

Universidade Federal de Campina Grande
Centro de Engenharia Elétrica e Informática
Coordenação de Pós-Graduação em Ciência da Computação

Explorando as relações entre os aspectos de
novidades musicais e as preferências pelos ouvintes

Andryw Marques Ramos

Dissertação submetida à Coordenação do Curso de Pós-Graduação em
Ciência da Computação da Universidade Federal de Campina Grande -
Campus I como parte dos requisitos necessários para obtenção do grau
de Mestre em Ciência da Computação.

Área de Concentração: Ciência da Computação
Linha de Pesquisa: Metodologias e técnicas da computação

Nazareno Ferreira de Andrade
(Orientador)

Campina Grande, Paraíba, Brasil
©Andryw Marques Ramos, 00/00/2014

Resumo

Seu resumo aqui

Abstract

Abstract Here

Agradecimentos

Agradecimentos

Conteúdo

1	Introdução	1
1.1	Organização do documento	3
2	Novidades e descobertas no âmbito musical	4
2.1	Novidade no consumo em geral	4
2.2	Incorporação de novidade em sistemas computacionais	4
2.3	Grupos de pessoas baseados no comportamento	5
2.4	Nossas contribuições	6
3	Características da novidade	7
4	Dados	10
4.1	Last.FM	10
4.2	Ouvinte	11
4.2.1	Timeline	11
4.2.2	Histórico do usuário	12
4.2.3	Filtros	13
4.3	Metadados dos Artistas	14
5	Modelos	17
5.1	Perfil musical do ouvinte	17
5.1.1	Ecleticidade	20
5.2	Características das novidades	21
5.2.1	Familiaridade	21
5.2.2	Popularidade	22

5.3	Preferências	22
6	Preferências dos ouvintes para diferentes aspectos de novidades	24
6.1	Preferências gerais	24
6.2	Preferências individuais	25
7	Grupos de usuários para diferentes aspectos de novidade	27
7.1	Conjunto de sujeitos	27
7.2	Dados que caracterizam sujeitos	28
7.3	Algoritmo de agrupamento	29
7.4	Escolha do número de grupos	29
7.5	Grupos	30
7.6	Discussão dos grupos encontrados	32
7.7	Discussão dos resultados	33
8	Comparação das novidades com os artistas conhecidos	34
8.1	Características dos artistas conhecidos	34
8.2	Comparação entre relações das preferências e aspectos das novidades e dos artistas conhecidos	35
8.3	Grupos de ouvintes baseados na diferença das relações entre preferências e aspectos das novidades e dos artistas conhecidos	38
9	Conclusão	41
10	Trabalhos Futuros	42
A	Apêndices	44

Lista de Símbolos

Lista de Figuras

4.1	Linha do tempo utilizada no trabalho	12
5.1	Curva de popularidade dos artistas do LAST.FM. A curva é enviesada. . . .	22
7.1	Número de grupo X Distância média dentro dos grupos. O joelho do gráfico está em torno da configuração com 7 grupos.	30
7.2	Centróides dos 7 grupos encontrados na análise. As métricas estão normalizadas pelo z-score, onde zero representa a média de todos os ouvintes, e a unidade de variação é um desvio padrão, para cada métrica. No eixo vertical, <i>fam</i> significa a familiaridade, <i>pop</i> significa popularidade, <i>AT</i> significa atenção total.	31
8.1	Número de grupo X Distância média dentro dos grupos.	39
8.2	Centróides dos 6 grupos encontrados na análise da diferença das correlações	39

Lista de Tabelas

4.1	Tags do artista Michael Jackson, junto com o valor normalizado de cada uma.	15
4.2	Número de ouvintes (popularidade) de alguns artistas no LAST.FM	16
6.1	Correlação (Coeficiente de Spearman) entre aspectos da novidade e preferências, analisando todas as novidades juntas	25
8.1	Correlações calculadas para artistas com novidade e artistas conhecidos. . .	35
8.2	Teste-T pareado entre as correlações dos aspectos e preferências dos ouvintes para artistas com novidade e artistas conhecidos.	36
8.3	Correlação (Coeficiente de Kendall) entre correlações calculadas para novidades (linhas) e artistas conhecidos (colunas)	37

Capítulo 1

Introdução

A procura e descoberta de novas músicas e artistas é um aspecto importante no consumo musical. Maddi [2] argumenta que consumidores em geral possuem um *impulso interno*, que tem como finalidade descobrir novas experiências, afim de criar novos sentimentos e emoções.

E este consumo musical mudou nos últimos anos. Serviços de streaming como, Spotify, Youtube, Soundcloud, rádios online como a do Last.Fm, até mesmo os sites de compra de música digital, como Itunes, Beatport, possibilitam o acesso a uma grande variedade de músicas. Isso facilita o acesso a músicas e artistas não escutados antes pelo ouvinte, as chamadas novidades. Porém, como há uma grande quantidade de novidades, encontrar aquelas que sejam relevantes acaba sendo uma tarefa custosa. Muitas vezes o ouvinte acaba não descobrindo novidades que seriam de seu interesse.

Para resolver este problema, tanto sistemas comerciais como a academia apresentam soluções, principalmente incorporando a sistemas de recomendação o domínio da *novidade*. Mas boa parte destas abordagens tratam novidade de forma unidimensional. Podemos caracterizar uma novidade com diferentes dimensões, ou aspectos, como a familiaridade e a popularidade. Por exemplo, um ouvinte pode preferir novidades similares, ou familiares, a músicas que ele costuma escutar, mas preferir novidades não populares, e vice-versa.

Com o intuito de expandir o entendimento sobre as novidades, conduzimos uma análise sobre o impacto dos diferentes aspectos das novidades para as preferências de ouvintes de música. Para auxiliar esta análise, procuramos responder as seguintes perguntas de pesquisa:

1. Há alguma relação geral entre algum aspecto das novidades e as preferências dos ou-

vintes?

2. Individualmente, os usuários preferem algum aspecto das novidades?
3. Existem grupos de ouvintes baseados nas preferências pelos aspectos das novidades comparadas com seus hábitos musicais?
4. As relações entre as preferências dos ouvintes e os aspectos das novidades são as mesmas que as relações entre as preferências dos ouvintes e os aspectos das músicas já conhecidas?

A primeira pergunta tenta descobrir se todos os ouvintes preferem algum aspecto específico da novidade. O objetivo foi encontrar respostas como: "Os ouvintes no geral preferem novidades familiares a outras músicas escutadas anteriormente (um ouvinte de rock prefere novidades de rock a novidades de rap)", "Os ouvintes no geral preferem novidades menos populares" ou "não existe uma preferência no geral". A segunda pergunta é uma especificação da primeira, em um âmbito individual. O objetivo foi encontrar respostas como: "75% dos ouvintes possuem preferência por algum aspecto, sendo que 15% preferem novidades familiares, 30% preferem não-familiares, etc."

A terceira pergunta tenta encontrar grupos de ouvintes, baseados nessas preferências pelos aspectos das novidades, junto com algumas características dos hábitos musicais do ouvinte. Encontrar usuários que compartilhem as mesmas características possibilitam que ferramentas, como recomendadores, os tratem de forma diferente.

Já a quarta pergunta verifica se as preferências dos ouvintes pelos aspectos das novidades são semelhantes às preferências dos ouvintes pelos aspectos dos artistas os quais os ouvintes já tinham escutados anteriormente - as conhecidas. Relacionamos estas preferências no mesmo período de tempo para entender se o comportamento dos ouvintes é o mesmo para ambos os tipos de artistas escutados ou se há alguma diferença.

Para responder as perguntas de pesquisa, coletamos dados históricos referentes à escuta de música de usuários do Last.FM, junto com metadados que caracterizam os artistas. No nosso estudo, as novidades e as conhecidas são artistas. Com os dados históricos, conseguimos identificar as novidades, as conhecidas e as preferências dos ouvintes. Com os metadados, conseguimos identificar os seus aspectos - familiaridade e popularidade.

Descobrimos que não há uma correlação geral entre a familiaridade ou popularidade e as preferências das novidades pelo ouvinte. Porém, individualmente, boa parte dos ouvintes preferem um e/ou outro aspecto das novidades.

Como boa parte dos ouvintes preferem algum aspecto da novidade, conseguimos identificar 7 grupos de ouvintes baseados na em 5 características: relação entre a familiaridade e as preferências, a relação entre a popularidade e as preferências, a ecleticidade do ouvinte, a popularidade dos artistas escutados pelo ouvinte e a proporção de novidades que ele escutou, no período observado.

Já no âmbito das novidades e conhecidas, os ouvintes possuem diferentes preferências pelos aspectos da primeira, comparados com os aspectos da segunda. No geral, a preferência por artistas familiares e populares é maior para as conhecidas que para as novidades.

Com nossos resultados, indicamos a necessidade do tratamento multi-dimensional das novidades. Claramente os ouvintes possuem preferências diferentes para diferentes aspectos das novidades. Isso pode ajudar no aperfeiçoamento de sistemas de recomendação de novidades musicais. Seria necessário um recomendador personalizado, não apenas baseado no histórico de execução do ouvinte, mas também nas preferências ou pela familiaridade ou pela popularidade das novidades.

A descoberta dos grupos pode permitir que os sistemas desenvolvam soluções diferentes para usuários diferentes. Por exemplo, construir interfaces diferentes para cada grupo, algoritmos diferentes de recomendação ou até direcionamento diferente de notícias.

