

## SEGMENTAÇÃO VIA MACHINE LEARNING: PROPOSTA DE CLUSTERIZAÇÃO DE CONSUMIDORES DO E-COMMERCE DE UMA EMPRESA MULTINACIONAL DO VAREJO ESPORTIVO

A. A. FALQUETO<sup>1</sup>, L. C. CEZAR<sup>2</sup>,  
Universidade Federal de Viçosa<sup>1,2</sup>  
ORCID ID: <https://orcid.org/0000-0002-5114-0332><sup>1</sup>  
alice.falqueto@ufv.br<sup>1</sup>

Artigo submetido em 19/02/2021 e aceito em 29/11/2021

DOI: 10.15628/holos.2021.12032

### RESUMO

O objetivo desse artigo é apresentar uma proposta de segmentação da base de consumidores do e-commerce de uma empresa multinacional do varejo esportivo, a partir da clusterização de dados via Machine Learning. Para isso, foi realizado um estudo quantitativo com dados de 526.686 clientes do e-commerce de uma empresa multinacional que atua no Brasil nesse setor. Os dados foram analisados pela análise de cluster, utilizando a

metodologia proposta por Jain, Murty e Flynn (1999). A partir da segmentação atual, limitada ao valor gasto pelo cliente, a nova proposta de segmentação, construída a partir do algoritmo K-means considerou novas variáveis como o número de pedidos nos doze meses anteriores, e seus respectivos tempos de inatividade. O uso desse algoritmo via Machine Learning se mostrou satisfatório, visto que foi possível obter três segmentos válidos que se diferenciavam da segmentação adotada atualmente.

**PALAVRAS-CHAVE:** Segmentação; Aprendizado de Máquina; Clusterização; E-commerce; Varejo

## SEGMENTATION FROM MACHINE LEARNING: CLUSTERING PROPOSAL FOR E-COMMERCE CONSUMERS OF A MULTINATIONAL SPORTS RETAIL COMPANY

### ABSTRACT

The objective of this article is to present a proposal for segmentation of the e-commerce consumer base of a multinational sports retail company, from the clustering of data via Machine Learning. For this, a quantitative study was conducted with data from 526,686 e-commerce customers of a multinational company that operates in Brazil in this sector. The data were analyzed by cluster analysis, using the methodology proposed by Jain, Murty and Flynn (1999). From the current

segmentation, limited to the amount spent by the client, the new segmentation proposal, constructed from the K-means algorithm, considered new variables such as the number of orders in the previous twelve months, and their respective downtime. The use of this algorithm via Machine Learning proved to be satisfactory since it was possible to obtain three valid segments that differed from the segmentation currently adopted.

**KEYWORDS:** Segmentation; Machine Learning; Clustering; E-commerce; Retail

## 1 INTRODUÇÃO

Nos últimos anos observou-se uma grande migração dos negócios para os meios digitais. A praticidade em realizar compras online sem sair de casa, e a possibilidade de comparar preços com facilidade em diferentes lojas tem atraído cada vez mais consumidores para as plataformas de *e-commerce* (Turban & King, 2004). Essa migração em massa para o meio digital foi acompanhada por mudanças em várias áreas da gestão empresarial (Saura, 2021). Não somente os consumidores passaram a ter mais informação sobre as empresas, mas também os negócios começaram a se interessar cada vez mais pelos dados dos seus consumidores, com o intuito de desenvolver novas estratégias de vendas e divulgação (Luz, 2020).

Uma das estratégias que evoluiu com o crescimento do digital foi a segmentação de consumidores. Com um grande volume de potenciais clientes acessando suas plataformas online, as empresas entenderam que deveriam cada vez mais usar diferentes estratégias para atrair e servir melhor a cada tipo de consumidor, de modo a obter vantagem competitiva em relação às escolhas da concorrência (Hooley, Piercy & Nicoulaud, 2011). Diferentes formas de trabalhar com os dados de consumidores vem sendo desenvolvido pelas organizações, com o objetivo de oferecer experiências cada vez mais personalizadas. A segmentação especificamente se propõe a agrupar consumidores com características semelhantes, identificar padrões e entender o comportamento de compra dos clientes, de modo a otimizar as ações para atrair e reter melhor o público-alvo de cada organização (Madeira, Silveira & Toledo, 2015).

Com um grande crescimento no volume dos dados e na velocidade de disseminação da informação na era digital, a estatística associada à tecnologia tem sido amplamente utilizada para desenvolver soluções analíticas (Samuel, 2017). Uma área que tem se destacado e ganhado cada vez mais espaço nas companhias é o *Machine Learning* (ML), que utiliza técnicas computacionais para aprendizado e a construção de sistemas que podem se organizar e adquirir conhecimento de maneira automatizada, com base na solução de problemas anteriores (Monard & Baranauskas, 2005). Dentre as técnicas do ML mais utilizadas para segmentação de consumidores está a clusterização, que consiste em agrupar conjuntos de clientes com características semelhantes com base no seu histórico, e pode fornecer informações valiosas para as companhias sobre o comportamento de compra dos seus consumidores (Anitha & Patil, 2019).

A clusterização de dados para agrupar pessoas tem sido amplamente utilizada no varejo (Fildes et al., 2019). As empresas estão buscando entender melhor, detectar semelhanças e diferenças nos dados dos clientes em todos os aspectos. Prever os comportamentos, propor melhores opções e oportunidades aos clientes tornou-se muito importante para o engajamento dos consumidores (Dogan, Ayçin & Bulut, 2018). No varejo de vestuário e calçados não é diferente. Os profissionais de marketing de moda devem ter sempre em mente as diferenças entre os clientes ao planejar suas estratégias de promoção e comunicação, para identificar quais clientes são valiosos para a empresa, quais clientes precisam de atividades promocionais, etc. (Kaur & Anand, 2018).



Mesmo com o tema ganhando importância em estudos internacionais, no Brasil a discussão sobre segmentação ainda é tímida. Em poucos estudos brasileiros esse assunto é tratado como tema central e a maioria não traz implicações gerenciais confiáveis para serem utilizadas nas organizações, pois não discutem a importância e a eficiência da segmentação na sua vinculação com a estratégia empresarial (Souza & Freitas, 2016).

Desse modo, entendendo as vantagens de uma segmentação de consumidores adequada e a importância de usar a tecnologia como aliada para realização de análises mais robustas, bem como a necessidade de se desenvolver mais trabalhos sobre segmentação no Brasil, o presente trabalho se propõe a responder o seguinte problema: *como pode ser realizada a segmentação da base de consumidores do e-commerce de uma empresa do varejo esportivo a partir da clusterização de dados em Machine Learning?* No intuito de responder tal questão foi conduzida uma pesquisa em uma multinacional do varejo esportivo<sup>1</sup>.

