

# Revisão do Projeto

**Disciplina: Trabalho de Conclusão de Curso I – BCC**

Caro orientando,

segue abaixo a tabela de cálculo da média das notas obtidas no Pré-Projeto e Projeto, as DUAS revisões do seu projeto contendo a avaliação do professor “avaliador” e professor “TCC1”. Lembro que os ajustes indicados nestas revisões não precisam ser feitos no projeto, mas sim quando levarem o conteúdo do projeto para o artigo (se for o caso). Este material contendo todo o histórico das revisões é encaminhado para o professor de TCC2.

Atenciosamente,

| Nome                | PreProjeto |   |   |      |           |    |   |      |      |      | Ori.  | Banca | Projeto   |      |   |   |    |       |    |   |   |    | Média |     |
|---------------------|------------|---|---|------|-----------|----|---|------|------|------|-------|-------|-----------|------|---|---|----|-------|----|---|---|----|-------|-----|
|                     | TCC1       |   |   |      | Avaliador |    |   |      | TCC1 |      |       |       | Avaliador |      |   |   |    |       |    |   |   |    |       |     |
|                     | A          | P | N | Nota | A         | P  | N | Nota | A    | P    |       |       | N         | Nota | A | P | N  | Nota  |    |   |   |    |       |     |
| HenriqueJoseWilbert | 15         | 1 | 0 | 16   | 9,79      | 11 | 4 | 0    | 15   | 9,11 | 10,00 | 10,00 | 10,00     | 16   | 0 | 0 | 16 | 10,00 | 14 | 1 | 0 | 15 | 9,78  | 9,7 |

|                                      |               |                      |
|--------------------------------------|---------------|----------------------|
| CURSO DE CIÊNCIA DA COMPUTAÇÃO – TCC |               |                      |
| ( ) PRÉ-PROJETO                      | ( X ) PROJETO | ANO/SEMESTRE: 2021/2 |

## UTILIZAÇÃO DE CLUSTERIZAÇÃO PARA AUXÍLIO EM TOMADA DE DECISÃO A PARTIR DE DADOS DE VAREJO

Henrique José Wilbert

Prof. Aurélio Faustino Hoppe – Orientador

Prof. Christian Daniel Falaster – Coorientador

### 1 INTRODUÇÃO

Com a evolução da tecnologia de informação a partir dos anos 80 e início dos anos 90, várias grandes empresas adotaram sistemas de gerenciamento na forma de softwares Enterprise Resource Planning (ERP) (RASHID; HOSSAIN; PATRICK, 2001, p.2). Estes softwares auxiliam em suas rotinas à nível operacional, seja no controle do estoque, fiscal, financeiro, transacional e até recursos humanos. A partir disso, alcançou-se um patamar de eficiência nunca concebido, visto que registros antes realizados em papel e caneta, passaram a ser produzidos automaticamente. Ainda segundo os autores, em paralelo a informatização desses processos, houve também um crescimento da quantidade de dados armazenados referentes à produtos, clientes, transações, gastos e receitas.

Diante deste contexto, avançaram-se também as táticas de marketing direto, como por exemplo, o envio de catálogos por correio, até ofertas altamente objetivas para indivíduos selecionados, cujas informações transacionais estavam presentes na base de dados. Percebeu-se então que o foco das relações empresa-cliente não está em clientes novos, e sim em clientes já existentes nas bases de dados, visto que o custo para adquirir um cliente novo através de publicidade é muito maior que o custo de alimentar uma relação já existente (PETRISON; BLATTBERG; WANG, 1997, p. 119).

Segundo Reinartz, Thomas e Kumar (2005, p.77), quando empresas tratam os gastos entre aquisição e retenção de clientes, destinar menos recursos para a retenção impactará em uma lucratividade menor a longo prazo, comparando-se a investimentos menores em aquisição de clientes. Ainda segundo os autores, no conceito de relações de retenção, atribui-se grande ênfase à lealdade e lucratividade de um cliente, sendo lealdade a tendência de o cliente comprar e a lucratividade, a medida geral de quanto lucro um cliente traz à empresa através de suas compras.

De acordo Nguyen, Sherif e Newby (2007, p.114) com o avanço da gerência das relações com clientes foram abertas novas vias pelas quais sua lealdade e lucratividade pode ser cultivada, atraindo uma crescente demanda por parte de empresas, visto que a adoção destes meios permite que as organizações melhorem seu serviço ao consumidor, consequentemente gerando renda. Com isso, diferentes ferramentas acabam sendo utilizadas, como sistemas de recomendação que, geralmente em ramos e-commerce, levam em conta várias características pertinentes ao comportamento do cliente, construindo um perfil próprio que será utilizado para realizar a recomendação de um produto que talvez seja de seu interesse. Outra ferramenta pertinente à lucros e lealdade é a segmentação, que visa separar uma única e confusa massa de clientes em segmentos homogêneos em termos de comportamento, permitindo o desenvolvimento de campanhas e estratégias de marketing especializadas à cada grupo de acordo com suas características (TSIPTSIS; CHORIANOPOULOS, 2009, p.4).

Em relação a segmentação de clientes, algumas métricas tornam-se relevantes nos contextos aos quais estão inseridas. Segundo Kumar (2008, p. 29), o modelo *Recency-Frequency-Monetary* (RFM), é utilizado em empresas de venda por catálogo, enquanto empresas de *high-tech* tendem a usar *Share of Wallet* (SOW) para implementar suas estratégias de marketing. Já o modelo *Past Customer Value* (PCV), geralmente é utilizado em empresas de serviços financeiros. Dentre os modelos citados, o RFM é o que possui maior facilidade de aplicação em diversas áreas de comércio, varejo e supermercados, visto que são necessários apenas os dados transacionais (vendas) dos clientes, dos quais são obtidos os atributos de Recência (R), Frequência (F) e Monetário (M).

A partir desses dados, segundo Tsipstis e Chorianopoulos (2009, p.335), é possível detectar bons clientes a partir das melhores pontuações de RFM. Se o cliente efetuou uma compra recentemente, seu atributo R será alto. Caso ele compre muitas vezes ao longo de um determinado período, seu atributo F será maior. E, por fim, caso seus gastos totais forem significativos, terá um atributo M alto. Ao categorizar o cliente dentro destas três características, é possível obter uma hierarquia de importância, tendo os clientes que possuem valores RFM altos no topo, e clientes que possuem valores baixos na base. Apesar destas

vantagens, o modelo padrão original é um tanto quanto arbitrário, segmentando os clientes em quintis, cinco grupos com 20% dos clientes, não atentando-se às nuances e todas as interpretações que a base de clientes pode possuir. Além disso, o método também pode produzir uma grande quantidade de grupos, que por muitas vezes, não representam significativamente os clientes de um estabelecimento, e caso o método de quintis seja utilizado, 125 grupos serão criados. Outro ponto a se observar é a variada gama de interpretações que os atributos RFM podem ter em relação aos tipos de atividades dos estabelecimentos, sendo necessário a adaptação do modelo para cada empresa.

Diante deste cenário, este trabalho propõe a criação de um artefato computacional que utilize o modelo RFM em conjunto com diferentes algoritmos de clusterização ao invés de quintis para segmentar clientes. Também será extraído de maneira automática as informações de múltiplas bases de dados (atacado, varejo e comércio), visando adequar-se dinamicamente em relação a eventuais diferenças de comportamento dos clientes.

### 1.1 OBJETIVOS

O objetivo deste trabalho é disponibilizar um artefato computacional de auxílio à segmentação de clientes a partir de múltiplas bases de dados utilizando o modelo *Recency-Frequency-Monetary* (RFM).

Os objetivos específicos são:

- implementar e testar diferentes algoritmos de clusterização;
- disponibilizar um mecanismo de visualização dos agrupamentos;
- avaliar a qualidade dos clusters em relação à sua separação e homogeneidade.

Comentado [AS1]: No objetivo geral você não menciona isso.

## 2 TRABALHOS CORRELATOS

Nesta seção serão apresentados os trabalhos similares ao proposto. Na seção 2.1 é apresentado o trabalho de Gustriansyah, Suhandi e Antony (2020), que consiste na aplicação do modelo *Recency Frequency Monetary* (RFM) para clusterização de produtos utilizando K-means, tendo como foco otimizar o número de clusters através de índices de validação. A seção 2.2 descreve o modelo de segmentação denominado *Length, Recency, Frequency, Monetary and Periodicity* (LRFMP) proposto por Peker, Kocyigit e Eren (2017), ao qual considera a longevidade e a periodicidade. Por fim, na seção 2.3 detalha-se o modelo R+FM, sendo uma versão modificada do modelo original que foi aplicada em clientes de uma empresa de e-commerce (TAVAKOLI *et al.*, 2018).

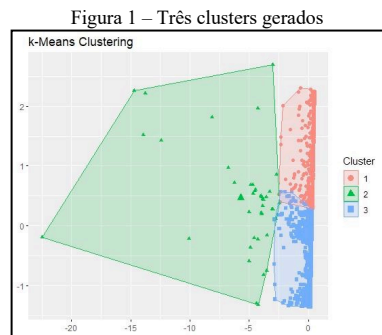
### 2.1 CLUSTERING OPTIMIZATION IN RFM ANALYSIS BASED ON K-MEANS

Gustriansyah, Suhandi e Antony (2020) utilizaram o algoritmo de clusterização K-means para agrupar 2.043 produtos de uma farmácia visando otimizar o manuseio de estoque. Foram utilizadas três características de acordo com o modelo *Recency Frequency Monetary* (RFM) para a separação dos produtos, levando em consideração dados transacionais capturados num período de um ano. O atributo recência classificou os produtos através da última venda realizada num intervalo de 1 a 364 dias. A frequência estabelece a quantidade de transações em que o produto ocorreu, variando num intervalo de 1 a 14.872. Já o atributo monetário, refere-se ao valor total proveniente das vendas acumuladas do produto, sendo definido num intervalo entre 1.250 e 1.151.952.500 Rupias Indonésias (Rp.).

Após a atribuição de valores RFM aos produtos, foram utilizados oito índices de validação do melhor número de clusters: *Elbow Method* (EM), que calcula a variação intra cluster conforme são aumentados os clusters e conclui que o melhor número é aquele que está no cotovelo (elbow) da curva. *Silhouette Index* (SI) que resulta em uma nota de -1 a 1 que indica a quão adequada é a classificação de um objeto dentro de um cluster em comparação aos outros. *Calinski-Harabasz Index* (CHI) que também mede a adequação da quantidade de clusters levando em conta a dispersão entre e intra clusters. *Davies-Bouldin Index* (DBI) que calcula as similaridades entre clusters levando em conta as distâncias e tamanhos dos clusters, quanto menor este índice melhor a separação entre os clusters. *Ratowski Index* (RI) que é baseado na média da soma dos quadrados dos dados entre clusters e a soma total dos quadrados de cada dado dentro de um cluster, dentre as quantidades calculadas escolhe-se a que obtém um maior índice. *Hubert Index* (HI) que é um método visual que indica a quantidade preferida através de um pico no gráfico e é calculado pelo coeficiente de correlação entre matrizes de distância. *Ball-Hall Index* (BHI) definido pela média da distância dos itens com os respectivos centroides do cluster, onde no gráfico o ponto de quantidade de clusters com maior diferença do anterior é sugerido. *Krzanowski-Lai Index* (KLI), que propõe índices internos definidos pelas diferenças entre matrizes de dispersão, e aponta a melhor quantidade de clusters pelo maior número gerado ao realizar a equação com quantidade k. Constatou-se que a maioria deles indicou que o melhor número de clusters seria 3, com base nas condições de interpretação de cada índice explicadas anteriormente.

Segundo Gustriansyah, Suhandi e Antony (2020) nos testes para verificar os clusters gerados, utilizou-se a equação de variância (R). Sendo “R” o valor da divisão entre a distância média dos dados cluster (distância intra-cluster) pela distância média dos dados em outros clusters (inter-cluster). O valor médio alcançado para R foi de 0.19113, sendo que quanto mais próximo de zero, maior a similaridade entre os membros dentro de cada cluster.

Gustriansyah, Suhandi e Antony (2020) utilizaram o software R Programming para gerar a clusterização, resultando na visualização demonstrada na Figura 7, tendo dois clusters com densidade maior e um cluster mais disperso. É possível observar também que o cluster em verde possui uma maior variância entre os próprios dados, enquanto os outros dois clusters possuem uma menor diferença interna.



Fonte: Gustriansyah, Suhandi e Antony (2020).

Segundo Gustriansyah, Suhandi e Antony (2020) também foram adquiridos os valores médios RFM para cada cluster, conforme apresenta a Tabela 1, sendo que o cluster 3 possui a maior média dos três atributos, e o cluster 1 possui a menor média dos três. É possível identificar um intervalo entre os valores médios de cada atributo, indicando uma diferença significativa inter-cluster.

Tabela 1 – Valores médios de RFM de cada cluster

| Cluster | Recency  | Frequency  | Monetary<br>(in thousands) |
|---------|----------|------------|----------------------------|
| 1       | 75.8167  | 3,436.744  | 3,089,608                  |
| 2       | 224.3947 | 13,013.333 | 76,920,847                 |
| 3       | 331.9681 | 107.418    | 286,927,000                |

Fonte: Gustriansyah, Suhandi e Antony (2020).

Gustriansyah, Suhandi e Antony (2020) concluem que o método gerou clusters com alta similaridade em relação aos dados existentes, apresentado uma segmentação mais objetiva quando comparado ao modelo RFM tradicional no qual os dados são divididos igualmente em cinco segmentos (20% dos dados para a cada segmento). Além disso, os autores também sugerem como extensões a utilização de outros métodos para a comparação, como *Particle Swarm Optimization* (PSO), que é um método computacional que otimiza soluções para uma equação de uma certa medida de qualidade, medido de (centroide que são parte do conjunto de dados) ou *maximizing-expectancy*, que é um método iterativo para encontrar estimativas de parâmetros para modelos estatísticos com variáveis não observadas.

## 2.2 LRFMP MODEL FOR CUSTOMER SEGMENTATION IN THE GROCERY RETAIL INDUSTRY: A CASE STUDY

Peker, Kocyigit e Eren (2017) propuseram o modelo *Length, Recency, Frequency, Monetary* (LRFMP) denominado *Length, Recency, Frequency, Monetary and Periodicity* (LRFMP) para classificar dados reais de 16.024 clientes de mercados de uma franquia na Turquia. Para isso, utilizou-se o algoritmo K-means para segmentar os clientes e três índices de validação de clusters para a otimização das suas quantidades, *Silhouette Index* (SI), *Calinski-Harabasz Index* (CHI) e *Davies-Bouldin Index* (DBI). Após a segmentação dos dados, verificou-se estratégias de gerenciamento e relações com os clientes para aumentar a lucratividade, como tratamento preferencial para clientes importantes, implementação de cartões fidelidade para aumentar a frequência de compra de clientes não costumam comprar com frequência, promoções voltadas para clientes incertos com sua escolha de local de compra, dentre outras estratégias.

Primeiramente, Peker, Kocyigit e Eren (2017) adaptaram o modelo LRFM, incluindo o parâmetro *periodicity* (periodicidade), pois a análise dos dados foi realizada a partir do histórico de compras em supermercados, que são estabelecimentos com alto número de visitas, tornando importante a regularidade nos padrões de visita e compra. Peker, Kocyigit e Eren (2017) definem a periodicidade como a regularidade das visitas de um determinado cliente, e é definida pelo desvio padrão dos seus tempos inter-visita (quantia de dias entre duas visitas consecutivas). Se um cliente possui valores baixos de periodicidade, significa que este realiza visitas ou compras em intervalos fixos, podendo caracterizá-lo como cliente regular. Além disso, os autores também modificaram o atributo de recência, transformando-o na média das diferenças entre a data das três últimas compras e a data atual, ao invés da simples diferença entre a data da última compra e a data atual estabelecida no modelo RFM padrão.

Após adquirir os atributos LRFMP dos dados transacionais dos clientes, Peker, Kocyigit e Eren (2017) aplicaram um método de normalização simples nos dados, considerando o intervalo de 0 e 1. Esta normalização foi feita pois os valores LRFMP variam em relação ao intervalo e escala, fato que poderia afetar negativamente a análise dos clusters.

Antes de aplicar a clusterização, Peker, Kocyigit e Eren (2017) utilizaram três índices para validação da quantidade possível de clusters: *Silhouette Index* (SI) que resulta em uma nota de -1 a 1 que indica o quão adequada é a classificação de um objeto dentro de um cluster em comparação aos outros, quanto maior o valor, melhor. *Calinski-Harabasz Index* (CHI) que mede a adequação da quantidade de clusters levando em conta a dispersão entre e intra clusters, um valor alto é preferido. *Davies-Bouldin Index* (DBI) que calcula as similaridades entre clusters levando em conta as distâncias e tamanhos dos clusters, quanto menor este índice melhor será a separação entre os clusters. A partir deles, Peker, Kocyigit e Eren (2017) executaram o algoritmo K-means variando o k de 2 a 9, e os resultados destas iterações foram avaliadas utilizando os três índices. Com base nos resultados, decidiu-se utilizar um número de 5 clusters, pois 2 dos 3 índices sugerem 5 como sendo a quantidade ideal.

Peker, Kocyigit e Eren (2017) utilizaram uma base de dados de uma franquia de mercados que possui mais de dez lojas na cidade de Antália na Turquia. Os dados são compostos por cerca de dois milhões de transações de 16.024 clientes num período de dois anos. Foram removidos os clientes com menos que três compras. Além disso, os autores removeram dados duplicados, transações com valores faltantes assim como, agregaram as compras dentro de um mesmo dia. Depois dessas operações, a quantidade de clientes caiu para 10.471, sendo aplicado na sequência o K-means. A Tabela 6 demonstra a quantidade de clientes nos clusters, os valores médios de LRFMP para cada cluster. Já na última coluna, aplicou-se uma técnica no qual o atributo do cluster recebe uma seta para cima (↑) caso o seu valor for maior que a média do atributo dos outros clusters, e uma seta para baixo (↓), caso seu valor for menor que a média.

Tabela 2 – Valores médios dos clusters

| Cluster | Sample size | Average <i>L</i> | Average <i>R</i> | Average <i>F</i> | Average <i>M</i> | Average <i>P</i> | LRFMP Scores   |
|---------|-------------|------------------|------------------|------------------|------------------|------------------|--|
| 1       | 538         | 633.29           | 39.67            | 175.24           | 24.32            | 4.99             | <i>L</i> ↑ <i>R</i> ↓ <i>F</i> ↑ <i>M</i> ↓ <i>P</i> ↓ |
| 2       | 4,681       | 564.50           | 90.19            | 33.44            | 31.73            | 31.49            | <i>L</i> ↑ <i>R</i> ↓ <i>F</i> ↑ <i>M</i> ↓ <i>P</i> ↓ |
| 3       | 1,091       | 482.17           | 301.18           | 5.33             | 34.21            | 159.32           | <i>L</i> ↑ <i>R</i> ↑ <i>F</i> ↓ <i>M</i> ↓ <i>P</i> ↑ |
| 4       | 818         | 374.01           | 220.24           | 11.85            | 104.18           | 45.81            | <i>L</i> ↓ <i>R</i> ↑ <i>F</i> ↓ <i>M</i> ↑ <i>P</i> ↑ |
| 5       | 3,343       | 173.70           | 399.79           | 10.34            | 30.14            | 27.85            | <i>L</i> ↓ <i>R</i> ↑ <i>F</i> ↓ <i>M</i> ↓ <i>P</i> ↓ |
| Average |             | 419.81           | 218.59           | 28.74            | 36.76            | 43.41            |  |

Fonte: Peker, Kocyigit e Eren (2017).

A partir destes resultados, Peker, Kocyigit e Eren (2017) descreveram as características dos grupos. O grupo 1 representa clientes leais de alta contribuição que, apesar de comporem a menor parcela dos clientes (5,14%), possuem a maior contribuição total entre os grupos. Também é possível observar, que este grupo possui a menor periodicidade média de todos, caracterizando estes clientes como regulares. O grupo 2, representando a maior parcela dos clientes (44,70%) foi classificado como clientes leais de baixa contribuição pois apesar de visitar mais frequentemente as lojas, não possuem tanta contribuição quanto o grupo 1. O grupo 3, com tamanho de 10,42%, foi classificado como clientes incertos, pois possui o atributo de longevidade alto e recência também alta, significando que são clientes com longa história de compra, porém sem muitas compras recentes. Vale notar que este grupo possui o maior valor de periodicidade de todos os grupos, caracterizando-o como um grupo de clientes sem rotina de compra definida. O grupo 4 e 5 foram classificados como clientes perdidos, visto que possuem poucas compras recentes, baixa frequência, e baixa longevidade, denotando um cliente que tem uma pouca interação com a franquia. O grupo 4, contendo uma pequena parcela de 7,81% dos clientes, gasta consideravelmente mais, logo foi classificado como contribuição alta, e o 5, cuja parcela é 31,93%, classificado como contribuição baixa.

A partir desta classificação, Peker, Kocyigit e Eren (2017) estabeleceram estratégias para cada grupo de clientes, como tratamento especial (vagas de estacionamento preferencial, presentes de aniversário, filas preferenciais) para clientes do grupo 1, de maneira a não perder a relação leal com a loja. Para os grupos 3, 4 e 5 cuja frequência é baixa, foi sugerida a adoção de programas de cartão fidelidade para aumentar a frequência deles. Para clientes incertos como no grupo 3, aplicou-se descontos e promoções, de maneira a incentivar os clientes, supostamente sensíveis aos preços, a voltar sua atenção à franquia. Para clientes perdidos, sugeriu-se uma análise mais profunda sobre o motivo da perda, como análise de feedback, inferência de motivos, dentre outros.

Por fim, Peker, Kocyigit e Eren (2017) concluem que o estudo contribuiu com a proposta de um novo modelo RFM, que possibilita uma análise mais profunda que o seu modelo original, visto que as características utilizadas e modificadas permitem uma melhor definição do comportamento de cada cliente. Outra contribuição foi a adição do atributo periodicidade (P) no modelo que, ao contrário do modelo RFM padrão, permite identificar se os clientes de um grupo variam em sua rotina de compras. Outra melhoria apontada é a modificação do atributo de recência, que uma vez calculado como uma média, permite uma caracterização mais precisa que o atributo R comumente utilizado. Uma das limitações destacadas pelos autores é a localidade do estudo, pois foi realizado somente com dados originários de uma cidade, sendo que o comportamento de clientes pode variar de acordo com as diferentes localidades onde é feita a análise. A partir disso, sugere-se uma análise mais ampla contemplando outros locais. Outra sugestão feita por Peker, Kocyigit e Eren (2017) é a adição de novos atributos ao modelo, como a quantidade de produtos comprados, quantidade de produtos percebíveis e não percebíveis comprados, a fim de promover uma interpretação mais profunda do comportamento.

### 2.3 CUSTOMER SEGMENTATION AND STRATEGY DEVELOPMENT BASED ON USER BEHAVIOR ANALYSIS, RFM MODEL AND DATA MINING TECHNIQUES: A CASE STUDY

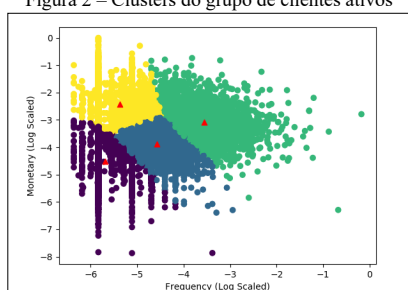
Tavakoli *et al.* (2018) desenvolveram o modelo RFM, denominado “R+FM”, sendo utilizado em conjunto com o algoritmo de clusterização K-means para segmentar 3 milhões de clientes da maior empresa de *E-commerce* do Oriente Médio. Além disso, o modelo de segmentação foi comparado com o utilizado pela empresa, sendo aplicado em uma campanha de Short Message Service (SMS) focada em aumentar os ganhos de cada segmento.

Tavakoli *et al.* (2018) defendem a utilização de um novo modelo de caracterização de clientes, argumentando que o modelo ideal necessita adaptar-se às mudanças de comportamento dos clientes, possuir certa independência de supervisão, levando em consideração a similaridade dos comportamentos dos clientes e a relação entre os atributos Frequência e Monetário. Para isso, construiu-se uma variante do modelo Recency, Monetary, Frequency (RFM) denominado de R+FM, que possui o atributo de recência separado dos demais, utilizando uma segmentação à parte do modelo FM. Os autores separaram os clientes em 3 grupos: os que compraram recentemente (cuja última compra foi dentro de 90 dias), denominados de ativos, clientes que compraram em um passado recente (cuja última compra foi entre 90 e 365 dias), denominados de expirando e por fim, os clientes que não compraram por um longo tempo (cuja última compra foi a mais de 365 dias), denominados de expirados. Para o atributo de frequência, Tavakoli *et al.* (2018) atentaram-se especialmente com a data da primeira compra, pois acreditam que a frequência tem uma importância maior conforme sua recência. Logo, definiram a frequência como a quantidade de compras dividida pela quantidade de dias desde a primeira compra, utilizando também uma função exponencial de decaimento, que efetivamente atribui um peso maior para anos mais recentes, sendo cada ano duas vezes mais pesado que o ano anterior. Como atributo monetário, estabeleceu-se a média dos valores das compras de um cliente, visto que um valor de soma total de compras, segundo os autores, estaria encorajando duas vezes os clientes.

Após a definição do modelo, Tavakoli *et al.* (2018), balancearam a relação entre frequência e monetário, criando a quarta característica que é definida pela combinação linear dos dois atributos, que nada mais é que a soma de cada atributo ponderada pelo peso de cada um. No tratamento dos dados, utilizou-se a técnica de remoção de *outliers* (clientes que não se encaixam no padrão normal) que não se encontram dentro dos intervalos interquartis, que são os intervalos que possuem os dados que pertencem à tendência média do conjunto de dados em geral. Também foram escalados os atributos de frequência e monetário para que seus intervalos sejam iguais, sendo aplicada a normalização *min-max*, que transforma os valores para estarem dentro do intervalo entre 1 e 0. Como os dados monetários e de frequência tratados possuem uma característica de cauda longa, fenômeno estatístico onde os dados são distribuídos de forma decrescente, foi aplicada uma transformação logarítmica para normalizar a distribuição, visto que a quantidade de valores baixos é muito alta, podendo atrapalhar a análise.

Como existem duas segmentações (R e FM), Tavakoli *et al.* (2018) estabeleceram segmentos FM para cada segmento R, resultando nos seguintes grupos: para clientes ativos existem os grupos de alto valor, médio valor com alto monetário, médio valor com alta frequência e baixo valor. Para clientes que estão expirando existem os grupos de alto, médio e baixo valores, sendo aplicado também para clientes expirados. Estes grupos foram definidos por Tavakoli *et al.* (2018) com ajuda da empresa de *E-commerce* Digikala. A partir disso, aplicou-se o K-means com  $k=4$  para o grupo de clientes ativos,  $k=3$  para o grupo de clientes expirando e  $k=3$  para o grupo de clientes expirados, resultando em um total de 10 clusters. Na Figura 8 são identificados os clusters gerados somente a partir do grupo de clientes ativos, organizados em um gráfico de valor monetário por frequência. Sendo o cluster de cor verde composto pelos clientes de alto valor, o cluster de cor amarela composto pelos clientes de médio valor com alto monetário, o cluster de cor azul composto pelos clientes de médio valor com alta frequência, e por fim, o cluster de cor roxa composto pelos clientes de baixo valor.

Figura 2 – Clusters do grupo de clientes ativos



Fonte: Tavakoli *et al.* (2018).

Após a geração dos grupos, Tavakoli *et al.* (2018) discorrem sobre possíveis estratégias para cada segmento, sugerindo um maior foco em clientes ativos com valor médio e baixo, bem como a manutenção de clientes já valiosos. Os autores também enfatizam a importância em recuperar os clientes do grupo expirando, cuja chance de retorno não é tão baixa quanto o grupo expirado, que por si só requer uma estratégia especial de reengajamento dos clientes à empresa.

Além da elaboração de estratégias, Tavakoli *et al.* (2018) implementaram uma campanha de SMS focada somente no segmento de clientes ativos (recência abaixo de 90 dias), pois a empresa já tinha realizado outras campanhas em clientes ativos anteriormente. Nesta campanha cada cliente foi presenteado com um *Voucher* condizente com o segmento ao qual o cliente pertencia. Para clientes ativos com valor alto foi oferecido um desconto de 10 por cento com um desconto máximo de até \$20, com o objetivo de manter a lealdade destes clientes. Para clientes ativos de valor médio com alto valor monetário, foi oferecido um desconto de 10 por cento com um desconto máximo de até \$10, o valor foi menor pois o objetivo era aumentar frequência de compra destes clientes, que em tese já gastam bastante. Para clientes ativos de valor médio com alta frequência foi oferecido um desconto de 10 por cento com um desconto máximo de até \$20 para vendas que custem mais que \$50 (que é o valor médio gasto por este segmento), incentivando assim uma compra de maior valor. Por fim, para clientes ativos de baixo valor, foi oferecido um desconto de 10 por cento com um desconto máximo de até \$20, com o objetivo de converter estes clientes em mais leais.

Para a análise da campanha, Tavakoli *et al.* (2018) selecionaram aleatoriamente 20% dos clientes de cada segmento para compor um grupo de controle, cujos *Vouchers* não foram enviados. Este grupo de controle foi comparado com os outros grupos da campanha para obter um valor de referência do aumento do valor monetário após sua conclusão. Os resultados alcançados por Tavakoli *et al.* (2018) podem ser observados no Quadro 1, ao qual percebe-se um aumento de \$14,30 na média monetária dos clientes ativos com alta frequência, mais do que o aumento de \$3,20 sofrido pelo grupo de controle. É possível também observar que a média monetária de todos os grupos alvo da campanha aumentou consideravelmente, enquanto o grupo de controle aumentou pouco ou até diminuiu, indicando uma efetividade no objetivo da campanha.