Por fim, o comportamento diferenciado do ouvinte para os aspectos das novidades e para os aspectos das conhecidas corrobora que a novidade é um âmbito especial no consumo de música, precisando ser tratada de maneira específica.

1.1 Organização do documento

/* Será escrita no final (Ima) */

Capítulo 2

Novidades e descobertas no âmbito musical

2.1 Novidade no consumo em geral

Vários trabalhos estudam o comportamento de novidade e diversidade no consumo das pessoas. Os consumidores em geral podem possuir comportamento com o intuito de manter a consistência de suas escolhas ou buscar a variedade / novidade, em diferentes âmbitos (Fishbach et. al). Devido a existência desta faceta adicional do consumo das pessoas, construtores de sistemas computacionais, como sistemas de recomendação, estão cada vez mais preocupados em incorporá-la nestes sistemas.

2.2 Incorporação de novidade em sistemas computacionais

Muitos dos algoritmos utilizados em sistemas de recomendação de itens e produtos estão interessados em aumentar a acurácia e precisão dos resultados. O objetivo principal é maximizar a quantidade de itens que o consumidor irá gostar. Esses algoritmos, como a filtragem colaborativa e filtragem baseada em conteúdo, recomendam itens já consumidos ou parecidos com os que já foram consumidos. Esta forma de recomendação, apesar de aumentar a acurácia dos resultados de treinamento, pode gerar um problema, já que negligencia o fator da novidade / diversidade, o qual está incorporado no comportamento das pessoas. Recomendar "mais do mesmo" deixa os usuários do sistema entediados (Zhang et. al, Herlocker

et. al, Morgan et. al).

Relacionado com essa novidade/diversidade em sistemas computacionais, Vargas et al. propõem um modelo que relaciona usuários, itens e novidade. Existe uma diferença entre a descoberta (o usuário conhece um item, que deixa de ser novidade), a relevância (item de interesse do usuário) e a escolha (quando o usuário seleciona um item relevante). Além disso, eles apontam que em geral as soluções envolvendo novidades de itens são apresentadas em dois modelos: o modelo baseado em popularidade e o baseado na similaridade de itens previamente expostos. Este tipo de visão da novidade como dois modelos também é corroborado por Beloggin et. al.

O modelo baseado em popularidade define que a popularidade do item está relacionada à sua descoberta pelas pessoas. Quanto menos popular um item, menos ele foi descoberto pelas pessoas, possuindo uma maior probabilidade de ser novidade para a maioria das mesmas (O' Celma). Vários trabalhos utilizam métricas relacionadas com a popularidade para detectar o quanto os algoritmos estudados expõem novidades (O' Celma, Ziegler et. al, Vargas et. al, Beloggin et. al, Zhang et. al).

O modelo baseado em similaridade define que há uma maior probabilidade de um item ser novidade para um usuário se este não for similar a outros itens descobertos e escolhidos pelo mesmo. Boa parte das abordagens baseadas neste tipo de modelo definem classes de itens baseadas na similaridade dos mesmos. Assim, é recomendado para o usuário itens de classes que não são similares a classes anteriormente escolhidas pelo usuário (Ziegler et. al, Nakatsuji et al, Zhang et. al).

2.3 Grupos de pessoas baseados no comportamento

Outro tema estudado no âmbito de novidades é a existência de diferentes grupos com diferentes preferências para novidades. Por exemplo, Munson e Resnick descobriram, em um conjunto de usuários online, subgrupos baseados nas preferências por novas opiniões: apreciadores de diversidade, aversos a desafios e procuradores de suporte. Os apreciadores de diversidade são usuários que se interessam tanto por opiniões similares a suas quanto desafiadores. Eles não se satisfazem com apenas opiniões similares. Já os aversos a desafios se satisfazem mais com opiniões semelhantes, diminuindo a satisfação se lerem opiniões de-

safiadoras. Já os procuradores de suporte se satisfazem com um certo número de opiniões semelhantes, que suportam seu ponto de vista, sendo indiferentes a demais opiniões conflitantes. Este resultado mostra que diferentes pessoas possuem diferentes comportamentos frente a novidades (no estudo, novas opiniões).

2.4 Nossas contribuições

Apesar de trabalhos passados utilizarem conceitos de novidade e diversidade, tanto de itens musicais quanto de itens no geral, não foram encontrados estudos que relacionem diferentes aspectos das novidades com a relevância das mesmas para os usuários. Como apontado na Seção 2.2, a novidade pode ser modelada em no mínimo dois aspectos. Além disso, O' Celma comenta que é importante sabermos a relevância de novidades para os usuários, para que tenhamos um conhecimento mais completo do comportamento do usuário. Assim, nós unimos os aspectos formalizados por Vagas e aderessados em vários estudos com a ideia de relevância, que O' Celma corrobora a importância. Com essa junção de conceitos podemos responder que aspectos das novidades são relevantes para ouvintes musicais.

Além disso, outro resultado do nosso trabalho segue a linha de trabalhos como o do Munson. Descobrimos diferentes grupos de ouvintes baseados nas preferências dos mesmos pelos aspectos das novidades. Mesmo em âmbitos diferentes podemos notar que novidade não pode ser tratado de forma única para todos os indivíduos. Há uma necessidade de tratamento específico para cada grupo, em um sistema computacional como sistema de recomendação, por exemplo.

Capítulo 3

Características da novidade

Para a construção dos experimentos, primeiro foi definido o conceito de novidade que iríamos trabalhar, junto com o das suas características. Foi importante esta definição inicial pois os termos utilizados na pesquisa (novidade, familiaridade, popularidade, relevância) são termos gerais, que podem possuir mais de um significado, não tendo um consenso da literatura.

1. Novidade / Item com novidade

Novidade é o conceito central deste trabalho. Um item com novidade é um item que não foi acessado pela pessoa anteriormente. No âmbito musical, itens podem ser músicas, artistas e álbuns, e as pessoas que escutam esses itens são ouvintes. Mais especificamente, tratamos as novidades como artistas que não foram escutados anteriormente pelo ouvinte. Por exemplo, se em algum momento o ouvinte João escutou o artista Eminem pela primeira vez, ele deixou de ser uma novidade para ser um artista conhecido. Antes desse momento Eminem era considerado uma novidade para João.

2. Item conhecido

Um item conhecido é o oposto da novidade. Assim, é um item que já foi acessado anteriormente pela pessoa. No nosso trabalho, um artista conhecido é um artista que já foi escutado anteriormente pelo ouvinte. No Capítulo 8 é mostrado o resultado da comparação entre as preferências dos ouvintes pelas características dos artistas com novidade e as preferências do mesmo pelas características dos artistas conhecidos.

3. Familiaridade

Muitos trabalhos /* **refs (lma)** */ caracterizam uma novidade baseada na similaridade do item acessado pela pessoa em relação a outros itens acessados anteriormente pela mesma. Trazendo para o âmbito musical, rotulamos este tipo de característica como familiaridade. A familiaridade de um artista para um ouvinte reflete o quanto este ouvinte foi exposto a outros artistas que têm descritores semelhantes aos do artista escutado. Um descritor é um símbolo que descreve / caracteriza um artista. No âmbito deste trabalho, descritores são termos que descrevem um artista, como gênero musical (pop, forró), localização (latina, brasileira, britânica), humor (animada, depressiva), entre outros.

Além da similaridade dos descritores, nós levamos em conta o quanto estes artistas similares ao artista em questão foram escutados pelo ouvinte. Isso porque a familiaridade de um artista é influenciado também pelo quanto o ouvinte escutou artistas similares a ele(Hargreaves /* **Re (lma)** */). Por exemplo, se Maria tem hábito de escutar muitos artistas pop e poucos artistas de rock, Britney Spears (cantora pop) é mais familiar a ela que Evanescence (banda de rock). Assim, a familiaridade está relacionada com a similaridade entre os descritores do artista e os descritores dos artistas do perfil do ouvinte junto com a influência desses artistas no perfil do ouvinte.

4. Popularidade

Outra característica bem relacionada com itens com novidade na literatura é a popularidade /* **ref (lma)** */. Neste trabalho definimos a popularidade como sendo o quanto de pessoas já escutaram o artista em questão. Por exemplo, Michael Jackson, artista que muitas pessoas de todo o mundo já escutaram, é mais popular que Rapadura, um rapper brasileiro que foi escutado apenas por um nicho específico de pessoas.

5. Preferência / Relevância

Com os aspectos familiaridade e popularidade, estudamos quais destas características das novidades são relevantes para o ouvinte. Em outras palavras, qual a preferência dos ouvintes por esses aspectos. Utilizamos dois conceitos de preferência:

- (a) Atenção total

É o quanto de atenção que um ouvinte deu para o artista em um período especificado. Para isso utilizamos a quantidade de músicas do artista que o ouvinte escutou no período. Quanto mais músicas do artista, mais atenção o ouvinte devotou ao artista.

(b) Período de atenção

É o período que o ouvinte devotou de atenção ao artista. No nosso trabalho, a unidade de tempo é uma semana. Assim, quanto mais semanas, de um período definido de tempo, o ouvinte escutou alguma música do artista em questão, maior o período de atenção do ouvinte.