A empresa escolhida para o estudo atua no segmento de calçados, vestuário e equipamentos esportivos. Este estudo teve como foco os consumidores do comércio eletrônico da companhia. Tal investigação se faz necessária, visto que, atualmente, uma das segmentações adotadas pela companhia se baseia no valor gasto por um cliente nos doze meses anteriores, com o intuito de recompensar consumidores mais fidelizados, por meio de ações de marketing específicas para esse segmento e concessão de benefícios. Entretanto, essa classificação acaba sendo não tão benéfica a consumidores que compram produtos com um valor de venda menor, mesmo que esses consumidores tenham uma boa frequência de compra, e com consumidores que compraram pela primeira vez há pouco tempo. Assim, o atual formato acaba por privilegiar consumidores que compraram poucas vezes, mas que adquiriram produtos com um valor de venda mais alto, o que justifica a necessidade de propor uma nova segmentação para a base de clientes do e-commerce da companhia.

## 2 CLUSTERIZAÇÃO DE CONSUMIDORES VIA MACHINE LEARNING

Em conjunto com a evolução do comércio eletrônico e das relações entre empresas e consumidores, também ocorreu nos últimos anos uma aceleração na disseminação da informação (Luz, 2020). Samuel (2017) aponta que houve um crescimento dramático nos volumes de dados, na velocidade, na complexidade e na imprevisibilidade das informações que aumentaram os desafios para os profissionais de todas as áreas. Nesse cenário, o uso da matemática, estatística e tecnologia foram fundamentais para criar soluções analíticas para área de negócios (Yin & Fernandez, 2020).

O grande volume de dados requisitados e administrados pelas organizações na atual era do *Big Data*, tem levado ao uso cada vez mais ampliado da concepção de *business intelligence* como meio interativo para tomada de decisões, de forma mais precisas e, embasadas nas mudanças constantes das preferências de consumo e decisões de compra do público-alvo (Kachamas et al., 2019; Yin & Fernandez, 2020). A análise desse intenso volume de dados sobre consumidores, tem invocado técnicas cada vez mais robustas de inteligência artificial (IA), como modelos que tentam

<sup>1</sup> No intuito de resguardar as políticas estratégicas da organização, seu nome não será divulgado.



identificar padrões para prever comportamentos futuros (Farrokhi et al., 2021; Fuentes et al., 2021; Qannari, 2017).

De acordo com Antonopoulos et al. (2020, p.1) “os métodos de IA podem ser usados para enfrentar vários desafios, como selecionar o conjunto ideal de consumidores para responder, aprender seus atributos e preferências, preços dinâmicos, programação e controle de dispositivos”. A IA tenta simular a inteligência humana para resolução de problemas cada vez mais complexos, no intuito de minimizar vieses subjetivos na interpretação dos dados (Fuentes et al., 2021). Para Ahani et al. (2019) o uso de IA em áreas como a segmentação de consumidores tem se tornado urgente, visto que as segmentações tradicionalmente utilizadas não permitem prever com certo grau de acurácia as mudanças ocorridas nas preferências de consumo, que levam os consumidores a manifestarem suas satisfações e insatisfações por diferentes caminhos (como por exemplo as redes sociais), demandando a criação de algoritmos híbridos que dialoguem com diferentes bases de dados. Nesse sentido, Saura (2021) destaca o uso de técnicas de Aprendizado de Máquina ou *Machine Learning* (ML) para extrair percepções acionáveis e identificar padrões para a construção de estratégias organizacionais.

O ML pode ser vista como uma ramificação da IA, cujo intuito é o desenvolvimento de técnicas computacionais para aprendizado, e a construção de sistemas que possam se organizar e adquirir conhecimento de maneira automatizada. Esse tipo de sistema pressupõe que um computador pode tomar decisões com base em experiências e na solução de problemas anteriores (Monard & Baranauskas, 2005). Inúmeras pesquisas têm evidenciado o papel da ML para resolução de problemas como por exemplo a sofisticação nas análises sensoriais (Fuentes et al., 2021; Qannari, 2017); a previsão de demanda no setor de energia (Antonopoulos et al., 2020) e no varejo (Fildes et al., 2019); as melhorias no processo de CRM (Chagas et al., 2020) e na segmentação de consumidores (Ahani et al., 2019; Farrokhi et al., 2021; Kachamas et al., 2019); o uso para o aperfeiçoamento das técnicas de marketing digital (Saura, 2021); dentre inúmeras outras. Todas compartilham a percepção da importância da ML para melhoria na compreensão dos dados, buscando aperfeiçoar a *performance* da organização a partir de decisões mais assertivas.

O ML pode ser classificado em diversas categorias, como aprendizagem supervisionada, aprendizagem não supervisionada, aprendizagem por reforço, aprendizagem ativa, aprendizagem híbrida, aprendizagem profunda e aprendizagem semi supervisionada (Samuel, 2017). De acordo com Lima, Machado e Lopes (2015), as três principais classificações e suas definições são: 1) Supervisionada, cujo objetivo principal é que o sistema aprenda a prever valores; 2) Não supervisionada, em que o sistema aprende a encontrar padrões e instâncias semelhantes; e 3) Aprendizagem por reforço, em que o sistema de ML aprende a melhorar com base nos resultados.

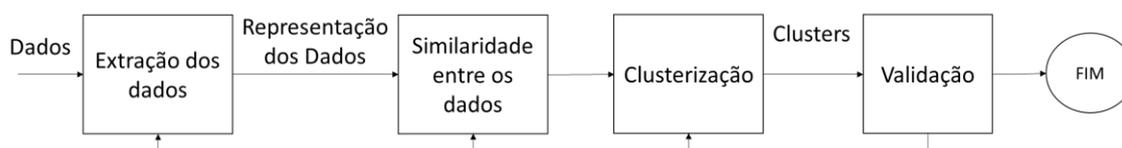
Um das possíveis aplicações do aprendizado não supervisionado é a segmentação de consumidores por meio da formação de grupos, que podem ser chamados de *clusters* (Saura, 2021). Um *cluster* é entendido como um grupo significativo de objetos que possuem características comuns, e são úteis para explorar dados, identificar padrões e fazer análises, oferecendo uma forma de compreender e extrair informações relevantes de grandes conjuntos de dados (Dogan, Ayçin & Bulut, 2018). A segmentação auxilia as empresas a alocarem seus recursos de marketing

de maneira eficaz, possibilitando que grupos específicos de cliente recebam a devida atenção, subsidiando assim a fidelização, por meio da construção de estratégias de longo prazo para a gestão de relacionamentos com os clientes (CRM) (Anitha & Patil, 2019; Chagas et al., 2020).

Existem diversos métodos de clusterização disponíveis atualmente e a principal diferença entre eles é abordagem em alguns aspectos relevantes, que influenciam na qualidade da divisão dos *clusters*, como a representação dos dados, necessidade de ajuste de parâmetros iniciais e medida de similaridade (Jain, Murty & Flynn, 1999). Além disso, muitos algoritmos apresentam dificuldades em alguns aspectos, como por exemplo a impossibilidade de detectar *clusters* de tamanhos e formas diferentes (Saura, 2021). Por outro lado, alguns algoritmos que têm capacidade de formar grupos variados são muito sensíveis a ruídos, e podem acabar fazendo a clusterização de uma forma equivocada (Oliveira, 2008).