Quadro 1 – Dados da média monetária do grupo de controle e grupo de campanha

| Recency | Segment                          | Average Monetary (USD) |                |                 |                |
|---------|----------------------------------|------------------------|----------------|-----------------|----------------|
|         |                                  | Control Users          |                | Campaign Users  |                |
|         |                                  | Before Campaign        | After Campaign | Before Campaign | After Campaign |
| Active  | High Value                       | 74.2                   | 73.8           | 88.2            | 89.2           |
|         | Medium Value with High Monetary  | 100.2                  | 97.6           | 104.6           | 105.2          |
|         | Medium Value with High Frequency | 32                     | 35.2           | 35.4            | 49.7           |
|         | Low Value                        | 50.7                   | 53.2           | 56.4            | 65.2           |

Fonte: Tavakoli *et al.* (2018).

Tavakoli *et al.* (2018) concluem que houve uma melhora no desempenho da campanha lançada em comparação com as anteriores, indicando ainda, que elas obtinham uma taxa de compra de 0,1 por cento, sendo que a campanha lançada para validação do modelo obteve uma taxa de 1 por cento, cerca de dez vezes mais efetivo. Os autores justificam esta melhora ao processo de segmentação do modelo, resultando em clusters mais significativos, facilitando a aplicação de vouchers específicos. Por fim, Tavakoli *et al.* (2018) sugerem o melhoramento da definição do atributo de recência de forma que seja mais útil ao time de marketing. Também recomendam calcular o Customer Lifetime Value (CLV), que atribui o valor vitalício à cada segmento e cliente, de forma a quantificar o valor que um cliente pode proporcionar à empresa.

### 3 PROPOSTA DO PROTÓTIPO

Essa seção visa apresentar a justificativa para a elaboração deste trabalho, os requisitos que serão seguidos e a metodologia que será utilizada.

#### 3.1 JUSTIFICATIVA

No Quadro 2 é apresentado um comparativo entre os trabalhos correlatos. As linhas representam as características relevantes e as colunas representam os trabalhos.

Quadro 2 – Comparativo entre os trabalhos correlatos

| Características                           | Correlatos | Gustriansyah, Suhandi e Antony (2020)   | Peker, Kocyigit e Eren (2017)                         | Tavakoli <i>et al.</i> (2018)                      |
|---|------------|---|---|--|
| Alvo da clusterização                     |            | Produtos                                | Clientes  | Clientes   |
| Modelo utilizado                          |            | RFM                                     | LRFMP   | R+FM   |
| Objetivo da segmentação                   |            | Gerenciamento de estoque                | Gerenciamento das relações com cliente                | Gerenciamento das relações com cliente             |
| Algoritmo de clusterização utilizado      |            | K-means                                 | K-means   | K-means  |
| Foco metodológico                         |            | Otimização de k com diferentes métricas | Formulação de um modelo novo e análise dos resultados | Formulação de um modelo novo e campanha de ofertas |
| Número de dados (clientes/produtos)       |            | 2.043                                   | 16.024  | ~3.000.000   |
| Quantidade de índices para validação de k |            | 8                                       | 3   | -  |
| Quantidade de clusters gerados            |            | 3                                       | 5   | 10   |
| Inferências sobre os dados                |            | -                                       | Sim   | Sim  |

Fonte: elaborado pelo autor.

A partir do Quadro 2, pode-se observar que Gustriansyah, Suhandi e Antony (2020) agruparam produtos de uma base de dados utilizando o modelo RFM padrão. Já Peker, Kocyigit e Eren (2017) optaram pelo desenvolvimento de um modelo novo, considerando a periodicidade (LRFMP). Tavakoli *et al.* (2018) também desenvolveram um novo modelo, ao qual a característica recência foi modificada e separada (R+FM).

Gustriansyah, Suhandi e Antony (2020) tinham como objetivo melhorar o gerenciamento de estoque, prezando por uma segmentação mais conclusiva sobre os produtos, visto que o modelo RFM padrão define segmentos arbitrariamente sem adequar-se às peculiaridades dos dados, enquanto o modelo aplicado através de k-means alcançou uma segmentação com dados altamente similares em cada cluster. Por outro lado, Peker, Kocyigit e Eren (2017) e Tavakoli *et al.* (2018) objetivavam o gerenciamento das relações com os clientes através de estratégias focadas em segmentos, visando aumentar a renda que eles fornecem à empresa. Todos os autores utilizaram o algoritmo K-means, por ser confiável e amplamente difundido. Vale ressaltar que no trabalho de Gustriansyah, Suhandi e Antony (2020), o algoritmo teve um foco metodológico maior, visto que foram utilizados 8 índices de validação para k clusters, visando otimizar a organização dos segmentos.



A quantidade de dados segmentados variou bastante entre os três trabalhos devido aos diferentes contextos de aplicação. Gustriansyah, Suhandi e Antony (2020) tinham 2.043 produtos na base de dados para segmentar, resultando em 3 clusters. Já Peker, Kocyigit e Eren (2017) possuíam o registro de 16.024 clientes de uma rede de padarias, sendo especificados 5 segmentos, obtidos através de uma análise por três índices de validação (Silhouette, Calinski-Harabasz e Davies-Bouldin). Por fim, Tavakoli *et al.* (2018) agruparam dados de 3 milhões de clientes pertencentes à base de dados de um e-commerce do Oriente Médio, resultando em 10 clusters, sendo 3 pertencentes à característica de recência, e os outros 7 distribuídos entre as características de frequência e monetária. Ressalta-se que Tavakoli *et al.* (2018) testaram o modelo em produção, montando uma campanha que focava no segmento de clientes ativos, visando primariamente aumentar os lucros da empresa, utilizando também um grupo de controle e comparação de renda antes e depois da campanha.

Gustriansyah, Suhandi e Antony (2020) demonstraram a possibilidade da aplicação de RFM fora do uso convencional de segmentação de clientes, e adquiriram clusters com uma variância média de 0.19113. Além disso, os autores sugeriram outras formas de comparação de dados, como Particle Swarm Optimization (PSO), medioides ou até *maximizing-expectancy*. Peker, Kocyigit e Eren (2017) segmentaram clientes de uma rede de mercados na Turquia em “clientes leais de alta contribuição”, “clientes leais de baixa contribuição”, “clientes incertos”, “clientes perdidos de alto gasto” e “clientes perdidos de baixo gasto”. Desta maneira, os autores providenciaram visões e estratégias (promoções, ofertas, regalias) de aumento de renda sobre os comportamentos dos clientes, porém limitaram-se a aplicar em um segmento específico de mercado. Por fim, Tavakoli *et al.* (2018) agruparam clientes de uma empresa de e-commerce com base em sua recência, resultando em clientes “Ativos”, “Expirando” e “Expirados”, e destes segmentos, sucessivamente separados em grupos de “Alto”, “Médio” e “Baixo” valores, validando posteriormente a segmentação através de uma campanha de ofertas para os clientes do grupo “Ativos”.

Todos os trabalhos aqui citados procuraram implementar o modelo RFM num contexto de clusterização por K-means, alterando o modelo e o manejo dos dados de acordo com cada categoria, seja ele produto ou cliente, varejo ou mercado. Com isso, criaram-se atributos e foram modificados alguns já existentes para atender às especificidades de cada contexto, visto que todos os trabalhos focaram em uma só base de dados, inevitavelmente adequando-se às mesmas.

Desta forma, este trabalho demonstra ser relevante, pois almeja aplicar o modelo RFM em conjunto com vários algoritmos de clusterização para realizar a segmentação a partir de atributos que realcem o comportamento dos clientes, disponibilizando um artefato computacional que se adequa à vários contextos (mercado, comércio, varejo etc.), utilizando várias bases de dados reais para testar a validade dos algoritmos utilizados. Inicialmente o cenário de análise será o varejo e, conforme os resultados alcançados, serão testados em outras bases/cenários. Vislumbra-se utilizar três índices (Silhouette, Calinski-Harabasz e Davies-Bouldin) para validação da qualidade dos clusters em relação à sua separação e homogeneidade. Além disso, deseja-se obter clusters significativos e coerentes com cada segmento de mercado aplicado. Outra contribuição deste trabalho refere-se ao âmbito comercial, com a geração de informações sobre as similaridades de clientes de cada segmento de mercado, podendo auxiliar gestores e administradores de empresas a obter uma visão crítica sobre os comportamentos de clientes ao longo das diferentes bases de dados, podendo também denotar características comuns a todos. Outra relevância seria a utilização deste trabalho em ambiente acadêmico, visto que serão aplicados diferentes algoritmos de clusterização, podendo providenciar informações sobre seus desempenhos e qualidade de agrupamento, além de ser aplicados processos de obtenção, limpeza e transformação de dados.

**Comentado [AS2]:** Acho que este verbo não condiz com o que você quer dizer aqui.

### 3.2 REQUISITOS PRINCIPAIS DO PROBLEMA A SER TRABALHADO

O artefato computacional a ser desenvolvido deverá:

- d) adquirir os dados transacionais de clientes a partir de um banco de dados (Requisito Funcional - RF);
- e) extrair dos clientes as características (recência, frequência e monetária) utilizadas no modelo RFM (RF);
- f) filtrar os clientes sem quantidade de compras relevantes (RF);
- g) normalizar os dados para evitar disparidades nas escalas dos dados, principalmente no atributo monetário (RF);
- h) aplicar três índices de validação (Silhouette, Calinski-Harabasz e Davies-Bouldin) para verificar a qualidade dos clusters (RF);
- i) apresentar em um gráfico 3D os clientes, com sua localização definida pela pontuação do cliente nas características RFM (RF);
- j) utilizar algoritmos de clusterização tais como hierárquicos, K-means e *Density-Based Spatial*

- Clustering of Applications With Noise* (DBSCAN) (RF);
- k) utilizar a linguagem Python para o desenvolvimento (Requisito Não Funcional - RNF);
  - l) utilizar o ambiente de desenvolvimento Jupyter Notebook (RNF);
  - m) utilizar o banco de dados PostgreSQL para ler os dados das bases utilizadas (RNF).

### 3.3 METODOLOGIA

O trabalho será desenvolvido observando as seguintes etapas:

- n) levantamento bibliográfico: pesquisar trabalhos relacionados e estudar sobre o modelo RFM e suas aplicações, algoritmos de clusterização, métodos de tratamento de dados e índices de validação;
- o) seleção de bases de dados: obter bases de dados de usuários cedidas pela empresa Intelidata Informática, desenvolvedora de software de gestão comercial. Serão selecionadas conforme sua adequação ao objetivo do trabalho, variando em tamanho e segmento de mercado;
- p) definição das características do modelo RFM: definir os atributos utilizados para caracterizar os clientes no modelo RFM;
- q) definição de métricas do modelo RFM: definir as métricas para mensuração e atribuição de pontuação de cada característica no modelo RFM;
- r) definição dos algoritmos de clusterização: pesquisar e escolher o algoritmo de clusterização que realizará o agrupamento das características RFM;
- s) implementação: implementar o artefato computacional de segmentação levando em consideração as etapas (b) até (e), utilizando a linguagem Python;
- t) análise dos clusters: avaliar a qualidade dos clusters gerados a partir dos diferentes algoritmos de clusterização e seus comportamentos em múltiplas bases de dados, aplicando índices de validação e apresentando os agrupamentos na forma de gráficos, de maneira que a sua separação e homogeneidade possam ser observadas.

As etapas serão realizadas nos períodos relacionados no Quadro 3.

Quadro 3 – Cronograma de atividades a serem realizadas

| etapas / quinzenas                          | 2022 |   |      |   |      |   |      |   |      |   |
|---|------|---|------|---|------|---|------|---|------|---|
|   | fev. |   | mar. |   | abr. |   | maio |   | jun. |   |
|   | 1    | 2 | 1    | 2 | 1    | 2 | 1    | 2 | 1    | 2 |
| levantamento bibliográfico                  |      |   |      |   |      |   |      |   |      |   |
| seleção de bases de dados                   |      |   |      |   |      |   |      |   |      |   |
| definição das características do modelo RFM |      |   |      |   |      |   |      |   |      |   |
| definição de métricas do modelo RFM         |      |   |      |   |      |   |      |   |      |   |
| definição dos algoritmos de clusterização   |      |   |      |   |      |   |      |   |      |   |
| implementação                               |      |   |      |   |      |   |      |   |      |   |
| análise dos clusters                        |      |   |      |   |      |   |      |   |      |   |

Fonte: elaborado pelo autor.

## 4 REVISÃO BIBLIOGRÁFICA

Nesta seção são abordados os assuntos que servirão de base para a realização deste trabalho. A seção 4.1 discorre sobre clustering. A seção 4.2 apresenta o modelo RFM. Por fim, a seção 4.3 aborda os índices de validação de clusters.

### 4.1 CLUSTERING

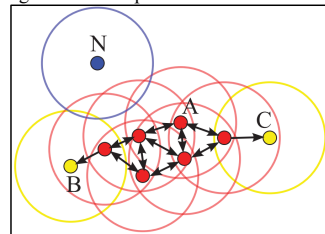
Para Cherkassky e Mulier (2007), clustering se trata do problema de separar um conjunto de dados em grupos chamados de “clusters” baseado em alguma medida de similaridade. O objetivo é encontrar um conjunto de clusters dos quais as amostras dentro dos mesmos são mais similares entre si do que quando comparadas com amostras de outros clusters. A análise destes clusters consiste no fato de que a medida de similaridade entre eles é escolhida subjetivamente baseado na sua habilidade de criar clusters interessantes ao analista. Cherkassky e Mulier (2007) classificam os algoritmos em dois tipos principais: (i) Hierárquicos, seguindo uma estrutura de árvores; (ii) Particionais, que geram clusters a partir de sucessivas secções, cujos métodos são identificados em dois grupos: particionais onde cada dado é atribuído à um e somente à um cluster e particionais que podem pertencer à vários clusters.

Segundo Schubert *et al.* (2017), o algoritmo *Density-Based Spatial Clustering of Applications With Noise* (DBSCAN) foi publicado em 1996, sendo comumente utilizado com sucesso em várias aplicações do mundo real. Schubert *et al.* (2017) apontam que o modelo DBSCAN utiliza uma estimativa de densidade mínima simples, baseada no número mínimo de vizinhos (*minPts*) dentro de um raio  $\epsilon$ , cuja medida de

Comentado [AS3]: Sugestão: de um

distância é arbitrária. A partir disso, objetos com mais que o número mínimo de vizinhos dentro do raio (incluindo o ponto analisado) são considerados como pontos centrais. Todos os vizinhos de um ponto central que estejam dentro do raio  $\epsilon$  são considerados parte do mesmo cluster que o ponto central. Se qualquer um desses vizinhos for também um ponto central, suas vizinhanças são incluídas no cluster. Além disso, pontos não centrais deste cluster são chamados de pontos-limite. Pontos que não são alcançáveis por qualquer ponto central são considerados ruído e não pertencem a nenhum cluster. A Figura 3 exemplifica o funcionamento do DBSCAN. O parâmetro  $minPts$  é 4, o raio  $\epsilon$  é denotado pelos círculos, N é um ponto ruído, A são pontos centrais e B e C são pontos-limite.

Figura 3 – Exemplo do modelo DBSCAN



Fonte: Schubert *et al.* (2017).

Zhao *et al.* (2005) ressaltam que a maioria dos algoritmos de *Hierarchical Clustering* são aglomerativos, cujos objetos são inicialmente designados para um cluster. A partir disso, clusters são repetidamente combinados até que a árvore de clusters total seja formada, ao qual seleciona-se uma altura de corte para obter os clusters desejados. Apesar disso, os autores descrevem que algoritmos particionais também podem ser utilizados para obter soluções hierárquicas através de sequências de bissecções repetidas. Segundo Zhao *et al.* (2005), na abordagem hierárquica aglomerativa, o parâmetro-chave é o método usado para determinar o par de clusters a serem unidos a cada execução do algoritmo. Na maioria das abordagens o par mais similar é selecionado, podendo ser feito de várias maneiras tradicionais como a *single-link*, *complete-link* e média grupal (ou *Unweighted Pair Group Method with Arithmetic mean - UPGMA*). O esquema *single-link* mede a similaridade de dois clusters através dos pontos mais próximos entre eles. Já o *complete-link* mede inversamente, de maneira que a similaridade é obtida através dos pontos mais distantes entre os dois clusters. Por fim, o método UPGMA leva em consideração a distância média entre elementos de cada cluster.

De acordo com Ghosh e Kumar (2013), o K-means é um método de particionamento mais aplicado para analisar dados, separando os objetos em clusters mutuamente exclusivos (K) de maneira que os objetos de cada cluster fiquem tão perto entre si quanto possível, porém tão longe quanto possível de objetos em clusters diferentes. Cada cluster possui um ponto central (centroide), cuja localização é obtida através da média da localização de todos os pontos pertencentes ao cluster. Segundo Ghosh e Kumar (2013), o algoritmo baseia-se na constante atualização da posição dos centroides e recálculo dos pontos mais próximos, sendo que inicialmente os k-centroides são aleatoriamente distribuídos no espaço. O algoritmo acaba sua execução quando nenhum ponto muda de cluster ou nenhum centroide se move.

#### 4.2 MODELO RFM

Segundo Hughes (2011), o modelo RFM é “Um meio antigo e altamente preditivo de determinar quem irá responder e comprar. Um método de codificar clientes existentes. Usado para prever resposta, tamanho médio de pedido, e outros fatores”. Este modelo categoriza geralmente clientes através das características de Recência (R), Frequência (F) e Monetária (M). As métricas utilizadas para medir tais características podem variar, porém geralmente classificam recência como a quantidade de dias desde a última compra, frequência como a quantidade de compras dentro de um determinado período, e monetária como o total acumulado de todas as vendas realizadas para um cliente.

No estudo realizado por Verhoef *et al.* (2003), cerca de 90% das empresas questionadas sobre a aplicação métodos de segmentação (como o RFM) afirmam que possuíam como objetivo a seleção de alvos, ou seja, encontrar o segmento de clientes que mais se identificam com a empresa, e 64,4% citaram como objetivo o tratamento diferencial de clientes, com promoções, preços e ofertas especiais. Verhoef *et al.* (2003) ainda evidenciam que este modelo é geralmente utilizado em marketing direto, onde a comunicação acontece diretamente entre a empresa e o consumidor, realizada através de mídias sociais, e-mail, mensagens SMS ou até pelo correio.

Comentado [AS4]: coloca o significado da sigla aqui

Christy *et al.* (2021) indicam que o propósito da pontuação RFM é projetar comportamentos futuros, ajudando nas decisões de segmentação posteriores. Para realizar essa projeção, é importante transformar este comportamento em números que podem ser utilizados ao longo do tempo. Um método comum de atribuição de pontuação seria utilizando quintis, sendo este um dos métodos mais utilizados. Ele consiste em organizar os clientes em ordem decrescente (melhor para pior), dividi-los em 5 grupos de tamanhos iguais, os melhores recebem a pontuação de 5 e os piores 1. Os melhores clientes então, teriam a pontuação total de 5,5,5.

**Comentado [AS5]:** define aqui o que cada valor 5 significa

Segundo Tsoy e Shchekoldin (2016), o método de quintis possui alguns problemas, sendo um deles a dificuldade de detectar nuances dentro de cada quintil, muitas vezes agrupando clientes com comportamentos muito diferentes. Isto se aplica geralmente nos quintis do topo, onde há um número bem pequeno de clientes muito importantes, e um número grande de clientes menos importantes, deixando claro uma heterogeneidade presente dentro do próprio grupo em questão. Outras vezes, nos quintis baixos, os autores afirmam que este método tende arbitrariamente a separar clientes com comportamentos parecidos.

De acordo com Safari, Safari e Montazer (2016), um dos processos mais importantes da aplicação do modelo RFM é a designação do peso dos atributos seguida da ponderação, onde é multiplicado o peso desejado pelo próprio valor do atributo. As fórmulas mais utilizadas tendem a priorizar recência acima de todos, seguido de frequência e por último monetário. Desta maneira, elevando clientes recentes à uma prioridade maior que os outros. Segundo os autores, esta designação pode variar de acordo com o contexto de cada empresa que emprega o modelo RFM, requerendo conhecimento prévio do ramo de negócios do estabelecimento para a correta atribuição de pesos.

#### 4.3 ÍNDICES DE VALIDAÇÃO

Para Hämäläinen, Jauhiainen e Kärkkäinen (2017), uma validação de clusters considera a qualidade do resultado de um algoritmo de cluster, tentando achar a separação que melhor se encaixa com a natureza dos dados. O número de clusters dado como parâmetro para vários algoritmos deveria ser decidido com base na estrutura natural dos dados, porém não há solução clara sobre o melhor número de clusters.

Hämäläinen, Jauhiainen e Kärkkäinen (2017) destacam que os índices de validação medem a tangibilidade do objetivo, sendo esta, uma alta similaridade entre os dados dentro de um cluster e uma alta diferença entre os dados de clusters diferentes. Essas medidas são chamadas de separação intra e inter cluster, respectivamente. Uma boa medida de separação intra-cluster possui números baixos, já uma boa medida de separação inter cluster possui números altos. Os autores ainda destacam que nenhum índice de validação é perfeito para qualquer contexto, alguns índices são mais adequados à diferentes tipos de dados, logo recomenda-se utilizar múltiplos índices para a análise de clusters. Arbelaitz *et al.* (2013) concluem em seu estudo através de uma análise estatística, que dos 30 índices pesquisados, 10 provam ser recomendáveis para utilização. No topo desta lista, encontram-se os índices Silhouette, Calinski-Harabasz e Davies-Bouldin.

De acordo com Rousseeuw (1987), para gerar o índice de Silhouette de um dado, são necessárias apenas duas coisas: os clusters obtidos e o conjunto das distâncias entre todos os dados observados, sendo calculado para cada  $i$  o seu respectivo índice Silhouette  $s(i)$ . Calcula-se também a média de dissimilaridade das distâncias de  $i$  com o resto dos dados do cluster de  $i$ , denotado por  $a(i)$ . No passo seguinte, é obtido o valor mínimo entre as distâncias de  $i$  e qualquer outro cluster (é descoberto então o cluster vizinho de  $i$ , ou seja, o cluster com que  $i$  mais se encaixaria caso não estivesse em seu cluster original), denotado por  $b(i)$ . Este processo pode ser resumido pela Equação 1, que resulta em um número entre -1 e 1, sendo -1 uma categorização ruim do objeto  $i$  (não condizente com seu cluster atual) e sendo 1 uma categorização ótima. Para obter a qualidade da clusterização em geral, é obtida a média de  $s(i)$  para todos os objetos  $i$  do conjunto de dados.

$$s(i) = \frac{b(i) - a(i)}{\max\{a(i), b(i)\}} \quad (1)$$

Para Calinski e Harabasz (1974), o seu índice Variance Rate Criterion (VRC), demonstrado na Equação 2, considera: a soma dos quadrados entre grupos (Between Group Sum of Squares - BGSS) que retrata a variância entre clusters levando em conta a distância de seus centroides até o centroide global; a soma dos quadrados dentro grupos (Within Group Sum of Squares - WGSS) que retrata a variância dentro de clusters levando em conta as distâncias dos pontos em um cluster até o seu centroide. Considera-se também o número de observações/dados ( $n$ ) e o número de clusters ( $k$ ). Quando este índice é utilizado, procura-se maximizar o resultado conforme o valor de  $k$  é aumentado.

$$VRC = \frac{BGSS}{k-1} / \frac{WGSS}{n-k} \quad (2)$$

Davies e Bouldin (1979) denotam que o objetivo de seu índice é definir uma medida de separação de clusters  $R(S_i, S_j, M_{ij})$  que permita a computação da similaridade média de cada cluster com o seu cluster mais similar (vizinho), o valor mais baixo possível seria o resultado ideal. Com  $S_i$  sendo a medida de dispersão do cluster  $i$ ,  $S_j$  sendo a medida de dispersão do cluster  $j$  e  $M_{ij}$  sendo a distância entre os clusters  $i$  e  $j$ , conforme Equação 3.

$$R_{ij} \equiv \frac{S_i + S_j}{M_{ij}} \quad R \equiv \frac{1}{N} \sum_{i=1}^N R_i \quad (3)$$

De acordo com Davies e Bouldin (1979), primeiro obtêm-se  $R_{ij}$  de todos os clusters, isto é, a razão de distâncias inter e intra-cluster entre o cluster  $i$  e  $j$ . Após isso, obtêm-se  $R_i$  (o valor mais alto de  $R_{ij}$ ) identificando para cada cluster, o cluster vizinho ao qual ele mais se assemelha. Por fim, é calculado o índice em si ( $R$ ), sendo este, a soma total das similaridades de  $n$  clusters com seus vizinhos mais próximos.

### REFERÊNCIAS

- ARBELAITZ, Olatz *et al.* An extensive comparative study of cluster validity indices. **Pattern Recognition**, [S.l.], v. 46, n. 1, p. 243-256, jan. 2013. Elsevier BV. <http://dx.doi.org/10.1016/j.patcog.2012.07.021>.
- CALINSKI, T.; HARABASZ, J. A dendrite method for cluster analysis. **Communications In Statistics - Theory And Methods**, [S.l.], v. 3, n. 1, p. 1-27, 1974. Informa UK Limited. <http://dx.doi.org/10.1080/03610927408827101>.
- CHERKASSKY, Vladimir S.; MULIER, Filip. Methods for data reduction and dimensionality reduction. In: CHERKASSKY, Vladimir S.; MULIER, Filip. **Learning from data: concepts, theory, and methods**. 2. ed. Hoboken: Ieee Press, 2007. Cap. 6, p. 191
- CHRISTY, A. Joy *et al.* RFM ranking – An effective approach to customer segmentation. **Journal Of King Saud University - Computer And Information Sciences**, [S.l.], v. 33, n. 10, p. 1251-1257, dez. 2021. Elsevier BV. <http://dx.doi.org/10.1016/j.jksuci.2018.09.004>.
- DAVIES, David L.; BOULDIN, Donald W. A Cluster Separation Measure. **IEEE Transactions on Pattern Analysis And Machine Intelligence**, [S.l.], v. -1, n. 2, p. 224-227, abr. 1979. Institute of Electrical and Electronics Engineers (IEEE). <http://dx.doi.org/10.1109/tpami.1979.4766909>.
- GHOSH, Soumi; KUMAR, Sanjay. Comparative Analysis of K-Means and Fuzzy C-Means Algorithms. **International Journal of Advanced Computer Science And Applications**, [S.l.], v. 4, n. 4, p. 35-39, 2013. The Science and Information Organization. <http://dx.doi.org/10.14569/ijacsa.2013.040406>.
- GUSTRIANSYAH, Rendra; SUHANDI, Nazori; ANTONY, Fery. Clustering optimization in RFM analysis Based on k-Means. **Indonesian Journal of Electrical Engineering And Computer Science**, [S.l.], v. 18, n. 1, p. 470-477, abr. 2020. Mensal.
- HÄMÄLÄINEN, Joonas; JAUHAINEN, Susanne; KÄRKKÄINEN, Tommi. Comparison of Internal Clustering Validation Indices for Prototype-Based Clustering. **Algorithms**, [S.l.], v. 10, n. 3, p. 105-119, 6 set. 2017. MDPI AG. <http://dx.doi.org/10.3390/a10030105>.
- HUGHES, Arthur M. **Strategic Database Marketing 4e: the masterplan for starting and managing a profitable, customer-based marketing program**. 4. ed. [S.l.]: McGraw-Hill, 2011. 608 p.
- KUMAR, Vijay. **Managing Customers for Profit: strategies to increase profits and build loyalty**. Upper Saddle River: Pearson Prentice Hall, 2008. 296 p.
- NGUYEN, Thuyuyen H.; SHERIF, Joseph S.; NEWBY, Michael. **Strategies for successful CRM implementation**. Information Management & Computer Security, [S.l.], v. 15, n. 2, p. 102-115, maio 2007.
- PEKER, Serhat; KOCYIGIT, Altan; EREN, P. Erhan. LRFMP model for customer segmentation in the grocery retail industry: a case study. **Marketing Intelligence & Planning**, [S.l.], v. 35, n. 4, p. 544-559, 6 maio 2017. Emerald. <http://dx.doi.org/10.1108/mip-11-2016-0210>.
- PETRISON, Lisa A.; BLATTBERG, Robert C.; WANG, Paul. Database marketing: past, present, and future. **Journal Of Direct Marketing**, [S.l.], v. 11, n. 4, p. 109-125, mar. 1997. Wiley. [http://dx.doi.org/10.1002/\(sici\)1522-7138\(199723\)11:43.0.co;2-g](http://dx.doi.org/10.1002/(sici)1522-7138(199723)11:43.0.co;2-g).
- RASHID, Mohammad A.; HOSSAIN, Liaquat; PATRICK, Jon David. The Evolution of ERP Systems: a historical perspective. In: NAH, Fiona Fui-Hoon. **Enterprise Resource Planning: solutions and management**. Hershey: Irm Press, 2001. p. 35-50.
- REINARTZ, Werner; THOMAS, Jacquelyn S.; KUMAR, V. Balancing Acquisition and Retention Resources to Maximize Customer Profitability. **Journal Of Marketing**, [S.l.], v. 69, n. 1, p. 63-79, jan. 2005.

