, propõe-se o Método de Análise Textual em Músicas Brasileiras (ATMBR).

É importante reforçar que o presente estudo se situa no âmbito da Ciência da Informação, com ênfase na Gestão da Informação e na subárea de Organização e Recuperação da Informação (ORI). Enquanto a Gestão da Informação estrutura processos sistemáticos para a produção e o uso da informação, visando à tomada de decisão e à construção do conhecimento, a ORI desempenha um papel central ao fornecer os métodos e ferramentas necessários para o acesso eficiente e qualificado a esses recursos informacionais (CHOO, 2003).

No contexto específico da Recuperação da Informação Musical, a música é conceituada como um objeto informacional complexo, não estruturado e multidimensional. Diferentemente de documentos textuais convencionais, os registros musicais integram camadas sonoras, simbólicas e contextuais, exigindo estratégias de representação que abranjam desde a melodia e o ritmo até o contexto histórico-cultural, conforme apontam os desafios examinados por Downie (2003).

A integração entre os princípios da Gestão da Informação e os modelos de recuperação musical promove o desenvolvimento de sistemas mais eficazes e interoperáveis. Essa abordagem não apenas favorece a preservação e o acesso ao patrimônio musical, mas também reforça o papel estratégico da Gestão da Informação na mediação entre a produção cultural, os sistemas tecnológicos e a apropriação do conhecimento pelo usuário (Downie et al., 2010).

1.1 JUSTIFICATIVA

A necessidade de rotulagem automática em larga escala é uma premissa para a eficiência dos modernos sistemas de recomendação de músicas (Turnbull et al., 2022), Florido; Raittz, 2018). Contudo, os métodos tradicionais de mineração de texto e aprendizagem de máquina, embora eficazes, frequentemente esbarram em alto custo computacional e de tempo quando aplicados a milhões de letras de música (Gotham; Bemman; Vatulkin, 2025).

Esta proposta de investigação justifica-se por abordar essa lacuna metodológica. O diferencial desta pesquisa é a aplicação de algoritmos de bioinformática na análise textual de letras de música. A escolha motiva-se pela característica intrínseca desses algoritmos, originalmente projetados para o sequenciamento genético, de processar massivamente dados sequenciais com eficiência e velocidade, de 10 a 100 vezes superior, que os métodos tradicionais de alinhamento textual não oferecem (De Pierri et al., 2020)

Ao adaptar essas ferramentas da biologia computacional para a mineração de texto, busca-se uma metodologia singular, propondo um método de baixo custo computacional para rotular automaticamente um grande volume de músicas. Este estudo contribui, portanto, em duas frentes principais:

- a) Contribuição Metodológica: a proposição do método ATMBR, que combina mineração de texto com algoritmos de bioinformática para classificar sentimentos, emoções, gêneros e estilos em letras de música de forma eficiente.
- b) Contribuição de Dados: a disponibilização de uma nova base de dados de músicas brasileiras, com corpus curado e rotulado, servindo como recurso valioso e "padrão-ouro¹" para futuras pesquisas na área.

1.2 OBJETIVOS

Sêneca (2017), um dos representantes do estoicismo na Roma antiga afirmava: “Quando se navega sem destino, nenhum vento é favorável”. Assim, toda pesquisa precisa estabelecer objetivos que norteiem o estudo, servindo como estratégia para dimensionar e traçar os caminhos necessários para alcançar os resultados propostos.

¹ Padrão-Ouro – termo designado para um subconjunto de dados curado e validado manualmente por especialistas, servindo como verdade fundamental (*ground truth*) para aferir a acurácia dos algoritmos

O objetivo geral deste estudo é definir uma metodologia para análise textual de letras de músicas e para a rotulagem musical no contexto de aprendizado de máquina em *corpora* textuais de músicas brasileiras e aplicar métodos robustos de algoritmos de mineração de texto combinados com algoritmos de bioinformática capazes de processar massivamente dados da mineração de texto, valendo-se de sua característica intrínseca de eficiência proveniente do sequenciamento genético.

Relaciona-se a seguir os objetivos específicos para atingir o objetivo geral:

- Pesquisar as referências teóricas sobre análise textual e rotulação de músicas, mineração de texto, algoritmos de bioinformática e sua integração com métodos de aprendizagem de máquina;
- Analisar e comparar métodos e algoritmos de mineração de texto associados a algoritmos de bioinformática para análise textual e rotulação de músicas;
- Realizar pré-processamento aplicando métodos de mineração de texto, variando e verificando parâmetros para obtenção de uma base de dados curada e vetorizada;
- Construir uma base de dados de músicas brasileiras, *corpora* musical curada e vetorizada, para métodos de aprendizagem de máquina e rotulação de músicas;
- Desenvolver modelos de AM supervisionada a partir de *corpora* curados e vetorizados para rotular automaticamente emoções e gêneros musicais.

1.3 METODOLOGIA DE PESQUISA

A metodologia de pesquisa adotada possui caráter predominantemente qualitativo, com a finalidade de analisar os resultados obtidos no método na busca por assuntos relacionados ao tema proposto. As bases indexadas de trabalhos científicos consultadas foram direcionadas aos temas por meio de palavras-chave, o que resultou em quantitativos, conforme apresentado no Quadro 1. O referencial teórico baseou-se na seleção e na leitura de fontes primárias (artigos e teses), secundárias (capítulos de livros, artigos de revisão e relatórios) e terciárias (bibliotecas e resumos).

Para isso, foram consultados títulos e autores que abordam temáticas relacionadas ao estudo, cujo teor se mostrou fundamental para o desenvolvimento desta proposta. Após a seleção e a verificação das publicações, restaram, ao final do processo, 57 estudos utilizados para compor o referencial teórico e o estado da arte.

Quadro 1- Palavras-chave utilizadas para pesquisa e quantitativos

Tema	Palavras-chave Operador Booleano	Período Ano	Bases		Seleção total
			CAPES Encontrado /Seleção	Dimensions Encontrado /Seleção	
Recuperação de Informação Musical em letras de música	recuperação AND informação AND musical AND letras OR música	2019 a 2025	21 / 2	5 / 1	3
Textual analysis of song lyrics, text mining, and machine learning	((lyrics) AND ((text mining) OR (machine learning)))	2019 a 2025	10 / 8	12295 / 240	248
Song lyric, textual music, textual machine learning, text mining, AI	((lyric OR song) OR ((textual OR music) AND analysis)) AND ((machine AND learning) OR (text AND mining) OR (AI))	2019 a 2025	10/2	4005/616	618
Sentiment emotion song lyric textual machine learning text mining AI	((sentiment OR emotion) OR ((textual OR music) AND analysis)) AND ((machine AND learning) OR (text AND mining) OR (AI))	2019 a 2025		1400/502	503
Total de publicações selecionadas para pré-processamento					1372

Fonte: O autor (2025).

É importante salientar que, nas ferramentas de busca da CAPES e do Dimensions, foram aplicadas restrições de pesquisa, como publicações dos últimos sete anos (2019 a 2025) e estudos situados nas áreas de aprendizagem de máquina e de inteligência artificial.

1.4 ORGANIZAÇÃO DO ESTUDO

Além da Introdução e objetivos apresentados no **Capítulo 1**, esta proposta de tese está organizada da seguinte forma: (i) **Capítulo 2**, que apresenta os fundamentos teóricos aplicados ao método, o estado da arte e os trabalhos correlatos; (ii) **Capítulo 3**, que descreve o método proposto de ATMBR e os passos que compõem sua aplicação; (iii) **Capítulo 4**, que traz a aplicação do experimento piloto do método; (iv) **Capítulo 5**, no qual são discutidos os resultados do experimento inicial; e (v) **Capítulo 6**, que contém as conclusões e as perspectivas de trabalhos futuros. Ao final, encontram-se as referências e os apêndices.

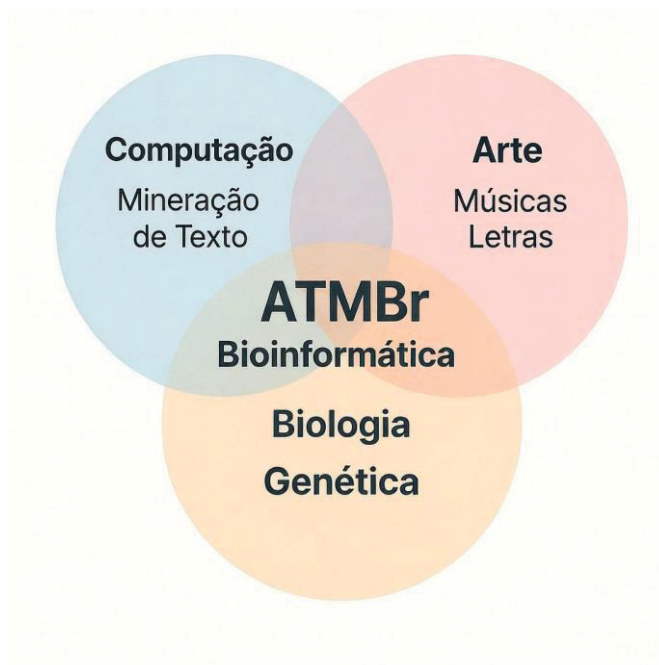
2 FUNDAMENTAÇÃO TEÓRICA

Dada a natureza interdisciplinar das linhas de pesquisa do Programa de Pós-graduação em Gestão da Informação, no qual se insere esta tese, cabe contextualizar as áreas de conhecimento e as disciplinas abrangidas pelo trabalho, bem como a relação de interdisciplinaridade entre elas.

A primeira delas é a arte e sua manifestação, por meio da música, como forma de expressão cultural da humanidade. A segunda corresponde à aplicação da computação, mais especificamente de algoritmos de aprendizagem de máquina para o reconhecimento de padrões e a recuperação de informação musical. A terceira abrange a biologia, considerada aqui como o aspecto inovador, por meio da utilização de algoritmos de bioinformática, com o intuito de reduzir significativamente o custo computacional no processamento de dados por meio de algoritmos de AM.

Na Figura 2, apresenta-se um diagrama que identifica as disciplinas abordadas e as áreas comuns subjacentes ao método proposto.

Figura 2 – Áreas do conhecimento envolvidas no método proposto.



Fonte: O Autor (2025).

Para Japiassu (1994), a interdisciplinaridade caracteriza-se pela intensidade das trocas entre especialistas e pelo grau de integração estabelecido em um projeto de pesquisa. Segundo o autor, é o próprio objeto de pesquisa que gera a necessidade

científica de interdisciplinaridade. A principal característica da interdisciplinaridade é combinar os resultados de várias disciplinas, tomando deles emprestados desenhos conceituais de análise para integrá-los após a comparação e a avaliação.

Segundo Piaget (1964), o termo interdisciplinar deve ser usado quando há colaboração entre diferentes disciplinas ou entre diferentes departamentos de uma mesma ciência, levando a uma interação em sentido estrito, a uma certa reciprocidade na troca e a um enriquecimento mútuo em toda a linha. Essa colaboração reforça a ideia de que o conhecimento avança quando diferentes áreas dialogam e compartilham suas estruturas conceituais.

Os desafios fundamentais e os problemas de alcance global são negligenciados pelas disciplinas científicas individuais. Eles são abordados apenas no campo da filosofia, mas não recebem contribuições significativas das ciências. Nessas circunstâncias, as mentes moldadas pelas disciplinas perdem sua capacidade natural de contextualizar o conhecimento, assim como de integrá-lo em seus contextos mais amplos (Morin, 2013). A diminuição da percepção do todo leva a uma diminuição da responsabilidade, ou seja, cada pessoa tende a se sentir responsável apenas por sua área de especialização, assim como há uma diminuição da solidariedade, isto é, cada pessoa não sente mais os laços com seus concidadãos (Morin, 2013).

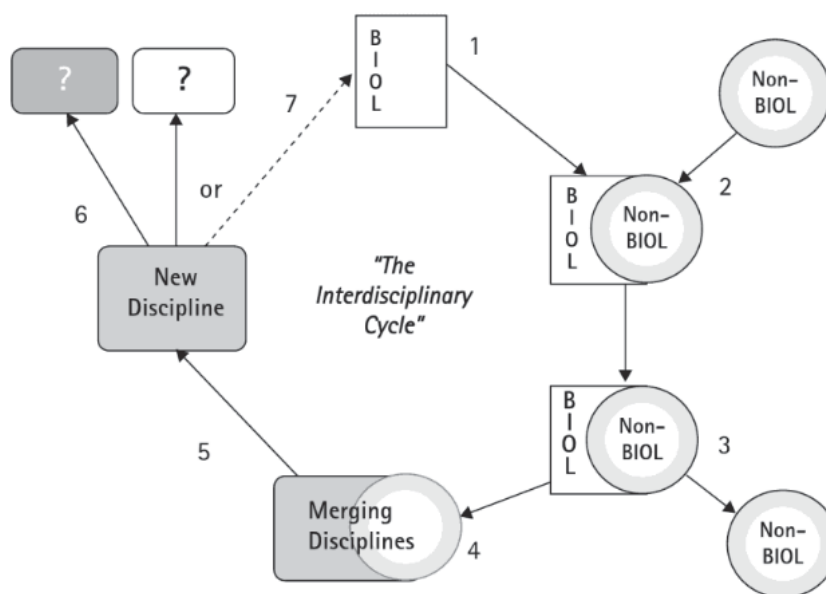
Em geral, o futuro pertence à pesquisa interdisciplinar, mas, na prática, tal pesquisa é muitas vezes difícil de organizar devido à ausência de entendimento mútuo, deliberadamente fomentada. O primeiro objetivo da pesquisa deve, portanto, identificar esferas nas quais sejam possíveis comparações entre as tendências de evolução e desenvolvimento nas ciências humanas, a fim de incentivar a colaboração e os intercâmbios interdisciplinares, ou, simplesmente, fortalecer a pesquisa em cada ramo por meio das comparações obtidas. Deve-se reconhecer que o problema é muito mais complicado no ser humano do que nas ciências naturais (Fazenda, 2012).

Segundo Burggren et al. (2017), na Figura 3, apresenta-se o ciclo da interdisciplinaridade nas Ciências Biológicas e em outras disciplinas que vivenciam a interdisciplinaridade e como esta se cruza (1). Uma disciplina autônoma, como a biologia experimental, frequentemente influencia as disciplinas não biológicas (2). Essa mistura de ideias e técnicas pode ser, em última análise, apenas passageira (3) ou resultar em uma verdadeira fusão das disciplinas (4). Essa nova disciplina, formada

pela fusão da biologia e da não-biologia (5), pode eventualmente fragmentar-se em novas disciplinas ou persistir (6) para retornar ao ciclo interdisciplinar.

O trabalho interdisciplinar nas ciências biológicas pode ser desafiador. Comunicar-se com colaboradores de outras disciplinas requer reaprender conceitos disciplinares dependentes, adotar um novo vocabulário e comprometer-se com novas abordagens.

Figura 3 - Ciclo da Interdisciplinaridade



Fonte: (Burggren et al., 2017)

Mesmo quando concluída com sucesso, a ciência interdisciplinar pode não ser totalmente apreciada por avaliadores conservadores ou mais tradicionalmente inclinados (Burggren et al., 2017). Apesar dessas limitações, a interdisciplinaridade nas ciências biológicas vem crescendo, impulsionada por um espectro de motivações que vai desde a curiosidade intelectual desenfreada até soluções práticas para problemas médicos e de engenharia. Claramente, no futuro, as ciências biológicas continuarão a operar dentro de um ciclo interdisciplinar, gerando novas subdisciplinas, como aconteceu com a bioinformática, mudando o tecido da própria biologia (Burggren et al., 2017).

A ciência não é um processo linear de acúmulo de conhecimento, mas sim um processo de mudança revolucionária. Kuhn define uma revolução científica como uma mudança fundamental na estrutura conceitual de uma disciplina científica,

frequentemente acompanhada de alterações nos métodos de pesquisa, nos valores e nos paradigmas científicos (Kuhn; Hacking, 2012).

A utilização de métodos de pesquisa e de paradigmas científicos aplicados à música contribui significativamente para o avanço do conhecimento musical. Esse conhecimento pode ser empregado para aprofundar a compreensão da música, criar interpretações e impulsionar tanto o ensino quanto o aprendizado musical. Dentro desse contexto, na sequência, abordam-se os referenciais teóricos deste estudo, como as disciplinas de Música, com enfoque na recuperação de informação musical; Aprendizagem de Máquina; Mineração de Texto; e Bioinformática, que compõem o presente estudo.

2.1 MÚSICA

Beethoven, trocando cartas com Bettina, escreveu certa vez (Sullivan, 1936):

"Fale com Goethe sobre mim. Diga a ele para ouvir minhas sinfonias e ele dirá que estou certo em dizer que a música é a única entrada incorpórea para o mundo superior do conhecimento, que compreende a humanidade, mas que a humanidade não pode compreender"

Estima-se que a exposição intensa à música seja uma forma de enriquecimento que exerce efeito semelhante nas regiões do hipocampo do cérebro. Alguns estudos mostram o poder da música sobre quem a ouve e quem a toca. Atualmente, pesquisas qualitativas e quantitativas são realizadas pela American Music Therapy Association (AMTA)¹ para explicar os diferentes efeitos da música em pessoas doentes e de diferentes idades, tornando o poder terapêutico da música específico para cada doença (Côrte; Lodovici Neto, 2009).

A música é composta por três elementos básicos: harmonia, melodia e ritmo. Harmonia é um conjunto agradável de notas musicais entoadas simultaneamente (acordes). A melodia é tocada separadamente, especialmente ao cantar. O ritmo é a noção métrico-temporal da música, ou seja, a batida e o tempo de duração (Pilhofer; Day, 2019).

¹ AMTA - Fundada em 1998, resultante da fusão entre a Associação Nacional de Musicoterapia (fundada em 1950) e a Associação Americana de Musicoterapia (fundada em 1971).

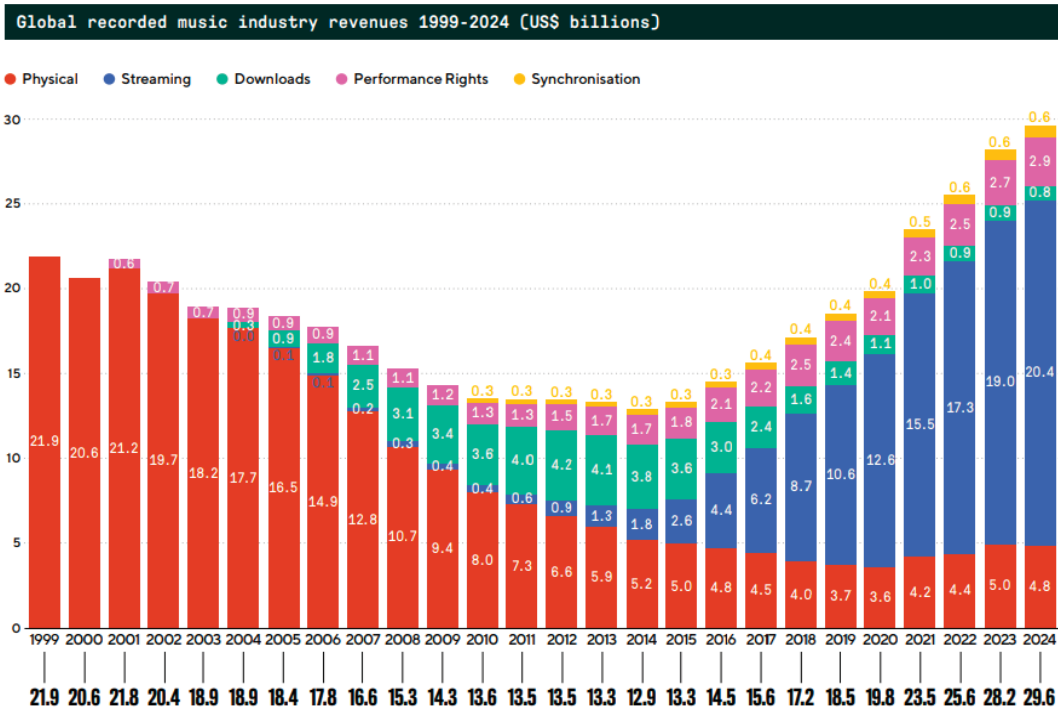
2.1.1 Músicas Digitais

Músicas digitais são uma forma de música representada por dados. Esses dados podem ser armazenados em um computador ou em outro dispositivo digital e reproduzidos por meio de um reprodutor de música digital.

Ela é diferente da música analógica, que é representada por ondas sonoras. A música analógica é gravada em um suporte físico, como um disco de vinil ou uma fita cassete, e reproduzida por meio de um dispositivo físico. A música digital apresenta algumas vantagens em relação à música analógica: é mais durável, pois não sofre danos físicos em sua mídia, e tem maior portabilidade, pois pode ser armazenada em um computador ou em outro dispositivo digital. Além disso, pode ser facilmente compartilhada com outras pessoas por meio de e-mail, sites de compartilhamento de arquivos e redes sociais (Scherzinger, 2019; Strawn; Pohlmann, 1986)

Conforme relatório anual da Federação Internacional da Indústria Fonográfica (IFPI)¹ de 2025 (Figura 4), observa-se um crescimento sensível, a partir de 2014, no volume de faturamento da indústria fonográfica.

Figura 4 – Faturamento mundial da indústria fonográfica (1999 a 2024)



Fonte: IFPI - Relatório de 2025

¹ IFPI é uma das vozes da indústria fonográfica mundial, representando mais de 8.000 gravadoras em todo o mundo. Promovem o valor da música gravada, fazem campanha pelos direitos dos produtores de discos e expandem o comércio de música gravada em todo o mundo.

Cabe ressaltar também que o faturamento de vendas de músicas no formato de *streaming* continua a ser o principal motor de receita, ultrapassando um marco significativo, uma vez que as receitas de *streaming* atingiram US\$ 20,4 bilhões. O *streaming* agora representa 69% de toda a receita da indústria musical, e as receitas de assinaturas pagas cresceram 9,5%. O número global de utilizadores com contas de subscrição paga atingiu os 752 milhões (Industry, 2025).

Uma preocupação central do relatório é o impacto da inteligência artificial (IA) generativa. A IFPI apela aos legisladores para que protejam a música e a criatividade humana, garantindo que a IA seja usada para "apoiar e amplificar" a criatividade, e não para "substituí-la" (Industry, 2025).

A música digital é um dos tipos de dados mais importantes disponíveis na Internet. Existem diversos estudos e métodos para a análise de conteúdo de áudio e de letras, que empregam diferentes características e técnicas. A identificação e seleção das músicas em uma base com milhares de músicas de múltiplas origens e gêneros, utilizando processos manuais, requer amplos conhecimentos culturais, sociais e históricos da música (Pachet et al., 2005).

Figura 5 – Variedade de rótulos em músicas



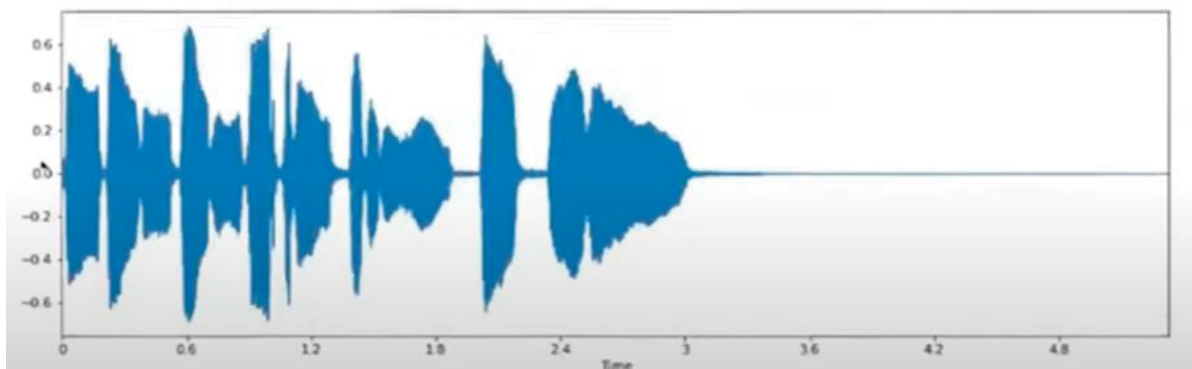
Fonte: O autor (2023).

Constitui desafio a aplicação de métodos de AM para integrar, de forma eficiente, informações provenientes de diferentes fontes, a fim de realizar a rotulação automática de músicas. As informações a serem integradas são diversas, como se observa na Figura 5, e incluem rótulos de gêneros musicais, estilos, emoções, ambientes, instrumentos e outros aspectos. Este é um problema de difícil solução, e uma das muitas dificuldades é a diversidade de categorias de rótulos relevantes para rotular músicas (Turnbull et al., 2009). Além disso, a combinação e a fusão desses rótulos permitirão que as músicas contidas em bases ou repositórios de larga escala, que, em muitos casos, apresentam deficiências no preenchimento dos rótulos, sejam rotuladas automaticamente.

2.1.2 Características da Música Digital

O som, assim como a música, é formado por ondas com diferentes comprimentos, amplitudes e durações ao longo de um determinado período. No exemplo da Figura 6, apresenta-se um breve trecho de uma música executada em flauta. A quantidade de ondas formadas por segundo é bem elevada e essa medida é denominada hertz.

Figura 6 – Amostra de trecho de uma música executada em flauta.



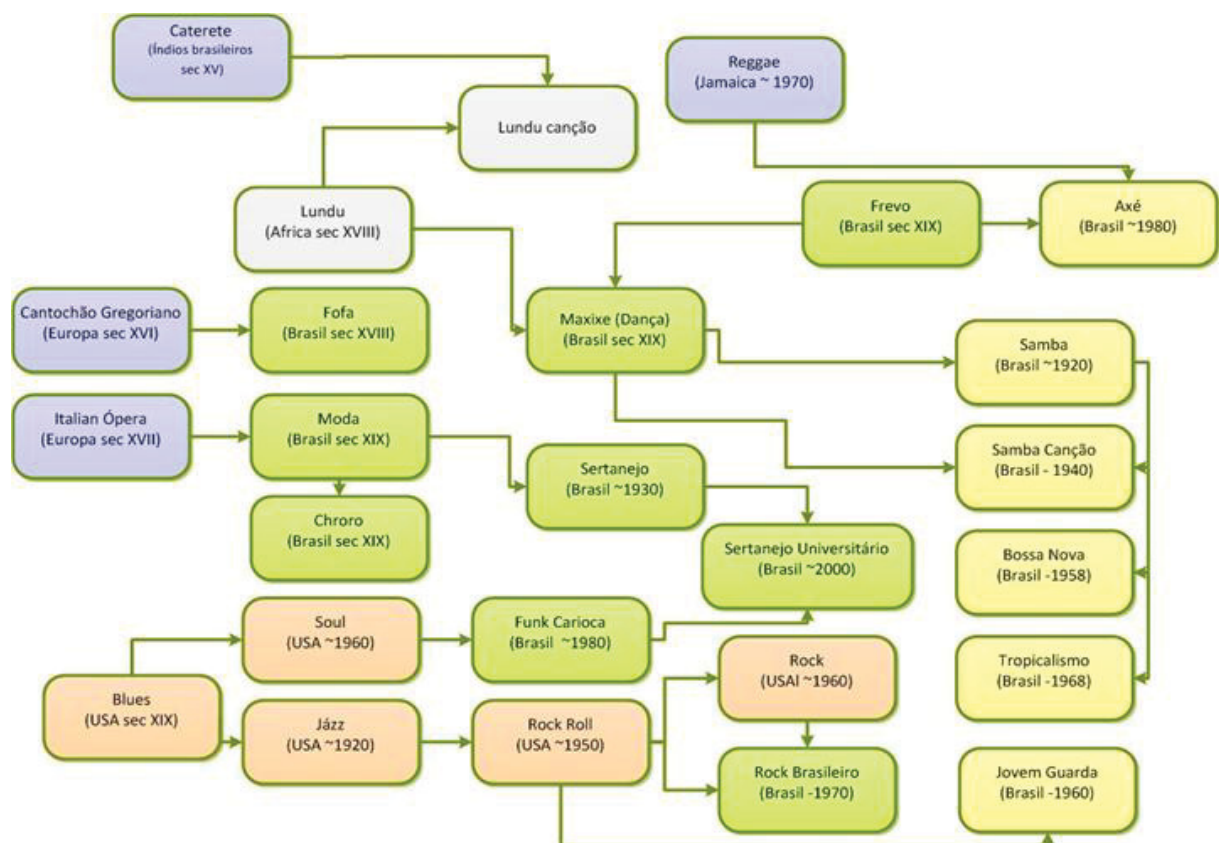
FONTE: Labrosa (2023).

Por este motivo, a impressão é compactada. A partir deste sinal de áudio é possível extrair várias características da música, como timbre, altura, compasso e outras propriedades. Os descritores, características ou atributos do sinal de áudio utilizados na maioria dos trabalhos que fazem a recuperação de informação musical (RIM), foram originalmente propostos por (Tzanetakis; Cook, 2000)

2.1.3 Música Popular Brasileira

A música popular brasileira (MPB) teve sua origem durante o Brasil Colônia, marcado pela chegada dos primeiros colonizadores portugueses. A interação entre a cultura portuguesa, as tradições dos indígenas locais e a influência africana, aliada ao sincretismo musical, propiciou o surgimento da MPB. Caldas (2010), em sua obra, realiza um estudo sobre as origens e a história da MPB, permitindo definir a estrutura e a taxonomia dos vários gêneros da música brasileira no período que vai do descobrimento do Brasil, em 1500, até os dias atuais, conforme apresentado na Figura 7. Interessante observar que o momento atual da MPB trouxe várias influências de gêneros musicais estrangeiros, absorvidas ao longo do tempo, precisamente a partir dos anos 50.

Figura 7 – Estrutura da Música Popular Brasileira.



Fonte: (Caldas, 2010)

A partir do diagrama apresentado, torna-se evidente a presença de influências estrangeiras, que se iniciam com o blues, jazz, soul, funk, reggae e rock na MPB. Essas influências, por conseguinte, permeiam tanto as letras quanto as melodias,

constituindo extensões significativas, conforme a Figura 7. O artista, pelo caráter peculiar da sua criatividade, sente-se livre para expressar sua mensagem e seus sentimentos em outra língua além do idioma nativo (Caldas, 2010).

2.1.4 Importâncias dos rótulos

A forma mais simples e direta de seleção de músicas em repertórios pessoais e em ambientes de *streaming* de música normalmente ocorre com base no título, no artista e/ou no gênero musical, mas as emoções e/ou sentimentos do ouvinte no momento de apreciar determinada música também influenciam sua escolha. Por este motivo, faz-se necessário abordar também os sentimentos e emoções humanas e como identificá-los no *corpus* musical para auxiliar no processo de recomendação (Schedl et al., 2022 ;Velankar; Kulkarni, 2023)

2.2 EMOÇÕES E SENTIMENTOS HUMANOS

A psicologia define emoção e sentimento como experiências humanas complexas que estão intimamente ligadas, mas também apresentam diferenças importantes. Segundo Damásio (1996), as emoções são respostas fisiológicas e comportamentais a estímulos internos ou externos. Elas são mediadas por circuitos neurais no tronco encefálico e nas amígdalas cerebrais. As emoções são geralmente rápidas e intensas e podem ser expressas por meio de alterações na expressão facial, na linguagem corporal, na frequência cardíaca, na respiração e na atividade hormonal.

Conforme Camras et al. (1981) e Ekman (2003), as principais diferenças entre emoção e sentimento, segundo a psicologia são:

- a) duração: as emoções são mais curtas que os sentimentos;
- b) intensidade: as emoções são mais intensas que os sentimentos;
- c) estímulo: as emoções são uma resposta a um estímulo, enquanto os sentimentos podem ser induzidos por um estímulo ou por fatores internos;
- d) componente cognitivo: os sentimentos têm um componente cognitivo mais forte que as emoções; e
- e) consciência: as emoções geralmente são conscientes, enquanto os sentimentos podem ser inconscientes.

As emoções geralmente são divididas em duas categorias principais: emoções básicas e complexas (Camras; Plutchik; Kellerman, 1981; Ekman, 2003). Emoções

básicas são aquelas compartilhadas por todas as pessoas, independentemente de sua cultura ou contexto social. As emoções básicas incluem alegria, tristeza, raiva, medo e nojo (Ekman, 2003). Emoções complexas são aquelas influenciadas por fatores cognitivos, como pensamentos, crenças e memórias. As emoções complexas incluem amor, orgulho, ciúme, inveja, gratidão, desesperança, culpa e vergonha (Ekman, 2003).

Além da classificação citada, básicas e complexas, as emoções podem ser geralmente divididas em duas categorias principais: sentimentos positivos e negativos. Emoções positivas são aquelas associadas a uma sensação de bem-estar, como felicidade, amor e satisfação, já as emoções negativas são aquelas associadas a uma sensação de mal-estar, como tristeza, raiva e medo (Camras; Plutchik; Kellerman, 1981; Damásio, 1996).

É importante compreender as emoções e os sentimentos para entender melhor nossas próprias experiências e as dos outros. As emoções e os sentimentos podem afetar nosso comportamento, nossos pensamentos e nosso bem-estar físico e mental. A psicologia estuda as emoções e os sentimentos para entender como eles funcionam, como são influenciados por fatores internos e externos e como podem ser usados para melhorar nossa vida (Damásio, 1996).

2.3 SISTEMAS DE RECOMENDAÇÃO MUSICAL (SRM)

Os SRM têm desempenhado um papel essencial no acesso à música digital, permitindo que os usuários descubram novas faixas e artistas. Esses sistemas fazem uso de diversas técnicas de AM e análise de dados para fornecer sugestões personalizadas de músicas (Afchar et al., 2022). Apesar dos avanços recentes, os SRM ainda não são perfeitos e podem gerar recomendações insatisfatórias. Isso ocorre porque os gostos e necessidades musicais dos usuários são complexos e influenciados por uma variedade de fatores, como preferências pessoais, o ambiente e o contexto, entre outros.

Os SRM tradicionais geralmente concentram-se em interações usuário-item, como histórico de reprodução, avaliações e rótulos (Schedl, 2019). No entanto, essas informações não são suficientes para capturar a riqueza e a complexidade dos gostos musicais dos usuários. Para fornecer recomendações mais precisas e relevantes, os SRM devem considerar os aspectos intrínsecos, extrínsecos e contextuais dos

ouvintes. Os aspectos intrínsecos referem-se às preferências pessoais do usuário, como o gênero, o estilo e o clima. Os aspectos extrínsecos estão relacionados ao contexto social e cultural do usuário, como a localização, a idade e o gênero. Já os aspectos contextuais referem-se ao momento específico em que o usuário está recebendo a recomendação, como o humor, a atividade e o evento (Schedl et al., 2018).

Outro desafio importante dos SRM é o aspecto do "problema da partida a frio", *cold-start problem*, que surge quando um sistema tenta recomendar músicas recém-adicionadas que ainda não têm histórico de interação do usuário (Okada; Tan; Kamioka, 2021)

Os avanços recentes das redes neurais profundas e do Deep Learning (DL) abriram novas possibilidades para aprimorar os sistemas de recomendação musical. Espera-se que esses sistemas se tornem cada vez mais precisos e personalizados para as necessidades dos usuários (Fessahaye et al., 2019). Os SRM são usados hoje por uma variedade de serviços de música, incluindo *Spotify*, *Apple Music*, *Deezer*, *Youtube Music* e *Amazon Music*¹. Esses serviços usam os SRM para recomendar músicas para os usuários baseado em seus dados de histórico de escuta (Velankar; Kulkarni, 2023).

Existem basicamente dois tipos principais de SRM: com filtro colaborativo e com filtro baseado em conteúdo. O filtro colaborativo usa os históricos de escuta dos usuários para recomendar músicas de outros com gostos semelhantes. E o filtro baseado em conteúdo usa informações sobre as músicas, como gênero, artista e título, para recomendar músicas semelhantes às que o usuário já ouviu. Os SRM são uma tecnologia em rápido desenvolvimento, tornando-se cada vez mais sofisticados por meio da aplicação da IA, e são usados por uma variedade de serviços de música. Além disso, eles também serão usados por diversos outros serviços, como serviços de vídeo, de livros e de jogos. Os SRM têm o potencial de revolucionar a forma como consumimos conteúdo; ademais, podem auxiliar na descoberta de novos conteúdos de interesse e promover economia de tempo (Schedl et al., 2022).

A análise textual e rotulação de músicas propostas neste trabalho demandam o uso de metodologias de mineração de texto. Essas metodologias são essenciais para

¹ As empresas relacionadas fornecem serviços de streaming musical na internet: Spotify®, Apple® Music®, Deezer®, Amazon Music® e YouTube Music®.

a compreensão do trabalho, pois permitem identificar padrões e relações nos dados textuais, o que é fundamental para definir rótulos adequados às músicas.

2.4 MINERAÇÃO DE TEXTO

A mineração de texto tem um viés de contribuição diferenciado para o desenvolvimento de ferramentas voltadas à descoberta de relações de termos que dificilmente uma consulta com lógica booleana pura consegue trazer em bases indexadas (Lowe et al., 2018). A mineração de texto traz novas informações e conhecimentos ocultos por meio da extração automática de dados de diferentes recursos escritos. Essencialmente, trata-se da relação entre as informações extraídas para formar novos fatos ou hipóteses a serem exploradas por meios mais convencionais de experimentação.

A diferença da mineração de texto da pesquisa na Web é que, na pesquisa tradicional, busca-se algo conhecido, normalmente escrito por outra pessoa. A dificuldade está em deixar de lado o material irrelevante e em encontrar as informações pertinentes às necessidades. Já na mineração de texto, o objetivo é descobrir informações até então desconhecidas, algo que ninguém conhece e ainda não publicado (Hearst, 2003).

Uma das principais tendências na mineração de texto nos últimos anos é o uso de abordagens baseadas em aprendizado de máquina. Essas abordagens envolvem o treinamento de algoritmos de aprendizado de máquina para identificar padrões e relações em dados textuais. Isso tem permitido a criação de modelos preditivos e classificadores eficientes para tarefas como extração de informações, classificação de sentimentos e sumarização de textos (Shruti; Priyanka, 2021).

2.4.1 Aplicação da mineração de texto.

A mineração de texto, na prática, consiste na execução de processos sequenciais e interativos que transformam ou organizam uma grande quantidade de informações e documentos em uma estrutura sistematizada, permitindo sua utilização de forma inteligente e eficiente. As etapas que compõem a mineração de texto podem ser definidas pelos seguintes passos (Kamran Kowsari et al., 2021; Larocca Neto et al., 2000):

- pré-processamento de dados textuais;

- processamento de texto; e
- pós-processamento pós-mineração de texto.

Naturalmente, a redução do número de palavras dos documentos na mineração de texto é uma exigência não apenas na representação vetorial, mas também em qualquer outra representação textual. É importante também que o método de pré-processamento seja robusto, isto é, capaz de lidar com textos que contenham ruídos, erros gramaticais e erros de digitação.

2.4.2 Métodos de Mineração de Texto

O campo da mineração de texto se concentra na extração de padrões, tendências e conhecimento útil a partir de grandes volumes de dados textuais não estruturados (Souza et al., 2017). Para que os métodos e algoritmos clássicos de mineração possam ser comumente aplicados, é indispensável uma etapa de pré-processamento (Shruti; Priyanka, 2021). Esta etapa é fundamental, pois transforma o texto bruto, ruidoso e não estruturado em uma representação estruturada, limpa e normalizada, que pode ser compreendida pelos algoritmos.

Essas técnicas de pré-processamento, já consolidadas na literatura nacional e internacional, são amplamente adotadas e incluem várias etapas fundamentais (Larocca Neto et al., 2000 ; De Lucca; Nunes, 2002):

- a) *Case Folding*, que consiste em converter todos os caracteres de um documento para um mesmo formato, caixa alta (maiúsculas) ou caixa baixa (minúsculas);
- b) *Stemming*, que consiste em converter cada palavra em seu radical, eliminando flexões verbais, sufixos e prefixos. Algoritmos empregados para isso geralmente incorporam uma grande quantidade de conhecimentos linguísticos, de modo que dependem do idioma.
- c) *Stopwords*, palavras de parada que ocorrem com frequência em um documento. Como são muito comuns, possuem pouca informação sobre o conteúdo do texto, sendo geralmente removidas da representação;
- d) Representação N-grama, que é uma alternativa para a contenção e remoção de palavras de parada, sendo um N-grama uma fatia de N caracteres de uma sequência mais longa; e

e) Lematização, que representa as palavras por meio do infinitivo dos verbos e do masculino singular de substantivos e adjetivos.

Uma técnica amplamente aplicada à mineração de texto e ao processamento de linguagem natural é o TF-IDF (*Term Frequency – Inverse Document Frequency*). Serve para medir a importância de uma palavra (ou termo) em um documento em relação a um corpus. É uma técnica estatística que atribui pesos às palavras com base na frequência local (no documento) e na raridade global (no corpus). É simples, interpretável, eficiente e ainda hoje é uma das bases mais importantes da mineração de texto, mesmo com o avanço de modelos neurais mais sofisticados (Salton; McGill, 1987).

A mineração de texto com aplicação em letras de músicas é uma área de pesquisa que busca identificar padrões e relações nos dados textuais dessas letras. Esses dados podem ser as letras, informações sobre artistas e músicas, ou até mesmo comentários dos usuários (Saluja; Jain; Yadav, 2019). Ela também fornece meios para transformar grandes quantidades de dados textuais brutos em informações estruturadas e significativas, essenciais para alimentar e aprimorar os algoritmos de aprendizagem de máquina, que discutiremos a seguir.

2.5 APRENDIZAGEM DE MÁQUINA

A aprendizagem de máquina (ou *Machine Learning*, em inglês) é um ramo da IA que se dedica a desenvolver algoritmos capazes de aprender com dados e de tomar decisões ou fazer previsões sem serem explicitamente programados para cada tarefa. O conceito de aprendizagem de máquina teve suas raízes na década de 1940, quando pesquisadores começaram a explorar como as máquinas poderiam aprender com experiências passadas. No final dos anos 1950 e 1960, surgiram as primeiras abordagens de aprendizagem de máquina, como o *Perceptron* de Frank Rosenblatt, um modelo de rede neural artificial (Buchanan, 2005). Contudo, durante as décadas de 1970 e 1980, ocorreu o “inverno da IA”, período em que o progresso da área estagnou devido às limitações computacionais e aos desafios teóricos (Buchanan, 2005).