Para assegurar a qualidade da clusterização, Jain, Murty e Flynn (1999) propõem que o processo siga algumas etapas, como apresentado na Figura 1

Figura 1.



**Figura 1: Etapas do Processo de Clusterização**

Fonte: Adaptado de Jain, Murty e Flynn (1999)

A partir do exposto na Figura 1 é preciso destacar que na etapa de Extração dos Dados são identificadas as variáveis de interesse a partir do conjunto de dados inicial, sendo feita a formatação para que o algoritmo consiga processar. Na etapa de Similaridade entre os dados é adotada uma medida para que a proximidade entre dois valores possa ser quantificada. Na etapa de Clusterização é escolhido o algoritmo para fazer o agrupamento e é obtido a divisão do conjunto de dados inicial em *clusters*. Por fim, na Validação, os *clusters* são avaliados por meio da comparação com outros algoritmos ou com o uso de índices estatísticos. Oliveira (2008) complementa que, caso na etapa de validação seja encontrado algum problema torna-se necessário redefinir os atributos ou medidas de similaridade definidos anteriormente e, posteriormente, refazer todo o processo.

Apesar de existir uma infinidade de algoritmos de clusterização encontrados na literatura, Ali (2020) destaca que cada um se adapta a diferentes tipos de problemas e tem suas próprias regras por trás do cálculo dos *clusters*. Para cada classificação existem inúmeros algoritmos com diversas variações. Nesse trabalho, foram explorados os algoritmos de clusterização por particionamento, uma das técnicas de agrupamento mais conhecidas em que os dados são divididos uma única vez em um número determinando de *clusters* (Jain, Murty, & Flynn, 1999). De acordo com Oliveira (2008), essa técnica pode ser vantajosa quando a quantidade de dados é grande, tendo em vista que outras técnicas, como a clusterização hierárquica, dividem os dados

gradualmente, obtendo muitas partições, dificultando assim, o armazenamento de todas as possibilidades de divisão.

Segundo Oliveira (2008, p.13) o problema da clusterização por particionamento pode ser definido como: “dado um conjunto de  $n$  dados caracterizados por  $d$  atributos cada, determine uma participação do conjunto inicial em  $K$  clusters”. O objetivo é obter *clusters* diferentes com o máximo de similaridade possível entre os elementos dentro de um *cluster*, e o mínimo de similaridade entre *clusters* diferentes. A escolha de  $K$  depende do problema e interfere na eficiência do algoritmo, por isso é interessante se atentar a várias possíveis partições, de modo a otimizar a clusterização.

### 3 METODOLOGIA

O processo de realização da pesquisa foi norteado predominantemente pela natureza quantitativa e descritiva. Em relação aos meios, a pesquisa pode ser classificada como pesquisa documental uma vez que foram utilizados para a composição da base de dados, relatórios internos que ainda não receberam tratamento analítico (Malhotra, 2012). Os relatórios da empresa, classificados aqui como dados secundários, foram coletados do banco de dados de pedidos, filtrados a partir do histórico de vendas e comportamento de compra dos consumidores, considerando o período de outubro de 2019 a setembro de 2020. Para serem utilizados, os dados brutos precisaram ser processados com o uso da Linguagem de Consulta Estruturada (SQL). Os dados utilizados para o estudo não são públicos. Para coleta, um dos autores desse estudo, obteve autorização para utilização dos dados por estagiar na companhia no período de realização da investigação.

A base de consumidores da empresa é de mais de 1,3 milhões de clientes específicos do *e-commerce*. Foram excluídos dessa população os consumidores que não realizaram nenhum pedido no período analisado e, os consumidores que não poderiam ser avaliados em todos os atributos definidos como relevantes para a investigação. A partir dessa estratificação foram utilizados nesse estudo, dados de 526.686 clientes. Com essa seleção, pode-se caracterizar o processo de definição da amostra como uma amostragem estratificada que, segundo Malhotra (2012) permite ao pesquisador dividir a população em subpopulações a partir da definição de critérios específicos.

A empresa escolhida para o estudo é uma multinacional que atua no segmento de calçados, vestuário e equipamentos esportivos, e tem filiais e subsidiárias espalhadas em 52 países. De acordo com a consultoria especializada em marketing esportivo Sports Value (2018), a companhia faturou cerca de US\$34,4 bilhões no exercício de 2017 globalmente, e cresceu cerca de 111% de 2002 a 2012. No Brasil, a empresa vende diretamente seus produtos aos consumidores por meio das 30 lojas próprias e por meio do seu *e-commerce*, além de fornecer os produtos da sua marca para as maiores redes de varejo esportivo do país. O comércio eletrônico da companhia funciona no Brasil desde 2013, representando assim um canal estratégico da empresa para a gestão do relacionamento com seus clientes.

Para análise dos dados foi utilizada a técnica de análise de *cluster*, que é uma técnica de análise multivariada de dados, cujo objetivo é agrupar dados de acordo com as similaridades entre



eles (Malhotra, 2012). Para operacionalização da análise, seguiu-se a metodologia proposta por Jain, Murty e Flynn (1999), estruturada em quatro etapas como descrito na sequência.

### 3.1 Pré-processamento dos dados

Essa etapa consistiu na realização de uma exploração inicial dos dados para entender a melhor forma de realização da extração das informações que foram utilizadas na clusterização. O intuito dessa análise inicial é verificar a qualidade dos dados, garantindo que sejam identificados clientes únicos, valores nulos, valores negativos e remoção de *outliers*. Para isso, foi feita uma extração do banco de dados da companhia utilizando o software SQL Server.

Após o pré-processamento dos dados, a segmentação utilizada atualmente foi caracterizada e analisada, por meio da obtenção de estatísticas e gráficos. Os valores observados foram o percentual de consumidores em cada segmento e a média de cada um dos atributos definidos para a análise. O cálculo dessas estatísticas iniciais tornou-se importante para a comparação com os *clusters* obtidos.

### 3.2 Medida de similaridade

Para realizar o particionamento dos dados, é necessário utilizar uma medida para medir a proximidade de dois valores. De acordo com Oliveira (2008), a maneira mais comum de calcular a similaridade entre dados representados em duas os três dimensões é a partir da Distância Euclidiana, que é calculada pela equação (1), em que  $x_i$  e  $x_j$  representam os dois dados cuja distância se deseja medir e,  $d$  é o número de atributos do dado. Essa foi a métrica utilizada na aplicação do algoritmo.

$$d(x_i, x_j) = \sqrt{\sum_{k=1}^d (|x_{ik} - x_{jk}|)^2} \quad (1)$$

### 3.3 Clusterização e segmentação dos consumidores

A clusterização foi realizada com o algoritmo K-Means (ou K-médias), um dos mais utilizados para segmentação de clientes segundo Chagas et al. (2020). Essa técnica de partição encontra o número de *clusters* especificado pelo usuário e cada um é representado por seu centroide. A principal vantagem desse algoritmo é a simplicidade que permite que ele seja computacionalmente mais rápido e tenha um bom desempenho em grandes conjuntos de dados. A simplicidade do K-Means também tem outra vantagem que é a exigência de apenas um parâmetro de entrada "K", o que ajuda a diminuir classificações incorretas de dados (Anitha & Patil, 2019). Para analisar o algoritmo de clusterização foi utilizado o software RapidMiner.