Excluído: eee

Excluído: On

Excluído: Of

Excluído: Of

- ROUSSEEUW, Peter J. Silhouettes: a graphical aid to the interpretation and validation of cluster analysis. **Journal Of Computational And Applied Mathematics**, [S.l.], v. 20, n. 0, p. 53-65, nov. 1987. Elsevier BV. [http://dx.doi.org/10.1016/0377-0427\(87\)90125-7](http://dx.doi.org/10.1016/0377-0427(87)90125-7).
- SAFARI, Fariba; SAFARI, Narges; MONTAZER, Gholam Ali. Customer lifetime value determination based on RFM model. **Marketing Intelligence & Planning**, [S.l.], v. 34, n. 4, p. 446-461, 6 jun. 2016. Emerald. <http://dx.doi.org/10.1108/mip-03-2015-0060>.
- SCHUBERT, Erich *et al.* DBSCAN Revisited, Revisited: why and how you should (still) use DBSCAN. **Acm Transactions On Database Systems**, [S.l.], v. 42, n. 3, p. 1-21, 24 ago. 2017. Association for Computing Machinery (ACM). <http://dx.doi.org/10.1145/3068335>.
- TAVAKOLI, Mohammadreza *et al.* Customer Segmentation and Strategy Development Based on User Behavior Analysis, RFM Model and Data Mining Techniques: a case study. In: 2018 IEEE 15TH INTERNATIONAL CONFERENCE ON E-BUSINESS ENGINEERING (ICEBE), 15., 2018, Xiam. **Proceedings [...]**. [S.l.]: Ieee, 2018. p. 119-126.
- TSIPTSIS, Konstantinos K.; CHORIANOPOULOS, Antonios. **Data Mining Techniques in CRM: inside customer segmentation**. Chichester: John Wiley & Sons, 2009. 374 p.
- TSOY, Marina E.; SHCHEKOLDIN, Vladislav Yu. RFM-analysis as a tool for segmentation of high-tech products' consumers. In: 2016 13TH INTERNATIONAL SCIENTIFIC-TECHNICAL CONFERENCE ON ACTUAL PROBLEMS OF ELECTRONICS INSTRUMENT ENGINEERING (APEIE), 13., 2016, Novosibirsk. **2016 13th International Scientific-Technical Conference on Actual Problems of Electronics Instrument Engineering (APEIE)**. [S.l.]: Ieee, 2016. p. 290-293.
- VERHOEF, Peter C *et al.* The commercial use of segmentation and predictive modeling techniques for database marketing in the Netherlands. **Decision Support Systems**, [S.l.], v. 34, n. 4, p. 471-481, mar. 2003.
- ZHAO, Ying *et al.* Hierarchical Clustering Algorithms for Document Datasets. **Data Mining And Knowledge Discovery**, [S.l.], v. 10, n. 2, p. 141-168, mar. 2005. Springer Science and Business Media LLC. <http://dx.doi.org/10.1007/s10618-005-0361-3>.

## FORMULÁRIO DE AVALIAÇÃO BCC – PROFESSOR AVALIADOR

Avaliador(a): Andreza Sartori

Atenção: quando o avaliador marcar algum item como atende parcialmente ou não atende, deve obrigatoriamente indicar os motivos no texto, para que o aluno saiba o porquê da avaliação.

| ASPECTOS AVALIADOS     |  | Atende | atende parcialmente | não atende |
|------------------------|--|--------|---------------------|------------|
| ASPECTOS TÉCNICOS      | 1. INTRODUÇÃO<br>O tema de pesquisa está devidamente contextualizado/delimitado?   | X      |                     |            |
|                        | O problema está claramente formulado?  | X      |                     |            |
|                        | 2. OBJETIVOS<br>O objetivo principal está claramente definido e é passível de ser alcançado?   |        | X                   |            |
|                        | Os objetivos específicos são coerentes com o objetivo principal?   | X      |                     |            |
|                        | 3. TRABALHOS CORRELATOS<br>São apresentados trabalhos correlatos, bem como descritas as principais funcionalidades e os pontos fortes e fracos?                          | X      |                     |            |
|                        | 4. JUSTIFICATIVA<br>Foi apresentado e discutido um quadro relacionando os trabalhos correlatos e suas principais funcionalidades com a proposta apresentada?             | X      |                     |            |
|                        | São apresentados argumentos científicos, técnicos ou metodológicos que justificam a proposta?  | X      |                     |            |
|                        | São apresentadas as contribuições teóricas, práticas ou sociais que justificam a proposta?   | X      |                     |            |
|                        | 5. REQUISITOS PRINCIPAIS DO PROBLEMA A SER TRABALHADO<br>Os requisitos funcionais e não funcionais foram claramente descritos?   | X      |                     |            |
|                        | 6. METODOLOGIA<br>Foram relacionadas todas as etapas necessárias para o desenvolvimento do TCC?  | X      |                     |            |
|                        | Os métodos, recursos e o cronograma estão devidamente apresentados e são compatíveis com a metodologia proposta?   | X      |                     |            |
|                        | 7. REVISÃO BIBLIOGRÁFICA (atenção para a diferença de conteúdo entre projeto e pré-projeto)<br>Os assuntos apresentados são suficientes e têm relação com o tema do TCC? | X      |                     |            |
| ASPECTOS METODOLÓGICOS | As referências contemplam adequadamente os assuntos abordados (são indicadas obras atualizadas e as mais importantes da área)?   | X      |                     |            |
|                        | 8. LINGUAGEM USADA (redação)<br>O texto completo é coerente e redigido corretamente em língua portuguesa, usando linguagem formal/científica?                            | X      |                     |            |
|                        | A exposição do assunto é ordenada (as ideias estão bem encadeadas e a linguagem utilizada é clara)?  | X      |                     |            |

O projeto de TCC será reprovado se:

- qualquer um dos itens tiver resposta NÃO ATENDE;
- pelo menos 4 (**quatro**) itens dos **ASPECTOS TÉCNICOS** tiverem resposta ATENDE PARCIALMENTE; ou
- pelo menos 4 (**quatro**) itens dos **ASPECTOS METODOLÓGICOS** tiverem resposta ATENDE PARCIALMENTE.

**PARECER:** ( x ) APROVADO ( ) REPROVADO

|                                      |               |                      |
|--------------------------------------|---------------|----------------------|
| CURSO DE CIÊNCIA DA COMPUTAÇÃO – TCC |               |                      |
| ( ) PRÉ-PROJETO                      | ( X ) PROJETO | ANO/SEMESTRE: 2021/2 |

## UTILIZAÇÃO DE CLUSTERIZAÇÃO PARA AUXÍLIO EM TOMADA DE DECISÃO A PARTIR DE DADOS DE VAREJO

Henrique José Wilbert

Prof. Aurélio Faustino Hoppe – Orientador

Prof. Christian Daniel Falaster – Coorientador

### 1 INTRODUÇÃO

Com a evolução da tecnologia de informação a partir dos anos 80 e início dos anos 90, várias grandes empresas adotaram sistemas de gerenciamento na forma de softwares Enterprise Resource Planning (ERP) (RASHID; HOSSAIN; PATRICK, 2001, p.2). Estes softwares auxiliam em suas rotinas à nível operacional, seja no controle do estoque, fiscal, financeiro, transacional e até recursos humanos. A partir disso, alcançou-se um patamar de eficiência nunca concebido, visto que registros antes realizados em papel e caneta, passaram a ser produzidos automaticamente. Ainda segundo os autores, em paralelo a informatização desses processos, houve também um crescimento da quantidade de dados armazenados referentes à produtos, clientes, transações, gastos e receitas.

Diante deste contexto, avançaram-se também as táticas de marketing direto, como por exemplo, o envio de catálogos por correio, até ofertas altamente objetivas para indivíduos selecionados, cujas informações transacionais estavam presentes na base de dados. Percebeu-se então que o foco das relações empresa-cliente não está em clientes novos, e sim em clientes já existentes nas bases de dados, visto que o custo para adquirir um cliente novo através de publicidade é muito maior que o custo de alimentar uma relação já existente (PETRISON; BLATTBERG; WANG, 1997, p. 119).

Segundo Reinartz, Thomas e Kumar (2005, p.77), quando empresas tratam os gastos entre aquisição e retenção de clientes, destinar menos recursos para a retenção impactará em uma lucratividade menor a longo prazo, comparando-se a investimentos menores em aquisição de clientes. Ainda segundo os autores, no conceito de relações de retenção, atribui-se grande ênfase à lealdade e lucratividade de um cliente, sendo lealdade a tendência de o cliente comprar e a lucratividade, a medida geral de quanto lucro um cliente traz à empresa através de suas compras.

De acordo Nguyen, Sherif e Newby (2007, p.114) com o avanço da gerência das relações com clientes foram abertas novas vias pelas quais sua lealdade e lucratividade pode ser cultivada, atraindo uma crescente demanda por parte de empresas, visto que a adoção destes meios permite que as organizações melhorem seu serviço ao consumidor, consequentemente gerando renda. Com isso, diferentes ferramentas acabam sendo utilizadas, como sistemas de recomendação que, geralmente em ramos e-commerce, levam em conta várias características pertinentes ao comportamento do cliente, construindo um perfil próprio que será utilizado para realizar a recomendação de um produto que talvez seja de seu interesse. Outra ferramenta pertinente à lucros e lealdade é a segmentação, que visa separar uma única e confusa massa de clientes em segmentos homogêneos em termos de comportamento, permitindo o desenvolvimento de campanhas e estratégias de marketing especializadas à cada grupo de acordo com suas características (TSIPTSIS; CHORIANOPOULOS, 2009, p.4).

Em relação a segmentação de clientes, algumas métricas tornam-se relevantes nos contextos aos quais estão inseridas. Segundo Kumar (2008, p. 29), o modelo *Recency-Frequency-Monetary* (RFM), é utilizado em empresas de venda por catálogo, enquanto empresas de *high-tech* tendem a usar *Share of Wallet* (SOW) para implementar suas estratégias de marketing. Já o modelo *Past Customer Value* (PCV), geralmente é utilizado em empresas de serviços financeiros. Dentre os modelos citados, o RFM é o que possui maior facilidade de aplicação em diversas áreas de comércio, varejo e supermercados, visto que são necessários apenas os dados transacionais (vendas) dos clientes, dos quais são obtidos os atributos de Recência (R), Frequência (F) e Monetário (M).

A partir desses dados, segundo Tsipstis e Chorianopoulos (2009, p.335), é possível detectar bons clientes a partir das melhores pontuações de RFM. Se o cliente efetuou uma compra recentemente, seu atributo R será alto. Caso ele compre muitas vezes ao longo de um determinado período, seu atributo F será maior. E, por fim, caso seus gastos totais forem significativos, terá um atributo M alto. Ao categorizar o cliente dentro destas três características, é possível obter uma hierarquia de importância, tendo os clientes que possuem valores RFM altos no topo, e clientes que possuem valores baixos na base. Apesar destas



vantagens, o modelo padrão original é um tanto quanto arbitrário, segmentando os clientes em quintis, cinco grupos com 20% dos clientes, não atentando-se às nuances e todas as interpretações que a base de clientes pode possuir. Além disso, o método também pode produzir uma grande quantidade de grupos, que por muitas vezes, não representam significativamente os clientes de um estabelecimento, e caso o método de quintis seja utilizado, 125 grupos serão criados. Outro ponto a se observar é a variada gama de interpretações que os atributos RFM podem ter em relação aos tipos de atividades dos estabelecimentos, sendo necessário a adaptação do modelo para cada empresa.

Diante deste cenário, este trabalho propõe a criação de um artefato computacional que utilize o modelo RFM em conjunto com diferentes algoritmos de clusterização ao invés de quintis para segmentar clientes. Também será extraído de maneira automática as informações de múltiplas bases de dados (atacado, varejo e comércio), visando adequar-se dinamicamente em relação a eventuais diferenças de comportamento dos clientes.

## 1.1 OBJETIVOS

O objetivo deste trabalho é disponibilizar um artefato computacional de auxílio à segmentação de clientes a partir de múltiplas bases de dados utilizando o modelo RFM.

Os objetivos específicos são:

- implementar e testar diferentes algoritmos de clusterização;
- disponibilizar um mecanismo de visualização dos agrupamentos;
- avaliar a qualidade dos clusters em relação à sua separação e homogeneidade.

## 2 TRABALHOS CORRELATOS

Nesta seção serão apresentados os trabalhos similares ao proposto. Na seção 2.1 é apresentado o trabalho de Gustriansyah, Suhandi e Antony (2020), que consiste na aplicação do modelo *Recency Frequency Monetary* (RFM) para clusterização de produtos utilizando K-means, tendo como foco otimizar o número de clusters através de índices de validação. A seção 2.2 descreve o modelo de segmentação denominado *Length, Recency, Frequency, Monetary and Periodicity* (LRFMP) proposto por Peker, Kocyigit e Eren (2017), ao qual considera a longevidade e a periodicidade. Por fim, na seção 2.3 detalha-se o modelo R+FM, sendo uma versão modificada do modelo original que foi aplicada em clientes de uma empresa de e-commerce (TAVAKOLI *et al.*, 2018).

### 2.1 CLUSTERING OPTIMIZATION IN RFM ANALYSIS BASED ON K-MEANS

Gustriansyah, Suhandi e Antony (2020) utilizaram o algoritmo de clusterização K-means para agrupar 2.043 produtos de uma farmácia visando otimizar o manuseio de estoque. Foram utilizadas três características de acordo com o modelo *Recency Frequency Monetary* (RFM) para a separação dos produtos, levando em consideração dados transacionais capturados num período de um ano. O atributo recência classificou os produtos através da última venda realizada num intervalo de 1 a 364 dias. A frequência estabelece a quantidade de transações em que o produto ocorreu, variando num intervalo de 1 a 14.872. Já o atributo monetário, refere-se ao valor total proveniente das vendas acumuladas do produto, sendo definido num intervalo entre 1.250 e 1.151.952.500 Rupias Indonésias (Rp.).

Após a atribuição de valores RFM aos produtos, foram utilizados oito índices de validação do melhor número de *clusters*: *Elbow Method* (EM), que calcula a variação intra cluster conforme são aumentados os clusters e conclui que o melhor número é aquele que está no cotovelo (elbow) da curva. *Silhouette Index* (SI) que resulta em uma nota de -1 a 1 que indica a quão adequada é a classificação de um objeto dentro de um cluster em comparação aos outros. *Calinski-Harabasz Index* (CHI) que também mede a adequação da quantidade de clusters levando em conta a dispersão entre e intra clusters. *Davies-Bouldin Index* (DBI) que calcula as similaridades entre clusters levando em conta as distâncias e tamanhos dos clusters, quanto menor este índice melhor a separação entre os clusters. *Ratkowski Index* (RI) que é baseado na média da soma dos quadrados dos dados entre clusters e a soma total dos quadrados de cada dado dentro de um cluster, dentre as quantidades calculadas escolhe-se a que obtém um maior índice. *Hubert Index* (HI) que é um método visual que indica a quantidade preferida através de um pico no gráfico e é calculado pelo coeficiente de correlação entre matrizes de distância. *Ball-Hall Index* (BHI) definido pela média da distância dos itens com os respectivos centroides do cluster, onde no gráfico o ponto de quantidade de clusters com maior diferença do anterior é sugerido. *Krzanowski-Lai Index* (KLI), que propõe índices internos definidos pelas diferenças entre matrizes de dispersão, e aponta a melhor quantidade de clusters pelo maior número gerado ao realizar a equação com quantidade k. Constatou-se que a maioria deles indicou que o melhor número de clusters seria 3, com base nas condições de interpretação de cada índice explicadas anteriormente.

Segundo Gustriansyah, Suhandi e Antony (2020) nos testes para verificar os clusters gerados, utilizou-se a equação de variância (R). Sendo “R” o valor da divisão entre a distância média dos dados cluster (distância intra-cluster) pela distância média dos dados em outros clusters (inter-cluster). O valor médio alcançado para R foi de 0.19113, sendo que quanto mais próximo de zero, maior a similaridade entre os membros dentro de cada cluster.

Gustriansyah, Suhandi e Antony (2020) utilizaram o software R Programming para gerar a clusterização, resultando na visualização demonstrada na Figura 7, tendo dois clusters com densidade maior e um cluster mais disperso. É possível observar também que o cluster em verde possui uma maior variância entre os próprios dados, enquanto os outros dois clusters possuem uma menor diferença interna.

Figura 4 – Três clusters gerados



Fonte: Gustriansyah, Suhandi e Antony (2020).

Segundo Gustriansyah, Suhandi e Antony (2020) também foram adquiridos os valores médios RFM para cada cluster, conforme apresenta a Tabela 1, sendo que o cluster 3 possui a maior média dos três atributos, e o cluster 1 possui a menor média dos três. É possível identificar um intervalo entre os valores médios de cada atributo, indicando uma diferença significativa inter-cluster.

Tabela 3 – Valores médios de RFM de cada cluster

| Cluster | Recency  | Frequency  | Monetary<br>(in thousands) |
|---------|----------|------------|----------------------------|
| 1       | 75.8167  | 3,436.744  | 3,089,608                  |
| 2       | 224.3947 | 13,013.333 | 76,920,847                 |
| 3       | 331.9681 | 107.418    | 286,927,000                |

Fonte: Gustriansyah, Suhandi e Antony (2020).

Gustriansyah, Suhandi e Antony (2020) concluem que o método gerou clusters com alta similaridade em relação aos dados existentes, apresentado uma segmentação mais objetiva quando comparado ao modelo RFM tradicional no qual os dados são divididos igualmente em cinco segmentos (20% dos dados para a cada segmento). Além disso, os autores também sugerem como extensões a utilização de outros métodos para a comparação, como *Particle Swarm Optimization* (PSO), que é um método computacional que otimiza soluções para uma equação de uma certa medida de qualidade, medido de (centroide que são parte do conjunto de dados) ou *maximizing-expectancy*, que é um método iterativo para encontrar estimativas de parâmetros para modelos estatísticos com variáveis não observadas.

## 2.2 LRFMP MODEL FOR CUSTOMER SEGMENTATION IN THE GROCERY RETAIL INDUSTRY: A CASE STUDY

Peker, Kocyigit e Eren (2017) propuseram o modelo *Length, Recency, Frequency, Monetary* (LRFMP) denominado *Length, Recency, Frequency, Monetary and Periodicity* (LRFMP) para classificar dados reais de 16.024 clientes de mercados de uma franquia na Turquia. Para isso, utilizou-se o algoritmo K-means para segmentar os clientes e três índices de validação de clusters para a otimização das suas quantidades, *Silhouette Index* (SI), *Calinski-Harabasz Index* (CHI) e *Davies-Bouldin Index* (DBI). Após a segmentação dos dados, verificou-se estratégias de gerenciamento e relações com os clientes para aumentar a lucratividade, como tratamento preferencial para clientes importantes, implementação de cartões fidelidade para aumentar a frequência de compra de clientes não costumam comprar com frequência, promoções voltadas para clientes incertos com sua escolha de local de compra, dentre outras estratégias.

Primeiramente, Peker, Kocyigit e Eren (2017) adaptaram o modelo LRFM, incluindo o parâmetro *periodicity* (periodicidade), pois a análise dos dados foi realizada a partir do histórico de compras em supermercados, que são estabelecimentos com alto número de visitas, tornando importante a regularidade nos padrões de visita e compra. Peker, Kocyigit e Eren (2017) definem a periodicidade como a regularidade das visitas de um determinado cliente, e é definida pelo desvio padrão dos seus tempos inter-visita (quantia de dias entre duas visitas consecutivas). Se um cliente possui valores baixos de periodicidade, significa que este realiza visitas ou compras em intervalos fixos, podendo caracterizá-lo como cliente regular. Além disso, os autores também modificaram o atributo de recência, transformando-o na média das diferenças entre a data das três últimas compras e a data atual, ao invés da simples diferença entre a data da última compra e a data atual estabelecida no modelo RFM padrão.

Após adquirir os atributos LRFMP dos dados transacionais dos clientes, Peker, Kocyigit e Eren (2017) aplicaram um método de normalização simples nos dados, considerando o intervalo de 0 e 1. Esta normalização foi feita pois os valores LRFMP variam em relação ao intervalo e escala, fato que poderia afetar negativamente a análise dos clusters.

Antes de aplicar a clusterização, Peker, Kocyigit e Eren (2017) utilizaram três índices para validação da quantidade possível de clusters: *Silhouette Index* (SI) que resulta em uma nota de -1 a 1 que indica o quão adequada é a classificação de um objeto dentro de um cluster em comparação aos outros, quanto maior o valor, melhor. *Calinski-Harabasz Index* (CHI) que mede a adequação da quantidade de clusters levando em conta a dispersão entre e intra clusters, um valor alto é preferido. *Davies-Bouldin Index* (DBI) que calcula as similaridades entre clusters levando em conta as distâncias e tamanhos dos clusters, quanto menor este índice melhor será a separação entre os clusters. A partir deles, Peker, Kocyigit e Eren (2017) executaram o algoritmo K-means variando o k de 2 a 9, e os resultados destas iterações foram avaliadas utilizando os três índices. Com base nos resultados, decidiu-se utilizar um número de 5 clusters, pois 2 dos 3 índices sugerem 5 como sendo a quantidade ideal.

Peker, Kocyigit e Eren (2017) utilizaram uma base de dados de uma franquia de mercados que possui mais de dez lojas na cidade de Antália na Turquia. Os dados são compostos por cerca de dois milhões de transações de 16.024 clientes num período de dois anos. Foram removidos os clientes com menos que três compras. Além disso, os autores removeram dados duplicados, transações com valores faltantes assim como, agregaram as compras dentro de um mesmo dia. Depois dessas operações, a quantidade de clientes caiu para 10.471, sendo aplicado na sequência o K-means. A Tabela 6 demonstra a quantidade de clientes nos clusters, os valores médios de LRFMP para cada cluster. Já na última coluna, aplicou-se uma técnica no qual o atributo do cluster recebe uma seta para cima (↑) caso o seu valor for maior que a média do atributo dos outros clusters, e uma seta para baixo (↓), caso seu valor for menor que a média.

Tabela 4 – Valores médios dos clusters

| Cluster | Sample size | Average <i>L</i> | Average <i>R</i> | Average <i>F</i> | Average <i>M</i> | Average <i>P</i> | LRFMP Scores   |
|---------|-------------|------------------|------------------|------------------|------------------|------------------|--|
| 1       | 538         | 633.29           | 39.67            | 175.24           | 24.32            | 4.99             | <i>L</i> ↑ <i>R</i> ↓ <i>F</i> ↑ <i>M</i> ↓ <i>P</i> ↓ |
| 2       | 4,681       | 564.50           | 90.19            | 33.44            | 31.73            | 31.49            | <i>L</i> ↑ <i>R</i> ↓ <i>F</i> ↑ <i>M</i> ↓ <i>P</i> ↓ |
| 3       | 1,091       | 482.17           | 301.18           | 5.33             | 34.21            | 159.32           | <i>L</i> ↑ <i>R</i> ↑ <i>F</i> ↓ <i>M</i> ↓ <i>P</i> ↑ |
| 4       | 818         | 374.01           | 220.24           | 11.85            | 104.18           | 45.81            | <i>L</i> ↓ <i>R</i> ↑ <i>F</i> ↓ <i>M</i> ↑ <i>P</i> ↑ |
| 5       | 3,343       | 173.70           | 399.79           | 10.34            | 30.14            | 27.85            | <i>L</i> ↓ <i>R</i> ↑ <i>F</i> ↓ <i>M</i> ↓ <i>P</i> ↓ |
| Average |             | 419.81           | 218.59           | 28.74            | 36.76            | 43.41            |  |

Fonte: Peker, Kocyigit e Eren (2017).

A partir destes resultados, Peker, Kocyigit e Eren (2017) descreveram as características dos grupos. O grupo 1 representa clientes leais de alta contribuição que, apesar de comporem a menor parcela dos clientes (5,14%), possuem a maior contribuição total entre os grupos. Também é possível observar, que este grupo possui a menor periodicidade média de todos, caracterizando estes clientes como regulares. O grupo 2, representando a maior parcela dos clientes (44,70%) foi classificado como clientes leais de baixa contribuição pois apesar de visitar mais frequentemente as lojas, não possuem tanta contribuição quanto o grupo 1. O grupo 3, com tamanho de 10,42%, foi classificado como clientes incertos, pois possui o atributo de longevidade alto e recência também alta, significando que são clientes com longa história de compra, porém sem muitas compras recentes. Vale notar que este grupo possui o maior valor de periodicidade de todos os grupos, caracterizando-o como um grupo de clientes sem rotina de compra definida. O grupo 4 e 5 foram classificados como clientes perdidos, visto que possuem poucas compras recentes, baixa frequência, e baixa longevidade, denotando um cliente que tem uma pouca interação com a franquia. O grupo 4, contendo uma pequena parcela de 7,81% dos clientes, gasta consideravelmente mais, logo foi classificado como contribuição alta, e o 5, cuja parcela é 31,93%, classificado como contribuição baixa.

A partir desta classificação, Peker, Kocyigit e Eren (2017) estabeleceram estratégias para cada grupo de clientes, como tratamento especial (vagas de estacionamento preferencial, presentes de aniversário, filas preferenciais) para clientes do grupo 1, de maneira a não perder a relação leal com a loja. Para os grupos 3, 4 e 5 cuja frequência é baixa, foi sugerida a adoção de programas de cartão fidelidade para aumentar a frequência deles. Para clientes incertos como no grupo 3, aplicou-se descontos e promoções, de maneira a incentivar os clientes, supostamente sensíveis aos preços, a voltar sua atenção à franquia. Para clientes perdidos, sugeriu-se uma análise mais profunda sobre o motivo da perda, como análise de feedback, inferência de motivos, dentre outros.

Por fim, Peker, Kocyigit e Eren (2017) concluem que o estudo contribuiu com a proposta de um novo modelo RFM, que possibilita uma análise mais profunda que o seu modelo original, visto que as características utilizadas e modificadas permitem uma melhor definição do comportamento de cada cliente. Outra contribuição foi a adição do atributo periodicidade (P) no modelo que, ao contrário do modelo RFM padrão, permite identificar se os clientes de um grupo variam em sua rotina de compras. Outra melhoria apontada é a modificação do atributo de recência, que uma vez calculado como uma média, permite uma caracterização mais precisa que o atributo R comumente utilizado. Uma das limitações destacadas pelos autores é a localidade do estudo, pois foi realizado somente com dados originários de uma cidade, sendo que o comportamento de clientes pode variar de acordo com as diferentes localidades onde é feita a análise. A partir disso, sugere-se uma análise mais ampla contemplando outros locais. Outra sugestão feita por Peker, Kocyigit e Eren (2017) é a adição de novos atributos ao modelo, como a quantidade de produtos comprados, quantidade de produtos percebíveis e não percebíveis comprados, a fim de promover uma interpretação mais profunda do comportamento.

### 2.3 CUSTOMER SEGMENTATION AND STRATEGY DEVELOPMENT BASED ON USER BEHAVIOR ANALYSIS, RFM MODEL AND DATA MINING TECHNIQUES: A CASE STUDY

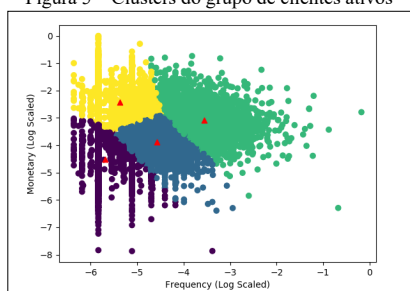
Tavakoli *et al.* (2018) desenvolveram o modelo RFM, denominado “R+FM”, sendo utilizado em conjunto com o algoritmo de clusterização K-means para segmentar 3 milhões de clientes da maior empresa de *E-commerce* do Oriente Médio. Além disso, o modelo de segmentação foi comparado com o utilizado pela empresa, sendo aplicado em uma campanha de Short Message Service (SMS) focada em aumentar os ganhos de cada segmento.

Tavakoli *et al.* (2018) defendem a utilização de um novo modelo de caracterização de clientes, argumentando que o modelo ideal necessita adaptar-se às mudanças de comportamento dos clientes, possuir certa independência de supervisão, levando em consideração a similaridade dos comportamentos dos clientes e a relação entre os atributos Frequência e Monetário. Para isso, construiu-se uma variante do modelo Recency, Monetary, Frequency (RFM) denominado de R+FM, que possui o atributo de recência separado dos demais, utilizando uma segmentação à parte do modelo FM. Os autores separaram os clientes em 3 grupos: os que compraram recentemente (cuja última compra foi dentro de 90 dias), denominados de ativos, clientes que compraram em um passado recente (cuja última compra foi entre 90 e 365 dias), denominados de expirando e por fim, os clientes que não compraram por um longo tempo (cuja última compra foi a mais de 365 dias), denominados de expirados. Para o atributo de frequência, Tavakoli *et al.* (2018) atentaram-se especialmente com a data da primeira compra, pois acreditam que a frequência tem uma importância maior conforme sua recência. Logo, definiram a frequência como a quantidade de compras dividida pela quantidade de dias desde a primeira compra, utilizando também uma função exponencial de decaimento, que efetivamente atribui um peso maior para anos mais recentes, sendo cada ano duas vezes mais pesado que o ano anterior. Como atributo monetário, estabeleceu-se a média dos valores das compras de um cliente, visto que um valor de soma total de compras, segundo os autores, estaria encorajando duas vezes os clientes.

Após a definição do modelo, Tavakoli *et al.* (2018), balancearam a relação entre frequência e monetário, criando a quarta característica que é definida pela combinação linear dos dois atributos, que nada mais é que a soma de cada atributo ponderada pelo peso de cada um. No tratamento dos dados, utilizou-se a técnica de remoção de *outliers* (clientes que não se encaixam no padrão normal) que não se encontram dentro dos intervalos interquartis, que são os intervalos que possuem os dados que pertencem à tendência média do conjunto de dados em geral. Também foram escalados os atributos de frequência e monetário para que seus intervalos sejam iguais, sendo aplicada a normalização *min-max*, que transforma os valores para estarem dentro do intervalo entre 1 e 0. Como os dados monetários e de frequência tratados possuem uma característica de cauda longa, fenômeno estatístico onde os dados são distribuídos de forma decrescente, foi aplicada uma transformação logarítmica para normalizar a distribuição, visto que a quantidade de valores baixos é muito alta, podendo atrapalhar a análise.

Como existem duas segmentações (R e FM), Tavakoli *et al.* (2018) estabeleceram segmentos FM para cada segmento R, resultando nos seguintes grupos: para clientes ativos existem os grupos de alto valor, médio valor com alto monetário, médio valor com alta frequência e baixo valor. Para clientes que estão expirando existem os grupos de alto, médio e baixo valores, sendo aplicado também para clientes expirados. Estes grupos foram definidos por Tavakoli *et al.* (2018) com ajuda da empresa de *E-commerce* Digikala. A partir disso, aplicou-se o K-means com  $k=4$  para o grupo de clientes ativos,  $k=3$  para o grupo de clientes expirando e  $k=3$  para o grupo de clientes expirados, resultando em um total de 10 clusters. Na Figura 8 são identificados os clusters gerados somente a partir do grupo de clientes ativos, organizados em um gráfico de valor monetário por frequência. Sendo o cluster de cor verde composto pelos clientes de alto valor, o cluster de cor amarela composto pelos clientes de médio valor com alto monetário, o cluster de cor azul composto pelos clientes de médio valor com alta frequência, e por fim, o cluster de cor roxa composto pelos clientes de baixo valor.