Capítulo 4

Dados

Após a descrição das características das novidades (e conhecidas) utilizadas no nosso estudo, esta seção descreve os dados que foram utilizados na pesquisa. Podemos dividir os dados em 2 partes: a primeira é representada pelo histórico musical dos sujeitos a serem analisados e a segunda pelos metadados dos artistas escutados pelos sujeitos. Os sujeitos dos experimentos representam os ouvintes. O histórico musical foi utilizado para identificar as novidades, as conhecidas, e as preferências dos ouvintes por ambas. Já os metadados foram utilizados para identificar os aspectos das novidades/conhecidas - a familiaridade e a popularidade. Os dados foram coletados da plataforma do LAST.FM.

4.1 Last.FM

O Last.FM é uma rede social musical que tem como principal característica o *Scrobbling*, um serviço que permite registrar o histórico de músicas escutadas pelos usuários. Além disso, o site fornece recursos como: serviço de rádio online, recomendador de novidades, tabelas com detalhes do histórico de execução do usuário, informações sobre artistas, turnês e possibilidade de criação de fóruns, entre outros.

O Last.FM fornece uma API ¹(Application Programming Interface - conjunto de rotinas fornecidas por um software para que aplicativos acessem suas funcionalidades) que permite o acesso a dados presentes no site. É possível coletar informações dos usuários, histórico de escuta dos usuários e informações sobre as músicas / álbuns / artistas. Para nossos experi-

¹www.lastfm.com.br/api

mentos, nós coletamos 2 tipos de dados: o primeiro consiste num conjunto de usuários, junto com seu histórico de escuta, e o segundo em metadados dos artistas escutados. Os usuários do LAST.FM foram os sujeitos da pesquisa, e como dito anteriormente, representam os ouvintes. Já os artistas que o usuário nunca escutou anteriormente são as novidades, enquanto os que ele já escutou são as conhecidas.

4.2 Ouvinte

Com o intuito de estudar os artistas escutados pelos ouvintes, foram coletados dados acerca de um conjunto de usuários do LAST.FM. A coleta deste conjunto de usuários foi feita a partir do procedimento de *SnowBall Sampling* [1], iniciada pelo perfil do autor e sendo expandida pela coleta dos vizinhos musicais. Vizinho musical é um conceito utilizado no LAST.FM, onde duas pessoas são vizinhas se possuem gostos musicais parecidos. Foi coletado um conjunto de 100 mil usuários.

Após a coleta dos usuários, o próximo passo foi coletar o histórico de escuta dos mesmos. Para identificar as novidades, o histórico do usuário foi dividido em períodos, que serão detalhados na Subseção 4.2.1.

4.2.1 Timeline

Os dados referentes ao histórico de cada sujeito foram coletados no período entre a primeira vez que o usuário escutou alguma música no LAST.FM e Agosto de 2013. Este período foi dividido em duas partes, como pode-se ver na Figura 4.1: *Histórico Inicial* do sujeito e o *Período de Experimento*. O Histórico Inicial contempla o período desde a primeira música que o sujeito escutou no LAST.FM até Agosto de 2012, enquanto o Período de Experimento engloba o período entre Agosto de 2012 e Agosto de 2013 (um ano no total). Além dessa divisão, especificamos os seis primeiros meses do Período de Experimento como *Período de Observação*.

Com esta divisão, foram identificadas quais as novidades escutadas pelo usuário. Os artistas escutados pelo usuário no Período de Observação que não foram escutados no Histórico Inicial são consideradas novidades. Já os artistas que foram previamente escutados no Histórico Inicial são considerados como artistas conhecidos. Não consideramos o Período

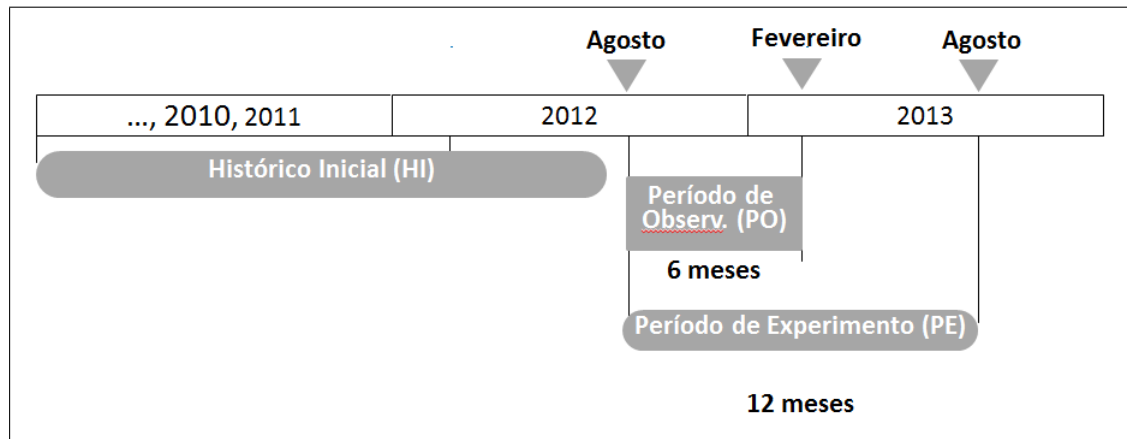


Figura 4.1: Linha do tempo utilizada no trabalho

de Experimento todo para evitar viés no cálculo das características das novidades. Uma novidade a que um sujeito foi exposto no começo do Período de Experimento tem maior probabilidade de ser escutada mais vezes a uma novidade que o sujeito foi exposto no final do Período de Experimento. Assim, identificamos como novidades os artistas escutados no Período de Observação, e levamos em conta as métricas referente a elas durante todo o Período de Experimento.

4.2.2 Histórico do usuário

Do Histórico Total do ouvinte, coletamos todos os artistas que ele escutou desde a entrada do usuário no LAST.FM até Agosto de 2013, junto com o total de execuções das músicas do artista. O método da API utilizado foi *getTopArtists*. Para o Período de Experimento, fizemos dois tipos de coleta. A primeira, utilizando o *getTopArtists* dos 12 meses, coletamos todos os artistas escutados, junto com o número de execuções das músicas de cada. O intuito dessa primeira coleta foi identificar os artistas do Histórico Inicial. A segunda parte consiste nos artistas escutados em cada semana deste período, junto com o número de execuções em cada semana. O método utilizado foi o *getWeeklyArtistChart*. Já esta segunda coleta foi realizada com o intuito de obter, além do número de execuções das músicas de cada artista, o número de semanas que o usuário escutou cada artista.

Com os dados do Histórico Total e do Período de Experimento conseguimos identificar o Histórico Inicial e as novidades. Os artistas do Histórico Inicial são os artistas que o ouvinte

não escutou apenas no Período de Experimento. Ou seja, artistas com número de execuções no Histórico Total do ouvinte maior que no Período do Experimento. Já as novidades são os artistas com o mesmo número de execuções no Histórico Total e no Período de Experimento.

Após a coleta e definição de cada período, foi realizado uma filtragem nos dados, que será descrita na Seção 4.2.3.

4.2.3 Filtros

Como o LAST.FM é uma rede social, formada por diferentes tipos de usuários, com diferentes hábitos musicais, foi preciso fazer uma filtragem nos sujeitos, para selecionar os adequados aos propósitos dos experimentos. Abaixo estão as características que os sujeitos precisavam ter para serem selecionados, junto com a maneira de filtragem utilizada.

1. Possuir alta atividade de escuta no período de Histórico Inicial.

Filtro: Exclusão de usuários que tenham escutado menos de 100 artistas no período de Histórico Inicial.

2. Possuam altos hábitos de escuta no Período de Experimento.

Filtro: Exclusão de usuários que escutaram menos de 100 músicas por semana em pelo menos 1/4 das semanas do Período de Experimento.

3. Foram expostos a um número de novidades que permitam a investigação de relações entre as características das novidades e as suas preferências

Filtro: Exclusão de usuários que possuam menos de 10 novidades.

4. Possuam números irrealistas de execuções: foi detectado que alguns usuários possuíam um número muito grande de execuções musicais. Alguns, por exemplo, tiveram uma média de mais de uma música por minuto, o que na realidade é impraticável. Uma explicação para esse fato seria a criação de robôs que trocassem a música assim que o sistema contabilizasse um *Scrobble*.

Filtro: Exclusão de usuários que tiveram uma média de execuções maior que 16 horas de execuções por dia, no Período de Experimento. Como as pessoas dormem em média 8 horas por dia, um ouvinte que passe o dia todo, enquanto acordado, escutando

música, escutaria 16 horas de música por dia. Suponto que uma música tem em média 4 minutos, foram excluídos os usuários que tiveram média maior que 240 músicas por dia ($\frac{16hrs \times 60min}{4min/musica} = 240musicas/dia$)

5. Não utilizem majoritariamente a rádio do LAST.FM: um dos objetivos é identificar as preferências dos usuários. Assim, é importante que a maior parte dos artistas escutados pelo usuário sejam escolhidos por ele, e não por uma rádio.

Exclusão de usuários que não escutaram nenhum artista mais de 15x na semana, em mais de 1/4 das semanas do período de observação.

O processo de filtragem resultou em uma amostra de 10.207 sujeitos.

4.3 Metadados dos Artistas

Com o intuito de calcular os aspectos das novidades, onde as novidades são artistas escutados pelo usuário, foram coletados dois tipos de metadados referentes aos artistas: as *tags* que descrevem o artista e a popularidade do artista no LAST.FM. O primeiro foi utilizado para calcular a familiaridade de uma novidade para um usuário, como também para construir um modelo de perfil do usuário e calcular a ecleticidade do mesmo. O segundo foi utilizado para representar a popularidade da novidade. Todos os metadados foram coletados da API do Last.fm no dia 01 de Setembro de 2013.