### 3.4 Validação

Essa etapa consiste em determinar a relevância e qualidade dos resultados, por meio de índices estatísticos. Para isso, foi utilizada a divisão previamente estabelecida do conjunto de

dados de acordo com segmentação atual, como critério externo de comparação. O intuito é que os *clusters* tenham o tamanho aproximado da segmentação atual, para que estejam alinhados com as ações estratégicas da empresa. Todavia, para que a nova divisão fosse considerada válida, ela deveria se diferenciar de alguma forma da divisão atual. Nesse sentido foram utilizados testes de hipóteses para comparar os atributos nas duas segmentações e determinar se a nova segmentação se diferenciava da antiga. Para se obter inferências quando as variáveis independentes são de natureza não métrica, utilizam-se testes não paramétricos. Como se desejava comparar a diferença na posição de duas populações diferentes, foi utilizado o teste U de Mann-Whitney (Malhotra, 2012). Para aplicação do teste foi utilizado o software Minitab.

## 4 RESULTADOS E DISCUSSÃO

A seguir são apresentados os resultados encontrados para cada passo da metodologia aplicada. Os resultados foram divididos nos seguintes tópicos: 4.1 Avaliação da segmentação vigente; 4.2 Construção da base de dados; 4.3 Clusterização da base de dados; e 4.4 Comparação dos resultados com a segmentação vigente. Cada item é discutido na sequência.

### 4.1 Avaliação da segmentação vigente

Atualmente, a segmentação utilizada pela empresa é baseada no comportamento de compra do consumidor nos últimos doze meses. Ela foi criada com o intuito de valorizar os consumidores que se cadastraram como membros da companhia e se mostraram fidelizados à marca, oferecendo descontos especiais em produtos, acesso antecipado e exclusivo a alguns lançamentos, participação em eventos, planos de exercício personalizados e serviços especiais no momento da compra. Essa segmentação é baseada no valor gasto em dólares por um cliente nos doze meses anteriores e os divide em 3 grupos conforme apresentado na Tabela 1:

**Tabela 1: Composição dos segmentos de acordo com a segmentação utilizada atualmente**

Segmento	Definição - Valor gasto nos últimos 12 meses	% do Total de Compradores
Desenvolvimento em escala	Menos \$100	75%
Potencial de Crescimento	Entre \$100 e \$300	21%
Alto Valor	Mais de \$300	4%

Fonte: Elaborado pelos autores (2021).

Cada um dos segmentos apresentados na Tabela 1 gozam de benefícios específicos como apresentado no Quadro 1:

**Quadro 1: Benefícios oferecidos pela companhia a cada segmento**

	SEGMENTO			
	Desenvolvimento em Escala	Potencial de Crescimento	Alto Valor	Não membros
Produtos exclusivos	X	X	X	

Acesso antecipado a lançamentos			X
Cupom de aniversário	X	X	X
Frete grátis (com valor mínimo no pedido)	X	X	X
Acesso antecipado a promoções			X
Acesso aos aplicativos de atividade física	X	X	X
Acesso a eventos esportivos da companhia	X	X	X
Recompensa pela prática de atividade física	X	X	X
Condições especiais para trocas e devoluções	X	X	X

Fonte: Elaborado pelos autores (2021).

O grupo de “Não membros” representa os consumidores que compram de forma anônima e não possuem cadastro no e-commerce da companhia, sendo que esse grupo não recebe nenhum dos benefícios concedidos aos membros cadastrados. Os grupos “Desenvolvimento em Escala” e “Potencial de Crescimento”, recebem atualmente os mesmos benefícios e o grupo “Alto Valor” recebe duas vantagens a mais: o acesso antecipado a lançamentos e promoções. Além dos benefícios, há uma diferenciação nas ações de marketing para cada um dos grupos, de acordo com suas características de compra. Por exemplo, consumidores do grupo “Desenvolvimento em Escala” são atingidos por campanhas mais voltadas para produtos de Futebol, enquanto “Consumidores de Alto Valor” por consumirem mais produtos da linha Casual, recebem comunicações específicas desses produtos.

O grande problema da classificação utilizada pela companhia atualmente é que ela acaba sendo não tão benéfica aos consumidores que compram produtos com um valor de venda menor, visto que, mesmo que esses consumidores tenham uma boa frequência de compra, podem acabar não atingindo o valor de \$300 dólares em um período de 12 meses. Este valor quando convertido para reais representa uma alta quantia. Da mesma forma, essa segmentação pode acabar privilegiando consumidores que realizaram apenas uma compra nos doze meses anteriores, mas com um produto de valor maior.

Outro problema com a classificação feita dessa forma, é que ela abrange um período de doze meses, sem considerar se os consumidores estão há muito tempo sem comprar ou interagir de alguma forma com a marca, ou mesmo se são novos consumidores (adquiridos recentemente), que certamente possuem um comportamento diferente dos demais.

## 4.2 Construção da base de dados

Para averiguação mais detalhada da segmentação atual, foram inseridos novos parâmetros, considerando as médias para um período de doze meses: 1) Valor gasto no período (Demanda); 2) Número de pedidos (Pedidos) e; 3) Tempo de inatividade. Para a data da última atividade, foram consideradas os seguintes tipos de interação: 1) Compra, 2) Navegação no site, 3) Uso dos

aplicativos para prática de atividade física, e 4) Abertura dos e-mails promocionais enviados pela companhia, escolhendo-se a data mais recente entre todas as interações para a análise. Esses fatores foram escolhidos para levar em conta não só o valor de compra dos produtos, mas também se os consumidores são frequentes e se não estão há muito tempo sem interagir com a marca.

A partir dessa definição realizou-se a segmentação dos consumidores de acordo com os critérios utilizados hoje, ou seja, baseada somente no valor gasto pelos clientes em dólares nos últimos doze meses de acordo com a nomenclatura já utilizada pela empresa. Para cada um dos atributos, foi calculada a média dos segmentos. Os resultados estão apresentados na Tabela 2.

**Tabela 2: Média dos atributos para cada um dos grupos com base na segmentação atual**

Segmento	Demanda	Pedidos	Tempo de Inatividade
Desenvolvimento em escala	\$45	1.19	146
Potencial de Crescimento	\$159	2.33	123
Alto Valor	\$628	7.19	79

Fonte: Elaborado pelos autores (2021).