Figura 5 – Clusters do grupo de clientes ativos



Fonte: Tavakoli *et al.* (2018).

Após a geração dos grupos, Tavakoli *et al.* (2018) discorrem sobre possíveis estratégias para cada segmento, sugerindo um maior foco em clientes ativos com valor médio e baixo, bem como a manutenção de clientes já valiosos. Os autores também enfatizam a importância em recuperar os clientes do grupo expirando, cuja chance de retorno não é tão baixa quanto o grupo expirado, que por si só requer uma estratégia especial de reengajamento dos clientes à empresa.

Além da elaboração de estratégias, Tavakoli *et al.* (2018) implementaram uma campanha de SMS focada somente no segmento de clientes ativos (recência abaixo de 90 dias), pois a empresa já tinha realizado outras campanhas em clientes ativos anteriormente. Nesta campanha cada cliente foi presenteado com um *Voucher* condizente com o segmento ao qual o cliente pertencia. Para clientes ativos com valor alto foi oferecido um desconto de 10 por cento com um desconto máximo de até \$20, com o objetivo de manter a lealdade destes clientes. Para clientes ativos de valor médio com alto valor monetário, foi oferecido um desconto de 10 por cento com um desconto máximo de até \$10, o valor foi menor pois o objetivo era aumentar frequência de compra destes clientes, que em tese já gastam bastante. Para clientes ativos de valor médio com alta frequência foi oferecido um desconto de 10 por cento com um desconto máximo de até \$20 para vendas que custem mais que \$50 (que é o valor médio gasto por este segmento), incentivando assim uma compra de maior valor. Por fim, para clientes ativos de baixo valor, foi oferecido um desconto de 10 por cento com um desconto máximo de até \$20, com o objetivo de converter estes clientes em mais leais.

Para a análise da campanha, Tavakoli *et al.* (2018) selecionaram aleatoriamente 20% dos clientes de cada segmento para compor um grupo de controle, cujos *Vouchers* não foram enviados. Este grupo de controle foi comparado com os outros grupos da campanha para obter um valor de referência do aumento do valor monetário após sua conclusão. Os resultados alcançados por Tavakoli *et al.* (2018) podem ser observados no Quadro 1, ao qual percebe-se um aumento de \$14,30 na média monetária dos clientes ativos com alta frequência, mais do que o aumento de \$3,20 sofrido pelo grupo de controle. É possível também observar que a média monetária de todos os grupos alvo da campanha aumentou consideravelmente, enquanto o grupo de controle aumentou pouco ou até diminuiu, indicando uma efetividade no objetivo da campanha.

Quadro 1 – Dados da média monetária do grupo de controle e grupo de campanha

| Recency | Segment<br>Monetary and Frequency | Average Monetary (USD) |                |                 |                |
|---------|-----------------------------------|------------------------|----------------|-----------------|----------------|
|         |                                   | Control Users          |                | Campaign Users  |                |
|         |                                   | Before Campaign        | After Campaign | Before Campaign | After Campaign |
| Active  | High Value                        | 74.2                   | 73.8           | 88.2            | 89.2           |
|         | Medium Value with High Monetary   | 100.2                  | 97.6           | 104.6           | 105.2          |
|         | Medium Value with High Frequency  | 32                     | 35.2           | 35.4            | 49.7           |
|         | Low Value                         | 50.7                   | 53.2           | 56.4            | 65.2           |

Fonte: Tavakoli *et al.* (2018).

Tavakoli *et al.* (2018) concluem que houve uma melhora no desempenho da campanha lançada em comparação com as anteriores, indicando ainda, que elas obtinham uma taxa de compra de 0,1 por cento, sendo que a campanha lançada para validação do modelo obteve uma taxa de 1 por cento, cerca de dez vezes mais efetivo. Os autores justificam esta melhora ao processo de segmentação do modelo, resultando em clusters mais significativos, facilitando a aplicação de vouchers específicos. Por fim, Tavakoli *et al.* (2018) sugerem o melhoramento da definição do atributo de recência de forma que seja mais útil ao time de marketing. Também recomendam calcular o Customer Lifetime Value (CLV), que atribui o valor vitalício à cada segmento e cliente, de forma a quantificar o valor que um cliente pode proporcionar à empresa.

### 3 PROPOSTA DO PROTÓTIPO

Essa seção visa apresentar a justificativa para a elaboração deste trabalho, os requisitos que serão seguidos e a metodologia que será utilizada.

#### 3.1 JUSTIFICATIVA

No Quadro 2 é apresentado um comparativo entre os trabalhos correlatos. As linhas representam as características relevantes e as colunas representam os trabalhos.

Quadro 2 – Comparativo entre os trabalhos correlatos

| Características                           | Correlatos<br>Gustriansyah,<br>Suhandi e<br>Antony (2020) | Peker, Kocyigit e<br>Eren (2017)                      | Tavakoli <i>et al.</i><br>(2018)                   |
|---|---|---|--|
| Alvo da clusterização                     | Produtos  | Clientes  | Clientes   |
| Modelo utilizado                          | RFM   | LRFMP   | R+FM   |
| Objetivo da segmentação                   | Gerenciamento de estoque                                  | Gerenciamento das relações com cliente                | Gerenciamento das relações com cliente             |
| Algoritmo de clusterização utilizado      | K-means   | K-means   | K-means  |
| Foco metodológico                         | Otimização de k com diferentes métricas                   | Formulação de um modelo novo e análise dos resultados | Formulação de um modelo novo e campanha de ofertas |
| Número de dados (clientes/produtos)       | 2.043   | 16.024  | ~3.000.000   |
| Quantidade de índices para validação de k | 8   | 3   | -  |
| Quantidade de clusters gerados            | 3   | 5   | 10   |
| Inferências sobre os dados                | -   | Sim   | Sim  |

Fonte: elaborado pelo autor.

A partir do Quadro 2, pode-se observar que Gustriansyah, Suhandi e Antony (2020) agruparam produtos de uma base de dados utilizando o modelo RFM padrão. Já Peker, Kocyigit e Eren (2017) optaram pelo desenvolvimento de um modelo novo, considerando a periodicidade (LRFMP). Tavakoli *et al.* (2018) também desenvolveram um novo modelo, ao qual a característica recência foi modificada e separada (R+FM).

Gustriansyah, Suhandi e Antony (2020) tinham como objetivo melhorar o gerenciamento de estoque, prezando por uma segmentação mais conclusiva sobre os produtos, visto que o modelo RFM padrão define segmentos arbitrariamente sem adequar-se às peculiaridades dos dados, enquanto o modelo aplicado através de k-means alcançou uma segmentação com dados altamente similares em cada cluster. Por outro lado, Peker, Kocyigit e Eren (2017) e Tavakoli *et al.* (2018) objetivavam o gerenciamento das relações com os clientes através de estratégias focadas em segmentos, visando aumentar a renda que eles fornecem à empresa. Todos os autores utilizaram o algoritmo K-means, por ser confiável e amplamente difundido. Vale ressaltar que no trabalho de Gustriansyah, Suhandi e Antony (2020), o algoritmo teve um foco metodológico maior, visto que foram utilizados 8 índices de validação para k clusters, visando otimizar a organização dos segmentos.

A quantidade de dados segmentados variou bastante entre os três trabalhos devido aos diferentes contextos de aplicação. Gustriansyah, Suhandi e Antony (2020) tinham 2.043 produtos na base de dados para segmentar, resultando em 3 clusters. Já Peker, Kocyigit e Eren (2017) possuíam o registro de 16.024 clientes de uma rede de padarias, sendo especificados 5 segmentos, obtidos através de uma análise por três índices de validação (Silhouette, Calinski-Harabasz e Davies-Bouldin). Por fim, Tavakoli *et al.* (2018) agruparam dados de 3 milhões de clientes pertencentes à base de dados de um e-commerce do Oriente Médio, resultando em 10 clusters, sendo 3 pertencentes à característica de recência, e os outros 7 distribuídos entre as características de frequência e monetária. Ressalta-se que Tavakoli *et al.* (2018) testaram o modelo em produção, montando uma campanha que focava no segmento de clientes ativos, visando primariamente aumentar os lucros da empresa, utilizando também um grupo de controle e comparação de renda antes e depois da campanha.

Gustriansyah, Suhandi e Antony (2020) demonstraram a possibilidade da aplicação de RFM fora do uso convencional de segmentação de clientes, e adquiriram clusters com uma variância média de 0.19113. Além disso, os autores sugeriram outras formas de comparação de dados, como Particle Swarm Optimization (PSO), medioides ou até *maximizing-expectancy*. Peker, Kocyigit e Eren (2017) segmentaram clientes de uma rede de mercados na Turquia em “clientes leais de alta contribuição”, “clientes leais de baixa contribuição”, “clientes incertos”, “clientes perdidos de alto gasto” e “clientes perdidos de baixo gasto”. Desta maneira, os autores providenciaram visões e estratégias (promoções, ofertas, regalias) de aumento de renda sobre os comportamentos dos clientes, porém limitaram-se a aplicar em um segmento específico de mercado. Por fim, Tavakoli *et al.* (2018) agruparam clientes de uma empresa de e-commerce com base em sua recência, resultando em clientes “Ativos”, “Expirando” e “Expirados”, e destes segmentos, sucessivamente separados em grupos de “Alto”, “Médio” e “Baixo” valores, validando posteriormente a segmentação através de uma campanha de ofertas para os clientes do grupo “Ativos”.

Todos os trabalhos aqui citados procuraram implementar o modelo RFM num contexto de clusterização por K-means, alterando o modelo e o manejo dos dados de acordo com cada categoria, seja ele produto ou cliente, varejo ou mercado. Com isso, criaram-se atributos e foram modificados alguns já existentes para atender às especificidades de cada contexto, visto que todos os trabalhos focaram em uma só base de dados, inevitavelmente adequando-se às mesmas.

Desta forma, este trabalho demonstra ser relevante, pois almeja aplicar o modelo RFM em conjunto com vários algoritmos de clusterização para realizar a segmentação a partir de atributos que realcem o comportamento dos clientes, disponibilizando um artefato computacional que se adeque à vários contextos (mercado, comércio, varejo etc.), utilizando várias bases de dados reais para testar a validade dos algoritmos utilizados. Inicialmente o cenário de análise será o varejo e, conforme os resultados alcançados, serão testados em outras bases/cenários. Vislumbra-se utilizar três índices (Silhouette, Calinski-Harabasz e Davies-Bouldin) para validação da qualidade dos clusters em relação à sua separação e homogeneidade. Além disso, deseja-se obter clusters significativos e coerentes com cada segmento de mercado aplicado. Outra contribuição deste trabalho refere-se ao âmbito comercial, com a geração de informações sobre as similaridades de clientes de cada segmento de mercado, podendo auxiliar gestores e administradores de empresas a obter uma visão crítica sobre os comportamentos de clientes ao longo das diferentes bases de dados, podendo também denotar características comuns a todos. Outra relevância seria a utilização deste trabalho em ambiente acadêmico, visto que serão aplicados diferentes algoritmos de clusterização, podendo providenciar informações sobre seus desempenhos e qualidade de agrupamento, além de ser aplicados processos de obtenção, limpeza e transformação de dados.

### 3.2 REQUISITOS PRINCIPAIS DO PROBLEMA A SER TRABALHADO

O artefato computacional a ser desenvolvido deverá:

- a) adquirir os dados transacionais de clientes a partir de um banco de dados (Requisito Funcional - RF);
- b) extrair dos clientes as características (recência, frequência e monetária) utilizadas no modelo RFM (RF);
- c) filtrar os clientes sem quantidade de compras relevantes (RF);
- d) normalizar os dados para evitar disparidades nas escalas dos dados, principalmente no atributo monetário (RF);
- e) aplicar três índices de validação (Silhouette, Calinski-Harabasz e Davies-Bouldin) para verificar a qualidade dos clusters (RF);
- f) apresentar em um gráfico 3D os clientes, com sua localização definida pela pontuação do cliente nas características RFM (RF);
- g) utilizar algoritmos de clusterização tais como hierárquicos, K-means e *Density-Based Spatial*

- Clustering of Applications With Noise* (DBSCAN) (RF);
- h) utilizar a linguagem Python para o desenvolvimento (Requisito Não Funcional - RNF);
  - i) utilizar o ambiente de desenvolvimento Jupyter Notebook (RNF);
  - j) utilizar o banco de dados PostgreSQL para ler os dados das bases utilizadas (RNF).

### 3.3 METODOLOGIA

O trabalho será desenvolvido observando as seguintes etapas:

- a) levantamento bibliográfico: pesquisar trabalhos relacionados e estudar sobre o modelo RFM e suas aplicações, algoritmos de clusterização, métodos de tratamento de dados e índices de validação;
- b) seleção de bases de dados: obter bases de dados de usuários cedidas pela empresa Intelidata Informática, desenvolvedora de software de gestão comercial. Serão selecionadas conforme sua adequação ao objetivo do trabalho, variando em tamanho e segmento de mercado;
- c) definição das características do modelo RFM: definir os atributos utilizados para caracterizar os clientes no modelo RFM;
- d) definição de métricas do modelo RFM: definir as métricas para mensuração e atribuição de pontuação de cada característica no modelo RFM;
- e) definição dos algoritmos de clusterização: pesquisar e escolher o algoritmo de clusterização que realizará o agrupamento das características RFM;
- f) implementação: implementar o artefato computacional de segmentação levando em consideração as etapas (b) até (e), utilizando a linguagem Python;
- g) análise dos clusters: avaliar a qualidade dos clusters gerados a partir dos diferentes algoritmos de clusterização e seus comportamentos em múltiplas bases de dados, aplicando índices de validação e apresentando os agrupamentos na forma de gráficos, de maneira que a sua separação e homogeneidade possam ser observadas.

As etapas serão realizadas nos períodos relacionados no Quadro 3.

Quadro 3 – Cronograma de atividades a serem realizadas

| etapas / quinzenas                          | 2022 |   |      |   |      |   |      |   |      |   |
|---|------|---|------|---|------|---|------|---|------|---|
|   | fev. |   | mar. |   | abr. |   | maio |   | jun. |   |
|   | 1    | 2 | 1    | 2 | 1    | 2 | 1    | 2 | 1    | 2 |
| levantamento bibliográfico                  |      |   |      |   |      |   |      |   |      |   |
| seleção de bases de dados                   |      |   |      |   |      |   |      |   |      |   |
| definição das características do modelo RFM |      |   |      |   |      |   |      |   |      |   |
| definição de métricas do modelo RFM         |      |   |      |   |      |   |      |   |      |   |
| definição dos algoritmos de clusterização   |      |   |      |   |      |   |      |   |      |   |
| implementação                               |      |   |      |   |      |   |      |   |      |   |
| análise dos clusters                        |      |   |      |   |      |   |      |   |      |   |

Fonte: elaborado pelo autor.

## 4 REVISÃO BIBLIOGRÁFICA

Nesta seção são abordados os assuntos que servirão de base para a realização deste trabalho. A seção 4.1 discorre sobre clustering. A seção 4.2 apresenta o modelo RFM. Por fim, a seção 4.3 aborda os índices de validação de clusters.

### 4.1 CLUSTERING

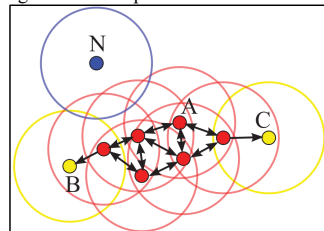
Para Cherkassky e Mulier (2007), clustering se trata do problema de separar um conjunto de dados em grupos chamados de “clusters” baseado em alguma medida de similaridade. O objetivo é encontrar um conjunto de clusters dos quais as amostras dentro dos mesmos são mais similares entre si do que quando comparadas com amostras de outros clusters. A análise destes clusters consiste no fato de que a medida de similaridade entre eles é escolhida subjetivamente baseado na sua habilidade de criar clusters interessantes ao analista. Cherkassky e Mulier (2007) classificam os algoritmos em dois tipos principais: (i) Hierárquicos, seguindo uma estrutura de árvores; (ii) Particionais, que geram clusters a partir de sucessivas secções, cujos métodos são identificados em dois grupos: particionais onde cada dado é atribuído à um e somente à um cluster e particionais que podem pertencer à vários clusters.

Segundo Schubert *et al.* (2017), o algoritmo *Density-Based Spatial Clustering of Applications With Noise* (DBSCAN) foi publicado em 1996, sendo comumente utilizado com sucesso em várias aplicações do mundo real. Schubert *et al.* (2017) apontam que o modelo DBSCAN utiliza uma estimativa de densidade mínima simples, baseada no número mínimo de vizinhos (*minPts*) dentro de um raio  $\epsilon$ , cuja medida de



distância é arbitrária. A partir disso, objetos com mais que o número mínimo de vizinhos dentro do raio (incluindo o ponto analisado) são considerados como pontos centrais. Todos os vizinhos de um ponto central que estejam dentro do raio  $\epsilon$  são considerados parte do mesmo cluster que o ponto central. Se qualquer um desses vizinhos for também um ponto central, suas vizinhanças são incluídas no cluster. Além disso, pontos não centrais deste cluster são chamados de pontos-limite. Pontos que não são alcançáveis por qualquer ponto central são considerados ruído e não pertencem a nenhum cluster. A Figura 3 exemplifica o funcionamento do DBSCAN. O parâmetro  $minPts$  é 4, o raio  $\epsilon$  é denotado pelos círculos, N é um ponto ruído, A são pontos centrais e B e C são pontos-limite.

Figura 6 – Exemplo do modelo DBSCAN



Fonte: Schubert *et al.* (2017).

Zhao *et al.* (2005) ressaltam que a maioria dos algoritmos de *Hierarchical Clustering* são aglomerativos, cujos objetos são inicialmente designados para um cluster. A partir disso, clusters são repetidamente combinados até que a árvore de clusters total seja formada, ao qual seleciona-se uma altura de corte para obter os clusters desejados. Apesar disso, os autores descrevem que algoritmos particionais também podem ser utilizados para obter soluções hierárquicas através de sequências de bissecções repetidas. Segundo Zhao *et al.* (2005), na abordagem hierárquica aglomerativa, o parâmetro-chave é o método usado para determinar o par de clusters a serem unidos a cada execução do algoritmo. Na maioria das abordagens o par mais similar é selecionado, podendo ser feito de várias maneiras tradicionais como a *single-link*, *complete-link* e média grupal (ou *Unweighted Pair Group Method with Arithmetic mean - UPGMA*). O esquema *single-link* mede a similaridade de dois clusters através dos pontos mais próximos entre eles. Já o *complete-link* mede inversamente, de maneira que a similaridade é obtida através dos pontos mais distantes entre os dois clusters. Por fim, o método UPGMA leva em consideração a distância média entre elementos de cada cluster.

De acordo com Ghosh e Kumar (2013), o K-means é um método de particionamento mais aplicado para analisar dados, separando os objetos em clusters mutuamente exclusivos (K) de maneira que os objetos de cada cluster fiquem tão perto entre si quanto possível, porém tão longe quanto possível de objetos em clusters diferentes. Cada cluster possui um ponto central (centroide), cuja localização é obtida através da média da localização de todos os pontos pertencentes ao cluster. Segundo Ghosh e Kumar (2013), o algoritmo baseia-se na constante atualização da posição dos centroides e recálculo dos pontos mais próximos, sendo que inicialmente os k-centroides são aleatoriamente distribuídos no espaço. O algoritmo acaba sua execução quando nenhum ponto muda de cluster ou nenhum centroide se move.

#### 4.2 MODELO RFM

Segundo Hughes (2011), o modelo RFM é “Um meio antigo e altamente preditivo de determinar quem irá responder e comprar. Um método de codificar clientes existentes. Usado para prever resposta, tamanho médio de pedido, e outros fatores”. Este modelo categoriza geralmente clientes através das características de Recência (R), Frequência (F) e Monetária (M). As métricas utilizadas para medir tais características podem variar, porém geralmente classificam recência como a quantidade de dias desde a última compra, frequência como a quantidade de compras dentro de um determinado período, e monetária como o total acumulado de todas as vendas realizadas para um cliente.

No estudo realizado por Verhoef *et al.* (2003), cerca de 90% das empresas questionadas sobre a aplicação métodos de segmentação (como o RFM) afirmam que possuíam como objetivo a seleção de alvos, ou seja, encontrar o segmento de clientes que mais se identificam com a empresa, e 64,4% citaram como objetivo o tratamento diferencial de clientes, com promoções, preços e ofertas especiais. Verhoef *et al.* (2003) ainda evidenciam que este modelo é geralmente utilizado em marketing direto, onde a comunicação acontece diretamente entre a empresa e o consumidor, realizada através de mídias sociais, e-mail, mensagens SMS ou até pelo correio.

Comentado [DSdR6]: Fonte courier new

Comentado [DSdR7]: Fonte courier new

Comentado [DSdR8]: Fonte courier new

Christy *et al.* (2021) indicam que o propósito da pontuação RFM é projetar comportamentos futuros, ajudando nas decisões de segmentação posteriores. Para realizar essa projeção, é importante transformar este comportamento em números que podem ser utilizados ao longo do tempo. Um método comum de atribuição de pontuação seria utilizando quintis, sendo este um dos métodos mais utilizados. Ele consiste em organizar os clientes em ordem decrescente (melhor para pior), dividi-los em 5 grupos de tamanhos iguais, os melhores recebem a pontuação de 5 e os piores 1. Os melhores clientes então, teriam a pontuação total 5,5,5.

Segundo Tsoy e Shchekoldin (2016), o método de quintis possui alguns problemas, sendo um deles a dificuldade de detectar nuances dentro de cada quintil, muitas vezes agrupando clientes com comportamentos muito diferentes. Isto se aplica geralmente nos quintis do topo, onde há um número bem pequeno de clientes muito importantes, e um número grande de clientes menos importantes, deixando claro uma heterogeneidade presente dentro do próprio grupo em questão. Outras vezes, nos quintis baixos, os autores afirmam que este método tende arbitrariamente a separar clientes com comportamentos parecidos.

De acordo com Safari, Safari e Montazer (2016), um dos processos mais importantes da aplicação do modelo RFM é a designação do peso dos atributos seguida da ponderação, onde é multiplicado o peso desejado pelo próprio valor do atributo. As fórmulas mais utilizadas tendem a priorizar recência acima de todos, seguido de frequência e por último monetário. Desta maneira, elevando clientes recentes à uma prioridade maior que os outros. Segundo os autores, esta designação pode variar de acordo com o contexto de cada empresa que emprega o modelo RFM, requerendo conhecimento prévio do ramo de negócios do estabelecimento para a correta atribuição de pesos.

#### 4.3 ÍNDICES DE VALIDAÇÃO

Para Hämäläinen, Jauhiainen e Kärkkäinen (2017), uma validação de clusters considera a qualidade do resultado de um algoritmo de cluster, tentando achar a separação que melhor se encaixa com a natureza dos dados. O número de clusters dado como parâmetro para vários algoritmos deveria ser decidido com base na estrutura natural dos dados, porém não há solução clara sobre o melhor número de clusters.

Hämäläinen, Jauhiainen e Kärkkäinen (2017) destacam que os índices de validação medem a tangibilidade do objetivo, sendo esta, uma alta similaridade entre os dados dentro de um cluster e uma alta diferença entre os dados de clusters diferentes. Essas medidas são chamadas de separação intra e inter cluster, respectivamente. Uma boa medida de separação intra-cluster possui números baixos, já uma boa medida de separação inter cluster possui números altos. Os autores ainda destacam que nenhum índice de validação é perfeito para qualquer contexto, alguns índices são mais adequados à diferentes tipos de dados, logo recomenda-se utilizar múltiplos índices para a análise de clusters. Arbelaitz *et al.* (2013) concluem em seu estudo através de uma análise estatística, que dos 30 índices pesquisados, 10 provam ser recomendáveis para utilização. No topo desta lista, encontram-se os índices Silhouette, Calinski-Harabasz e Davies-Bouldin.

De acordo com Rousseeuw (1987), para gerar o índice de Silhouette de um dado, são necessárias apenas duas coisas: os clusters obtidos e o conjunto das distâncias entre todos os dados observados, sendo calculado para cada  $i$  o seu respectivo índice Silhouette  $s(i)$ . Calcula-se também a média de dissimilaridade das distâncias de  $i$  com o resto dos dados do cluster de  $i$ , denotado por  $a(i)$ . No passo seguinte, é obtido o valor mínimo entre as distâncias de  $i$  e qualquer outro cluster (é descoberto então o cluster vizinho de  $i$ , ou seja, o cluster com que  $i$  mais se encaixaria caso não estivesse em seu cluster original), denotado por  $b(i)$ . Este processo pode ser resumido pela Equação 1, que resulta em um número entre -1 e 1, sendo -1 uma categorização ruim do objeto  $i$  (não condizente com seu cluster atual) e sendo 1 uma categorização ótima. Para obter a qualidade da clusterização em geral, é obtida a média de  $s(i)$  para todos os objetos  $i$  do conjunto de dados.

$$s(i) = \frac{b(i) - a(i)}{\max\{a(i), b(i)\}} \quad (1)$$

Para Calinski e Harabasz (1974), o seu índice Variance Rate Criterion (VRC), demonstrado na Equação 2, considera: a soma dos quadrados entre grupos (Between Group Sum of Squares - BGSS) que retrata a variância entre clusters levando em conta a distância de seus centroides até o centroide global; a soma dos quadrados dentro grupos (Within Group Sum of Squares - WGSS) que retrata a variância dentro de clusters levando em conta as distâncias dos pontos em um cluster até o seu centroide. Considera-se também o número de observações/dados ( $n$ ) e o número de clusters ( $k$ ). Quando este índice é utilizado, procura-se maximizar o resultado conforme o valor de  $k$  é aumentado.

$$VRC = \frac{BGSS}{k-1} / \frac{WGSS}{n-k} \quad (2)$$

Davies e Bouldin (1979) denotam que o objetivo de seu índice é definir uma medida de separação de clusters  $R(S_i, S_j, M_{ij})$  que permita a computação da similaridade média de cada cluster com o seu cluster mais similar (vizinho), o valor mais baixo possível seria o resultado ideal. Com  $S_i$  sendo a medida de dispersão do cluster  $i$ ,  $S_j$  sendo a medida de dispersão do cluster  $j$  e  $M_{ij}$  sendo a distância entre os clusters  $i$  e  $j$ , conforme Equação 3.

$$R_{ij} \equiv \frac{S_i + S_j}{M_{ij}} \quad R \equiv \frac{1}{N} \sum_{i=1}^N R_i \quad (3)$$

De acordo com Davies e Bouldin (1979), primeiro obtêm-se  $R_{ij}$  de todos os clusters, isto é, a razão de distâncias inter e intra-cluster entre o cluster  $i$  e  $j$ . Após isso, obtêm-se  $R_i$  (o valor mais alto de  $R_{ij}$ ) identificando para cada cluster, o cluster vizinho ao qual ele mais se assemelha. Por fim, é calculado o índice em si ( $R$ ), sendo este, a soma total das similaridades de  $n$  clusters com seus vizinhos mais próximos.

### REFERÊNCIAS

- ARBELAITZ, Olatz *et al.* An extensive comparative study of cluster validity indices. **Pattern Recognition**, [S.l.], v. 46, n. 1, p. 243-256, jan. 2013. Elsevier BV. <http://dx.doi.org/10.1016/j.patcog.2012.07.021>.
- CALINSKI, T.; HARABASZ, J. A dendrite method for cluster analysis. **Communications In Statistics - Theory And Methods**, [S.l.], v. 3, n. 1, p. 1-27, 1974. Informa UK Limited. <http://dx.doi.org/10.1080/03610927408827101>.
- CHERKASSKY, Vladimir S.; MULIER, Filip. Methods for data reduction and dimensionality reduction. In: CHERKASSKY, Vladimir S.; MULIER, Filip. **Learning from data: concepts, theory, and methods**. 2. ed. Hoboken: Ieee Press, 2007. Cap. 6, p. 191
- CHRISTY, A. Joy *et al.* RFM ranking – An effective approach to customer segmentation. **Journal Of King Saud University - Computer And Information Sciences**, [S.l.], v. 33, n. 10, p. 1251-1257, dez. 2021. Elsevier BV. <http://dx.doi.org/10.1016/j.jksuci.2018.09.004>.
- DAVIES, David L.; BOULDIN, Donald W. A Cluster Separation Measure. **Ieee Transactions On Pattern Analysis And Machine Intelligence**, [S.l.], v. -1, n. 2, p. 224-227, abr. 1979. Institute of Electrical and Electronics Engineers (IEEE). <http://dx.doi.org/10.1109/tpami.1979.4766909>.
- GHOSH, Soumi; KUMAR, Sanjay. Comparative Analysis of K-Means and Fuzzy C-Means Algorithms. **International Journal Of Advanced Computer Science And Applications**, [S.l.], v. 4, n. 4, p. 35-39, 2013. The Science and Information Organization. <http://dx.doi.org/10.14569/ijacsa.2013.040406>.
- GUSTRIANSYAH, Rendra; SUHANDI, Nazori; ANTONY, Fery. Clustering optimization in RFM analysis Based on k-Means. **Indonesian Journal Of Electrical Engineering And Computer Science**, [S.l.], v. 18, n. 1, p. 470-477, abr. 2020. Mensal.
- HÄMÄLÄINEN, Joonas; JAUHAINEN, Susanne; KÄRKKÄINEN, Tommi. Comparison of Internal Clustering Validation Indices for Prototype-Based Clustering. **Algorithms**, [S.l.], v. 10, n. 3, p. 105-119, 6 set. 2017. MDPI AG. <http://dx.doi.org/10.3390/a10030105>.
- HUGHES, Arthur M. **Strategic Database Marketing 4e: the masterplan for starting and managing a profitable, customer-based marketing program**. 4. ed. [S.l.]: McGraw-Hill, 2011. 608 p.
- KUMAR, Vijay. **Managing Customers for Profit: strategies to increase profits and build loyalty**. Upper Saddle River: Pearson Prentice Hall, 2008. 296 p.
- NGUYEN, Thuyuyen H.; SHERIF, Joseph S.; NEWBY, Michael. **Strategies for successful CRM implementation**. Information Management & Computer Security, [S.l.], v. 15, n. 2, p. 102-115, maio 2007.
- PEKER, Serhat; KOCYIGIT, Altan; EREN, P. Erhan. LRFMP model for customer segmentation in the grocery retail industry: a case study. **Marketing Intelligence & Planning**, [S.l.], v. 35, n. 4, p. 544-559, 6 maio 2017. Emerald. <http://dx.doi.org/10.1108/mip-11-2016-0210>.
- PETRISON, Lisa A.; BLATTBERG, Robert C.; WANG, Paul. Database marketing: past, present, and future. **Journal Of Direct Marketing**, [S.l.], v. 11, n. 4, p. 109-125, mar. 1997. Wiley. [http://dx.doi.org/10.1002/\(sici\)1522-7138\(199723\)11:43.0.co;2-g](http://dx.doi.org/10.1002/(sici)1522-7138(199723)11:43.0.co;2-g).
- RASHID, Mohammad A.; HOSSAIN, Liaquat; PATRICK, Jon David. The Evolution of ERP Systems: a historical perspective. In: NAH, Fiona Fui-Hoon. **Enterprise Resource Planning: solutions and management**. Hershey: Irm Press, 2001. p. 35-50.
- REINARTZ, Werner; THOMAS, Jacquelyn S.; KUMAR, V. Balancing Acquisition and Retention Resources to Maximize Customer Profitability. **Journal Of Marketing**, [S.l.], v. 69, n. 1, p. 63-79, jan. 2005.