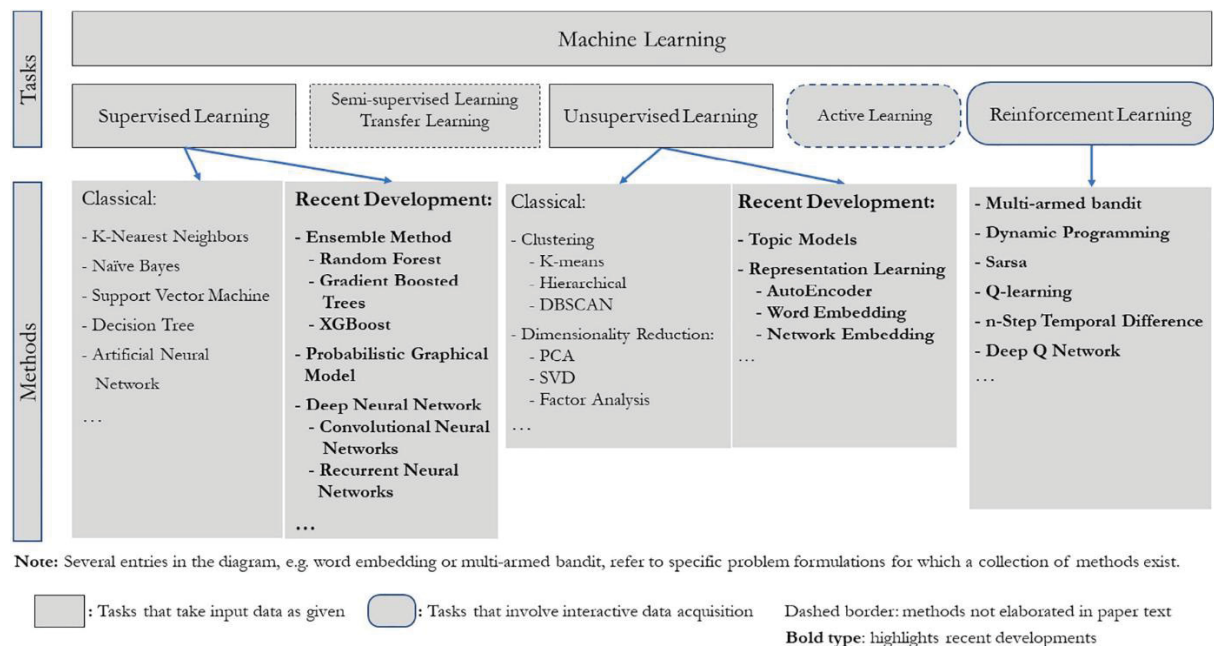
A virada ocorreu nos anos 1990, quando o aumento do poder computacional e a disponibilidade de grandes conjuntos de dados impulsionaram o avanço do campo. Algoritmos baseados em aprendizado supervisionado, como as máquinas de vetores

de suporte (SVM) e as redes neurais profundas, começaram a se destacar. A partir do século XXI, com a explosão do acesso à internet, a quantidade de dados disponíveis cresceu exponencialmente, catalisando ainda mais o desenvolvimento da área e possibilitando a adoção generalizada de soluções de AM em diversos setores (Ma; Sun, 2020).

2.5.1 Tipos de aprendizagem de máquina

A classificação canônica dos paradigmas da aprendizagem de máquina, popularizada por Mitchell (1997), utiliza como critério o tipo de “experiência” ou “*feedback*” que o algoritmo recebe durante o treino. Esta abordagem divide o campo nos três tipos principais:

Figura 8 – Tarefas e Métodos de AM.



Fonte: (Ma; Sun, 2020)

- Aprendizado Supervisionado:** os algoritmos são treinados em um conjunto de dados rotulados, em que cada exemplo possui uma entrada e a saída esperada correspondente. O objetivo é fazer com que o modelo aprenda a mapear as entradas para as saídas corretas, de modo que, ao se deparar com novos dados, ele possa fazer previsões precisas.
- Aprendizado Não Supervisionado:** diferentemente do aprendizado supervisionado, o algoritmo é treinado com um conjunto de dados não rotulado. O objetivo é encontrar padrões, estruturas ou grupos naturais nos

dados sem qualquer orientação explícita. Isso pode incluir técnicas de clusterização, redução de dimensionalidade e associação; e

- c) Aprendizado por Reforço: o algoritmo aprende por meio de interações com um ambiente. Ele recebe *feedback* na forma de recompensas ou punições após cada ação e busca aprender uma estratégia para maximizar a recompensa total ao longo do tempo.

A Figura 8, apresenta os algoritmos principais de aprendizagem supervisionada e não supervisionada (Ma e Sun,2020).

2.5.2 Pesquisas Recentes

A área de AM é altamente dinâmica e tem registrado inúmeros avanços recentes. Algumas das pesquisas mais notáveis incluem:

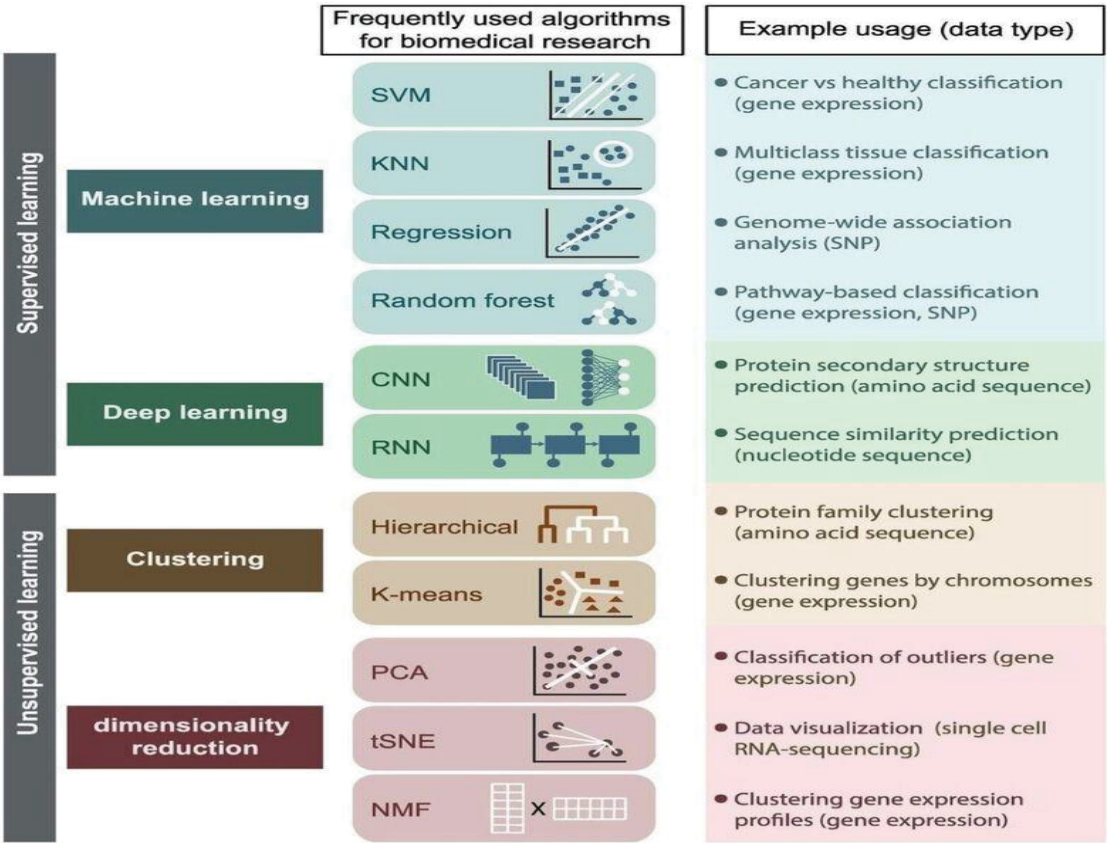
- a) Redes Neurais Profundas (*Deep Learning*): as redes neurais profundas têm sido fundamentais para o avanço em visão computacional, processamento de linguagem natural e outros domínios. Pesquisas contínuas têm focado em melhorar a eficiência e a interpretabilidade dessas redes, bem como em desenvolver arquiteturas inovadoras.
- b) Aprendizado por Transferência: essa área busca aplicar conhecimentos adquiridos em uma tarefa para melhorar o desempenho em outra tarefa relacionada. É especialmente útil quando os dados de treinamento são escassos, permitindo que modelos pré-treinados sejam ajustados para tarefas específicas;
- c) Aprendizado Federado: concentra-se em desenvolver algoritmos capazes de aprender modelos em dispositivos distribuídos (como *smartphones*) sem precisar compartilhar dados centralmente, abordando preocupações com privacidade e segurança;
- d) Aprendizado Automático com Menos Dados (*Few-Shot Learning*): nessa linha de pesquisa, o objetivo é treinar modelos capazes de generalizar e aprender com poucos exemplos de treinamento, tornando os algoritmos mais eficientes e flexíveis; e
- e) Ética em AM: com o aumento do uso de algoritmos em aplicações críticas, a pesquisa em ética em AM tem se concentrado em questões como o viés algorítmico, a transparência, a justiça e a responsabilidade.

Essas são apenas algumas das muitas áreas de pesquisa em AM que estão sendo exploradas atualmente. Com o rápido avanço da tecnologia e a contínua colaboração entre pesquisadores e a indústria, espera-se que essa área continue a crescer e a impactar significativamente diversos setores da sociedade. Ela vem sendo cada vez mais usada em algoritmos de bioinformática, uma área da ciência que estuda a coleta, o armazenamento, a análise e, principalmente, a interpretação de dados biológicos.

2.6 ALGORITMOS DE BIOINFORMÁTICA

A bioinformática é uma área multidisciplinar que combina conceitos da biologia, da ciência da computação e da estatística para analisar e interpretar grandes volumes de dados biológicos. Os avanços tecnológicos na biologia molecular, como o sequenciamento do DNA, geraram uma grande quantidade de informações genéticas e moleculares que precisavam ser processadas e interpretadas de forma eficiente. Isso levou ao desenvolvimento de algoritmos específicos para lidar com essas tarefas desafiadoras (Mariano et al., 2019).

Figura 9 – Inteligência Artificial aplicada à Bioinformática



Fonte: (Ma; Sun, 2020)

Na década de 1980, com o Projeto Genoma Humano em andamento, houve uma demanda crescente por métodos computacionais para analisar sequências de DNA. Nessa época, surgiram os primeiros algoritmos para alinhar sequências e identificar similaridades entre elas. Desde então, a bioinformática tornou-se uma área essencial na pesquisa biomédica e tem impulsionado avanços significativos na compreensão da genética e da biologia molecular (Wang et al., 2004). Na Figura 9, apresentam-se os principais algoritmos de IA, suas respectivas áreas da bioinformática e suas aplicações. Um exemplo de uso de cada algoritmo, com os respectivos dados de entrada, está indicado à direita.

Para melhor entendimento, descrevem-se as abreviações utilizadas: SVM, máquinas de vetores de suporte; KNN, K vizinhos mais próximos; CNN, redes neurais convolucionais; RNN, redes neurais recorrentes; PCA, análise de componentes principais; t-SNE, incorporação de vizinhos estocásticos t-distribuídos; NMF, fatoração de matriz não negativa. Estes são os algoritmos principais da IA com as respectivas áreas da bioinformática e suas aplicações (Auslander; Gussow; Koonin, 2021)

2.6.1 Tipos de Algoritmos de Bioinformática

Os algoritmos de bioinformática abrangem uma ampla gama de aplicações, e alguns dos principais tipos são (Van Ryssen, 2006):

- a) Alinhamento de Sequências: usado para comparar sequências de DNA, RNA ou proteínas, identificando regiões conservadas e semelhanças funcionais. O alinhamento é fundamental para identificar genes, estudar a evolução molecular e compreender a estrutura e a função de proteínas. Entre os métodos mais utilizados estão o alinhamento global (*Needleman-Wunsch*) e o alinhamento local (*Smith-Waterman*);
- b) Montagem de Genomas: lida com o desafio de montar fragmentos de sequências de DNA obtidos por técnicas de sequenciamento, reconstituindo o genoma completo. A montagem de genomas é fundamental para compreender a estrutura genética dos organismos e identificar variações genéticas. O método de sobreposição e montagem (*Overlap-Layout-Consensus*) é amplamente utilizado;
- c) Anotação de Genomas: envolve identificar e descrever genes e elementos funcionais de um genoma, como exons, íntrons e regiões regulatórias.

Algoritmos de anotação atribuem funções a sequências e auxiliam no entendimento da organização genômica;

- d) Predição de Estrutura de Proteínas: busca prever a estrutura tridimensional de proteínas a partir de suas sequências de aminoácidos. Métodos como homologia comparativa e modelagem por *threading* são amplamente utilizados e
- e) Filogenia: constroem árvores evolutivas que representam as relações entre espécies ou sequências. Auxiliam no entendimento da evolução das espécies e no rastreamento da disseminação de doenças. O método de máxima verossimilhança é um dos mais utilizados.

2.6.2 Pesquisas Mais Recentes em Algoritmos de Bioinformática

A bioinformática está em constante evolução, impulsionada pelos avanços tecnológicos e pelas novas descobertas biológicas. Algumas das pesquisas mais recentes incluem:

- a) Bioinformática de Nova Geração: campo relativamente novo que se concentra no desenvolvimento de ferramentas e técnicas para analisar e interpretar dados biológicos gerados por sequenciamento de próxima geração (NGS). As técnicas de NGS permitem o sequenciamento de grandes quantidades de DNA, RNA ou proteínas em curto prazo. Isso levou a um aumento significativo na geração de dados biológicos, que podem ser difíceis de analisar e interpretar (Horner et al., 2009);
- b) Aprendizado de Máquina em Bioinformática: é uma subárea da IA que se concentra no desenvolvimento de algoritmos capazes de aprender com dados e de tomar decisões sem serem programados explicitamente. As aplicações de aprendizado de máquina em bioinformática são diversas e incluem (Libbrecht; Noble, 2015):
 - Predição de estrutura de proteínas: o aprendizado de máquina pode ser usado para prever a estrutura 3D de proteínas a partir de sua sequência de aminoácidos, útil para o desenvolvimento de novos medicamentos e tratamentos;

- Anotação funcional de genes: pode prever a função de genes a partir da sequência de DNA, auxiliando na compreensão do papel dos genes em doenças e no desenvolvimento de novos tratamentos;
 - Identificação de marcadores genéticos: pode identificar marcadores associados a doenças ou características específicas, úteis para o diagnóstico e o tratamento dessas doenças. Essas abordagens são usadas para predição de estrutura de proteínas, anotação funcional de genes, identificação de marcadores genéticos e outras tarefas.
- c) Bioinformática de Populações e Medicina Personalizada: campo relativamente novo que se concentra no estudo da variação genética entre populações e indivíduos. A medicina personalizada visa fornecer tratamentos adaptados a cada paciente, com base em sua composição genética. As pesquisas recentes nessa área buscam desenvolver algoritmos para identificar genes associados a doenças, prever respostas a tratamentos e realizar estudos de associação genômica, com base em técnicas de aprendizado de máquina e análise de dados complexos (Ginsburg; Willard, 2009);
- d) Biologia de Sistemas e Redes Biológicas: campo interdisciplinar voltado à compreensão do comportamento de sistemas biológicos complexos, como células, tecidos e órgãos. A bioinformática fornece ferramentas e técnicas para analisar e interpretar grandes volumes de dados, essenciais para estudos de biologia de sistemas. Os algoritmos de análise de redes são importantes para estudos de biologia de sistemas. As redes biológicas são representações matemáticas das interações entre componentes biológicos, como genes, proteínas e moléculas, e os algoritmos de análise de redes podem ser usados para identificar módulos funcionais, que são conjuntos de componentes que interagem entre si para realizar uma função específica. Esses módulos podem revelar propriedades emergentes em sistemas biológicos, que não podem ser explicadas apenas pelo comportamento dos componentes individuais (Hidalgo et al., 2009).
- e) Bioinformática Estrutural: campo da bioinformática que se concentra no estudo da estrutura e da dinâmica de macromoléculas biológicas, como proteínas e ácidos nucleicos. A estrutura determina a função dessas

moléculas, o que torna a bioinformática estrutural essencial. Algoritmos avançados de simulação e modelagem permitem gerar modelos estruturais e dinâmicos a partir de informações limitadas, como sequências de aminoácidos ou de nucleotídeos. Esses modelos ajudam a investigar a relação entre estrutura e função biológica, além de prever estruturas ainda não estudadas experimentalmente (Levitt, 2001).

Essas pesquisas e avanços em algoritmos de bioinformática têm desempenhado um papel fundamental na aceleração do conhecimento biomédico, abrindo caminho para novas terapias, diagnósticos mais precisos e uma melhor compreensão da vida em nível molecular. A interação contínua entre a ciência da computação e a biologia continuará a impulsionar o progresso nesta área vital da pesquisa científica (Gollery, 2005).

2.6.3 Algoritmo de Bioinformática SWeeP

O algoritmo de bioinformática SWeeP (sigla de *Spaced Words Projection*) é um modelo matemático e computacional inovador, desenvolvido para representar, de forma vetorial, grandes conjuntos de dados de sequências biológicas em vetores compactos. Originalmente, o SWeeP foi dedicado à vetorização em formato de sequência biológica (FASTA). O método transforma qualquer conjunto de sequências de aminoácidos em um único vetor, que pode representar, por exemplo, todas as proteínas de um organismo.

a) Finalidade do Algoritmo SWeeP

A principal finalidade do SWeeP é atender às exigências do *Big Data* na área biológica, especialmente quando a análise e a comparação de sequências em grande escala se tornam inviáveis ou extremamente demoradas com os métodos tradicionais.

As aplicações principais são:

- **Vetorização e Redução de Dimensionalidade:** o SWeeP representa sequências de proteínas em vetores relativamente pequenos, preservando a comparabilidade entre elas. Ele utiliza a projeção de conjuntos de k-Mers espaçados (palavras espaçadas) em uma base *quase-ortonormal, orientada aleatoriamente, para criar um vetor de alta dimensão (HDV), que é posteriormente projetado em um vetor de menor dimensão (redução de*

dimensionalidade).

- **Análise de Genomas e Proteomas Completos:** o SWeeP foi utilizado para gerar representações vetoriais de todos os proteomas bacterianos completos disponíveis no NCBI (mais de 10.324 proteomas). Essa representação é consistente e pode ser utilizada em análises de aprendizado de máquina.
- **Suporte à Taxonomia e Filogenia:** trata-se de uma ferramenta robusta para apoiar a análise e a discussão de modelos taxonômicos. É utilizada para a construção de árvores filogenéticas, como demonstrado na geração da maior árvore de proteomas bacterianos já construída; e
- **Mineração de Textos (MT):** o SWeeP também foi implementado no pacote BioTEX para vetorizar textos codificados em formato de sequência biológica (BSF), permitindo a aplicação de ferramentas de bioinformática em estratégias de MT.

b) Vantagens do SWeeP

O SWeeP oferece várias vantagens, especialmente em contextos de *Big Data* e análise filogenética, devido à sua natureza de método livre de alinhamento, ver Quadro 2:

Comparativo com Outras Ferramentas Similares

O SWeeP é um método livre de alinhamento, o que o coloca em contraste com os métodos tradicionais baseados em alinhamento e também, em comparação com outros métodos livres de alinhamento baseados em k-Mers.

Comparação com Métodos de Alinhamento

Métodos tradicionais de alinhamento de sequências, como o BLAST, são considerados o padrão atual, mas apresentam limitações:

- **Velocidade e Viabilidade:** o alinhamento de grandes conjuntos de dados, como genomas completos, requer muito tempo e, muitas vezes, torna-se inviável, podendo levar horas ou dias. O SWeeP oferece maior agilidade ao realizar tarefas que o alinhamento não consegue executar com eficiência em escala de *Big Data*.

Quadro 2 - Vantagens do uso do algoritmo SWeeP

Vantagem	Detalhes e Suporte nas Fontes
Agilidade e Rapidez	Permite a construção rápida e sensível de árvores filogenéticas compactas. Por ser um método livre de alinhamento, acelera significativamente a comparação de dados. Por exemplo, a construção de uma árvore filogenética a partir de 10.324 genomas bacterianos completos levou apenas algumas horas.
Eficiência Computacional	Reduz significativamente os custos computacionais. A sua capacidade de redução de dimensionalidade torna a Análise de Componentes Principais (PCA) viável para grandes volumes de sequências, necessitando de apenas 0,027 GB para 600 coordenadas, em contraste com os 190 GB estimados para a matriz completa.
Qualidade e Robustez da Análise	Foi capaz de classificar corretamente todas as instâncias em um teste de aprendizado de máquina supervisionado para os gêneros bacterianos <i>Corynebacterium</i> , <i>Klebsiella</i> e <i>Escherichia</i> . Os resultados obtidos são robustos e consistentes com os de outros estudos de filogenia.
Versatilidade	Possui múltiplas aplicações, podendo ser usado em <i>machine learning</i> , PCA, comparação de sequências, comparação de genomas completos, construção de árvores filogenéticas e análise taxonômica global. É ajustável, permitindo várias projeções e a alteração de k-Mers (máscaras) para se adequar aos dados a serem minerados.
Integração com MT	Combinado com o BioTEX, permite o manuseio de grandes quantidades de texto de forma rápida e robusta, possibilitando o uso do arsenal de algoritmos de bioinformática para análise de textos.

Fonte: Adaptado de (De Pierri et al., 2020)

- **Identidade de Sequência:** o alinhamento pode ser problemático quando a identidade de sequência é baixa.
- **Coerência Filogenética:** em um estudo comparativo com 41 mitogenomas, a árvore filogenética gerada pelo SWeeP (via *Neighbor-Joining*) foi considerada mais coerente do que a gerada pelo Clustal Omega.

Comparação com Outros Métodos Livres de Alinhamento

Embora existam outros métodos livres de alinhamento para análise comparativa de genomas completos, como Prot-SpaM, Kmacs e BioVec, o SWeeP se diferencia em diversos aspectos:

- **Velocidade:** o SWeeP demonstrou ser significativamente mais rápido.

Em um conjunto de dados mitocondriais, o SWeeP foi até 100 vezes mais rápido do que o Kmacs e cerca de 10 vezes mais rápido do que o BioVec e o Prot-SpaM;

- Redução de Dimensionalidade: O SWeeP permite reduzir a dimensionalidade de seus vetores, o que constitui um diferencial importante, já que as ferramentas mencionadas (Kmacs, BioVec, Prot-SpaM) não dispõem desse recurso. A redução de dimensionalidade é essencial para viabilizar análises, como a PCA, em grandes conjuntos de dados.

Em resumo, a abordagem do SWeeP para representar sequências biológicas em espaços vetoriais compactos amplia os recursos matemáticos e computacionais disponíveis para a análise e mineração de sequências biológicas e de textos científicos, tornando esses processos mais eficientes e acessíveis.

2.6.4 Método de mineração de texto com algoritmos de bioinformática

A MT é uma técnica que permite extrair informações úteis de grandes volumes de texto. Quando aplicada à bioinformática, ela desempenha um papel importante na análise de dados genômicos, proteômicos e de outras informações biológicas. A AM, por sua vez, é uma abordagem computacional que permite aos sistemas aprenderem com os dados e melhorarem seu desempenho ao longo do tempo (Cohen K. Bretonnel; Hunter, 2013).

A combinação dessas duas técnicas pode trazer avanços significativos na análise de dados biomoleculares. Analogamente ao que ocorre com os textos de modo geral, ele vale para os dados biológicos. Com a grande quantidade de informações geradas pela evolução das técnicas laboratoriais, surgiram muitas ferramentas e estratégias computacionais para tratamento e conformação à bioinformática (Machado et al., 2022). Na bioinformática, comumente empregam-se dados biológicos representados por cadeias de caracteres denominadas arquivos FASTA. As sequências biológicas são representadas nesses arquivos por um conjunto de letras (Ayyildiz; Piazza, 2019).

Um conjunto de caracteres representa sequências nos arquivos FASTA, traduzindo-as em informação, de modo semelhante ao que ocorre em textos escritos em linguagem natural, que apenas representam dados distintos. Dentro dessa

perspectiva, é interessante identificar uma estratégia que utilize textos em linguagem natural, codificados em formato baseado na representação FASTA. Assim, determinados métodos de bioinformática tornam-se aplicáveis a textos escritos em linguagem natural (Machado et al., 2022).

2.7 ESTADO DA ARTE DA RECUPERAÇÃO DE INFORMAÇÃO MUSICAL EM LETRAS DE MÚSICAS

O estado da arte da recuperação de informação das músicas, especificamente em suas letras, conforme explorado em artigos selecionados, evidencia um campo de pesquisa em rápida evolução, marcado pela adoção de uma variedade de metodologias de AM e de Processamento de Linguagem Natural (PLN). Esta seção resume as principais aplicações, avanços e tendências nessa área, destacando as contribuições de cada estudo.

A análise e a RIM aplicadas às letras das músicas constituem um importante campo de pesquisa interdisciplinar que abrange a música, a linguística, a ciência da computação e a psicologia. As letras, como forma de expressão artística, contêm ricas camadas de significado e emoção, constituindo um campo fértil para análise por meio de técnicas avançadas de AM e PLN.

As principais aplicações e métodos identificados nas publicações selecionadas, conforme a seção 1.3, para este estudo, são classificados e apresentados no Quadro 3, a seguir.

2.7.1 Aplicações atuais na recuperação de informação textual em músicas

Encontram-se na literatura várias aplicações de recuperação de informação musical em letras de músicas, e as mais importantes e comuns são:

- Classificação de gênero e estilo: o uso de letras pode ser útil para a classificação de gênero. Embora pareça que os áudios das músicas sejam mais úteis do que as letras para essa tarefa, devido à alta dimensionalidade dos dados de áudio e à baixa dimensionalidade das letras, é possível classificar músicas em gêneros apenas com base em suas letras. Nesta linha de pesquisa, o tema gênero e estilo corresponde a 14% dos trabalhos, de acordo com a seleção das publicações contidas no QUADRO 2 para o tema estado da arte (BOGHRATI; BERGER, 2023; DA SILVA; SILVA; MARCACINI, 2020);

Quadro 3 - Recuperação de informação musical em letras de músicas

Publicação	Autores	Ano	Aplicação	Método / Ferramenta
Quantifying Cultural Change: Gender Bias in Music	Baghdati, Berger	2023	Analisando o viés de gênero na música ao longo do tempo	PLN Processamento de Linguagem Natural e AM
Tracking Emotions from Song Lyrics	Jo, Kim	2023	Analisando emoções em letras de K-pop	MINERAÇÃO DE TEXTO análise de frequência de morfemas, modelagem de tópicos estruturais
Data Science Approach to Compare the Lyrics	Rosebaugh , Shamir	2022	Análise dos estilos, legibilidade e sentimento dos compositores nas letras	Análise quantitativa de texto, ML básico supervisionado
More Than Words	Preniqi et al.	2022	Ligar preferências musicais a valores morais	Análise de texto, abordagens de regressão
Sentiment Classification of English and Hindi Music Lyrics	Sumith et al.	2022	Sentiment classification in English and Hindi lyrics	AM Supervisionada (Random Forest, Naive Bayes, SVM, AdaBoost)
Music TM-Dataset for Joint Representation	Zeng et al.	2021	Aprendizagem de representação para recuperação intermodal	Recuperação intermodal em música
Using ML Analysis to Interpret Music Emotion and Lyric Features	Xu et al.	2021	Relação entre letra e emoção musical	AM supervisionada Random Forest LIWC, regressão
4MuLA	Silva et al.	2020	Classificação de gêneros, similaridade musical e de artistas, popularidade	Regressão multitarefa, multimodal, multilíngue
Familiar Feelings	Lloyd, Casey	2020	Reconhecimento de Emoções Musicais (MER)	Randon Forest, análise de importância de recursos
Conditional LSTM-GAN for Melody Generation	Yi Yu et al.	2019	Geração de melodia a partir de letras	LSTM-GAN (Memória de Curto Prazo - Rede Adversarial Gerativa)
L,M&A	Saluja et al.	2019	Algoritmos de recomendação de músicas, análise de sentimento de letras	Mineração de texto Tidytext
Opinion Mining and Classification of Music Lyrics	Ahuja, Sangal	2018	Análise de sentimento em letras de músicas em inglês	AM supervisionado, marcação de PDV, WorldNet
Music Mood Classification	Chauhan et al.	2016	Deteção de humor em músicas hindi	LDA Alocação latente de Dirichlet, unigrama, termo frequência
Unsupervised Tagging of Spanish Lyrics	Parra, León	2013	Agrupando músicas com emoções semelhantes	Clustering, Bag Of Words (BOW), técnicas de PDV
Lyrics-based Emotion Classification	Kim, Kwon	2011	Classificação de emoções a partir de letras	Seleção de características por análise sintática parcial

Fonte Autor (2023)

- Geração de letras de música: algoritmos de geração de texto podem ser treinados com base em letras de músicas existentes para gerar letras originais em um estilo específico. Isso tem aplicações na composição musical e na publicidade, e essa aplicação corresponde a 3% da seleção;
- Análise de emoção, humor e sentimento: técnicas de aprendizado de máquina permitem determinar a polaridade emocional de cada letra (positiva, negativa ou neutra). Isso é útil para entender o conteúdo emocional das músicas e até mesmo prever a recepção do público (Brattico et al., 2011; Chauhan; Chauhan, 2017; Ekman, 2003 ; Jo; Kim, 2022 ; Kim; Kwon, 2011; May; Casey,2020; Saluja et al., 2019 ; Da Silva et al.,2020 e Srivastava et al.,2022; Sumith et al. ,2022; Xu et al. 2021). A maioria dos trabalhos concentra-se nesta aplicação, que compreende 79% da seleção e
- Outros temas, como a geração de melodias e a análise linguística, com verificação da semântica das letras, permitem uma compreensão mais aprofundada de sua diversidade temática, correspondendo a 3% da seleção.

2.7.2 Algoritmos de aprendizagem de máquina para análise textual

Na análise de RIM em letras de músicas, o uso de métodos de AM e MT, bem como de vários algoritmos e técnicas, pode ser empregado. A escolha dos algoritmos depende dos objetivos específicos da análise e das características dos dados disponíveis. A partir da pesquisa das publicações sobre o tema do estudo, foi possível verificar, além das aplicações, técnicas e metodologias recentes, a identificação dos trabalhos atuais que estão relacionados:

- a) Classificação de Texto: algoritmos de classificação, como Nãive Bayes (NB), Máquinas de Vetores de Suporte (MVS), Florestas Aleatórias (FA) e Redes Neurais (RN), podem ser usados para classificar letras de músicas em diferentes categorias, como gêneros musicais ou estilos (Ahuja; Sangal 2018; Kim; Kwon 2011; May; Casey 2020; e Sumith et al. 2022; Xu et al. 2021);

- b) **Análise de Sentimento:** algoritmos de análise de sentimento, incluindo classificação binária (positiva/negativa), análise de emoções discretas ou de polaridade contínua, podem ser aplicados para determinar o sentimento das letras. Algoritmos como o VADER (*Valence Aware Dictionary and sEntiment Reasoner*) ou modelos de aprendizado profundo são comuns para essa tarefa (Hutto; Gilbert, 2014);
- c) **Agrupamento de Texto:** algoritmos como o K-Means ou DBSCAN podem ser usados para agrupar letras semelhantes com base em características linguísticas ou temáticas, permitindo a identificação de padrões e tendências (De Pierri et al., 2020; Hotho; Nürnberger; Paaß, 2005; Parra; León, 2013);
- d) **Redes Neurais Recorrentes e *Short-Term Memory*:** para tarefas que envolvem sequências de palavras, as *RNN* e as *Long Short-Term Memory* (LSTM) são utilizadas para capturar dependências temporais nas letras (Srivastava et al., 2022);
- e) **Análise de Tópicos:** algoritmos como *Latent Dirichlet Allocation* (LDA) ou *Non-Negative Matrix Factorization* (NMF) podem ser aplicados para identificar tópicos predominantes em conjuntos de letras de músicas (Chauhan; Chauhan, 2017);
- f) **Redes Neurais Convolucionais (CNN):** quando as letras de músicas são tratadas como imagens de texto, as CNN podem ser usadas para identificar padrões visuais, especialmente quando a formatação visual é relevante (Blaszke; Kostek, 2022);
- g) **Processamento de Linguagem Natural (PLN):** campo da IA que visa compreender, interpretar e manipular a linguagem humana por meio de máquinas. Um aspecto importante do PLN é a representação de palavras e frases de modo que os computadores possam processá-las de forma eficiente. Nesse contexto, técnicas de *word embedding*, como BERT, FastText, GloVe e Word2Vec, desempenham um papel fundamental por serem modelos pré-treinados. Essas técnicas podem ser usadas em tarefas de análise textual em letras de músicas, como geração de texto, tradução automática e resumo (Doh et al., 2022) e
- h) **Aprendizado Profundo (*Deep Learning*):** redes neurais profundas, incluindo modelos de atenção e *Transformers*, têm se mostrado eficazes em tarefas

complexas de processamento de texto, como tradução automática, geração de texto e análise de sentimentos (Jia, 2022).

A escolha do algoritmo depende das necessidades específicas do projeto de análise textual de letras de músicas e do conjunto de dados disponível. Cada algoritmo apresenta vantagens e limitações, e a seleção adequada deve ocorrer de acordo com os objetivos da pesquisa e com as características das letras de músicas em questão.

2.8 MÉTODOS DE INCORPORAÇÃO DE PALAVRAS

Embora a MT tenha dado passos largos na solução dos problemas que envolvem as questões de classificação de documentos, extração de informações, sumarização automática de textos, análise de tendências e temas, conversão de texto para fala, descoberta de conhecimento e outras aplicações, essa disciplina passa para um novo patamar com recursos da área de processamento de linguagem natural (Gupta et al., 2020).

Cabe reforçar que, para a metodologia proposta, o uso de tecnologias recentes, como word embeddings e a incorporação de palavras, será aplicado neste trabalho, e as abordagens descritas serão empregadas para a aplicação do conceito.

2.8.1 Modelo Word2Vec

O modelo Word2Vec é projetado para transformar palavras em vetores numéricos de alta dimensão, criando representações que capturam relações semânticas e sintáticas complexas entre elas. O Word2Vec destaca-se por sua notável capacidade de detectar semelhanças e analogias entre palavras, com base em seus contextos de uso. Opera-se principalmente com dois modelos de arquitetura. (Mikolov et al., 2013):

a) CBOW (*Continuous Bag-of-Words*): esse modelo prevê uma palavra-alvo a partir das palavras de contexto ao redor. Por exemplo, dadas as palavras ao redor, o modelo tentará prever a palavra faltante em uma frase.

b) *Skip-gram*: funciona de forma inversa ao CBOW, prevendo as palavras de contexto a partir de uma palavra-alvo. Esse modelo é particularmente eficaz para trabalhar com conjuntos de dados menores e para capturar palavras mais raras.

Ambos os modelos são treinados com redes neurais e aprendem os pesos (que se tornam os vetores de palavras), ajustando-se para prever palavras com base em

suas vizinhanças. Uma vez treinado, o Word2Vec mapeia cada palavra para um espaço vetorial, de modo que palavras com contextos semelhantes se localizam próximas umas das outras nesse espaço, refletindo seus significados e relações semânticas (Mikolov et al., 2013).

A principal inovação do Word2Vec é sua eficiência e eficácia na captura de nuances semânticas, tornando-o uma ferramenta fundamental em diversas aplicações de PLN, desde análise de sentimentos até sistemas de recomendação e tradução automática (Mikolov et al., 2013).

2.8.2 Modelo GloVe

O GloVe (*Global Vectors for Word Representation*), desenvolvido por Pennington (2014), é uma técnica para criar vetores de palavras (incorporação) que capturam significados complexos e relações entre palavras com base em sua coocorrência em um grande *corpus* de textos (Pennington; Socher; Manning, 2014).

Características-chave do GloVe:

a) Baseado na Coocorrência de Palavras:

- Diferentemente do Word2Vec, que usa contexto local das palavras (janelas de palavras adjacentes), o GloVe foca na coocorrência global de pares de palavras em todo o *corpus*.
- O modelo constrói uma grande matriz que registra quantas vezes cada par de palavras coocorre no *corpus*, capturando as frequências de ocorrência de cada palavra.

b) Representação Vetorial:

- O GloVe utiliza métodos de decomposição de matrizes para transformar a matriz de coocorrência em vetores de palavras.
- Cada palavra é representada por um vetor de alta dimensão que encapsula seu significado e uso, com base na frequência com que aparece em conjunto com outras palavras.

c) Captura de Relações Semânticas e Sintáticas

- Os vetores resultantes representam palavras, de modo que as relações semânticas e sintáticas se refletem na proximidade geométrica entre eles.

- Isso significa que palavras com significados semelhantes ou usos relacionados tendem a estar mais próximas umas das outras no espaço vetorial.

d) Aplicações Diversas:

- O GloVe é utilizado em várias tarefas de PLN, como análise de sentimentos, tradução automática, identificação de entidades nomeadas e outras aplicações.
- Ele é eficaz em captar nuances de significado e em relacionar palavras de forma complexa, tornando-o valioso para sistemas que dependem de uma compreensão profunda da linguagem.

e) Eficiência e Escalabilidade:

- O modelo é eficiente em termos computacionais e escalável para grandes conjuntos de dados, o que é importante para trabalhar com grandes *corpora* de texto.

Em resumo, GloVe é uma abordagem inovadora para a criação de incorporação de palavras, diferenciando-se por sua ênfase na análise global de coocorrência de palavras em larga escala, proporcionando uma rica captura de relações semânticas e sintáticas em dados textuais (Pennington et al., 2014).

2.8.3 Modelo BERT

O BERT (*Bidirectional Encoder Representations from Transformers*), desenvolvido por Devlin (2019), é um modelo que revolucionou a forma como as máquinas de inteligência artificial compreendem a linguagem humana. As suas principais características e funções são (Devlin et al., 2019):

a) Bidirecionalidade

- Ao contrário de modelos anteriores, que analisavam o texto linearmente (da esquerda para a direita ou vice-versa), o BERT processa o texto de forma bidirecional. Isso significa que ele considera o contexto completo de uma palavra em uma frase, analisando as palavras antes e depois dela.

b) Baseado em *Transformers*

- O BERT utiliza a arquitetura de *transformers*, especificamente a parte do *encoder*, para processar palavras em contexto. Os *transformers* são

modelos de atenção que focam em diferentes partes de uma sentença para compreender seu significado completo.

c) Pré-treinamento e Ajuste Fino

- O modelo é pré-treinado em dois tipos de tarefas não supervisionadas: previsão de linguagem mascarada (*Masked Language Model*) e previsão de sentenças próximas (*Next Sentence Prediction*). Isso permite que ele desenvolva uma compreensão profunda da linguagem.
- Após o pré-treinamento, o BERT pode ser ajustado finamente com dados adicionais para tarefas específicas de PLN, como a classificação de texto, a análise de sentimentos, o reconhecimento de entidades nomeadas, entre outras.

d) Generalização e Desempenho

- O BERT alcançou desempenho de ponta em diversas tarefas de PLN, demonstrando uma compreensão excepcional das nuances linguísticas.
- Devido ao seu pré-treinamento extensivo, o BERT apresenta forte capacidade de generalização para diferentes aplicações, exigindo menor quantidade de dados durante o ajuste fino em tarefas específicas.

e) Impacto

- Desde seu lançamento, o BERT tornou-se um dos modelos de referência em PLN, influenciando o desenvolvimento de inúmeros outros modelos baseados em *transformers*.

Em resumo, o BERT representa um grande avanço na PLN, oferecendo uma abordagem poderosa e flexível para compreender a linguagem humana, sobretudo por sua capacidade de processamento bidirecional e pelo uso eficiente da arquitetura de *transformers*. É importante reforçar que as aplicações mais comuns e frequentemente exploradas para subsidiar sistemas de recomendação e rotulação continuam sendo objeto de investigação, sobretudo em razão da ausência de métodos computacionais de baixo custo para a recuperação de informação, e, no contexto musical, essa realidade não é diferente.

Na seção 3, discorre-se sobre o método proposto de Análise Textual de Músicas Brasileiras (ATMBR), baseado nas tecnologias de MT e AM, e nos recursos recentes de processamento de linguagem natural apresentados nesta seção.

3 METODOLOGIA PARA ANÁLISE TEXTUAL DE MÚSICAS

Nesta seção, descrevem-se o método proposto para a análise textual de músicas brasileiras e a forma como se busca solucionar o problema de rotulação das músicas a partir de suas letras e dos diferentes níveis de informação semântica nelas presentes. Inicialmente, apresenta-se uma visão geral do método de Análise Textual da Música Brasileira (ATMBr). Na subseção seguinte, serão descritos a origem das fontes de informação dos corpora musicais e, em seguida, os processos de pré-processamento, cura e normalização das informações. Em continuidade, aborda-se o processo de extração de *embeddings* das letras de músicas, o pré-treinamento dos modelos de linguagem, a análise dos mapas de correlação e a geração dos modelos destinados à rotulação automática. Na sequência, apresentam-se as métricas utilizadas para a avaliação e validação dos resultados. Ao final, apresentam-se as considerações conclusivas do método proposto.

3.1 VISÃO GERAL DO MÉTODO ATMBr

Conforme o problema descrito no capítulo inicial do presente estudo, o método proposto empregará uma solução de MT com aprendizagem de máquina e a aplicação de métodos de processamento de linguagem natural, combinados com o uso de algoritmos de bioinformática. Ao final, serão gerados modelos de classificadores supervisionados para analisar, classificar e rotular músicas com base nas características extraídas das letras das músicas vetorizadas.

O método ATM tem como base quatro macroprocessos (Figura 10), que abordam os seguintes aspectos:

- a) pré-processamento, cura e normalização do *corpora* musical;
- b) vetorização do *corpora* musical aplicando o algoritmo SWeeP;
- c) geração de modelos de aprendizagem de máquina; e
- d) rotulação aplicando classificadores supervisionados

Serão detalhados na continuidade, com as particularidades de cada módulo e os processos que a compõem. O método foi concebido para processar um *corpus* de letras de músicas, contendo, basicamente, um código de referência, título, compositor e/ou artista, ano de lançamento e *links* para a música e a letra. Ao final do processo,

gera-se um dendrograma que permite analisar todo o *corpus* com base nos termos mais significativos e frequentes.

Figura 10 - Diagrama do funcionamento do método ATM

Fonte: O autor (2025).

Para o experimento-piloto, foi inicialmente obtido o *corpora* de letras de músicas brasileiras, com a classificação de gêneros, proveniente do sítio de serviço de músicas Vagalume¹. Após baixar os dados de um repositório, foram realizadas algumas operações importantes antes da fase de pré-processamento. O *corpora* musical passou por uma etapa intensa de cura, ou seja, uma limpeza prévia com o intuito de

remover títulos repetidos, bem como aqueles sem dados completos ou inconsistentes. Conforme apresentado no Quadro 4, verificam-se os quantitativos do *corpora* utilizado.

Quadro 4 – Informações do conjunto de dados extraído do *streaming* Vagalume

Descrição	Quantitativo
<i>Corpus</i> musical quantidade	138.850
<i>Corpus</i> pré-processados (curado)	80.490
<i>Corpus</i> utilizado para o experimento piloto	5.997
Relações de palavras geradas pelo <i>corpora</i>	7.676

Fonte: O Autor (2025).

O *corpora* das letras de músicas contém os seguintes metadados (Quadro 5), que serão utilizados como entradas para o método, sendo o título e a letra da música as principais fontes no processo.