Tags são palavras (ou conjunto de palavras), como *rock*, *rap* e *pop*, associadas a um recurso, como músicas, álbuns e artistas. No Last.FM os usuários podem marcar cada um dos recursos com alguma tag, caracterizando-as *tags sociais*. Estas tags podem representar gêneros musicais (rock, samba), localização (brasil, nordeste, germany, west coast), humor (sad, chill, happy), opinião (love, favorite), referência pessoal (seen live, i own it), entre outros. Como as tags podem ser de vários tipos (não apenas gênero musical), elas podem ser consideradas *descritores* das músicas. Ao longo do texto, a palavra *descriptor* será utilizada como um termo que descreve uma música / artista.

Para cada artista foram coletadas as tags atribuídas a ele pelos usuários, junto com a popularidade de cada tag. Esta popularidade está relacionada à quantidade de vezes que a tag foi atribuída para o artista específico, pelos usuários do Last.FM. A popularidade da tag

fornecida pelo LAST.FM é normalizada, onde a tag mais atribuída possui valor igual a 100 e as outras tags possuem valores proporcionais, de acordo com a frequência de atribuição de cada uma. Formalmente, seja A o conjunto de artistas, e T o conjunto de tags. Seja $fa : A \times T \rightarrow R$ uma função que denote a frequência absoluta que uma tag $t \in T$ foi atribuída a um artista $a \in A$. O valor normalizado da tag t , representado pela função $f : A \times T \rightarrow R$ é representado pela equação 4.1.

$$f(t, a) = \frac{fa(t, a)}{\max_{x \in T}(fa(x, a))} \times 100 \quad (4.1)$$

A tabela 4.1 apresenta as 5 tags com maior valor do artista Michael Jackson. Pode-se ver que a tag *pop* foi a mais atribuída para Michael Jackson, possuindo valor 100. O método utilizado da API do Last.fm foi o *artist.gettoptags*.

Tag	Valor
pop	100
80s	49
dance	40
soul	35
funk	32

Tabela 4.1: Tags do artista Michael Jackson, junto com o valor normalizado de cada uma.

Como as tags são associadas pelos usuários do Last.FM, há problemas relacionadas a esse processo /* **REFERENCIA (Ima)** */. Usuários podem atribuir tags que não condizem com a realidade, podem errar na escrita da tag, etc. Para utilizar tags que realmente descrevam o artista, foi realizado um processo de filtragem. De cada artista foram consideradas no máximo quatro tags que tenham valor mínimo de 30 (onde 100 é a tag mais atribuída àquele artista). Além disso, foram eliminadas manualmente as tags com conotação pessoal, como *seen live* (vi ao vivo), *favorite* (favorito).

Sobre a popularidade do artista, foram coletados o número de usuários do LAST.FM que escutaram cada artista. Como dito anteriormente, esses valores foram coletado no dia 01 de Setembro de 2013. Como o valor da popularidade pode atingir ordens bem diferentes, com o intuito de estudar a magnitude de cada valor, foram utilizados nos cálculos o logaritmo na

base 10 do valor da popularidade. A tabela 4.2 mostra exemplos de popularidade de alguns artistas. O método utilizado da API foi o *artist.getinfo*.

Artista	Número de ouvintes (popularidade)
Michael Jackson	2.998.428
The Beatles	3.177.625
Red Hot Chili Peppers	4.032.453
Eminem	3.756.890
Chico Buarque	314.584

Tabela 4.2: Número de ouvintes (popularidade) de alguns artistas no LAST.FM

Capítulo 5

Modelos

Após a coleta e filtragem dos dados, utilizamos 3 conceitos para representar os artistas escutados pelos sujeitos:

1. Modelo de perfil do ouvinte: foi construindo o modelo do perfil do ouvinte, com o intuito de: gerar uma representação visual do que foi escutado pelo ouvinte; viabilizar o cálculo da familiaridade de um artista para um ouvinte; e gerar a métrica de ecleticidade, que foi utilizada na criação dos grupos de ouvintes (Capítulo 7).
2. Características das novidades: foram modeladas as características da novidade a serem utilizadas nos experimentos - familiaridade e popularidade.
3. Métricas de relevância: por fim, foram modeladas duas métricas que refletem a preferência do ouvinte para um artista, ou a relevância deste artista para o ouvinte, durante um período de tempo - a atenção total e o período de atenção.

5.1 Perfil musical do ouvinte

Um perfil de um usuário é um modelo, geralmente apresentado em forma de grafo ou tabela, que representa características de um determinado usuário sobre determinado tema. Assim, o perfil musical de um ouvinte é uma representação das músicas ou artistas que ele escuta.

A necessidade de uma representação do perfil musical do ouvinte surgiu primeiramente para calcular a familiaridade de um artista para o ouvinte. Como definida no Capítulo 3, a familiaridade de um artista para um ouvinte está relacionada com a semelhança dos descritores

do artista com os descritores dos artistas escutados previamente pelo ouvinte, considerando a influência individual destes descritores no histórico musical do ouvinte. Para realizar o cálculo, é necessário um modelo do perfil do ouvinte que sumarie os descritores das músicas por ele escutadas no decorrer do tempo, para comparar com os descritores da novidade escutada.

Uma forma de representação intuitiva do perfil de uma pessoa seria o conjunto de gêneros de artistas que essa pessoa escuta/escutou. Tipicamente, ao perguntar a ouvintes qual seu perfil musical, respostas como estas surgem: "Meu perfil é formado por artistas de Rock", "Meu perfil é formado por artista de Forró e artistas de Pagode". Uma representação formal desse conceito se adequaria à nossa necessidade, pois gêneros musicais também são considerados como descritores.

Assim, o nosso modelo de perfil de ouvinte é formado pelo conjunto de grupos de artistas mais representativos do histórico do ouvinte, onde cada grupo é formado por artistas semelhantes de acordo com os seus descritores. Com isso, um sujeito pode ter um perfil formado por um grupo de artistas de Rock, um grupo de artistas de Samba, um grupo de artistas de Forró, e assim por diante. Os artistas do histórico musical utilizados para a construção do perfil foram os que não são novidades e que possuísem número de execução no histórico do usuário maior que a média do número de execuções total dos artistas do histórico do usuário.

Para a construção deste conjunto de grupos, foi utilizado um algoritmo de agrupamento hierárquico aglomerativo **REFERENCIA (Ima)** **/*. Este tipo de algoritmo inicializa cada elemento (em nosso caso cada artista) em um grupo, e a cada passo, ele une os dois grupos mais próximos (similares). Desta maneira, é necessário definir uma medida de distância, ou dissimilaridade, entre os grupos.

Na maior parte dos métodos utilizados no agrupamento hierárquico aglomerativo, a medida de distância entre grupos pode ser gerada a partir de uma métrica de distância entre os pares de elementos e um critério de união que especifica quais grupos unir em cada passo, em função desta distância. Então, foi definido como métrica de distância entre pares de artistas o complemento da similaridade do cosseno entre os vetores de tags dos artistas e como critério de união (*linkage criterion*) o agrupamento de união pela média (o *average linkage clustering*).

Para calcular a similaridade do cosseno **REFERENCIA (Ima)** **/*, os artistas foram

representados por vetores, onde o vetor é formado pelas tags atribuídas a cada artista (e que passaram pelo processo de filtragem). Já o valor de cada tag é o valor normalizado do número de vezes que a tag foi atribuída ao artista. Tanto a filtragem quanto valor de cada tag foram descritos na Seção 4.3. Assim, seja A o conjunto de artistas, e T o conjunto de tags, seja $f : A \times T \rightarrow R$ a função que denote a frequência que uma tag $t \in T$ foi atribuída a um artista $a \in A$, então o vetor que representa um artista $a \in A$ é:

$$\vec{a} := (f(a, t_1), f(a, t_2), \dots, f(a, t_{|T|})) \quad (5.1)$$

A similaridade do cosseno entre dois artistas é definida pela equação 5.2. Como o algoritmo aglomerativo hierárquico requer uma medida de distância, e não de similaridade, foi calculado o complemento da similaridade do cosseno (5.3).

$$\cos(\vec{a}, \vec{a}') := \frac{\langle \vec{a}, \vec{a}' \rangle}{\|\vec{a}\| \|\vec{a}'\|} \quad (5.2)$$

$$\text{dis}(\vec{a}, \vec{a}') := 1 - \cos(\vec{a}, \vec{a}') \quad (5.3)$$

O *average linkage clustering* /* **REFERENCIA (lma)** */ é um método de união de grupos baseado na média das distâncias entre cada par de elementos de cada grupo. A distância entre dois grupos é definida pela equação 5.4. Seja X e Y grupos, $x \in X$ um artista do grupo X e $y \in Y$ um artista do grupo Y . A distância $d(X, Y)$ entre os grupos X e Y é definida pela média das distâncias de todos os pares $x \in X$ e $y \in Y$.

$$d(X, Y) := \frac{1}{|X||Y|} \sum_{x \in X} \sum_{y \in Y} \text{dis}(x, y) \quad (5.4)$$

Após a definição da distância entre grupos, o algoritmo de agrupamento foi aplicado. Como é um método aglomerativo hierárquico, o algoritmo inicia cada artista dentro de um grupo separado. Em cada etapa os grupos mais próximos vão sendo aglutinados, até chegar em 1 grupo com todos os artistas. Para selecionar o número de grupos de um ouvinte, o algoritmo foi interrompido no momento em que a distância mínima entre 2 grupos fosse igual a 0,30.

Os perfis obtidos tiveram média de 33,28 grupos (onde cada grupo possui pelo menos 2 artistas) e desvio padrão de 20,3. A figura /* **FIGURA (lma)** */ representa o histograma do

número de grupos. Um exemplo de perfil está representado na figura /* **FIGURA (lma)** */. Como comparação, a figura /* **FIGURA (lma)** */ possui o grupo de artistas utilizados no agrupamento.

A criação do perfil, além de auxiliar na visualização do gosto musical do usuário, evidenciado nas figuras /* **FIGURAS (lma)** */, e do cálculo da familiaridade (Seção 5.2.1), faz parte do cálculo da ecleticidade.

5.1.1 Ecleticidade

A ecleticidade representa o quão eclético musicalmente um ouvinte é - o quão diferente são os grupos de artistas que ele escuta. Ou seja, um ouvinte com alta ecleticidade é um que escuta muitos estilos diferentes de música. Esta métrica foi utilizada para conhecer melhor os hábitos dos ouvintes, e foi utilizada na geração dos grupos de ouvintes baseados nas preferências pelos aspectos das novidades comparadas com seus hábitos musicais, descritos no Capítulo 7.

Inicialmente foi pensado em utilizar o número de grupos do perfil do ouvinte como critério de ecleticidade. Quanto mais grupos o ouvinte possuir no perfil, mais eclético ele seria. Porém, dois ouvintes podem possuir o mesmo número de grupos mas um ouvinte pode possuir no perfil grupos mais similares (como um grupo de indie rock e outro de british rock) e outro possuir menos similares (como um grupo de forró e outro de indie rock).

Uma alternativa a essa abordagem seria contabilizar o quanto de diferença ou diversidade cada grupo adiciona ao perfil. Quanto mais diversidade cada grupo proporcionar, mais eclético o ouvinte é. Para isso, recomendamos novamente a um algoritmo de agrupamento hierárquico, porém agora utilizando os grupos do perfil. A cada passo calculamos e armazenamos a distância entre os dois grupos que foram unidos. Por fim, definimos a ecleticidade como a soma de todas estas distâncias. Formalmente, seja $P := \{C_1, \dots, C_n\}$ o perfil do ouvinte, formado pelos grupos do perfil C_i . Seja $X^{(p)} := \{X_1^{(p)}, \dots, X_n^{(p)}\}$ o conjunto de grupos no passo p do algoritmo hierárquico, onde $X^{(1)} = P$ e $X^{(j)}$, onde $j > 1$, grupos do algoritmo hierárquico criados a partir do conjunto inicial $X^{(1)}$. Além disso, seja $d(X_k^{(p)}, X_l^{(p)})$ a distância entre os grupos $X_k^{(p)} \in X^{(p)}$ e $X_l^{(p)} \in X^{(p)}$. Então, a ecleticidade do ouvinte com perfil P é calculada como: 5.5.

$$e(P) = \sum_{i=1}^{|x^{(p)}|} \min(d(X_k^{(p)}, X_l^{(p)})) \quad (5.5)$$

A distância $d(X_k^{(p)}, X_l^{(p)})$ entre os grupos do algoritmo hierárquico foi calculada utilizando o *average linkage method* (equação 5.4). Como o *average linkage method* depende da distância entre cada par de elemento (onde cada elemento é um grupo $C_i \in P$), foi definido como distância entre dois grupos $C_i \in P$ e $C_j \in P$ o complemento da similaridade do cosseno (equação 5.3) entre os vetores c_i e c_j que representam os centróides de C_i e C_j .

5.2 Características das novidades

Para caracterizar as novidades multidimensionalmente, utilizamos dois aspectos: a familiaridade e a popularidade. Esta seção descreve o cálculo de cada aspecto.

5.2.1 Familiaridade

Como definido no Capítulo 3, a familiaridade de um artista a para um ouvinte o reflete o quanto este ouvinte foi exposto a outros artistas que têm descritores semelhantes aos do artista a . A familiaridade está relacionada com a similaridade entre os descritores do artista e os descritores dos artistas do perfil do ouvinte, considerando com a influência individual desses artistas.

Formalmente, seja $P := \{C_1, \dots, C_n\}$ o perfil do ouvinte, formado pelos grupos de artistas C_i . Seja \vec{c}_i o centróide do grupo C_i , e p_i a influência do grupo C_i no perfil do ouvinte, definida como a proporção de todas as execuções de músicas pelo ouvinte que são as músicas cujos artistas estão em C_i . Assim, a familiaridade entre um artista a e o perfil do ouvinte P é o valor máximo da similaridade entre o artista a e algum grupo C_i multiplicado pela influência p_i de C_i :

$$fam(a, P) = \max(\cos(\vec{a}, \vec{c}_i) \times p_i) \quad (5.6)$$

A influência p_i de C_i é representado pela média do total de execuções das músicas dos artistas presentes no grupo. Quanto mais vezes os artistas do grupo i foram escutados, mais influentes os descritores deste grupo são para o ouvinte.

5.2.2 Popularidade

O segundo aspecto da novidade estudado foi a popularidade. Para calcular a popularidade, utilizamos o logaritmo na base 10 da número de ouvintes do artista no LAST.FM. O logaritmo foi utilizado pois a distribuição da popularidade dos artistas é enviesada (5.1).



Figura 5.1: Curva de popularidade dos artistas do LAST.FM. A curva é enviesada.

5.3 Preferências

Para mensurar o quanto o ouvinte preferiu a novidade, foram utilizados duas métricas: a atenção total e o período de atenção. Como novidades podem ser descobertas em todo o Período de Experimento, alguns destes artistas possuem uma janela de tempo no experimento menor (artistas escutadas no final do Período de Experimento). Para contornar esse problema, utilizamos duas soluções. Primeiro, utilizamos como denominador no cálculo das métricas o número de semanas da Janela de Tempo de exposição à novidade, que vai da primeira semana que foi escutada a novidade até o fim do Período de Experimento. Segundo, como mencionado na Seção 4.2.1, apenas as novidades descobertas no Período de Observação foram consideradas na análise, mas todo o Período de Experimento foi utilizado para cálculo das métricas. Isso dá a cada novidade um mínimo de 6 meses de coleta de dados, que limita um possível viés para Janelas de Tempo pequenas.

A atenção total representa a atenção que o ouvinte deu para o artista no Período de Experimento. A atenção total do ouvinte para o artista é representada pelo total de número de

execuções de músicas do artista que ele escutou no Período de Experimento, dividido pelo número de semanas de sua Janela de Tempo.

Já o período de atenção o tempo que o ouvinte deu atenção ao artista. Assim, é o número de semanas que o ouvinte escutou o artista dividido pelo número de semanas de sua Janela de Tempo.

Capítulo 6

Preferências dos ouvintes para diferentes aspectos de novidades

Após a coleta e filtragem dos dados, e da modelagem dos conceitos, partimos para responder as perguntas de pesquisa. Este capítulo abrange as duas primeiras perguntas, abordando as preferências dos ouvintes pelos aspectos - familiaridade e popularidade - da novidade, de forma geral e de forma individual.

6.1 Preferências gerais

Este trabalho visa entender como diferentes aspectos de novidades podem influenciar as preferências de ouvintes por elas. Para esta pesquisa, a primeira pergunta levantada foi: *Há uma correlação geral entre algum aspecto da novidade e as preferências dos ouvintes?* Para responder essa pergunta, calculamos a correlação entre cada aspecto da novidade - familiaridade e popularidade - e cada métrica de preferência - atenção total e período de atenção, de todas as novidades de todos os ouvintes juntas.

Com todas métricas das novidades em mãos, foi utilizado o método de correlação não-paramétrico de Spearman. O resultado gerado por este método pode variar de -1 a 1. Quanto mais próximo de 1, mais as variáveis estão correlacionadas positivamente - se uma cresce/decresce a outra cresce/decresce. Quanto mais próximo de -1, mais as variáveis estão correlacionadas negativamente - se uma cresce a outra decresce, e vice-versa. Se o valor estiver próximo a 0 não há correlação.

Dimensões	F	P
Período de atenção (PdA)	0,08	0,07
Familiarity (F)	-	0,08

Tabela 6.1: Correlação (Coeficiente de Spearman) entre aspectos da novidade e preferências, analisando todas as novidades juntas

A tabela 6.1 mostra o coeficiente de Spearman para cada par de aspecto da novidade / preferência. Como pode-se ver, todos os valores encontrados da correlação são próximos de zero. Podemos concluir que no geral não existe uma correlação entre aspectos e preferências pela novidade, para todos os ouvintes juntos. Por exemplo, os ouvintes em geral não preferem novidades familiares, ou no geral não preferem novidades não-familiares. Nós levantamos duas hipóteses para explicar esse resultado:

1. Diferentes ouvintes possuem diferentes preferências para os aspectos das novidades. Nesta hipótese, cogitamos que diferentes ouvintes possuem diferentes preferências musicais. Assim, existem ouvintes que preferem novidades populares, outros preferem novidades não populares, etc. Colocando todos estes ouvintes juntos, a correlação geral vai ser próxima a zero.
2. Individualmente, os ouvintes não preferem um aspecto a outro das novidades. Cada ouvinte pode preferir, por exemplo, tanto novidades familiares quanto não familiares, fazendo com que a correlação entre a preferência e o aspecto das novidades seja próxima a zero.

6.2 Preferências individuais

O resultado da primeira pergunta de pesquisa e estas hipóteses nos levam à segunda pergunta de pesquisa: individualmente, os ouvintes preferem algum aspecto de novidade? Para saber se os ouvintes possuem alguma correlação entre os aspectos e as preferências das novidades, calculamos as correlações para cada sujeito individualmente. Cerca de 74% dos sujeitos possuem alguma correlação com valor maior que 0,15 ou menor que -0,15, e cerca de 26%

possuem alguma correlação maior que 0,3 ou menor que -0,3. Desta maneira, individualmente, boa parte dos ouvintes possuem alguma correlação entre algum aspecto e alguma preferência da novidade. A figura /* **FIGURA (lma)** */ mostra a distribuição acumulada dos valores das correlações para cada par aspecto / preferência.

Analisando estes dois resultados juntos, pode-se concluir que, apesar de não existir uma tendência geral quanto às preferências dos sujeitos para os diferentes aspectos das novidades, a maior parte dos sujeitos possuem alguma preferência para algum aspecto de novidade no seu comportamento. Isso sugere a presença de diferentes tipos de ouvintes nos nossos dados. Para identificar estes tipos, foi utilizado um algoritmo de agrupamento nos sujeitos, que será discutido no próximo capítulo.

Capítulo 7

Grupos de usuários para diferentes aspectos de novidade

De acordo com os resultados do Capítulo 6, há a evidência de diferentes tipos de sujeitos nos nossos dados, de acordo com as preferências pelos aspectos das novidades. Isto nos leva a terceira pergunta de pesquisa: Existem grupos de ouvintes relacionados com as preferências pelos aspectos das novidades baseadas no perfil? Para identificar estes grupos, foi realizada um agrupamento nos dados que caracterizam os sujeitos. Esta análise será descrita neste capítulo.

7.1 Conjunto de sujeitos

Primeiramente, a análise foi realizada com os dados dos mesmos sujeitos utilizados na análise das preferências dos ouvintes para os aspectos das novidades (Capítulo 6). Depois, para estudar especificamente os sujeitos com alguma correlação entre preferência e aspectos da novidade, foi feita a análise com dois subconjuntos: o primeiro subconjunto consiste em sujeitos com valor de correlação entre algum aspecto e preferência da novidade maior que 0,15 ou menor que -0,15; o segundo consiste em sujeitos com valor de correlação maior que 0,3 ou menor que -0,3. Os resultados encontrados nas três análises foram similares. Portanto, mostraremos apenas os resultados da primeira análise, enquanto os outros resultados estão no apêndice /* X (lma) */.

7.2 Dados que caracterizam sujeitos

O objetivo da análise atual é identificar os grupos de ouvintes baseados nas preferências pelos aspectos das novidades comparadas com seus hábitos musicais. Foram escolhidas 5 métricas de caracterização dos sujeitos, que podem ser divididas em dois grupos: métricas relacionadas com novidades e métricas relacionadas com os hábitos musicais.

1. **Relacionadas com novidades:** Métricas que caracterizam os sujeitos a partir de suas preferências pelos aspectos das novidades.

- *Correlação entre a familiaridade das novidades e a atenção total dedicada a elas durante o experimento.*
- *Correlação entre popularidade das novidades e a atenção total dedicada a elas durante o experimento.*

**Calculamos a correlação (coeficiente de Spearman) entre a atenção total e o período de atenção, e o valor encontrado foi de 0,71. Como ambas as variáveis estão correlacionadas, decidimos utilizar no algoritmo de agrupamento apenas as correlações que envolvem a atenção total.*

2. **Relacionadas com os hábitos musicais:** Métricas que caracterizam os hábitos musicais do sujeito e que estão relacionadas com novidades.

- *Ecleticidade* A ecleticidade representa o quão diferente os artistas do perfil do ouvinte são, de acordo com seus descritores. Sua definição e cálculo foi feito na Subseção 5.1.1. A ecleticidade está relacionada com a familiaridade da novidade, pois quanto mais eclético um sujeito for, maior a probabilidade dele ser familiar a vários tipos de novidades.
- *Popularidade média dos artistas do perfil do ouvinte* Esta métrica está relacionada com a popularidade das novidades escutadas pelo ouvinte.
- *Proporção de novidades escutadas pelo ouvinte no Período de Observação* Com esta métrica pode-se identificar se o ouvinte possuiu o hábito de escutar muitas ou poucas novidades, no Período de Observação.

7.3 Algoritmo de agrupamento

Para fazer o agrupamento dos sujeitos, foi utilizado o algoritmo de agrupamento aglomerativo hierárquico. Como descrito na Seção 5.1, este tipo de algoritmo necessita de uma métrica de dissimilaridade entre os pares de sujeitos e um critério de união que especifica quais grupos unir em cada passo.

Para calcular a dissimilaridade entre os sujeitos, primeiro foram calculadas as métricas descritas na Seção 7.2. Após o cálculo, estes dados foram normalizados, baseados no Z-Score. Então, a dissimilaridade entre dois sujeitos foi calculada a partir da distância euclidiana, onde cada sujeito é representado por um vetor contendo as 5 métricas normalizadas. Formalmente, seja S o conjunto de vetores com dados normalizados que representam cada sujeito; e $x \in S$ e $y \in S$ vetores de dois sujeitos de S . x_i representa a posição i do vetor x , no caso uma das 5 métricas que caracterizam os sujeitos. A dissimilaridade entre os sujeitos, representada pela dissimilaridade $disSuj(x, y)$ dos vetores que os representam é dado pela equação 7.1. Como critério de união, foi utilizado o método Ward /* **REFERENCIA (Ima)** */.

$$disSuj(x, y) := \sqrt{\sum_{i=1}^5 (x_i - y_i)^2}$$

(7.1)

7.4 Escolha do número de grupos

O método de agrupamento hierárquico não expõe explicitamente o número de grupos resultantes. A cada passo, o algoritmo une dois grupos, até que todos os sujeitos estejam em um só grupo. Uma abordagem para definir a melhor configuração de número de grupos é o método do joelho /* **REF (Ima)** */, ao plotar um gráfico onde o eixo X é o número de grupos e o Y um critério de avaliação. A figura 7.1 mostra a distância média dentro dos grupos para cada configuração de número de grupos. O método de joelho determina escolher uma configuração de grupos que não adicionem muita heterogeneidade, evidenciado a partir

da curvatura máxima do gráfico (joelho). Pela figura, a distância média dentro dos grupos começa a aumentar vertiginosamente nas configurações com menos de 7 grupos. Assim, analisamos as configuração com 6, 7 e 8 grupos.

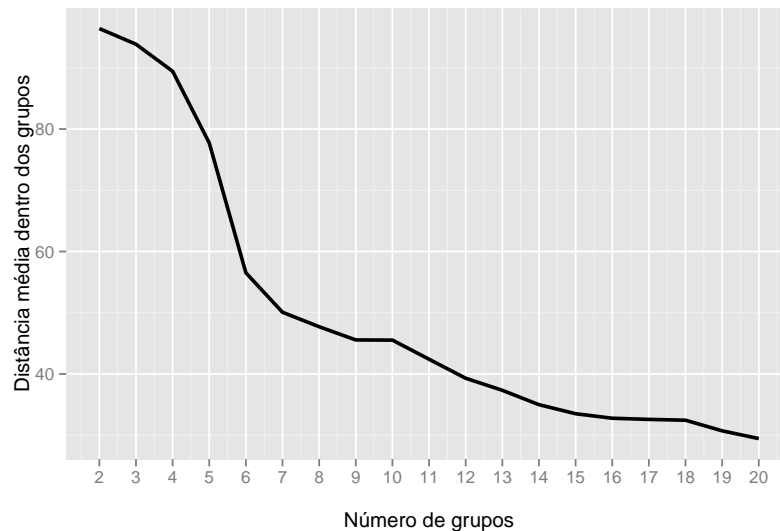


Figura 7.1: Número de grupo X Distância média dentro dos grupos. O joelho do gráfico está em torno da configuração com 7 grupos.

Comparando os centróides dos grupos para as configurações 6, 7 e 8, os centróides dos grupos 7 e 8 são semelhantes. /* **figura Z (lma)** */ Desta maneira, escolhemos a solução com 7 grupos a mais representativa para nossos estudos.

7.5 Grupos

Após a escolha de 7 grupos de ouvintes, utilizamos o centróide de cada grupo para analisar as suas principais características. A figura 7.2 representa os valores dos centróides normalizados. Podemos dividir os grupos em dois tipos: o primeiro, onde as características que se destacam são as relacionadas com as preferências pelos aspectos das novidades e o segundo, onde as características que se destacam são as relacionadas com os hábitos musicais dos ouvintes. De acordo com cada centróide, rotulamos os grupos da seguinte forma:

1. Grupos marcados pelas preferências pelos aspectos das novidades

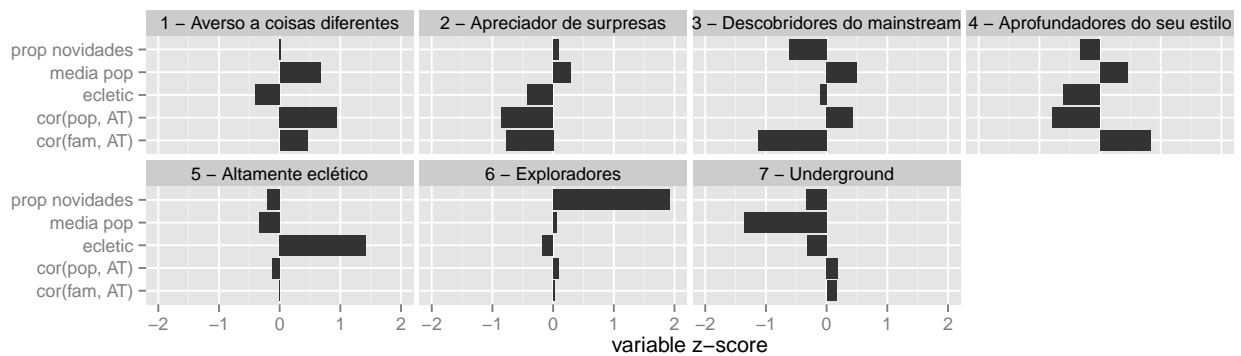


Figura 7.2: Centróides dos 7 grupos encontrados na análise. As métricas estão normalizadas pelo z-score, onde zero representa a média de todos os ouvintes, e a unidade de variação é um desvio padrão, para cada métrica. No eixo vertical, *fam* significa a familiaridade, *pop* significa popularidade, *AT* significa atenção total.

- (a) Averso a coisas diferentes (ou acomodado) [total de ouvintes: 2317 (20%)]:
Maior grupo com característica marcante pelas preferências pelos aspectos das novidades, formado por ouvintes que preferem novidades familiares e populares, além de possuírem hábitos musicais marcados por artistas populares.
 - (b) Apreciador de surpresas [total de ouvintes: 1859 (17%)]: Ouvintes preferem novidades não-familiares e pouco populares.
 - (c) Descobridores do mainstream [total de ouvintes: 1022 (8%)]: Ouvintes que preferem novidades não-familiares e populares.
 - (d) Aprofundadores do seu estilo [total de ouvintes: 1467 (14%)]: Ouvintes que preferem novidades pouco populares e familiares, além de possuir pouca ecleticidade.
2. Grupos marcados pelas características dos hábitos musicais
- (a) Altamente eclético (2) [total de ouvintes: 2456 (21%)]: Maior grupo de todos, formado por ouvintes que possuem alta ecleticidade
 - (b) Exploradores (3) [total de ouvintes: 1047 (9%)]: Ouvintes que possuem alta proporção de novidades escutadas durante o Período de Observação
 - (c) Underground (4) [total de ouvintes: 1281 (11%)]: Ouvintes com hábito musical

marcado por artistas pouco populares.

7.6 Discussão dos grupos encontrados

Encontramos 4 grupos que possuem características principais pelas preferências pelos aspectos das novidades. Coincidentemente, foram encontrados grupos com todas as combinações possíveis de preferências pelos aspectos.

O maior destes 4 grupos é o que chamamos *Averso a coisas diferentes*. É um grupo de ouvintes que preferem novidades familiares e populares, além de possuírem hábitos musicais marcados por artistas populares, pouca ecleticidade e proporção mediana de novidades escutadas. É um tipo de ouvinte que não procura expandir seu perfil musical, preferindo escutar o que está na mídia do que ele habitualmente já escuta.

Os outros 3 grupos dos marcados pelas preferências pelos aspectos das novidades possuem uma distribuição mais homogênea do número de ouvintes. Opostos aos *Aversos a coisas diferentes*, os *Apreciadores de surpresas* preferem novidades não familiares e não populares, além de possuir pouca ecleticidade. Desta maneira, estes ouvintes normalmente escutam artistas bem parecidos, mas tentam aumentar esse leque de artistas do perfil preferindo novidades não familiares e não populares. Eles preferem surpresas, artistas diferentes do que já escutaram.

Os *Descobridores do mainstream* preferem novidades populares, que estão na mídia, mesmo não sendo familiares. Uma possível explicação seria ouvintes que escutam o que está nas paradas das rádios, sem se importar se são parecidos com o que ele escutava antes ou não. Já os *Aprofundadores do seu estilo* preferem novidades familiares e pouco populares, além de possuírem pouca ecleticidade. Esse tipo de ouvinte são fechados no seu nicho de estilos musicais, e ou não preferem o que está na mídia destes estilos, ou já escutaram tudo que está na mídia e agora querem expandir para artistas não populares destes estilos.

Analisando os 3 grupos restantes, o *Altamente eclético* é o grupo de ouvintes com maior ecleticidade, comparando com os demais. Nota-se que neste grupo existem diferentes sujeitos, por que as outras métricas são próximas a zero. O grupo de *Exploradores*, formado por ouvintes com alta proporção de novidades escutadas no período de observação, também possui essa característica do *Altamente eclético*. Por fim, o grupo *Underground* é formado

por ouvintes que têm hábito de escutar artistas pouco populares no geral.

7.7 Discussão dos resultados

Capítulo 8

Comparação das novidades com os artistas conhecidos

Até este momento trabalhamos apenas com novidades. Porém, será que as preferências dos ouvintes pelas novidades são as mesmas que pelos artistas conhecidos? Caso esta hipótese seja verdade, mostraríamos que o comportamento dos ouvintes para estes dois âmbitos, novidades e não-novidades, seriam similares, podendo estender os resultados das novidades para os itens conhecidos e vice-versa. Para responder esta pergunta expandimos nossos experimentos para englobar também os artistas conhecidos, tentando responder a seguinte pergunta de pesquisa: *As relações entre as preferências dos ouvintes e os aspectos das novidades são as mesmas que as relações entre as preferências dos ouvintes e os aspectos dos artistas já conhecidos?*

8.1 Características dos artistas conhecidos

Como dito na Seção 4.2.1, os artistas conhecidos são os artistas escutados no Período de Observação que já foram escutados previamente no Histórico Inicial do ouvinte. Para eles foram calculadas 3 das 4 métricas que foram calculadas para os artistas com novidade (Capítulo 5): familiaridade e popularidade (aspectos) e total de atenção (preferência).

Para calcular as relações entre os aspectos dos artistas conhecidos e as preferências pelos ouvintes, utilizamos a mesma metodologia descrita no Capítulo 6. Utilizamos o método de correlação não-paramétrico de Spearman entre cada par de aspecto / preferência.

8.2 Comparação entre relações das preferências e aspectos das novidades e dos artistas conhecidos

Para guiar nossos experimentos, decidimos especificar a primeira pergunta de pesquisa, que é mais genérica. Assim, tentamos responder duas perguntas de pesquisa:

1. *As relações entre as preferências dos ouvintes e os aspectos das novidades são significativamente diferentes das relações entre as preferências dos ouvintes e os aspectos dos artistas já conhecidas?* Se as relações comparadas forem significativamente iguais, poderíamos estender os resultados das novidades para o comportamento geral do ouvinte. Se forem diferentes, veremos a importância da análise separada de itens com novidades e itens conhecidos.

Para cada um dos 10.207 sujeitos descritos na Seção 4.2.3 foram calculadas as correlações descritas no Capítulo 6 e na Seção 8.1. Desta maneira cada ouvinte possui pares de correlações (Tabela 8.1), onde um valor do par é calculado utilizando artistas com novidade e o outro os artistas conhecidos. Por exemplo, cada ouvinte possui um valor de correlação entre a familiaridade e o total de atenção para as novidades e um valor de correlação entre a familiaridade e o total de atenção para os artistas conhecidos. Desta maneira, para responder a pergunta exposta no parágrafo anterior, utilizamos um teste-T pareado, para cada par de correlação novidade / artista conhecido correspondente de todos os ouvintes. O intuito desta análise seria descobrir se, por exemplo, o valor da correlação entre a familiaridade e o total de atenção para as novidades e para os artistas conhecidos são significativamente diferentes.

Artistas com novidade	Artistas conhecidos
cor(familiaridade, atenção total)	cor(familiaridade, atenção total)
cor(popularidade, atenção total)	cor(popularidade, atenção total)

Tabela 8.1: Correlações calculadas para artistas com novidade e artistas conhecidos.

A tabela 8.2 mostra o p-valor e o intervalo de confiança da média das diferenças para cada par de aspecto / preferência. Em ambas métricas comparadas, a probabilidade da média das diferenças entre os valores para os artistas com novidade e os artistas

conhecidos ser igual a zero é baixíssima ($p\text{-valor} < 2,2e-16$). Desta maneira, podemos afirmar que as correlações no geral são diferentes. Ou seja, as relações dos ouvintes entre as preferências e os aspectos, para artistas com novidade e artistas conhecidos, são diferentes.

Métricas	P-Valor	Intervalo de Confiança
Familiaridade / Período de atenção	$< 2,2e-16$	[-0,14 -0,13]
Popularidade / Período de atenção	$< 2,2e-16$	[-0,04 -0,03]

Tabela 8.2: Teste-T pareado entre as correlações dos aspectos e preferências dos ouvintes para artistas com novidade e artistas conhecidos.

Outro fato que podemos observar da tabela é que o intervalo de confiança, com 95% de confiança, para as diferenças das correlações entre familiaridade e atenção total para artistas com novidade e artistas conhecidos é [-0,14 -0,13], sugerindo que as preferências pelos artistas conhecidos familiares são maiores que as preferências pelos artistas com novidade familiares. Uma possível explicação seria que os artistas conhecidos mais familiares são os artistas que o ouvinte mais escutou no seu histórico. Assim, naturalmente a preferência por eles é bem maior.

Já o intervalo de confiança, com 95% de confiança, para as diferenças das correlações entre popularidade e atenção total para artistas com novidade e artistas conhecidos é [-0,04 -0,03]. Também sugere que as preferências pelos artistas mais populares conhecidos são maiores que pelos artistas mais populares com novidade, apesar de que este intervalo de confiança está mais próximo de zero que o intervalo de confiança anterior.

No geral, podemos concluir que as preferências dos ouvintes por artistas mais familiares e mais populares é menor se os artistas forem novidades do que se eles forem artistas conhecidos. É como se a preferência por artistas conhecidos já esteja consolidada e que ao escutar novidades os ouvintes tendem a explorar mais, escutando no mesmo período de tempo mais artistas menos familiares e menos populares que os artistas conhecidos.

2. Existe correlação entre as relações das preferências pelos aspectos dos artistas com novidade e as relações das preferências pelos aspectos dos artistas conhecidos? Pró-

ximo passo foi verificar se existiam correlações entre as relações calculadas para os artistas com novidade e para os artistas conhecidos. Será que quanto maior a correlação de um ouvinte entre familiaridade e atenção total para artistas com novidade, por exemplo, maior a correlação entre a familiaridade e atenção total para artistas conhecidos? Para responder a pergunta de pesquisa, utilizamos o método de correlação não-paramétrico de Kendall.

Novidades X Artistas conhecidos	cor(fam., aten. tot.)	cor(pop. , aten. tot.)
cor(familiaridade , atenção total)	0,05	-0.02
cor(popularidade , atenção total)	-0.04	0,09

Tabela 8.3: Correlação (Coeficiente de Kendall) entre correlações calculadas para novidades (linhas) e artistas conhecidos (colunas)

Como podemos observar na Tabela 8.3, a correlação é baixa em todos os casos. Não podemos afirmar, por exemplo, que dado dois ouvintes, se um possuir maior correlação entre familiaridade e período de atenção para novidades que o outro ouvinte ele também terá maior correlação entre familiaridade e período de atenção para artistas conhecidos, e vice-versa.

Esta falta de correlação geral pode ser ocasionada pela diversidade de comportamentos para novidades e artistas conhecidos dos ouvintes. Levantamos a hipótese de que, por exemplo, existem ouvintes que possuem uma preferência mais forte por novidades familiares que por artistas conhecidos familiares, enquanto outros não. Hipótese semelhante a esta foi exposta no Capítulo 6, que nos levou a criar perfis de ouvintes baseados nas relações entre preferências e os aspectos das novidades (Capítulo 7). Desta maneira, fizemos uma investigação parecida para testar a hipótese.

8.3 Grupos de ouvintes baseados na diferença das relações entre preferências e aspectos das novidades e dos artistas conhecidos

Como as correlações encontradas na segunda pergunta de pesquisa da seção anterior foram baixas, seguindo a metodologia dos Capítulos 6 / 7 tentamos encontrar perfis de ouvintes, com o intuito de responder a seguinte pergunta: *Existem diferentes perfis de ouvintes baseados nas diferenças das relações preferências / aspectos das novidades e dos itens conhecidos?*.

Primeiro passo foi calcular duas diferenças, para cada ouvinte: primeiro, a diferença entre as correlações da familiaridade e atenção total das novidades e dos artistas conhecidos (linha 1 da Tabela 8.1) e segundo, a diferença entre as correlações da popularidade e atenção total das novidades e dos artistas conhecidos (linha 2 da Tabela 8.1) .

Com estas diferenças em mãos utilizamos o algoritmo de agrupamento aglomerativo hierárquico. A escolha do número de grupos foi baseada no método do joelho. Na imagem 8.1 podemos observar um início de joelho entre as configurações com 4 e 6 grupos. Podemos notar pela tabela /* **FALTA (Ima)** */ que as configurações de 5 e 6 grupos possuem uma mudança maior que a mudança entre as configurações 6 e 7. Assim, como adicionar um grupo na configuração 6 não muda muito, escolhemos a configuração de 6 grupos.

Após a escolha dos 6 grupos de ouvintes, utilizamos o centróide de cada grupo para analisar as suas principais características. A figura 8.2 representa os valores dos centróides.

Podemos dividir os grupos em quatro tipos:

1. Tipo I - Preferências similares para novidades e itens conhecidos: Composto pelo grupo 1, o centroide deste evidencia que as diferenças entre as correlações, dos dois pares de preferência / aspecto, calculadas para novidades e itens conhecidos, são próximas a zero.
2. Tipo II - Correlações para itens conhecidos maiores que para novidades: O segundo tipo é composto pelos grupos 2 e 3, onde as correlações, dos dois pares de preferência / aspecto, são maiores para os itens conhecidos que para as novidades, seja em maior (grupo 3) ou menor (grupo 6) grau.

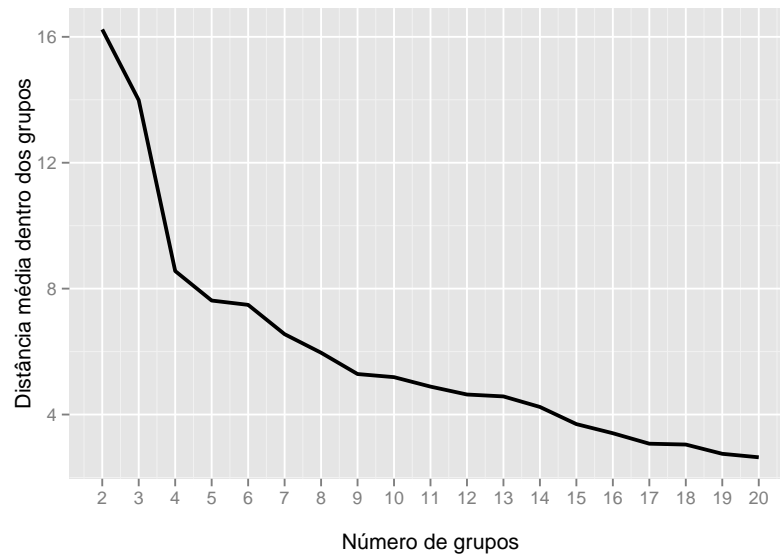


Figura 8.1: Número de grupo X Distância média dentro dos grupos.

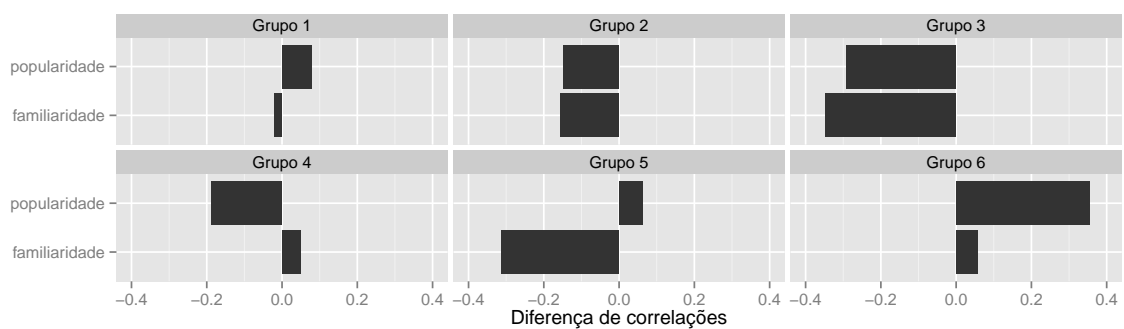


Figura 8.2: Centróides dos 6 grupos encontrados na análise da diferença das correlações

3. Tipo III - Correlação para itens conhecidos maior que para novidades em apenas um aspecto: composto pelos grupos 4 e 5, apenas um par preferência / aspecto possui correlação para itens conhecidos maior que para novidades, enquanto que o outro par é próximo a zero. O grupo 4 é formado no geral por ouvintes que possuem correlações entre preferência e familiaridade para itens conhecidos maiores que para novidades. Já o grupo 5 é formado no geral por ouvintes que possuem correlações entre preferência e popularidade para itens conhecidos maiores que para novidades.
4. Tipo IV - Correlação das preferências e popularidade para novidades maior que para itens conhecidos: por último, este tipo de ouvinte possui correlação entre preferência e popularidade maior para novidades que para artistas conhecidos. Composto pelo grupo 6.

Boa parte dos grupos (4 dos 6) são formados no geral por ouvintes que possuem pelo menos correlação de algum par preferência / aspecto maior para itens conhecidos que para novidades. Isso corrobora os resultados encontrados na Seção 8.2, que mostra que no geral os ouvintes possuem correlações maiores para itens conhecidos que para novidades. Cerca de 23% dos ouvintes estão no grupo de comportamento parecido para novidades e itens conhecidos. Assim, apenas 1/4 dos ouvintes poderíamos tentar prever seu comportamento para novidades baseados no comportamento de itens conhecidos, enquanto os demais ouvintes precisam de um tratamento diferenciado para estes dois tipos de comportamentos.

Capítulo 9

Conclusão

Capítulo 10

Trabalhos Futuros

Bibliografia

- [1] Leo A. Goodman. Snowball sampling. *The Annals of Mathematical Statistics*, 32(1):148–170, 03 1961.
- [2] S. R. Maddi. *Theories of cognitive consistency*, chapter The pursuit of consistency and variety. Rand McNally, Chicago, 1968.

Apêndice A

Apêndices