ROUSSEEUW, Peter J. Silhouettes: a graphical aid to the interpretation and validation of cluster analysis. **Journal Of Computational And Applied Mathematics**, [S.l.], v. 20, n. 0, p. 53-65, nov. 1987. Elsevier BV. [http://dx.doi.org/10.1016/0377-0427\(87\)90125-7](http://dx.doi.org/10.1016/0377-0427(87)90125-7).

SAFARI, Fariba; SAFARI, Narges; MONTAZER, Gholam Ali. Customer lifetime value determination based on RFM model. **Marketing Intelligence & Planning**, [S.l.], v. 34, n. 4, p. 446-461, 6 jun. 2016. Emerald. <http://dx.doi.org/10.1108/mip-03-2015-0060>.

SCHUBERT, Erich *et al.* DBSCAN Revisited, Revisited: why and how you should (still) use DBSCAN. **Acm Transactions On Database Systems**, [S.l.], v. 42, n. 3, p. 1-21, 24 ago. 2017. Association for Computing Machinery (ACM). <http://dx.doi.org/10.1145/3068335>.

TAVAKOLI, Mohammadreza *et al.* Customer Segmentation and Strategy Development Based on User Behavior Analysis, RFM Model and Data Mining Techniques: a case study. In: 2018 IEEE 15TH INTERNATIONAL CONFERENCE ON E-BUSINESS ENGINEERING (ICEBE), 15., 2018, Xiam. **Proceedings [...]**. [S.l.]: Ieee, 2018. p. 119-126.

TSIPTSIS, Konstantinos K.; CHORIANOPOULOS, Antonios. **Data Mining Techniques in CRM: inside customer segmentation**. Chichester: John Wiley & Sons, 2009. 374 p.

TSOY, Marina E.; SHCHEKOLDIN, Vladislav Yu. RFM-analysis as a tool for segmentation of high-tech products' consumers. In: 2016 13TH INTERNATIONAL SCIENTIFIC-TECHNICAL CONFERENCE ON ACTUAL PROBLEMS OF ELECTRONICS INSTRUMENT ENGINEERING (APEIE), 13., 2016, Novosibirsk. **2016 13th International Scientific-Technical Conference on Actual Problems of Electronics Instrument Engineering (APEIE)**. [S.l.]: Ieee, 2016. p. 290-293.

VERHOEF, Peter C *et al.* The commercial use of segmentation and predictive modeling techniques for database marketing in the Netherlands. **Decision Support Systems**, [S.l.], v. 34, n. 4, p. 471-481, mar. 2003.

ZHAO, Ying *et al.* Hierarchical Clustering Algorithms for Document Datasets. **Data Mining And Knowledge Discovery**, [S.l.], v. 10, n. 2, p. 141-168, mar. 2005. Springer Science and Business Media LLC. <http://dx.doi.org/10.1007/s10618-005-0361-3>.

# FORMULÁRIO DE AVALIAÇÃO BCC – PROFESSOR TCC I

Avaliador(a): Dalton Solano dos Reis

| ASPECTOS AVALIADOS     |  | atende | atende parcialmente | não atende |
|------------------------|--|--------|---------------------|------------|
| ASPECTOS TÉCNICOS      | 1. INTRODUÇÃO<br>O tema de pesquisa está devidamente contextualizado/delimitado?   | X      |                     |            |
|                        | O problema está claramente formulado?  | X      |                     |            |
|                        | 2. OBJETIVOS<br>O objetivo principal está claramente definido e é passível de ser alcançado?   | X      |                     |            |
|                        | Os objetivos específicos são coerentes com o objetivo principal?   | X      |                     |            |
|                        | 3. JUSTIFICATIVA<br>São apresentados argumentos científicos, técnicos ou metodológicos que justificam a proposta?  | X      |                     |            |
|                        | São apresentadas as contribuições teóricas, práticas ou sociais que justificam a proposta?   | X      |                     |            |
|                        | 4. METODOLOGIA<br>Foram relacionadas todas as etapas necessárias para o desenvolvimento do TCC?  | X      |                     |            |
|                        | Os métodos, recursos e o cronograma estão devidamente apresentados?  | X      |                     |            |
| ASPECTOS METODOLÓGICOS | 5. REVISÃO BIBLIOGRÁFICA (atenção para a diferença de conteúdo entre projeto e pré-projeto)<br>Os assuntos apresentados são suficientes e têm relação com o tema do TCC? | X      |                     |            |
|                        | 6. LINGUAGEM USADA (redação)<br>O texto completo é coerente e redigido corretamente em língua portuguesa, usando linguagem formal/científica?                            | X      |                     |            |
|                        | A exposição do assunto é ordenada (as ideias estão bem encadeadas e a linguagem utilizada é clara)?  | X      |                     |            |
|                        | 7. ORGANIZAÇÃO E APRESENTAÇÃO GRÁFICA DO TEXTO<br>A organização e apresentação dos capítulos, seções, subseções e parágrafos estão de acordo com o modelo estabelecido?  | X      |                     |            |
|                        | 8. ILUSTRAÇÕES (figuras, quadros, tabelas)<br>As ilustrações são legíveis e obedecem às normas da ABNT?  | X      |                     |            |
|                        | 9. REFERÊNCIAS E CITAÇÕES<br>As referências obedecem às normas da ABNT?  | X      |                     |            |
|                        | As citações obedecem às normas da ABNT?  | X      |                     |            |
|                        | Todos os documentos citados foram referenciados e vice-versa, isto é, as citações e referências são consistentes?  | X      |                     |            |

O projeto de TCC será reprovado se:

- qualquer um dos itens tiver resposta NÃO ATENDE;
- pelo menos 4 (quatro) itens dos **ASPECTOS TÉCNICOS** tiverem resposta ATENDE PARCIALMENTE; ou
- pelo menos 4 (quatro) itens dos **ASPECTOS METODOLÓGICOS** tiverem resposta ATENDE PARCIALMENTE.

**PARECER:** ( X ) APROVADO ( ) REPROVADO

# Revisão do Pré-projeto

**Disciplina: Trabalho de Conclusão de Curso I – BCC**

Caro orientando,

segue abaixo o Termo de Compromisso, as DUAS revisões do seu pré-projeto contendo a avaliação do professor “avaliador” e professor “TCC1”, junto com as avaliações da defesa na banca de qualificação. É muito importante que revise com cuidado e discuta possíveis dúvidas decorrente das revisões com o seu professor orientador, e com o professor de TCC1. Sempre procure fazer todos os ajustes solicitados, até mesmo o menores detalhes, pois todos são importantes e irão refletir na sua nota nesta disciplina. Mas, caso o professor orientador julgue que algumas anotações das revisões não devam ser feitas, ou mesmo que sejam feitas de forma diferente a solicitada pelo revisor, anexe ao final do seu projeto a ficha “Projeto: Observações – Professor Orientador” disponível no material da disciplina, e justifique o motivo.

Lembrem que agora o limite de páginas do projeto é no máximo 12 (doze) páginas. E que a seção de “Revisão Bibliográfica” deve ser complementada.

Atenciosamente,

UNIVERSIDADE REGIONAL DE BLUMENAU  
CENTRO DE CIÊNCIAS EXATAS E NATURAIS  
CURSO DE CIÊNCIA DA COMPUTAÇÃO – BACHARELADO

**TERMO DE COMPROMISSO**

|   |  |
|---|--|
| <b>I – IDENTIFICAÇÃO DO ALUNO</b>   |  |
| Nome:   | Henrique José Wilbert  |
| CV Lattes:  | <a href="http://lattes.cnpq.br/4610551855601588">http://lattes.cnpq.br/4610551855601588</a>  |
| E-mail:   | <a href="mailto:hwilbert@furb.br">hwilbert@furb.br</a>   |
| Telefone:   | (47)99201-6054   |
| <b>II – IDENTIFICAÇÃO DO TRABALHO</b>   |  |
| Título provisório:  | UTILIZAÇÃO DE CLUSTERIZAÇÃO PARA AUXÍLIO EM TOMADA DE DECISÃO A PARTIR DE DADOS DE VAREJO  |
| Orientador:   | Aurélio Faustino Hoppe   |
| Coorientador (se houver):   | Christian Daniel Falaster  |
| Linha de Pesquisa:  | <input type="checkbox"/> Tecnologias aplicadas à informática na educação<br><input checked="" type="checkbox"/> Tecnologias aplicadas ao desenvolvimento de sistemas |
| <b>III – COMPROMISSO DE REALIZAÇÃO DO TCC</b>   |  |
| Eu (aluno),   | Henrique José Wilbert  |
| comprometo-me a realizar o trabalho proposto no semestre 2022-1, de acordo com as normas e os prazos determinados pela FURB, conforme previsto na resolução nº.20/2016. |  |
| Assinatura:   | NÃO É NECESSÁRIO – Encaminhar por mail ao orientador   |
| <b>IV – COMPROMISSO DE ORIENTAÇÃO</b>   |  |
| Eu (orientador),  | Aurélio Faustino Hoppe   |
| comprometo-me a orientar o trabalho proposto no semestre 2022-1, de acordo com as normas e os prazos determinados pela FURB, conforme previsto na resolução nº.20/2016. |  |
| Assinatura:   | NÃO É NECESSÁRIO – Encaminhar por mail ao professor de TCC I   |

Blumenau, 10 de agosto de 2021

| CURSO DE CIÊNCIA DA COMPUTAÇÃO – TCC |             |                      |
|--------------------------------------|-------------|----------------------|
| ( X ) PRÉ-PROJETO                    | ( ) PROJETO | ANO/SEMESTRE: 2021/2 |

## UTILIZAÇÃO DE CLUSTERIZAÇÃO PARA AUXÍLIO EM TOMADA DE DECISÃO A PARTIR DE DADOS DE VAREJO

Henrique José Wilbert

Prof. Aurélio Faustino Hoppe – Orientador

Prof. Christian Daniel Falaster – Coorientador

### 1 INTRODUÇÃO

Com a evolução da tecnologia de informação a partir dos anos 80 e início dos anos 90, várias grandes empresas adotaram sistemas de gerenciamento na forma de softwares Enterprise Resource Planning (ERP) (RASHID; HOSSAIN; PATRICK, 2001, p.2). Estes softwares auxiliam em suas rotinas à nível operacional, seja no controle do estoque, fiscal, financeiro, transacional e até recursos humanos. A partir disso, alcançou-se um patamar de eficiência nunca concebido, visto que registros antes realizados em papel e caneta, passaram a ser produzidos automaticamente. Ainda segundo os autores, em paralelo a informatização desses processos, houve também um crescimento da quantidade de dados armazenados referentes a produtos, clientes, transações, gastos e receitas.

**Excluído:** destes

Diante deste contexto, avançaram-se também as táticas de marketing direto, como por exemplo, o envio de catálogos por correio, até ofertas altamente objetivas, focando em indivíduos selecionados, cujas informações transacionais estavam presentes na base de dados. Percebeu-se que o foco das relações empresa-cliente não está em clientes novos, e sim em clientes já existentes nas bases de dados, visto que o custo para adquirir um cliente novo através de publicidade é muito maior que o custo de alimentar uma relação já existente (PETRISON; BLATTBERG; WANG, 1997, p. 119, tradução nossa).

Segundo Reinartz, Thomas e Kumar (2005, p.77), quando empresas tratam os gastos entre aquisição e retenção de clientes, destinar menos recursos para a retenção impactará em uma lucratividade menor a longo prazo, comparando-se a investimentos menores em aquisição de clientes. Ainda segundo os autores, no conceito de relações de retenção, atribui-se grande ênfase à lealdade e lucratividade de um cliente, sendo lealdade a tendência de o cliente comprar e a lucratividade, a medida geral de quanto lucro um cliente traz à empresa através de suas compras.

**Excluído:** do

**Comentado [AS9]:** Tem a questão da indicação também. Um cliente feliz indica outros clientes.

De acordo Nguyen, Sherif e Newby (2007, p.114) com o avanço da gerência das relações com clientes foram abertas novas vias pelas quais sua lealdade e lucratividade pode ser cultivada, atraindo uma crescente demanda por parte de empresas, visto que a adoção destes meios permite que as organizações melhorarem seu serviço ao consumidor, consequentemente gerando renda. Com isso, diferentes ferramentas acabam sendo utilizadas, como sistemas de recomendação que, geralmente em ramos e-commerce, levam em conta várias características pertinentes ao comportamento do cliente, construindo um perfil próprio que será utilizado para realizar a recomendação de um produto que talvez seja de seu interesse. Outra ferramenta pertinente à lucros e lealdade é a segmentação, que visa separar uma única e confusa massa de clientes em segmentos homogêneos em termos de comportamento, permitindo o desenvolvimento de campanhas e estratégias de marketing especializadas à cada grupo de acordo com suas características (TSIPTSIS; CHORIANOPOULOS, 2009, p.4).

Em relação a segmentação de clientes, algumas métricas tornam-se relevantes nos contextos aos quais estão inseridas. Segundo Kumar (2008, p.29), o modelo Recency-Frequency-Monetary (RFM), é utilizado em empresas de venda por catálogo, enquanto empresas de high-tech tendem a usar Share of Wallet (SOW) para implementar suas estratégias de marketing. Já o modelo Past Customer Value (PCV), geralmente é utilizado em empresas de serviços financeiros. Dentre os modelos citados, o RFM é o que possui maior facilidade de aplicação em diversas áreas de comércio, varejo e supermercados, visto que são necessários apenas os dados transacionais dos clientes (vendas), dos quais são obtidos os atributos de Recência (R), Frequência (F) e Monetário (M).

**Excluído:** recência

**Excluído:** frequência

**Excluído:** monetário

A partir desses dados, segundo Tsiptsis e Chorianopoulos (2009, p.335), é possível detectar bons clientes a partir das melhores pontuações de RFM. Se o cliente efetuou uma compra recentemente, seu atributo R será alto. Caso ele compre muitas vezes ao longo de um determinado período, seu atributo F será maior. E, por fim, caso seus gastos totais forem significativos, terá um atributo M alto. Ao categorizar o cliente dentro destas três características, é possível obter uma hierarquia de importância, tendo os clientes que possuem valores RFM altos no topo, e clientes que possuem valores baixos na base. Apesar destas



vantagens, o modelo padrão original é um tanto quanto arbitrário, segmentando os clientes em quintis, cinco grupos com 20% dos clientes, não atentando-se às nuances e todas as interpretações que a base de clientes pode possuir. Além disso, o método também pode produzir uma grande quantidade de grupos, que por muitas vezes, não representam significativamente os clientes de um estabelecimento, e caso o método de quintis seja utilizado, 125 grupos serão criados. Outro ponto a se observar é a variada gama de interpretações que os atributos RFM podem ter em relação aos tipos de atividades dos estabelecimentos, sendo necessário a adaptação do modelo para cada empresa.

Diante deste cenário, este trabalho propõe a criação de um artefato computacional que utilize o modelo RFM em conjunto com diferentes algoritmos de clusterização ao invés de quintis para segmentar clientes, extraindo de maneira automática as informações de bases de dados, sendo aplicado ao contexto de diversas empresas de varejo, atacado e comércio, visando adequar-se dinamicamente às suas eventuais diferenças de comportamento nos clientes.

## 1.1 OBJETIVOS

O objetivo deste trabalho é desenvolver um artefato computacional de auxílio à segmentação de clientes a partir de múltiplas bases de dados utilizando o modelo RFM.

Os objetivos específicos são:

- implementar e testar diferentes algoritmos de clusterização;
- disponibilizar um mecanismo de visualização dos agrupamentos;
- avaliar a qualidade em relação ao agrupamento dos clientes.

## 2 TRABALHOS CORRELATOS

Neste capítulo serão apresentados os trabalhos similares ao proposto. Na seção 2.1 é apresentado o trabalho de Gustriansyah, Suhandi e Antony (2020), que consiste na aplicação do modelo *Recency Frequency Monetary* (RFM) para clusterização de produtos utilizando K-means, tendo como foco otimizar o número de clusters através de índices de validação. A seção 2.2 descreve o modelo de segmentação denominado *Length, Recency, Frequency, Monetary and Periodicity* (LRFMP) proposto por Peker, Kocyigit e Eren (2017), ao qual considera a longevidade e a periodicidade. Por fim, na seção 2.3 detalha-se o modelo R+FM, sendo uma versão modificada do modelo original que foi aplicada em clientes de uma empresa de e-commerce (TAVAKOLI et al., 2018).

### 2.1 CLUSTERING OPTIMIZATION IN RFM ANALYSIS BASED ON K-MEANS

Gustriansyah, Suhandi e Antony (2020) utilizaram o algoritmo de clusterização K-means para agrupar 2.043 produtos de uma farmácia visando otimizar o manuseio de estoque. Foram utilizadas três características de acordo com o modelo *Recency Frequency Monetary* (RFM) para a separação dos produtos, levando em consideração dados transacionais capturados num período de um ano. O atributo recência classificou os produtos através da última venda realizada num intervalo de 1 a 364 dias. A frequência estabelece a quantidade de transações em que o produto ocorreu, variando num intervalo de 1 a 14.872. Já o atributo monetário, refere-se ao valor total proveniente das vendas acumuladas do produto, sendo definido num intervalo entre 1250 e 1.151.952.500 Rupias Indonésias (Rp.).

Após a atribuição de valores RFM aos produtos, foram utilizados oito índices de validação do melhor número de clusters: *Elbow Method* (EM), que calcula a variação intra-cluster conforme são aumentados os clusters e conclui que o melhor número é aquele que está no cotovelo (elbow) da curva. *Silhouette Index* (SI) que resulta em uma nota de -1 a 1 que indica a quão adequada é a classificação de um objeto dentro de um cluster em comparação aos outros. *Calinski-Harabasz Index* (CHI) que também mede a adequação da quantidade de clusters levando em conta a dispersão entre e intra clusters. *Davies-Bouldin Index* (DBI) que calcula as similaridades entre clusters levando em conta as distâncias e tamanhos dos clusters, quanto menor este índice melhor a separação entre os clusters. *Ratowski Index* (RI) que é baseado na média da soma dos quadrados dos dados entre clusters e a soma total dos quadrados de cada dado dentro de um cluster, dentre as quantidades calculadas escolhe-se a que obtém um maior índice. *Hubert Index* (HI) que é um método visual que indica a quantidade preferida através de um pico no gráfico e é calculado pelo coeficiente de correlação entre matrizes de distância. *Ball-Hall Index* (BHI) definido pela média da distância dos itens com os respectivos centroides do cluster, onde no gráfico o ponto de quantidade de clusters com maior diferença do anterior é sugerido. *Krzanowski-Lai Index* (KLI), que propõe índices internos definidos pelas diferenças entre matrizes de dispersão, e aponta a melhor quantidade de clusters pelo maior número gerado ao realizar a equação com quantidade k. Constatou-se que a maioria deles indicou que o melhor número de clusters seria 3, com base nas condições de interpretação de cada índice explicadas anteriormente.

Comentado [AS10]: Poderias dividir esta frase em 2

Comentado [AS11]: Desenvolver é metodologia.

Faltou falar do uso da clusterização aqui

Comentado [AS12]: Do que?

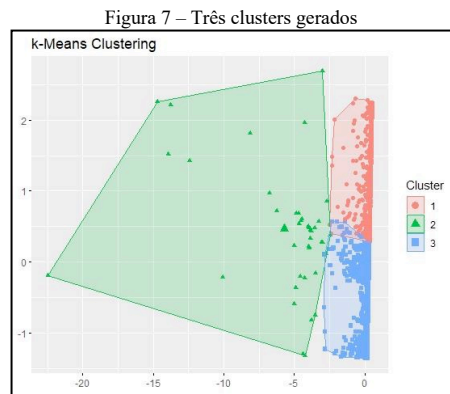
Excluído: final

Formatado: Fonte: Itálico

Segundo Gustriansyah, Suhandi e Antony (2020) nos testes para verificar os clusters gerados, utilizou-se a equação de variância (R). Sendo “R” o valor da divisão entre a distância média dos dados cluster (distância intra-cluster) pela distância média dos dados em outros clusters (inter-cluster). O valor médio alcançado para R foi de 0.19113, sendo que quanto mais próximo de zero, maior a similaridade entre os membros dentro de cada cluster.

Gustriansyah, Suhandi e Antony (2020) utilizaram o software R Programming para gerar a clusterização, resultando na visualização demonstrada na Figura 7, tendo dois clusters com densidade maior e um cluster mais disperso. É possível observar também que o cluster em verde possui uma maior variância entre os próprios dados, enquanto os outros dois clusters possuem uma menor diferença interna.

**Excluído:** Também se observa



Fonte: Gustriansyah, Suhandi e Antony (2020).

Segundo Gustriansyah, Suhandi e Antony (2020) também foram adquiridos os valores médios RFM para cada cluster, conforme apresenta a Tabela 1, sendo que o cluster 3 possui a maior média dos três atributos, e o cluster 1 possui a menor média dos três. É possível identificar um intervalo entre os valores médios de cada atributo, indicando uma diferença significativa inter-cluster.

**Excluído:** número

Tabela 5 – Valores médios de RFM de cada cluster

| Cluster | Recency  | Frequency  | Monetary<br>(in thousands) |
|---------|----------|------------|----------------------------|
| 1       | 75.8167  | 3,436.744  | 3,089,608                  |
| 2       | 224.3947 | 13,013.333 | 76,920,847                 |
| 3       | 331.9681 | 107.418    | 286,927,000                |

Fonte: Gustriansyah, Suhandi e Antony (2020).

Gustriansyah, Suhandi e Antony (2020) concluem que o método gerou clusters com alta similaridade em relação aos dados existentes, apresentado uma segmentação mais objetiva quando comparado ao modelo RFM tradicional no qual os dados são divididos igualmente em cinco segmentos (20% dos dados para a cada segmento). Além disso, os autores também sugerem como extensões a utilização de outros métodos para a comparação, como *Particle Swarm Optimization* (PSO), que é um método computacional que otimiza soluções para uma equação de uma certa medida de qualidade, medianoide (centroides que são parte do conjunto de dados) ou *maximizing-expectancy*, que é um método iterativo para encontrar estimativas de parâmetros para modelos estatísticos com variáveis não observadas.

## 2.2 LRFMP MODEL FOR CUSTOMER SEGMENTATION IN THE GROCERY RETAIL INDUSTRY: A CASE STUDY

Peker, Kocyigit e Eren (2017) propuseram o modelo *Length, Recency, Frequency, Monetary* (LRFM) denominado *Length, Recency, Frequency, Monetary and Periodicity* (LRFMP) para classificar dados reais de 16.024 clientes de mercados de uma franquia na Turquia. Para isso, utilizou-se o algoritmo K-means para segmentar os clientes e três índices de validação de clusters para a otimização das suas quantidades, *Silhouette Index* (SI), *Calinski-Harabasz Index* (CHI) e *Davies-Bouldin Index* (DBI). Após a segmentação dos dados, verificou-se estratégias de gerenciamento e relações com os clientes para aumentar a lucratividade, como tratamento preferencial para clientes importantes, implementação de cartões

fidelidade para aumentar a frequência de compra de clientes não costumam comprar com frequência, promoções voltadas para clientes incertos com sua escolha de local de compra, dentre outras estratégias.

Primeiramente, Peker, Kocyigit e Eren (2017) adaptaram o modelo LRFM, incluindo o parâmetro *periodicity* (periodicidade), pois a análise dos dados foi realizada a partir do histórico de compras em supermercados, que são estabelecimentos com alto número de visitas, tornando importante a regularidade nos padrões de visita e compra. Peker, Kocyigit e Eren (2017) definem a periodicidade como a regularidade das visitas de um determinado cliente. Sendo atribuída como o desvio padrão dos tempos inter-visita do cliente (quantia de dias entre duas visitas consecutivas). Se um cliente possui valores baixos de periodicidade, significa que este realiza visitas ou compras em intervalos fixos, podendo caracterizá-lo como cliente regular. Além disso, os autores também modificaram o atributo de recência, transformando-o na média das diferenças entre a data das três últimas compras e a data atual, ao invés da simples diferença entre a data da última compra e a data atual estabelecida no modelo RFM padrão.

**Comentado [AS13]:** Deve-se evitar iniciar a frase com gerúndio. Gerúndio complementa alguma ideia.

Após adquirir os atributos LRFMP dos dados transacionais dos clientes, Peker, Kocyigit e Eren (2017) aplicaram um método de normalização simples nos dados, considerando o intervalo de 0 e 1. Esta normalização foi feita pois os valores LRFMP variam em relação ao intervalo e escala, fato que poderia afetar negativamente a análise dos clusters.

Antes de aplicar a clusterização, Peker, Kocyigit e Eren (2017) utilizaram três índices para validação da quantidade possível de clusters: *Silhouette Index* (SI) que resulta em uma nota de -1 a 1 que indica o quão adequada é a classificação de um objeto dentro de um cluster em comparação aos outros, quanto maior o valor, melhor. *Calinski-Harabasz Index* (CHI) que mede a adequação da quantidade de clusters levando em conta a dispersão entre e intra clusters, um valor alto é preferido. *Davies-Bouldin Index* (DBI) que calcula as similaridades entre clusters levando em conta as distâncias e tamanhos dos clusters, quanto menor este índice melhor será a separação entre os clusters. A partir deles, Peker, Kocyigit e Eren (2017) executaram o algoritmo K-means variando o k de 2 a 9, e os resultados destas iterações foram avaliadas utilizando os três índices. Com base nos resultados, decidiu-se utilizar um número de 5 clusters, pois 2 dos 3 índices sugerem 5 como sendo a quantidade ideal.

Peker, Kocyigit e Eren (2017) utilizaram uma base de dados de uma franquia de mercados que possui mais de dez lojas na cidade de Antália na Turquia. Os dados são compostos por cerca de dois milhões de transações de 16.024 clientes num período de dois anos. Foram removidos os clientes com menos que três compras. Além disso, os autores removeram dados duplicados, transações com valores faltantes assim como, agregaram as compras dentro de um mesmo dia. Depois dessas operações, a quantidade de clientes caiu para 10.471, sendo aplicado na sequência o K-means. A Tabela 6 demonstra a quantidade de clientes nos clusters, os valores médios de LRFMP para cada cluster. Já na última coluna, aplicou-se uma técnica no qual o atributo do cluster recebe uma seta para cima (↑) caso o seu valor for maior que a média do atributo dos outros clusters, e uma seta para baixo (↓), caso seu valor for menor que a média.

**Excluído:** também

Tabela 6 – Valores médios dos clusters

| Cluster | Sample size | Average L | Average R | Average F | Average M | Average P | LRFMP Scores   |
|---------|-------------|-----------|-----------|-----------|-----------|-----------|----------------|
| 1       | 538         | 633.29    | 39.67     | 175.24    | 24.32     | 4.99      | L↑ R↓ F↑ M↓ P↓ |
| 2       | 4,681       | 564.50    | 90.19     | 33.44     | 31.73     | 31.49     | L↑ R↓ F↑ M↓ P↓ |
| 3       | 1,091       | 482.17    | 301.18    | 5.33      | 34.21     | 159.32    | L↑ R↑ F↓ M↓ P↑ |
| 4       | 818         | 374.01    | 220.24    | 11.85     | 104.18    | 45.81     | L↓ R↑ F↓ M↑ P↑ |
| 5       | 3,343       | 173.70    | 399.79    | 10.34     | 30.14     | 27.85     | L↓ R↑ F↓ M↓ P↓ |
| Average |             | 419.81    | 218.59    | 28.74     | 36.76     | 43.41     |                |

Fonte: Peker, Kocyigit e Eren (2017).

A partir destes resultados, Peker, Kocyigit e Eren (2017) descreveram as características dos grupos. O grupo 1 representa clientes leais de alta contribuição que, apesar de comporem a menor parcela dos clientes (5,14%), possuem a maior contribuição total entre os grupos. Também é possível observar, que este grupo possui a menor periodicidade média de todos, caracterizando estes clientes como regulares. O grupo 2, representando a maior parcela dos clientes (44,70%) foi classificado como clientes leais de baixa contribuição pois apesar de visitar mais frequentemente as lojas, não possuem tanta contribuição quanto o grupo 1. O grupo 3, com tamanho de 10,42%, foi classificado como clientes incertos, pois possui o atributo de longevidade alto e recência também alta, significando que são clientes com longa história de compra, porém sem muitas compras recentes, vale notar que este grupo possui o maior valor de periodicidade de todos os grupos, caracterizando-o como um grupo de clientes sem rotina de compra definida. O grupo 4 e 5 foram classificados como clientes perdidos, visto que possuem poucas compras recentes, baixa frequência, e baixa longevidade, denotando um cliente que tem uma pouca interação com a franquia. O

grupo 4, contendo uma pequena parcela de 7,81% dos clientes, gasta consideravelmente mais, logo foi classificado como contribuição alta, e o 5, cuja parcela é 31,93%, classificado como contribuição baixa.