Quadro 5 – Descrição dos metadados utilizados

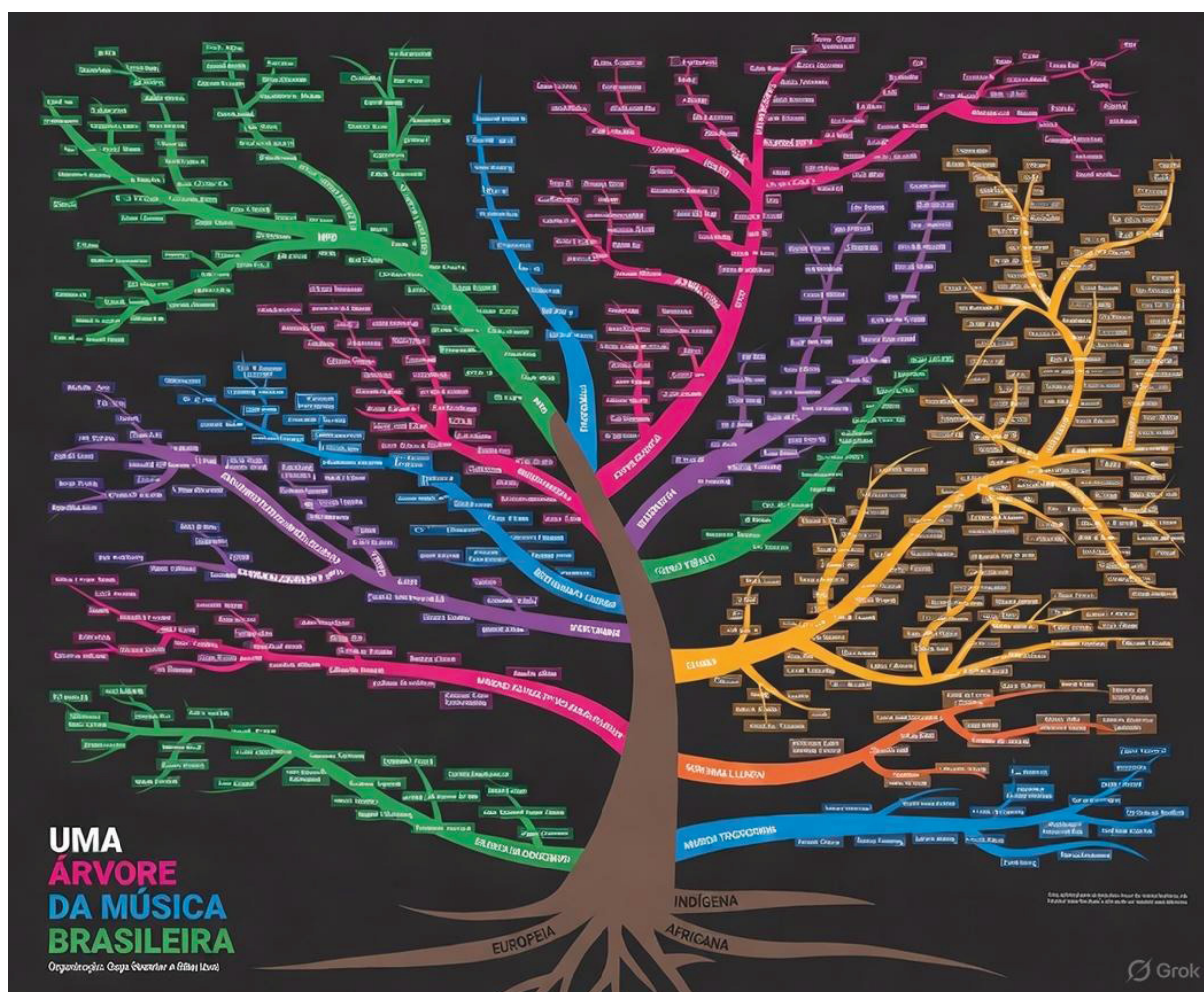
Descrição do Metadados	característica
Código de referência	Numérico
Título da música	Literal
Compositor / artista	Literal
Gênero musical	Literal
Ano de lançamento da música	Numérico
Link da letra da música para visualização	Literal
Letra da música	Literal

Fonte: O Autor (2025).

É importante pontuar que o *corpora* das letras de músicas apresenta a seguinte distribuição entre catorze gêneros musicais brasileiros (ver Quadro 5). Os gêneros descritos são específicos e particulares à origem dos dados, o sítio Vagalume. A abordagem e o estudo dos gêneros musicais brasileiros são temas amplamente

discutidos na literatura, como evidencia Caldas (2010). Internacionalmente, também se configura como objeto de intensa pesquisa, principalmente devido à sua aplicação em sistemas de recomendação. Outra visão que retrata, em sua gênese, o gênero musical brasileiro e suas raízes é apresentada pelos autores Mori e Stroeter (2020) em sua obra “Uma Árvore da Música Brasileira” (Figura 11).

Figura 11 – Árvore dos gêneros musicais brasileiros



Fonte: Mori e Stroeter (2020)

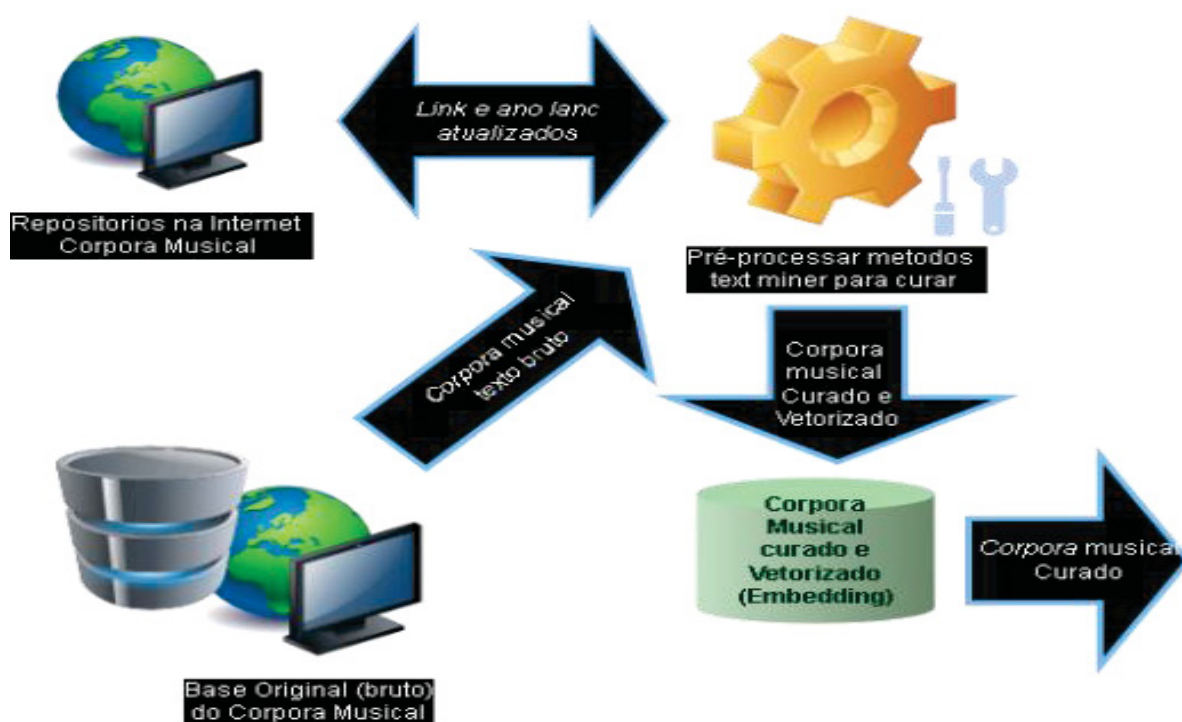
3.1.2 Pré-processamento para a cura¹ e normalização do *corpora* musical

A redução do número de termos de um *corpora* textual em MT é uma exigência não apenas na representação vetorial, mas também para qualquer outro tipo de representação de texto. A importância do método de pré-processamento se fundamentou na necessidade de robustez, para que fosse capaz de lidar com letras

¹ Cura - Significa o cuidado curatorial completo e contínuo que transforma dados digitais brutos de alta qualidade, confiáveis e reutilizáveis por outros pesquisadores ou sistemas no futuro.

de músicas e textos com alto índice de ruído, contendo erros gramaticais e de digitação. O pré-processamento consumiu cerca de 70% do tempo total de aplicação do método. Embora sejam processos simples, a grande quantidade de títulos musicais gerou um volume expressivo de ruídos que se multiplicaram ao longo das etapas. É importante lembrar que o *corpora* musical proveniente do sítio Vagalume tem a proposta de ser alimentado por usuários do sítio, o que, na literatura, é definido como método de colaboração de informações (*crowdsourcing*). Esse fato acabou por gerar uma série de ruídos, dificultando a captura precisa da letra da música. O aprofundamento deste processo, bem como sua relação com ontologias e a representação do conhecimento, resultou na publicação de Florido, De Paula Pinto e Raittz (2025).

Figura 12 – Pré-processamento e cura do *corpora* musical



Fonte: O Autor (2025).

Os processos aplicados (Figura 12) usualmente foram os descritos na seção 2.4.2 do referencial teórico: *case folding*, *stemming*, *stopwords* e representação de N-gramas com sentimento positivo e negativo. Cabe reforçar, nesses passos, as principais dificuldades encontradas no pré-processamento do *corpora*:

- O crescimento rápido do número de músicas e letras faz com que títulos se repitam frequentemente, mesmo quando as letras são diferentes, sendo necessário verificar, além do título, se a letra também é repetida;
- No meio artístico, é muito comum vários intérpretes regravarem composições famosas e consagradas que não sejam de sua autoria, gerando várias repetições da mesma letra;
- Na seleção do *corpus* musical deste estudo, constatou-se a presença de termos em inglês e em espanhol, em alguns casos, na totalidade da obra. Também apareceram palavras em italiano e até em hebraico, especialmente em letras do gênero gospel; e
- Por algum problema na origem dos dados, foram encontrados no *corpora* vários títulos de artistas não brasileiros. Como o cunho do trabalho é o estudo da música popular brasileira, foram descartadas as letras de artistas e/ou compositores em outras línguas, por fora do escopo.

3.1.1 Aplicação do algoritmo de bioinformática SWeeP

As ferramentas e algoritmos de bioinformática disponíveis não foram desenvolvidos para processar textos em linguagem natural. Para utilizá-los, é necessário empregar o *Biological Sequence Format* (BSF), o que exige a criação de um *parser*¹ capaz de converter textos em linguagem natural para esse formato. Isso permite que o material textual seja processado por diversas ferramentas de bioinformática, incluindo as voltadas à comparação de sequências, aliadas à agilidade e à robustez das tecnologias de vetorização.

Para isso, utiliza-se um pacote que implementa estratégias de MT usando ferramentas de bioinformática e que fornece recursos de AMINOcode para codificar texto em BSF (Machado et al., 2022). Esse pacote oferece duas funções principais para converter textos convencionais em um formato biológico válido:

¹ *Parser* - ferramenta ou componente de software que analisa dados de entrada para construir uma representação estrutural, geralmente seguindo as regras de uma gramática formal. Em computação e programação, os *parsers* são comumente usados para processar e interpretar textos, transformando-os em uma estrutura de dados mais facilmente manipulável por um programa.

- O AMINOcode consiste em substituir caracteres de texto por letras que representam aminoácidos, conforme uma lista codificada. A Tabela 1, no Anexo I, mostra as regras de substituição de caracteres no AMINOcode (Machado et al., 2022)
- O DNAbits realiza a conversão de caracteres de texto dividindo cada byte em quatro pares de bits. O *coding* mantém as informações de acordo com o Código Padrão Americano para Intercâmbio de Informações (ASCII), substituindo cada par de bits por A, C, G e T, conforme a regra predefinida na Tabela 1 no ANEXO C, (Machado et al., 2022).

Como exemplo, no título “recuperação de informação musical”, o primeiro passo consiste em realizar o pré-processamento, aplicando as operações básicas da MT: *case folding*, *stemming*, remoção de *stopwords* e lematização. Em seguida, os termos resultantes são convertidos com o AMINOcode para possibilitar o processamento por algoritmos de bioinformática. No Quadro 5, apresentam-se as conversões realizadas, que culminam na cadeia de aminoácidos, a serem posteriormente processadas pelo algoritmo. A aplicação do AMINOcode resultou na seguinte sequência de aminoácidos (Machado et al., 2022).

“MYVSYICYSYINFYORMYATYIYONYSRYETRYIYEVYAL”.

Quadro 5 – Exemplo de Conversão de termos para pesquisa

Passos	Termo 1					Termo 2										Termo 3											
1.Texto bruto	m	u	s	i	c		i	n	f	o	r	m	a	t	i	o	n		r	e	t	r	i	e	v	a	l
2.Pré-processado	m	u	s	i	c	spc	i	n	f	o	r	m	a	t	i	o	n	spc	r	e	t	r	i	e	v	a	l
3.AMINOcode	m	vv	s	vi	c	vs	vi	n	f	vo	r	m	va	t	vi	vo	n	vs	r	ve	t	r	vi	ve	v	va	l

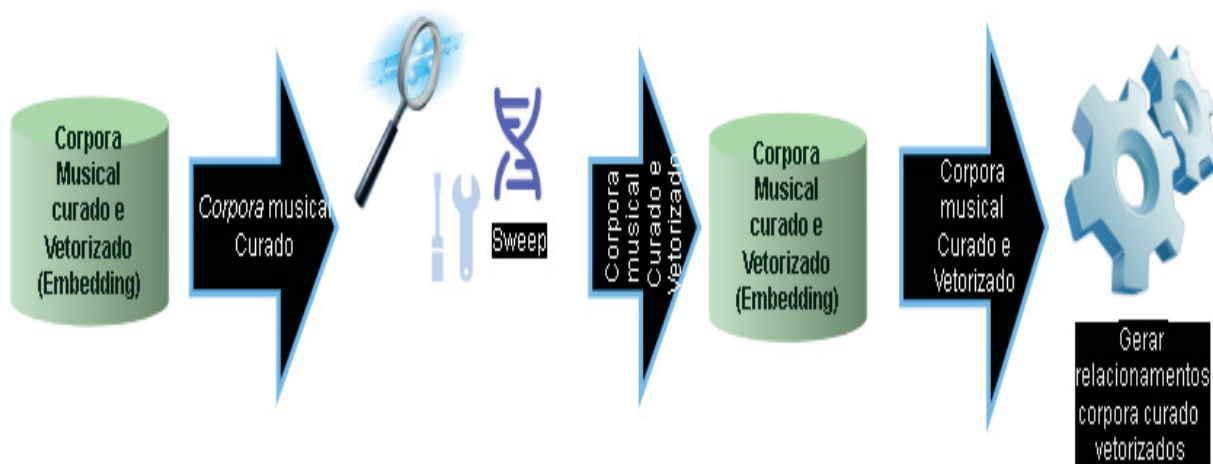
Fonte: Machado et al. (2022)

Na continuação, os cabeçalhos das sequências foram convertidos para o formato BSF com o AMINOcode, resultando em um arquivo FASTA com as informações textuais codificadas. Para as análises de conteúdo desses cabeçalhos, as sequências foram posteriormente convertidas em vetores pela função Sweep (De Pierri et al., 2020). Em seguida, os grupos de cabeçalhos foram definidos por meio de métodos de agrupamento aglomerativo e, finalmente, as sequências pertencentes a cada *cluster* foram alinhadas separadamente utilizando o *Ômega Clustal* (Sievers; Higgins, 2018). Após o alinhamento, a sequência de consenso foi obtida e decodificada a partir do aminoácido correspondente.

De acordo com o diagrama apresentado na Figura 11, o ponto fundamental do método é a aplicação do algoritmo SWeeP, responsável por gerar o *corpora* musical vetorizado (*embeddings*) das letras de músicas. Mais detalhadamente, conforme ilustrado na Figura 13, os seguintes passos são realizados:

- Após o pré-processamento do *corpora* musical, o *corpora* curado, além de ter o conteúdo das letras tokenizado, é submetido à técnica TF-IDF, para medir a importância de um termo em um documento em relação a uma coleção de documentos (Salton; McGill, 1987). Esses valores são então incorporados ao *corpora curado*.
- Em seguida, ainda utilizando o *corpora* curado, as letras e caracteres de cada palavra são convertidos por uma função que produz o formato BSF/FASTA. Cada caractere é convertido em um tipo de aminoácido; ver exemplo no Quadro 4, resultando em um grande vetor de aminoácidos associado ao corpus;
- Na sequência, este grande vetor de aminoácidos é transformado em números também e submetido ao algoritmo SWeeP para processar e gerar um *corpora* musical vetorizado denominado (CMV);
- Após essa fase, é gerada uma base das relações entre o modelo de linguagem do português do Brasil e entre todas as palavras que agora chamamos de termos, que retornam ao seu formato original.

Figura 13 - Diagrama da ATMBR recorte processo SWeeP



Fonte: O Autor (2025)

- a) Relações das músicas com outras, a partir dos termos mais significativos de uma para com as outras (Figura 14), pela contextualização dos termos em todo o *corpora*;
- b) Relação dos termos e as 20 músicas mais similares (Figura 15) com o termo;

Figura 15 – Consulta das 20 músicas mais relacionadas com o termo “saudade”

SAUDADE - Documentos relacionados ao Word

consulta de pesquisa

Número dos vizinhos

Procurar

Ajuda

Total de inscrições: 20 Entradas selecionadas: 20

Exportar CSV

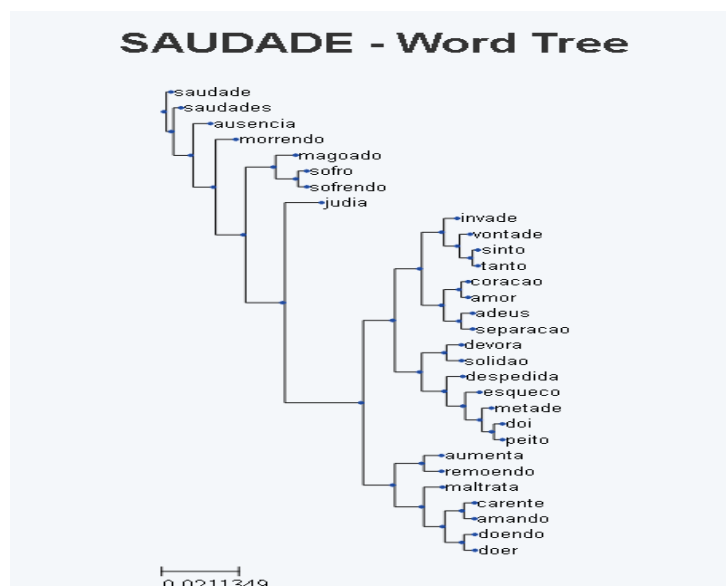
Limpar filtros

#	Classificação	Semelhança	Título + Letra	Artista	Gênero	Ano	Link	EU IA
0	1	2	3	4	5	6	7	8
71276	0	0,9964	Ta Pintando Uma Saudade cadecadecade voce morre de te ver cadecadecade voce morrendo de de te ver te lembrar me faz sofrer te esquecer eu nao consigo viver por viver por voce tao esquecido com os olhos rasos dagua sufocando a minha magoa nesse peito dolorido cadecadecade voce morrendo de de te ver cadecadecade voce morrendo de te ver chora viola sentida encostada no meu peito consolando a minha vida a desse amor desfeito solidao e desatino por ai sem ter destino voce vivendo desse jeito cadecadecade voce morrendo de te ver cadecadecade voce morrendo de de te ver cadecadecade voce morrendo de de te ver cadecadecade voce morrendo de de te ver	Otavio Augusto e Gabriel	sertanejo	1990	Link	1127641
60022	1	0,9956	Nao Posso Acreditar o tempo passou e em mim tão restou a dorsaudade as lembranças desse amor que o tempo não levousaudade voce deixou marcas no meu coração se foi tudo um sonhoou então desilusao volta e me que seu amor ainda e minha volta e me que voce nao me esqueceu que ainda me ainda me suuhzaninhaaa	Raça Negra	pagode	1996	Link	1086316
60021	2	0,9956	Pensando Em Você pensandouuh pensandoaahh ja não sinto o teu cheiro e o teu colo meu amor lembra que foi verdadeiro que pena tudo acabou as tardes ficaram vazias e a noite que eu ia te ver lembra aquela chuva fina que saudade de voce vou lembrar voce e o seu corpo inteiro que vontade de te ver aonde você está pensandouuh pensandoaahh	Estakazero	forro	2007	Link	10/12/2017
66119	3	0,9955	Lembrança de Você A Lembrança de voce doi doi a saudade no peito volte logo amor nao me deixe sofrer nao me faça chorar coracao no peito batendo sem jeito danos por alguem que nao quer voltar voce foi embora nao me disse nada nao me deu adeus deixando eu sofrer de tristeza e saudade foi pura maldade o que voce me fez	Trio Nordestino	forro	1984	Link	1018546
62812	4	0,9955	A Saudade Bateu Valeu a saudade bateu doeu a saudade bateu a saudade bateu doeu o amor e como um passarinho gosta de voltar ao velho ninho coração que não se cansa guarda na saudade a esperança	Angela Maria	guarda	2002	Link	1135505

Fonte: O Autor (2025)

- c) Dendrograma (Figura 16) das relações dos termos e suas respectivas distâncias dentro de cada título musical;

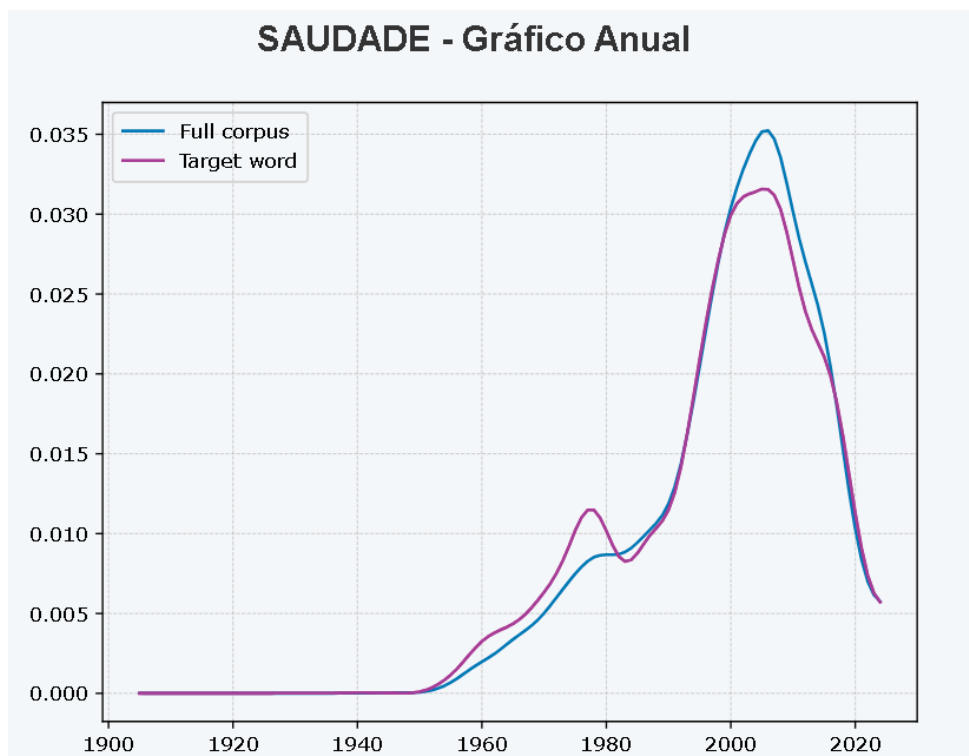
Figura 16 – Consulta do termo “Saudade” e as distâncias entre os termos



Fonte: O Autor (2025)

d) Gráfico das frequências temporais dos termos (Figura 17).

Figura 17 - Frequência do termo “saudades” na linha do tempo



Fonte: O Autor (2025)

O entendimento e o uso das informações musicais recuperadas ficarão mais claros com os experimentos realizados e descritos na seção 4.

3.1.2 Geração de modelos do *corpora* processado

Nesta fase do método, foi possível gerar modelos com as informações disponibilizadas e com aplicação de métodos e algoritmos de AM supervisionada. Para a geração de modelos de classificação supervisionada, optou-se por classificadores específicos e amplamente utilizados para esse fim. Com base nos *embeddings* gerados a partir dos *corpora* musicais, realizou-se um estudo comparativo para identificar o classificador supervisionado com maior acurácia média. Os classificadores selecionados foram: MLP (redes neurais densas), SVM (máquinas de suporte vectorial) e *Random Forest* (florestas aleatórias).

O critério estabelecido no método ATMBR para seleção das amostras das letras de músicas, com o fim de gerar os modelos, baseou-se em dois aspectos principais:

- a) a seleção do gênero musical, por fazer parte do *corpora* desde a origem dos dados, embora não seja um dado unânime, o gênero musical é uma premissa necessária para a análise por um classificador supervisionado; e
- b) seleção do termo e seu enquadramento em determinados padrões, como emoções e sentimentos, permitindo também sua aplicação com os classificadores mencionados anteriormente.

Após a seleção das amostras, recuperou-se o *corpora* antes do pré-processamento, com o propósito de identificar os agrupamentos tratados no item “a”, cuja classe será utilizada para que o classificador de AM supervisionada produza o modelo.

3.1.3 Métricas de avaliação

As métricas de avaliação adotadas para o método são a precisão e a acurácia, descritas no Apêndice E deste estudo.

- a) Precisão (*precision*): é o número de classificações positivas corretas (Verdadeiros Positivos) dividido pelo número de exemplos classificados na mesma classe (a soma dos resultados de Verdadeiros Positivos e Falsos Positivos).
- b) Exatidão / Acurácia: é o número de classificações corretas (a soma de Verdadeiros Positivos e Verdadeiros Negativos) dividido pelo número total de classificações (a soma de todos os itens).

3.1.4 Validação do método ATMBR

A validação do método será realizada a partir de um conjunto de amostras do *corpora* vetorizado (*embeddings*), composto por títulos que contenham termos ou rótulos para validação com classificadores supervisionados. A validação será de duas maneiras:

- Aplicação de classificadores supervisionados MLP para identificação de gêneros musicais e de classificadores de análise de sentimentos com o algoritmo VADER para comparação com o resultado da análise textual gerado pela plataforma.; e
- Elaboração de uma avaliação da plataforma Ritmo Brasil por meio de um questionário aplicado a um público de amostra estratificada, disponibilizando um vídeo para explicação da finalidade da plataforma e instruções de como avaliar a plataforma na sequência, aplica-se o método

ATMBr descrito e apresentam-se os resultados obtidos em função dos experimentos realizados, bem como o processo de validação piloto conduzido do com usuários.

3.2 CONSIDERAÇÕES DO MÉTODO ATMBr

Alguns desafios observados no pré-processamento e no experimento da análise textual de músicas merecem destaque:

- a) Qualidade dos dados: é essencial para o sucesso dos modelos. Letras de músicas frequentemente apresentam erros ortográficos, gírias e linguagem figurada, o que pode ser desafiador para os algoritmos.
- b) Interpretabilidade: modelos de AM, como redes neurais profundas, podem ser difíceis de interpretar, o que pode tornar complicado entender como as decisões são tomadas; e
- c) Viés cultural do gênero musical: alguns gêneros musicais utilizam termos de outras línguas, distintos do português. Embora o objetivo principal seja a música brasileira, o viés artístico dos compositores é influenciado por outras línguas, como o inglês, o espanhol, o francês e o hebraico.

A análise textual das letras de músicas, quando combinada com métodos de aprendizagem de máquina, amplia o conhecimento sobre música, linguagem e cultura. Essa abordagem possibilita a automatização de tarefas de análise, mas também apresenta desafios relacionados à qualidade dos dados e à cultura artística.

Os desafios colocados no processo de pré-processamento do *corpora* musical e a transformação do seu conteúdo contribuíram para uma publicação sobre o tema e para a consequente definição de uma ontologia musical, ver Apêndice A (Florido; De Paula Pinto; Raittz, 2025).

4 EXPERIMENTO DO MÉTODO ATMBR

Com o intuito de validar o método, foi realizado um experimento inicial com o corpus musical em sua totalidade, visando avaliar os resultados obtidos pelo método proposto. O resultado é traduzido pelos dados e informações finais obtidos, bem como pela geração da plataforma “Ritmo Brasil”, cuja importância se estende ao âmbito da pesquisa e do estudo da música.

4.1 Vieses de análises do experimento

Inúmeros são os aspectos proporcionados pela investigação, mas pela exiguidade de tempo cabe relatar e analisar alguns aspectos importantes do âmbito da pesquisa que foi abordado a partir das informações geradas pela plataforma Ritmo Brasil:

- a) verificação da correlação dos termos com as músicas;
- b) aplicação de classificadores supervisionados para identificação dos gêneros musicais;
- c) análise das emoções pelos termos gerados;
- d) dimensão criativa das letras;
- e) memória histórica; e
- f) regionalismo contido nas letras das músicas.

4.2 Ambiente de desenvolvimento computacional

O ambiente de desenvolvimento utilizado para a implantação do método foi o MatLab, um *framework* de programação e desenvolvimento proprietário. Ferramenta de grande capacidade usada para resolver uma ampla gama de problemas numéricos e científicos. É dotada de uma vasta biblioteca para aplicação de algoritmos de mineração de texto, classificadores e estatísticos.

Além desta ferramenta, utilizou-se também a linguagem de programação e o *framework* Python intensamente aplicados no desenvolvimento de algoritmos e sistemas no meio científico. Detém também uma ampla biblioteca de algoritmos de aprendizagem de máquina, redes neurais, classificadores supervisionados e não supervisionados, além de recursos estatísticos.

4.3 Seleção do *corpora* musical para o experimento

A etapa que mais tempo levou antes da realização dos experimentos propriamente ditos citados na seção 4.1 foi o pré-processamento dos dados, com a cura do *corpora* musical bruto.

O *corpora* musical bruto originalmente continha 138.850 títulos. Após o pré-processamento e a curadoria dos dados, a amostra utilizada no experimento foi reduzida para 80.490 títulos, o que representa 58,3% do total inicial, ou seja, uma perda de 41,7%. Essa redução significativa deve-se, principalmente, à origem dos dados: a utilização de uma plataforma de *crowdsourcing*, alimentada por usuários, para o repositório de letras de músicas e para a base de dados. Esse fator trouxe diversos desafios, especialmente relacionados à qualidade e à consistência das informações inseridas pelos próprios usuários, tais como:

- Ambiguidade e Variação: Nomes de artistas inseridos de formas diferentes (ex.: "Chiquinha Gonzaga", "Francisca Edviges Neves Gonzaga");
- Informações Incompletas: Dados como ano de lançamento e autoria frequentemente equivocados ou ausentes;
- Subjetividade: A classificação de gênero musical é subjetiva e varia entre os colaboradores;
- Erros de Linguagem: Presença massiva de erros de ortografia, acentuação e pontuação;
- Ruído Textual: Presença de termos como "Intro", "Refrão", "Solo", cifras e acordes (ex.: "Am", "G", "C") misturados ao texto da letra;
- Duplicatas e Versões: Múltiplas versões da mesma letra com pequenas variações.

A análise da distribuição dos títulos musicais, conforme apresentado na Tabela 1, evidencia a predominância de determinados gêneros no corpus musical curado. O gênero gospel destaca-se como o mais representativo, abrangendo 23,1% do total de títulos, seguido pelo sertanejo, com 18,7%, e pelo MPB, com 13,9%. Juntos, esses três gêneros concentram 55,6% do acervo analisado, o que demonstra sua expressiva relevância no cenário musical brasileiro contemporâneo. Essa concentração sugere não apenas tendências de produção e consumo musical, mas também possíveis

vieses na composição da base de dados utilizada, que devem ser considerados em análises subsequentes.

Tabela 1 – Quantitativos de títulos musicais, termos versus gênero

Gênero	Quantidade títulos	Percentual (%)	Quantitativo de termos
GOSPEL	18.554	23,1%	23.513
SERTANEJO	15.019	18,7%	20.075
MPB	11.188	13,9%	26.252
FORRO	6.761	8,4%	14.296
ROCK	5.307	6,6%	13.469
PAGODE	5.026	6,2%	9.607
SAMBA	4.330	5,4%	13.190
AXE	2.898	3,6%	8.965
INFANTIL	2.608	3,2%	8.823
POP	2.417	3,0%	6.997
BOSSA-NOVA	1.855	2,3%	6.632
VELHA-GUARDA	1.649	2,0%	6.267
FUNK-CARIOCA	1.601	2,0%	5.168
JOVEM-GUARDA	1.278	1,6%	4.049
Todos os GENEROS	80.491	100,0%	75.955

Fonte: O Autor (2025).

É importante destacar que a soma de todos os termos (167.303) por gênero musical não corresponde ao total geral de termos de todos os gêneros (75.955). Isso ocorre porque o processamento foi realizado separadamente para cada gênero, ou seja, cada gênero musical gerou um *corpora* musical de termos específicos, resultando em *corpora* distintos para cada categoria. A partir da definição do *corpora* curado, este foi submetido ao processamento com os algoritmos descritos na seção 3.1.3, incluindo a aplicação do algoritmo de bioinformática SWeeP para a geração dos *corpora* vetorizados (*embeddings*).

4.4 Plataforma de letras de músicas Ritmo Brasil

O principal resultado do método ATMBR é a plataforma Ritmo Brasil, configurada como um ambiente de consulta e pesquisa voltado a músicos e apreciadores de música popular brasileira.

A plataforma está hospedada provisoriamente em um domínio na UFPR/SEPT para validar os conceitos e disponibilizar o acesso aos usuários de música. Como forma de validar a plataforma, foi elaborada uma pesquisa de opinião, cujo resultado é relatado na seção 5.1.

A plataforma Ritmo Brasil (Figura 18) possui uma página inicial com opções de acesso aos vários gêneros musicais, e o primeiro ícone corresponde ao conjunto de todas as músicas. Nessa opção, o tempo de acesso depende do dispositivo utilizado e pode haver certa demora devido ao uso de memória aleatória.

Figura 18 – Plataforma RITMO BRASIL



Fonte: Autor (2025)

Ao escolher a opção de um gênero ou de todos os gêneros, abaixo de cada ícone há duas opções: "Músicas" e "Termos". Na opção "Músicas", abre-se uma tela

com as informações básicas dos títulos, que contém, como diferencial, o link para as 20 músicas mais relacionadas ao título escolhido (Figura 19).

Figura 19 – Tela de consulta opção “Músicas”

Brazilian Lyrics BOSSA NOVA - Texts

Search query: Neighbors number:

[Search](#) [Help](#)

Total entries: 1855 Selected entries: 1855 [Export CSV](#) [Clear Filters](#)

#	Title	Artist	Gender	Year	Related Lyrics	Link	ID
0	1	2	3	4	5	6	7
0	In a Bar	Familia da Musica	bossa-nova	1962	Related Lyrics	Link	1006095
1	Night and Day	Agostinho dos Santos	bossa-nova	1969	Related Lyrics	Link	1004704
2	Yesterday	Agostinho dos Santos	bossa-nova	1967	Related Lyrics	Link	1004779
3	Derniere Valse	Lisa Ono	bossa-nova	2008	Related Lyrics	Link	1006471
4	Brasil Com P	Maria Rita	bossa-nova	2007	Related Lyrics	Link	1006569
5	Elegia Desesperada	Vinicius de Moraes	bossa-nova	1977	Related Lyrics	Link	1007430
6	O Desespero da Piedade	Vinicius de Moraes	bossa-nova	1980	Related Lyrics	Link	1007460
7	Rosa de Hiroshima	Vinicius de Moraes	bossa-nova	1972	Related Lyrics	Link	1007492

Fonte: Autor (2025)

Na opção “Termos”, abre-se uma tela com todos os termos presentes em todas as músicas, e a página contém links para as 20 músicas relacionadas ao termo selecionado (Figura 20).

Figura 20 - Tela de consulta opção “Termos”

Letra brasileira de BOSSA NOVA - Palavras

consulta de pesquisa: Número dos vizinhos:

[Procurar](#) [Ajuda](#)

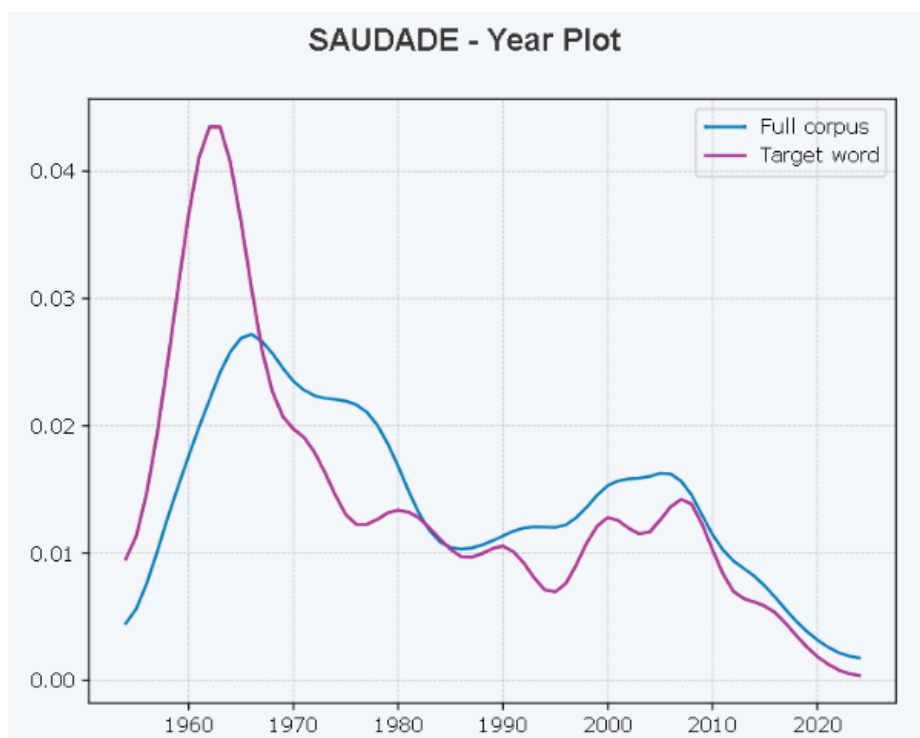
Total de inscrições: 6632 Entradas selecionadas: 6632 [Exportar CSV](#) [Limpar filtros](#)

#	Palavra	Ocasionalmente.	Palavras relacionadas	Letras relacionadas	Árvore de Palavras	Gráfico anual
0	1	2	3	4	5	6
1183	E	1717	e, o, que, de, a, nao, eu, se, meu, me	Letras relacionadas	Árvore de Palavras	Gráfico anual
1182	QUE	1655	que, e, o, de, a, nao, eu, meu, me, se	Letras relacionadas	Árvore de Palavras	Gráfico anual
1180	DE	1588	de, a, o, e, que, nao, se, eu, me, meu	Letras relacionadas	Árvore de Palavras	Gráfico anual
1184	O	1584	o, e, de, a, que, nao, se, eu, meu, me	Letras relacionadas	Árvore de Palavras	Gráfico anual
1181	UM	1572	a, de, o, e, que, nao, se, eu, me, meu	Letras relacionadas	Árvore de Palavras	Gráfico anual
1179	NAO	1281	nao, que, o, e, de, eu, a, se, me, so	Letras relacionadas	Árvore de Palavras	Gráfico anual
1177	UE	1161	eu, meu, que, não, eu, e, o, de, então, a	Letras relacionadas	Árvore de Palavras	Gráfico anual
1161	SE	1047	se, não, a, sem, de, em, o, quem, e, por	Letras relacionadas	Árvore de Palavras	Gráfico anual
5625	UM	992	um, do, não, uma, se, em, com, a, de, mas	Letras relacionadas	Árvore de Palavras	Gráfico anual
5626	FAZER	953	fazer, não, um, em, da, se, na, seu, a, de	Letras relacionadas	Árvore de Palavras	Gráfico anual
1178	MEU	915	meu, eu, me, que, o, nao, e, de, a, so	Letras relacionadas	Árvore de Palavras	Gráfico anual
856	PRA	890	pra, se, não, que, e, de, o, quem, a, então	Letras relacionadas	Árvore de Palavras	Gráfico anual
1176	MEU	870	eu, meu, eu, não, que, o, de, se, e, a	Letras relacionadas	Árvore de Palavras	Gráfico anual
1203	AMOR	867	amor, eu, então, que, amar, meu, tao, e, o,	Letras relacionadas	Árvore de Palavras	Gráfico anual

Fonte: Autor (2025)

Além desses *links*, há também acesso ao dendrograma dos termos mais próximos e ao *link* para os anos (décadas) em que o termo aparece nas letras das músicas (Figura 21).

Figura 21- Incidência da termo Saudade durante as décadas 60 a 2020



Fonte: o autor (2025)

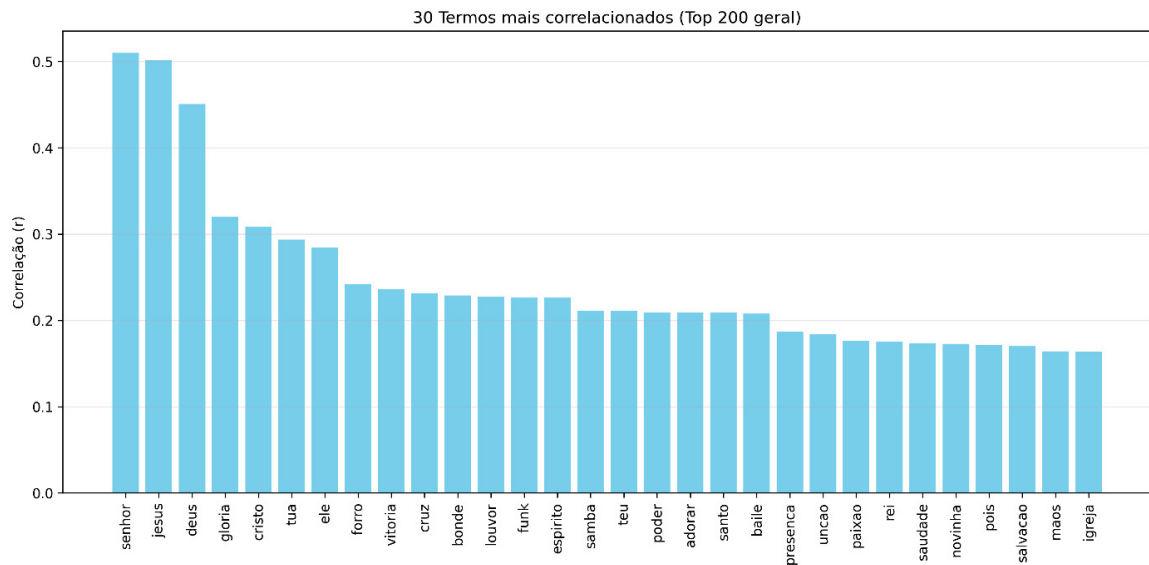
Os detalhes de ordenação das colunas e as formas de consulta nas páginas estão descritos no botão “*help*” da página. Como forma de instruir os usuários, há também um manual de consulta (ver Apêndice F). Mais detalhes sobre o uso e as informações retornadas por cada *link* serão descritos na sequência desta seção.

4.5 Correlação dos termos com as músicas

A partir da totalidade de termos (palavras) de todas as músicas, realizou-se uma análise por meio da correlação de *Pearson*, com o intuito de verificar a incidência e as ocorrências desses termos.

De acordo com a Figura 22, o gráfico de barras evidencia a forte correlação entre termos religiosos e a forte presença do gênero gospel.

Figura 22 – Correlação de termos versus títulos



Fonte: Autor (2025)

4.6 Análise da dimensão emoção no experimento

Para o experimento em questão foi adotado um protocolo de procedimentos de padrão aos termos que exprimem emoções básicas definidas por (Ekman, 2003). Esse protocolo permitirá comparações com outras emoções e com experimentos que vierem a seguir.

Figura 23 – Consulta com o termo “ALEGRIA” em todo o *corpora*.

Letras brasileiras MPB - Palavras						
alegria			Número dos vizinhos			
<div> <div>Procurar</div> <div>Ajuda</div> </div>						
Total de inscrições: 26252 Entradas selecionadas: 48			<div> <div>Exportar CSV</div> <div>Limpar filtros</div> </div>			
#	Palavra	Ocasionalmente.	Palavras relacionadas	Letras relacionadas	Árvore de Palavras	Gráfico anual
0	1	2	3	4	5	6
18619	ALEGRIA	622	alegria, alegre, milagre, fantasia, abre, viva, riso, festa, melodia, negra	Letras relacionadas	Árvore de Palavras	Gráfico anual
18620	ALEGRE	102	alegre, alegria, festa, norte, viva, terra, fonte, povo, milagre, porto	Letras relacionadas	Árvore de Palavras	Gráfico anual
12364	ALEGRIA S	28	alegria s, alegria, tristezas, alegre, tristes, coracoes, ventos, animais, nossas, harmonia	Letras relacionadas	Árvore de Palavras	Gráfico anual
20803	ALEGRAR	27	alegrar, alegria, manha, verao, traz, abre, dia, doce, mar, luz	Letras relacionadas	Árvore de Palavras	Gráfico anual
12365	TRISTEZAS	24	tristezas, tristeza, tristeza, alegria s, tristezas, distante, nossas, estas, vazias, mistérios	Letras relacionadas	Árvore de Palavras	Gráfico anual
20804	ALEGRA	16	alegra, alegria, trazendo, preparando, alegrar, alegre, abre, energia, montanha, pressa	Letras relacionadas	Árvore de Palavras	Gráfico anual
20923	SENSACIONAL	14	sensacional, tradição, carnaval, continente, fantasia, bahia, alegria, folia, feriado, praça	Letras relacionadas	Árvore de Palavras	Gráfico anual
9134	ABENCOADO	13	abençoado, chamada, rainha, maria, bonito, natal, alegria, bahia, tereza, melodia	Letras relacionadas	Árvore de Palavras	Gráfico anual
20894	PALAVRAO	11	palavrao, palhao, levanta, maravilhoso, espalhado, fechado, fantasia, massa, alegria, atento	Letras relacionadas	Árvore de Palavras	Gráfico anual
9885	ROMÊNIA	9	romania, marcha, tomara, alegria, praça, maria, maceio, canta, bonita, palhao	Letras relacionadas	Árvore de Palavras	Gráfico anual
2467	COROS	6	choras, revoadas, choveu, alegria, chuva, vindas, chorando, embalar, chorou, pranto	Letras relacionadas	Árvore de Palavras	Gráfico anual
19024	CONVIVIO	6	convívio, contente, alegria, alcanco, forte, humor, sorte, sinto, mordida, calendário	Letras relacionadas	Árvore de Palavras	Gráfico anual
19077	ARTES	6	artes, alegria s, alegre, união, fotografia, alegria, lares, setenta, meiodia, obra	Letras relacionadas	Árvore de Palavras	Gráfico anual
20899	CAITARA	0	caitara, caitara, caitara, caitara, caitara, caitara, caitara, caitara, caitara, caitara	Letras relacionadas	Árvore de Palavras	Gráfico anual

Fonte: O Autor (2025).