Pela Tabela 2 é possível observar a necessidade de propor uma nova segmentação, por evidenciarem que dois atributos importantes – Pedidos e Tempo de inatividade – são ignorados pela segmentação vigente. Além disso é possível observar que os valores realmente gastos no período de um ano estão distantes da proposta de segmentação.

### 4.3 Clusterização da base de dados

Após a construção da tabela com as informações selecionadas, os dados foram importados para o software RapidMiner para realizar a clusterização. O algoritmo foi aplicado, utilizando a Distância Euclidiana como medida de similaridade entre os dados. Foram testadas as segmentações com o parâmetro de entrada K para os valores 3, 4, 5 e 6 e, a partir da observação dos grupos formados, foi definido que o número de *clusters* que melhor se ajustava a segmentação utilizada atualmente foi para K = 4. Em todos os outros valores de K testados, o algoritmo gerava grupos com um número de dados muito pequenos que não atendiam o critério de substancialidade, definido por Madeira, Silveira e Toledo (2015), ou seja, os segmentos não eram grandes o suficiente para formar 3 segmentos e justificar sua ativação.

Como a segmentação deveria estar alinhada com a estratégia da companhia, foram formados 4 grupos (*cluster* 0, 1, 2 e 3). O *cluster* 0 no entanto, reuniu somente 569 compradores na amostra apresentada, o que não justifica sua ativação de forma separada. Dessa forma, esse *cluster* foi reagrupado com o *Cluster* 1, pois esses dois *clusters* são os que apresentaram a média de pedidos significativamente maior do que os demais *clusters*. Para a comparação com a segmentação atual, os *clusters* obtidos foram definidos com os mesmos nomes utilizados pela companhia.

Após o reagrupamento em três *clusters*, foi realizada a validação dos segmentos obtidos. Para que a divisão encontrada fosse considerada válida, ela deveria se diferenciar de alguma forma da divisão atual. Para entender se haviam diferenças, foi realizado o teste de Mann-Whitney para

duas amostras, com o intuito de comparar as medianas dos atributos em cada um dos grupos nas duas segmentações. Assim, para cada mediana a ser comparada foram definidas as seguintes hipóteses:

*Hipótese nula*  $H_0: \eta_1 - \eta_2 = 0$  e *Hipótese alternativa*  $H_1: \eta_1 - \eta_2 \neq 0$ ,

Sendo que:  $\eta_1$ : mediana da segmentação atual e  $\eta_2$ : mediana da segmentação proposta

A confiança atingida pelo teste foi de 95%. Os resultados obtidos para o *p-Valor* em cada uma das comparações podem ser visualizados na Tabela 3.

**Tabela 3: p-Valor obtido no Teste de Mann-Whitney para duas amostras em cada uma das médias comparadas**

Segmento	Demanda	Pedidos	Tempo de inatividade
Desenvolvimento em escala	0,000	0,848	0,000
Potencial de Crescimento	0,000	0,000	0,000
Alto Valor	0,000	0,000	0,000

Fonte: Elaborado pelos autores (2021).

A hipótese nula afirma que a diferença entre as medianas para as duas segmentações é 0. Quando o valor de *p* é menor do que o nível de significância de 0,05, a hipótese nula é rejeitada e conclui-se que a diferença entre as medianas da população é estatisticamente significativa. Ou seja, como pode ser observado na Tabela 3, com exceção da mediana de Pedidos para o segmento de “Desenvolvimento em Escala” que apresentou *p-Valor* de 0,848, todas as outras medianas da segmentação proposta apresentaram *p-Valor* de 0,000 e, portanto, se diferenciam das medianas da segmentação atual. Dessa forma, é possível afirmar que os grupos obtidos são válidos, pois são diferentes dos anteriores em pelo menos um dos atributos.

#### 4.4 Comparação dos resultados com a segmentação vigente

O primeiro ponto de comparação foi o tamanho dos grupos em cada um dos segmentos. A Tabela 4 mostra a composição dos grupos em cada uma das segmentações.

**Tabela 4: Composição dos grupos na Segmentação Atual x Segmentação Proposta**

Segmento	Segmentação Atual	Segmentação Proposta
Desenvolvimento em escala	75%	62%
Potencial de Crescimento	21%	34%
Alto Valor	4%	4%

Fonte: Elaborado pelos autores (2021).

Como supracitado, na divisão atual, por considerar somente valores em dólares para o agrupamento dos consumidores, há uma grande restrição dos grupos considerados mais fidelizados: “Alto Valor”, com somente 4% dos consumidores e “Potencial de Crescimento”, com 21%.

A nova proposta de segmentação manteve o tamanho do segmento de “Alto Valor”, com 4% dos consumidores. Mas uma diferença pode ser percebida nos outros dois segmentos. Há um crescimento de 13 pontos percentuais da base de consumidores classificados como “Potencial de Crescimento”, de 21% para 34%. O grupo “Desenvolvimento em Escala” cai na mesma proporção, diminuindo sua participação de 75% para 62%, ainda assim constituindo a maioria dos consumidores. Isso mostra que a nova segmentação tende a ser menos restritiva ao classificar os consumidores em grupos mais fidelizados. O segundo ponto de comparação foi a média dos atributos avaliados em cada um dos grupos nas duas propostas de segmentação. Os valores obtidos para as médias estão dispostos lado a lado na Tabela 5 para facilitar a comparação.

**Tabela 5: Média dos atributos para cada um dos segmentos**

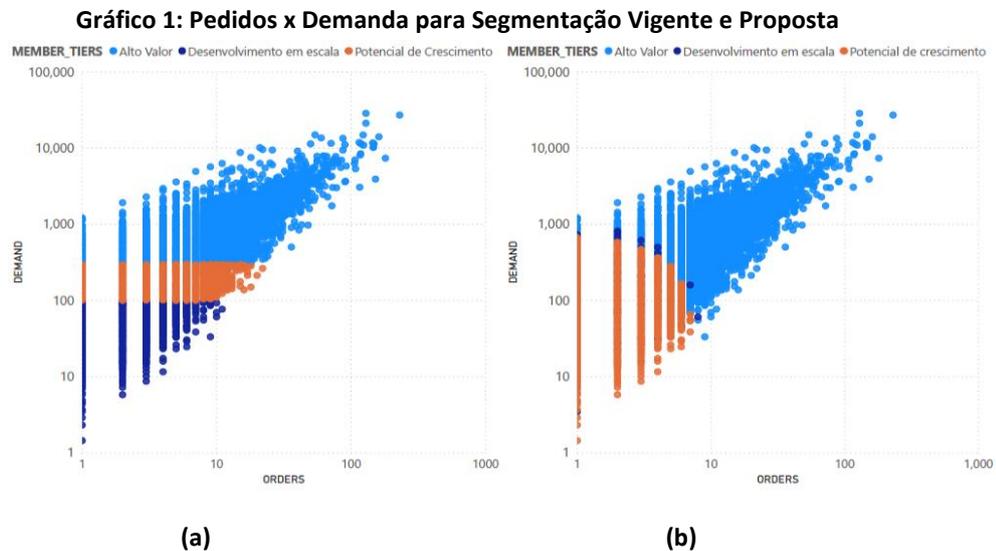
Atributo	Demanda		Pedidos		Tempo de Inatividade	
	Atual	Proposta	Atual	Proposta	Atual	Proposta
Desenvolvimento em escala	\$45	\$65	1.19	1.22	146	250
Potencial de Crescimento	\$159	\$77	2.33	1.52	123	83
Alto Valor	\$628	\$639	7.19	8.52	79	64

Fonte: Elaborado pelos autores (2021).

Para o segmento de “Alto Valor” há um crescimento de apenas 1,8% na média de demanda, mostrando que esse atributo não foi muito afetado pela nova segmentação nesse grupo. Já em “Potencial de Crescimento” há uma redução de 51,6% na média de demanda, enquanto para os consumidores classificados como “Desenvolvimento em Escala” há um crescimento de 44% na média da demanda. Tais dados evidenciam que houve uma redistribuição entre os consumidores que estavam nesses dois grupos anteriormente, redimensionando-os melhor em relação ao valor médio gasto no período, em face dos demais atributos considerados na análise. O mesmo aconteceu para o atributo de Pedidos. Enquanto há crescimento na média de pedidos para “Desenvolvimento em Escala” de 2,5%, há uma queda de 34,7% nesse valor para “Potencial de Crescimento”. No caso dos pedidos, também há o crescimento da média desse atributo em 18,5% para o segmento de “Alto Valor”.

Essa redução nas médias de Demanda e Pedidos para “Potencial de Crescimento” e aumento nas médias para “Desenvolvimento em Escala” podem ser explicados pelo último atributo considerado para análise: o Tempo de Inatividade. Esse atributo foi o que apresentou as maiores variações na média, para “Desenvolvimento em Escala” em que se observa um aumento de 71,2%. Para “Potencial de Crescimento”, há uma redução de 32,5% e, para “Alto Valor”, há uma redução de 19%. Tais resultados evidenciam que a classificação “Potencial de Crescimento” passou a englobar provavelmente consumidores cadastrados recentemente, o que explica a redução na média de Demanda e Pedidos para esse grupo, visto que, como são novos, podem não ter feito muitos pedidos ainda. Da mesma forma, os consumidores em “Desenvolvimento em Escala”, apesar do aumento na média de Demanda e Pedidos, tiveram um aumento principalmente no tempo em que não interagem com a companhia, mostrando que esse atributo foi essencial para classificar alguns consumidores nesse grupo.

Por fim, também foi comparada a distribuição dos consumidores nos segmentos. Para isso foram elaborados gráficos de dispersão combinando aos atributos analisados. O Gráfico 1 mostra a distribuição de Pedidos x Demanda para os três segmentos antes e depois



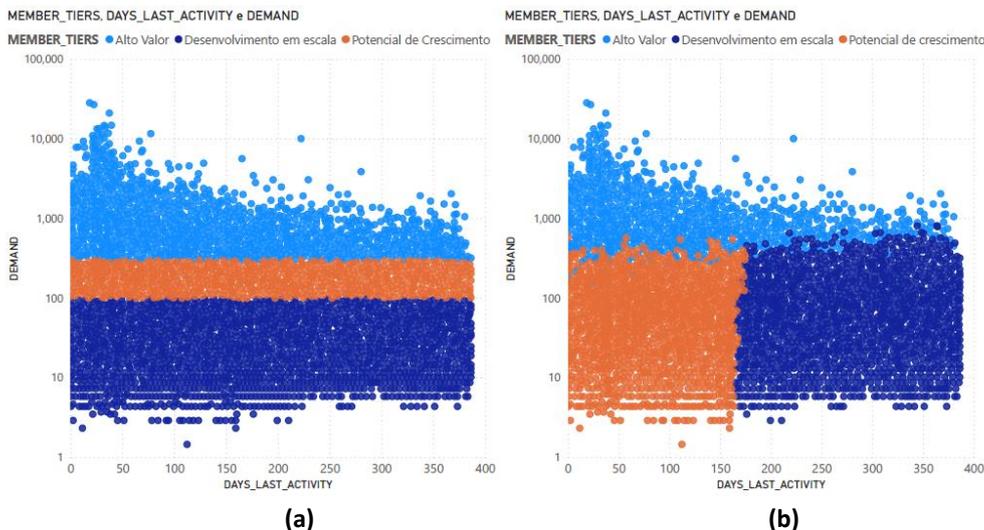
Fonte: Elaborado pelos autores (2020).

O Gráfico 1 (a) representa a segmentação atual e, como visto anteriormente, o único atributo considerado para a segmentação é a demanda, evidenciado pelas divisões claramente definidas no eixo da Demanda. Como consequência disso é possível observar consumidores classificados como “Alto Valor” que realizaram somente um pedido no período considerado, assim como consumidores classificados como “Potencial de Crescimento” e “Desenvolvimento em Escala” com mais de dez pedidos.

Já o Gráfico 1 (b) representa a segmentação proposta pela técnica de *machine learning*. Nesse caso é possível observar que nenhum dos dois atributos plotados definem de forma predominante os segmentos. Também é possível observar uma correção parcial dos problemas apontados na segmentação anterior. Há uma redução considerável no grupo de “Alto Valor” com somente um pedido. Também não é possível visualizar mais consumidores classificados como “Potencial de Crescimento” e “Desenvolvimento em Escala” com mais de dez pedidos realizados. Outra característica bastante evidente na nova segmentação é uma grande sobreposição desses dois últimos segmentos no que se refere a pedidos e demanda.

A segunda visualização plotada traz uma visão de Tempo de Inatividade x Demanda e pode ser encontrada no Gráfico 2.

**Gráfico 2: Tempo de inatividade x Demanda para a Segmentação Vigente e Proposta**



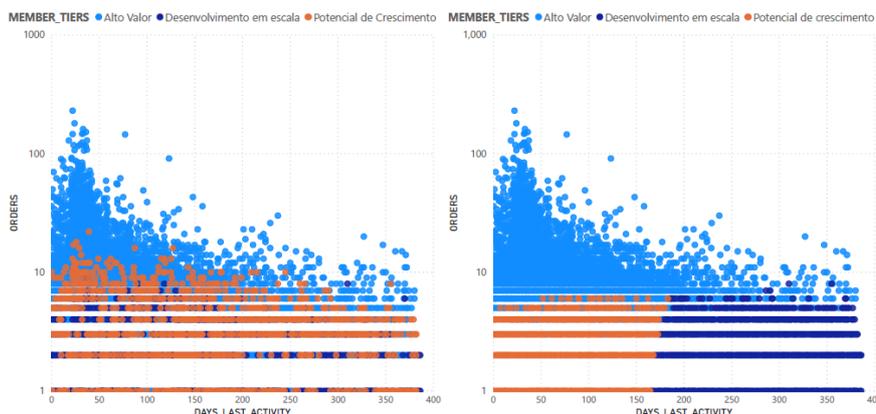
Fonte: Elaborado pelos autores (2021).

O Gráfico 2 (a) traz a visão da segmentação atual e evidencia que apenas o atributo da demanda é considerado pela clara divisão no eixo Y, sem considerar se os consumidores estão inativos há muito tempo ou não. Por outro lado, o Gráfico 2 (a) deixa claro que esse atributo passou a ser considerado, principalmente nos segmentos de “Potencial de Crescimento” e “Desenvolvimento em Escala”. Como visto no Gráfico 2 (b), há uma sobreposição de pedidos e demanda nesses dois segmentos e aqui fica claro que a principal diferenciação entre esses grupos passou a ser a recência da atividade.

No Gráfico 2 (b) também pode ser observado que para o segmento de “Alto Valor”, alguns consumidores ainda foram incluídos nesse grupo, mesmo com muito tempo de inatividade. Nesse caso, o atributo de Demanda teve mais relevância na classificação, visto que foram consumidores que gastaram um valor próximo a \$1000 dólares.

Por fim, a última visualização plotada para as duas segmentações é apresentada no Gráfico 3 e traz a visão de Tempo de Inatividade x Pedidos.

**Gráfico 3: Tempo de inatividade x Pedidos para a Segmentação Vigente e Proposta**



(a)

(b)

Fonte: Elaborado pelos autores (2021).

Nessa visualização não é possível identificar no Gráfico 3 (a) uma divisão clara para segmentação atual (Gráfico (b)), pois ele não tem a representação do atributo demanda. Pelo contrário, esse gráfico mostra uma sobreposição dos segmentos, principalmente “Potencial de Crescimento” e “Desenvolvimento em Escala”. Também ficam evidenciadas características da segmentação atual já observadas anteriormente, como a irrelevância do tempo de inatividade para o agrupamento e consumidores com um alto número de pedidos que não são classificados como “Alto Valor”.

Já no Gráfico 3 (b) podem ser observadas novamente as características já descritas para a nova segmentação proposta. Para os segmentos de “Desenvolvimento em Escala” e “Potencial de Crescimento” há uma divisão mais evidente no tempo de inatividade que diferencia esses dois grupos. Na nova proposição todos os consumidores com mais de dez pedidos também passam a ser classificados como de “Alto Valor”.

## 5 CONCLUSÃO

O intuito desse estudo foi responder como pode ser realizada a segmentação da base de consumidores do *e-commerce* de uma empresa do varejo esportivo, a partir da clusterização de dados em *machine learning*. Para isso, foi desenvolvida uma pesquisa a partir do *e-commerce* de uma multinacional do varejo esportivo que atua no Brasil. O estudo consistiu na avaliação da segmentação utilizada atualmente e na proposição de uma nova segmentação utilizando um algoritmo para formação de *clusters*.

A partir da análise, foi possível compreender que a segmentação utilizada atualmente pela companhia, que se baseia somente no valor gasto pelos consumidores nos doze meses anteriores, pode não ser a melhor opção, por considerar apenas um atributo no agrupamento dos clientes. A análise revelou consumidores que estavam há muito tempo inativos em todos os grupos, e consumidores com um alto número de pedidos, que estavam designados em grupos considerados menos fidelizados. Desse modo, a nova proposta de segmentação englobou outros atributos: o número de pedidos nos doze meses anteriores e o tempo de inatividade dos consumidores, que, como definido por Hooley, Piercy e Nicoulaud (2011) são variáveis relacionadas ao comportamento do cliente.

Para a nova segmentação foi utilizado o algoritmo de clusterização K-means, uma técnica de particionamento que permite agrupar elementos dentro de um conjunto de dados que possuem maior similaridade entre si. O uso desse algoritmo de *Machine Learning* se mostrou satisfatório, visto que foi possível obter três segmentos válidos que se diferenciavam da segmentação adotada atualmente. Também foi possível observar uma melhora nos atributos avaliados, visto que a nova segmentação proposta passou a considerar também, o número de pedidos e o tempo de inatividade para o agrupamento dos consumidores, além da demanda que já era utilizada anteriormente.



É importante ressaltar que a nova segmentação se mostrou válida e vantajosa em comparação à segmentação atual para os dados históricos utilizados no estudo. Entretanto, essa segmentação não foi testada de forma gerencial para confirmar se o novo agrupamento melhoraria a ativação dos consumidores como mecanismo de reposta às ações comerciais da companhia. Dessa forma, como complemento desse estudo, sugere-se a realização da testagem dos segmentos obtidos, com estratégias definidas de comunicação para cada grupo de acordo com suas características. Além disso, sugere-se a realização de análises descritivas mais detalhadas de cada grupo. Apesar da segmentação basear-se no comportamento de compra, os *clusters* obtidos ainda podem ser descritos, por exemplo, em função de suas características demográficas e preferências por determinada categoria de produto. A caracterização desses grupos pode tornar os segmentos mais acessíveis para a construção de estratégias de marketing.

A presente investigação avança na discussão sobre o uso de técnicas de *Machine Learning* na área de negócios (em especial na área de marketing), como alternativa para o aperfeiçoamento na segmentação de clientes. O estudo contribui ainda para a literatura sobre *e-commerce*, marketing de varejo e clusterização de dados, por trazer à luz, a partir da proposta de segmentação apresentada, possibilidades de questionamentos e reflexões a respeito do uso de grandes volumes de dados de clientes pelas organizações. Além disso o estudo desperta atenção para as possibilidades de compreensão de grupos sociais específicos e seus respectivos comportamentos de compra, a partir da construção de robustos bancos de dados que reúnam tais atributos pelas organizações do varejo.

## 6 REFERÊNCIAS

- Ahani, A., Nilashi, M., Ibrahim, O., Sanzogni, L., & Weaven, S. (2019). Market segmentation and travel choice prediction in Spa hotels through TripAdvisor's online reviews. *International Journal of Hospitality Management*, 80, 52–77. <https://doi.org/10.1016/j.ijhm.2019.01.003>.
- Ali, M. (2020). *How to implement Clustering in Power BI using PyCaret*. Acesso em 17 de outubro de 2020, disponível em: <https://towardsdatascience.com/how-to-implement-clustering-in-power-bi-using-pycaret-4b5e34b1405b>.
- Anitha, P., & Patil, M. M. (2019). RFM model for customer purchase behavior using K-Means algorithm. *Journal of King Saud University - Computer and Information Sciences*. <https://doi.org/10.1016/j.jksuci.2019.12.011>.
- Antonopoulos, I., Robu, V., Couraud, B., Kirli, D., Norbu, S., Kiprakis, A., Flynn, D., Elizondo-Gonzalez, S., & Wattam, S. (2020). Artificial intelligence and machine learning approaches to energy demand-side response: a systematic review. *Renewable and Sustainable Energy Reviews*, 130, 109899. <https://doi.org/10.1016/j.rser.2020.109899>.
- Chagas, B. N. R., Viana, J., Reinhold, O., Lobato, F. M. F., Jacob, A. F. L., & Alt, R. (2020). A literature review of the current applications of machine learning and their practical implications. *Web Intelligence*, 18(1), 69–83. <https://doi.org/10.3233/WEB-200429>.



- Dogan, O., Ayçin, E., & Bulut, Z. A. (2018). Customer Segmentation by Using RFM Model and Clustering Methods: A Case Study in Retail Industry. *International Journal of Contemporary Economics and Administrative Sciences*, 8(1), 1-19. <http://www.ijceas.com/index.php/ijceas/article/view/174>.
- Farrokhi, A., Farahbakhsh, R., Rezazadeh, J., & Minerva, R. (2021). Application of Internet of Things and artificial intelligence for smart fitness: A survey. *Computer Networks*, 189, 107859. <https://doi.org/10.1016/j.comnet.2021.107859>.
- Fildes, R., Ma, S., & Kolassa, S. (2019). Retail forecasting: Research and practice. *International Journal of Forecasting*. <https://doi.org/10.1016/j.ijforecast.2019.06.004>.
- Fuentes, S., Tongson, E., & Gonzalez Viejo, C. (2021). Novel digital technologies implemented in sensory science and consumer perception. *Current Opinion in Food Science*, 41, 99–106. <https://doi.org/10.1016/j.cofs.2021.03.014>.
- Kachamas, P., Akkaradamrongrat, S., Sinthupinyo, S., & Chandrachai, A. (2019). Application of Artificial Intelligent in the Prediction of Consumer Behavior from Facebook Posts Analysis. *International Journal of Machine Learning and Computing*, 9(1), 91–97. <https://doi.org/10.18178/ijmlc.2019.9.1.770>.
- Hooley, G., Piercy, N. F., & Nicoulaud, B. (2011). *Estratégia de marketing e posicionamento competitivo* (4ª ed.). (L. Pauleti, & S. Midori, Trans.) São Paulo: Pearson Prentice Hall.
- Jain, A. K., Murty, M. N., & Flynn, P. J. (1999). Data clustering. *ACM Computing Surveys*, 31(3), 264–323. <https://doi.org/10.1145/331499.331504>.
- Kaur, H., & Anand, S. (2018). Segmenting Generation Y using the Big Five personality traits: understanding differences in fashion consciousness, status consumption and materialism. *Young Consumers*, 19(4), 382–401. <https://doi.org/10.1108/YC-03-2018-00788>.
- Lima, B. V. A. de, Machado, V. P., & Lopes, L. A. (2015). Aprendizado de máquina para rotulação automática de usuários de uma rede social acadêmica. *Revista Eletrônica de Sistemas de Informação*, 14(1), 4. <https://doi.org/10.21529/RESI.2015.1401004>.
- Luz, V. V. (2020). *Comportamento do consumidor na era digital* (1ª ed.). Curitiba: Contentus.
- Madeira, A. B., Silveira, J. G., & Toledo, L. A. (2015). Marketing Segmentation: Your Role For Diversity in Dynamical Systems. *Revista Eletrônica de Gestão Organizacional*, 13(1), 71-78. <https://periodicos.ufpe.br/revistas/gestaoorg/article/download/22025/18453>.
- Malhotra, N. K. (2012). *Pesquisa de marketing: uma orientação aplicada* (6 ed.). Porto Alegre: Bookman.
- Monard, M. C., & Baranauskas, J. A. (2005). Conceitos sobre aprendizado de máquina. In: Rezende, S. O. (Ed.). *Sistemas inteligentes: fundamentos e aplicações* (pp. 39-56). Barueri: Manole.



- Oliveira, T. S. (2008). *Clusterização de dados utilizando técnicas de redes complexas e computação bioinspirada*. Dissertação (Mestrado em Ciências de Computação e Matemática Computacional) - Universidade de São Paulo, São Carlos.
- Qannari, E. M. (2017). Sensometrics approaches in sensory and consumer research. *Current Opinion in Food Science*, 15, 8–13. <https://doi.org/10.1016/j.cofs.2017.04.001>.
- Samuel, J. (2017). Information Token Driven Machine Learning For Electronic Markets: Performance Effects In Behavioral Financial Big Data Analytics. *Journal of Information Systems and Technology Management*, 14(3), 371–384. <https://doi.org/10.4301/S1807-17752017000300005>.
- Saura, J. R. (2021). Using Data Sciences in Digital Marketing: Framework, methods, and performance metrics. *Journal of Innovation & Knowledge*, 6(2), 92–102. <https://doi.org/10.1016/j.jik.2020.08.001>.
- Souza, L. L. F., & Freitas, A. A. F. (2016). Revisão da produção científica brasileira em segmentação de mercado. *Revista de Ciências Da Administração*, 96–108. <https://doi.org/10.5007/2175-8077.2016v18n45p96>.
- Sports Value. (2018). *A competição global das marcas de material esportivo*. Sports Value. Acesso em 20 de novembro de 2020, disponível em: <https://www.sportsvalue.com.br/wp-content/uploads/2018/08/SportsValue-Empresas-de-material-esportivo-2018-1.pdf>.
- Turban, E., & King, D. (2004). *Comércio Eletrônico: Estratégia e Gestão* (1ª ed.). São Paulo: Prentice Hall.
- Yin, J., & Fernandez, V. (2020). A systematic review on business analytics. *Journal of Industrial Engineering and Management*, 13(2), 283. <https://doi.org/10.3926/jiem.3030>.

#### COMO CITAR ESTE ARTIGO:

Falqueto, A. A., & Cezar, L. C. (2022). SEGMENTAÇÃO VIA MACHINE LEARNING: PROPOSTA DE CLUSTERIZAÇÃO DE CONSUMIDORES DO E-COMMERCE DE UMA EMPRESA MULTINACIONAL DO VAREJO ESPORTIVO. *HOLOS*, 4. Recuperado de <https://www2.ifrn.edu.br/ojs/index.php/HOLOS/article/view/12032>

#### SOBRE OS AUTORES

##### A. A. FALQUETO

Graduanda em Engenharia de Produção. Universidade Federal de Viçosa. E-mail: [alice.falqueto@ufv.br](mailto:alice.falqueto@ufv.br)  
ORCID ID: <https://orcid.org/0000-0002-5114-0332>

##### L. C. CEZAR

Doutor em Administração pela Universidade Federal do Espírito Santo. Professor do Departamento de



Administração e Contabilidade (DAD). Universidade Federal de Viçosa (UFV). E-mail: layon.cezar@ufv.br  
ORCID ID: <https://orcid.org/0000-0003-2062-4593>

**Editor(a) Responsável:** Prof. Dr. Miler Franco D'anjour

**Pareceristas Ad Hoc:** Richard Medeiros Araújo e Roberto Rodney Ferreira Junior



**Recebido:** 19 de fevereiro de 2021

**Aceito:** 29 de novembro de 2021

**Publicado:** 28 de dezembro de 2022