A partir desta classificação, Peker, Kocyigit e Eren (2017) estabeleceram estratégias para cada grupo de clientes, como tratamento especial (vagas de estacionamento preferencial, presentes de aniversário, filas preferenciais) para clientes do grupo 1, de maneira a não perder a relação leal com a loja. Para os grupos 3, 4 e 5 cuja frequência é baixa, foi sugerida a adoção de programas de cartão fidelidade para aumentar a frequência deles. Para clientes incertos como no grupo 3, aplicou-se descontos e promoções, de maneira a incentivar os clientes, supostamente sensíveis aos preços, a voltar sua atenção à franquia. Para clientes perdidos, sugeriu-se uma análise mais profunda sobre o motivo da perda, como análise de feedback, inferência de motivos, dentre outros.

Por fim, Peker, Kocyigit e Eren (2017) concluem que o estudo contribuiu com a proposta de um novo modelo RFM, que possibilita uma análise mais profunda que o seu modelo original, visto que as características utilizadas e modificadas permitem uma melhor definição do comportamento de cada cliente. Outra contribuição foi a adição do atributo periodicidade (P) no modelo, que, ao contrário do modelo RFM padrão, permite identificar se os clientes de um grupo variam em sua rotina de compras. Outra melhoria apontada é a modificação do atributo de recência, que uma vez calculado como uma média, permite uma caracterização mais precisa que o atributo R comumente utilizado. Uma das limitações destacadas pelos autores é a localidade do estudo, pois foi realizado somente com dados originários de uma cidade, sendo que o comportamento de clientes pode variar de acordo com as diferentes localidades onde é feita a análise. A partir disso, sugere-se uma análise mais ampla contemplando outros locais. Outra sugestão feita por Peker, Kocyigit e Eren (2017) é a adição de novos atributos ao modelo, como a quantidade de produtos comprados, quantidade de produtos perecíveis e não perecíveis comprados, a fim de promover uma interpretação mais profunda do comportamento.

### 2.3 CUSTOMER SEGMENTATION AND STRATEGY DEVELOPMENT BASED ON USER BEHAVIOR ANALYSIS, RFM MODEL AND DATA MINING TECHNIQUES: A CASE STUDY

Tavakoli *et al.* (2018) desenvolveram o modelo RFM, denominado “R+FM”, sendo utilizado em conjunto com o algoritmo de clusterização K-means para segmentar 3 milhões de clientes da maior empresa de *E-commerce* do Oriente Médio. Além disso, o modelo de segmentação foi comparado com o utilizado pela empresa, sendo aplicado em uma campanha de Short Message Service (SMS) focada em aumentar os ganhos de cada segmento.

Tavakoli *et al.* (2018) defendem a utilização de um novo modelo de caracterização de clientes, argumentando que o modelo ideal necessita adaptar-se às mudanças de comportamento dos clientes, possuir certa independência de supervisão, levando em consideração a similaridade dos comportamentos dos clientes e a relação entre os atributos Frequência e Monetário. Para isso, construiu-se uma variante do modelo Recency, Monetary, Frequency (RFM) denominado de R+FM, que possui o atributo de recência separado dos demais, utilizando uma segmentação à parte do modelo FM. Os autores separaram os clientes em 3 grupos: os que compraram recentemente (cuja última compra foi dentro de 90 dias), denominados de ativos, clientes que compraram em um passado recente (cuja última compra foi entre 90 e 365 dias), denominados de expirando e por fim, os clientes que não compraram por um longo tempo (cuja última compra foi a mais de 365 dias), denominados de expirados. Para o atributo de frequência, Tavakoli *et al.* (2018) atentaram-se especialmente com a data da primeira compra, pois acreditam que a frequência tem uma importância maior conforme sua recência. Logo, definiram a frequência como a quantidade de compras dividida pela quantidade de dias desde a primeira compra, utilizando também uma função exponencial de decaimento, que efetivamente atribui um peso maior para anos mais recentes, sendo cada ano duas vezes mais pesado que o ano anterior. Como atributo monetário, estabeleceu-se a média dos valores das compras de um cliente, visto que um valor de soma total de compras, segundo os autores, estaria encorajando duas vezes os clientes.

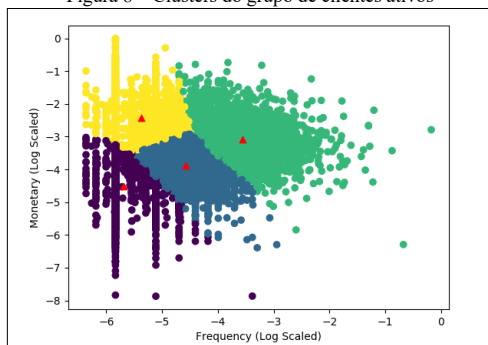
Após a definição do modelo, Tavakoli *et al.* (2018), balancearam a relação entre frequência e monetário, criando a quarta característica que é definida pela combinação linear dos dois atributos, que nada mais é que a soma de cada atributo ponderada pelo peso de cada um. No tratamento dos dados, utilizou-se a técnica de remoção de *outliers* (clientes que não se encaixam no padrão normal) que não se encontram dentro dos intervalos interquartis, que são os intervalos que possuem os dados que pertencem à tendência média do conjunto de dados em geral. Também foram escalados os atributos de frequência e monetário para que seus intervalos sejam iguais, sendo aplicada a normalização *min-max*, que transforma os valores para estarem dentro do intervalo entre 1 e 0. Como os dados monetários e de frequência tratados possuem uma característica de cauda longa, fenômeno estatístico onde os dados são distribuídos de forma

Excluído: i

decrecente, foi aplicada uma transformação logarítmica para normalizar a distribuição, visto que a quantidade de valores baixos é muito alta, podendo atrapalhar a análise.

Como existem duas segmentações (R e FM), Tavakoli *et al.* (2018) estabeleceram segmentos FM para cada segmento R, resultando nos seguintes grupos: para clientes ativos existem os grupos de alto valor, médio valor com alto monetário, médio valor com alta frequência e baixo valor. Para clientes que estão expirando existem os grupos de alto, médio e baixo valores, sendo aplicado também para clientes expirados. Estes grupos foram definidos por Tavakoli *et al.* (2018) com ajuda da empresa de *E-commerce* Digikala. A partir disso, aplicou-se o K-means com  $k=4$  para o grupo de clientes ativos,  $k=3$  para o grupo de clientes expirando e  $k=3$  para o grupo de clientes expirados, resultando em um total de 10 clusters. Na Figura 8 são identificados os clusters gerados somente a partir do grupo de clientes ativos, organizados em um gráfico de valor monetário por frequência. Sendo o cluster de cor verde composto pelos clientes de alto valor, o cluster de cor amarela composto pelos clientes de médio valor com alto monetário, o cluster de cor azul composto pelos clientes de médio valor com alta frequência, e por fim, o cluster de cor roxa composto pelos clientes de baixo valor.

Figura 8 – Clusters do grupo de clientes ativos



Fonte: Tavakoli *et al.* (2018).

Após a geração dos grupos, Tavakoli *et al.* (2018) discorrem sobre possíveis estratégias para cada segmento. Sugerindo um maior foco em clientes ativos com valor médio e baixo, bem como a manutenção de clientes já valiosos. Os autores também enfatizam a importância em recuperar os clientes do grupo expirando, cuja chance de retorno não é tão baixa quanto o grupo expirado, que por si só requer uma estratégia especial de reengajamento dos clientes à empresa.

Além da elaboração de estratégias, Tavakoli *et al.* (2018) implementaram uma campanha de SMS focada somente no segmento de clientes ativos (recência abaixo de 90 dias), pois a empresa já tinha realizado outras campanhas em clientes ativos anteriormente. Nesta campanha cada cliente foi apresentado com um *Voucher* condizente com o segmento ao qual o cliente pertencia. Para clientes ativos com valor alto foi oferecido um desconto de 10 por cento com um desconto máximo de até \$20, com o objetivo de manter a lealdade destes clientes. Para clientes ativos de valor médio com alto valor monetário, foi oferecido um desconto de 10 por cento com um desconto máximo de até \$10, o valor foi menor pois o objetivo era aumentar frequência de compra destes clientes, que em tese já gastam bastante. Para clientes ativos de valor médio com alta frequência foi oferecido um desconto de 10 por cento com um desconto máximo de até \$20 para vendas que custem mais que \$50 (que é o valor médio gasto por este segmento), incentivando assim uma compra de maior valor. Por fim, para clientes ativos de baixo valor, foi oferecido um desconto de 10 por cento com um desconto máximo de até \$20, com o objetivo de converter estes clientes em mais leais.

Para a análise da campanha, Tavakoli *et al.* (2018) selecionaram aleatoriamente 20% dos clientes de cada segmento para compor um grupo de controle, cujos *Vouchers* não foram enviados. Este grupo de controle foi comparado com os outros grupos da campanha para obter um valor de referência do aumento do valor monetário após sua conclusão. Os resultados alcançados por Tavakoli *et al.* (2018) podem ser observados no Quadro 1, ao qual percebe-se um aumento de \$14,30 na média monetária dos clientes ativos com alta frequência, mais do que o aumento de \$3,20 sofrido pelo grupo de controle. É possível também observar que a média monetária de todos os grupos alvo da campanha aumentou consideravelmente, enquanto o grupo de controle aumentou pouco ou até diminuiu, indicando uma efetividade no objetivo da campanha.

**Comentado [AS14]:** Deve-se evitar iniciar a frase com gerúndio. Gerúndio complementa alguma ideia.

Podes unir esta frase com a anterior.

Quadro 1 – Dados da média monetária do grupo de controle e grupo de campanha

| Recency | Segment                          | Average Monetary (USD) |                |                 |                |
|---------|----------------------------------|------------------------|----------------|-----------------|----------------|
|         |                                  | Control Users          |                | Campaign Users  |                |
|         |                                  | Before Campaign        | After Campaign | Before Campaign | After Campaign |
| Active  | High Value                       | 74.2                   | 73.8           | 88.2            | 89.2           |
|         | Medium Value with High Monetary  | 100.2                  | 97.6           | 104.6           | 105.2          |
|         | Medium Value with High Frequency | 32                     | 35.2           | 35.4            | 49.7           |
|         | Low Value                        | 50.7                   | 53.2           | 56.4            | 65.2           |

Fonte: Tavakoli *et al.* (2018).

Tavakoli *et al.* (2018) concluem que houve uma melhora no desempenho da campanha lançada em comparação com as anteriores, indicando ainda, que elas obtinham uma taxa de compra de 0,1 por cento, sendo que a campanha lançada para validação do modelo obteve uma taxa de 1 por cento, cerca de dez vezes mais efetivo. Os autores justificam esta melhora ao processo de segmentação do modelo, resultando em clusters mais significativos, facilitando a aplicação de vouchers específicos. Por fim, Tavakoli *et al.* (2018) sugerem o melhoramento da definição do atributo de recência de forma que seja mais útil ao time de marketing. Também recomendam calcular o Customer Lifetime Value (CLV), que atribui o valor vitalício à cada segmento e cliente, de forma a quantificar o valor que um cliente pode proporcionar à empresa.

### 3 PROPOSTA DO PROTÓTIPO

Esse capítulo visa apresentar a justificativa para a elaboração deste trabalho, os requisitos que serão seguidos e a metodologia que será utilizada. Será apresentada também uma breve revisão bibliográfica das principais áreas de estudo que serão exploradas, bem como os principais termos utilizados.

#### 3.1 JUSTIFICATIVA

No Quadro 2 é apresentado um comparativo entre os trabalhos correlatos. As linhas representam as características relevantes e as colunas representam os trabalhos.

Quadro 2 – Comparativo entre os trabalhos correlatos

| Características                           | Correlatos | Gustriansyah, Suhandi e Antony (2020)   | Peker, Kocyigit e Eren (2017)                         | Tavakoli <i>et al.</i> (2018)                      |
|---|------------|---|---|--|
| Alvo da clusterização                     |            | Produtos                                | Clientes  | Clientes   |
| Modelo utilizado                          |            | RFM                                     | LRFMP   | R+FM   |
| Objetivo da segmentação                   |            | Gerenciamento de estoque                | Gerenciamento das relações com cliente                | Gerenciamento das relações com cliente             |
| Algoritmo de clusterização utilizado      |            | K-means                                 | K-means   | K-means  |
| Foco metodológico                         |            | Otimização de k com diferentes métricas | Formulação de um modelo novo e análise dos resultados | Formulação de um modelo novo e campanha de ofertas |
| Número de dados (clientes/produtos)       |            | 2.043                                   | 16.024  | ~3.000.000   |
| Quantidade de índices para validação de k |            | 8                                       | 3   | -  |
| Quantidade de clusters gerados            |            | 3                                       | 5   | 10   |
| Inferências sobre os dados                |            | -                                       | Sim   | Sim  |

Fonte: elaborado pelo autor.

A partir do Quadro 2, pode-se observar que Gustriansyah, Suhandi e Antony (2020) clusterizaram produtos de uma base de dados utilizando o modelo RFM padrão. Já Peker, Kocyigit e Eren (2017) optaram pelo desenvolvimento de um modelo novo, considerando a periodicidade (LRFMP). Tavakoli *et al.* (2018) também desenvolveram um novo modelo, ao qual a característica recência foi modificada e separada (R+FM).

Gustriansyah, Suhandi e Antony (2020) tinham como objetivo melhorar o gerenciamento de estoque, prezando por uma segmentação mais conclusiva sobre os produtos, visto que o modelo RFM padrão define segmentos arbitrariamente sem adequar-se às peculiaridades dos dados, enquanto o modelo aplicado através de k-means alcançou uma segmentação com dados altamente similares em cada cluster. Por outro lado, Peker, Kocyigit e Eren (2017) e Tavakoli *et al.* (2018) objetivavam o gerenciamento das relações com os clientes através de estratégias focadas em segmentos, visando aumentar a renda que eles fornecem à empresa. Todos os autores utilizaram o algoritmo K-means, por ser confiável e amplamente difundido. Vale ressaltar que no trabalho de Gustriansyah, Suhandi e Antony (2020), o algoritmo teve um

foco metodológico maior, visto que foram utilizados 8 índices de validação para k clusters, visando otimizar a organização dos segmentos.

A quantidade de dados segmentados variou bastante entre os três trabalhos devido aos diferentes contextos de aplicação. Gustriansyah, Suhandi e Antony (2020) tinham 2.043 produtos na base de dados para segmentar, resultando em 3 clusters. Já Peker, Kocyigit e Eren (2017) possuíam o registro de 16.024 clientes de uma rede de padarias, sendo especificados 5 segmentos, obtidos através de uma análise por três índices de validação (Silhouette, Calinski-Harabasz e Davies-Bouldin). Por fim, Tavakoli *et al.* (2018) agruparam dados de 3 milhões de clientes pertencentes à base de dados de um e-commerce do Oriente Médio, resultando em 10 clusters, sendo 3 pertencentes à característica de recência, e os outros 7 distribuídos entre as características de frequência e monetária. Ressalta que Tavakoli *et al.* (2018) testaram o modelo em produção, montando uma campanha que focava no segmento de clientes ativos, visando primariamente aumentar os lucros da empresa, utilizando também um grupo de controle e comparação de renda antes e depois da campanha.

Excluído: E

Comentado [AS15]: Quem?

Gustriansyah, Suhandi e Antony (2020) demonstraram a possibilidade da aplicação de RFM fora do uso convencional de segmentação de clientes, e adquiriram clusters com uma variância média de 0.19113. Além disso, os autores sugeriram outras formas de comparação de dados, como Particle Swarm Optimization (PSO), medioides ou até *maximizing-expectancy*. Peker, Kocyigit e Eren (2017) segmentaram clientes de uma rede de mercados na Turquia em “clientes leais de alta contribuição”, “clientes leais de baixa contribuição”, “clientes incertos”, “clientes perdidos de alto gasto” e “clientes perdidos de baixo gasto”. Desta maneira, os autores providenciaram visões e estratégias (promoções, ofertas, regalias) de aumento de renda sobre os comportamentos dos clientes, porém limitaram-se a aplicar em um segmento específico de mercado. Por fim, Tavakoli *et al.* (2018) agruparam clientes de uma empresa de e-commerce com base em sua recência, resultando em clientes “Ativos”, “Expirando” e “Expirados”, e destes segmentos, sucessivamente separados em grupos de “Alto”, “Médio” e “Baixo” valores, validando posteriormente a segmentação através de uma campanha de ofertas para os clientes do grupo “Ativos”.

Excluído: E

Todos os trabalhos aqui citados procuraram implementar o modelo RFM num contexto de clusterização por K-means, alterando o modelo e o manejo dos dados de acordo com cada categoria, seja ele produto ou cliente, varejo ou mercado. Com isso, criaram-se atributos e foram modificados alguns já existentes para atender às especificidades de cada contexto, visto que todos os trabalhos focaram em uma só base de dados, inevitavelmente adequando-se às mesmas.

Desta forma, este trabalho demonstra ser relevante, pois almeja aplicar o modelo RFM em conjunto com vários algoritmos de clusterização em forma de um artefato computacional que se adeque à vários contextos (mercado, comércio, varejo etc.), utilizando várias bases de dados reais para testar a validade dos algoritmos utilizados. Vislumbra-se utilizar três índices para validação da qualidade dos clusters (Silhouette, Calinski-Harabasz e Davies-Bouldin). Além disso, deseja-se obter clusters significativos e coerentes com cada segmento de mercado aplicado. Outra contribuição deste trabalho seria no âmbito comercial, com a geração de informações sobre as similaridades de clientes de cada segmento de mercado, podendo auxiliar gestores e administradores de empresas a obter uma visão crítica sobre os comportamentos de clientes ao longo das diferentes bases de dados, podendo também denotar características comuns a todos. Outra relevância seria a utilização deste trabalho em ambiente acadêmico, visto que serão aplicados diferentes algoritmos de clusterização, podendo providenciar informações sobre seus desempenhos e qualidade de agrupamento, além de ser aplicados processos de obtenção, limpeza e transformação de dados.

Comentado [AS16]: Como?

Comentado [AS17]: Isso é possível? Não é um escopo muito amplo para um TCC?

Formatado: Fonte: (Padrão) Times New Roman, 10 pt

### 3.2 REQUISITOS PRINCIPAIS DO PROBLEMA A SER TRABALHADO

O artefato computacional a ser desenvolvido deverá:

- adquirir os dados transacionais de clientes a partir de um banco de dados (Requisito Funcional - RF);
- extrair dos clientes as características (recência, frequência e monetária) utilizadas no modelo RFM (RF);
- filtrar os clientes sem quantidade de compras relevantes (RF);
- normalizar os dados para evitar disparidades nas escalas dos dados, principalmente no atributo monetário (RF);
- aplicar três índices de validação (Silhouette, Calinski-Harabasz e Davies-Bouldin) para verificar a qualidade dos clusters (RF);
- apresentar em um gráfico 3D os clientes, com sua localização definida pela pontuação do cliente nas características RFM (RF);
- utilizar algoritmos de clusterização tais como K-means, mean-shift e DBSCAN (RF);
- utilizar a linguagem Python para o desenvolvimento (Requisito Não Funcional - RNF);

Comentado [AS18]: Coloque o significado da sigla



- i) utilizar o ambiente de desenvolvimento Jupyter Notebook (RNF);
- j) utilizar o banco de dados PostgreSQL para ler os dados das bases utilizadas (RNF).

### 3.3 METODOLOGIA

O trabalho será desenvolvido observando as seguintes etapas:

- a) levantamento bibliográfico: pesquisar trabalhos relacionados e estudar sobre o modelo RFM e suas aplicações, algoritmos de clusterização, métodos de tratamento de dados e índices de validação;
- b) seleção de bases de dados: obter bases de dados de usuários cedidas pela empresa Intelidata Informática, desenvolvedora de software de gestão comercial. Serão selecionadas conforme sua adequação ao objetivo do trabalho, variando em tamanho e segmento de mercado;
- c) definição das características do modelo RFM: definir os atributos utilizados para caracterizar os clientes no modelo RFM;
- d) definição de métricas do modelo RFM: definir as métricas para mensuração e atribuição de pontuação de cada característica no modelo RFM;
- e) definição dos algoritmos de clusterização: pesquisar e escolher o algoritmo de clusterização que realizará o agrupamento das características RFM;
- f) implementação: implementar o artefato computacional de segmentação levando em consideração as etapas (b) até (e), utilizando a linguagem Python;
- g) análise dos clusters: avaliar a qualidade dos clusters gerados a partir dos diferentes algoritmos de clusterização e seus comportamentos em múltiplas bases de dados, aplicando índices de validação e apresentando-os na forma de gráficos.

As etapas serão realizadas nos períodos relacionados no Quadro 3.

Quadro 3 – Cronograma de atividades a serem realizadas

| etapas / quinzenas                          | 2022 |   |      |   |      |   |      |   |      |   |
|---|------|---|------|---|------|---|------|---|------|---|
|   | fev. |   | mar. |   | abr. |   | maio |   | jun. |   |
|   | 1    | 2 | 1    | 2 | 1    | 2 | 1    | 2 | 1    | 2 |
| levantamento bibliográfico                  |      |   |      |   |      |   |      |   |      |   |
| seleção de bases de dados                   |      |   |      |   |      |   |      |   |      |   |
| definição das características do modelo RFM |      |   |      |   |      |   |      |   |      |   |
| definição de métricas do modelo RFM         |      |   |      |   |      |   |      |   |      |   |
| definição dos algoritmos de clusterização   |      |   |      |   |      |   |      |   |      |   |
| implementação                               |      |   |      |   |      |   |      |   |      |   |
| análise dos clusters                        |      |   |      |   |      |   |      |   |      |   |

Fonte: elaborado pelo autor.

## 4 REVISÃO BIBLIOGRÁFICA

Neste capítulo serão brevemente aprofundados os assuntos que servirão de base para a realização deste trabalho. Serão tratados os temas de clustering, modelo RFM e índices de validação de clusters.

Para Cherkassky e Mulier (2007), clustering se trata do problema de separar um conjunto de dados em grupos chamados de “clusters” baseado em alguma medida de similaridade. O objetivo é encontrar um conjunto de clusters dos quais as amostras dentro dos mesmos são mais similares entre si do que quando comparadas com amostras de outros clusters. Existem vários algoritmos para clusterizar dados, especializando-se em situações específicas, como o *Density-based spatial clustering of applications with noise* (DBSCAN), que é um algoritmo que agrupa dados baseado em sua densidade, utilizando conceitos de pontos núcleo, pontos vizinhos e ruído. Outro método utilizado é o *Hierarchical Clustering*, que constrói uma hierarquia de clusters, podendo aplicar abordagens de cima para baixo, onde as observações pertencem a um cluster que é dividido em outros menores, ou a abordagem de baixo para cima, onde cada dado inicia como um cluster, e é então agrupado conforme o algoritmo é executado. Por fim, um dos algoritmos mais utilizados é o K-means, cuja orientação é focada em centroides, ou seja, seus clusters são representados por um ponto central (não necessariamente fazendo parte do conjunto de dados) que possui a menor distância possível entre si mesmo e o resto dos dados do cluster.

Segundo Hughes (2011), o modelo RFM é “Um meio antigo e altamente preditivo de determinar quem irá responder e comprar. Um método de codificar clientes existentes. Usado para prever resposta, tamanho médio de pedido, e outros fatores”. Este modelo categoriza geralmente clientes através das características de recência (R), frequência (F) e monetária (M). As métricas utilizadas para medir tais características podem variar, porém geralmente classificam recência como a quantidade de dias desde a

Comentado [AS19]: redundante

Comentado [AS20]: 1ª letra maiúscula



última compra, frequência como a quantidade de compras dentro de um determinado período, e monetária como o total acumulado de todas as vendas realizadas para um cliente. Este modelo é muito utilizado em marketing direto, onde os meios de comunicação são diretamente entre a empresa e o consumidor, realizado através de mídias sociais, e-mail, mensagens SMS ou até pelo correio. No estudo realizado por Verhoef *et al.* (2003), cerca de 90% das empresas questionadas sobre a aplicação métodos de segmentação (como o RFM) afirmam que possuíam como objetivo a seleção de alvos, ou seja, encontrar o segmento de clientes que mais se identificam com a empresa, e 64,4% citaram como objetivo o tratamento diferencial de clientes, com promoções, preços e ofertas especiais.

O índice de Silhouette é utilizado no conceito de clusterização, para analisar a qualidade de um dado localizado em determinado cluster, levando em conta a distância média entre clusters. Com o cálculo deste índice, valores próximos de 1 para um determinado dado em um cluster são considerados bons, valores perto de 0 indicam que o dado está entre clusters e caso o valor seja próximo de -1 significa que provavelmente o dado está no cluster errado. O índice de Calinski-Harabasz é definido como a razão de duas somas calculadas entre todos os clusters: a soma da dispersão intra-clusters e soma da dispersão inter-clusters. Quanto maior seu valor, melhor a performance da clusterização observada. Por fim, o índice de Davies-Bouldin indica a similaridade média entre clusters, levando em conta sua distância e tamanho. Um valor baixo para este índice indica uma melhor separação entre os clusters.

#### REFERÊNCIAS

CHERKASSKY, Vladimir S.; MULIER, Filip. Methods for data reduction and dimensionality reduction. In: CHERKASSKY, Vladimir S.; MULIER, Filip. **Learning from data: concepts, theory, and methods**. 2. ed. Hoboken: Ieee Press, 2007. Cap. 6, p. 191

GUSTRIANSYAH, Rendra; SUHANDI, Nazori; ANTONY, Fery. Clustering optimization in RFM analysis Based on k-Means. **Indonesian Journal Of Electrical Engineering And Computer Science**, [S.l.], v. 18, n. 1, p. 470-477, abr. 2020. Mensal. Disponível em: <http://ijeecs.iaescore.com/index.php/IJEECS/article/view/20264>. Acesso em: 02 set. 2021

HUGHES, Arthur Middleton. **Strategic Database Marketing 4e: the masterplan for starting and managing a profitable, customer-based marketing program**. 4. ed. [S.l.]: McGraw-Hill, 2011. 608 p.

KUMAR, V. **Managing Customers for Profit: strategies to increase profits and build loyalty**. Upper Saddle River: Pearson Prentice Hall, 2008. 296 p.

NGUYEN, Thuyuyen H.; SHERIF, Joseph S.; NEWBY, Michael. Strategies for successful CRM implementation. **Information Management & Computer Security**, [S.l.], v. 15, n. 2, p. 102-115, maio 2007. Disponível em: <https://www.emerald.com/insight/content/doi/10.1108/09685220710748001/full/html?journalCode=imcs>. Acesso em: 26 set. 2021.

PEKER, Serhat; KOCYIGIT, Altan; EREN, P. Erhan. LRFMP model for customer segmentation in the grocery retail industry: a case study. **Marketing Intelligence & Planning**, [S.l.], v. 35, n. 4, p. 544-559, 6 maio 2017. Emerald. <http://dx.doi.org/10.1108/mip-11-2016-0210>. Disponível em: <https://www.emerald.com/insight/content/doi/10.1108/MIP-11-2016-0210/full/html>. Acesso em: 07 set. 2021.

PETRISON, Lisa A.; BLATTBERG, Robert C.; WANG, Paul. Database marketing: past, present, and future. **Journal Of Direct Marketing**, [S.l.], v. 11, n. 4, p. 109-125, mar. 1997. Wiley. [http://dx.doi.org/10.1002/\(sici\)1522-7138\(199723\)11:43.0.co;2-g](http://dx.doi.org/10.1002/(sici)1522-7138(199723)11:43.0.co;2-g). Disponível em: [https://onlinelibrary.wiley.com/doi/10.1002/\(SICI\)1522-7138\(199723\)11:4%3C109::AID-DIR12%3E3.0.CO;2-G](https://onlinelibrary.wiley.com/doi/10.1002/(SICI)1522-7138(199723)11:4%3C109::AID-DIR12%3E3.0.CO;2-G). Acesso em: 19 set. 2021.

RASHID, Mohammad A.; HOSSAIN, Liaquat; PATRICK, Jon David. The Evolution of ERP Systems: a historical perspective. In: NAH, Fiona Fui-Hoon. **Enterprise Resource Planning: solutions and management**. Hershey: Irm Press, 2001. p. 35-50. Disponível em: <https://books.google.com.br/books?id=qBcJwDWk4ioC&lpg=PR1&ots=9MrXoQhaRL&dq=Enterprise%20Resource%20Planning%3A%20Solutions%20and%20Management&lr&hl=pt-BR&pg=PR1#v=onepage&q=Enterprise%20Resource%20Planning%20Solutions%20and%20Management&f=false>. Acesso em: 19 set. 2021.

REINARTZ, Werner; THOMAS, Jacquelyn S.; KUMAR, V. Balancing Acquisition and Retention Resources to Maximize Customer Profitability. **Journal Of Marketing**, [S.l.], v. 69, n. 1, p. 63-79, jan. 2005.

Comentado [AS21]: Fonte?

Comentado [AS22]: Não está de acordo com a norma

TAVAKOLI, Mohammadreza *et al.* Customer Segmentation and Strategy Development Based on User Behavior Analysis, RFM Model and Data Mining Techniques: a case study. In: 2018 IEEE 15TH INTERNATIONAL CONFERENCE ON E-BUSINESS ENGINEERING (ICEBE), 15., 2018, Xiam. **Proceedings** [...] . [S.l.]: Ieee, 2018. p. 119-126. Disponível em: [https://www.researchgate.net/publication/330027350\\_Customer\\_Segmentation\\_and\\_Strategy\\_Development\\_Based\\_on\\_User\\_Behavior\\_Analysis\\_RF\\_Model\\_and\\_Data\\_Mining\\_Techniques\\_A\\_Case\\_Study](https://www.researchgate.net/publication/330027350_Customer_Segmentation_and_Strategy_Development_Based_on_User_Behavior_Analysis_RF_Model_and_Data_Mining_Techniques_A_Case_Study). Acesso em: 11 set. 2021.

TSIPTSIS, Konstantinos K.; CHORIANOPOULOS, Antonios. Data Mining Techniques in CRM: inside customer segmentation. Chichester: John Wiley & Sons, 2009. 374 p.