- a) Os termos definidos para o experimento no âmbito de emoções foram basicamente selecionados: alegria, tristeza, medo e raiva;
- b) Para cada termo são geradas uma lista de vinte músicas ou mais relacionadas com o termo; e
- c) Foram considerados apenas os termos similares mencionados no item “a”, como as emoções.

Para o termo “ALEGRIA”, como uma das emoções relacionadas à validação do experimento, verificaram-se na plataforma Ritmo Brasil os seguintes termos relacionados: alegrar, alegrias, alegre, alegrar, alegremente, alegrando, euforia, alegrou e harmonia (ver Figura 23).

Figura 24 - Termos relacionados com o termo principal

#	Palavra	Ocasionalmente.	Palavras relacionadas
0	1	2	3
62166	ALEGRIA	5231	alegria, alegrar, alegrias, alegre, alegrar, alegremente, alegrando, euforia, alegrou, harmonia
62164	ALEGRIA S	152	alegrias, alegria, alegre, alegremente, alegrar, alegrar, alegres, euforia, risos, melancolia

Fonte: O Autor (2025).

No *link* “Letras Relacionadas” estão listadas as músicas associadas ao termo “ALEGRIA”. Cabe reforçar que, na pesquisa do termo “ALEGRIA” na página textual das músicas (ver Figura 25), o termo em questão aparece destacado em amarelo nas letras selecionadas pelo algoritmo.

Além do termo principal, devem-se considerar também as outras palavras relacionadas (ver Figura 23), o que indica que a música contém o termo emocional “ALEGRIA”, reforçado por essas ocorrências. Esses termos derivam dos embeddings e são selecionados com base nas distâncias euclidianas a partir do termo em foco.

Outra informação importante é a análise do dendrograma gerado para cada termo do corpora musical (Figura 26), na qual se observam os termos e seus relacionamentos com as respectivas distâncias. Existem outros termos relacionados com o termo principal “ALEGRIA”, que, embora não estejam listados na Figura 23, também são relevantes a serem considerados.

Título + Letra	Artista	Gênero
3	4	5
<p>Moda-Rue modular modular da áfrica eu venho sou oriundo de la meus meus antepassados falaram que um vais crescer e levar em outros terras tudo o que existe na africa la bem distante cantando toda a alegria reinante cantaram os bantus sudaneses cantaram angola e congo e</p>	Roberto Ribeiro	samba
<p>Aleluim Desconhecido quem brinca a quem dançara quem saíra de aleluim quem tocara quem chorara quem gemera seu bandolim quem brincara quem dança quem saíra de aleluim quem to quem chorara quem gemera seu bandolim quem tera encamado a sua alegria quem tera ensaiado a sua folia quem tera o trabalho louvado quem tera o destino sagrado pra fazer minha gente feliz novamente quem brincara quem dançara quem saíra de aleluim quem tocara quem chorara quem gemera seu bandolim quem brincara quem dançara quem saíra de aleluim quem tocara quem chorara quem gemera seu bandolim quem tera a osadia quem tem o palhao quem tera seu encanto seu niso tao facil quem tera o officio bonito quem tera esse vicio benedito pra fazer minha gente novamente feliz quem brincara quem dançara quem saíra de aleluim quem tocara quem chorara quem gemera seu bandolim quem brincara quem dançara quem saíra de aleluim quem tocara quem chorara quem gemera seu bandolim</p>	Ivan Lins	mpb
<p>A Fórmula Mágica a sua fórmula e mágica chegou a magia para todo e pra sua alegria a fórmula mágica se transformou solo refrão issoissoisso e a sua fórmula que um vai mas eu sabiaeu sabia sabiaeu que um vinia a sua fórmula mágica fórmula magicafórmula magicaso poderia ser da sua contafórmula magicafórmula magicotoda horatodo a todo tempo</p>	Marina Lima	mpb
<p>Alegria do Senhor e a nossa força aleluia simbora alegria do senhor alegria do senhor e a nossa força alegria do senhor alegria do senhor e a nossa força alegria alegria alegria e a nossa força alegria alegria alegria alegria e a nossa força alegrias do senhor alegria do senhor e a nossa força alegres no senhor o deus da nossa salvação alegravos no senhor o deus da nossa salvação alegravos no senhor o deus da nossa salvação quem e Jesus ele eaa minha segurança ele e a minha fortaleza seu amor me enche de alegria nele sempre resultou! ele eaa minha segurança ele e a minha fortaleza seu amor me enche de alegria nele sempre resultarei alegria do senhor alegria do senhor e a nossa força alegria do senhor e a nossa força alegria alegria alegria alegria e a nossa força e a nossa força e a nossa força alegria alegria alegria e a nossa força alegradores no senhor o deus da nossa salvação alegravos no senhor o deus da nossa salvação quem e Jesus ele eaa minha segurança ele e a minha fortaleza seu amor me enche de alegria nele sempre resultou! ele eaa minha segurança ele e a minha fortaleza seu amor me enche de alegria nele sempre resultarei alegria do senhor alegria do senhor e a minha fortaleza seu amor me enche de alegria nele sempre resultou na palma da mão pra esquentar a galera segura e um vai não pode deixar cair seguro oooo</p>	Crianças do Rei	evangelho
<p>Origem da Felicidade a alegria e o alimento da alma alegria e nossa grande inspiração alegria recompensa os sacrifícios alegria soberana decisão alegria recompensa os sacrifícios alegria soberana decisão tempera o teu medo com a esperança inclina essa balança a teu favor alegria e a acreditar que o samba e o criador a alegria e a acreditar que o samba e o criador</p>	Mart'halia	samba
<p>Casa de Deus alegreime quando me disse vamos a na não existe tristeza não existe solidao na so existe alegria amor e comunhao eu vou cantar de alegria eu vou gritar de alegria eu vou pular de alegria com meus irmãos vou celebrar eu vou cantar de alegria eu vou gritar de alegria eu vou pular de alegria e a Jesus vou celebrar eu amo a</p>	Regis Danese	evangelho
<p>Tempo de Festa este e um este e um tempo de luvor pra celebrar aquele que primeiro nos amou transformou nosso choro em riso nos deu novas vestes de luvor pra celebrar aquele que primeiro nos amou nos jogou do império das trevas e nos deu perdão e paz arrancou todas as feridas nos fez felizes demais festa alegria e uma dança de celebração ao único digno Jesus seu nome e Jesus festa alegria e um povo que se reúne aqui do rei do rei dos reis seu nome e Jesus cantamos de alegria dançamos de alegria pulamos de alegria gritamos de alegria festa alegria e uma dança de celebração ao único digno Jesus seu nome e Jesus festa alegria e um povo que se reúne aqui do rei do rei dos reis cantamos de alegria dançamos de alegria pulamos de alegria gritamos de alegria</p>	Diante do Trono	infantil
<p>Nada de Tristeza viva a alegria viva a alegria a tristeza matou a muitos quem a ela se entregou ergue o teu braço e louva o teu senhor Jesus e a alegria Jesus e o amor viva a alegria</p>	Salette Ferreira	evangelho
<p>Hosana nos reunimos aqui para adorar o nosso rei pra celebrar sua fidelidade ele prometeu que entre nós sempre estará vamos declarar bem vindo aqui senhor glorioso salvador festa alegria e uma dança de celebração ao único digno Jesus seu nome e Jesus festa alegria e um povo que se reúne aqui do rei do rei dos reis cantamos de alegria dançamos de alegria pulamos de alegria gritamos de alegria seu nome e Jesus</p>	Diante do Trono	infantil

Verifica-se, por exemplo, que, na lista das vinte músicas, o termo “EUFORIA”, que aparece no dendrograma, concentra-se intensamente em duas músicas, sendo, portanto, um termo de peso considerável nessas obras.

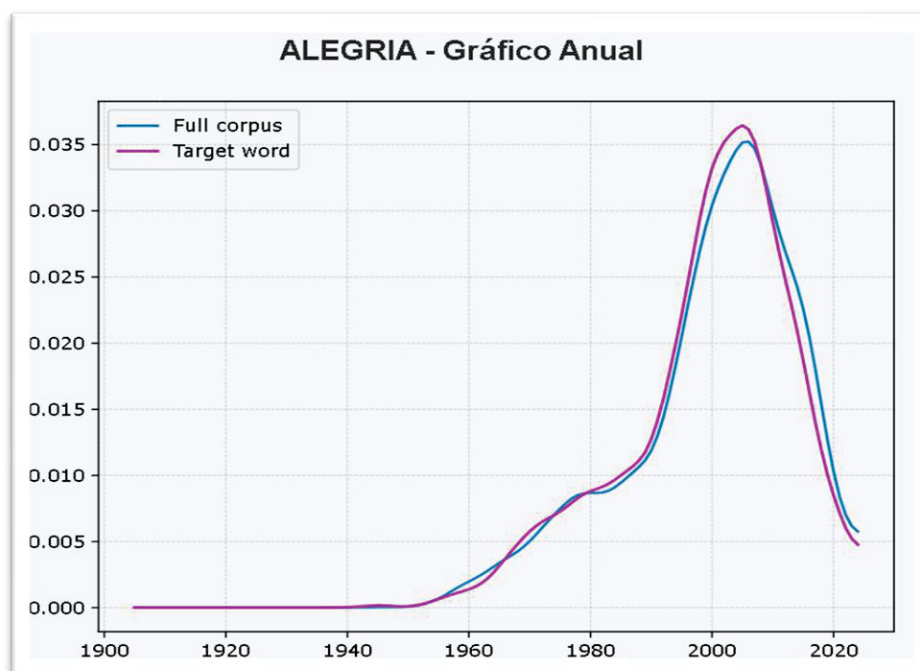
ALEGRIA - Árvore de Palavras

alegria
alegrar
alegra
alegre
alegremente
alegrias
alegrando
alegrou
alegres
palmas
dancas
transmitindo
saltar
guia
cante
cantar
canto
cantando
raiou
calmaria
dia
fantasia
melodia
riso
melancolia
euforia
negra
festejar
harmonia
energia

0.0357335

Uma última informação disponível na plataforma é o gráfico dos anos em que os títulos com o termo “ALEGRIA” foram usados nas músicas ao longo das décadas. Logicamente, a grande concentração de lançamentos das músicas presentes no *corpora* musical ocorre no período de 1980 a 2020, conforme a Figura 27.

Figura 27 – Termo “ALEGRIA” aplicado às músicas ao longo das décadas.



Fonte: O Autor (2025).

No experimento, as análises também foram realizadas para o termo “TRISTEZA”, relacionado à emoção. Na sequência, foram selecionadas 220 músicas para a emoção “ALEGRIA” e 240 para “TRISTEZA”, totalizando um *corpora* de teste de 460 músicas, com a finalidade de classificar as emoções utilizando o algoritmo de análise de sentimento VADER.

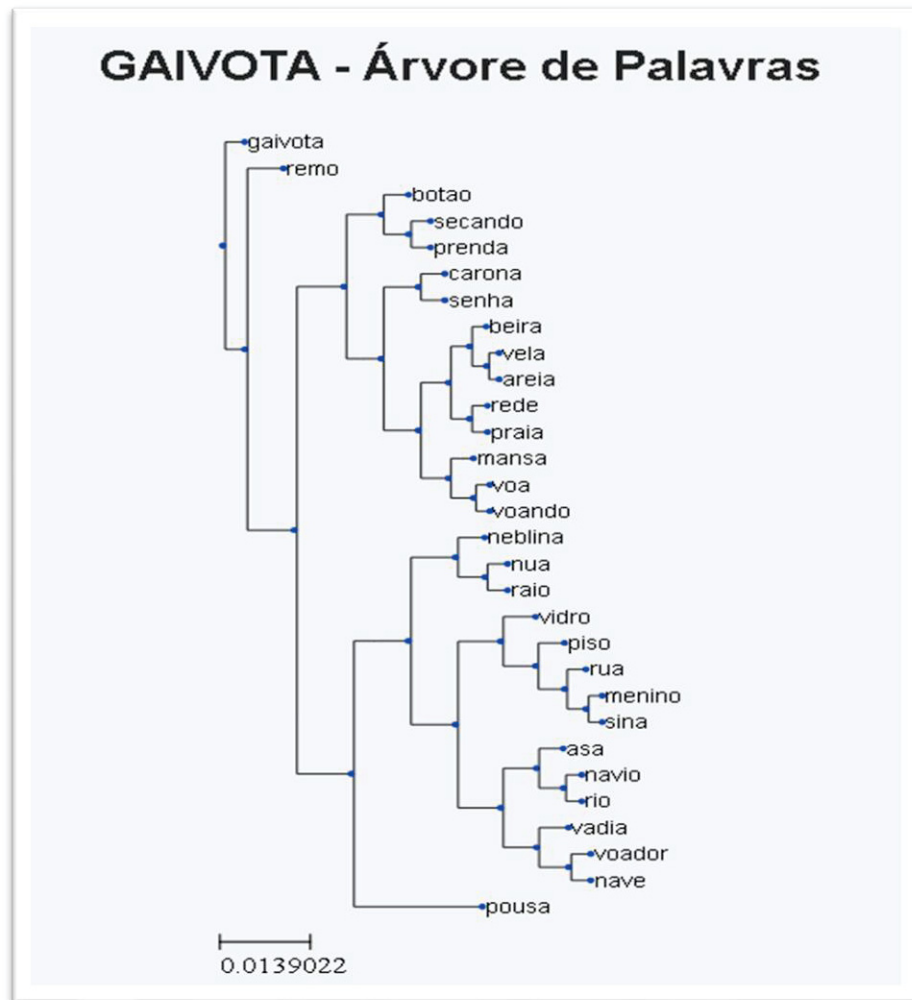
4.7 Análise da dimensão criativa no experimento

Outra abordagem de análise, realizada com *corpora* musicais, consistiu na escolha de um termo, por exemplo, “GAIVOTA” ou “GAIVOTAS”. Neste caso, não se refere a emoções, mas a um estudo empírico de criação.

A partir dos termos relacionados com “GAIVOTA” e “GAIVOTAS” contidos nos dendrogramas de árvore, selecionaram-se os seguintes termos: gaivota, voa, voando, neblina, vela, remo, beira, rede, secando, azuis, luar, colorindo, lua, estrelas, luas, dourado, azul, sereno, ver Figura 28. Com os termos associados, procedeu-se a um

experimento solicitando à IA ChatGpt para elaborar uma letra de poesia com o tema "Gaivota" usando os termos associados.

Figura 28 – Dendrograma do termo GAIVOTA



FONTE: O autor (2025)

A partir do algoritmo interno de PLN foi gerada a poesia no Quadro 6. Após isso, foi solicitada, ainda, a criação de uma melodia na plataforma de IA Suno para a referida poesia, com as cifras da melodia sugeridas pela IA ChatGpt.

Quadro 6 - Produção criativa artificial dos termos relacionados com “GAIVOTA”

<p>Título: Gaivota Luar Link para a obra: https://suno.com/s/JjQt2luFrVKy8J7L (Estrofe 1)</p> <p>Pela beira do mar eu te vi voando, gaivota branca no céu se soltando, cortando a neblina que vem devagar, voa, voa leve, sem olhar pra trás. (Estrofe 2)</p> <p>Na rede secando, o vestido azul, guardo o teu cheiro que o vento levou. O remo descansa, a vela dobrou, o barco parado onde o tempo parou.</p> <p>(Refrão)</p> <p>Azuis do teu olhar ainda estão em mim, sereno horizonte que nunca tem fim. O luar colorindo o mar de prata e dor, a lua e as estrelas chorando amor. Um dourado no céu quando a lua surgiu, duas luas no mar, uma delas partiu.</p> <p>[Ponte]</p> <p>Gaivota que voa, me leva contigo, pra longe da beira, pra longe do umbigo. Gaivota que voa, me leva contigo, pra longe da dor que ficou.</p> <p>[Refrão Final]</p> <p>Azuis do teu olhar ainda estão em mim, sereno horizonte que nunca tem fim. O luar colorindo o que o amor desfez, gaivota que voa... me leva outra vez.</p>	<p>Gaivota Luar</p> <p>[Tom: Am] [Estrofe 1] Am7 Am7/G Pela beira do mar eu te vi voando, F7M E7(13) gaivota branca no céu se soltando, Am7 D7(9) cortando a neblina que vem devagar, Gm6 A7(9-) voa, voa leve, sem olhar pra trás. [Estrofe 2] Dm7 G7(13) Na rede secando, o vestido azul, C7M(9) C7M(9)/B guardo o teu cheiro que o vento levou. Am7 Em7 A7(9) O remo descansa, a vela dobrou, Dm7 G7(13) C7M F7M E7(13) Am7 o barco parado onde o tempo parou.</p> <p>[Refrão] C7M Am7 Azuis do teu olhar ainda estão em mim, Dm7 G7(13) C7M sereno horizonte que nunca tem fim. C7M Am7 O luar colorindo o mar de prata e dor, Bm7(b5) E7(b13) Am7 a lua e as estrelas chorando amor. F7M Em7 A7(9) Um dourado no céu quando a lua surgiu, Dm7 G7(13) C7M F7M E7(13) Am7 duas luas no mar, uma delas partiu.</p> <p>[Ponte] F7M G7(13) Gaivota que voa, me leva contigo, C7M Am7 pra longe da beira, pra longe do umbigo. Dm7 G7(13) Gaivota que voa, me leva contigo, C7M E7(13) pra longe da dor que ficou.</p> <p>[Refrão Final] C7M Am7 Azuis do teu olhar ainda estão em mim, Dm7 G7(13) C7M sereno horizonte que nunca tem fim. C7M Am7 O luar colorindo o que o amor desfez, Bm7(b5) E7(b13) Am7 gaivota que voa... me leva outra vez.</p>
---	--

Cabe agora validar com um músico se há sentido na obra gerada com base nos termos do algoritmo ATMBR e na melodia sugerida pela IA Suno.

4.8 Análise da dimensão memória histórica indígena no experimento

A partir do termo relacionado com “TUPI”, verificou-se a relação de termos mais próximos: tupi, guarani, índia, tupã, Amazonas, Luanda, índio, baiano, tropical, Brasil. Na Figura 29, é importante observar que, além da seleção da música diretamente pela palavra “TUPI”, há várias ocorrências na coluna 3, totalizando 32 links de música na consulta. Isto permite analisar 640 músicas com o termo "TUPI".

Figura 29 – Pesquisa da memória histórica, termo indígena “TUPI”

Letras brasileiras - Palavras						
tupi,			Número dos vizinhos			
Procurar Ajuda						
Total de inscrições: 75955 Entradas selecionadas: 32			Exportar CSV Limpar filtros			
#	Palavra	Ocasionalmente.	Palavras relacionadas	Letras relacionadas	Árvore de Palavras	Gráfico anual
0	1	2	3	4	5	6
18081	TUPI	36	tupi, guarani, índia, tupa, Amazonas, Luanda, índio, baiano, tropical, Brasil	Letras relacionadas	Árvore de Palavras	Gráfico anual
12827	TUPA	31	tupa, tupi, pajé, tropical, Luis, Barbosa, Amazonas, pavão, pereira, xang	Letras relacionadas	Árvore de Palavras	Gráfico anual
18082	GUARANI	26	guarani, tupi, Amazonas, índia, índio, Zumbi, Paraguai, Cabral, palmares, pavão	Letras relacionadas	Árvore de Palavras	Gráfico anual
10637	CURUMIM	21	curumim, índia, Uirapuru, tupi, Bemtevi, tupa, Luzia, Passarim, Aladim, guarani	Letras relacionadas	Árvore de Palavras	Gráfico anual
10638	UIRAPURU	19	uirapuru, curumim, Bumbameuboi, guarani, tupi, índia, tabuleiro, pinheiro, maracatu, pajé	Letras relacionadas	Árvore de Palavras	Gráfico anual
18080	TUPINAMBA	14	tupinamba, Zumbi, palmares, tupi, África, africano, guarani, índio, africana, umbanda	Letras relacionadas	Árvore de Palavras	Gráfico anual
16424	SELVAS	12	selvas, índia, quintais, raras, negros, tupi, guetos, tropicais, praias, febril	Letras relacionadas	Árvore de Palavras	Gráfico anual
5271	CAPIBARIBE	11	capibaribe, canavial, Pernambuco, guarani, bumba, iracema, azulão, Amazonas, tupi, maracatu	Letras relacionadas	Árvore de Palavras	Gráfico anual
5243	TUPINIQUIM	10	tupiniquim, brasileira, Paraíba, tabuleiro, saci, tupi, banana, índia, oba, banana	Letras relacionadas	Árvore de Palavras	Gráfico anual
53033	LIMEIRA	9	limeira, umbanda, tupi, Parati, resina, guarani, índia, curio, saci, Zumbi	Letras relacionadas	Árvore de Palavras	Gráfico anual

FONTE: Autor (2025)

Na Figura 30, abaixo, são apresentadas letras de músicas relacionadas ao termo TUPI. Elas estão bem contextualizadas em relação ao termo, estabelecendo toda a relação com a questão histórica do Brasil, com os costumes e com os povos indígenas. Um professor de história pode fazer um passeio pelo tema apenas com músicas relacionadas.

Figura 30 – Memória histórica, lista de músicas relacionadas “TUPI”

Total de inscrições: 20 Entradas selecionadas: 2

Exportar CSV Limpar filtros

#	Classificação	Semelhança	Título + Letra	Artista	Gênero	Ano	Link	EU IA
0	1	2	3	4	5	6	7	8
79686	0	0,9934	Querelas do Brasil o Brasil não conhece o brasil o brasil nunca foi ao brasil tapir jabuti ilianha ali alaude piau ururau aqui atauda piacarioca porecrameca jobim akarone jobimacu oh oh oh perere camara tororo olerere piriri ratata karate olara o brasil não merece o brasil o brasil caçando o brasil jereba saci caandrades cunhas anranha aranha sertoes guimaraes bachianas aguas imanonaima ariranboia na aura das maos do jobimacu oh oh oh jerere sarara cururu olerere blablaba bafafa sururu olara do brasil sos ao brasil do brasil sos ao brasil do brasil sos ao brasil tinhorao urutu sucuni o jobim sabia bermevi cabucu cordovil cavambi olerere madeira olana e bangu olara cascadura água santa acari olerere ipanema e nova iguaçu olara do brasil sos ao brasil do brasil sos ao brasil	Elis Regina	bossa nova	1978	Link	1005652
79687	1	0,9926	In(d)ignacao eu fiquei indignado ele ficou indignado a massa indignada duro de tao indignado a nossa indignacao e uma mosca sem asas não ultrapassa as janelas de nossas casas indignacao indigna indigna nacao indignacao indigna indigna inacao lazzo matumbe araketu a bahia indignada carlos cachaca morenqueraivo meirelles o samba indignado vila dias papagaio cafezal pendura saia pau comeu tomas santa martao morro indignado ramiro lucio flavio e escadinha o crime indignado e jaizinho la na ponta indignados sata e seus asseclas imigrantes da asinia boca branca os inocentes e os leões da lagoinha indignados jaguarao oiapoque e guaicurus a zona indignada ginga mao branca negrinhos de sinha a capoeira indignada gaviões galoucura máfia azul jovem mancha verdeflamante independente a massa indignada	Vagabunda	pedra	1992	Link	1100161
79683	2	0,9920	Coração do Brasil são paulo tem fama e todo mundo que e a maior potência do nosso pais o seu interior conheça eu quis viajar em todo estado e muito feliz campinas barretos e mogi das cruz brotas e atibaia santos e queluz indústria paulista que tudo produz e a estrela que brilha e o mundo seduz talbete jauba urubanan ribeirão preto itu parmitir pirassununga e jabuticaba paraguacu sorocaba e pinhar praju olimpia e aracatuba são roque pompeia rio preto angatuba de ranchania marilha ubatuba de araraquara e caraguatatuba jundiá lorena prudente rio claro santo anastacio abare santo amaro franca ourinhos bragança e amparo botucatu que meus pais se realizados são carlos limeira lensei jambeiro piracicaba bufet e cruzeiro de casa branca piquete barreiro itapetininga bernardo e mineiro santa cruz do rio pardo são sebastião caçapava itapulis e promissão batatais de belouro e cubatao são jose dos campos e outras povoacao o estado de sao paulo e de uma nacao tanto em riqueza e em populacao por isso que eu digo com muita razão do rico brasil sao paulo e o coração	Tiao Carreiro e Carreirinho	sertanejo	1972	Link	1132115

FONTE: Autor (2025)

4.9 Análise da dimensão do regionalismo cultural no experimento

A partir dos termos relacionados com “MINEIRO”, conforme a Figura 31, verificou-se que os termos mais próximos de “MINEIRO” são: violeiro, caipira, minas, tiao, puxando, galo e mato. Termos bem característicos do ambiente do estado de Minas Gerais.

Figura 31 – Pesquisa na dimensão regionalismo termo “MINEIRO”

Brazilian Lyrics - Words

MINEIRO Neighbors number

Search Help

Total entries: 75955 Selected entries: 61

Export CSV Clear Filters

#	Word	Occ.	Related Words	Related Lyrics	Word Tree	Year Plot
0	1	2	3	4	5	6
17959	MINAS	328	minas, gerais, paulo, vila, paulista, zona, martins, mineiro, selva, bala	Related Lyrics	Word Tree	Year Plot
3982	VIOLEIRO	220	violeiro, viola, caipira, mato, mineiro, gado, galo, roca, madeira, boiada	Related Lyrics	Word Tree	Year Plot
11868	MINEIRO	125	mineiro, violeiro, caipira, ootas, minas, tiao, puxando, paulista, galo, mato	Related Lyrics	Word Tree	Year Plot

Fonte: Autor (2025)

Os dados da Figura 32 ilustram a relação entre as letras de música e o termo “MINEIRO”. As canções refletem a riqueza do regionalismo brasileiro, abordando costumes e o ambiente típico de Minas Gerais. Esses elementos sugerem que tais músicas constituem um recurso valioso para aulas de Geografia, facilitando a assimilação da cultura mineira pelos estudantes.

Figura 32 – Regionalismo, lista de músicas relacionadas com “MINEIRO”

MINEIRO - Word Related Documents

MIN

Neighbors number

Search

Help

Total entries: 20 Selected entries: 6

Export CSV

Clear Filters

#	Rank	Similarity	Title + Lyrics	Artist	Gender	Year	Link	ID
0	1	2	3	4	5	6	7	8
79398	15	0.9931	Foi-se O Que Era Doce me descaderei de tanto xaxa no bobo de noivado da do ribamar vi quando cheguei as moca de la cozinhando uns inhame com os de arrevira buzanza de flor chulapa de mel e a covanca soprando um sussurro descido do ceu tinha gago anao gente de azar com a espinhela caida pedindo pro inhame estala jabacule virge espetacular assunto assim as veis e calar mas desque eu provei do bobo eu roxo pra comentar sanfona guitarra batuque berreiro e veja voce o viradesvira o caminho da roca e o balance inhame e bobo frango asado cuscus e maracuja pucanga cobreiro retreta jarguete e tamandua foguete beijando as estrela e as moca la ze pingulim chico do pincel paqueraro lazinha que era muie de xexeu serafim tres pema resolveu chia pois muie nao e farinha que vai pra onde venta deuse um sururu de saculeja tudo dando e levando enquanto sem se mancar pedro gargarejo com a mao no manjar preparava um caldinho pra noiva gargareja jabacule acudi um que tava no chao tomei uma no ouvido de adevorve o pirao foi um cimiterio foi um carnaval de paixoes confundidas quem e que tira a moral pra ser semvergonha basta ser decente e quem vende saude possivelmente e doente foise o que era doce ninguem quer contar quanto macho afinouse na festa do ribamar	Joao Bosco	samba	1974	Link	1104626
7929	17	0.9931	Quem nasceu nasceu quem nao nasceu nao nascera saida frente do trem bala pa tu nao se machucar quem nao nasceu nao nascera e a reliquia vermelho e nos que ta saudades eternas do da saudades do bruninho ele era meu chara saudades do caverinha do tuninho e do caju o da vila ideal e o gordo do saudade do dance e do andresinho moral que saudades do coco la vigario geral com saudades do cucu saudade do bozinho era o terror do zinco saudades do tiaginho o da o a la do salgueiro o orlando jogador nego bruto e o mineiro liberdade pro maluco liberdade pro marcinho poca russa benemerito alda sombra e o babinho o beira mar maido sapinho da provi la do barbanti da ilha liberdade para rose marcelo chara o dinho o patrick e o charlei brow liberdade isaia que o nosso general rapaziada que no problema na rua guerriando e o marretao bolado e o fabiano o piloto no mandela o papai la no manguinho no engenho e o lacai e o mano kinho rei davi e o bacalhau quadilha do chapadao do cajero biscoito e o macarrao o julin da o naiba e o paulinho no falete fugueteiro e o bonde do paulinzinho la furquim e o loco estilo colombiano e o menor e o gordao la no morro do turano e o bonde do jogador favela do arara maguino alex luthor lambari e o polegar quem nao nasceu nao nascera saida frente do trem bala pa tu nao se machucar quem nao nasceu nao nascera e a reliquia vermelho e nos que ta	Mc G3	funk-carioca	1999	Link	1021387
8556	18	0.9931	Vou Com Gas vou contente pra minas gerais vou contente pra minas gerais vou contente pra minas gerais vou contente pra minas gerais chegando la vou rever minha mineira mineira mineira vou contente pra minas gerais vou contente pra minas gerais vou contente pra minas gerais vou contente pra minas gerais chegando la vou	Tim Maia	mpb	1979	Link	1079146

Fonte: Autor (2025)

O experimento que analisa a dimensão do regionalismo cultural nas letras de músicas revela-se bastante rico, uma vez que o Brasil, em sua diversidade, apresenta expressões musicais regionais que vão do Sul, passando pelo Sudeste, Nordeste e Centro-Oeste, até alcançar o Norte, compondo um mosaico cultural amplo e profundamente significativo

4.10 Instrumento para avaliação da plataforma RITMO BRASIL

Após a geração da plataforma Ritmo Brasil, resultante da aplicação do método ATMBR e da disponibilização para os usuários na internet. Elaborou-se um instrumento de avaliação da plataforma, destinado a músicos e estudiosos da música brasileira popular.

Enviou-se um texto de explicação para um grupo seletivo de músicos e apreciadores da música, sobre o preenchimento do instrumento eletrônico de pesquisa, ver o Apêndice C, baseado em duas seções:

- Seção 1 - foram feitas perguntas dirigidas a respeito da idade e área de formação na graduação, formação musical e preferência do gênero musical;

- b) Seção 2 - perguntas objetivas, com escala Likert, para avaliar a usabilidade da plataforma e as opções de consulta, em termos quantitativos, quanto a gêneros e títulos.
- c) Seção 3 – Avaliação da consulta que destaca termos ou palavras mais significativos nas letras, incluindo contagens de ocorrências, palavras relacionadas, árvores de palavras, gráficos por ano e as 20 músicas mais relacionadas ao termo selecionado.
- Avaliar a relevância dessa funcionalidade, consulta de palavras em músicas, para buscar inspiração em composições de letras;
 - Os gráficos de distribuição de lançamentos de músicas por ano são úteis para identificar as músicas (ou tendências e contextos)?
 - De modo geral, a plataforma é útil para a composição de novas letras musicais brasileiras?
 - Sugerir alguma funcionalidade e/ou dados a serem acrescentados à plataforma, comente abaixo.
 - Opinar sobre a utilidade da plataforma como subsídio a estudos de música ou como fonte de inspiração para composições e/ou produções musicais?

A realização dos experimentos, fundamentada no corpus musical previamente curado e vetorizado, viabilizou a execução das etapas descritas anteriormente. Com base nesses procedimentos, a análise dos resultados é apresentada no capítulo 5 a seguir.

5 APRESENTAÇÃO DOS RESULTADOS

5.1 Avaliação da plataforma Ritmo Brasil

O instrumento de avaliação foi aplicado a músicos e a apreciadores da música popular brasileira, com formação escolar variada. Embora a maioria possua algum conhecimento musical, seja pela aprendizagem de instrumentos, pela participação em bandas ou pela atuação profissional, vale destacar a participação de integrantes da banda Sabotage, a Banda Mais Bonita da cidade, e de um integrante do grupo Mutantes. O período de aplicação do instrumento ocorreu entre 11 e 24 de novembro de 2025, resultando em respostas válidas, posteriormente tabuladas e analisadas. A partir dessas respostas eletrônicas, observaram-se os seguintes resultados.

5.1.1 Visão Geral da Amostra

Para esta análise, os respondentes foram segmentados em dois grupos distintos, com base na autodeclaração de "Experiência com a música" e "Experiência profissional".

Quadro 7 – Amostra de respondentes

Grupo	Respondentes
A - Somente Ouvintes "Sou apenas ouvinte" e/ou não possuem formação instrumental/teórica declarada	7 (23%).
B - Músicos, Compositores e Estudantes de Música: Músicos profissionais e compositores, pessoas com formação instrumental ou com graduação ou especialização na área.	23 (77%)
Total de respondentes	30 (100%)

Fonte: o Autor (2025)

A maioria dos participantes da pesquisa demonstrou possuir conhecimentos musicais e alguma experiência prévia com a música, totalizando 23 respondentes, o que representa 77% do total de respostas obtidas.

5.1.2 Correlações entre Formação Escolar e Gostos Musicais

A correlação entre a área de formação e os gêneros preferidos levou às seguintes conclusões:

- Exatas/Engenharias: o grupo mais numeroso, com preferência marcante pelo rock BR e pela MPB. Também foi o grupo que mais citou a "Velha Guarda", especialmente entre os respondentes com mais de 60 anos.
- Artes / Música: grupo com perfil mais técnico. A preferência recai sobre gêneros ricos em harmonia e letra, como bossa nova, MPB e samba.
- Humanas / Sociais Aplicadas: demonstram apreço pela narrativa lírica. Preferência dominante por MPB e samba, com abertura para o rock BR e sertanejo, entre os mais jovens.
- Biológicas / Médicas: perfil mais eclético, com base sólida em Bossa Nova e MPB.

De maneira geral, o MPB aparece como o gênero unificador entre todas as formações, enquanto o rock BR é a preferência mais evidente entre os respondentes de exatas.

5.1.3 Preferência de gêneros e relação com a faixa etária

A análise demonstra claramente o efeito geracional sobre o gosto musical:

- 18 a 29 anos: grupo mais diverso. Embora consumam MPB e rock, são os únicos a citar consistentemente o funk carioca, o forró, o axé e o sertanejo. Mostram abertura a ritmos populares contemporâneos e regionais.
- 30 a 49 anos: há hegemonia do rock BR e do POP. Esta faixa etária viveu o auge das bandas de rock nacional, e isso se reflete nas respostas, com o rock sendo quase unânime como preferência, ao lado da MPB.
- 50 a 59 anos: trata-se de uma faixa de transição. O rock começa a ceder espaço a gêneros mais tradicionais, como a MPB e o samba.
- 60 anos ou mais: nesta faixa etária, o gosto está mais consolidado na tradição brasileira, com predominância absoluta de Bossa Nova, Samba e Velha Guarda. O interesse por gêneros de "alta energia" (Funk, Rock) é residual ou inexistente.

5.1.4 O quanto é útil a plataforma para os músicos

A utilidade foi avaliada como alta, com nota média de aproximadamente 9/10, porém, mais direcionada à etapa de letra/lírica.

- Ferramenta de Referência: músicos veem grande utilidade em ter um banco de dados centralizado. "Um prato cheio pra quem quiser pesquisar" (respondente 15).
- Apoio ao Bloqueio Criativo: considerado muito útil no início da composição, quando o músico precisa de um norte temático ou de rimas (respondente 11).
- Validação de Mercado: "auxilia a compreender o que já foi feito". Alguns destacam a importância de verificar se determinada temática já está saturada (respondente 25).
- Limitação: a utilidade diminui quando o músico busca apoio melódico ou harmônico, pois a plataforma (nesta versão) é focada exclusivamente no texto.

5.1.5 Possibilidades a serem exploradas pelos músicos

Com base nas respostas, identificam-se os seguintes casos de uso prático:

- Engenharia reversa de sucessos: analisar quais palavras são mais recorrentes em músicas de sucesso de um gênero específico (ex.: "o que faz um sertanejo ser sertanejo em termos de texto?");
- Expansão de vocabulário: usar a busca para encontrar sinônimos ou palavras correlatas que fujam do lugar-comum, enriquecendo a poesia da canção;
- Análise temporal (Zeitgeist¹): compositores podem explorar como o uso de determinadas palavras mudou ao longo das décadas, criando músicas que soem "vintage" ou "modernas" propositalmente; e
- Pesquisa de títulos: verificar se o título proposto para uma nova música já existe e em que contexto foi utilizado.

¹ Zeitgeist - Do alemão "espírito do tempo". O conceito denota o conjunto de ideias, crenças e sensibilidades que caracterizam um período histórico. No contexto desta análise, refere-se à capacidade das letras de música de capturar e refletir as tensões sociais, as gírias e as preocupações coletivas vigentes no momento de sua composição.

5.1.6 Recursos que deveriam ser adicionados

As sugestões dos respondentes foram técnicas, com foco em usabilidade (UX) e aspectos musicológicos:

- Idioma: tradução da interface para o português. O inglês foi repetidamente citado (respondentes 11 e 28) como uma barreira ao uso fluido.
- Busca inteligente (NLP):
 - aceitar termos sem acento ou com grafia aproximada (ex.: "intenções" "intenções");
 - permitir buscas por frases ou trechos, não apenas palavras isoladas; e
 - *autocomplete* na barra de pesquisa, para guiar o usuário.
- Funcionalidades musicais:
 - dicionário de rimas para sugerir palavras que rimam com o termo pesquisado;
 - Inclusão de BPM, tonalidade e harmonia para análises como "Palavras tristes ocorrem mais em tons menores?", e *links* externos, permitindo integração direta com *Spotify* e *YouTube* para ouvir músicas encontradas.
- Navegação: botão de "voltar" mais visível e paginação dos resultados para evitar lentidão.

5.1.7 Finalidade da Plataforma e Uso para Pesquisas

Análise textual da questão 2.10. Esta seção avalia a percepção de valor da plataforma, diferenciando as visões por perfil de usuário.

a) Percepção dos "somente ouvintes"

Este grupo tende a avaliar a plataforma sob a ótica da curiosidade cultural e da experiência de navegação, com menor foco na aplicação técnica.

- Finalidade percebida: fonte de conhecimento histórico e de compreensão das tendências culturais.
- Análise das respostas:

- Respondente 16 destaca o potencial para "análises textuais, lexicais e temáticas aprofundadas sobre a evolução dos estilos";
- Respondente 36 ressalta a capacidade de identificar "movimentos, temas e palavras relativas a diferentes períodos históricos", funcionando como espelho da sociedade da época;
- Críticas / limitações: alguns relataram dificuldade de compreensão inicial (respondente 35 considerou "não intuitivo"). O respondente 20 levantou uma preocupação mais filosófica: o receio de que a ferramenta possa: "induzir o compositor a criar não com sentimento, apenas produzir", vendo a IA como um risco de mecanização da arte.

b) Percepção dos músicos, compositores e estudantes

- A visão deste grupo foi mais pragmática. A plataforma é validada como ferramenta de trabalho (o "martelo e o cinzel" do compositor) e de pesquisa musicológica.
- Finalidade percebida: apoio à composição (principalmente na letra), verificação de originalidade e estudo do vocabulário característico de cada gênero.
- Análise das respostas:
 - Respondente 9 define a finalidade como mapeamento de estilo: "mostra um caminho interessante, onde a maioria dos compositores transita, dificilmente poderemos fugir dos termos apresentados";
 - Respondente 23, estudante de Letras, descreve a plataforma como "uma carta na manga" para verificar associações entre palavras.
 - Os respondentes 25 e 11 destacam a utilidade "ex-ante" (inspiração inicial) e "ex-post" (verificação), especialmente para evitar o plágio ou a repetição excessiva.
 - Ressalva importante: os respondentes 11 e 26 reforçam que a plataforma é excelente para a letra, mas limitada para melodia e harmonia, lembrando que a canção é um conjunto. Assim, sua finalidade principal permanece o "subsídio textual".

Em suma, a plataforma Ritmo Brasil atinge seu objetivo acadêmico ao demonstrar a aplicabilidade da IA na análise textual de música. Contudo, para sua evolução como produto de informação, conforme a avaliação com os músicos, recomenda-se a adição de recursos semânticos avançados, como um dicionário de rimas, bem como melhorias na interface. A plataforma não é apenas um banco de dados, mas também um espelho quantitativo da subjetividade brasileira, oferecendo métricas para o que antes era apenas intuitivo.

5.2 Resultados dos experimentos da análise de sentimentos.

Com o intuito de verificar e validar o uso do método ATMBR na plataforma Ritmo Brasil, na perspectiva da consulta de termos de sentimentos, foi realizado um experimento comparativo aplicando o método VADER para análise direta dos termos contidos nas músicas e a classificação supervisionada (MLP, SVM e Random Forest) utilizando o *corpora* vetorizado (*embedding*).

Ressalta-se que o comparativo apresentou dados interessantes sobre as emoções presentes nas letras das músicas. Foram selecionadas músicas associadas às emoções “alegria”, “tristeza”, “medo” e “raiva”, considerando não apenas as palavras exatas, mas também suas derivações, como “alegres” e “alegremente”, preservando o mesmo radical (Quadro 8). O mesmo procedimento foi aplicado aos demais sentimentos.

Quadro 8 – *Corpora* selecionado para os experimentos

Emoção	Seleção dos termos	Quantidade Letras musicas
Alegria	alegria, alegrar, alegremente, alegre, alegre, alegrarei, alegres, alegrou, alegrando, alegre, alegrias	219
Tristeza	entristece, entristecer, entristeceu, entristecido, triste, tristes, tristezas, tristonho, tristemente, tristeza, tristonha, tristonhos	80
Medo	medo, medonha, medonho, medos, amedrontar	80
Raiva	enraivecida, raiva, raivas, raivinha, raivoso	40
	Total de letras de músicas analisadas	659

Fonte: Autor (2025)

A plataforma permite exportar as letras das músicas no formato “csv” e, a partir destes arquivos, foi gerado o corpus musical utilizado na comparação, conforme apresentado no QUADRO 8.

A Tabela 2 apresenta os resultados da classificação do *corpora* das letras e das emoções básicas gerados pelo método ATMBR e submetidos ao algoritmo VADER. As músicas foram categorizadas nas seguintes emoções: alegre, triste, raiva e medo.

Tabela 2 – Resultado aplicação do algoritmo VADER

Emoção	Precision	Recall	F1-Score	Total Real
ALEGRIA	0.518	0.659	0.580	88
TRISTEZA	0.479	0.495	0.486	91
RAIVA	0.288	0.211	0.244	90
MEDO	0.398	0.385	0.391	91

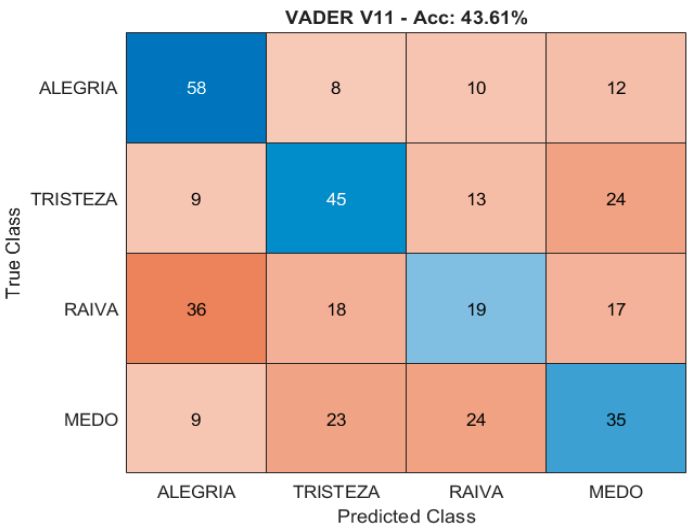
Fonte: Autor (2025)

Acurácia geral: 43.61% (157/360) e Macro-F1 Média: 0.425

Os modelos baseados exclusivamente em léxico (VADER-pt), mesmo com léxico extenso e traduzido, não ultrapassaram 43,61% de acurácia geral e apresentaram desempenho inferior nas classes minoritárias (Raiva e Medo).

A matriz de confusão (ver a Figura 32) apresenta erros acentuados entre as classes.

Figura 33 – Matriz de confusão, análise de emoções VADER



Fonte: Autor (2025)

O modelo com classificadores supervisionados (MLP, SVM e Random Forest) é bem superior ao do primeiro experimento, que utilizou apenas o algoritmo VADER.

Os resultados alcançados são de nível de estado da arte e apresentam-se como solução para a classificação robusta de emoções (Alegre, Triste, Raiva e Medo) em letras de música em português brasileiro.

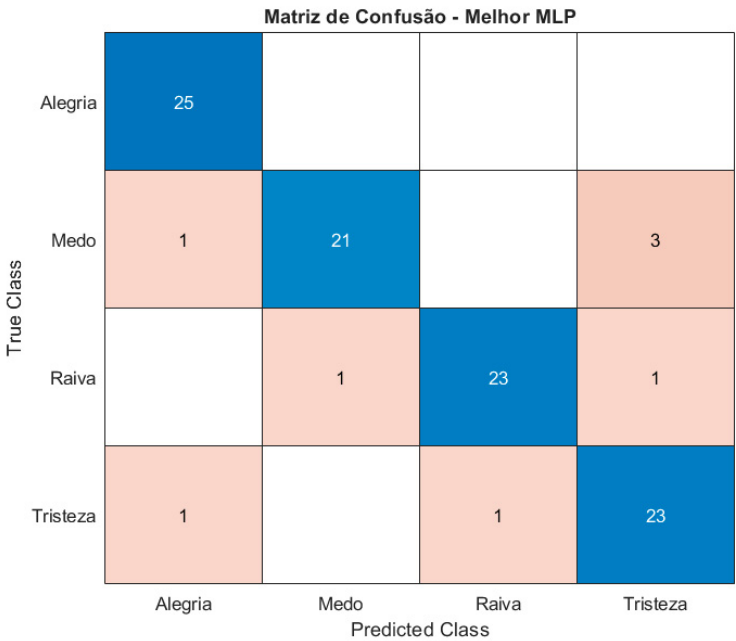
Tabela 3 – Resultados comparativos algoritmos supervisionados

Método	PCA	Configuração	Acurácia	F1-macro	Precisão
MLP	34	Hidden=377	92.0%	0.9230	0.9261
MLP	34	Hidden= 89	91.0%	0.9132	0.9165
MLP	89	Hidden=233	90.0%	0.9054	0.9108
Random Forest	55	Trees=500	89.0%	0.9025	0.9153
SVM	34	RBF (C=10, auto)	89.0%	0.8915	0.8930

Fonte: Autor (2025)

É sensível e, claro, à abordagem do uso de classificadores supervisionados com o corpus vetorizado (embeddings semânticos). O embedding ATMBR foi treinado em grande volume de texto em português brasileiro, capturando nuances semânticas, contexto, ironia e expressões idiomáticas presentes nas letras de músicas.

Figura 34 – Matriz de confusão da MLP (PCA=34, HD =377)



Fonte: Autor (2025)

Mesmo com redução drástica da dimensionalidade (PCA com 34 componentes), o MLP manteve 92% de acurácia. Neste caso, o espaço semântico

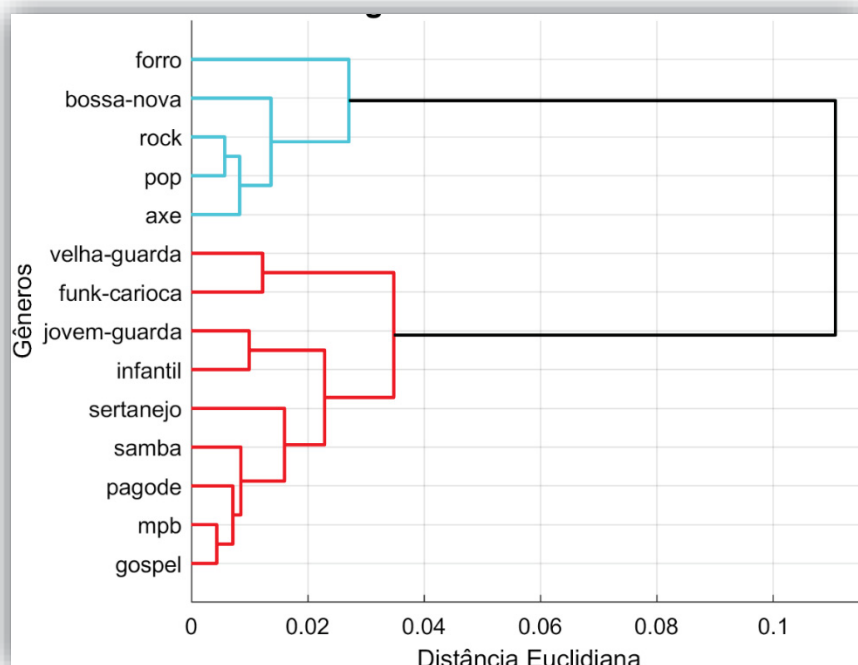
aprendido é extremamente informativo para a tarefa. A matriz de confusão do melhor MLP (Figura 33) mostra uma separação quase perfeita: Alegria: 25/25 corretas; Tristeza: 23/25; Medo: 23/25; Raiva: 21/25. Apenas 6 erros em 100 amostras de teste, comportamento próximo ao estado da arte para a tarefa.

Para sistemas de identificação automática de emoção em letras de músicas brasileiras, recomenda-se fortemente o uso de representações distribuídas (embeddings) treinadas em português e de classificadores supervisionados, como *MLPs* ou *Random Forests*. Métodos puramente léxicos, como o VADER, mesmo com adaptações, permanecem muito aquém do desempenho alcançável com o aprendizado de representações contextuais.

5.3 Definição do dendrograma dos gêneros musicais

A partir dos embeddings gerados para o corpus de letras de músicas, aplicou-se um classificador supervisionado do tipo *Perceptron* Multicamadas (MLP). Como as classes correspondentes aos gêneros musicais de cada letra já eram conhecidas, o modelo foi treinado supervisionado.

Figura 35 – Dendrograma dos gêneros musicais



Fonte: Autor (2025)

Após a obtenção das representações vetoriais (embeddings) e a predição/confirmação dos gêneros pelo classificador, tornou-se possível construir um

dendrograma (Figura 35) que revela a estrutura hierárquica de similaridade entre os diferentes gêneros musicais, com base exclusivamente nas características textuais das letras.

A análise do dendrograma de agrupamento hierárquico, baseada em embeddings textuais, revelou uma dicotomia clara na música brasileira. Identificaram-se dois macro grupos: um de influência Pop/Rock internacional (agregando Rock, Pop, Bossa Nova, Axé e Forró, este último devido à similaridade linguística contemporânea) e um segundo, mais coeso, de matriz nacional e afro-brasileira (Samba, Pagode, MPB, Sertanejo e Funk).

Internamente, destacam-se a forte coesão dos núcleos Samba, Pagode e MPB, e o isolamento do gênero Gospel, cujo léxico religioso o torna o maior *outlier* da amostra. Conclui-se que o conteúdo semântico das letras é suficiente para recuperar, de forma não supervisionada, a genealogia sócio-histórica dos gêneros, distinguindo com precisão as raízes culturais das letras de músicas brasileiras.

5.4 Resultados da dimensão descritiva

O aspecto das análises da historicidade e do regionalismo brasileiro, verificado no *corpora* musical gerado pelo método e nas listas de músicas, foi apenas um aspecto explorado sob a ótica educacional.

Na dimensão descritiva, inúmeras possibilidades de abordagem educacional se delineiam, e os músicos, devido ao seu conhecimento específico, percebem os detalhes contextuais dos termos empregados nas músicas. Dessa forma, torna-se possível explorar mais profundamente a dimensão da criação, dos gêneros, da temporalidade e de outros enfoques relacionados.

6 CONSIDERAÇÕES FINAIS

A interdisciplinaridade entre áreas de conhecimento, como a computação, as artes e a biologia, presente neste projeto em específico, tem como meta a construção de metodologias e ferramentas, visando desenvolver mais rapidamente a recuperação de informações musicais. Embora produza resultados satisfatórios para o propósito a que se destina, em muitos casos, a interdisciplinaridade com outras áreas das ciências não é amplamente explorada, e este caso não é exceção.

Esta pesquisa contribui ao associar a disciplina de mineração de texto e de aprendizado de máquina à recuperação de informações textuais de música, utilizando algoritmos da bioinformática. Essa abordagem auxilia não apenas na área específica, mas também beneficia outras disciplinas que buscam novas perspectivas. Na ciência da informação, o volume de dados e informações cresce exponencialmente e esta situação carece de novas tecnologias para atender às demandas.

A metodologia proposta, com suas funcionalidades fundamentadas em métodos de bioinformática e com custo computacional reduzido em comparação com outras ferramentas, visa atender a uma interseção entre a música e a ciência da informação. Outro aspecto relevante é o destaque para as pesquisas sobre *corpora* musicais brasileiros, frequentemente negligenciadas nos estudos e pesquisas internacionais. Esta pesquisa também apresenta e compartilha um *corpora* genuinamente brasileiro para futuras pesquisas na área.

Como etapas futuras para a conclusão do método ATMBR, é necessário realizar experimentos mais aprofundados com o *corpora* musical vetorizado, aprimorar o pré-processamento, desenvolver modelos multimodais com sinal de áudio, desenvolver novos métodos de rotulação mais específicos e encaminhar uma publicação que aborde os resultados alcançados.

REFERÊNCIAS

- AFCHAR, Darius et al. Explainability in music recommender systems. **AI Magazine**, v. 43, n. 2, 2022.
- AHUJA, Mahesh; SANGAL, A. L. Opinion Mining and Classification of Music Lyrics Using Supervised Learning Algorithms. In: **2018 First International Conference on Secure Cyber Computing and Communication (ICSCCC)**, 2018. p. 223-227.
- AUSLANDER, Noam; GUSSOW, Ayal B.; KOONIN, Eugene V. Incorporating machine learning into established bioinformatics frameworks. **International Journal of Molecular Sciences**, 2021.
- AYYILDIZ, Dilara; PIAZZA, Silvano. Introduction to Bioinformatics. In: **Methods in Molecular Biology**. New York, NY: Humana Press, 2019. v. 1986, p. 1-15. DOI: https://doi.org/10.1007/978-1-4939-9442-7_1.
- BAIRD, Amee; SAMSON, Séverine. **Memory for music in Alzheimer's disease: Unforgettable?** **Neuropsychology Review**, v. 19, n. 1, p. 85-101, 2009. DOI: <https://doi.org/10.1007/s11065-009-9085-2>.
- BLASZKE, Maciej; KOSTEK, Bożena. Musical Instrument Identification Using Deep Learning Approach. **Sensors**, v. 22, n. 8, p. 3033, 15 abr. 2022. DOI: <https://doi.org/10.3390/s22083033>.
- BOGHRATI, Reihane; BERGER, Jonah. Quantifying Cultural Change: Gender Bias in Music. **Journal of Experimental Psychology: General**, v. 152, n. 2, p. 314–329, 2023. DOI: <https://doi.org/10.1037/xge0001412>
- BRATTICO, Elvira et al. A functional MRI study of happy and sad emotions in music with and without lyrics. **Frontiers in Psychology**, v. 2, art. 308, dez. 2011. DOI: <https://doi.org/10.3389/fpsyg.2011.00308>.
- BUCHANAN, Bruce G. A (very) brief history of artificial intelligence. **AI Magazine**, v. 26, n. 4, 2005. DOI: <https://doi.org/10.1609/aimag.v26i4.1848>
- BURGGREN, Warren et al. **The Oxford Handbook of Interdisciplinarity**. Oxford. Oxford University Press, 2017. DOI: <https://doi.org/10.1093/oxfordhb/9780198733522.013.9>
- CAFÉ, Lígia Maria Arruda; BARROS, Camila Monteiro de. Abordagens metodológicas das pesquisas sobre organização da informação musical. **Revista Brasileira de Biblioteconomia e Documentação**, v. 14, n. 3, p. 304–323, jul. 2018. Disponível em: <https://rbbd.febab.org.br/rbbd/article/view/1126> . Acesso em: 12 dez. 2023.
- CALDAS, W. **Iniciação à música popular brasileira**. São Paulo: Amarylis, 2010.
- CAMRAS, Linda; PLUTCHIK, Robert; KELLERMAN, Henry. Emotion: Theory, Research, and Experience. Vol. 1. Theories of Emotion. **The American Journal of Psychology**, v. 94, n. 2, 1981. DOI: <https://doi.org/10.2307/1422757>
- CHANG, Hong Yi; HUANG, Shih Chang; WU, Jia Hao. A personalized music recommendation system based on electroencephalography feedback. **Multimedia Tools and Applications**, v. 76, n. 19, 2017. DOI: <https://doi.org/10.1007/s11042-015-3202-4>

CHAUHAN, Swati; CHAUHAN, Prachi. Music mood classification based on lyrical analysis of Hindi songs using Latent Dirichlet Allocation. In: International Conference on Information Technology, 2017. DOI: <https://doi.org/10.1109/INCITE.2016.7857593>

CHOO, Chun Wei. **A organização do conhecimento: como as organizações usam as informações para criar significado, construir conhecimento e tomar decisões**. Tradução: Eliana Rocha. São Paulo: Editora Senac, 2003. ISBN: 85-7359-341-5 / 978-85-7359-341-9

COHEN, K. Bretonnel; HUNTER, Lawrence E. Chapter 16: Text Mining for Translational Bioinformatics. **PLOS Computational Biology**, v. 9, n. 4, p. e1003044, abr. 2013. DOI: <https://doi.org/10.1371/journal.pcbi.1003044>.

CÔRTE, Beltrina; LODOVICI NETO, Pedro. A musicoterapia na doença de Parkinson. **Ciência & Saúde Coletiva**, v. 14, n. 6, 2009.

DA SILVA, Angelo Cesar Mendes; SILVA, Diego Furtado; MARCACINI, Ricardo Marcondes. 4MuLA: A Multitask, Multimodal, and Multilingual Dataset of Music Lyrics and Audio Features. In: **Proceedings of the Brazilian Symposium on Multimedia and the Web (WebMedia '20)**, New York: ACM, 2020. p. 145-148. DOI: <https://doi.org/10.1145/3428658.3431089>.

DAMÁSIO, António R. **O erro de Descartes: emoção, razão e o cérebro humano**. Tradução de Dora Vicente e Georgina Segurado. São Paulo: Companhia das Letras, 1996.

DE LUCCA, J. L.; NUNES, Maria das Graças Volpe. **Lematização versus Stemming**. São Carlos - SP, 2002.

DE PIERRI, Camilla Reginatto et al. SWeeP: representing large biological sequence datasets in compact vectors. **Scientific Reports**, v. 10, n. 1, 2020. DOI: <https://doi.org/10.1038/s41598-019-55627-4>.

DEVLIN, Jacob; CHANG, Ming-Wei, LEE, Kenton; TOUTANOVA, Kristina. BERT: Pre-training deep bidirectional transformers for language understanding. In: **Proceedings of NAACL-HLT 2019**, Minneapolis, Minnesota, USA, 2019. p. 4171–4186. DOI: <https://doi.org/10.18653/v1/N19-1423>.

DOH, S.; WON, M.; CHOI, K.; NAM, J. Toward Universal Text-To-Music Retrieval. In: **2023 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)**, Rhodes Island, Grécia, 2023. Pag 1-5. DOI: <https://doi.org/10.1109/ICASSP49357.2023.10094670>.

DOWNIE, J. Stephen. Music information retrieval. **Annual Review of Information Science and Technology**, v. 37, n. 1, p. 295–340, 2003. DOI: <https://doi.org/10.1002/aris.1440370108>

DOWNIE, J. Stephen et al. The music information retrieval evaluation eXchange: Some observations and insights. **Studies in Computational Intelligence**, v. 274, p. 93–115, 2010. DOI: https://doi.org/10.1007/978-3-642-11674-2_5

EKMAN, Paul. **Emotions revealed: recognizing faces and feelings to improve communication and emotional life**. New York: Times Books, 2003.

FAZENDA, Ivani Catarina Arantes. Interdisciplinaridade Transdisciplinaridade: visões culturais e epistemológicas e as condições de produção. **Revista Interdisciplinaridade**, v. 1, n. 2, 2012. Disponível em: <https://revistainterdisciplinaridade.com/article/view/2> Acesso em: 12/12/23

FESSAHAYE, Ferdos et al. T-recsys: A novel music recommendation system using deep learning. In: **2019 IEEE international conference on consumer electronics (ICCE)**. IEEE, 2019. p. 1-6.

FLORIDO, Irapuru H.; RAITTZ, Roberto T. Hybrid method for automatic music labeling. In: **2018 XLIV Latin American Computer Conference (CLEI)**, 1-5 out. 2018, São Paulo, Brasil. Anais [...]. IEEE, 2018. p. [indicar páginas]. ISBN 978-1-7281-0437-9. DOI: <https://doi.org/10.1109/CLEI.2018.00038>

FLORIDO, Irapuru H. et al. Estratégias de metacognição na literatura e música sob a ótica da andragogia. **IF-Sophia: Revista eletrônica de investigações Filosófica, Científica e Tecnológica**, v. 8, n. 25, 15 jul. 2023.

FLORIDO, Irapuru H.; DE PAULA PINTO, José S.; RAITTZ, Roberto T.; MACHADO, Diogo J. Ontologia de Letras de Músicas Brasileiras Alimentada por Crowdsourcing. **The Advances in Knowledge Representation (AKR)**, nov. 2025.

GINSBURG, Geoffrey S.; WILLARD, Huntington F. Genomic and personalized medicine: foundations and applications. **Translational Research**, v. 154, n. 6, p. 277–287, dez. 2009. Disponível em: [https://www.translationalres.com/article/S1931-5244\(09\)00274-6/fulltext](https://www.translationalres.com/article/S1931-5244(09)00274-6/fulltext). DOI: 10.1016/j.trsl.2009.09.005. Acesso em: 12 dez. 2023.

GOLLERY, Martin. Bioinformatics: Sequence and Genome Analysis, David W. Mount. **Clinical Chemistry**, v. 51, n. 11, p. 2219, nov. 2005. DOI: 10.1373/clinchem.2005.053850.

GOTHAM, Mark; BEMMAN, Brian; VATOLKIN, Igor. Towards an ‘Everything Corpus’: A Framework and Guidelines for the Curation of More Comprehensive Multimodal Music Data. **Transactions of the International Society for Music Information Retrieval**, v. 8, n. 1, p. 70–92, 5 maio 2025. DOI: 10.5334/tismir.228.

GUPTA, Aaryan; DENGRE, Vinya; KHERUWALA, Hamza Abubakar; SHAH, Manan. Comprehensive review of text-mining applications in finance. **Financial Innovation**, v. 6, art. 39, 2020. DOI: 10.1186/s40854-020-00205-1.

HEARST, Marti. **What is Text Mining?** 2022. Disponível em: <https://people.ischool.berkeley.edu/~hearst/text-mining.html>. Acesso em: 10 ago. 2022.

HORNER, David Stephen; PAVESI, Giulio; CASTRIGNANÒ, Tiziana; D’ONORIO DE MEO, Paolo; LIUNI, Sabino; SAMMETH, Michael; PICARDI, Ernesto; PESOLE, Graziano. Bioinformatics approaches for genomics and post genomics applications of next-generation sequencing. **Briefings in Bioinformatics**, v. 11, n. 2, p. 181–197, mar. 2010. DOI: 10.1093/bib/bbp046.

HOTH, Andreas; NÜRNBERGER, Andreas; PAASS, Gerhard. A Brief Survey of Text Mining. **Journal for Language Technology and Computational Linguistics**, v. 20, n. 1, p. 19–62, jul. 2005. DOI: 10.21248/jlcl.20.2005.68.

HUTTO, C. J.; GILBERT, Eric. VADER: A Parsimonious Rule-Based Model for Sentiment Analysis of Social Media Text. In: **Eighth International AAI Conference on Weblogs and Social Media (ICWSM)**, 2014. p. 216–225. DOI: 10.1609/icwsml.v8i1.14550.

IEGER-RAITZ, Rosangela et al. What are we learning with Yoga? Mapping the scientific literature on Yoga using a vector-text-mining approach. **PLOS ONE**, v. 20, n. 5, e0322791, maio 2025. DOI: <https://doi.org/10.1371/journal.pone.0322791>.

INTERNATIONAL FEDERATION OF THE PHONOGRAPHIC INDUSTRY. **Global Music Report 2025**. London: IFPI, 2025. Disponível em: https://ifpi-website-cms.s3.eu-west-2.amazonaws.com/GMR_2025_State_of_the_Industry_Final_83665b84be.pdf. Acesso em: 10 dez. 2025.

JAPIASSU, Hilton. **A questão da interdisciplinaridade**. Porto Alegre: Secretaria Municipal de Educação, jul. 1994. Disponível em: <http://smeduquedecaxias.rj.gov.br>.

JIA, Xiaoguang. Music Emotion Classification Method Based on Deep Learning and Improved Attention Mechanism. **Computational Intelligence and Neuroscience**, v. 2022, Article ID 5181899, 2022. DOI: <https://doi.org/10.1155/2022/5181899>.

JO, Wonkwang; KIM, M. Justin. Tracking emotions from song lyrics: Analyzing 30 years of K-pop hits. **Emotion**, v. 23, n. 6, p. 1658–1669, 2022. DOI: <https://doi.org/10.1037/emo0001185>.

JOYCE, John. Pandora and the Music Genome Project. **Scientific Computing**, v. 23, set. 2006. Disponível em: https://www.researchgate.net/publication/295343382_Pandora_and_the_music_genome_project_Song_structure_analysis_tools_facilitate_new_music_discovery.

KIM, Minho; KWON, Hyuk Chul. Lyrics-based emotion classification using feature selection by partial syntactic analysis. In: **23rd IEEE International Conference on Tools with Artificial Intelligence (ICTAI)**, 2011. DOI: 10.1109/ICTAI.2011.174.

KUHN, Thomas S.; HACKING, Ian. **The Structure of Scientific Revolutions**. 4. ed. Chicago: University of Chicago Press, 2012.

LAROCCA NETO, Joel et al. Generating text summaries through the relative importance of topics. In: MONARD, M. C.; SICHMAN, J. S. (org.). **IBERAMIA-SBIA 2000: 7th Ibero-American Conference on Artificial Intelligence and 15th Brazilian Symposium on Artificial Intelligence**. Lecture Notes in Artificial Intelligence, v. 1952, p. 300–309, 2000. DOI: https://doi.org/10.1007/3-540-44399-1_31.

LEVITT, Michael. The birth of computational structural biology. **Nature Structural Biology**, v. 8, n. 5, p. 392–393, 2001. DOI: <https://doi.org/10.1038/87545>. Acesso em: 18 dez. 2023.

LIBBRECHT, Maxwell W.; NOBLE, William Stafford. Machine learning applications in genetics and genomics. **Nature Reviews Genetics**, v. 16, p. 321–332, 2015. DOI: <https://doi.org/10.1038/nrg3920>. Acesso em: 18 dez. 2025.

LOWE, M. Sara et al. The boolean is dead, long live the boolean! natural language versus boolean searching in introductory undergraduate instruction. **College and Research Libraries**, v. 79, n. 4, 2018. DOI: <https://doi.org/10.5860/crl.79.4.496>.

MA, Liye; SUN, Baohong. Machine learning and AI in marketing – Connecting computing power to human insights. **International Journal of Research in Marketing**, v. 37, n. 3, p. 481–504, 2020. DOI: <https://doi.org/10.1016/j.ijresmar.2020.04.005>.

MACHADO, Diogo de Jesus Soares et al. Biotext: Exploiting Biological-Text Format for Text Mining. **bioRxiv**, 2022. DOI: <https://doi.org/10.1101/2021.04.08.439078>.

MARIANO, Diego et al. Introducing Programming Skills for Life Science Students. **Biochemistry and Molecular Biology Education**, v. 47, n. 3, p. 296–304, 2019. DOI: <https://doi.org/10.1002/bmb.21228>.

MAY, Lloyd; CASEY, Michael. Familiar Feelings: Listener-Rated Familiarity in Music Emotion Recognition. In: **21st International Society for Music Information Retrieval Conference (ISMIR)**, 2020. Disponível em: <https://archives.ismir.net/ismir2020/paper/000019.pdf>. Acesso em: 18 dez. 2023.

MIKOLOV, Tomas et al. Efficient estimation of word representations in vector space. In: **International Conference on Learning Representations (ICLR)**, 2013. Disponível em: <https://arxiv.org/abs/1301.3781>.

MORI, Elisa; STROETER, Guga. **Uma árvore da música brasileira**. São Paulo: Edições Sesc, 2020.

MORIN, Edgard. **Os sete saberes necessários à educação do futuro**. 2. ed. São Paulo: Cortez, 2013.

OKADA, Keisuke; TAN, Phan Xuan; KAMIOKA, Eiji. Five-Factor Musical Preference Prediction for Solving New User Cold-Start Problem in Content-Based Music Recommender System. In: **2021 International Conference on Technologies and Applications of Artificial Intelligence (TAAI)**, 2021. DOI: <https://doi.org/10.1109/TAAI53890.2021.00041>.

ORAMAS, S.; GOUYON, F.; HOGAN, S.; LANDAU, C.; EHMANN, A. MGPHot: A Dataset of Musicological Annotations for Popular Music (1958-2022). **Transactions of the International Society for Music Information Retrieval**, v. 8, n. 1, p. 108–120, 2025. DOI: <https://doi.org/10.5334/tismir.229>.

PACHET, François et al. Editorial metadata in electronic music distribution systems: Between universalism and isolationism. **Journal of New Music Research**, v. 34, n. 2, p. 103–116, 2005. DOI: <https://doi.org/10.1080/09298210500124218>.

PARRA, Fabio Leonardo; LEÓN, Elizabeth. Unsupervised tagging of Spanish lyrics dataset using clustering. In: **14th International Society for Music Information Retrieval Conference (ISMIR)**, 2013. Disponível em: <https://archives.ismir.net/ismir2013/paper/000019.pdf>.

PENNINGTON, Jeffrey; SOCHER, Richard; MANNING, Christopher D. GloVe: Global vectors for word representation. In: **Proceedings of the 2014 Conference on Empirical Methods in Natural Language Processing (EMNLP)**, Doha, Qatar, 2014. p. 1532–1543. DOI: <https://doi.org/10.3115/v1/D14-1162>.

PERALTA, Elsa. Clara Bertrand Cabral - Patrimônio Cultural Imaterial: Convenção da Unesco e seus Contextos. **Midas**, n. 2, 2013. DOI: <https://doi.org/10.4000/midas.334>.

PIAGET, Jean. Classification of disciplines and interdisciplinary connexions. **International Social Science Journal**, v. 16, n. 4, p. 553–570, 1964.

POWERS, D. M. W. Evaluation: from precision, recall and f-measure to roc, informedness, markedness & correlation. **Journal of Machine Learning Technologies**, v. 2, n. 1, p. 37–63, 2011. Disponível em: https://www.researchgate.net/publication/220766891_Evaluation_From_Precision_R

ecall_and_F-Measure_to_ROC_Informedness_Markedness_Correlation. Acesso em: 18 dez. 2023.

REN, Fujii; BRACEWELL, David B. Advanced Information Retrieval. **Electronic Notes in Theoretical Computer Science**, v. 225, p. 57–72, 2009. DOI: <https://doi.org/10.1016/j.entcs.2008.12.057>.

SALTON, Gerard; MCGILL, Michael J. **Introduction to modern information retrieval**. New York: McGraw-Hill, 1987.

SALUJA, Vasu; JAIN, Minni; YADAV, Prakarsh. L,M&A: An algorithm for music lyrics mining and sentiment analysis. In: **2019 10th International Conference on Computing, Communication and Networking Technologies (ICCCNT), 2019**. DOI: <https://doi.org/10.1109/ICCCNT45670.2019.8944747>

SCHEDL, Markus et al. Current challenges and visions in music recommender systems research. **International Journal of Multimedia Information Retrieval**, v. 7, n. 2, p. 95–116, 2018. DOI: <https://doi.org/10.1007/s13735-018-0154-2>.

SCHEDL, Markus. Deep Learning in Music Recommendation Systems. **Frontiers in Applied Mathematics and Statistics**, v. 5, art. 44, 2019. DOI: <https://doi.org/10.3389/fams.2019.00044>.

SCHERZINGER, Martin. Toward a history of digital music: New technologies, business practices and intellectual property regimes. In: COOK, Nicholas; CLARKE, Eric F.; LEWIS, Daniel; RODGERS, Tara (ed.). **The Cambridge Companion to Music in Digital Culture**. Cambridge: Cambridge University Press, 2019. p. 46–66. DOI: <https://doi.org/10.1017/9781108325906.004>.

SÊNECA, Lucius Annaeus. **Cartas de um estoico**, volume 1: um guia para a vida feliz. São Paulo: [Editora], 2017.

SHRUTI, Chandrayan; PRIYANKA, Bamne. A brief survey of Text Mining and its applications. **International Journal of Emerging Trends in Engineering Research**, v. 9, n. 8, p. 1043–1048, 2021. DOI: <https://doi.org/10.30534/ijeter/2021/25982021>.

SIEVERS, Fabian; HIGGINS, Desmond G. Clustal Omega for making accurate alignments of many protein sequences. **Protein Science**, v. 27, n. 1, p. 135–145, 2018. DOI: <https://doi.org/10.1002/pro.3290>.

SILVEIRA, Marcus Bernardes de Oliveira. A tradição em dois escopos: patrimônio cultural e sambas de roda. **Antíteses**, v. 8, n. 16, p. 279–299, 18 jan. 2016. DOI: <https://doi.org/10.5433/1984-3356.2016v8n16p279>.

SONG, Jonah. The Evolution and Impact of Streaming Services: Changing the Media Landscape. **Global Media Journal**, v. 22, n. 72, 2024.

SORDO, Mohamed et al. Inferring Semantic Facets of a Music Folksonomy with Wikipedia. **Journal of New Music Research**, v. 42, n. 4, p. 346–363, dez. 2013. DOI: <https://doi.org/10.1080/09298215.2013.848035>.

SOUZA, Bruno A. et al. Lematização versus Stemming. **Proceedings of the 22nd Brazilian Symposium on Multimedia and the Web**, v. 48, n. 2, p. 1043–1048, 2017. DOI: <https://doi.org/10.30534/ijeter/2021/25982021>.

SRIVASTAVA, Abhishek et al. Melody Generation from Lyrics Using Three Branch Conditional LSTM-GAN. In: **2022 International Conference on Innovations in**

Computational Intelligence and Computer Vision (ICICV), 2022. DOI: <https://doi.org/10.1109/ICICV54921.2022.9786042>.

STRAWN, John; POHLMANN, Ken C. Principles of Digital Audio. **Computer Music Journal**, v. 10, n. 3, p. 86–87, 1986. DOI: <https://doi.org/10.2307/3680212>.

SULLIVAN, J. W. N. **Beethoven: His Spiritual Development**. [S. l.]: Victor Gollancz Ltd, 1936.

SUMITH, N. et al. Sentiment Classification of English and Hindi Music Lyrics Using Supervised Machine Learning Algorithms. In: **2022 2nd International Conference on Advance Computing and Innovative Technologies in Engineering (ICACITE)**, 2022. DOI: <https://doi.org/10.1109/ICACITE53722.2022.9823722>.

THOMA, Myriam V. et al. The Effect of Music on the Human Stress Response. **PLoS ONE**, v. 8, n. 8, e70156, 2013. DOI: <https://doi.org/10.1371/journal.pone.0070156>.

MITCHELL, Tom Michael. **Machine Learning**. New York: McGraw Hill, 1997.

TURNBULL, Douglas R. et al. Combining audio content and social context for semantic music discovery. In: **Proceedings of the 32nd International ACM SIGIR Conference on Research and Development in Information Retrieval (SIGIR)**, 2009. DOI: <https://doi.org/10.1145/1571941.1571992>.

TURNBULL, Douglas R. et al. Exploring Popularity Bias in Music Recommendation Models and Commercial Streaming Services. **arXiv preprint**, 19 ago. 2022. Disponível em: <https://arxiv.org/abs/2208.09332>. Acesso em: 18 dez. 2023.

TZANETAKIS, George; COOK, Perry. **MARSYAS: a framework for audio analysis. Organised Sound**, v. 4, n. 3, p. 169–175, 16 dez. 2000. DOI: <https://doi.org/10.1017/S1355771800003135>. Acesso em: 18 dez. 2023.

VAN RYSEN, Stefaan. An Introduction to Bioinformatics Algorithms by Neil C. Jones and Pavel A. Pevzner. **Leonardo**, v. 39, n. 5, p. 505–506, 2006. DOI: <https://doi.org/10.1162/leon.2006.39.5.505>.

VELANKAR, Makarand; KULKARNI, Parag. Music Recommendation Systems: Overview and Challenges. In: **Signals and Communication Technology**, [s. l.]: Springer, 2020. DOI: https://doi.org/10.1007/978-981-15-6315-9_13.

WANG, Zhuo; CHEN, Yazhu; LI, Yixue. A brief review of computational gene prediction methods. **Genomics, Proteomics & Bioinformatics**, v. 2, n. 4, p. 216–221, 2004. DOI: [https://doi.org/10.1016/S1672-0229\(04\)02024-9](https://doi.org/10.1016/S1672-0229(04)02024-9)

WEBSTER, Jack. Taste in the platform age: music streaming services and new forms of class distinction. **Information Communication and Society**, v. 23, n. 13, p. 1902–1918, 2020. DOI: <https://doi.org/10.1080/1369118X.2019.1670224>.

XU, Liang et al. Using machine learning analysis to interpret the relationship between music emotion and lyric features. **PeerJ Computer Science**, v. 7, e785, 2021. DOI: <https://doi.org/10.7717/peerj-cs.785>.

Apêndice A – Artigo aceito na revista *Advances Knowledge Representation* (AKR)

Ontologia de Letras de Músicas Brasileiras Alimentada por Crowdsourcing

Crowdsourcing Brazilian Song Lyrics Ontology

Resumo: A música, desde os primórdios da civilização, desempenha um papel fundamental como forma de entretenimento e tem relevância significativa para a ciência da informação, especialmente na área de recuperação de informação musical. Embora existam diversos trabalhos na literatura, principalmente em língua inglesa, observa-se uma carência de estudos voltados à Música Popular Brasileira (MPB). O objetivo desta pesquisa é propor e desenvolver uma ontologia de domínio para letras da MPB a partir de uma base de dados de larga escala, originada de sítios brasileiros, e contextualizar os desafios e oportunidades das informações provenientes do método de crowdsourcing aplicado nesses sítios. O método de formalização do conhecimento envolve o uso da linguagem Ontology Web Language (OWL) e a aplicação do guia Ontology Development 101 para o desenvolvimento de ontologias. São ainda analisados os aspectos relacionados ao uso de crowdsourcing em plataformas colaborativas online para preencher as instâncias da ontologia, com foco na garantia da qualidade das informações. A ontologia proposta abrange elementos essenciais, como título, artista, gênero, letra, compositor e ano de lançamento, com ênfase na letra, a fim de possibilitar a extração de termos-chave para análises futuras de grandes volumes de músicas. Adicionalmente, o artigo apresenta os resultados detalhados do pré-processamento de milhares de letras da MPB, identificando problemas como erros linguísticos, atribuições incorretas de autoria e a presença de cifras no texto lírico. Conclui-se que a ontologia é um artefato eficaz para estruturar o conhecimento musical, e que o uso do crowdsourcing para sua instanciação, embora promissor, apresenta desafios significativos de qualidade de dados que foram identificados no estudo.

Abstract: Music, since the dawn of civilization, plays a fundamental role as a form of entertainment and holds significant relevance for information science, especially in the area of music information retrieval. Although there are several works in the literature, primarily in English, there is a noticeable lack of studies focused on Brazilian Popular Music (MPB). The objective of this research is to propose and develop a domain ontology for MPB lyrics based on a large-scale database originated from Brazilian websites, and to contextualize the challenges and opportunities presented by information derived from the crowdsourcing method applied by these websites. The method for knowledge formalization involves using the Ontology Web Language (OWL) and following the Ontology Development 101 guide. Furthermore, aspects of crowdsourcing on collaborative online platforms for populating ontology instances are analyzed, with a focus on ensuring information quality. The proposed ontology encompasses essential elements such as title, artist, genre, lyrics, composer, and release year, with an emphasis on the lyrics to enable the extraction of keywords for future analyses of large volumes of songs. Additionally, the paper presents detailed results of the pre-processing of thousands of MPB lyrics, identifying problems such as linguistic errors, incorrect authorship attribution, and the presence of chords in the lyrical text. It is concluded that the ontology is an effective artifact for structuring musical knowledge, and that the use of crowdsourcing for its instantiation, while promising, presents significant data quality challenges that were identified in the study.

1. Introdução

A música, enquanto manifestação cultural fundamental da experiência humana, possui relevância intrínseca em âmbito social, emocional e econômico. Sua importância ultrapassa o mero entretenimento, estabelecendo-se como um poderoso meio de comunicação e expressão. O mercado global de música demonstra uma expressiva vitalidade, como evidenciado pelo "Global Music Report 2024" da International Federation of the Phonographic Industry (IFPI), referente a 2023, que registrou um crescimento robusto de 10,2%, atingindo US\$ 28,6 bilhões, impulsionado principalmente pelo streaming, responsável por 67,3% da receita total.

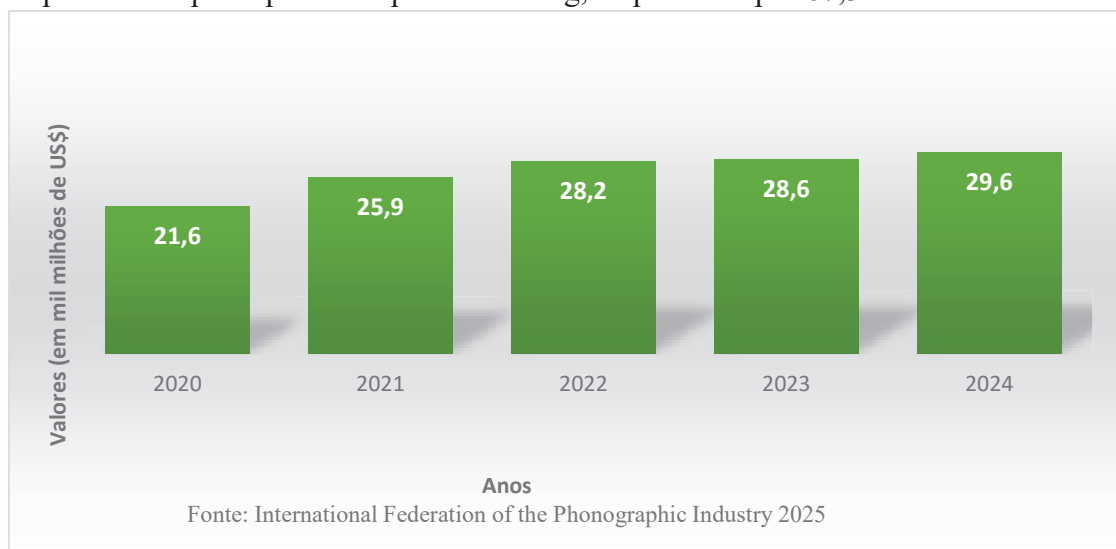


Figura 1- Faturamento dos 5 últimos anos da Indústria da Música Mundial

A tendência de crescimento e expansão, refletida pela crescente digitalização e consumo musical em diversas plataformas, projeta-se para o futuro, conforme indicado no relatório "Global Music Report 2025" (com dados de 2024), ver Figura 1. Nesse contexto, o entretenimento musical não apenas enriquece o tecido cultural das sociedades, mas também dinamiza uma complexa cadeia produtiva.

Entretanto, a organização e a recuperação eficiente do vasto volume de informações associadas às obras musicais, como metadados, letras e contextos históricos e culturais, representam um desafio crescente. Uma ontologia é definida como uma especificação formal e explícita de uma conceitualização compartilhada, possibilitando a representação do conhecimento de um domínio de forma compreensível tanto para humanos quanto para máquinas (Gruber, 1993; Rashid et al., 2018). Para a formalização da ontologia de letras de músicas aqui proposta, optou-se pela Web Ontology Language (OWL), uma linguagem padrão recomendada pelo W3C que oferece elevado poder expressivo para descrever classes, propriedades e indivíduos (W3C OWL Working Group, 2012).

Considerando o imenso volume de dados musicais, especialmente no contexto da rica produção brasileira, o crowdsourcing surge como uma estratégia viável para alimentar e enriquecer a ontologia. Essa abordagem, baseada na colaboração em massa, permite que uma ampla comunidade de usuários contribua com dados, conhecimentos e validações, viabilizando a construção de bases de conhecimento abrangentes.

Diante deste cenário, a presente pesquisa busca responder à seguinte pergunta: Como estruturar formalmente o conhecimento sobre letras da MPB por meio de uma ontologia de

domínio e quais são os desafios e as oportunidades de utilizar o crowdsourcing para seu populamento?

Dessa forma, o presente artigo tem como objetivo delinear uma proposta de ontologia de domínio fundamental para letras de músicas, detalhando seus componentes essenciais (título, artista, gênero, letra, compositor, ano de lançamento, termos principais) e discutindo a aplicação do crowdsourcing nessa população, sem deixar de analisar criticamente os desafios inerentes a essa metodologia.

As ontologias musicais representam um arcabouço valioso para aprofundar estudos sobre as dimensões culturais, regionais e temporais da música, além de consolidar um acervo de conhecimento relevante para as áreas de musicologia, gestão do conhecimento e ciência da informação.

2. Referencial Teórico

A aplicação de ontologias na área musical demanda alguns aspectos essenciais para a associação a métodos de colaboração de informações (crowdsourcing). A seguir, serão abordados os conceitos desta temática para contextualizar a proposta da ontologia.

2.1 Método da Revisão de Literatura

A construção deste trabalho foi fundamentada em uma breve revisão da literatura, e a busca por referenciais teóricos foi realizada em bases de dados acadêmicas, como Periódicos Capes, Google Scholar e Dimensions. Foram utilizados os seguintes descritores para a busca nos repositórios indexados: "ontologia", "letras de música", "crowdsourcing", "music information retrieval" e "music ontology".

O recorte temporal estabelecido foi do período dos últimos 15 anos, incluindo trabalhos semanais clássicos e os idiomas considerados foram o português e inglês.

2.2 Ontologias e Recuperação de informação

Ao longo da história do desenvolvimento da inteligência artificial de nível humano, diferentes paradigmas competiram pela predominância. A inteligência artificial simbólica (IAS) prevaleceu durante grande parte do século XX, mas atualmente o paradigma conexionista, baseado no aprendizado de máquina com redes neurais profundas, vem ganhando destaque. Contudo, ambos os paradigmas apresentam vantagens e limitações, e um grande desafio atual na área é promover a integração entre eles. Um pilar central do paradigma simbólico é que a inteligência resulta da manipulação de representações composicionais abstratas, com elementos que representam objetos e suas inter-relações. Se essa premissa estiver correta, um objetivo essencial para o aprendizado profundo é desenvolver arquiteturas capazes de detectar objetos e relações em dados brutos e aprender a representá-los de maneira eficaz para processamentos futuros. (Garnelo & Shanahan, 2019).

O uso recente de "ontologias" na inteligência artificial simbólica (IAS) reflete um retorno evidente à visão aristotélica do mundo: ao empregar e compartilhar ontologias, assume-se que os sistemas de informação têm categorizações diretas e naturais da realidade, que, para serem compartilhadas, precisam, no mínimo, ser traduzíveis em termos de seus conceitos (Cassapo, 2004). Em ciência da computação e sistemas de informação, uma ontologia é uma especificação formal e explícita de uma conceitualização compartilhada (Gruber, 1993). Noy e McGuinness (2001) definem uma Ontologia como “uma descrição explícita formal de conceitos num domínio do discurso, classes, propriedades de cada conceito descrevendo características, atributos dos conceitos e restrições sobre as propriedades” (McGuinness, 2017).

Em termos mais simples, ontologia é um modelo de dados que representa um conjunto de conceitos em um domínio e as relações entre eles. As ontologias são utilizadas para promover um entendimento comum de uma área de conhecimento, permitindo a interoperabilidade entre sistemas e a realização de inferências lógicas. Elas fornecem um vocabulário compartilhado e uma estrutura semântica que podem ser utilizados para modelar um domínio, especificando os tipos de objetos, conceitos, suas propriedades e relações (Gruber, 1993; Guarino & Giaretta, 1995).

2.3 Linguagem de Representação de Conhecimento (OWL)

A Web Semântica é uma tendência para o futuro da internet, na qual as informações possuem significados explícitos, facilitando o processamento e a integração automática de dados por máquinas. Ela se baseia na capacidade do XML de criar esquemas de marcação personalizados e na abordagem flexível do Resource Description Framework (RDF) para a representação de dados. O primeiro nível acima do RDF necessário para a Web Semântica é uma linguagem de ontologia capaz de descrever formalmente o significado da terminologia empregada em documentos web. Para que as máquinas realizem tarefas de raciocínio úteis nesses documentos, a linguagem deve ir além da semântica básica do RDF Schema (McGuinness et al., 2004).

A OWL (Web Ontology Language) é uma família de linguagens de representação de conhecimento para a criação de ontologias, recomendada pelo W3C (W3C OWL Working Group, 2012). Construída sobre o RDF, a OWL oferece uma sintaxe formal e uma semântica rica para a descrição de classes (conjuntos de indivíduos), propriedades (relações entre indivíduos e atributos dos indivíduos) e indivíduos (objetos do domínio). A OWL permite a definição de restrições, cardinalidade, características das propriedades (como transitividade e simetria) e relações complexas entre classes (como disjunção e equivalência), o que a torna uma ferramenta poderosa para a modelagem semântica e para o raciocínio automatizado sobre o conhecimento do domínio.

2.4 Músicas e Crowdsourcing

O crowdsourcing consiste em obter serviços, ideias ou conteúdo por meio da solicitação de contribuições de um grande grupo de pessoas, especialmente de uma comunidade online (Howe, 2006). Este modelo baseia-se na premissa de que a inteligência coletiva pode superar a de especialistas individuais em determinadas tarefas, especialmente aquelas que envolvem grande volume de dados ou subjetividade na interpretação.

Conforme Estellés-Arolas e González-Ladrón-de-Guevara (2012), crowdsourcing é uma atividade online participativa em que uma entidade (e.g., empresa, indivíduo) propõe, via chamada aberta, a uma multidão heterogênea a realização voluntária de uma tarefa, resultando em benefícios mútuos.

De acordo com a revisão sistemática de Estellés-Arolas e González-Ladrón-de-Guevara (2012), foram identificados oito características essenciais (*differentia specifica*) que definem qualquer iniciativa de crowdsourcing:

1. Multidão (Crowd): Um grupo grande e diverso de indivíduos com diferentes níveis de conhecimento, que não necessariamente se conhecem entre si.
2. Tarefa: Uma atividade com um objetivo claro, de complexidade e modularidade variáveis, que pode envolver trabalho, dinheiro (como no crowdfunding), conhecimento ou experiência.
3. Recompensa para a Multidão: Satisfação de necessidades, como recompensa econômica, reconhecimento social, autoestima ou desenvolvimento de habilidades (Maslow, 1943).

4. Iniciador (Crowdsourcer): Pode ser uma empresa, instituição, organização sem fins lucrativos ou indivíduo que propõe a tarefa.
5. Benefício do Iniciador: O crowdsourcer obtém a solução de um problema, conhecimento, ideias ou valor agregado a partir do trabalho da multidão.
6. Tipo de Processo: Um processo participativo, distribuído e online.
7. Chamada Aberta: Uma convocação flexível que pode ser totalmente aberta ou restrita a uma comunidade com conhecimentos específicos.
8. Meio Utilizado: A internet, frequentemente associada às tecnologias da Web 2.0.

No contexto da catalogação musical e enriquecimento de ontologias, o crowdsourcing se propõe a ser uma estratégia eficaz para a coleta e validação de metadados, transcrição de letras, identificação de temas, anotação de emoções e classificação de gêneros, aproveitando o conhecimento distribuído de fãs, músicos e pesquisadores.

2.5 Ontologias e Letras de Músicas

A aplicação de ontologias no domínio musical tem se mostrado uma área em expansão para pesquisa e desenvolvimento, especialmente no que tange à organização, recuperação e análise de informações complexas associadas às obras musicais. Quando o foco se volta para as letras das músicas, as ontologias oferecem um potencial significativo para identificar as camadas de significado, contexto e estrutura.

A interoperabilidade dos metadados musicais é um desafio considerável, e as ontologias emergem como soluções para lidar com a heterogeneidade desses dados, permitindo o alinhamento, a integração e o acesso unificado a diversos conjuntos de dados (de Berardinis et al., 2023). Em seu trabalho, De Berardinis (2023) propõe um modelo semântico flexível para metadados relacionados a artistas, composições, performances e gravações, o qual é altamente relevante para o estudo de letras ao permitir a integração de metadados contextuais que informam o significado lírico, como gênero, intenção do compositor ou período histórico (de Berardinis et al., 2023). Sua capacidade de considerar a proveniência das informações é crucial para lidar com múltiplas interpretações do significado de uma letra.

A Music Ontology, apresentada por Raimond et al (2007), é um framework formal amplamente utilizado, servindo como base para informações editoriais, culturais e acústicas (Raimond et al., 2007). Para letras, ela permite vincular informações sobre o compositor, o período histórico ou o gênero. No entanto, sua generalidade pode ser uma limitação para capturar nuances específicas de letras em contextos culturais diversos.

No trabalho de Proutskova et al. (2020), destaca-se a necessidade de extensões para uma Ethno-Music-Ontology, importante para letras em tradições musicais não ocidentais, que frequentemente carregam significados culturais, espirituais ou históricos específicos que modelos padrão podem não abranger adequadamente, como estruturas poéticas tradicionais ou o papel da letra em rituais (Proutskova et al., 2020).

2.6 Vantagens do Uso de Ontologias para Letras de Músicas

No contexto do uso de ontologias para a geração de conhecimento e a recuperação de informação musical, reforçam-se as vantagens do uso da ontologia, de modo a facilitar, compartilhar e estudar todo um corpus musical da música popular brasileira que não se encontra disponível em repositórios e na web. Em seguida, ressaltamos as vantagens do uso de uma ontologia, dirigida a esta área de conhecimento.

- Contextualização Semântica Aprofundada: Ontologias permitem vincular letras a um rico conjunto de metadados (artista, álbum, ano, contexto de criação, performance),

enriquecendo a análise de seu significado, a Music Meta Ontology exemplifica essa capacidade (de Berardinis et al., 2023).

- **Modelagem de Aspectos Culturais e Estilísticos:** É possível modelar como as letras refletem e são influenciadas por contextos culturais, gêneros e estilos específicos (relevância da Ethno-Music-Ontology e da Genre Ontology Learning, como discutido por (Schreiber, 2016). A Genre Ontology Learning, ao aprender ontologias de gêneros a partir de etiquetas coletivas, pode ajudar a mapear relações semânticas entre letras e gêneros, identificando temas e estilos linguísticos característicos.

- **Integração com Elementos Sonoros e Performativos:** Embora algumas ontologias se concentram primariamente em áudio, elas podem indiretamente beneficiar a análise de letras. A Audio Feature Ontology (Allik et al., 2016) e a Audio Effects Ontology (Wilmering et al., 2013), juntamente com o Studio Ontology Framework (Fazekas & Sandler, 2011), podem ajudar a entender como características acústicas e efeitos de áudio são usados para reforçar o impacto emocional ou destacar partes específicas das letras. A Ontologia de Instrumentos Musicais (Kolozali et al., 2011) também pode ser relevante ao conectar o timbre de instrumentos a conotações simbólicas presentes nas letras.

- **Análise Semântica de Alto Nível e Raciocínio:** Métodos como o proposto em "Predicting High-level Music Semantics Using Social Tags via Ontology-based Reasoning" são diretamente aplicáveis às letras para prever humor, tema ou uso com base em etiquetas sociais, enriquecendo a base de conhecimento por meio de léxicos como o WordNet (Wang et al., 2010).

- **Suporte à Educação e Anotação Detalhada:** Ontologias como a Musical Forms and Structures Ontology (MFSO) e a Musical Performance Ontology (MPO) (Sébastien et al., 2013) podem estruturar anotações que conectam o texto lírico à forma musical e à performance, descrevendo como a estrutura de uma canção (verso-refrão) reflete a narrativa da letra ou como a entrega vocal afeta a interpretação do texto.

- **Descoberta e Recomendação Aprimoradas:** Ao estruturar o conhecimento sobre o conteúdo lírico, ontologias podem aprimorar sistemas de recomendação musical que considerem temas, emoções ou estilos linguísticos semelhantes nas letras.

2.7 Desafios no Uso de Ontologias para Letras de Músicas

Assim como existem os bônus da aplicação de determinados métodos em contrapartida identificamos também alguns aspectos de desafios as serem superados no contexto do uso da ontologia:

- **Captura da Complexidade e Subjetividade Lírica:** Letras de músicas frequentemente contêm ambiguidades, metáforas, ironia e uma vasta gama de emoções sutis. Modelar essa complexidade e a subjetividade da interpretação é um desafio significativo para uma representação ontológica formal. **Necessidade de Extensões Específicas:** Muitas ontologias musicais existentes, embora úteis, não foram primariamente concebidas com foco profundo nas letras. A Music Ontology, por exemplo, é abrangente, mas pode necessitar de extensões para detalhar aspectos específicos da semântica, da estrutura e do simbolismo lírico. Ontologias focadas em áudio também requerem módulos ou extensões para se conectarem significativamente ao texto (Raimond et al., 2007);

- **Variação Cultural e Linguística:** O significado e o uso de certos termos, temas ou dispositivos estilísticos nas letras podem variar enormemente entre diferentes culturas e línguas. Desenvolver ontologias que sejam sensíveis a essas variações (como proposto pela Ethno-Music-Ontology) é crucial, mas complexo (Proutskova et al., 2020);

- **Integração de Múltiplas Dimensões:** Uma análise completa de letras pode requerer a integração de informações de diversas fontes e naturezas: o texto em si, o contexto do artista e da obra, a performance, os elementos sonoros da música e a recepção pelo público. Unificar essas dimensões num modelo ontológico coerente é uma tarefa complexa;

- Limitações de Ontologias de Gênero Musicais: Como aponta a pesquisa em Genre Ontology Learning, ontologias de gênero existentes podem apresentar grande desconexão (Schreiber, 2016). As letras podem, de fato, oferecer um meio de conectar subgêneros ou identificar influências cruzadas que as classificações formais de gênero não capturam.

- Escassez de Ontologias Dedicadas a Dispositivos Literários em Letras: Embora a análise de dispositivos literários seja comum em estudos de letras, faltam ontologias robustas que formalizem o conhecimento dessas Figuras de linguagem e de seu impacto específico no contexto musical.

Em síntese, com base nas referências consultadas, resume-se a situação: para uma análise abrangente das letras, seria necessário integrar diferentes perspectivas ontológicas (metadados, aspectos culturais, sonoros, performativos, semântica textual) num framework unificado, possivelmente com extensões específicas para modelar a estrutura e o simbolismo intrínseco às letras. O desafio reside em criar modelos que sejam, ao mesmo tempo, expressivos o suficiente para capturar a riqueza das letras e flexíveis para se adaptarem a diferentes contextos e necessidades de análise.

A Ontologia de Música Popular Brasileira (MpbDomPubl), proposta neste trabalho, é classificada como uma ontologia de domínio, pois visa representar o conhecimento específico de um domínio particular: as letras de músicas da MPB. Diferentemente de uma ontologia de alto nível, seu propósito não é definir conceitos genéricos. Por essa razão, nesta etapa da pesquisa, optou-se por não alinhar à ontologia de alto nível consolidada, como a BFO. No entanto, reconhece-se a importância dessa prática, e tal alinhamento é considerado um passo fundamental para trabalhos futuros, visando garantir maior consistência e interoperabilidade.

3. Metodologia

O desenvolvimento da ontologia foi fundamentado no método *Ontology Development 101* (McGuinness et al., 2017), seguindo os passos recomendados no guia citado. O protocolo para a construção da MpbDomPubl segue as etapas descritas e consolidadas para o desenvolvimento de ontologias baseadas em OWL.

3.1 Definição do Escopo e Domínio:

- Domínio: Letras de músicas, com foco inicial na Música Popular Brasileira (MPB) de domínio público.
- Escopo: Representar informações factuais sobre músicas (título, artista, compositor, ano de lançamento, gênero) e o conteúdo textual das letras, permitindo a extração e representação dos termos principais.
- Questões de Competência: A ontologia deve ser capaz de responder a perguntas como: "Quais músicas de Chiquinha Gonzaga foram lançadas nos anos 30 e falam sobre 'carnaval'?" "Quais artistas cantam sobre temas políticos no gênero MPB?".

3.2 Enumeração de Termos Importantes:

Identificar os termos-chave do domínio: Título, Artista, Gênero Musical, Letra, Termo (Palavra), Ano de Lançamento.

3.3 Definição das Restrições das Propriedades

Especificar cardinalidade: uma música deve ter pelo menos um título, tipos de dados, domínios e imagens das propriedades.

3.4 Definição Formal das Classes e Propriedades

A seguir, detalham-se as principais classes da MpbDomPubl, seguindo o padrão de definição por gênero próximo e diferença específica ("X é um Y que Z").

3.4.1 Definição Formal das Classes

- Música: É uma entidade (ou obra) que se caracteriza por possuir um título, uma letra, um ou mais intérpretes e compositores, um gênero musical e um ano de lançamento.
- Pessoa: É um agente que está envolvido na criação ou performance de uma música.
- Gênero Musical: É uma categoria que classifica a obra musical (ex.: Choro, Samba, sertanejo, etc.)

A formalização na linguagem OWL (Ontology Web Language) foi realizada por meio de axiomas que definem as relações (Propriedades de Objeto) entre classes. Por exemplo, a propriedade `temArtista` tem como Domínio a classe Música e como Imagem a classe Pessoa (ou Artista, dependendo de como você definiu o domínio direto), o que garante a consistência e a coerência do modelo.

3.4.2 Definição Formal das Propriedades

Descreve-se a seguir as relações e atributos aplicados a ontologia:

- Propriedades de Objeto (são as relações):
 - o `temGenero` (domínio: Música, imagem: Gênero Musical)
 - o `temArtista` (domínio: Música, imagem: Pessoa)
- Propriedades de Dados (Atributos):
 - o `título` (domínio: Música, tipo: `xsd:string`)
 - o `anoDeLançamento` (domínio: Música, tipo: `xsd:integer`)
 - o `artista` (domínio: Pessoa, tipo: `xsd:string`)
 - o `gênero` (domínio: Gênero Musical, tipo: `xsd:string`)
 - o `letra` (domínio: Musica, tipo: `xsd:string`)
 - o `termo` (domínio: Musica, tipo: `xsd:string`)

3.5 Criação dos Indivíduos (Instâncias)

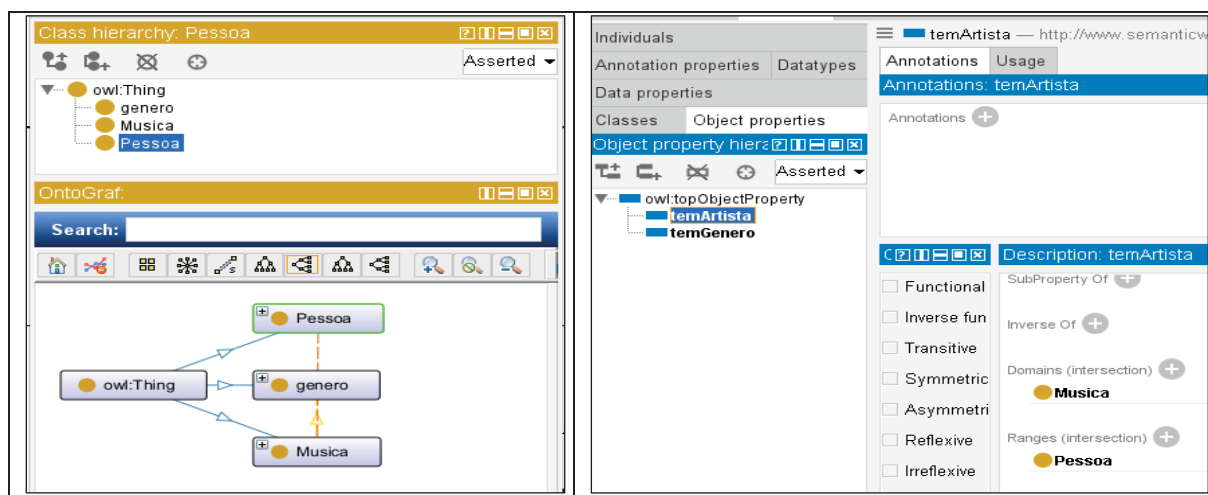
Esta etapa seria massivamente importada de bases de dados disponíveis ou realizada via crowdsourcing, com os usuários da plataforma musical inserindo dados sobre músicas específicas, que seriam mapeados para as classes e propriedades da ontologia.

3.6 Iteração e Avaliação

A ontologia deve ser refinada e avaliada quanto à sua consistência, completude e capacidade de responder às questões de competência.

3.7 Detalhes da Implementação

A implementação da ontologia foi realizada na ferramenta Protégé. As Figuras 2 e 3 detalham os principais aspectos da implementação.



(a) Hierarquia das Classes

(b) Propriedade do Objeto

Figura 2 – Ontologia proposta no framework Protégé.

A hierarquia de classes (Figura 2a) demonstra a especialização de conceitos. As relações entre as classes foram definidas por meio de Propriedades de Objeto, com domínios e imagens estritamente definidos (Figura 2b).

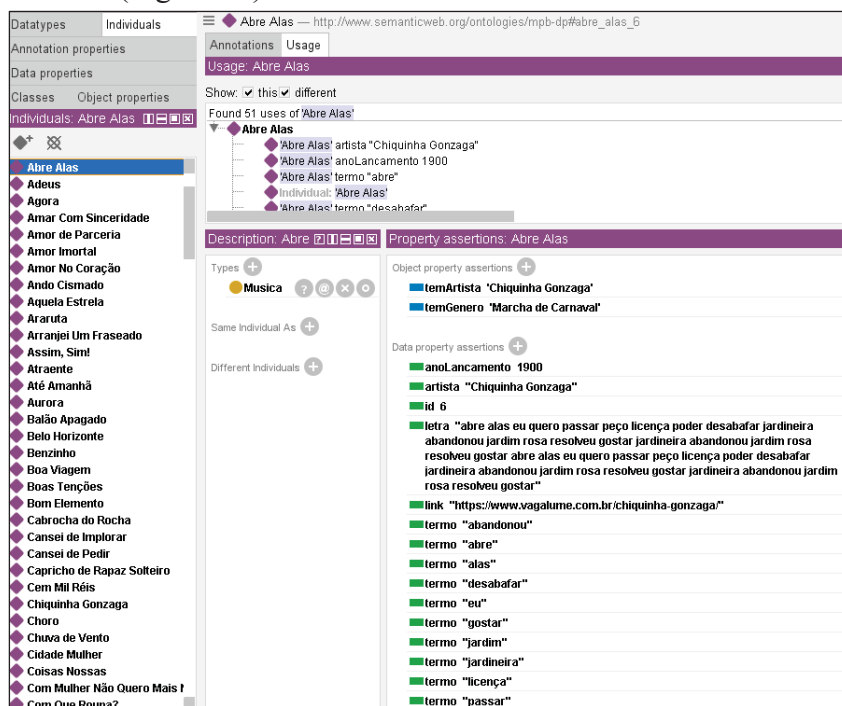


Figura 3 – Exemplo da instanciação de uma música.

Os atributos e restrições lógicas, como as de cardinalidade, foram aplicados para garantir a integridade do modelo. Na Figura 3, temos um exemplo de instauração de uma música da artista “Chiquinha Gonzaga”, intitulada “Abre Alas”, que apresenta também as propriedades das músicas e, adicionalmente, os termos mais significativos presentes na letra da música.

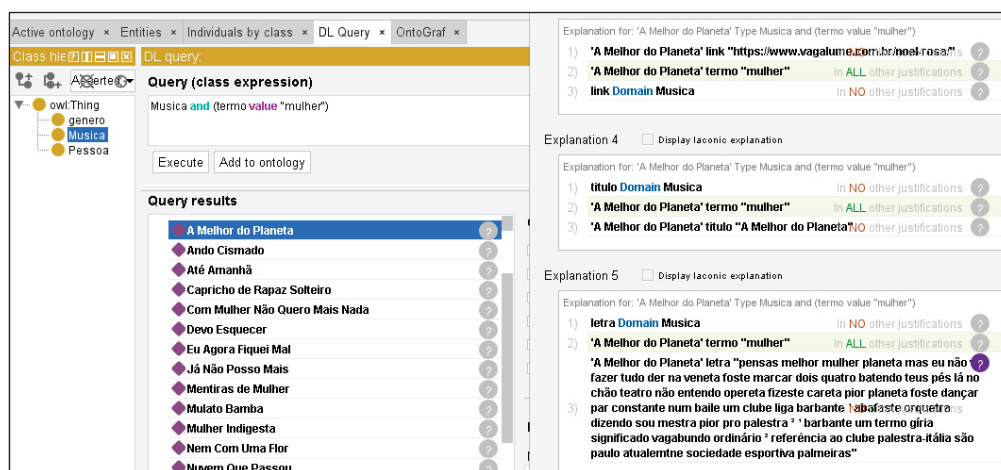


Figura 4 – Exemplo da inferência na ontologia.

Na Figura 4, a ferramenta Protege apresenta uma consulta a ontologia para identificar todas as músicas que continham o termo "mulher" em suas letras, com o comando “Musica and (termo value ‘mulher’””. Na tela abaixo, apresentam-se as músicas associadas ao termo “mulher” e, ao lado, o descritivo em que o termo foi encontrado, por exemplo, no título e/ou na letra.

A ontologia, formalizada na linguagem OWL conforme especificado no Anexo I, foi importada na ferramenta Protégé (McGuinness et al., 2017). Este processo resultou na geração do mapa conceitual e as relações entre os conceitos definidos.

5. Validação da Ontologia e Aderência aos Princípios FAIR

A validação de uma ontologia é uma etapa fundamental para garantir sua qualidade, consistência e utilidade. Para o presente trabalho, foi adotada uma abordagem de validação multifacetada, combinando testes de competência, verificação de consistência lógica e a publicação do artefato seguindo os princípios de dados FAIR (Findable, Accessible, Interoperable, and Reusable).

5.1 Validação por Questões de Competência e Inferência

Uma das formas de avaliar a ontologia é verificar sua capacidade de responder às questões de competência definidas na fase de escopo. Conforme demonstrado na Figura 4, a ontologia foi submetida a consultas (queries) na ferramenta Protégé para testar sua capacidade de inferência e recuperação de conhecimento. Um exemplo foi a execução de uma consulta para identificar todas as músicas que continham o termo "mulher" em suas letras, o que demonstrou que o modelo é capaz de realizar inferências para as quais foi projetado.

5.2 Validação de Consistência Lógica

A consistência lógica da ontologia MpbDomPubl foi verificada por meio de raciocinadores lógicos (reasoners), como o HermiT, integrados à ferramenta Protégé. A execução do raciocinador permitiu validar a coerência dos axiomas, a hierarquia de classes e as restrições definidas, sem identificar inconsistências ou contradições no modelo. Esse procedimento garante que a ontologia é sintaticamente correta e logicamente sólida.

5.3 Publicação e Aderência aos Princípios FAIR

Visando garantir que a ontologia seja encontrável (Findable), acessível (Accessible), interoperável (Interoperable) e reutilizável (Reusable), ela foi publicada em um repositório

online. A ontologia está disponível em seu formato OWL no repositório de dados Zenodo no seguinte endereço: MPBDomPubl - Ontologia da Música Popular Brasileira de Domínio Público

Para promover a reutilização, a ontologia foi disponibilizada sob a licença Creative Commons Attribution 4.0 International (CC BY 4.0), que permite o compartilhamento e a adaptação do material para qualquer finalidade, inclusive comercial, desde que o devido crédito seja atribuído.

6. Discussão e Análise dos Resultados

A ontologia destina-se a pesquisadores das áreas de musicologia e ciência da informação, que poderão acessá-la por meio de repositório público para realizar análises semânticas. O público contribuinte (crowd) seria composto por fãs de música, estudantes e pesquisadores, que alimentariam a base de dados por meio de uma futura plataforma colaborativa.

6.1 Ontologias e inferências.

A grande vantagem do uso de ontologias em determinado domínio de conhecimento, particularmente neste trabalho, é a possibilidade de pesquisar termos presentes nas letras das músicas. Além de permitir o estudo e a compreensão da semântica das obras relacionadas, traz conhecimentos implícitos às músicas. A Figura 4 ilustra como uma consulta simples permite identificar músicas que contêm o termo “mulher” em sua composição.

6.2. Desafios com o Crowdsourcing para Popular uma Ontologia

A utilização de uma plataforma de crowdsourcing para alimentar a ontologia apresenta desafios significativos, principalmente relacionados à qualidade e consistência dos dados inseridos pelos usuários:

- Ambiguidade e Variação: Nomes de artistas e compositores podem ser inseridos de formas diferentes (ex.: "Chiquinha Gonzaga", "Francisca Edviges Neves Gonzaga");
- Informações Incompletas ou Incorretas: Dados como ano de lançamento e autoria podem ser desconhecidos por muitos usuários ou inseridos de forma equivocada;
- Subjetividade na Classificação de Gênero: A definição do gênero musical de uma canção pode ser subjetiva e variar entre os colaboradores;
- Erros de Linguagem: Presença massiva de erros de ortografia, acentuação e pontuação. A inserção manual de letras de músicas invariavelmente leva a erros ortográficos e de pontuação;
- Erros de Autoria: Músicas atribuídas a artistas incorretos, um problema comum em fontes não oficiais;
- Ruído Textual: Presença de termos como "Intro", "Refrão", "Solo" e outros elementos estruturais da canção que não fazem parte da letra em si. Inserção de informações não pertinentes, como cifras e acordes (ex.: "Am", "G", "C"), misturadas ao texto da letra (ex.: Estribillo, 2x, 3x, Repetir); e
- Duplicatas e Versões: Existência de múltiplas versões da mesma letra, com pequenas variações.

6.3. Resultados da Análise de Pré-processamento de Letras de Músicas

Uma análise preliminar de um conjunto de dados com letras de músicas, obtido de duas fontes brasileiras que disponibilizam letras da música popular brasileira: “Vagalume” e “Letras e Músicas”. As letras das músicas nestes sítios, alimentados pelo processo de crowdsourcing,

revelaram a necessidade de um robusto processo de pré-processamento antes da extração de termos e da população da ontologia.

Conforme relatado na seção anterior, os principais problemas identificados são: ambiguidade e variação, informações incompletas ou incorretas, subjetividade na classificação de gênero: a definição do gênero musical de uma canção pode ser subjetiva e variar entre os colaboradores, erros de linguagem, erros de autoria, ruído textual, duplicatas e versões.

A resolução dessas complexidades demandou um investimento considerável de tempo e esforço para garantir que o corpus musical atingisse um nível de qualidade satisfatório, tornando-o adequado à aplicação dos algoritmos de transformação planejados. É importante salientar a reconhecida necessidade de desenvolver algoritmos específicos para a limpeza e normalização dos dados, bem como a implementação de mecanismos de validação na plataforma de crowdsourcing, visando mitigar a inserção de informações de baixa qualidade desde o início da coleta. Após a conclusão da fase de pré-processamento, o corpus final contabilizou 81 mil músicas, o que representa uma redução de aproximadamente 41,3% em relação ao quantitativo inicial de 138 mil letras.

É importante ressaltar que, de acordo com as leis de direitos autorais, as letras das músicas analisadas neste trabalho têm finalidade exclusivamente acadêmica e não comercial. Diante disso, as músicas contidas na ontologia proposta, ou seja, as instancias, referem-se apenas às músicas de domínio público, as obras de autoria dos artistas: Chiquinha Gonzaga, Ernesto Nazareth, Noel Rosa e Zequinha de Abreu.

7. Conclusão

Em suma, a proposta de uma ontologia para letras de músicas brasileiras, alimentada pela metodologia de crowdsourcing, configura-se como uma alternativa inovadora e de grande valia para a organização e análise do nosso vasto patrimônio musical. A utilização da linguagem OWL para estruturar semanticamente este conhecimento permite que as informações musicais sejam não apenas armazenadas de forma eficiente, mas também interpretadas e relacionadas por sistemas inteligentes. Embora os desafios associados à coleta colaborativa de dados e ao pré-processamento necessário das letras sejam consideráveis, os benefícios potenciais são substanciais.

A criação de um acervo musical com essa estrutura semântica abre novos caminhos para a realização de estudos aprofundados sobre a música brasileira em suas dimensões culturais, regionais, temporais e linguísticas. Além disso, estabelece uma base de conhecimento fundamental para o avanço de pesquisas em musicologia, ciência da computação, linguística computacional e humanidades digitais, fomentando novas formas de interação e um entendimento mais completo da riqueza da nossa música. Assim, a implementação desta ontologia para letras da MPB reforça a relevância da criação de um repositório específico e especializado que enriqueça futuras investigações na área de recuperação de informação musical.

Referencias

Allik, A., Fazekas, G., & Sandler, M. (2016). An ontology for audio features. Proceedings of the 17th International Society for Music Information Retrieval Conference, ISMIR 2016.

Cassapo, F. M. (2004). Uma Sociedade Multiagente para o Mapeamento Automático Inteligente de Competências em Ambiente de Colaboração. Pontifícia Universidade Católica do Paraná.

de Berardinis, J., Carriero, V. A., Meroño-Peñuela, A., Poltronieri, A., & Presutti, V. (2023). The Music Meta Ontology: A Flexible Semantic Model for the Interoperability of Music

Metadata. Proceedings of the 24th International Society for Music Information Retrieval Conference, 859–867. <https://doi.org/10.5281/zenodo.10265423>

Fazekas, G., & Sandler, M. B. (2011). The studio ontology framework. Proceedings of the 12th International Society for Music Information Retrieval Conference, ISMIR 2011.

Garnelo, M., & Shanahan, M. (2019). Reconciling deep learning with symbolic artificial intelligence: representing objects and relations. In *Current Opinion in Behavioral Sciences* (Vol. 29). <https://doi.org/10.1016/j.cobeha.2018.12.010>

Gruber, T. R. (1993). A translation approach to portable ontology specifications. *Knowledge Acquisition*, 5(2). <https://doi.org/10.1006/knac.1993.1008>

Guarino, N., & Giarretta, P. (1995). Ontologies and Knowledge Bases: Towards a Terminological Clarification. Towards Very Large Knowledge Bases. *Knowledge Building and Knowledge Sharing*, 1(9).

Howe, J. (2006). The Rise of Crowdsourcing. *Wired Magazine*, 14(06). <https://doi.org/10.1086/599595>

Industry, I. F. of the P. (2024). Global Music Report 2024.

Industry, I. F. of the P. (2025). Global Music Report 2025. https://ifpi-website-cms.s3.eu-west-2.amazonaws.com/GMR_2025_State_of_the_Industry_Final_83665b84be.pdf

Kolozali, S., Barthet, M., Fazekas, G., & Sandler, M. (2011). Knowledge representation issues in musical instrument ontology design. Proceedings of the 12th International Society for Music Information Retrieval Conference, ISMIR 2011.

Maslow, A. H. (1943). A theory of human motivation. *Psychological Review*, 50(4). <https://doi.org/10.1037/h0054346>

McGuinness, N. F. N., and D. L. (2017). Ontology Development 101: A Guide to Creating Your First Ontology. *Sustainability* (Switzerland), 9(12).

McGuinness, F. Deborah; McGuinness, L. van H., D. L.; Van Harmelen, F. (2004). OWL Web Ontology Language Overview. W3C recommendation. *W3C Recommendation*, 10(10).

Proutskova, P., Volk, A., Heidarian, P., & Fazekas, G. (2020). From music ontology towards ethno-music-ontology. Proceedings of the 21st International Society for Music Information Retrieval Conference, ISMIR 2020.

Raimond, Y., Abdallah, S. A., Sandler, M. B., & Giasson, F. (2007). Music Ontology. In S. Dixon, D. Bainbridge, & R. Typke (Eds.), *Proceedings of the 8th International Conference on Music Information Retrieval, {ISMIR} 2007*, Vienna, Austria, September 23-27, 2007 (pp. 417–422). Austrian Computer Society. http://ismir2007.ismir.net/proceedings/ISMIR2007%5C_p417%5C_raimond.pdf

Rashid, S. M., De Roure, D., & McGuinness, D. L. (2018). A music theory ontology. *ACM International Conference Proceeding Series*. <https://doi.org/10.1145/3243907.3243913>

Schreiber, H. (2016). Genre ontology learning: Comparing curated with crowd-sourced ontologies. Proceedings of the 17th International Society for Music Information Retrieval Conference, ISMIR 2016.

Sébastien, V., Sébastien, D., & Conruyt, N. (2013). Annotating works for music education: Propositions for musical forms and structures ontology and a musical performance ontology. Proceedings of the 14th International Society for Music Information Retrieval Conference, ISMIR 2013.

W3C OWL Working Group. (2012). OWL 2 Web Ontology Language Document Overview. OWL 2 Web Ontology Language, December.

Wang, J., Chen, X., Hu, Y., & Feng, T. (2010). Predicting high-level music semantics using social tags via ontology-based reasoning. Proceedings of the 11th International Society for Music Information Retrieval Conference, ISMIR 2010.

Wilmering, T., Fazekas, G., & Sandler, M. B. (2013). The audio effects ontology. Proceedings of the 14th International Society for Music Information Retrieval Conference, ISMIR 2013.

Anexo I

Descreve-se, abaixo, a especificação da ontologia das letras de músicas brasileiras de domínio público (MpbDomPubl). É importante salientar que, nesta transcrição em OWL, há apenas uma instância de música entre um total de 228. A ontologia completa pode ser acessada pelo link do repositório Zenodo MPBDomPubl - Ontologia da Música Popular Brasileira de Domínio Público

```
<?xml version="1.0"?>
<Ontology xmlns="http://www.w3.org/2002/07/owl#"
  xmlns:base="http://www.semanticweb.org/ontologies/mpb-dp"
  xmlns:rdf="http://www.w3.org/1999/02/22-rdf-syntax-ns#"
  xmlns:xml="http://www.w3.org/XML/1998/namespace"
  xmlns:xsd="http://www.w3.org/2001/XMLSchema#"
  xmlns:rdfs="http://www.w3.org/2000/01/rdf-schema#"
  ontologyIRI="http://www.semanticweb.org/ontologies/mpb-dp">
  <Prefix name="" IRI="http://www.semanticweb.org/ontologies/mpb-dp#"/>
  <Prefix name="mpb" IRI="http://www.semanticweb.org/ontologies/mpb-dp#"/>
  <Prefix name="owl" IRI="http://www.w3.org/2002/07/owl#"/>
  <Prefix name="rdf" IRI="http://www.w3.org/1999/02/22-rdf-syntax-ns#"/>
  <Prefix name="xml" IRI="http://www.w3.org/XML/1998/namespace"/>
  <Prefix name="xsd" IRI="http://www.w3.org/2001/XMLSchema#"/>
  <Prefix name="rdfs" IRI="http://www.w3.org/2000/01/rdf-schema#"/>
  <Declaration>
    <Class IRI="#Musica"/>
  </Declaration>
  <Declaration>
    <Class IRI="#Pessoa"/>
  </Declaration>
  <Declaration>
    <Class IRI="#genero"/>
  </Declaration>
  <Declaration>
    <ObjectProperty IRI="#temArtista"/>
  </Declaration>
  <Declaration>
    <ObjectProperty IRI="#temGenero"/>
  </Declaration>
  <Declaration>
    <DataProperty IRI="#anoLancamento"/>
  </Declaration>
  <Declaration>
    <DataProperty IRI="#artista"/>
  </Declaration>
  <Declaration>
    <DataProperty IRI="#id"/>
  </Declaration>
```

```

<Declaration>
  <DataProperty IRI="#letra"/>
</Declaration>
<Declaration>
  <DataProperty IRI="#link"/>
</Declaration>
<Declaration>
  <DataProperty IRI="#nomeArtista"/>
</Declaration>
<Declaration>
  <DataProperty IRI="#termo"/>
</Declaration>
<Declaration>
  <DataProperty IRI="#titulo"/>
</Declaration>
<Declaration>
  <NamedIndividual IRI="#abre_alas_6"/>
</Declaration>
<ClassAssertion>
  <Class IRI="#Musica"/>
  <NamedIndividual IRI="#abre_alas_6"/>
</ClassAssertion>
<ObjectProperty IRI="#temArtista"/>
  <NamedIndividual IRI="#abre_alas_6"/>
  <NamedIndividual IRI="#chiquinha_gonzaga"/>
</ObjectPropertyAssertion>
<ObjectPropertyAssertion>
  <ObjectProperty IRI="#temGenero"/>
  <NamedIndividual IRI="#abre_alas_6"/>
  <NamedIndividual IRI="#marcha_de_carnaval"/>
</ObjectPropertyAssertion>
<DataPropertyAssertion>
  <DataProperty IRI="#anoLancamento"/>
  <NamedIndividual IRI="#abre_alas_6"/>
  <Literal>
datatypeIRI="http://www.w3.org/2001/XMLSchema#integer">1900</Literal>
  </DataPropertyAssertion>
<DataPropertyAssertion>
  <DataProperty IRI="#artista"/>
  <NamedIndividual IRI="#abre_alas_6"/>
  <Literal>Chiquinha Gonzaga</Literal>
</DataPropertyAssertion>
<DataPropertyAssertion>
  <DataProperty IRI="#id"/>
  <NamedIndividual IRI="#abre_alas_6"/>
  <Literal>
datatypeIRI="http://www.w3.org/2001/XMLSchema#integer">6</Literal>
  </DataPropertyAssertion>
<DataPropertyAssertion>
  <DataProperty IRI="#letra"/>
  <NamedIndividual IRI="#abre_alas_6"/>
  <Literal>abre alas eu quero passar peço licença poder desabafar
jardineira abandonou jardim rosa resolveu gostar jardineira abandonou jardim
rosa resolveu gostar abre alas eu quero passar peço licença poder desabafar
jardineira abandonou jardim rosa resolveu gostar jardineira abandonou jardim
rosa resolveu gostar</Literal>
  </DataPropertyAssertion>
<DataPropertyAssertion>
  <DataProperty IRI="#link"/>
  <NamedIndividual IRI="#abre_alas_6"/>
  <Literal>https://www.vagalume.com.br/chiquinha-gonzaga/</Literal>

```

```

</DataPropertyAssertion>
<DataPropertyAssertion>
  <DataProperty IRI="#termo"/>
  <NamedIndividual IRI="#abre_alas_6"/>
  <Literal>abandonou</Literal>
</DataPropertyAssertion>
<DataPropertyAssertion>
  <DataProperty IRI="#termo"/>
  <NamedIndividual IRI="#abre_alas_6"/>
  <Literal>abre</Literal>
</DataPropertyAssertion>
<DataPropertyAssertion>
  <DataProperty IRI="#termo"/>
  <NamedIndividual IRI="#abre_alas_6"/>
  <Literal>alas</Literal>
</DataPropertyAssertion>
<DataPropertyAssertion>
  <DataProperty IRI="#termo"/>
  <NamedIndividual IRI="#abre_alas_6"/>
  <Literal>desabafar</Literal>
</DataPropertyAssertion>
<DataPropertyAssertion>
  <DataProperty IRI="#termo"/>
  <NamedIndividual IRI="#abre_alas_6"/>
  <Literal>eu</Literal>
</DataPropertyAssertion>
<DataPropertyAssertion>
  <DataProperty IRI="#termo"/>
  <NamedIndividual IRI="#abre_alas_6"/>
  <Literal>gostar</Literal>
</DataPropertyAssertion>
<DataPropertyAssertion>
  <DataProperty IRI="#termo"/>
  <NamedIndividual IRI="#abre_alas_6"/>
  <Literal>jardim</Literal>
</DataPropertyAssertion>
<DataPropertyAssertion>
  <DataProperty IRI="#termo"/>
  <NamedIndividual IRI="#abre_alas_6"/>
  <Literal>jardineira</Literal>
</DataPropertyAssertion>
<DataPropertyAssertion>
  <DataProperty IRI="#termo"/>
  <NamedIndividual IRI="#abre_alas_6"/>
  <Literal>licença</Literal>
</DataPropertyAssertion>
<DataPropertyAssertion>
  <DataProperty IRI="#termo"/>
  <NamedIndividual IRI="#abre_alas_6"/>
  <Literal>passar</Literal>
</DataPropertyAssertion>
<DataPropertyAssertion>
  <DataProperty IRI="#termo"/>
  <NamedIndividual IRI="#abre_alas_6"/>
  <Literal>peço</Literal>
</DataPropertyAssertion>
<DataPropertyAssertion>
  <DataProperty IRI="#termo"/>
  <NamedIndividual IRI="#abre_alas_6"/>
  <Literal>poder</Literal>
</DataPropertyAssertion>

```

```

<DataPropertyAssertion>
  <DataProperty IRI="#termo"/>
  <NamedIndividual IRI="#abre_alas_6"/>
  <Literal>quero</Literal>
</DataPropertyAssertion>
<DataPropertyAssertion>
  <DataProperty IRI="#termo"/>
  <NamedIndividual IRI="#abre_alas_6"/>
  <Literal>resolveu</Literal>
</DataPropertyAssertion>
<DataPropertyAssertion>
  <DataProperty IRI="#termo"/>
  <NamedIndividual IRI="#abre_alas_6"/>
  <Literal>rosa</Literal>
</DataPropertyAssertion>
<DataPropertyAssertion>
  <DataProperty IRI="#titulo"/>
  <NamedIndividual IRI="#abre_alas_6"/>
  <Literal>Abre Alas</Literal>
</DataPropertyAssertion>
<ObjectPropertyDomain>
  <ObjectProperty IRI="#temArtista"/>
  <Class IRI="#Musica"/>
</ObjectPropertyDomain>
<ObjectPropertyDomain>
  <ObjectProperty IRI="#temGenero"/>
  <Class IRI="#Musica"/>
</ObjectPropertyDomain>
<ObjectPropertyRange>
  <ObjectProperty IRI="#temArtista"/>
  <Class IRI="#Pessoa"/>
</ObjectPropertyRange>
<ObjectPropertyRange>
  <ObjectProperty IRI="#temGenero"/>
  <Class IRI="#genero"/>
</ObjectPropertyRange>
<DataPropertyDomain>
  <DataProperty IRI="#anoLancamento"/>
  <Class IRI="#Musica"/>
</DataPropertyDomain>
<DataPropertyDomain>
  <DataProperty IRI="#artista"/>
  <Class IRI="#Musica"/>
</DataPropertyDomain>
<DataPropertyDomain>
  <DataProperty IRI="#id"/>
  <Class IRI="#Musica"/>
</DataPropertyDomain>
<DataPropertyDomain>
  <DataProperty IRI="#letra"/>
  <Class IRI="#Musica"/>
</DataPropertyDomain>
<DataPropertyDomain>
  <DataProperty IRI="#link"/>
  <Class IRI="#Musica"/>
</DataPropertyDomain>
<DataPropertyDomain>
  <DataProperty IRI="#nomeArtista"/>
  <Class IRI="#Pessoa"/>
</DataPropertyDomain>
<DataPropertyDomain>

```



```

        <DataProperty IRI="#termo"/>
        <Class IRI="#Musica"/>
    </DataPropertyDomain>
    <DataPropertyDomain>
        <DataProperty IRI="#titulo"/>
        <Class IRI="#Musica"/>
    </DataPropertyDomain>
    <DataPropertyRange>
        <DataProperty IRI="#anoLancamento"/>
        <Datatype abbreviatedIRI="xsd:integer"/>
    </DataPropertyRange>
    <DataPropertyRange>
        <DataProperty IRI="#artista"/>
        <Datatype abbreviatedIRI="xsd:string"/>
    </DataPropertyRange>
    <DataPropertyRange>
        <DataProperty IRI="#id"/>
        <Datatype abbreviatedIRI="xsd:integer"/>
    </DataPropertyRange>
    <DataPropertyRange>
        <DataProperty IRI="#letra"/>
        <Datatype abbreviatedIRI="xsd:string"/>
    </DataPropertyRange>
    <DataPropertyRange>
        <DataProperty IRI="#link"/>
        <Datatype abbreviatedIRI="xsd:string"/>
    </DataPropertyRange>
    <DataPropertyRange>
        <DataProperty IRI="#nomeArtista"/>
        <Datatype abbreviatedIRI="xsd:string"/>
    </DataPropertyRange>
    <DataPropertyRange>
        <DataProperty IRI="#termo"/>
        <Datatype abbreviatedIRI="xsd:string"/>
    </DataPropertyRange>
    <DataPropertyRange>
        <DataProperty IRI="#titulo"/>
        <Datatype abbreviatedIRI="xsd:string"/>
    </DataPropertyRange>
    <AnnotationAssertion>
        <AnnotationProperty abbreviatedIRI="rdfs:label"/>
        <IRI>#abre_alas_6</IRI>
        <Literal>Abre Alas</Literal>
    </AnnotationAssertion>
</Ontology>

```

Apêndice B – Rotulação manual de músicas do conjunto de dados da MSD

Este experimento foi realizado como parte do projeto Análise e Recuperação de Informação Musical em Larga escala PIBIC 2014/2015 na PUC-PR pela estudante Ana Pavan no Programa de Pós-Graduação de Informática Aplicada.

Dentre os gêneros e subgêneros do conjunto de músicas da MSD foram selecionados 25 com a escolha aleatória de 10 músicas por grupo de gêneros, perfazendo um total de 250 músicas, ver QUADRO B.1 abaixo.

#	Gênero	Quant	Instrumentos Predominantes	Emoções Predominantes
1	BIGBAND	10	Bateria/Contrabaixo/Saxofone/Trompete	<u>Alegria/Paixão/Melancolia</u>
2	BLUES CONTEMPORARY	10	Guitarra/Bateria/Baixo/Vocal	<u>Melancolia/Alegria/Paixão</u>
3	COUNTRY Tradicional	10	Violão/Guitarra/Baixo/Vocal	<u>Alegria/Melancolia/Saudade</u>
4	DANCE	10	Baixo/Teclado/Vocal/Flauta	<u>Alegria/Melancolia/Saudade</u>
5	ELECTRONICA	10	Bateria/Baixo/Vocal/Flauta	<u>Melancolia/Alegria/Saudade</u>
6	FOLK INTERNATIONAL	10	Violão/Guitarra/Baixo/Vocal	<u>Alegria/Calma/Melancolia</u>
7	GOSPEL	10	Bateria/Baixo/Piano/Vocal	<u>Alegria/Calma/Saudade</u>
8	EXPERIMENTAL	10	Bateria/Baixo/Piano/Vocal	<u>Calma/Melancolia/Saudade</u>
9	GRUNGE EMO	10	Guitarra/Bateria/Baixo/Vocal	<u>Alegria/Melancolia/Raiva</u>
10	HIP HOP / RAP	10	Bateria/Teclado/Vocal/Flauta	<u>Melancolia/Paixão/Raiva</u>
11	JAZZ CLASSIC	10	Guitarra/Bateria/Baixo/Piano	<u>Alegria/Calma/Melancolia</u>
12	METAL ALTERNATIVE	10	Guitarra/Bateria/Baixo/Vocal	<u>Melancolia/Raiva/Saudade</u>
13	METAL DEATH	10	Guitarra/Bateria/Baixo/Vocal	<u>Melancolia/Raiva/Raiva</u>
14	METAL HEAVY	10	Guitarra/Bateria/Baixo/Vocal	<u>Alegria/Melancolia/Raiva</u>
15	POP CONTEMPORARY	10	Guitarra/Bateria/Vocal/Flauta	<u>Alegria/Paixão/Saudade</u>
16	INDIE POP	10	Guitarra/Bateria/Baixo/Vocal	<u>Melancolia/Alegria/Paixão</u>
17	POP LATIN	10	Violão/Bateria/Baixo/Vocal	<u>Alegria/Paixão/Calma</u>
18	PUNK	10	Guitarra/Bateria/Baixo/Vocal	<u>Melancolia/Raiva/Alegria</u>
19	REGGAE	10	Guitarra/Bateria/Baixo/Vocal	<u>Alegria/Melancolia/Paixão</u>
20	RNB SOUL	10	Guitarra/Bateria/Baixo/Vocal	<u>Alegria/Paixão/Melancolia</u>

21	ROCK ALTERNATIVE	10	Guitarra/Bateria/Baixo/Vocal	<u>Melancolia/Alegria/Saudade</u>
22	ROCK COLLEGE	10	Guitarra/Bateria/Baixo/Vocal	<u>Melancolia/Paixão/Saudade</u>
23	ROCK CONTEMPORARY	10	Violão/Bateria/Baixo/Vocal	<u>Melancolia/Paixão/Saudade</u>
24	ROCK HARD	10	Guitarra/Bateria/Baixo/Vocal	<u>Melancolia/Alegria/Raiva</u>
25	ROCK NEO PSYCHEDELIA	10	Guitarra/Bateria/Baixo/Vocal	<u>Alegria/Melancolia/Raiva</u>
	total	250		

QUADRO B.1 – Rotulação manual por gênero musical

Para o experimento de rotulação manual das músicas e composição do repositório de dados de músicas *Ground Truth* foram utilizados os seguintes termos classificados em dois grupos: instrumentos, ver QUADRO B.2, e emoções, ver QUADRO B.3, com a finalidade de padronizar o processo de rotulação.

Código	Instrumento
1	Pandeiro
2	Violão
3	Guitarra
4	Bateria
5	Baixo
6	Contrabaixo
7	Teclado
8	Saxofone
9	Trompete
10	Piano
11	Órgão
12	Vocal
13	Violino
14	Flauta
15	Tuba
16	Sintetizador
17	Instrumentos não comuns
18	Acordeão
19	Maracas

Código	Emoção
A	Alegria
M	Melancolia
NI	Não Identificado
S	Saudade
C	Calma
S	Saudade
P	Paixão
R	Raiva

QUADRO B.2 – códigos das emoções utilizados para rotulação

20	Castanhola
21	Bongo
22	Harpa
23	Gaita
QUADRO B.1 – códigos dos Instrumentos utilizados para rotulação	

Género	1	2	3	4	5	6	7	8	9	10
BIGBAND										
TRACK_ID	TRIXZYJ128F426F6EC'	TRIXZYJ128F426F7EC	TRIXZYJ128F426F8EC'	TRIXZYJ128F426F9EC'	TRIXZYJ128F426F10EC'	TRIXZYJ128F426F11EC'	TRIXZYJ128F426F12EC'	TRIXZYJ128F426F13EC'	TRIXZYJ128F426F14EC'	TRIXZYJ128F426F15EC'
TITLE	BOOGIE WOOGIE	GOODY GOODY	GOING TO TOWN	ROCK IN RHYTHM	MOOD INDIGO	ONE O'CLOCK JUMP	SUMMERTIME	APPOLO DAZE	SWEE' PEA	SOMEBODY LOVES ME
ARTIST	COUNT ORCHESTRA	BASIE BENNY GOODMAN	DUKE ELLINGTON	DUKE ELLINGTON	DUKE ELLINGTON	DUKE ELLINGTON	COUNT BASIE ORCHESTRA	COUNT BASIE ORCHESTRA	COUNT BASIE ORCHESTRA	COUNT ORCHESTRA
ALBUM	SINGLE	GOODY GOODY	RELEASE THE STARS	ROCK'IN IN RYTHM	MOOD INDIGO	ONE O'CLOCK JUMP	SUMMERTIME	SINGLE	SWEE' PEA	SINGLE
YEAR	1870	1936	2007	1931	1930	1937	1935	1930	1966	1924
COUNTRY	US	US	US	US	US	US	US	US	US	US
COMPOSER	COUNT BASIE	MATTY MALNECK	RUFUS WAINRIGHT	DUKE ELLINGTON	DUKE ELLINGTON	DUKE ELLINGTON	DUKE ELLINGTON	COUNT BASIE	TOMMY ROE	GEORGE GERSHWIN
INSTRUMENTS	4,6,8,10,12	4, 6, 8,9,12, 14	4,6,8,9	4,6,8,9,10,14	4,6,8,9,10,14	4,6,8,9,10,14,15	2,4,6,8,9,10,14	4,6,8,9,10,14	4, 6, 8, 9, 10, 14, 15	4,6,9,10
EMOTIONS	A, P	A, P	A	A	M	A,P	A,P	A	A,P	A,M
BLUES CONTEMPORARY										
TRACK_ID	TRIXZYJ128F426F16EC'	TRIXZYJ128F426F17EC'	TRIXZYJ128F426F18EC'	TRIXZYJ128F426F19EC'	TRUMJFH12903D056BD'	TRIUFDC128E079307B'	TRMLTWT128F92CBF07'	TRCPHZP128F426CC5A'	TRRWQPH128F9308EA8'	TRTTPKG12903CA2F42'
TITLE	MY LAST REGRET	MATCHBOX	I DON'T WANT TO SAY BEST WISHES	CANNED HEAT	SAVE ME	STILL RAININ'	BACK DOOR SLAM	EAT THE LUNCH YOU BROUGHT	TWO STEPS FROM THE END	MISSISSIPPI BLUES
ARTIST	ROBERT CRAY	JONNY LANG	DUKE ROBILLARD	RORY BLOCK	JANIVA MAGNESS	JONNY LANG	ROBERT CRAY	JANIVA MAGNESS	ROBERT CRAY	RORY BLOCK
ALBUM	TWENTY	SUN	BLUES FOR A SUNDAY MORNING	SYNCRONIZED	US AND THEM	WANDER THIS WORLD	TIME WILL TELL	BURY HIM AT CROSSROADS	TE TWENTY	SINGLE
YEAR	2007	1956	2005	1999	2005	1998	2003	2008	1970	1970
COUNTRY	US	US	US	ENGLAND	US	US	US	US	US	US
COMPOSER	ROBERT CRAY	CARL PERKINS	DUKE ROBILLARD	JAMIROQUAI	SHINEDOWN	JONNY LANG	ROBERT CRAY	JANIVA MAGNESS	ROBERT CRAY	RORY BLOCK
INSTRUMENTS	3,4,6,10,12	3,4,5,9,12	4,5,8,10,12	2,4,12	1,2,3,4,5,7,12	3,4,5,12	3,4,5,12	4,5,11,12	3,4,5,11,12	2,12
EMOTIONS	M,P	A,P	M,P	A,M	M	A,M	A,M	M,P	M	A,M
COUNTRY Tradicional										
TRACK_ID	TRCINKW12903CB1E90'	TRAQBWC128F426C0E4'	TRCWGV128F42635A0'	TRGKNFZ128F92FC080	TRGUFBV128F934366F'	TRUIMZU128E078985E'	TRLROHY12903CAFB30'	TRLBESE128F42591EE'	TRIXZYJ128F426F14EC'	TRWGUKO128F4263CG97'
TITLE	THE GIRL IN THE WOOD	EUGENIE OREGON	DESPERATELY	EASY COME EASY GO	PEACE IN THE VALLEY	TAKE ME BACK TO THE COUNTRY	SEAMANS BLUES	I CAN'T HELP IT	IF YOU ONLY KNEW	WHY SHOULD WE TRY ANYMORE
ARTIST	JIMMIE RODGERS	DOLLY PARTON	GEORGE STRAIT	GEORGE STRAIT	LORETTA LYNN	DOLLY PARTON	WILLIE NELSON	HANK WILLIAMS	TONNY RICE	HANK WILLIAMS
ALBUM	SINGLE	SINGLE	WRAPPED	EASY COME EASY GO	SINGLE	BEST OF COUNTRY	E	NA.NA..	COLD ON THE SHOULDER	NA.NA..
YEAR	1958	1972	1988	1993	1937	1998	2010	1972	1984	1972
COUNTRY	US	US	US	US	US	US	US	US	US	US
COMPOSER	JIMMIE RODGERS	DOLLY PARTON	BRUCE ROBINSON	AARON BARKER	THOMAS A. DORSEY	KAREN STALY	ERNEST TUBB	HANK WILLIAMS	LARRY RICE	HANK WILLIAMS
INSTRUMENTS	1,2,5,12	2,4,5,12	2,3,4,5,12	2,3,4,5,12	2,3,4,5,12	2,3,4,5,12	3,12	2,3,12	3,12	1,2,3,4,5,12
EMOTIONS	A,M,P	A,R,S	A,M,S	M,S	S	A	A	P,S	A,C	M,S

DANCE											
TRACK_ID	'TRMOWJ128E07818C8'	'TROCFHA128E07818C7'	'TRAMQNY128F14A03DC'	'TRBKRVVA128F92E6247'	'TRGFNSL128F92E6FAA'	'TRCOVJP12903CFCDF2'	'TRFYWXA12903CFCDF5'	'TRBGOGG128F931BD69'	TRJBQJ128F931BD66'	'TRFEKNV128E07818C0'	
TITLE	TOO LONG	AERODYNAMIC	BREAK MY HEART	CLEAR BLUE	SORROW HOME	PUESTA DEL SOL	KEEP ON RUNNING	HALO	FREEFALL	HARDER BETTER FASTER STRONGER	
ARTIST	DAFT PUNK	DAFT PUNK	ATB	MARKUS SCHULZ	MARKUS SCHULZ	JUDGE JULES	JUDGE JULES	CHRISTOPHER LAWRENCE	CHRISTOPHER LAWRENCE	DAFT PUNK	
ALBUM	DAFT CLUB	DAFT CLUB	ADDICTED TO MUSIC	ELECTRONIC ELEMENTS	WITHOUT YOU NEAR	PROVEN WORLDWIDE	PROVEN WORLDWIDE	ALL OR NOTHING	PHARMACY MUSIC	SINGLE	
YEAR	2003	2003	2003	2004	2005	2006	2006	2004	2004	2001	
COUNTRY	FRANCE	FRANCE	US	NETHERLANDS	NETHERLANDS	AUSTRALIA	AUSTRALIA	AUSTRALIA	AUSTRALIA	FRANCE	
COMPOSER	DAFT PUNK	DAFT PUNK	ATB	MARKUS SCHULZ	MARKUS SCHULZ	JUDGE JULES	JUDGE JULES	CHRISTOPHER LAWRENCE	CHRISTOPHER LAWRENCE	DAFT PUNK	
INSTRUMENTS	3,4,5,7,10,12,16	3,4,5,7,10,12,16	1,2,5,12,16	5,7,16	7,12,16	7,12,16	5,7,8,13,16	5,7,12,16	7,12,16	3,7,12,16	
EMOTIONS	M	A,M	M,S	M,S	M,S	A	P,S	A,M	A,M	A	
ELECTRONICA											
TRACK_ID	'TRBWWKE128E0784AD0'	'TRCVGOP128F4268656'	'TRBFTFC128F92FFAA4'	'TRMEELL128F427B0EE'	'TRRZZAE128E0798C75'	TRBXWUW128E0798C3F'	'TRFNHDM128F92C4F7D'	'TRBVVXD128F92C4F77'	'TRFOJGA128F14A0BE5'	'TRMWTYD128F14A0BDE'	
TITLE	FELT MOUNTAIN	CLOWNS	THE ONLY THING I KNOW	RIDE A WHITE HORSE	KAISER SALSEK	CRYSTAL CLEAR	TWO HEARTS IN ONE BODY	LEFTOVER	MOVE	DESTINATION EARTH	
ARTIST	GOLDFRAPP	GOLDFRAPP	GOTYE	GOLDFRAPP	GOLDIE	GOLDIE	MARASMA	MARASMA	NAOKI KENJI	NAOKI KENJI	
ALBUM	FELT MOUNTAIN	SEVENTH TREE	BOARDFACE	RIDE A WHITE HORSE	RING OF SATURN	RING OF SATURN	SINGLES	SINGLES	ECOUSTIC	ECOUSTIC	
YEAR	2000	2004	2003	2005	1998	1998	2000	2000	2004	2004	
COUNTRY	ENGLAND	ENGLAND	AUSTRALIA	ENGLAND	US	US	ITALY	ITALY	JAPAN	JAPAN	
COMPOSER	GOLDFRAPP	GOLDFRAPP	GOTYE	GOLDFRAPP	GOLDIE	GOLDIE	MARASMA	MARASMA	NAOKI KENJI	NAOKI KENJI	
INSTRUMENTS	4,5,10,11,12,16	2,5,12	3,4,5,12	5,7,12	3,4,5	3,4,5,7,8,12,16	5,12,16	5,10,12,16	4,6,16	1,5,7,16	
EMOTIONS	M	M,S	M,S	A,M	A	A,M,P	A,M	M,S	M,C	M,C	
FOLK INTERNATIONAL											
TRACK_ID	TRAMZUO128F421AE14'	TRBOGBX128F421AE09	TRFDAFZ128F4235573	TRAODJE128F149311E	TRCYGJP128F1493124	TRCPUMM128F4276D59	TREAYVD128F42801CF	TRAMUVB128F42819F3	TRBXQLG128F423846C	TRAGZPT128F93131D5	
TITLE	ROCK ME BABY	TULE TON SON TON	ROCKIN' CRADLE OF THE SOUTH	LISBOA AUSENTE	ROSA NEGRA	MAGIA DEL RITMO	AMIGO	I WISH YOU WELL	BLACK IS THE COLOUR	MI CAPONA	
ARTIST	CHUBBY CARRIER	CHUBBY CARRIER	CHUBBY CARRIER	ALA DOS NAMORADOS	ALA DOS NAMORADOS	GISPSY KINGS	GISPSY KINGS	CARA DILLON	CARA DILLON	PEPE PINTO	
ALBUM	DANCE ALL NIGHT	DANCE ALL NIGHT	WHO STOLE THE HOT SAUSE	COLCHON	ALA DOS NAMORADOS	SOMOS GRITANOS	ROOTS	AFTER THE MORNING	CARA DILLON	SUS 20 GRANDES EXITOS	
YEAR	1993	1993	1996	1994	2001	2001	2001	2006	2001	1930	
COUNTRY	US	US	US	PORTUGAL	PORTUGAL	FRANCE	FRANCE	UNITED KINGDOM	UNITED KINGDOM	SPAIN	
COMPOSER	CHUBBY CARRIER	CHUBBY CARRIER	CHUBBY CARRIER	ALA DOS NAMORADOS	ALA DOS NAMORADOS	GISPSY KINGS	GISPSY KINGS	CARA DILLON	CARA DILLON	PEPE PINTO	
INSTRUMENTS	3,4,5,12	3,4,5,8,18	3,4,5,8,18	2,5,12	2,5,12	2,3,4,5,12,20,21	2,3,4,5,12,20,21	2,12	6,9,10,12,13,14,22	2,12	
EMOTIONS	A	A	A	M, S,C	M,C	A,C	A,C	M,S	C	A,P,S	

GOSPEL											
TRACK_ID	TRDSFUR128F92F5B58	TRABOQY128F92F8B54	'TRMYKMU12903CA7826'	TRAONNE128F93037BC	TRCOXEU12903CC2717	'TRCSPAA128F92FA6A9'	'TRAE1RF128F14952B0'	'TRACPEN128F4293F01'	TRMYPBS128F4259E8D'	'TRAFFEOY128F4259E97'	
TITLE	KNOW GOD	SUPERMAN	MIND VS HEART	NOTHING REALLY CHANGES	I FEEL LIKE DYIN'	ONE WAY	WE'RE ALMOST THERE	STAY WITH ME	ROCKED ON THE DEEP	OUT OF MY DARKNESS	
ARTIST	VICKIE WINANS	VICKIE WINANS	NNEKA	LARRY NORMAN	LARRY NORMAN	LARRY NORMAN	SEVEN PLACES	SEVEN PLACES	GREATER VISION	LEGACY FIVE	
ALBUM	BRINGING IT ALL TOGETHER	BRINGING IT ALL TOGETHER	AL	UPON THIS ROCK	SOMETHING NEW UNDER THE SUN	STREET LEVEL	HEAR US SAY JESUS	GLOWING	QUARTETS	MONUMENTS	
YEAR	1990	1990	2008	1969	1980	1972	2004	2007	2003	2004	
COUNTRY	US	US	US	US	US	CANADA	US	US	US	US	
COMPOSER	VICKIE WINANS	VICKIE WINANS	NNEKA	LARRY NORMAN	LARRY NORMAN	LARRY NORMAN	SEVEN PLACES	SEVEN PLACES	EUGENE WRIGHT	WAYNE HAUN	
INSTRUMENTS	4,5,12,22	4,5,7,22	3,4,12	5,7,9,10,12	3,4,5,8,11	4,5,10,12	2,3,4,5,12	2,4,12	10,12	2,4,5,10,12	
EMOTIONS	A,C	A,S,C	A	A	A,M,S	A,C	A	A,C	A	A	
EXPERIMENTAL											
TRACK_ID	YNJM128E0781C42	'TRBKSZI128E07813CE'	'TRGFLCA128E0781C3C'	'TRBILEW128F92E7044'	'TRDJBTJ128F92E704D'	'TRKEWGSZ128F92E704D'	'TRAAGNS128E07824BD'	'TRBCDMC128F1452976'	TRMXDGA128F930E103	TRAASYQ128F930E0FE	
TITLE	KURUSHI	IT'S TIME FOR ACTION	ASK THE DRAGON	THE LAST GOOD TIME	BABY DOLL	OVER	COFFEE HOMEGROUND	A CORAL ROOM	MELANGE	DROWNING MAN	
ARTIST	YOKO ONO	YOKO ONO	YOKO ONO	OXBOW	OXBOW	OXBOW	KATE BUSH	KATE BUSH	GREENSLADE	GREENSLADE	
ALBUM	RISING	BLUEPRINT SUNRISE	FOR A	SERENAIDE IN RED	SERENAIDE IN RED	SERENAIDE IN RED	LIONHEART	AERIAL	SUNDANCE	GREENSLADE	
YEAR	1995	2001	1995	1997	1997	1997	1978	2005	1973	1973	
COUNTRY	US	US	US	US	US	US	ENGLAND	ENGLAND	ENGLAND	ENGLAND	
COMPOSER	YOKO ONO	YOKO ONO	YOKO ONO	MARIANNE FAITHFUL	MARIANNE FAITHFUL	MARIANNE FAITHFUL	KATE BUSH	KATE BUSH	TONY LAWSON	REEVES/DAVE LAWSON & GUERRA	
INSTRUMENTS	4,5,10,12	3,4,5,12,16	4,5,7,12	3,4,5,10,12,17	3,4,5,10,12,17	3,4,5,12	4,5,10,12,18	10,12	4,5,7,10,12	3,4,5,7,11,12	
EMOTIONS	M,S	NI	NI	M,R	M	M,C	A,C	M,S,C	M,C	A,M,C	
GRUNGE EMO											
TRACK_ID	TRPMDUU128F9320E24	TROBTVB128F9312AAA	TRCOOYB128E078ED95	TRDGXA128F42782F8	'TRHHKFG128F428A70B	TRRDFVG128F42893F7	TRJDTBE128F4289421	'TRYLLLG128F425E1D6	TRAWBOE128F92F2F46	TRADPIA128E078EE1B'	
TITLE	STAIN	SMELLS LIKE SPIRIT	COME AS YOU ARE	EVACUATION	DAUGHTER	DOLL	FOR ALL THE COWS	POX/AMERICAN EPJ	YOU ARE	HEART-SHAPED BOX	
ARTIST	NIRVANA	NIRVANA	NIRVAN A	PEARL JAM	PEARL JAM	FOO FIGHTERS	FOO FIGHTERS	SMASHING PUMPKINS	PEARL JAM	NIRVANA	
ALBUM	BLEW EP	NEVERMIND	LIVE TONIGHT OUT	BINAURAL	VS.	THE COLOUR AND THE SHAPE	FOO FIGHTERS	AMERICAN GOTHIC EP	RIOT ACT	IN UTERO	
YEAR	1989	1991	1994	2000	1993	1997	1992	2008	2002	1993	
COUNTRY	US	US	US	US	US	US	US	US	US	US	
COMPOSER	KURT KOBAIN	KURT KOBAIN	KURT KOBAIN	MIKE MCREADY	EDDIE VEDDER	DAVE GROHL	DAVE GROHL	SMASHING PUMPKINS	MATT CAMERON	KURT KOBAIN	
INSTRUMENTS	1,3,4,5,12	3,4,5,12	3,4,5,12	3,4,5,12	2,3,4,5,12	2,12	2,3,4,5,12	2,3,4,5,12	3,4,5,12	3,4,5,12	
EMOTIONS	M,R	A,R	C	R	A,M	M	A,M	M	M	M	

HIP HOP / RAP												
TRACK_ID	TRICWPS128F92C4A34	TRICWPS128F92C4A34	TRIVUXL128F9311BEC	TRATDKN128E078EDB3	TRBLASA12903CCAB0F	TRFKSWC128F9300987	TRTLKWA128F92FF219	TRTKTBZ128E078880B	TRUAGPE128F930C55F	TRWDGHI128E0785EA6		
TITLE	FOREVER YOUNG	WILL WORK FOR LOVE	I WONDER	THE WAY I AM	3 A.M.	MUST BE THE GANJA	GOD GAVE ME STYLE	I'M SUPPOSED TO DIE TONIGHT	99 SHIT	HARD TIMES		
ARTIST	USHER	USHER	KANYE WEST	KANYE WEST	EMINEM	EMINEM	EMINEM	50 CENT	50 CENT	KURTIS BLOW		
ALBUM	HERE I STAND	HERE I STAND	GRADUATION	THE GRATEST:CLASSIC JOINTZ, VOL 1-11	RELAPSE	RELAPSE	THE MASSACRE	THE MASSACRE	FREESTYLE B4 PAYSTYLE	KURTIS BLOW		
YEAR	2008	2008	2007	2000	2009	2009	2005	2005	2007	1980		
COUNTRY	US	US	US	US	US	US	US	US	US	US		
COMPOSER	MIKKEL STORLEER	MIKKEL STORLEER	THOMAS BANGHALTER	MARSHALL MATHERS	DON BLACK	DON BLACK	CURTIS JACKSON	MARSHALL MATHERS	CURTIS JACKSON	J.B. MOORDE		
INSTRUMENTS	4,7,10,12	3,4,7,10,12	7,12,16	5,7,12,16	4,10,12,16	6,12,16	12,16	5,12,16	4,5,7,12,16	3,4,5,12,16		
EMOTIONS	P	P	M	R	M	A	NI	M,S	M,R	A,M		
JAZZ CLASSIC												
TRACK_ID	TRAABHO12903D08576	TRAACPD128F931C3E7	TRBJOSO128F9342C69	TRCCMCP12903D08578	TRDYSKM128F9342C74	TRDLOJT128F9342C71	TRDMDDT128F145709A	TRBKGCC128F14582E7	TRBATZF128F146041B	TRFAHVA128F146710A		
TITLE	I KNOW YOU	MIDNIGHT SUN	SPIRIT	MIRAGE	SCUFFLE	MOOD SWINGS	OLD FOLKS	COOL BLUES	SONNYMOON FOR TWO	SOLID		
ARTIST	MIKE STERN	KITTY MARGOLIS	MIKE STERN	MIKE STERN	MIKE STERN	MIKE STERN	GRANT GREEN	GRANT GREEN	GRANT GREEN	GRANT GREEN		
ALBUM	THESE TIMES	EVOLUTION	VOICES	THESE TIMES	UPSIDE DOWNSIDE	UPSIDE DOWNSIDE	GRANSTAND	BORN TO BE BLUE	FIRST IMPRESSION	SOLID		
YEAR	2004	1993	2001	2004	1986	1986	1961	1962	2001	1964		
COUNTRY	US	US	US	US	US	US	US	US	US	US		
COMPOSER	MIKE STERN	JOHNNY MERCER	MIKE STERN	MIKE STERN	MIKE STERN	MIKE STERN	WILLARD ROBISON	CHARLIE PARKER	SONNY ROLLINS	SONNY ROLLINS		
INSTRUMENTS	4,5,7,12,22	3,4,5,6,12	3,4,5,10,12,16	3,4,5,12,16	3,4,5,12	2,3,4,5,10	2,3,4,5,10	2,3,4,5,10	4,5,10	4,5,8,9,10		
EMOTIONS	C	NDA	A,P	A,C	A	A,C	M,C	A,C	M,C	A,C		
METAL ALTERNATIVE												
TRACK_ID	TRBQVLO128F93260FD	TRDBEXY128F42A4ABC	TRXNHQT128F9326103	TRXASSM128F4241AF3	TRFHFCHU128EF3885ED	TRFRAEO128F4279C87	TRESBYB128E079207D	TRWUGQT128E0792647	TRFNEXE12903CB921E	TRAAYBP128E0792605		
TITLE	AGAIN	B.Y.O.B.	IN THE DARK	CASSIE	FLOWER	F**K THE SYSTEM	TAINTED LOVE	I PUT A SPELL ON YOU	BLESS THOSE	BOB		
ARTIST	FLYLEAF	SYSTEM OF A DOWN	FLYLEAF	FLYLEAF	SOUNDGARDEN	SYSTEM OF A DOWN	MARILYN MANSON	MARILYN MANSON	LIVING COLOUR	PRIMUS		
ALBUM	AGAIN	MESMERIZE	MEMENTO MORI	PAT SEALS	ULTRAMEGA OK	STEAL THIS ALBUM!	NOT ANOTHER TEEN MOVIE	SMELLS LIKE CHILDREN	ON A STAGE AT A WORLD CAFE LIVE	PORK SODA		
YEAR	2009	2005	2009	2004	1988	2005	2001	1993	2009	1993		
COUNTRY	US	US	US	US	US	US	US	US	US	US		
COMPOSER	FLYLEAF	DARON MALAKIAN	PAT SEALS	PAT SEALS	KIM THAYIL	DARON MALAKIAN	ED COBB	JAY HAWKINS	ANNIE BANDEZ	LARRY LALONDE		
INSTRUMENTS	3,4,5,12	3,4,5,12	3,4,5,12	3,4,5,12	3,4,5,12	3,4,5,12	4,5,7,12,16	2,3,4,5,12	3,4,5,12	3,4,5,6,12		
EMOTIONS	M,R,P	R	M,S	M,S	M,R	M,R	M,S	E	M,C	M		

METAL DEATH											
TRACK_ID	TRQVQLQ12903CA0AD5	TRGXQXZ128F92EEFD8	TRFEBZW128F92EEFE3	TRTMUCS128F93ZC8D0	TRTXDXO128F93ZC8D5	TRTWVTY128F934AE5C	TRZPELV128F934B4Z5	TRDAFY12903CF5BFB	TRDPQSM12903CDF2AE	TRDPMOZ128F14A1510	
TITLE	INTO THE BLACK SLUMBER	RISE OF THE LEVIATHAN	PROSTHETIC ERECTION	PAYBACK	SEE ME NOW	BLOODLANDS	SLAIN	INTRO	HEAVING EARTH	THE GAME	
ARTIST	ANNOTATIONS OF AN AUTOPSY	ANNOTATIONS OF AN AUTOPSY	ANNOTATIONS OF AN AUTOPSY	OBITUARY	OBITUARY	CANNIBAL CORPSE	CANNIBAL CORPSE	MORBID ANGEL	MORBID ANGEL	CREMATORY	
ALBUM	THE REING OF DARKNESS	BEFORE THE THRONE OF INFECTION	BEFORE THE THRONE OF INFECTION	DARKEST DAY	DARKEST DAY	CREATED TO KILL	THE WRETCHED SPAWN	BLESSED ARE THE SICK	FORMULAS FATAL TO THE FLESH	ACT SEVEN	
YEAR	2010	2008	2008	2009	2009	1995	2004	1991	1998	1999	
COUNTRY	ENGLAND	ENGLAND	ENGLAND	US	US	US	US	US	US	GERMANY	
COMPOSER	STEVE REGAN	ANNOTATIONS OF AN AUTOPSY	ANNOTATIONS OF AN AUTOPSY	TREVOR PERES	TREVOR PERES	ALEX WEBSTER	ALEX WEBSTER	MORBID ANGEL	MORBID ANGEL	CREMATORY	
INSTRUMENTS	3,4,5,12	3,4,5,12	3,4,5,12	3,4,5,12	3,4,5,12	3,4,5,12	3,4,5,12	3,4,5,12	3,4,5,12	3,4,5,12	
EMOTIONS	M,R	M,R	M,R	M,R	M,R	M,R	M,R	M,R	M,R	M,R	
METAL HEAVY											
TRACK_ID	TRDKLZ128F146838E	TRDWER128F42836ZF	TRFDEEV128F146F751	TRFJOJW128F92E4A46	TRJNEVF128F92E4A42	TRJVPZY128F92E4A48	TRJYGT12903CD3D5C	TRJKTUV128F931C3B2	TRJTACK128F4274196	TRKIVAO128F1467117	
TITLE	THE EDGE OF DARKNESS	PHANTOM OF THE OPERA	TWILIGHT ZONE	BACKS TO THE WALL	NEVER SURRENDER	BROKEN HEROES	IRON MAN	COMPUTER GOD	PRAY FOR BLOOD	PARANOID	
ARTIST	IRON MAIDEN	IRON MAIDEN	IRON MAIDEN	SAXON	SAXON	SAXON	BLACK SABBATH	BLACK SABBATH	MEGADETH	MEGADETH	
ALBUM	THE X FACTOR	LIVE AFTER DEATH	TWILIGHT ZONE	SAXON	DENIM AND LEATHER	THE SAXON CHRONICLES[VIDEO]	PARANOID	DEHUMANIZER	UNITED ABOMINATIONS	NAVITY TRIBUTE TO BLACK SABBATH	
YEAR	1995	1980	1981	1979	1981	2007	1970	1992	2007	1994	
COUNTRY	ENGLAND	ENGLAND	ENGLAND	ENGLAND	ENGLAND	ENGLAND	ENGLAND	ENGLAND	US	US	
COMPOSER	STEVE HARRIS	STEVE HARRIS	STEVE HARRIS	PAUL QUINN	SAXON	SAXON	OZZY OSBOURNE	GEEZER BUTLER	DAVE MUSTAINE	OZZY OSBOURNE	
INSTRUMENTS	3,4,5,12	3,4,5,12	3,4,5,12	3,4,5,12	3,4,5,12	3,4,5,12	3,4,5,12	3,4,5,12	3,4,5,12	3,4,5,12	
EMOTIONS	M	A	A,M	A,R	A,R	M	M,S	M,R	R	A,R	
POP CONTEMPORARY											
TRACK_ID	TRBBUHZ128F428F001	TRANRYD128E078F277	TRANRYD128E078F277	TRGTWMB128E078F2E5	TRJJXW1128F426BD83	TREKSNP128F92F33B3	TRWJUKVP128F92DCA7E	TRDBYVWP128F92FFC89	TRLHEVM128F9323439		
TITLE	4 MINUTES	MAMMA MIA	VOULEZ-VOUS	DANGING QUEEN	DOES YOUR MOTHER KNOW	THRILLER	BILLIE JEAN	WOMANIZER	IF YOU SEEK AMY	SOMEBODY LOVES ME	
ARTIST	MADONNA & JUSTIN TIMBERLAKE	ABBA	ABBA	ABBA	ABBA	MICHAEL JACKSON	MICHAEL JACKSON	BRITNEY SPEARS	BRITNEY SPEARS	BREAK THE ICE	
ALBUM	HARD CANDY	GREATEST HITS[VOGUE]	VOULEZ-VOUZ	ARRIVAL	ABBA LIVE[PROMO]	THRILLER	THRILLER	WOMANIZER	CIRCUS	BLACKOUT	
YEAR	2008	1975	1979	1976	1978	1982	1982	2008	2008	2007	
COUNTRY	US	SWEEDEN	SWEEDEN	SWEEDEN	SWEEDEN	US	US	US	US	US	
COMPOSER	JUSTIN TIMBERLAKE	BENNY ANDERSON	BENNY ANDERSON	BENNY ANDERSON	BENNY ANDERSON	MICHAEL JACKSON	MICHAEL JACKSON	NIKESHA BRISCOODE	MAX MARTIN	KERI HILSON	
INSTRUMENTS	9,12,16	3,4,5,7,12	3,4,5,7,12	3,4,5,7,12	3,4,5,7,12	3,4,5,12,16	3,4,5,12,16	12,16	12,16	12,16	
EMOTIONS	A,P,S	A,P	A,P	A,P	A,P	A	M,R	M,R,S	M,P	M,S	

INDIE POP										
TRACK_ID	TRAFUBQ128F92EB825	TRAZUOJ128F92EB81F	TRDQVAC128F92EB814	TREVJIS128F14A94E5	TRDQVOTB128F4261A05	TRPMGAI128F4261A03	TRABMBG128F426B9F1	TRABMF1128F1476F9F	TRBWDJH128F426E9FF	TREALF1128F147729E
TITLE	GIRLFRIEND IN A COMA	ASK	WHAT SHE SAID	SOME GIRLS ARE BIGGER THAN OTHERS	DARK VICENTE	EVERYBODY IS A STAR	SWEDEDEEDEE	WHEN THE WORLD IS RUNNING DOWN YOU MAKE THE BEST OF WAYS STILL AROUND	HE WAR	EVERY BREATH YOU TAKE
ARTIST	THE SMITHS	THE SMITHS	THE SMITHS	THE SMITHS	THE PASTELS	THE PASTELS	CAT POWER	THE POLICE	CAT POWER	THE POLICE
ALBUM	GIRLFRIEND IN A COMA	ASK	MEAT IS MURDER	THE QUEEN IS DEAD	THE LAST GREAT WILDERNESS	GEOGRAPHIC COMPILATION	THE COVERS RECORD	ZENYATTA MONDATT	YOU ARE FREE	SYNCHRONICITY
YEAR	1987	1986	1985	1986	2001	2001	2000	1980	2003	1983
COUNTRY	ENGLAND	ENGLAND	ENGLAND	ENGLAND	ENGLAND	ENGLAND	US	ENGLAND	US	ENGLAND
COMPOSER	MORRISSEY	MORRISSEY	MORRISSEY	MORRISSEY	KATRINA MITCHELL	SILVESTER STEWART	MICHAEL HURLEY	STING	CHAN MARSHALL	STING
INSTRUMENTS	2,3,4,5,12,22	2,3,4,5,7,12	3,4,5,12	3,4,5,12	3,4,5,12	3,4,5,14	2,5,12	3,4,5,12	2,3,4,5,12	3,4,5,6,12
EMOTIONS	M,S	A,P	M,P	M	M	M	M,C	A,M	A	A,M,P,C
POP LATIN										
TRACK_ID	TRYNUTX128F9309427	TRDNMMF128F9309430	TRHAVVT128F34217B	TRPHJVT128EF34217F	TRCWMTG128F4272F07	TRKALBA128F4272F08	TRICPPZ128F149591C	TRAVKCV128F149591D	TRKTLVF128F4279F4B	TRLRDOA128F149EC83
TITLE	LADY	INDIANAPOLIS	MEDICINE FOR MY SOUL	COMO YO	SIMPLESMENTE ESPIRITUAL	RETRATO	TWISTED	ONE NIGHT THING	PALOMA BLANCE(WHITE DOVE)	EL MURO
ARTIST	MENUDO	MENUDO	JUAN LUIS GUERRA	JUAN LUIS GUERRA	GIAN MARCO	GIAN MARCO	LUIS FONSI	LUIS FONSI	JULIO IGLESIAS	ENRIQUE IGLESIAS
ALBUM	POR AMOR	A TODO ROCK	LA LLAVE DE MI CORAZON	LA LLAVE DE MI CORAZON	A TIEMPO	A TIEMPO	FIGHT THE FEELING	FIGHT THE FEELING	HEY!	VIVIR
YEAR	1988	1988	2007	2007	2002	2002	2002	2002	1980	1987
COUNTRY	PUERTO RICO	PUERTO RICO	DOMINICAN REPUBLIC	DOMINICAN REPUBLIC	PERU	PERU	PUERTO RICO	PUERTO RICO	SPAIN	SPAIN
COMPOSER	JULIO SELIAS	JULIO SELIAS	JUAN LUIS GUERRA	JUAN LUIS GUERRA	GIAN MARCO	GIAN MARCO	LUIS FONSI	LUIS FONSI	JULIO IGLESIAS	ENRIQUE IGLESIAS
INSTRUMENTS	1,2,3,4,5,12,21	3,4,5,10,11,12,13	1,2,4,5,9,10,12,19	1,2,4,5,9,10,12,19	2,3,4,5,7,12	2,4,5,12	2,4,5,12	2,4,5,12	2,3,4,5,8,9,10,12,14,19	2,4,5,8,12
EMOTIONS	A,P	A	A,P	A,P,S	A,P	P,C	A,M,C	A,M,C	A,C	M,P,C
PUNK										
TRACK_ID	TRIFVND128F423EEDF	TRMKVAQ12903D0A693	TRWTOBI128F428E467	TRDCFFH128F428E46	TRCIFDQ128F428E45F	TRDPXRS128F92EAAFC	TRGYHUS128F92EAAF8	TRNSBLP128E0782864	TRMSBSM128E078287E	TRLKQCT128F932CE05
TITLE	FELL IN LOVE WITH A GIRL	SEVEN NATION ARMY	CATCH HELL BLUES	LITTLE CREAM SODA	BONE BROKE	RACE AGAINST MYSELF	NEOCON	MARS	CALL MR. LEE	FUCK FACISM
ARTIST	THE WHITE STRIPES	THE WHITE STRIPES	THE WHITE STRIPES	THE WHITE STRIPES	THE WHITE STRIPES	THE OFFSPRING	THE OFFSPRING	TELEVISION	TELEVISION	THE OPRESSED
ALBUM	WHITE BLOOD	UNDER GREAT WHITE NORTHERN LIGHTS	ICKY THUMP	ICKY THUMP	ICKY THUMP	RACE AGAINST MYSELF	NEOCON	TELEVISION	TELEVISION	THE OPRESSED
YEAR	2001	2010	2007	2007	2007	2003	2003	1992	1992	1996
COUNTRY	US	US	US	US	US	US	US	US	US	UNITED KINGDOM
COMPOSER	JACK WHITE/MEG WHITE	JACK WHITE/MEG WHITE	JACK WHITE	JACK WHITE	JACK WHITE	THE OFFSPRING	THE OFFSPRING	TELEVISION	TELEVISION	THE OPRESSED
INSTRUMENTS	3,4,12	3,4,12	3,4,12	3,4,12	3,4,12	3,4,5,12,16	3,4,5,17	3,4,5	3,4,5,12	3,4,5,12
EMOTIONS	A,P	M,R,C	M	A,M	M	A,M,R	M,R	M,R	M,R	R

REGGAE											
TRACK_ID	TRJWYJX128F930E050	TRCSQS128F930E025	TRPTBZF128F930E01E	TRAVMM12903CFAD62	TRAHRL1128F933D96F	TRALBHB12903CA2E52	TRBJWDK12903CE27D4	TRBKCNW128F9306BE6	TRFQYL128F42576F8	TRIEHD128F93365BD	
TITLE	KAYA	KEEP ON MOVING	DON'T ROCK MY BOAT	LIVELY UP YOURSELF	MINE YUH BUSINESS	FOLLY RANKING	MARIJUANA TREE	GIRLS I SEE	STREET LIFE	CHILD ABUSE	
ARTIST	BOB MARLEY	BOB MARLEY	BOB MARLEY	BOB MARLEY	JOHNNY OSBOURNE	JOHNNY OSBOURNE	JOHNNY OSBOURNE	BEENIE MAN	BEENIE MAN	LITTLE KIRK	
ALBUM	THE REGGAE	BEST OF BOB MARLEY	SOUL REBELS	BEST REGGAE	ROUGHER THAN MEN	GREEN LEAVES	NIGHTFALL SHOWCASE	HIGHLIGHTS	THE CRUSADERS	RAS REGGAE	
YEAR	1992	1989	1970	1994	1989	1980	1997	2000	1979	1991	
COUNTRY	JAMAICA	JAMAICA	JAMAICA	JAMAICA	UNITED KINGDOM	UNITED KINGDOM	UNITED KINGDOM	JAMAICA	JAMAICA	JAMAICA	
COMPOSER	GLEN ADAMS	BOB MARLEY	BOB MARLEY	BOB MARLEY	JOHNNY OSBOURNE	JOHNNY OSBOURNE	JOHNNY OSBOURNE	BEENIE MAN	BEENIE MAN	LITTLE KIRK	
INSTRUMENTS	1,3,4,5,8,12,16	1,2,4,5,12,21	1,3,5,12,19	3,4,5,8,9,12,19,21	4,5,7,12,16	3,4,5,12,21	3,4,5,12	4,5,12,16	3,4,5,12	4,5,10,12	
EMOTIONS	A,P	A	A,M	A	A	A,M	A,M	A,M,P	M	M	
RNB SOUL											
TRACK_ID	TRCSXLY128FR2EB4F4	TRPEEOA128F4295C28	TRZZUNB128F4263020	TRPNKJX128E078FF1B	TRVCDYT128F429810E	TRVCDYT128F429610E	TRRCQXU128E0799A11	TRFHVUH128F425F96E	TRFJWWU128E078D288	TRIEKHX128F429EE26	
TITLE	GOT TO BE THERE	I CAN'T HELP IT	IT'S THE FALLING IN LOVE	EVERYBODY'S SOMEBODY'S FOOL	FOREVER MORE	YOU MADE ME LOVE YOU	TELL YOUR HEART I LOVE YOU	KNOCKS ME OFF MY FEET	LET ME BE YOUR LOVEMAKER	SWEET WONDER	
ARTIST	MICHAEL JACKSON	MICHAEL JACKSON	MICHAEL JACKSON	CONNIE FRANCIS	R. KELLY	R. KELLY	STEVIE WONDER	STEVIE WONDER	BETTY WRIGHT	BETTY WRIGHT	
ALBUM	GOT TO BE THERE	OF THE WALL	OF THE WALL	THE VERY BEST OF CONNIE FRANCIS	CHOCOLATE FACTORY	CHOCOLATE FACTORY	THE COMPLETE WONDER	SONGS IN THE KEY OF LIFE	GOLDEN CLASSICS: CLEAN UP WOMAN	GOLDEN CLASSICS: CLEAN UP WOMAN	
YEAR	1972	1979	1979	1963	2003	2003	2005	1976	1988	1988	
COUNTRY	US	US	US	US	US	US	US	US	US	US	
COMPOSER	ELLIOT WILLENSKY	STEVIE WONDER	STEVIE WONDER	CONNIE FRANCIS	R. KELLY	R. KELLY	STEVIE WONDER	STEVIE WONDER	BETTY WRIGHT	BETTY WRIGHT	
INSTRUMENTS	3,4,5,10,12	3,4,5,10,12	3,5,4,12	3,5,6,12	2,3,4,5,12,16	3,4,5,12	3,4,5,12,16,23	4,5,7,12	3,4,5,12	1,3,4,5,12	
EMOTIONS	A,M,P	M,P	A,P	A,M,P	A,P	A,P	A,P	A,P	A,P	A,P	
ROCK ALTERNATIVE											
TRACK_ID	TRXHSSO128F425E6A4	TRAIBXQ128F425E6A5	TRYLCPG128F427CACA	TRXDWGM128F42474C	TRAYDXM128F9317358	TRNVWOK128F92DE8D0	TRINUNT128F9342A8E	TRSQEOW12903CF7C58	TRGOOGQ128F4284E69	TRRZTNO128F931132A	
TITLE	CHARMER	MY PARTY	HORSE TO WATER	SUMMER HIGH	NO YOU GIRLS	ULYSSES	LUCID DREAMS	STAY	I CAME AS A RAT	SPACEMAN	
ARTIST	KINGS OF LEON	KINGS OF LEON	REM	REM	FRANZ FERDINAND	FRANZ FERDINAND	FRANZ FERDINAND	FLY LEAF	MODEST MOUSE	THE KILLERS	
ALBUM	CHARMER	CHARMER	REVEAL	REVEAL	TONIGHT	ULYSSES	TONIGHT	MISSING	THE MOON & ANTARTICA	DAY & AGE	
YEAR	2007	2007	1980	1980	2009	2008	2009	2009	2000	2008	
COUNTRY	US	US	US	US	US	US	US	US	US	US	
COMPOSER	KINGS OF LEON	KINGS OF LEON	KINGS OF LEON	KINGS OF LEON	FRANZ FERDINAND	FRANZ FERDINAND	FRANZ FERDINAND	FLY LEAF	MODEST MOUSE	THE KILLERS	
INSTRUMENTS	3,4,5,12	3,4,5,12	3,4,5,12	2,3,4,5,7,12	3,4,5,12	3,4,5,12,16	3,4,5,12	2,12	2,3,4,5,12	3,4,5,12	
EMOTIONS	M,P	A,M	M,R	M	A,M,P	M	A,M,S	M,S	A,M,R	A,M,S	

ROCK COLLEGE										
TRACK_ID	TRCKY1M128E0782A2F	TRTYNQW128F4255339	TRENTGL128E0780C8E	TRLNWK128E0780CD8	TRNPPYW128F92CEAB0	TRQALM7128F147C48A	TRMTZ1E128C719698E	TRNVGNC128F9346A03	TRPNRDD128F425C31D	TRSUJZX128EF345A43
TITLE	THE ONE I LOVE	LOSING MY RELIGION	CLOCKS	YELLOW	ELECTION DAY	BUCK HILL	WITH OR WITHOUT YOU	REPULSION	MOUNTAIN MAN	LIGHTNING BULB
ARTIST	REM	REM	COLDPLAY	COLDPLAY	THE REPLACEMENTS	THE REPLACEMENTS	U2	DINOSAUR JR.	DINOSAUR JR.	DINOSAUR JR.
ALBUM	REMTV	REM LIVE	A PUSH OF BLOOD TO THE HEAD	LIVE	PLEASED TO MEET ME	HOOTENANNY	WITH OR WITHOUT YOU	DINOSAUR	DINOSAUR	BEYOND
YEAR	1987	2007	2002	2003	1987	1983	1987	1985	1985	2007
COUNTRY	US	US	US	US	US	US	US	US	US	US
COMPOSER	PETER BUCK	REM	CHRIS MARTIN	COLDPLAY	PAUL WESTERBER	ADAM CLAYTON	BONO	J MASCIS	J MASCIS	LOU BARLOW
INSTRUMENTS	3,4,5,12	3,4,5,12	3,4,5,12	3,4,5,12,14	3,4,5,12	1,3,4,5,10,12,13	3,4,5,12	3,4,5,12	3,4,5,12	3,4,5,12
EMOTIONS	M,P,S	M	A,P,C	M,P,S	M,R	M,P	A,M,S	A,M,S	R	M,R
TRACK_ID										
ROCK CONTEMPORARY										
TRACK_ID	TRADZQJ128F4271649	TRBFFPW128F425795F	TRBHZIX128E0786671	TROGAHB128F424E129	TRNSVGW128F424E127	TRKKTMM128F42B7815	TRVTFAM128E07928C4	TRAWLL128F92EDAD2	TRBHOVM128E0784F29	TRCLRSE128F92ECFEB
TITLE	INCOMPLETE	GIGGLING AGAIN FOR NO REASON	PRECIOUS ILLUSIONS	REAL WORLD	PUSH	MIAMI	CARRIAGE	EVERYTHING REMINDS ME OF YOU	SWEET PEA	SOMEBODY LOVES ME
ARTIST	ALANIS MORISSETTE	ALANIS MORISSETTE	ALANIS MORISSETTE	MATCHBOX TWENTY	MATCHBOX TWENTY	COUNTING CROWS	COUNTING CROWS	JEWEL	RUN 2 U	PERFECTLY CLEAR
ALBUM	FLAVORS OF ENTANGLEMENT	FLAVORS OF ENTANGLEMENT	PRECIOUS ILLUSIONS	REAL WORLD	PUSH	HARD CANDY	HARD CANDY	PERFECTLY CLEAR	SINGLE	PERFECTLY CLEAR
YEAR	2008	2008	2002	1998	2000	2002	2002	2003	2003	2008
COUNTRY	US	US	US	US	US	US	US	US	US	US
COMPOSER	ALANIS MORISSETTE	ALANIS MORISSETTE	ALANIS MORISSETTE	ROB THOMAS	ROB THOMAS	ADAM DURITZ	ADAM DURITZ	JOE FIRSTMAN	JEWEL	JEWEL
INSTRUMENTS	2,12	4,5,12,16	2,4,5,12	3,4,5,12	3,4,5,12	3,4,5,12	2,4,5,12	2,12	2,12,16	2,12
EMOTIONS	A,M,P,S	A,M,P,S	A,M,P,S	A,R	M,P	M,P	M,P	M,P,S	A,P	M,P,S,C
ROCK HARD										
TRACK_ID	TRJYVZH128F145BF8D	TRTAFYS128F14A51CF	TRKSBBWE128F92F1237	TRPPLTW128F145984E	TRJVVHWJ128F92D8B4F	TRULASX128F145BFD9	TRQFVEL128F930A83F	TRFNAUJ128F4298933	TRVUDCS128F931B0BB	TRKUUXY128F4270F37
TITLE	HEARTBREAKER	SMOKE ON THE WATER	CHILD IN TIME	HOLY MAN	GHOST SANDWICH	LOOKING FOR A STRANGER	FUCK YOU PAY ME	CUSTARD PIE	CHILDREN OF THE GRAVE	STAIRWAY TO HEAVEN
ARTIST	PAT BENATAR	DEEP PURPLE	DEEP PURPLE	DEEP PURPLE	ARE WEAPONS	PAT BENATAR	BRUCE HAMPTON	LED ZEPELLIN	BLACK SABBATH	LED ZEPELLIN
ALBUM	IN THE HEART OF THE NIGHT	LIVE IN LONDON	DEEP PURPLE IN ROCK	STORMBRINGER	ARE WEAPONS	GET NERVOUS	ONE RUINED LIFE OF A BRONZE TOURIST	PHYSICAL GRAFFITI	CHILDREN OF THE GRAVE	LED ZEPELLIN IV
YEAR	1979	1982	1970	1974	2000	1982	2003	1975	1976	1971
COUNTRY	US	US	US	US	US	US	US	ENGLAND	US	ENGLAND
COMPOSER	PAT BENATAR	RITCHIE BLACKMOND	RITCHIE BLACKMOND	DAVID COVERDALE	BRAIN F. MCPECK	FRANNE GOLDE	BRUCE HAMPTON	JIMMIE PAGE	GEEZER BUTLER	JIMMIE PAGE
INSTRUMENTS	3,4,5,12	3,4,5,7,12	3,4,5,7,12	3,4,5,12	3,4,12	3,4,5,12	3,4,5,12	3,4,5,12	3,4,5,7,12	3,4,5,12,14
EMOTIONS	M,P,R,S	M,R	M,R	A,M	A	A,P	R	A,P,R	M,R	A,M,S,C

ROCK PSYCHEDELIA	NEO												
TRACK_ID		TRAWAE128E078B013	TRSXLTB128E078B014	TRAFKYR12903CB4810	TRFPWUB128F14539CE	TRDHPGU128F4281235	TRQXCQY128F428122C	TRQZSCU128F4259CCE	TRYOPXW128F4259EAD	TRKSTEJ128F42810C9	TRWMV0B128F42810BF		
TITLE		ALL'S QUIET ON THE EASTER FRONT	THE KKK TOOK MY BABY AWAY	KNOW YOUR ENEMY	AMERICAN IDIOT	NOTHING SOMETHING	WE ARE ONE	THE SOUND OF SINNERS	SOMETHING ENGLAND	YOU ONLY LIVE ONCE	REPTILIA		
ARTIST		RAMONES	RAMONES	GREEN DAY	GREEN DAY	THE OFFSPRING	THE OFFSPRING	THE CLASH	THE CLASH	THE STROKES	THE STROKES		
ALBUM		PLEASANT DREAMS	PLEASANT DREAMS	21ST BREAKDOWN	AMERICAN IDIOT	IGNITION	IGNITION	SANDINISTRA	SANDINISTRA	FIRST IMPRESSIONS OF EARTH	ROOM OF FIRE		
YEAR		1981	1981	2009	2004	1993	1993	1980	1980	2006	2003		
COUNTRY		US	US	US	US	US	US	US	US	US	US		
COMPOSER		DEE DEE RAMONES	JOEY RAMONE	BILLIE JOE ARMSTRONG	BILLIE ARMSTRONG	JOE MARVIN FERGUSON	THE OFFSPRING	THE CLASH	JOE STRUMMER	JULIAN CASABLANCA	JULIAN CASABLANCA		
INSTRUMENTS		1,3,4,5,12,17	3,4,5,12	3,4,5,12	3,4,5,12	3,4,5,12	3,4,5,12	2,4,5,6,9,10	3,4,5,7,10,12	1,3,4,5,12	1,3,4,5,12		
EMOTIONS		A,M	A,M,P	R	A,R	A,M	A,R	A,M	A	A,M,P,S	A,M		

Apêndice C – Formulário de pesquisa de opinião da plataforma

Seção 1 de 3

Pesquisa de Opinião - Plataforma Ritmo Brasil

B *I* U  

Este questionário é parte integrante de uma pesquisa do Programa de Pós-Graduação em Gestão de Informação da Universidade Federal do Paraná (PPGGI/UFPR).

A pesquisa tem como foco a aplicação de inteligência artificial na análise textual de letras de músicas brasileiras.

Antes de responder ao questionário solicitamos acessar a plataforma [RITMO BRASIL](#) e/ou assistir o [vídeo explicativo](#) da plataforma onde foram pre-processadas e filtradas 80mil músicas com suas letras e links para sua audição e consulta.

Objetivo: O objetivo deste instrumento é avaliar e identificar o uso da plataforma "Ritmo Brasil", desenvolvida com algoritmos de inteligência artificial para a análise musical.

Termos de Participação: Solicitamos que responda ao questionário de forma objetiva e sincera.

Asseguramos que todas as informações coletadas são confidenciais, de acesso exclusivo do pesquisador e não serão analisadas de forma isolada.

Caso não concorde com as condições da pesquisa, poderá fechar o formulário e optar por não participar.

E-mail *

E-mail válido

A minha faixa etária encontra-se no intervalo: *

B *I* U  

- ☐ Menos de 18 anos
- ☐ 18 a 29 anos
- ☐ 30 a 39 anos
- ☐ 40 a 49 anos
- ☐ 50 a 59 anos
- ☐ 60 a 69 anos
- ☐ 70 ou mais

Qual sua formação escolar? *

- ☐ Ensino básico (1 a 9 ano)
- ☐ Ensino médio / técnico
- ☐ Ensino superior
- ☐ Pós-graduação (Especialização / Mestrado / Doutorado)

Qual sua principal área de formação, no caso de ter cursado um curso superior?

- ☐ Biológicas / Médicas
- ☐ Exatas / Engenharias
- ☐ Humanas / Sociais Aplicadas
- ☐ Artes

...

Qual sua experiência com a música? *

- ☐ Sou apenas ouvinte
- ☐ Estudei no ensino básico a disciplina de música
- ☐ Fiz curso de instrumentos musicais (violão, piano, guitarra, bateria, órgão, outros)
- ☐ Cursei uma graduação de música
- ☐ Fiz uma especialização em música
- ☐ Minha experiência com música é prática e independente

Tem experiência profissional na área da música? *

- ☐ Sim, sou músico(a)
- ☐ Sim, sou compositor(a)
- ☐ Sim, sou músico(a) e compositor(a)
- ☐ Não

Quais gêneros musicais da MPB são seus preferidos? *

- | | |
|---------------------------------------|---------------------------------------|
| <input type="checkbox"/> Axé | <input type="checkbox"/> POP |
| <input type="checkbox"/> Bossa Nova | <input type="checkbox"/> Rock BR |
| <input type="checkbox"/> Forró | <input type="checkbox"/> Samba |
| <input type="checkbox"/> Funk BR | <input type="checkbox"/> Sertanejo |
| <input type="checkbox"/> Gospel | <input type="checkbox"/> Velha Guarda |
| <input type="checkbox"/> Infantil | |
| <input type="checkbox"/> Jovem Guarda | |
| <input type="checkbox"/> MPB | |
| <input type="checkbox"/> Pagode | |

☐ Sim, sou músico(a)

☐ Sim, sou compositor(a)

☐ Sim, sou músico(a) e compositor(a)

☐ Não

☐ Axé

☐ Bossa Nova

☐ Forró

☐ Funk BR

☐ Gospel

☐ Infantil

☐ Jovem Guarda

☐ MPB

☐ Pagode

<> ...

Após o acesso solicitamos responder as questões seguintes classificando na escala de Likert entre 0 e 10 cada funcionalidade.

[illegible]

1 2 3 4 5 6 7 8 9 10

Muito Insuficiente ○ ○ ○ ○ ○ ○ ○ ○ ○ ○ Muito Satisfatório

Muito Insuficiente

Muito Satisfatório

Muito Insuficiente

Muito Satisfatório

Muito Insuficiente

Muito Satisfatório

Opção de Consulta por TERMOS / PALAVRAS



A consulta **TERMOS** que destaca termos ou palavras mais significativos nas letras, incluindo contagens de ocorrências, palavras relacionadas, árvores de palavras, gráficos por ano e as 20 músicas mais relacionadas ao termo selecionado.

2.6 Avaliar a relevância dessa funcionalidade, consulta de palavras em músicas, para buscar ^{*} inspiração em composições de letras.

	1	2	3	4	5	6	7	8	9	10	
Muito Insuficiente	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	Muito Satisfatório

...

2.7 Os gráficos de distribuição de lançamentos das músicas por ano são úteis para ^{*} identificar as músicas (ou tendências e contextos)?

	1	2	3	4	5	6	7	8	9	10	
Muito Insuficiente	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	Muito Satisfatório

2.8 De um modo geral, a plataforma é útil para o fim de composição de novas letras ^{*} musicais brasileiras?

	1	2	3	4	5	6	7	8	9	10	
Muito Insuficiente	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	Muito Satisfatório

2.9 Sugerir alguma funcionalidade e/ou dados a serem acrescentados na plataforma ^{*} comente abaixo.

Texto de resposta longa

Texto de resposta longa

...

2.10 Opinar sobre a utilidade da plataforma como subsídio a estudos da musica ou como ^{*} objeto de inspiração de composições e/ou produções musicais?

Texto de resposta longa

Apêndice D – Tabela de conversão entre AMNOcode e DNAbits.

Tabela C.1 - Conversão de caracteres código ASCII para AMINOcode e DNAbits

AMINOcode					
Letters	Minimal/ Expanded		Numbers	Expanded	Minimal
a	YA		0	YDA	YD
b	E		1	YDQ	
c	C		2	YDT	
d	D		3	YDH	
e	YE		4	YDF	
f	F		5	YDI	
g	G		6	YDS	
h	H		7	YDE	
i	YI		8	YDG	
j	I		9	YDN	
k	K				
l	L				
m	M				
n	N				
o	YQ				
p	P				
q	Q				
r	R				
s	S				
t	T				
u	YV				
v	V				
x	W				
z	A				
w	YW				
y	YY				
			Punctuation	Expanded	Minimal
			.	YPE	YP
			,	YPC	
			;	YPS	
			!	YPW	
			?	YPQ	
			:	YPT	
			Other	Minimal/Expanded	
			Space	YS	
			Any different	YK	
DNAbits					
Bits		Nucleotide Representation			
00		A			
10		C			
01		G			
11		T			

Apêndice E – Métricas de avaliação

A escolha das métricas de avaliação usadas para a rotulação de música pode-se observar no comparativo dos métodos, que são bem variadas. Mas de acordo com a natureza e características em comum como: fontes de informação do conjunto de dados, algoritmos a serem utilizados e similaridade do método a proposto serão utilizadas as seguintes métricas:

Precisão (precision): é o número de classificações positivas corretas (Verdadeiros Positivos), dividido pelo número de exemplos classificados na mesma classe (a soma dos resultados de Verdadeiros Positivos e Falsos Positivos). Deste modo V_p pode ser n_{11} , n_{22} , $n_{::}$, n_{cc} e F_p pode ser n_{21} , $n_{::}$, n_{c1} . Conforme a equação 1 (POWERS, 2011).

$$Precisao = \frac{V_p}{V_p + F_p} \quad (1)$$

Exatidão/Acurácia: é o número de classificações corretas (a soma dos resultados Verdadeiros Positivos e Verdadeiros Negativos) dividido pelo número total de classificações (a soma de todos os itens). Conforme representa a equação 2 (POWERS, 2011).

$$Exatidao = \frac{V_p + V_n}{(V_p + V_n + F_p + F_n)} \quad (2)$$

Revocação (recall): é o número de classificações positivas corretas (Verdadeiros Positivos), dividido pelo número real de exemplos da classe (a soma dos resultados de Verdadeiros Positivos e Falsos Negativos). Neste contexto, V_p segue a definição acima e F_n pode ser n_{12} , n_{13} , $n_{::}$, n_{1c} . Conforme a equação 3 (POWERS, 2011).

$$Revocacao = \frac{V_p}{V_p + F_n} \quad (3)$$

Apêndice F – Manual de instrução de consulta Ritmo Brasil

Manual de Instruções para Acesso e Consulta ao Site "Ritmo Brasil" de Letras de Músicas Brasileiras

Este manual fornece orientações detalhadas e estruturadas para acessar e utilizar o site "Ritmo Brasil", disponível em <http://200.236.3.21/index.html>. O recurso é projetado para profissionais da música — compositores, cantores, artistas e estudiosos, como musicólogos —, facilitando consultas a um corpus de letras de músicas brasileiras. As instruções seguem a sequência lógica de navegação, com referências às imagens fornecidas para ilustração. Para maior completude, as telas anexadas foram incorporadas por meio de descrições detalhadas, com base nas capturas de tela fornecidas, incluindo elementos visuais-chave, layouts e conteúdos exibidos. Essas descrições complementam as referências às imagens, permitindo uma compreensão visual mesmo em formatos textuais.

1. Acesso à Página Principal

- URL de Acesso Inicial: Acesse o site diretamente pelo link <http://200.236.3.21/index.html>.



- **Descrição da Interface:**
A página principal exibe o título "Ritmo Brasil" e o subtítulo "Navegue pela onda da Música Popular Brasileira". Ela apresenta uma grade de ícones que representam gêneros musicais, com a opção "Todos os GÊNEROS" no topo para consultas gerais. Os gêneros disponíveis incluem: axé, bossa nova, forró, funk carioca, gospel, infantil, jovem guarda, mpb, pagode, pop, rock, samba, sertanejo e velha guarda.
- **Descrição da Tela Anexada:**
Título principal "RITMO BRASIL".

Abaixo, uma grade de cartões quadrados ilustra cada gênero com ícones estilizados. Cada ícone apresenta o nome do gênero e, abaixo, links para "Músicas" e "Termos" em azul.

- **Procedimento Inicial:**

Selecione "Todos os GÊNEROS" para uma visão ampla do corpus ou clique em um gênero específico para restringir as consultas. Essa escolha determina o escopo das buscas subsequentes.

2. Opções de Consulta de cada Gênero

- Estrutura Geral: Abaixo de cada ícone de gênero (ou na seção "Todos os GÊNEROS"), há dois links principais para consultas:
- Músicas: Direciona para uma página de listagem de letras de músicas.

Exemplo para bossa-nova:

http://200.236.3.21/brazilian_lyrics_html_tm_bossa-nova/html_tm/TEXTS.html.



#	Título	Artista	Gênero	Ano	Letras relacionadas	Link	EU IA
0	Em um bar	Família da Música	bossa nova	1962	Letras relacionadas	Link	1006095
1	Noite e dia	Agostinho dos Santos	bossa nova	1969	Letras relacionadas	Link	1004704
2	Ontem	Agostinho dos Santos	bossa nova	1967	Letras relacionadas	Link	1004779
3	Última Valsa	Lisa Ono	bossa nova	2008	Letras relacionadas	Link	1006471
4	Brasil Com P	Maria Rita	bossa nova	2007	Letras relacionadas	Link	1006569
5	Elegia Desesperada	Vinício de Moraes	bossa nova	1977	Letras relacionadas	Link	1007430
6	O Desespero da Piedade	Vinício de Moraes	bossa nova	1980	Letras relacionadas	Link	1007460
7	Rosa de Hiroxima	Vinício de Moraes	bossa nova	1972	Letras relacionadas	Link	1007492
8	O desejado	Elizabeth Cardoso	bossa nova	1978	Letras relacionadas	Link	1005951
9	Doze Anos	Chico Buarque	bossa nova	1979	Letras relacionadas	Link	1005102
10	Querelas do Brasil	Elis Regina	bossa nova	1978	Letras relacionadas	Link	1005652
11	Oração Caribe	Agostinho dos Santos	bossa nova	1973	Letras relacionadas	Link	1004724
12	Amor de Mis Amores	Agostinho dos Santos	bossa nova	1973	Letras relacionadas	Link	1004604
13	Um Dia de Cao	Chico Buarque	bossa nova	1994	Letras relacionadas	Link	1005360
14	O Caçador de Esmeralda	Elis Regina	bossa nova	1973	Letras relacionadas	Link	1005612

- Termos: Direciona para uma página de análise de termos e de palavras significativas. Exemplo para bossa-nova:

http://200.236.3.21/brazilian_lyrics_html_tm_bossa-nova/html_tm/WORDS.html.

Voltar

Alternar tema

Em um bar - Documentos relacionados

consulta de pesquisa

Número dos vizinhos

Procurar

Ajuda

Total de inscrições: 20

Entradas selecionadas: 20

Exportar CSV

Limpar filtros

#	Classificação	Semelhança	Título + Letra	Artista	Gênero	Ano	Link	EU IA
0	1	2	3	4	5	6	7	8
			Num bar , um monte de bêbados bebendo cerveja, bebendo cerveja, esperando outra. O dono disse: "Você vai ter que repetir. Você só vai tomar a cerveja que sobrou." "Ai um cara deu um soco na mesa, fazendo canecas, garfos e tudo voarem. Jones, parado bem do meu lado, disse: "Você vai ter que dar um soco na sua testa." Até que eles beberam juntos, beberam seu melhor uísque, como amigos. Ah, sim, naquela época elas brigavam, só diziam "bons tempos", e fizeram as pazes. Sempre há tempo para mudar. Só porque fizemos algo estúpido quando éramos jovens, vamos deixar a estupidez para sempre. Sempre há tempo para mudar. "Você vai ter que tomar a cerveja que sobrou", disse o dono.	Família da Música	bossa nova	1962	Link	1006095
0	0	1,0000						
40	1	0,5501	Tete um milhão de estrelas um milhão dormindo um milhão de sentimentos vá embora outra voce sabe tudo ta bem mas eu lhe garanto que o mundo volta também hoje sou eu quem vai dizer voce que sabia demais nao viu que o tempo passou e fez de voce nunca mais nao chora nao pede segue em paz.	Wanda Sa	bossa nova	2002	Link	1007590
281	2	0,5419	Dez Leis que vai, deuses, bom, que vai, deuses, bom que vai, deuses, bom que lei, leis, bom jogo, lei, leis, bom, lei, leis, bom, jogo, lei lei e meu num gueto ao norte meu e meu num gueto ao norte meu e meu num gueto ao norte meu e meu são dez das leis mas um so rei são dez das leis mas um so rei são dez das leis mas um so rei	Marcos Valle	bossa nova	1970	Link	1006503

Recomendação: Escolha o link apropriado com base no objetivo: "Músicas" para buscas por composições específicas e "Termos" para análises linguísticas e de tendências.

3. Consulta pela Opção "Músicas"

- Descrição da Página: Ao selecionar "Músicas" em um gênero ou em "Todos os GÊNEROS", a página exibe uma tabela com todas as músicas disponíveis no escopo selecionado. Essa interface permite ordenação e filtragem.
- Colunas da Tabela: A tabela inclui as seguintes colunas:
 - Número sequencial.
 - ID: Identificador único da música.
 - Title: Título da música.
 - Artist: Artista ou compositor.
 - Gender: Gênero musical (pode ser redundante em consultas por gênero específico).
 - Year: Ano de lançamento ou de composição.
 - Related Lyrics: Link para músicas relacionadas.
 - Link: Link externo para o site de origem (ex.: Vagalume ou Letras.mus.br).
- Funcionalidades de Ordenação: Clique no cabeçalho de qualquer coluna para ordenar os dados em ordem ascendente ou descendente (alfabética para textos ou numérica para valores como ano ou ID).
- Aspecto Principal: Related Lyrics: Clique no link "Related Lyrics" de uma música específica para abrir uma nova página que lista as 20 músicas mais relacionadas, com base em termos compartilhados. Essa página inclui colunas como Rank, Similarity, Title + Lyrics, Artist, Gender, Year, Link e ID. A coluna "Similarity" indica o grau de proximidade semântica, auxiliando nas análises comparativas.

Voltar

Alternar tema

Em um bar - Documentos relacionados

consulta de pesquisa

Número dos vizinhos

Procurar

Ajuda

Total de inscrições: 20 Entradas selecionadas: 20

Exportar CSV

Limpar filtros

#	Classificação	Semelhança	Título + Letra	Artista	Gênero	Ano	Link	EU IA
0	1	2	3	4	5	6	7	8
0	0	1,0000	Num bar , um monte de bêbados bebendo cerveja, bebendo cerveja, esperando outra. O dono disse: "Você vai ter que repetir. Você só vai tomar a cerveja que sobrou." "Ai um cara deu um soco na mesa, fazendo canecas, garfos e tudo voarem. Jones, parado bem do meu lado, disse: "Você vai ter que dar um soco na sua testa." Até que eles beberam juntos, beberam seu melhor uísque, como amigos. Ah, sim, naquela época eles brigavam, só diziam "bons tempos", e fizeram as pazes. Sempre há tempo para mudar. Só porque fizemos algo estúpido quando éramos jovens, vamos deixar a estupidez para sempre. Sempre há tempo para mudar. "Você vai ter que tomar a cerveja que sobrou", disse o dono.	Família da Música	bossa nova	1962	Link	1006095
40	1	0,5501	Tete um milhão de estrelas um milhão dormindo um milhão de sentimentos vá embora outra voce sabe tudo ta bem mas eu lhe garanto que o mundo voltas tambem hoje sou eu quem vai dizer voce que sabia demais nao viu que o tempo passou e fez de voce nunca mais nao chora nao pede segue em paz	Wanda Sa	bossa nova	2002	Link	1007590
281	2	0,5419	Dez Leis que vai, deuses, bom, que vai, deuses, bom que vai, deuses, bom que lei, leis, bom jogo, lei, leis, bom, lei, leis, bom, jogo, lei lei e meu num gueto ao norte meu e meu num gueto ao norte meu e meu num gueto ao norte meu e meu sao dez das leis mas um so rei sao dez das leis mas um so rei sao repete ate o fim sao dez das leis mas um so rei	Marcos Valle	bossa nova	1970	Link	1006503

Descrição da Tela Anexada Documentos relacionados:

- A tabela lista 20 entradas com as colunas "", "Classificação", "Similaridade", "Título + Letra", "Artista", "Gênero", "Ano", "Link".
- Exemplos incluem títulos como "Num bar" (similaridade 1.0000), artista "Família da Música", gênero "bossa nova" e ano 1962. Há opções como "Procurar", "Ajuda", "Exportar CSV" e "Limpar filtros".
- Outros Links: A coluna "Link" direciona para fontes externas, como Vagalume e Letras.mus.br para detalhes adicionais.

4. Consulta pela Opção "Termos"

- Ao selecionar "Termos" em um gênero ou em "Todos os GÊNEROS", a página exibe uma tabela com termos (palavras) extraídos das letras, incluindo frequências e análises associadas.

Alternar tema

Letra brasileira de BOSSA NOVA - Palavras

consulta de pesquisa

Número dos vizinhos

Procurar

Ajuda

Total de inscrições: 6632 Entradas selecionadas: 6632

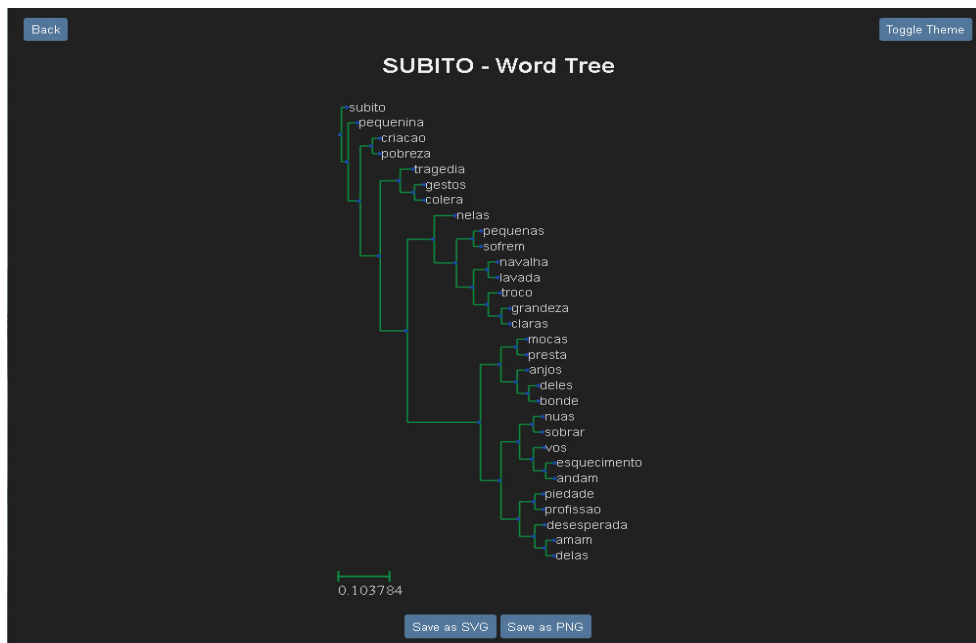
Exportar CSV

Limpar filtros

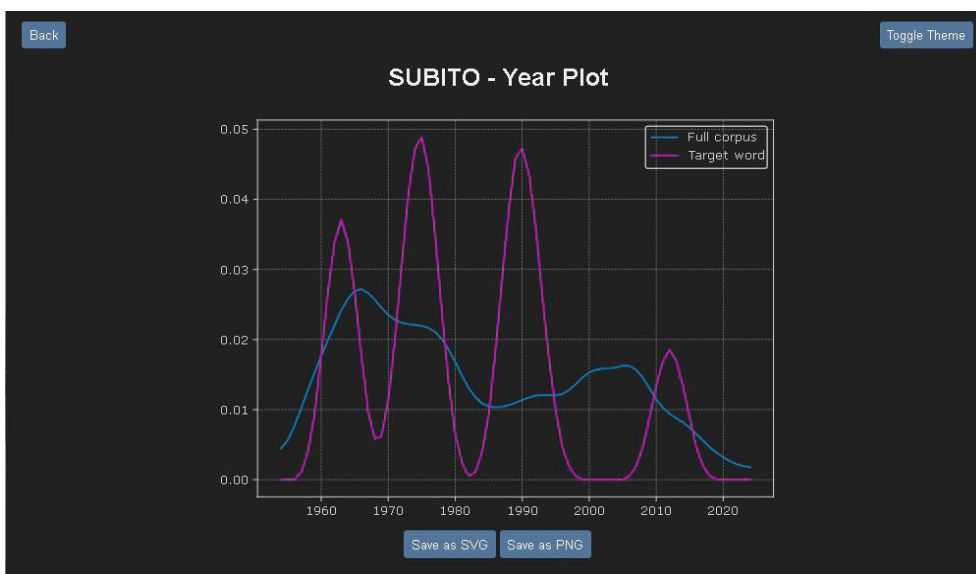
#	Palavra	Ocasionalmente.	Palavras relacionadas	Letras relacionadas	Árvore de Palavras	Gráfico anual
0	1	2	3	4	5	6
0	SUBITO	6	subito, pequenina, esquecimento, criação, andam, gestos, pobreza, nuas, amam, vos	Letras relacionadas	Árvore de Palavras	Gráfico anual
1	PEQUENINA	7	pequenina, criação, subito, esquecimento, pobreza, deles, nuas, andam, navalha, bonde	Letras relacionadas	Árvore de Palavras	Gráfico anual
2	CRIAÇÃO	6	criação, pobreza, pequenina, subito, gestos, colera, tragédia, andam, grandeza, suicídio	Letras relacionadas	Árvore de Palavras	Gráfico anual
3	VOS	5	vos, esquecimento, claras, nuas, deles, anjos, andam, subito, grandeza, pequenina	Letras relacionadas	Árvore de Palavras	Gráfico anual
4	ANJOS	6	anjos, deles, bonde, sobrar, pobres, vos, meninos, amam, espírito, presta	Letras relacionadas	Árvore de Palavras	Gráfico anual
5	DELES	5	deles, bonde, anjos, vos, navalha, grandeza, pequenina, claras, esquecimento, sobrar	Letras relacionadas	Árvore de Palavras	Gráfico anual

- Colunas da Tabela: A tabela inclui as seguintes colunas

A tela exibe "SUBITO - Word Tree", um diagrama em árvore verde sobre fundo escuro, ramificando-se a partir de "súbito" para palavras relacionadas, como "pequenina", "criação", "tragédia", "gestos", "nelas", etc., e terminando em "delas". Há opções para "Save as SVG" e "Save as PNG"



- Year Plot: Clique no link para visualizar um gráfico de linha que ilustra a frequência do termo ao longo dos anos, comparando-a ao corpus em geral. Isso destaca picos de uso, facilitando a análise de tendências temporais. O gráfico pode ser salvo em formato SVG ou PNG.



- Descrição da tela Year Plot: Apresenta o termo "SUBITO - Year Plot" e um gráfico de linhas com eixos de 1960 a 2020 (x) e de 0.00 a 0.05 (y). Duas linhas: azul para "Full corpus" e roxa para "Target word", com picos nos anos 1970-1980. As opções incluem "Save as SVG" e "Save as PNG".

5. Considerações Gerais para Consultas

As consultas gerais ("Todos os GÊNEROS") ou restritas por gênero permitem análises profundas, como tendências linguísticas, evoluções sociológicas e inspirações criativas nas letras de músicas brasileiras.

- Observações Importantes:

- a) Desempenho: Como as páginas são carregadas integralmente na memória RAM do dispositivo, consultas em "Todos os GÊNEROS" podem ser mais lentas do que em gêneros específicos, especialmente em ordenações ou filtros.
- b) Tradução: As páginas podem ser traduzidas automaticamente para o português no navegador (ex.: no Google Chrome, ative a opção de tradução integrada).
- c) Botão "Help": presente na maioria das páginas, oferece orientações sobre buscas e exportações.
- Ajuda: Barra de Pesquisa
Este projeto fornece um utilitário de busca em tabelas que permite aos usuários filtrar linhas com base em consultas complexas, com suporte a operadores lógicos (AND/OR) e a buscas específicas por coluna. Abaixo, segue uma explicação sobre como usar a funcionalidade de busca.

Características

- Pesquisa em todas as colunas: Insira uma consulta que procure todas as colunas da tabela.
- Pesquisa específica por coluna: Use colchetes ([index]) para limitar a pesquisa a uma coluna específica.
- Operadores lógicos: Combine termos com AND ou OR para consultas mais complexas.
- Termos correspondentes em destaque: Os termos encontrados estão destacados na tabela para maior clareza.
- Análise dinâmica de consultas: Suporta padrões de pesquisa avançados com expressões lógicas.
- Correspondência sem distinção entre maiúsculas e minúsculas: As pesquisas não distinguem as maiúsculas das minúsculas.
- Expansão de vizinhos: Inclua linhas acima e abaixo dos resultados correspondentes para fornecer contexto adicional.
- Como usar, siga estes passos simples para filtrar sua tabela:
- Digite a consulta: Insira uma consulta de pesquisa no campo de entrada, na sintaxe desejada.
- Definir vizinhos (opcional): Use o campo de entrada "Neighbors" para especificar quantas linhas acima e abaixo dos resultados correspondentes também devem ser exibidas
- Busca por gatilho: O filtro é aplicado imediatamente ao pressionar Enter.
- Resultados da revisão: As linhas correspondentes, juntamente com os vizinhos especificados, permanecerão visíveis e os termos correspondentes serão destacados.
- Sintaxe de Pesquisa
 - Pesquisa Básica: Digite uma palavra-chave na barra de pesquisa para filtrar as linhas que contêm esse termo em qualquer coluna. Exemplo: exemplo.
 - Pesquisa Específica por Coluna: Para pesquisar em uma coluna específica, use colchetes com o índice da coluna. Exemplo: exemplo [2] (busca "exemplo" na terceira coluna, índices começam em 0)

- Operadores Lógicos: Combine os termos de pesquisa usando AND ou OR (em maiúsculas). Exemplo: AND teste (linhas com ambos os termos). Exemplo: exemplo[1] OR teste[2] (linhas com "exemplo" na coluna 2 ou "teste" na coluna 3)
- Expansão do Bairro: Especifique o número de vizinhos no campo "Neighbors". Exemplo: Consulta: example; Vizinhos: 1 (exibe a linha correspondente, mais uma acima e outra abaixo).
- Parênteses para Agrupamento: Use parênteses para agrupar consultas complexas. Exemplo: (exemplo AND teste) OR amostra.

Este manual pode ser atualizado conforme as evoluções no site. Para suporte adicional, explore o botão "Help" ou contate o administrador do site se disponível.