VERHOEF, Peter C *et al.* The commercial use of segmentation and predictive modeling techniques for database marketing in the Netherlands. **Decision Support Systems**, [S.l.], v. 34, n. 4, p. 471-481, mar. 2003. Disponível em: <https://liacs.leidenuniv.nl/~puttenpwhvander/library/Others/segmpredmodel-hoekstra.pdf>. Acesso em: 19 set. 2021.

## FORMULÁRIO DE AVALIAÇÃO – PROFESSOR AVALIADOR

Avaliador(a): **Andreza Sartori**

Atenção: quando o avaliador marcar algum item como atende parcialmente ou não atende, deve obrigatoriamente indicar os motivos no texto, para que o aluno saiba o porquê da avaliação.

| ASPECTOS AVALIADOS <sup>1</sup> |  | Atende | atende parcialmente | não atende |
|---------------------------------|--|--------|---------------------|------------|
| ASPECTOS TÉCNICOS               | 1. INTRODUÇÃO<br>O tema de pesquisa está devidamente contextualizado/delimitado?   | X      |                     |            |
|                                 | O problema está claramente formulado?  | X      |                     |            |
|                                 | 2. OBJETIVOS<br>O objetivo principal está claramente definido e é passível de ser alcançado?   |        | X                   |            |
|                                 | Os objetivos específicos são coerentes com o objetivo principal?   |        | X                   |            |
|                                 | 3. TRABALHOS CORRELATOS<br>São apresentados trabalhos correlatos, bem como descritas as principais funcionalidades e os pontos fortes e fracos?                          | X      |                     |            |
|                                 | 4. JUSTIFICATIVA<br>Foi apresentado e discutido um quadro relacionando os trabalhos correlatos e suas principais funcionalidades com a proposta apresentada?             | X      |                     |            |
|                                 | São apresentados argumentos científicos, técnicos ou metodológicos que justificam a proposta?  |        | X                   |            |
|                                 | São apresentadas as contribuições teóricas, práticas ou sociais que justificam a proposta?   | X      |                     |            |
|                                 | 5. REQUISITOS PRINCIPAIS DO PROBLEMA A SER TRABALHADO<br>Os requisitos funcionais e não funcionais foram claramente descritos?   | X      |                     |            |
|                                 | 6. METODOLOGIA<br>Foram relacionadas todas as etapas necessárias para o desenvolvimento do TCC?  | X      |                     |            |
|                                 | Os métodos, recursos e o cronograma estão devidamente apresentados e são compatíveis com a metodologia proposta?   | X      |                     |            |
|                                 | 7. REVISÃO BIBLIOGRÁFICA (atenção para a diferença de conteúdo entre projeto e pré-projeto)<br>Os assuntos apresentados são suficientes e têm relação com o tema do TCC? | X      |                     |            |
|                                 | As referências contemplam adequadamente os assuntos abordados (são indicadas obras atualizadas e as mais importantes da área)?   |        | X                   |            |
| ASPECTOS METODOLÓGICOS          | 8. LINGUAGEM USADA (redação)<br>O texto completo é coerente e redigido corretamente em língua portuguesa, usando linguagem formal/científica?                            | X      |                     |            |
|                                 | A exposição do assunto é ordenada (as ideias estão bem encadeadas e a linguagem utilizada é clara)?  | X      |                     |            |



UNIVERSIDADE REGIONAL DE BLUMENAU  
CENTRO DE CIÊNCIAS EXATAS E NATURAIS  
DISCIPLINA: TRABALHO DE CONCLUSÃO DE CURSO I  
CURSO: CIÊNCIA DA COMPUTAÇÃO - BCC

**ATA DA DEFESA: BANCA DO PRÉ-PROJETO**

Venho, por meio deste, manifestar minha avaliação sobre a **apresentação** do Pré-Projeto de TCC

realizado pelo(a) acadêmico(a), Henrique Jose Wilbert no **SEGUNDO SEMESTRE DE 2021**, com o título UTILIZAÇÃO DE CLUSTERIZAÇÃO PARA AUXÍLIO EM TOMADA DE DECISÃO A PARTIR DE DADOS DE VAREJO, sob orientação do prof(a). Aurélio Faustino Hoppe.

A referida apresentação obteve a seguinte nota:

| Componente da Banca                                  | Nota<br>(de 0 a 10) |
|--|---------------------|
| <b>Professor(a) Avaliador(a):</b><br>Andreza Sartori | <b>10,0</b>         |

**ATENÇÃO.** A nota acima se refere somente a apresentação do pré-projeto e vai ser repassada para o aluno (orientando). Favor preencher os campos acima e enviar por e-mail ao professor de TCC1 ([dalton@furb.br](mailto:dalton@furb.br)). Não passar o arquivo com as anotações da revisão já enviado ao professor de TCC1 para o orientando e nem para o professor orientador. Após o professor de TCC1 receber esta ata preenchida, o professor de TCC1 vai disponibilizar para o orientando/orientador os arquivos com as revisões. Caso julgue necessário fazer mais alguma consideração relacionada ao pré-projeto ou a defesa, favor usar o espaço abaixo.

Observações da apresentação:

| CURSO DE CIÊNCIA DA COMPUTAÇÃO – TCC |             |                      |
|--------------------------------------|-------------|----------------------|
| ( X ) PRÉ-PROJETO                    | ( ) PROJETO | ANO/SEMESTRE: 2021/2 |

## UTILIZAÇÃO DE CLUSTERIZAÇÃO PARA AUXÍLIO EM TOMADA DE DECISÃO A PARTIR DE DADOS DE VAREJO

Henrique José Wilbert

Prof. Aurélio Faustino Hoppe – Orientador

Prof. Christian Daniel Falaster – Coorientador

### 1 INTRODUÇÃO

Com a evolução da tecnologia de informação a partir dos anos 80 e início dos anos 90, várias grandes empresas adotaram sistemas de gerenciamento na forma de softwares Enterprise Resource Planning (ERP) (RASHID; HOSSAIN; PATRICK, 2001, p.2). Estes softwares auxiliam em suas rotinas à nível operacional, seja no controle do estoque, fiscal, financeiro, transacional e até recursos humanos. A partir disso, alcançou-se um patamar de eficiência nunca concebido, visto que registros antes realizados em papel e caneta, passaram a ser produzidos automaticamente. Ainda segundo os autores, em paralelo a informatização destes processos, houve também um crescimento da quantidade de dados armazenados referentes à produtos, clientes, transações, gastos e receitas.

Comentado [DSdR23]: desses

Diante deste contexto, avançaram-se também as táticas de marketing direto, como por exemplo, o envio de catálogos por correio, até ofertas altamente objetivas, focando em indivíduos selecionados, cujas informações transacionais estavam presentes na base de dados. Percebeu-se que o foco das relações empresa-cliente não está em clientes novos, e sim em clientes já existentes nas bases de dados, visto que o custo para adquirir um cliente novo através de publicidade é muito maior que o custo de alimentar uma relação já existente (PETRISON; BLATTBERG; WANG, 1997, p. 119, tradução nossa).

Comentado [DSdR24]: Remover.  
Só em citação direta.

Segundo Reinartz, Thomas e Kumar (2005, p.77), quando empresas tratam os gastos entre aquisição e retenção de clientes, destinar menos recursos para a retenção impactará em uma lucratividade menor a longo prazo, comparando-se a investimentos menores em aquisição de clientes. Ainda segundo os autores, no conceito de relações de retenção, atribui-se grande ênfase à lealdade e lucratividade de um cliente, sendo lealdade a tendência do cliente comprar e a lucratividade, a medida geral de quanto lucro um cliente traz à empresa através de suas compras.

De acordo Nguyen, Sherif e Newby (2007, p.114) com o avanço da gerência das relações com clientes foram abertas novas vias pelas quais sua lealdade e lucratividade pode ser cultivada, atraindo uma crescente demanda por parte de empresas, visto que a adoção destes meios permite que as organizações melhorarem seu serviço ao consumidor, consequentemente gerando renda. Com isso, diferentes ferramentas acabam sendo utilizadas, como sistemas de recomendação que, geralmente em ramos e-commerce, levam em conta várias características pertinentes ao comportamento do cliente, construindo um perfil próprio que será utilizado para realizar a recomendação de um produto que talvez seja de seu interesse. Outra ferramenta pertinente à lucros e lealdade é a segmentação, que visa separar uma única e confusa massa de clientes em segmentos homogêneos em termos de comportamento, permitindo o desenvolvimento de campanhas e estratégias de marketing especializadas à cada grupo de acordo com suas características (TSIPTSIS; CHORIANOPOULOS, 2009, p.4).

Em relação a segmentação de clientes, algumas métricas tornam-se relevantes nos contextos aos quais estão inseridas. Segundo Kumar (2008, p.29), o modelo Recency-Frequency-Monetary (RFM), é utilizado em empresas de venda por catálogo, enquanto empresas de high-tech tendem a usar Share of Wallet (SOW) para implementar suas estratégias de marketing. Já o modelo Past Customer Value (PCV), geralmente é utilizado em empresas de serviços financeiros. Dentre os modelos citados, o RFM é o que possui maior facilidade de aplicação em diversas áreas de comércio, varejo e supermercados, visto que são necessários apenas os dados transacionais dos clientes (vendas), dos quais são obtidos os atributos de recência (R), frequência (F) e monetário (M).

Comentado [DSdR25]: Itálico.

Comentado [DSdR26]: Itálico.

A partir desses dados, segundo Tsipitsis e Chorianopoulos (2009, p.335), é possível detectar bons clientes a partir das melhores pontuações de RFM. Se o cliente efetuou uma compra recentemente, seu atributo R será alto. Caso ele compre muitas vezes ao longo de um determinado período, seu atributo F será maior. E, por fim, caso seus gastos totais forem significativos, terá um atributo M alto. Ao categorizar o cliente dentro destas três características, é possível obter uma hierarquia de importância, tendo os clientes que possuem valores RFM altos no topo, e clientes que possuem valores baixos na base. Apesar destas

Comentado [DSdR27]: Recência

Comentado [DSdR28]: Frequência

Comentado [DSdR29]: Monetário

vantagens, o modelo padrão original é um tanto quanto arbitrário, segmentando os clientes em quintis, cinco grupos com 20% dos clientes, não atentando-se às nuances e todas as interpretações que a base de clientes pode possuir. Além disso, o método também pode produzir uma grande quantidade de grupos, que por muitas vezes, não representam significativamente os clientes de um estabelecimento, e caso o método de quintis seja utilizado, 125 grupos serão criados. Outro ponto a se observar é a variada gama de interpretações que os atributos RFM podem ter em relação aos tipos de atividades dos estabelecimentos, sendo necessário a adaptação do modelo para cada empresa.

Diante deste cenário, este trabalho propõe a criação de um artefato computacional que utilize o modelo RFM em conjunto com diferentes algoritmos de clusterização ao invés de quintis para segmentar clientes, extraíndo de maneira automática as informações de bases de dados, sendo aplicado ao contexto de diversas empresas de varejo, atacado e comércio, visando adequar-se dinamicamente à suas eventuais diferenças de comportamento nos clientes.

## 1.1 OBJETIVOS

O objetivo deste trabalho é desenvolver um artefato computacional de auxílio à segmentação de clientes a partir de múltiplas bases de dados utilizando o modelo RFM.

Os objetivos específicos são:

- implementar e testar diferentes algoritmos de clusterização;
- disponibilizar um mecanismo de visualização dos agrupamentos;
- avaliar a qualidade em relação ao agrupamento dos clientes.

## 2 TRABALHOS CORRELATOS

Neste capítulo serão apresentados os trabalhos similares ao proposto. Na seção 2.1 é apresentado o trabalho de Gustriansyah, Suhandi e Antony (2020), que consiste na aplicação do modelo *Recency Frequency Monetary* (RFM) para clusterização de produtos utilizando K-means, tendo como foco otimizar o número de clusters através de índices de validação. A seção 2.2 descreve o modelo de segmentação denominado *Length, Recency, Frequency, Monetary and Periodicity* (LRFMP) proposto por Peker, Kocyigit e Eren (2017), ao qual considera a longevidade e a periodicidade. Por final, na seção 2.3 detalha-se o modelo R+FM, sendo uma versão modificada do modelo original que foi aplicada em clientes de uma empresa de e-commerce (TAVAKOLI et al., 2018).

### 2.1 CLUSTERING OPTIMIZATION IN RFM ANALYSIS BASED ON K-MEANS

Gustriansyah, Suhandi e Antony (2020) utilizaram o algoritmo de clusterização K-means para agrupar 2.043 produtos de uma farmácia visando otimizar o manuseio de estoque. Foram utilizadas três características de acordo com o modelo *Recency Frequency Monetary* (RFM) para a separação dos produtos, levando em consideração dados transacionais capturados num período de um ano. O atributo recência classificou os produtos através da última venda realizada num intervalo de 1 a 364 dias. A frequência estabelece a quantidade de transações em que o produto ocorreu, variando num intervalo de 1 a 14.872. Já o atributo monetário, refere-se ao valor total proveniente das vendas acumuladas do produto, sendo definido num intervalo entre 1250 e 1.151.952.500 Rupias Indonésias (Rp.).

Após a atribuição de valores RFM aos produtos, foram utilizados oito índices de validação do melhor número de clusters: *Elbow Method* (EM), que calcula a variação intra-cluster conforme são aumentados os clusters e conclui que o melhor número é aquele que está no cotovelo (elbow) da curva. *Silhouette Index* (SI) que resulta em uma nota de -1 a 1 que indica a quão adequada é a classificação de um objeto dentro de um cluster em comparação aos outros. *Calinski-Harabasz Index* (CHI) que também mede a adequação da quantidade de clusters levando em conta a dispersão entre e intra clusters. *Davies-Bouldin Index* (DBI) que calcula as similaridades entre clusters levando em conta as distâncias e tamanhos dos clusters, quanto menor este índice melhor a separação entre os clusters. *Ratkowski Index* (RI) que é baseado na média da soma dos quadrados dos dados entre clusters e a soma total dos quadrados de cada dado dentro de um cluster, dentre as quantidades calculadas escolhe-se a que obtém um maior índice. *Hubert Index* (HI) que é um método visual que indica a quantidade preferida através de um pico no gráfico e é calculado pelo coeficiente de correlação entre matrizes de distância. *Ball-Hall Index* (BHI) definido pela média da distância dos itens com os respectivos centroides do cluster, onde no gráfico o ponto de quantidade de clusters com maior diferença do anterior é sugerido. *Krzanowski-Lai Index* (KLI), que propõe índices internos definidos pelas diferenças entre matrizes de dispersão, e aponta a melhor quantidade de clusters pelo maior número gerado ao realizar a equação com quantidade k. Constatou-se que a maioria deles indicou que o melhor número de clusters seria 3, com base nas condições de interpretação de cada índice explicadas anteriormente.

**Comentado [DSdR30]:** Sei que ainda pode não se ter certeza o que será implementado. Mas sugiro utilizar algo mais expressivo, mesmo que depois com o desenvolvimento do TCC se altere.

**Comentado [DSdR31]:** Clientes. Também extraíndo

**Comentado [DSdR32]:** Menos genérico.

**Comentado [DSdR33]:** Nesta seção

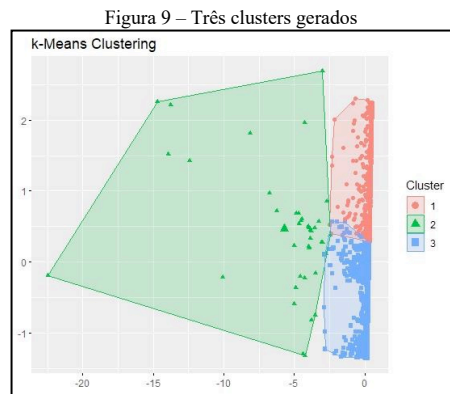
**Comentado [DSdR34]:** Espaço em branco.

**Comentado [DSdR35]:** fim

**Comentado [DSdR36]:** 1.250

Segundo Gustriansyah, Suhandi e Antony (2020) nos testes para verificar os clusters gerados, utilizou-se a equação de variância (R). Sendo “R” o valor da divisão entre a distância média dos dados cluster (distância intra-cluster) pela distância média dos dados em outros clusters (inter-cluster). O valor médio alcançado para R foi de 0.19113, sendo que quanto mais próximo de zero, maior a similaridade entre os membros dentro de cada cluster.

Gustriansyah, Suhandi e Antony (2020) utilizaram o software R Programming para gerar a clusterização, resultando na visualização demonstrada na Figura 7, tendo dois clusters com densidade maior e um cluster mais disperso. Também se observa que o cluster em verde possui uma maior variância entre os próprios dados, enquanto os outros dois clusters possuem uma menor diferença interna.



Fonte: Gustriansyah, Suhandi e Antony (2020).

Segundo Gustriansyah, Suhandi e Antony (2020) também foram adquiridos os valores médios RFM para cada cluster, conforme apresenta a Tabela 1, sendo que o cluster número 3 possui a maior média dos três atributos, e o cluster 1 possui a menor média dos três. É possível identificar um intervalo entre os valores médios de cada atributo, indicando uma diferença significativa inter-cluster.

Tabela 7 – Valores médios de RFM de cada cluster

| Cluster | Recency  | Frequency  | Monetary<br>(in thousands) |
|---------|----------|------------|----------------------------|
| 1       | 75.8167  | 3,436.744  | 3,089,608                  |
| 2       | 224.3947 | 13,013.333 | 76,920,847                 |
| 3       | 331.9681 | 107.418    | 286,927,000                |

Fonte: Gustriansyah, Suhandi e Antony (2020).

Gustriansyah, Suhandi e Antony (2020) concluem que o método gerou clusters com alta similaridade em relação aos dados existentes, apresentado uma segmentação mais objetiva quando comparado ao modelo RFM tradicional no qual os dados são divididos igualmente em cinco segmentos (20% dos dados para a cada segmento). Além disso, os autores também sugerem como extensões a utilização de outros métodos para a comparação, como *Particle Swarm Optimization* (PSO), que é um método computacional que otimiza soluções para uma equação de uma certa medida de qualidade, medianoide (centroide que são parte do conjunto de dados) ou *maximizing-expectancy*, que é um método iterativo para encontrar estimativas de parâmetros para modelos estatísticos com variáveis não observadas.

## 2.2 LRFMP MODEL FOR CUSTOMER SEGMENTATION IN THE GROCERY RETAIL INDUSTRY: A CASE STUDY

Peker, Kocyigit e Eren (2017) propuseram o modelo *Length, Recency, Frequency, Monetary* (LRFM) denominado *Length, Recency, Frequency, Monetary and Periodicity* (LRFMP) para classificar dados reais de 16.024 clientes de mercados de uma franquia na Turquia. Para isso, utilizou-se o algoritmo K-means para segmentar os clientes e três índices de validação de clusters para a otimização das suas quantidades, *Silhouette Index* (SI), *Calinski-Harabasz Index* (CHI) e *Davies-Bouldin Index* (DBI). Após a segmentação dos dados, verificou-se estratégias de gerenciamento e relações com os clientes para aumentar a lucratividade, como tratamento preferencial para clientes importantes, implementação de cartões

Comentado [DSdR37]: Remover.

fidelidade para aumentar a frequência de compra de clientes não costumam comprar com frequência, promoções voltadas para clientes incertos com sua escolha de local de compra, dentre outras estratégias.

Primeiramente, Peker, Kocyigit e Eren (2017) adaptaram o modelo LRFM, incluindo o parâmetro *periodicity* (periodicidade), pois a análise dos dados foi realizada a partir do histórico de compras em supermercados, que são estabelecimentos com alto número de visitas, tornando importante a regularidade nos padrões de visita e compra. Peker, Kocyigit e Eren (2017) definem a periodicidade como a regularidade das visitas de um determinado cliente. Sendo atribuída como o desvio padrão dos tempos inter-visita do cliente (quantia de dias entre duas visitas consecutivas). Se um cliente possui valores baixos de periodicidade, significa que este realiza visitas ou compras em intervalos fixos, podendo caracterizá-lo como cliente regular. Além disso, os autores também modificaram o atributo de recência, transformando-o na média das diferenças entre a data das três últimas compras e a data atual, ao invés da simples diferença entre a data da última compra e a data atual estabelecida no modelo RFM padrão.

**Comentado [DSdR38]:** Evitar iniciar frase c/ gerúndio.

Após adquirir os atributos LRFMP dos dados transacionais dos clientes, Peker, Kocyigit e Eren (2017) aplicaram um método de normalização simples nos dados, considerando o intervalo de 0 e 1. Esta normalização foi feita pois os valores LRFMP variam em relação ao intervalo e escala, fato que poderia afetar negativamente a análise dos clusters.

Antes de aplicar a clusterização, Peker, Kocyigit e Eren (2017) utilizaram três índices para validação da quantidade possível de clusters: *Silhouette Index* (SI) que resulta em uma nota de -1 a 1 que indica o quão adequada é a classificação de um objeto dentro de um cluster em comparação aos outros, quanto maior o valor, melhor. *Calinski-Harabasz Index* (CHI) que mede a adequação da quantidade de clusters levando em conta a dispersão entre e intra clusters, um valor alto é preferido. *Davies-Bouldin Index* (DBI) que calcula as similaridades entre clusters levando em conta as distâncias e tamanhos dos clusters, quanto menor este índice melhor será a separação entre os clusters. A partir deles, Peker, Kocyigit e Eren (2017) executaram o algoritmo K-means variando o k de 2 a 9, e os resultados destas iterações foram avaliadas utilizando os três índices. Com base nos resultados, decidiu-se utilizar um número de 5 clusters, pois 2 dos 3 índices sugerem 5 como sendo a quantidade ideal.

Peker, Kocyigit e Eren (2017) utilizaram uma base de dados de uma franquia de mercados que possui mais de dez lojas na cidade de Antália na Turquia. Os dados são compostos por cerca de dois milhões de transações de 16.024 clientes num período de dois anos. Foram removidos os clientes com menos que três compras. Além disso, os autores removeram dados duplicados, transações com valores faltantes assim como, também agregaram as compras dentro de um mesmo dia. Depois dessas operações, a quantidade de clientes caiu para 10.471, sendo aplicado na sequência o K-means. A Tabela 6 demonstra a quantidade de clientes nos clusters, os valores médios de LRFMP para cada cluster. Já na última coluna, aplicou-se uma técnica na qual o atributo do cluster recebe uma seta para cima (↑) caso o seu valor for maior que a média do atributo dos outros clusters, e uma seta para baixo (↓), caso seu valor for menor que a média.

Tabela 8 – Valores médios dos clusters

| Cluster | Sample size | Average L | Average R | Average F | Average M | Average P | LRFMP Scores   |
|---------|-------------|-----------|-----------|-----------|-----------|-----------|----------------|
| 1       | 538         | 633.29    | 39.67     | 175.24    | 24.32     | 4.99      | L↑ R↓ F↑ M↓ P↓ |
| 2       | 4,681       | 564.50    | 90.19     | 33.44     | 31.73     | 31.49     | L↑ R↓ F↑ M↓ P↓ |
| 3       | 1,091       | 482.17    | 301.18    | 5.33      | 34.21     | 159.32    | L↑ R↑ F↓ M↓ P↑ |
| 4       | 818         | 374.01    | 220.24    | 11.85     | 104.18    | 45.81     | L↓ R↑ F↓ M↑ P↑ |
| 5       | 3,343       | 173.70    | 399.79    | 10.34     | 30.14     | 27.85     | L↓ R↑ F↓ M↓ P↓ |
| Average |             | 419.81    | 218.59    | 28.74     | 36.76     | 43.41     |                |

Fonte: Peker, Kocyigit e Eren (2017).

A partir destes resultados, Peker, Kocyigit e Eren (2017) descreveram as características dos grupos. O grupo 1 representa clientes leais de alta contribuição que, apesar de comporem a menor parcela dos clientes (5,14%), possuem a maior contribuição total entre os grupos. Também é possível observar, que este grupo possui a menor periodicidade média de todos, caracterizando estes clientes como regulares. O grupo 2, representando a maior parcela dos clientes (44,70%) foi classificado como clientes leais de baixa contribuição pois apesar de visitar mais frequentemente as lojas, não possuem tanta contribuição quanto o grupo 1. O grupo 3, com tamanho de 10,42%, foi classificado como clientes incertos, pois possui o atributo de longevidade alto e recência também alta, significando que são clientes com longa história de compra, porém sem muitas compras recentes, vale notar que este grupo possui o maior valor de periodicidade de todos os grupos, caracterizando-o como um grupo de clientes sem rotina de compra definida. O grupo 4 e 5 foram classificados como clientes perdidos, visto que possuem poucas compras recentes, baixa frequência, e baixa longevidade, denotando um cliente que tem uma pouca interação com a franquia. O

**Comentado [DSdR39]:** Recentes. Vale



grupo 4, contendo uma pequena parcela de 7,81% dos clientes, gasta consideravelmente mais, logo foi classificado como contribuição alta, e o 5, cuja parcela é 31,93%, classificado como contribuição baixa.

A partir desta classificação, Peker, Kocyigit e Eren (2017) estabeleceram estratégias para cada grupo de clientes, como tratamento especial (vagas de estacionamento preferencial, presentes de aniversário, filas preferenciais) para clientes do grupo 1, de maneira a não perder a relação leal com a loja. Para os grupos 3, 4 e 5 cuja frequência é baixa, foi sugerida a adoção de programas de cartão fidelidade para aumentar a frequência deles. Para clientes incertos como no grupo 3, aplicou-se descontos e promoções, de maneira a incentivar os clientes, supostamente sensíveis aos preços, a voltar sua atenção à franquia. Para clientes perdidos, sugeriu-se uma análise mais profunda sobre o motivo da perda, como análise de feedback, inferência de motivos, dentre outros.

Por fim, Peker, Kocyigit e Eren (2017) concluem que o estudo contribuiu com a proposta de um novo modelo RFM, que possibilita uma análise mais profunda que o seu modelo original, visto que as características utilizadas e modificadas permitem uma melhor definição do comportamento de cada cliente. Outra contribuição foi a adição do atributo periodicidade (P) no modelo que, ao contrário do modelo RFM padrão, permite identificar se os clientes de um grupo variam em sua rotina de compras. Outra melhoria apontada é a modificação do atributo de recência, que uma vez calculado como uma média, permite uma caracterização mais precisa que o atributo R comumente utilizado. Uma das limitações destacadas pelos autores é a localidade do estudo, pois foi realizado somente com dados originários de uma cidade, sendo que o comportamento de clientes pode variar de acordo com as diferentes localidades onde é feita a análise. A partir disso, sugere-se uma análise mais ampla contemplando outros locais. Outra sugestão feita por Peker, Kocyigit e Eren (2017) é a adição de novos atributos ao modelo, como a quantidade de produtos comprados, quantidade de produtos perecíveis e não perecíveis comprados, a fim de promover uma interpretação mais profunda do comportamento.

### 2.3 CUSTOMER SEGMENTATION AND STRATEGY DEVELOPMENT BASED ON USER BEHAVIOR ANALYSIS, RFM MODEL AND DATA MINING TECHNIQUES: A CASE STUDY

Tavakoli *et al.* (2018) desenvolveram o modelo RFM, denominado “R+FM”, sendo utilizado em conjunto com o algoritmo de clusterização K-means para segmentar 3 milhões de clientes da maior empresa de *E-commerce* do Oriente Médio. Além disso, o modelo de segmentação foi comparado com o utilizado pela empresa, sendo aplicado em uma campanha de Short Message Service (SMS) focada em aumentar os ganhos de cada segmento.

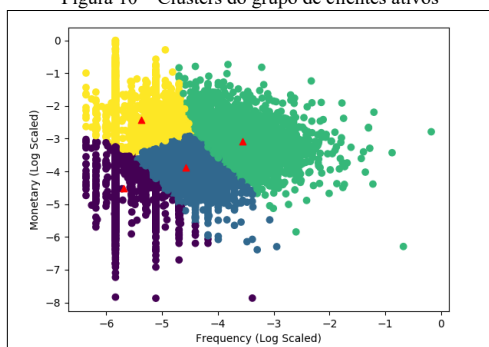
Tavakoli *et al.* (2018) defendem a utilização de um novo modelo de caracterização de clientes, argumentando que o modelo ideal necessita adaptar-se às mudanças de comportamento dos clientes, possuir certa independência de supervisão, levando em consideração a similaridade dos comportamentos dos clientes e a relação entre os atributos Frequência e Monetário. Para isso, construiu-se uma variante do modelo Recency, Monetary, Frequency (RFM) denominado de R+FM, que possui o atributo de recência separado dos demais, utilizando uma segmentação à parte do modelo FM. Os autores separaram os clientes em 3 grupos: os que compraram recentemente (cuja última compra foi dentro de 90 dias), denominados de ativos, clientes que compraram em um passado recente (cuja última compra foi entre 90 e 365 dias), denominados de expirando e por fim, os clientes que não compraram por um longo tempo (cuja última compra foi a mais de 365 dias), denominados de expirados. Para o atributo de frequência, Tavakoli *et al.* (2018) atentaram-se especialmente com a data da primeira compra, pois acreditam que a frequência tem uma importância maior conforme sua recência. Logo, definiram a frequência como a quantidade de compras dividida pela quantidade de dias desde a primeira compra, utilizando também uma função exponencial de decaimento, que efetivamente atribui um peso maior para anos mais recentes, sendo cada ano duas vezes mais pesado que o ano anterior. Como atributo monetário, estabeleceu-se a média dos valores das compras de um cliente, visto que um valor de soma total de compras, segundo os autores, estaria encorajando duas vezes os clientes.

Após a definição do modelo, Tavakoli *et al.* (2018), balancearam a relação entre frequência e monetário, criando a quarta característica que é definida pela combinação linear dos dois atributos, que nada mais é que a soma de cada atributo ponderada pelo peso de cada um. No tratamento dos dados, utilizou-se a técnica de remoção de *outliers* (clientes que não se encaixam no padrão normal) que não se encontram dentro dos intervalos interquartis, que são os intervalos que possuem os dados que pertencem à tendência média do conjunto de dados em geral. Também foram escalados os atributos de frequência e monetário para que seus intervalos sejam iguais, sendo aplicada a normalização *min-max*, que transforma os valores para estarem dentro do intervalo entre 1 e 0. Como os dados monetários e de frequência tratados possuem uma característica de cauda longa, fenômeno estatístico onde os dados são distribuídos de forma

decrecente, foi aplicada uma transformação logarítmica para normalizar a distribuição, visto que a quantidade de valores baixos é muito alta, podendo atrapalhar a análise.

Como existem duas segmentações (R e FM), Tavakoli *et al.* (2018) estabeleceram segmentos FM para cada segmento R, resultando nos seguintes grupos: para clientes ativos existem os grupos de alto valor, médio valor com alto monetário, médio valor com alta frequência e baixo valor. Para clientes que estão expirando existem os grupos de alto, médio e baixo valores, sendo aplicado também para clientes expirados. Estes grupos foram definidos por Tavakoli *et al.* (2018) com ajuda da empresa de *E-commerce* Digikala. A partir disso, aplicou-se o K-means com  $k=4$  para o grupo de clientes ativos,  $k=3$  para o grupo de clientes expirando e  $k=3$  para o grupo de clientes expirados, resultando em um total de 10 clusters. Na Figura 8 são identificados os clusters gerados somente a partir do grupo de clientes ativos, organizados em um gráfico de valor monetário por frequência. Sendo o cluster de cor verde composto pelos clientes de alto valor, o cluster de cor amarela composto pelos clientes de médio valor com alto monetário, o cluster de cor azul composto pelos clientes de médio valor com alta frequência, e por fim, o cluster de cor roxa composto pelos clientes de baixo valor.

Figura 10 – Clusters do grupo de clientes ativos



Fonte: Tavakoli *et al.* (2018).

Após a geração dos grupos, Tavakoli *et al.* (2018) discorrem sobre possíveis estratégias para cada segmento. Sugerindo um maior foco em clientes ativos com valor médio e baixo, bem como a manutenção de clientes já valiosos. Os autores também enfatizam a importância em recuperar os clientes do grupo expirando, cuja chance de retorno não é tão baixa quanto o grupo expirado, que por si só requer uma estratégia especial de reengajamento dos clientes à empresa.

Além da elaboração de estratégias, Tavakoli *et al.* (2018) implementaram uma campanha de SMS focada somente no segmento de clientes ativos (recência abaixo de 90 dias), pois a empresa já tinha realizado outras campanhas em clientes ativos anteriormente. Nesta campanha cada cliente foi apresentado com um *Voucher* condizente com o segmento ao qual o cliente pertencia. Para clientes ativos com valor alto foi oferecido um desconto de 10 por cento com um desconto máximo de até \$20, com o objetivo de manter a lealdade destes clientes. Para clientes ativos de valor médio com alto valor monetário, foi oferecido um desconto de 10 por cento com um desconto máximo de até \$10, o valor foi menor pois o objetivo era aumentar frequência de compra destes clientes, que em tese já gastam bastante. Para clientes ativos de valor médio com alta frequência foi oferecido um desconto de 10 por cento com um desconto máximo de até \$20 para vendas que custem mais que \$50 (que é o valor médio gasto por este segmento), incentivando assim uma compra de maior valor. Por fim, para clientes ativos de baixo valor, foi oferecido um desconto de 10 por cento com um desconto máximo de até \$20, com o objetivo de converter estes clientes em mais leais.

Para a análise da campanha, Tavakoli *et al.* (2018) selecionaram aleatoriamente 20% dos clientes de cada segmento para compor um grupo de controle, cujos *Vouchers* não foram enviados. Este grupo de controle foi comparado com os outros grupos da campanha para obter um valor de referência do aumento do valor monetário após sua conclusão. Os resultados alcançados por Tavakoli *et al.* (2018) podem ser observados no Quadro 1, ao qual percebe-se um aumento de \$14,30 na média monetária dos clientes ativos com alta frequência, mais do que o aumento de \$3,20 sofrido pelo grupo de controle. É possível também observar que a média monetária de todos os grupos alvo da campanha aumentou consideravelmente, enquanto o grupo de controle aumentou pouco ou até diminuiu, indicando uma efetividade no objetivo da campanha.

Comentado [DSdR40]: Evitar iniciar frase c/ gerúndio.

| Quadro 1 – Dados da média monetária do grupo de controle e grupo de campanha |                                  |                        |                |                 |                |
|--|----------------------------------|------------------------|----------------|-----------------|----------------|
| Recency  | Segment                          | Average Monetary (USD) |                |                 |                |
|  |                                  | Control Users          |                | Campaign Users  |                |
|  |                                  | Before Campaign        | After Campaign | Before Campaign | After Campaign |
| Active   | High Value                       | 74.2                   | 73.8           | 88.2            | 89.2           |
|  | Medium Value with High Monetary  | 100.2                  | 97.6           | 104.6           | 105.2          |
|  | Medium Value with High Frequency | 32                     | 35.2           | 35.4            | 49.7           |
|  | Low Value                        | 50.7                   | 53.2           | 56.4            | 65.2           |

Fonte: Tavakoli *et al.* (2018).

Tavakoli *et al.* (2018) concluem que houve uma melhora no desempenho da campanha lançada em comparação com as anteriores, indicando ainda, que elas obtinham uma taxa de compra de 0,1 por cento, sendo que a campanha lançada para validação do modelo obteve uma taxa de 1 por cento, cerca de dez vezes mais efetivo. Os autores justificam esta melhora ao processo de segmentação do modelo, resultando em clusters mais significativos, facilitando a aplicação de vouchers específicos. Por fim, Tavakoli *et al.* (2018) sugerem o melhoramento da definição do atributo de recência de forma que seja mais útil ao time de marketing. Também recomendam calcular o Customer Lifetime Value (CLV), que atribui o valor vitalício à cada segmento e cliente, de forma a quantificar o valor que um cliente pode proporcionar à empresa.

### 3 PROPOSTA DO PROTÓTIPO

Esse capítulo visa apresentar a justificativa para a elaboração deste trabalho, os requisitos que serão seguidos e a metodologia que será utilizada. Será apresentada também uma breve revisão bibliográfica das principais áreas de estudo que serão exploradas, bem como os principais termos utilizados.

#### 3.1 JUSTIFICATIVA

No Quadro 2 é apresentado um comparativo entre os trabalhos correlatos. As linhas representam as características relevantes e as colunas representam os trabalhos.

Quadro 2 – Comparativo entre os trabalhos correlatos

| Características                           | Correlatos | Gustriansyah, Suhandi e Antony (2020)   | Peker, Kocyigit e Eren (2017)                         | Tavakoli <i>et al.</i> (2018)                      |
|---|------------|---|---|--|
| Alvo da clusterização                     |            | Produtos                                | Clientes  | Clientes   |
| Modelo utilizado                          |            | RFM                                     | LRFMP   | R+FM   |
| Objetivo da segmentação                   |            | Gerenciamento de estoque                | Gerenciamento das relações com cliente                | Gerenciamento das relações com cliente             |
| Algoritmo de clusterização utilizado      |            | K-means                                 | K-means   | K-means  |
| Foco metodológico                         |            | Otimização de k com diferentes métricas | Formulação de um modelo novo e análise dos resultados | Formulação de um modelo novo e campanha de ofertas |
| Número de dados (clientes/produtos)       |            | 2.043                                   | 16.024  | ~3.000.000   |
| Quantidade de índices para validação de k |            | 8                                       | 3   | -  |
| Quantidade de clusters gerados            |            | 3                                       | 5   | 10   |
| Inferências sobre os dados                |            | -                                       | Sim   | Sim  |

Fonte: elaborado pelo autor.

A partir do Quadro 2, pode-se observar que Gustriansyah, Suhandi e Antony (2020) clusterizaram produtos de uma base de dados utilizando o modelo RFM padrão. Já Peker, Kocyigit e Eren (2017) optaram pelo desenvolvimento de um modelo novo, considerando a periodicidade (LRFMP). Tavakoli *et al.* (2018) também desenvolveram um novo modelo, ao qual a característica recência foi modificada e separada (R+FM).

Gustriansyah, Suhandi e Antony (2020) tinham como objetivo melhorar o gerenciamento de estoque, prezando por uma segmentação mais conclusiva sobre os produtos, visto que o modelo RFM padrão define segmentos arbitrariamente sem adequar-se às peculiaridades dos dados, enquanto o modelo aplicado através de k-means alcançou uma segmentação com dados altamente similares em cada cluster. Por outro lado, Peker, Kocyigit e Eren (2017) e Tavakoli *et al.* (2018) objetivavam o gerenciamento das relações com os clientes através de estratégias focadas em segmentos, visando aumentar a renda que eles fornecem à empresa. Todos os autores utilizaram o algoritmo K-means, por ser confiável e amplamente difundido. Vale ressaltar que no trabalho de Gustriansyah, Suhandi e Antony (2020), o algoritmo teve um

**Comentado [DSdR41]:** Essa seção

**Comentado [DSdR42]:** Remover. Preâmbulo da próxima seção.

foco metodológico maior, visto que foram utilizados 8 índices de validação para k clusters, visando otimizar a organização dos segmentos.

A quantidade de dados segmentados variou bastante entre os três trabalhos devido aos diferentes contextos de aplicação. Gustriansyah, Suhandi e Antony (2020) tinham 2.043 produtos na base de dados para segmentar, resultando em 3 clusters. Já Peker, Kocyigit e Eren (2017) possuíam o registro de 16.024 clientes de uma rede de padarias, sendo especificados 5 segmentos, obtidos através de uma análise por três índices de validação (Silhouette, Calinski-Harabasz e Davies-Bouldin). Por fim, Tavakoli *et al.* (2018) agruparam dados de 3 milhões de clientes pertencentes à base de dados de um E-commerce do Oriente Médio, resultando em 10 clusters, sendo 3 pertencentes à característica de recência, e os outros 7 distribuídos entre as características de frequência e monetária. Ressalta que Tavakoli *et al.* (2018) testaram o modelo em produção, montando uma campanha que focava no segmento de clientes ativos, visando primariamente aumentar os lucros da empresa, utilizando também um grupo de controle e comparação de renda antes e depois da campanha.

Gustriansyah, Suhandi e Antony (2020) demonstraram a possibilidade da aplicação de RFM fora do uso convencional de segmentação de clientes, e adquiriram clusters com uma variância média de 0.19113. Além disso, os autores sugeriram outras formas de comparação de dados, como Particle Swarm Optimization (PSO), medioides ou até *maximizing-expectancy*. Peker, Kocyigit e Eren (2017) segmentaram clientes de uma rede de mercados na Turquia em “clientes leais de alta contribuição”, “clientes leais de baixa contribuição”, “clientes incertos”, “clientes perdidos de alto gasto” e “clientes perdidos de baixo gasto”. Desta maneira, os autores providenciaram visões e estratégias (promoções, ofertas, regalias) de aumento de renda sobre os comportamentos dos clientes, porém limitaram-se a aplicar em um segmento específico de mercado. Por fim, Tavakoli *et al.* (2018) agruparam clientes de uma empresa de E-commerce com base em sua recência, resultando em clientes “Ativos”, “Expirando” e “Expirados”, e destes segmentos, sucessivamente separados em grupos de “Alto”, “Médio” e “Baixo” valores, validando posteriormente a segmentação através de uma campanha de ofertas para os clientes do grupo “Ativos”.

Todos os trabalhos aqui citados procuraram implementar o modelo RFM num contexto de clusterização por K-means, alterando o modelo e o manejo dos dados de acordo com cada categoria, seja ele produto ou cliente, varejo ou mercado. Com isso, criaram-se atributos e foram modificados alguns já existentes para atender às especificidades de cada contexto, visto que todos os trabalhos focaram em uma só base de dados, inevitavelmente adequando-se às mesmas.

Desta forma, este trabalho demonstra ser relevante, pois almeja aplicar o modelo RFM em conjunto com vários algoritmos de clusterização em forma de um artefato computacional que se adeque à vários contextos (mercado, comércio, varejo etc.), utilizando várias bases de dados reais para testar a validade dos algoritmos utilizados. ~~Vislumbra-se utilizar três índices para validação da qualidade dos clusters (Silhouette, Calinski-Harabasz e Davies-Bouldin). Além disso, deseja-se obter clusters significativos e coerentes com cada segmento de mercado aplicado. Outra contribuição deste trabalho seria no âmbito comercial, com a geração de informações sobre as similaridades de clientes de cada segmento de mercado, podendo auxiliar gestores e administradores de empresas a obter uma visão crítica sobre os comportamentos de clientes ao longo das diferentes bases de dados, podendo também denotar características comuns a todos. Outra relevância seria a utilização deste trabalho em ambiente acadêmico, visto que serão aplicados diferentes algoritmos de clusterização, podendo providenciar informações sobre seus desempenhos e qualidade de agrupamento, além de ser aplicados processos de obtenção, limpeza e transformação de dados.~~

Comentado [DSdR43]: Menos genérico.

Comentado [DSdR44]: Arrumar formato da fonte.

### 3.2 REQUISITOS PRINCIPAIS DO PROBLEMA A SER TRABALHADO

O artefato computacional a ser desenvolvido deverá:

- a) adquirir os dados transacionais de clientes a partir de um banco de dados (Requisito Funcional - RF);
- b) extrair dos clientes as características (recência, frequência e monetária) utilizadas no modelo RFM (RF);
- c) filtrar os clientes sem quantidade de compras relevantes (RF);
- d) normalizar os dados para evitar disparidades nas escalas dos dados, principalmente no atributo monetário (RF);
- e) aplicar três índices de validação (Silhouette, Calinski-Harabasz e Davies-Bouldin) para verificar a qualidade dos clusters (RF);
- f) apresentar em um gráfico 3D os clientes, com sua localização definida pela pontuação do cliente nas características RFM (RF);
- g) utilizar algoritmos de clusterização tais como K-means, mean-shift e DBSCAN (RF);

Comentado [DSdR45]: Menos genérico.

- h) utilizar a linguagem Python para o desenvolvimento (Requisito Não Funcional - RNF);
- i) utilizar o ambiente de desenvolvimento Jupyter Notebook (RNF);
- j) utilizar o banco de dados PostgreSQL para ler os dados das bases utilizadas (RNF).

### 3.3 METODOLOGIA

O trabalho será desenvolvido observando as seguintes etapas:

- a) levantamento bibliográfico: pesquisar trabalhos relacionados e estudar sobre o modelo RFM e suas aplicações, algoritmos de clusterização, métodos de tratamento de dados e índices de validação;
- b) seleção de bases de dados: obter bases de dados de usuários cedidas pela empresa Intelidata Informática, desenvolvedora de software de gestão comercial. Serão selecionadas conforme sua adequação ao objetivo do trabalho, variando em tamanho e segmento de mercado;
- c) definição das características do modelo RFM: definir os atributos utilizados para caracterizar os clientes no modelo RFM;
- d) definição de métricas do modelo RFM: definir as métricas para mensuração e atribuição de pontuação de cada característica no modelo RFM;
- e) definição dos algoritmos de clusterização: pesquisar e escolher o algoritmo de clusterização que realizará o agrupamento das características RFM;
- f) implementação: implementar o artefato computacional de segmentação levando em consideração as etapas (b) até (e), utilizando a linguagem Python;
- g) análise dos clusters: avaliar a qualidade dos clusters gerados a partir dos diferentes algoritmos de clusterização e seus comportamentos em múltiplas bases de dados, aplicando índices de validação e apresentando-os na forma de gráficos.

As etapas serão realizadas nos períodos relacionados no Quadro 3.

Quadro 3 – Cronograma de atividades a serem realizadas

| etapas / quinzenas                          | 2022 |   |      |   |      |   |      |   |      |   |
|---|------|---|------|---|------|---|------|---|------|---|
|   | fev. |   | mar. |   | abr. |   | maio |   | jun. |   |
|   | 1    | 2 | 1    | 2 | 1    | 2 | 1    | 2 | 1    | 2 |
| levantamento bibliográfico                  |      |   |      |   |      |   |      |   |      |   |
| seleção de bases de dados                   |      |   |      |   |      |   |      |   |      |   |
| definição das características do modelo RFM |      |   |      |   |      |   |      |   |      |   |
| definição de métricas do modelo RFM         |      |   |      |   |      |   |      |   |      |   |
| definição dos algoritmos de clusterização   |      |   |      |   |      |   |      |   |      |   |
| implementação                               |      |   |      |   |      |   |      |   |      |   |
| análise dos clusters                        |      |   |      |   |      |   |      |   |      |   |

Fonte: elaborado pelo autor.

### 4 REVISÃO BIBLIOGRÁFICA

Neste capítulo serão brevemente aprofundados os assuntos que servirão de base para a realização deste trabalho. Serão tratados os temas de clustering, modelo RFM e índices de validação de clusters.

Para Cherkassky e Mulier (2007), clustering se trata do problema de separar um conjunto de dados em grupos chamados de “clusters” baseado em alguma medida de similaridade. O objetivo é encontrar um conjunto de clusters dos quais as amostras dentro dos mesmos são mais similares entre si do que quando comparadas com amostras de outros clusters. Existem vários algoritmos para clusterizar dados, especializando-se em situações específicas, como o *Density-Based spatial clustering of applications with noise* (DBSCAN), que é um algoritmo que agrupa dados baseado em sua densidade, utilizando conceitos de pontos núcleo, pontos vizinhos e ruído. Outro método utilizado é o *Hierarchical Clustering*, que constrói uma hierarquia de clusters, podendo aplicar abordagens de cima para baixo, onde as observações parte de um cluster que é dividido em outros menores, ou a abordagem de baixo para cima, onde cada dado inicia como um cluster, e é então agrupado conforme o algoritmo é executado. Por fim, um dos algoritmos mais utilizados é o K-means, cuja orientação é focada em centroides, ou seja, seus clusters são representados por um ponto central (não necessariamente fazendo parte do conjunto de dados) que possui a menor distância possível entre si mesmo e o resto dos dados do cluster.

Segundo Hughes (2011), o modelo RFM é “Um meio antigo e altamente preditivo de determinar quem irá responder e comprar. Um método de codificar clientes existentes. Usado para prever resposta, tamanho médio de pedido, e outros fatores”. Este modelo categoriza geralmente clientes através das características de recência (R), frequência (F) e monetária (M). As métricas utilizadas para medir tais

Comentado [DSdR46]: Menos genérico.

Comentado [DSdR47]: Nesta seção

Comentado [DSdR48]: Maiúsculo.

Comentado [DSdR49]: Maiúsculo.

Comentado [DSdR50]: Maiúsculo.

Comentado [DSdR51]: Maiúsculo.

Comentado [DSdR52]: Maiúsculo.

características podem variar, porém geralmente classificam recência como a quantidade de dias desde a última compra, frequência como a quantidade de compras dentro de um determinado período, e monetária como o total acumulado de todas as vendas realizadas para um cliente. Este modelo é muito utilizado em marketing direto, onde os meios de comunicação são diretamente entre a empresa e o consumidor, realizado através de mídias sociais, e-mail, mensagens SMS ou até pelo correio. No estudo realizado por Verhoef *et al.* (2003), cerca de 90% das empresas questionadas sobre a aplicação métodos de segmentação (como o RFM) afirmam que possuíam como objetivo a seleção de alvos, ou seja, encontrar o segmento de clientes que mais se identificam com a empresa, e 64,4% citaram como objetivo o tratamento diferencial de clientes, com promoções, preços e ofertas especiais.

O índice de Silhouette é utilizado no conceito de clusterização, para analisar a qualidade de um dado localizado em determinado cluster, levando em conta a distância média entre clusters. Com o cálculo deste índice, valores próximos de 1 para um determinado dado em um cluster são considerados bons, valores perto de 0 indicam que o dado está entre clusters e caso o valor seja próximo de -1 significa que provavelmente o dado está no cluster errado. O índice de Calinski-Harabasz é definido como a razão de duas somas calculadas entre todos os clusters: a soma da dispersão intra-clusters e soma da dispersão inter-clusters. Quanto maior seu valor, melhor a performance da clusterização observada. Por fim, o índice de Davies-Bouldin indica a similaridade média entre clusters, levando em conta sua distância e tamanho. Um valor baixo para este índice indica uma melhor separação entre os clusters.

#### REFERÊNCIAS

CHERKASSKY, Vladimir S.; MULIER, Filip. Methods for data reduction and dimensionality reduction. In: CHERKASSKY, Vladimir S.; MULIER, Filip. **Learning from data: concepts, theory, and methods**. 2. ed. Hoboken: Ieee Press, 2007. Cap. 6, p. 191

GUSTRIANSYAH, Rendra; SUHANDI, Nazori; ANTONY, Fery. Clustering optimization in RFM analysis Based on k-Means. **Indonesian Journal Of Electrical Engineering And Computer Science**, [S. l.], v. 18, n. 1, p. 470-477, abr. 2020. Mensal. Disponível em: <http://ijeecs.iaescore.com/index.php/IJECS/article/view/20264>. Acesso em: 02 set. 2021

HUGHES, Arthur Middleton. **Strategic Database Marketing 4e: the masterplan for starting and managing a profitable, customer-based marketing program**. 4. ed. [S. l.]: McGraw-Hill, 2011. 608 p.

KUMAR, V. **Managing Customers for Profit: strategies to increase profits and build loyalty**. Upper Saddle River: Pearson Prentice Hall, 2008. 296 p.

NGUYEN, Thuyuyen H.; SHERIF, Joseph S.; NEWBY, Michael. Strategies for successful CRM implementation. **Information Management & Computer Security**, [S. l.], v. 15, n. 2, p. 102-115, maio 2007. Disponível em: <https://www.emerald.com/insight/content/doi/10.1108/09685220710748001/full/html?journalCode=imcs>. Acesso em: 26 set. 2021.

PEKER, Serhat; KOCYIGIT, Altan; EREN, P. Erhan. LRFMP model for customer segmentation in the grocery retail industry: a case study. **Marketing Intelligence & Planning**, [S. l.], v. 35, n. 4, p. 544-559, 6 maio 2017. Emerald. <http://dx.doi.org/10.1108/mip-11-2016-0210>. Disponível em: <https://www.emerald.com/insight/content/doi/10.1108/MIP-11-2016-0210/full/html>. Acesso em: 07 set. 2021.

PETRISON, Lisa A.; BLATTBERG, Robert C.; WANG, Paul. Database marketing: past, present, and future. **Journal Of Direct Marketing**, [S. l.], v. 11, n. 4, p. 109-125, mar. 1997. Wiley. [http://dx.doi.org/10.1002/\(sici\)1522-7138\(199723\)11:43.0.co;2-g](http://dx.doi.org/10.1002/(sici)1522-7138(199723)11:43.0.co;2-g). Disponível em: [https://onlinelibrary.wiley.com/doi/10.1002/\(SICI\)1522-7138\(199723\)11:4%3C109::AID-DIR12%3E3.0.CO;2-G](https://onlinelibrary.wiley.com/doi/10.1002/(SICI)1522-7138(199723)11:4%3C109::AID-DIR12%3E3.0.CO;2-G). Acesso em: 19 set. 2021.

RASHID, Mohammad A.; HOSSAIN, Liaquat; PATRICK, Jon David. The Evolution of ERP Systems: a historical perspective. In: NAH, Fiona Fui-Hoon. **Enterprise Resource Planning: solutions and management**. Hershey: Irm Press, 2001. p. 35-50. Disponível em: <https://books.google.com.br/books?id=qBcJwDWk4ioC&lpg=PR1&ots=9MrXoQhaRL&dq=Enterprise%20Resource%20Planning%3A%20Solutions%20and%20Management&lr&hl=pt-BR&pg=PR1#v=onepage&q=Enterprise%20Resource%20Planning:%20Solutions%20and%20Management&f=false>. Acesso em: 19 set. 2021.

REINARTZ, Werner; THOMAS, Jacquelyn S.; KUMAR, V. Balancing Acquisition and Retention Resources to Maximize Customer Profitability. **Journal Of Marketing**, [S. l.], v. 69, n. 1, p. 63-79, jan. 2005.

**Comentado [DSdR53]:** Arrumar estilo do formato das referências abaixo.

Usar estilo TF-REFERÊNCIAS ITEM  
Ex.: texto não justificado.

**Comentado [DSdR54]:** Norma ABNT: evento.

**Comentado [DSdR55]:** Abreviar.

**Comentado [DSdR56]:** Negrito.

TAVAKOLI, Mohammadreza *et al.* Customer Segmentation and Strategy Development Based on User Behavior Analysis, RFM Model and Data Mining Techniques: a case study. In: 2018 IEEE 15TH INTERNATIONAL CONFERENCE ON E-BUSINESS ENGINEERING (ICEBE), 15., 2018, Xiam. **Proceedings** [...]. [S.l.]: Ieee, 2018. p. 119-126. Disponível em: [https://www.researchgate.net/publication/330027350\\_Customer\\_Segmentation\\_and\\_Strategy\\_Development\\_Based\\_on\\_User\\_Behavior\\_Analysis\\_RF\\_Model\\_and\\_Data\\_Mining\\_Techniques\\_A\\_Case\\_Study](https://www.researchgate.net/publication/330027350_Customer_Segmentation_and_Strategy_Development_Based_on_User_Behavior_Analysis_RF_Model_and_Data_Mining_Techniques_A_Case_Study). Acesso em: 11 set. 2021.

TSIPTSIS, Konstantinos K.; CHORIANOPOULOS, Antonios. **Data Mining Techniques in CRM**: inside customer segmentation. Chichester: John Wiley & Sons, 2009. 374 p.

VERHOEF, Peter C *et al.* The commercial use of segmentation and predictive modeling techniques for database marketing in the Netherlands. **Decision Support Systems**, [S.l.], v. 34, n. 4, p. 471-481, mar. 2003. Disponível em: <https://liacs.leidenuniv.nl/~puttenpwhvander/library/Others/segmpredmodel-hoekstra.pdf>. Acesso em: 19 set. 2021.

**Comentado [DSdR57]:** Negrito.

# FORMULÁRIO DE AVALIAÇÃO BCC – PROFESSOR TCC I

Avaliador(a): Dalton Solano dos Reis

| ASPECTOS AVALIADOS <sup>1</sup> |  | atende | atende parcialmente | não atende |
|---------------------------------|--|--------|---------------------|------------|
| ASPECTOS TÉCNICOS               | 1. INTRODUÇÃO<br>O tema de pesquisa está devidamente contextualizado/delimitado?   | X      |                     |            |
|                                 | O problema está claramente formulado?  | X      |                     |            |
|                                 | 2. OBJETIVOS<br>O objetivo principal está claramente definido e é passível de ser alcançado?   | X      |                     |            |
|                                 | Os objetivos específicos são coerentes com o objetivo principal?   | X      |                     |            |
|                                 | 3. JUSTIFICATIVA<br>São apresentados argumentos científicos, técnicos ou metodológicos que justificam a proposta?  | X      |                     |            |
|                                 | São apresentadas as contribuições teóricas, práticas ou sociais que justificam a proposta?   | X      |                     |            |
|                                 | 4. METODOLOGIA<br>Foram relacionadas todas as etapas necessárias para o desenvolvimento do TCC?  | X      |                     |            |
|                                 | Os métodos, recursos e o cronograma estão devidamente apresentados?  | X      |                     |            |
| ASPECTOS METODOLÓGICOS          | 5. REVISÃO BIBLIOGRÁFICA (atenção para a diferença de conteúdo entre projeto e pré-projeto)<br>Os assuntos apresentados são suficientes e têm relação com o tema do TCC? | X      |                     |            |
|                                 | 6. LINGUAGEM USADA (redação)<br>O texto completo é coerente e redigido corretamente em língua portuguesa, usando linguagem formal/científica?                            | X      |                     |            |
|                                 | A exposição do assunto é ordenada (as ideias estão bem encadeadas e a linguagem utilizada é clara)?  | X      |                     |            |
|                                 | 7. ORGANIZAÇÃO E APRESENTAÇÃO GRÁFICA DO TEXTO<br>A organização e apresentação dos capítulos, seções, subseções e parágrafos estão de acordo com o modelo estabelecido?  | X      |                     |            |
|                                 | 8. ILUSTRAÇÕES (figuras, quadros, tabelas)<br>As ilustrações são legíveis e obedecem às normas da ABNT?  | X      |                     |            |
|                                 | 9. REFERÊNCIAS E CITAÇÕES<br>As referências obedecem às normas da ABNT?  |        | X                   |            |
|                                 | As citações obedecem às normas da ABNT?  | X      |                     |            |
|                                 | Todos os documentos citados foram referenciados e vice-versa, isto é, as citações e referências são consistentes?  | X      |                     |            |





UNIVERSIDADE REGIONAL DE BLUMENAU  
CENTRO DE CIÊNCIAS EXATAS E NATURAIS  
DISCIPLINA: TRABALHO DE CONCLUSÃO DE CURSO I  
CURSO: CIÊNCIA DA COMPUTAÇÃO - BCC

**ATA DA DEFESA: BANCA DO PRÉ-PROJETO**

Venho, por meio deste, manifestar minha avaliação sobre a **apresentação** do Pré-Projeto de TCC realizado pelo(a) acadêmico(a), Henrique Jose Wilbert no **SEGUNDO SEMESTRE DE 2021**, com o título UTILIZAÇÃO DE CLUSTERIZAÇÃO PARA AUXÍLIO EM TOMADA DE DECISÃO A PARTIR DE DADOS DE VAREJO.

A referida apresentação obteve a seguinte nota:

| Componente da Banca  | Nota<br>(de 0 a 10) |
|--|---------------------|
| <b>Professor(a) Orientador(a):</b><br>Aurélio Faustino Hoppe | <b>10,0</b>         |

A apresentação aconteceu em 25 / 10 / 2021 na sala de reunião virtual do MS-Teams, tendo início às 17 : 30 hs e foi encerrada às 18 : 00 hs.

**ATENÇÃO.** A nota acima se refere somente a apresentação do pré-projeto e vai ser repassada para o aluno (orientando). Favor preencher os campos acima e enviar por e-mail ao professor de TCC1 ([dalton@furb.br](mailto:dalton@furb.br)). Lembro que os arquivos com as anotações das revisões do professor de TCC1 e Avaliador serão enviados para o orientando e professor orientador após o professor de TCC1 receber esta ata preenchida. Caso julgue necessário fazer mais alguma consideração relacionada ao pré-projeto ou a defesa, favor usar o espaço abaixo.

Observações da apresentação: