

Uso de *Machine Learning* para entendimento das relações de consumo no modelo de negócio (*SaaS enabled Marketplace*) do Olist

RESUMO

Esta pesquisa visou identificar os fatores que impactam o review score dado por consumidores após compras via Olist Store em marketplaces, utilizando técnicas de Machine Learning, especificamente Random Forest e Feature Importance. A análise, baseada em um data frame com 10 categorias e 4308 amostras de treino (73,9% de score alto) e 1078 de teste (mesma proporção), obteve uma acurácia média de 79,1%. A variável `expected_diff`, que mede a diferença entre a data prevista e a real de entrega, foi a mais relevante para a predição do modelo, destacando a importância da entrega nas compras online, especialmente antes da COVID-19. As categorias de beleza/saúde (28,73%) e brinquedos (27,60%) mostraram maior necessidade de entregas rápidas e alta acurácia, enquanto informática/acessórios (22,47%) e móveis/decoração (23,21%) apresentaram melhor distribuição entre variáveis, mas com acurácias menores. A pesquisa conclui que um investimento robusto em logística é crucial para empresas que desejam obter vantagem competitiva no comércio eletrônico.

PALAVRAS-CHAVE: Comércio eletrônico. Técnicas de *Machine Learning*. Categorias de produtos. Previsão do *review score*.

Gabriela Amaral de Alencar Leite
Universidade Federal de São Paulo
(UNIFESP) São Paulo, São Paulo
gabriela.amaral.alencar@gmail.com

Luis Hernan Contreras Pinochet
Faculdade de Economia,
Administração, Contabilidade e
Atuária da Universidade de São
Paulo (FEAUSP), São Paulo, São
Paulo
luis.hernan@usp.br

Vanessa Itacaramby Pardim
Faculdade de Economia,
Administração, Contabilidade e
Atuária da Universidade de São
Paulo (FEAUSP), São Paulo, São
Paulo
vanessa.itacaramby@usp.br

Luciana Massaro Onusic
Universidade Federal de São Paulo
(UNIFESP) São Paulo, São Paulo
luciana.onusic@unifesp.br

INTRODUÇÃO

O empreendedorismo tem crescido em termos de importância em nações desenvolvidas e em desenvolvimento, tendo em vista seu impacto no desenvolvimento nacional e econômico (KUMAR et al., 2021; LIMA et al., 2021). Nesse contexto, de acordo com Pobee (2021), o comércio eletrônico ganha destaque, pois se constitui em uma ferramenta importante para os empreendedores por permitir, entre outros, 24 horas de atendimento, 7 dias por semana, promovendo a troca de informações em tempo real. Shehata e Montash (2020) destacam que pequenas empresas, incluindo empreendedores individuais, podem obter vantagem competitiva e estratégica ao adotar o comércio eletrônico. Isso ficou ainda mais evidente quando da eclosão da pandemia da COVID-19 (SCUTARIU et al., 2022).

Nesse contexto, a computação em nuvem se apresenta com grande potencial para transformar a maneira como as empresas de comércio eletrônico fazem negócios. De acordo com o *National Institute of Standards and Technology (NIST)*, os modelos atuais de serviços em nuvem são *Infrastructure-as-a-Service (IaaS)*, *Platform-as-a-Service (PaaS)* e *Software-as-a-Service (SaaS)*. Cada modelo representa um tipo diferente de serviço de nuvem para modelos de negócios que buscam a inovação (PINOCHET et al., 2021).

Inserindo a lógica destes modelos a uma empresa de comércio eletrônico é possível observar que a *IaaS* é usado para comprar infraestrutura de computador como um serviço sob demanda, como servidores, armazenamento, redes e sistemas operacionais. O *PaaS* fornece aos clientes uma plataforma de aplicativos pré-criada que pode ser usada, conforme necessário, para desenvolvimento em vez de investir em sua própria infraestrutura subjacente. Por fim, o *SaaS* pode ser usado para iniciar rapidamente um site de comércio eletrônico sem se preocupar com configurações de servidor, investimentos em equipamentos ou atualizações de *software* (WANG et al., 2024).

A adoção do comércio eletrônico e dos serviços de nuvem podem potencializar as oportunidades e, conseqüentemente, aumentar o desempenho das empresas no contexto da economia digital (SCUTARIU et al., 2022).

Para acompanhar o nível de satisfação, existem diferentes caminhos, mas o método *CSAT (Customer Satisfaction Score)*, traduzido como a nota atribuída a determinado serviço/produto indicando a satisfação do consumidor com determinada experiência (BIRKETT, 2021), se faz necessário para o sucesso de qualquer estratégia. Neste trabalho, o método *CSAT* é representado pela nota da experiência no *Olist Store*, maior loja de departamentos dentro de grandes *marketplaces* brasileiros, intermediando o consumo no meio digital.

Assim, o objetivo geral desta pesquisa é identificar os fatores que impactam no *review score* dado pelo consumidor após a realização da compra via *Olist Store*, em algum *marketplace*, por meio da análise de dados e aplicação do modelo de *Machine Learning* para prever a nota baseada em *CSAT*.

Com este estudo espera-se promover avanços nas técnicas analíticas de *Machine Learning* de *Random Forest* e *Feature Importance*, como forma de ranquear as principais variáveis que impactam na satisfação do cliente via *CSAT* (nota dada pelo usuário), e apresentar uma estratégia de comparação do

desempenho de diferentes categorias de produtos da venda no Olist Store na previsão do review dado pelo usuário por meio de técnicas de *Machine Learning*.

A Olist ocupa um importante papel de impulsionar o empreendedor de menor porte, com destaque para o contexto brasileiro, facilitando o acesso às primeiras páginas dos *marketplaces* (NIWAKTE, 2021).

A contribuição do estudo decorre da escolha pela aplicação de técnicas de *Machine Learning* em vista que os algoritmos de aprendizado de máquina têm obtido sucesso ao serem capazes de igualar e superar o desempenho humano em diversos campos, como tradução de idiomas, jogos de tabuleiro e carros inteligentes (MCCOY; AURET, 2019; ZAGHLOUL et al., 2024).

Essa pesquisa destaca a aplicação de técnicas de ciência de dados, especificamente o *Machine Learning*, para compreender os fatores de influência na pontuação de avaliação dada pelos consumidores após suas compras em um *marketplace* (ZAGHLOUL et al., 2024). As análises feitas revelam a importância da logística na satisfação do cliente e, conseqüentemente, na reputação do negócio. Essa abordagem demonstra a intersecção existente entre ciência, tecnologia e sociedade, em que avanços tecnológicos como o deste estudo são aplicados para resolver questões socioeconômicas, como a melhoria da experiência do consumidor e a competitividade das empresas no comércio eletrônico. Reconhecendo a relevância da logística na era digital, esse estudo ressalta a necessidade de as empresas entenderem e investirem em processos logísticos eficientes para se destacarem no mercado (LIMA et., 2021).

Em suma, a aplicação dessas técnicas pode trazer inúmeros benefícios às organizações, como aumentar a eficiência e a eficácia dos processos de auditoria (CHAN; VASARHELYI, 2011), ou prever, objetivo deste estudo, o comportamento de consumidores (PAÇO et al., 2018).

REFERENCIAL TEÓRICO

A influência do intermediário no comércio eletrônico

Um dos maiores avanços proporcionados pela internet é o comércio eletrônico, e as transações comerciais na rede, que podem ser definidas como qualquer transação comercial efetuada em algum ambiente online, entre diferentes partes, sendo as mais famosas: entre empresas, popularmente conhecido como B2B (*business-to-business*) e entre empresas e consumidores finais B2C (*business-to-consumer*) (Rodrigues, 2020).

Para explorar esse potencial de crescimento do comércio eletrônico, as plataformas online mudaram para o modelo de *marketplaces*. Isto se deve ao fato de que mesmo os varejistas online que tradicionalmente operavam como revendedores estão se movendo, cada vez mais, para um “*marketplace mode*” permitindo terceirizar transações diretamente com clientes (RATH et al., 2021).

Nesse contexto, tendo em vista a complexidade crescente do comércio eletrônico, surgem várias empresas, com destaque para a startups, que, segundo a Associação Brasileira de Startups (ABSTARTUPS, 2019), estão alavancando seus negócios com base em 3 modelos: *SaaS*, *Marketplace* e *SaaS enabled Marketplace*.

O primeiro, *SaaS (Software as a Service)*, automação usando *softwares*, baseado em algoritmos, para determinado processo que era feito manualmente.

Na sequência tem-se o *Marketplace*, cujo modelo consiste em conectar oferta e demanda, cobrando uma parcela por tal intermediação. Por fim, tem-se o *SaaS enabled Marketplace* que consiste na automatização de algum processo do usuário, além de conectar de alguma forma, a oferta e demanda, monetizando essa intermediação.

A Olist é considerada a maior loja de departamentos dentro dos maiores *marketplaces* brasileiros e ganhou mercado ajudando lojas físicas a venderem nesses *marketplaces* como Mercado Livre, Amazon e Submarino (Olist [s/d]). O modelo de negócio da Olist consiste em automatizar um processo que seria feito manualmente por meio de *softwares*, visto que o Olist centraliza todos os cadastros e unifica a gestão na plataforma única da empresa, com acesso exclusivo ao lojista (NIWATE, 2021).

Além disso, conecta a oferta com a demanda, pois se posiciona como uma loja única perante os consumidores no menu de busca dos *marketplaces*, mas é composta por todos os lojistas credenciados. A Figura 1 ilustra o ecossistema ao qual a Olist Store faz parte.

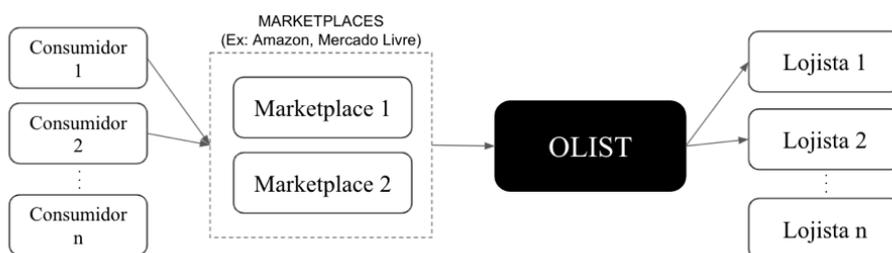


Figura 1: Fluxo de compras via Olist Store

Fonte: elaboração própria

É possível diferenciar a Olist dos demais *marketplaces*, tendo em vista que pode ser definida como um shopping virtual, ou seja, um site de *e-commerce* que reúne ofertas de produtos e serviços de diferentes vendedores (EUROMONITOR, 2018). Portanto, conforme observado na Figura 1, nas transações via Olist Store, são dois intermediários entre o consumidor e o lojista.

O modelo de negócios da empresa, auxilia, principalmente, pequenos e médios lojistas que enfrentam muitas vezes problemas por incompatibilidade tecnológica ao ingressarem nos aplicativos de comércio, por meio de automatização e monitoramento integrado.

Machine Learning e a compreensão de padrões

Nas últimas décadas, a tecnologia vem ganhando cada vez mais espaço na vida das pessoas, consequentemente vem sendo cada vez mais discutida, com *Machine Learning* não é diferente. O *Machine Learning* é um ramo da Inteligência

Artificial que tem como intuito explorar algoritmos que podem aprender e fazer predições com base nos dados de entrada e no passado (HUTTER et al., 2019).

Com o passar do tempo, se tornou mais necessário que as tecnologias fossem capazes de criar hipóteses baseadas nas experiências passadas fornecidas para resolver o problema (FACELI et al., 2011). Os algoritmos de *Machine Learning* podem ser categorizados, conforme o tipo de aprendizado, em 3 tipos, sendo eles: aprendizado supervisionado, aprendizado não-supervisionado e aprendizado por reforço (CASTLE, 2018; COSTI, 2020).

Para Silva (2021), o aprendizado supervisionado se refere aos algoritmos de *Machine Learning* que aprendem a “resposta correta” mapeando os *inputs* e *outputs* por meio de amostras inseridas no sistema, ou seja, o algoritmo entra em contato com a base de treino que possui a variável resposta, desta forma pode aprender, encontrar os padrões, criar uma regra interna e prever o resultado de uma nova amostra.

Assim, normalmente utilizado para a predição de eventos, o objetivo do aprendizado supervisionado é aprender uma função que, dada uma amostra de dados e resultados desejados, se aproxima melhor da relação entre entrada e saída observável nos dados (TSUNODA et al., 2020).

Já os não supervisionados, são treinados em dados não rotulados, ou seja, chega em respostas não presentes no banco de dados original, e devem determinar a importância do recurso independentemente, com base nos padrões inerentes à amostra (COSTI, 2020). Seu objetivo é procurar alguma estrutura em dados de amostra, agrupá-los em grupos de regras semelhantes e descobrir padrões ocultos. Métodos de agrupamento são usados para encontrar uma partição dos dados e classificar novos dados de entrada com uma regra de previsão (Jordan & Mitchell, 2015).

Por fim, a terceira categoria diz respeito ao aprendizado por reforço, cujo dados de treinamento indicam uma recompensa ou uma punição ao algoritmo com base nas metas estabelecidas. Na aprendizagem por tentativa e erro, usando recompensas e punições como *feedback*, os algoritmos encontram uma solução adequada maximizando as recompensas totais (JORDAN; MITCHELL, 2015; SILVA, 2021).

O presente estudo está alinhado na categoria do aprendizado de máquina supervisionado, visto que ele visa prever o campo de pontuação dada pelo cliente após a compra.

PROCEDIMENTOS METODOLÓGICOS

A pesquisa que se apresenta tem natureza preditiva, uma vez que consiste na aplicação de algoritmos para compreender a estrutura dos dados existentes e gerar regras de predição (Santos et al., 2019), e utilizou algoritmos de *Machine Learning* para analisar dados secundários obtidos de fontes abertas e de acesso público. Técnicas de *Machine Learning* estão cada vez sendo mais exploradas porque, quando os modelos são expostos a novos dados, eles são capazes de se adaptar independentemente, significando que é possível produzir, de forma rápida, modelos capazes de analisar dados maiores e mais complexos, e entregar resultados mais rápidos e precisos (MOHRI et al., 2018).

A escolha da base de dados

A base de dados que utilizada é composta por dados públicos da Olist, disponível no Kaggle, plataforma *online* gratuita fundada em 2010 e considerada uma grande comunidade de cientistas de dados com fóruns de discussão, desafios, competições e *datasets* públicos com grande quantidade de informações para exploração (BANACHEWICZ; MASSARON, 2022).

A base de dados do estudo é composta por dados temporais de pedidos da Olist Store, de 2016 a 2018 (dados mais atuais não foram disponibilizados pelo Kaggle ou empresa por serem considerados dados estratégicos e sensíveis no mercado), e informações de aproximadamente 100 mil transações. Composto por 9 bancos de dados, Figura 2, possui uma grande diversidade de tipos de informação com visibilidade de informações de itens do pedido, preço, frete, pagamento, localização do cliente, atributos do produto de diferentes setores, suas características do anúncio digital e algumas mais.

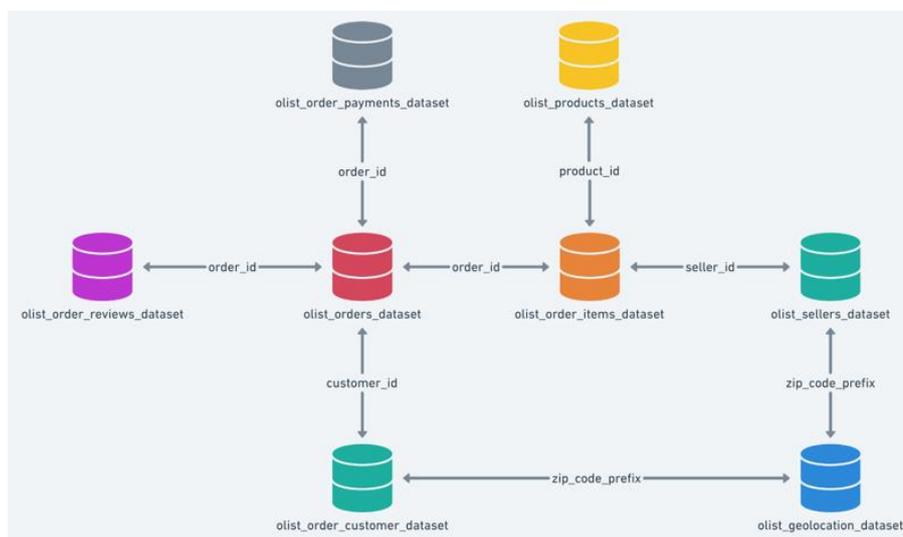


Figura 2: Esquema dos conjuntos de dados do estudo

Fonte: Kaggle

É importante ressaltar que os bancos de dados são compostos por algumas chaves primárias e únicas entre as bases, o que permite a união de todas, mas sem permitir qualquer tipo de identificação dos envolvidos, para respeitar a Lei Geral de Proteção de Dados Pessoais (LGPD).

O primeiro banco de dados, *olist_order_reviews_dataset*, contém informações sobre a nota de avaliação dada pelo comprador após a compra, em uma escala de 1 a 5, sendo auxiliar para a variável foco de estudo, acompanhada da data do evento e outras informações.

O segundo banco por sua vez, *olist_orders_dataset*, contém as informações da compra realizada por algum cliente. Por exemplo, data prevista divulgada ao cliente ao realizar a compra e a real data de chegada do pedido. Ao tirar a diferença entre essas duas datas, tem-se uma nova variável de interesse para o estudo.

O banco *olist_order_payments_dataset*, contém as informações sobre as formas de pagamento e valores de cada pedido realizado no período de 2016 a 2018, anos que compõe a base. Já o banco *olist_order_customer_dataset*, contém informações de identificação dos clientes que permitirão mapear todo o processo de compra, unificando aos demais conjuntos de dados como forma de pagamento, a data prevista de entrega e a data em que ele realmente chegou, chegando à avaliação e comentário dados ao final de todo o processo.

O banco *olist_geolocation_dataset* por sua vez, contém informações geográficas dos clientes como estados, CEPs das cidades brasileiras, entre outras informações. Já o banco *olist_sellers_dataset* fornece as informações correspondente aos vendedores que tiveram pedidos no período analisado, permitindo a identificação dos estados brasileiros e podendo ser associados com a demora ou rapidez na entrega.

Na sequência tem-se o banco *olist_order_items_dataset* que fornece as informações de cada item solicitado em cada pedido do período analisado, assim como valores do produto e do frete. Quando há mais de um item para o mesmo frete, o preço será dividido entre o número de itens que se adequem.

O último banco é o *olist_products_dataset* que contém informações sobre os produtos que foram comercializados no período analisado, permitindo a identificação das macrocategorias que a própria empresa utiliza, para buscar padrões de consumo por meio das aplicações da tecnologia à pesquisa.

Técnicas de pesquisa e Análise de Dados

O presente estudo utilizou o Google Colaboratory, também conhecido como Google Colab, como principal plataforma na nuvem para o desenvolvimento de toda a análise, desde o momento exploratório dos dados até o algoritmo de *Machine Learning* (GUNAWAN et al., 2020).

Por ser um ambiente totalmente *online* e diretamente no seu navegador, a plataforma nos proporciona conexão direta com a plataforma Kaggle, local onde o banco de dados foi disponibilizado, por meio de uma chave API única por usuário. As principais bibliotecas utilizadas no estudo são: NumPy, Pandas, Math, Matplotlib, Seaborn, Datetime e Sickit-Learn.

As bibliotecas de código surgiram no intuito de ajudar a economizar tempo e deixar a programação prática e menos extensa. Uma biblioteca é uma coleção de códigos pré-combinados que podem ser utilizados juntos de forma funcional para eliminar a necessidade de você digitar o código do zero, otimizando partes que são quase tidas como padrão para aquela aplicação em específico (TEAM, 2022). As principais bibliotecas utilizadas no estudo são:

1. Numpy: É uma biblioteca poderosa que auxilia a realização de cálculos em *arrays* multidimensionais, ou seja, dados multidimensionais. Possui módulos que auxiliam no armazenamento dos dados de treinamento e parâmetros do aprendizado de máquina (SANTIAGO JR, 2019).
2. Pandas: Essa é considerada uma das mais completas para análise e estruturação de dados (FIGUEIREDO, 2018). É possível identificar durante o código as ferramentas que foram utilizadas do Pandas pela repetição do código "pd." (p-d-) seguido da função que será utilizada.

3. Math: A biblioteca *math* pode auxiliar desde a função matemática mais simples até a mais complexa, assim como utilizar uma série de constantes matemáticas (SAMUEL, 2019).
4. Matplotlib: Uma das mais famosas bibliotecas no quesito de visualização de dados, é idealmente feito para conseguir fazer visualizações de forma muito simples (BARRETT et al., 2005), como um histograma em poucos comandos.
5. Seaborn: É outra ferramenta para auxiliar na visualização de dados e gráficos estatísticos. Oferece uma interface de alto nível com o Matplotlib e integra-se com as estruturas de dados do Pandas (WASKOM, 2021).
6. Datetime: biblioteca focada na manipulação de datas e horas.
7. Sickit-Learn: é uma biblioteca que integra uma ampla gama de algoritmos de *Machine Learning*, para problemas supervisionados e não supervisionados, com foco em ter uma linguagem acessível democratizando o acesso ao aprendizado de máquina (PEDREGOSA et al., 2011).

Após conexão e *download* dos bancos de dados no Google Colab, iniciou-se o processo de tratamento da base. Um dos estágios previstos na literatura é o *Data Cleaning*, que ajuda na manutenção dos dados de entrada para o modelo, visto que o modelo de *Machine Learning* não terá resultados positivos caso a entrada de dados não seja consistente (Tembusai et al., 2021). A primeira parte desta etapa se refere aos dados duplicados. Neste caso optou-se pela aplicação da função de *drop_duplicates()*, responsável por identificar todas as linhas que tivessem os mesmos dados em todas as colunas, evitando assim que houvesse um peso duplicado onde não há necessidade.

A segunda parte dessa etapa diz respeito aos dados faltantes. O mais comum na literatura é entender primeiro se a origem desses dados faltantes é por ser um dado que não existe, ou um dado que não foi gravado. Caso se trate de um dado que não existe, não faz sentido tentar prever.

Neste trabalho, em busca de um melhor resultado, houve um corte de todas as linhas que, por algum motivo, não tivessem todos os dados de todas as colunas preenchidos, permanecendo a frequência da coluna de menor valor por meio da função *dropna()*. Essa decisão se deu principalmente visto que mesmo nas colunas com tais dados faltantes, ainda há muitos dados de entrada, o que, portanto, não afetaria o modelo.

Um primeiro filtro feito na base de pedidos (*df_order*) é a continuidade apenas de compras que chegaram ao *status* de *delivered* (entregue, em português), para que sejam comparadas apenas compras que possuem já uma nota atribuída, pois tal avaliação apenas é enviada ao cliente após cumprimento de todas as etapas de compra até o recebimento.

O próximo passo no pré-processamento será a substituição de colunas antes com datas, pela diferença entre os eventos. As datas que não se enquadram nisso serão descartadas. As datas que foram consideradas são listadas na Tabela 1 com seus respectivos *dataframes* de origem.

Tabela 1 – colunas de datas consideradas

| Coluna | Dataframe de Origem |
|-------------------------------|---------------------|
| review_creation_date | df_reviews |
| review_answer_timestamp | df_reviews |
| order_purchase_timestamp | df_orders |
| order_approved_at | df_orders |
| order_delivered_carrier_date | df_orders |
| order_delivered_customer_date | df_orders |
| order_estimated_delivery_date | df_orders |

Fonte: elaboração própria

Para um melhor entendimento da seleção das datas, a Figura 3 traz o processo de compra com as respectivas etapas.

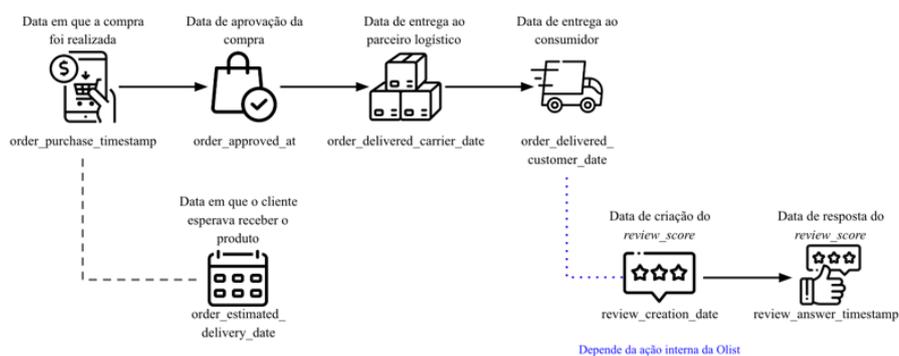


Figura 3: Fluxo das colunas de datas dentro do processo de compras do Olist
Fonte: Elaboração própria

A primeira coluna foi descartada e as próximas substituídas por um número inteiro como o número de dias que o pedido atrasou ou adiantou frente ao esperado pelo consumidor, tendo uma maior relevância do que ter duas colunas contendo datas distintas. A Figura 4 ilustra o processo de compras após as alterações de colunas.

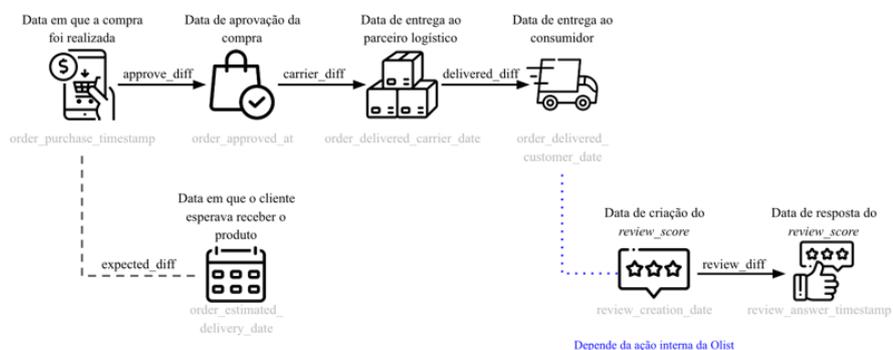


Figura 4: Fluxo das colunas de datas após alterações
Fonte: elaboração própria

Por conta da formatação presente no banco de dados do estudo, foi necessário primeiro converter as respectivas colunas para o formato *date time* e então feita a criação da nova coluna, contendo a operação mencionada. Após todas as análises iniciais com esse formato mais visual, houve uma nova conversão dos formatos resultantes da operação para *time delta*, números inteiros, para assim o modelo ser capaz de interpretá-los.

Outra etapa cumprida foi a criação de uma coluna chamada *'total_product_qtd'* que é responsável por apresentar quantos produtos são por compra, o que auxiliará em etapas seguintes.

A função *merge* também foi utilizada neste trabalho com a finalidade de unir *data frames* pelas suas chaves primárias (únicas). Antes de unir as bases, é preciso identificar quais dados realmente são interessantes para o objetivo do estudo, para assim escolher o melhor formato da aplicação (Estrella, 2020), por meio do requerimento *how* e da identificação por meio da Teoria dos Conjuntos, Figura 5.

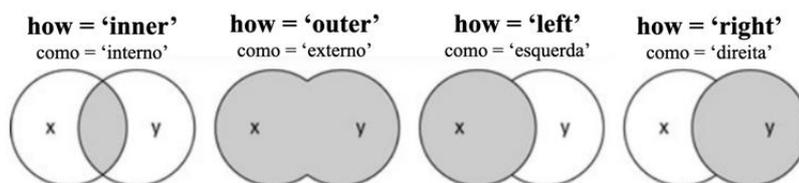


Figura 5: Teoria dos Conjuntos

Na aplicação da Teoria dos Conjuntos a junção de *data frames*, como mostra a Figura 5, inicialmente sempre será determinado pelo *df* que ocupará a posição da esquerda (exemplificado pela letra x), e qual *df* ocupará a posição da direita (exemplificado pela letra y).

Portanto, o *how* definirá o limite que o padrão deverá respeitar ao fazer o encontro dos dados, e no presente estudo utilizou-se o *left*. A função será utilizada em diferentes momentos da pesquisa, como no agrupamento do *df_orders* com *df_items*, mas tendo sua principal função em unir todos os bancos, que inicialmente eram separados, em uma única base que será utilizada para a análise final.

Uma das possíveis quebras presentes no banco de dados são as categorias, que vieram por padrão da base original, e são classificadas em 71 categorias, desde itens de casa até informática. Para tentar criar um recorte um pouco mais coerente, visto que as categorias podem ter comportamentos muito diferentes entre cada uma, buscou-se entender quais eram as 10 principais categorias, no quesito volume de pedidos, mais significativas, como mostra a Tabela 2.

Tabela 2 – 10 principais categorias

| Categoria | Número de pedidos |
|------------------------|--------------------------|
| beleza_saude | 954 |
| cama_mesa_banho | 922 |
| relogios_presentes | 843 |
| utilidades_domesticas | 693 |
| esporte_lazer | 629 |
| informatica_acessorios | 596 |
| moveis_decoracao | 513 |
| automotivo | 492 |
| telefonica | 385 |
| ferramentas_jardim | 284 |

Fonte: elaboração própria

O segundo filtro foi focado em apenas compras que tivessem 1 (um) item por pedido, visto que a variável dependente do estudo e algumas outras variáveis seguem o padrão de ser uma resposta por pedido, e não por item do pedido. À essa base final, deu-se o nome de *df_categorias*. Para uma melhor compreensão dos resultados, optou-se pela conversão das 5 classes (nota de 1 a 5) para 2, utilizando a média das notas para definir o limite, sendo a nova variável chamada de *score_class*. Sendo a classe 0 o *score* baixo (notas de 1 a 3) e a classe 1 o *score* alto (notas 4 e 5).

Por fim, as colunas finais presentes no banco de dados e que, portanto, são utilizadas para a aplicação dos modelos de *Machine Learning* estão descritas na Tabela 3, todas com o preenchimento de 5384 linhas:

Tabela 3 – colunas do banco de dados *df_categorias*

| Colunas | Descrição |
|---------------------|--|
| review_diff | Diferença de dias entre a criação do <i>review</i> da compra e o cliente respondê-la. |
| approve_diff | Diferença de dias entre a realização da compra e a aprovação dela. |
| carrier_diff | Diferença de dias entre a aprovação da compra e a entrega ao parceiro logístico. |
| delivered_diff | Diferença de dias entre a entrega ao parceiro logístico e a entrega ao consumidor. |
| expected_diff | Diferença de dias entre a data que o cliente esperava receber e a data em que recebeu. |
| price | Preço do pedido. |
| total_freight_value | Preço do frete. |
| product_name_lenght | Número de caracteres extraídos do nome do produto. |

| | |
|----------------------------|--|
| product_description_length | Número de caracteres extraídos da descrição do produto. |
| product_photos_qty | Número de fotos publicadas do produto. |
| product_weight_g | Peso do produto medido em gramas. |
| product_length_cm | Comprimento do produto medido em centímetros. |
| product_height_cm | Altura do produto medida em centímetros. |
| product_width_cm | Largura do produto medida em centímetros. |
| score_class | Conversão do <i>review score</i> em <i>Score Alto</i> e <i>Score Baixo</i> |

Fonte: elaboração própria

Técnicas de *Machine Learning*

No presente estudo, optou-se pela utilização do algoritmo de *Random Forest*, que é uma popular e poderosa técnica de *Machine Learning* (BREIMAN, 2001; SPEISER et al., 2019). Ele é feito pelo conjunto de muitas árvores de decisão individuais, que calculam entre si o resultado predito com partes randômicas da amostragem total, o que auxilia na redução do efeito de *overfitting* e melhora a generalização, e então escolhem pela maioria a classificação final (LOOSVELT et al., 2012). Os modelos de Árvore de Decisão são baseados em probabilidade em que cada *Tree* é uma *Decision Tree* individual, enquanto o *Random Forest* é o conjunto de muitas árvores. Portanto, o *Random Forest* geralmente fornece maior precisão em comparação com um único modelo de árvore de decisão, mantendo algumas das qualidades benéficas dos modelos de árvore (por exemplo, capacidade de interpretar relações entre preditores e resultado) (SPEISER; DURKALSKI; LEE, 2015).

Para iniciar a modelagem de *Machine Learning*, separou-se a base em duas sendo o *dataframe y_cat* contendo apenas a variável dependente (*score_class*) e o *dataframe x_cat* com todas as demais variáveis, que servirão para treino e teste do modelo. Em seguida, foi realizada a separação de ambos os *dataframes*, mas ainda sem haver dispersão da correspondência entre as duas tabelas, em treino e teste de forma aleatória e randômica.

Essa etapa tem a relevância de permitir que o modelo treine em uma base, o que permite ao pesquisador fazer ajustes, para então validar realmente sua performance rodando na base de teste. Assim, a base de teste se torna uma base desconhecida para o modelo, simulando a prática na vida real. A separação das compras foi feita em 80% para treino com uma amostragem de 4308, e 20% para teste com uma amostragem de 1078. O algoritmo *Random Forest* tende a ter uma melhor performance com a divisão selecionada (AVUÇLU; ELEN, 2020).

A etapa seguinte baseou-se em reduzir a variância das variáveis que seriam analisadas, por meio do método chamado *Z score* que é ilustrada pela função: $z\ score = \frac{(x-\mu)}{\sigma}$. Na expressão, x representa o valor que será normalizado em questão, como o valor da compra por exemplo, μ representa a média das variáveis e σ representa o desvio padrão (AHMED et al., 2019). A essa etapa dá-se o nome de *Feature Scaling*, e deve ser feita após a separação em treino e teste para que

uma não tenha influência na normalização da outra, afinal o intuito é principalmente que a base de teste nunca tenha sido vista pela máquina. Em seguida, iniciou-se o processo de aplicação do método do *Random Forest*.

Matriz de Confusão

A apresentação dos resultados é feita por meio da Matriz de Confusão, Figura 6, que é comumente utilizada para entender o comportamento de classificação de cada categoria em modelos supervisionados de classificação (HASNAIN et al., 2020). A matriz é composta por linhas e colunas, e apresenta duas possibilidades de resposta, *score* alto ou *score* baixo, a matriz terá dimensão 2x2. Para a análise dos resultados no estudo, foi criada uma Matriz de Confusão referente ao percentual de acerto de determinada categoria.

| | | Valor real | | |
|----------------|-------------|----------------------------|----------------------------|--------|
| | | Score Baixo | Score Alto | |
| Valor previsto | Score Baixo | Verdadeiro Positivo (VP %) | Falso Negativo (FN %) | = 100% |
| | Score Alto | Falso Positivo (FP %) | Verdadeiro Negativo (VN %) | = 100% |

Figura 6: Modelo prático da Matriz de Confusão no estudo

Cada valor pode ser descrito por:

- Verdadeiro Positivo (VP %): mostra quantos registros foram classificados como *score* alto e o *score* realmente era alto, dividido pelo total de vezes que o classificador previu *score* baixo, independentemente do valor real.
- Verdadeiro Negativo (VN %): mostra quantos registros foram classificados como *score* baixo e o *score* realmente era baixo, dividido pelo total de vezes que o classificador previu *score* alto, independentemente do valor real.
- Falso Positivo (FP %): mostra quantos registros foram classificados como *score* alto e o *score* era baixo, dividido pelo total de vezes que o classificador previu *score* alto, independentemente do valor real.
- Falso Negativo (FN %): mostra quantos registros foram classificados como negativos incorretamente, ou seja, a resposta do classificador foi que o *score* era baixo e o *score* era alto.

Os indicadores de performance descritos a seguir podem ser obtidos com os valores absolutos resultantes da Matriz de Confusão:

- Acurácia: responsável por mostrar a performance geral do modelo, ou seja, de todas as classificações requisitadas, quantas foram feitas corretamente. A acurácia é calculada pela expressão:

$$\text{Acurácia} = \frac{VP+VN}{VP+VN+FP+FN}$$
- Precisão: dentre todas as classificações de classe positivo, VP e FP, quantas estão corretas.
$$\text{Precisão} = \frac{VP}{VP+FP}$$

E por fim, utilizou-se o *Feature Importance*, que é um valor usado para ordenar as variáveis por importância de impacto na variável dependente (KANG; RYU, 2019), *score_class*, em uma escala de 0 a 100%. Esse *ranking* que pode auxiliar na seleção de variáveis para testes posteriores, auxiliando na performance e reduzindo os recursos necessários, se tornando mais rápido. Estas preocupações com performance computacional tendem a ser necessárias proporcionalmente ao tamanho da amostragem.

O processo todo descrito a partir da separação do *df_categorias* em treino e teste até a obtenção do *Feature Importance* foi repetido 10 vezes e então retirada as médias simples. x é representado pela saída do *Feature Importance*: $média = \frac{\sum x}{n}$. Após isso, repetiu-se o processo com foco nas 5 variáveis de maior relevância e bases separadas por categoria, focada nas 10 principais categorias em volume de compras.

RESULTADOS E DISCUSSÕES

A primeira aplicação foi com base no *df_categorias*, data frame contendo toda a base das 10 categorias juntas, com foco em entender por meio do *Feature Importance* quais variáveis tinham uma maior relevância, para então utilizá-las na análise por categoria. Com 4308 amostras no treino, sendo desses 73,9% (3184) da classe 1, ou seja, definidos como *score* alto, e 1078 no teste, sendo desses coincidentemente 73,9% (797) também da classe 1. Com essas amostragens, obteve-se uma acurácia média das 10 aplicações de 79,1%. Também tirou-se uma média simples da matriz de confusão de cada *loop*, tendo o resultado mostrado na Figura 7.

O método da *random forest* é amplamente utilizado devido à sua eficácia em lidar com uma variedade de problemas de classificação e regressão. Trata-se de uma metodologia especialmente adequada para conjuntos de dados complexos, com muitas variáveis e uma vasta quantidade de dados. Uma das vantagens do *random forest* é a sua capacidade de lidar bem com dados de diferentes tipos, sejam eles numéricos, categóricos ou uma combinação de ambos. Cada árvore nesse método é uma tomada de decisão independente, construída com uma amostra aleatória dos dados e utilizando apenas um subconjunto das características disponíveis (RIGATTI, 2017).

Por outro lado, a matriz de confusão é uma ferramenta de avaliação usada para avaliar o desempenho de modelos de classificação. Ela mostra a frequência com que cada classe real foi classificada correta ou incorretamente pelo modelo (MONARD & BARANAUSKAS, 2000).

A combinação do método *random forest* com a matriz de confusão oferece várias vantagens, dentre elas uma visão completa do desempenho do modelo, a identificação de padrões complexos e a detecção de desequilíbrios de classe. Em resumo, essa combinação é adequada para uma ampla gama de conjuntos de dados e problemas de classificação, proporcionando uma avaliação abrangente e detalhada do desempenho do modelo.

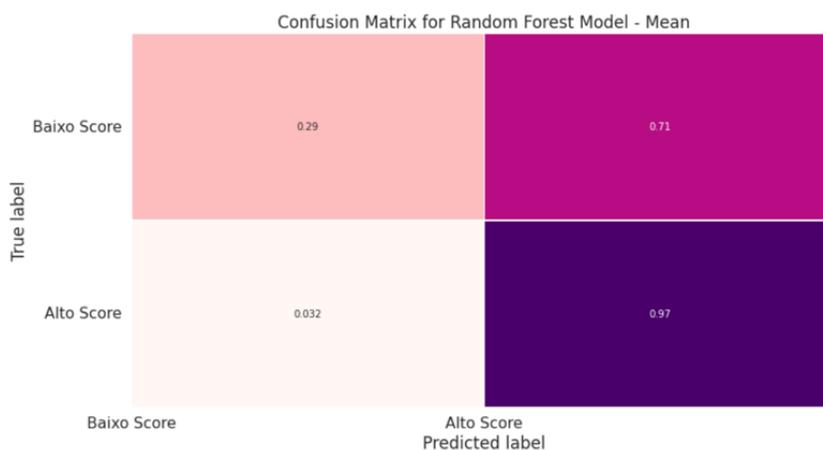


Figura 7: Matriz de Confusão com resultados médios do *Random Forest*

Outros resultados de performance são observados na Tabela 4.

Tabela 4 – Indicador de performance do algoritmo de *Machine Learning*

| Classe | Precisão |
|-------------|----------|
| Baixo Score | 0,76 |
| Alto Score | 0,79 |

Por fim, ao analisar o resultado de cada *looping* da aplicação do *Random Forest* para o *df_categorias*, concluiu-se que as 5 variáveis mais relevantes para a decisão da nota atribuída após todo o serviço são *expected_diff*, *delivered_diff*, *review_diff*, *carrier_diff* e *approve_diff*. Os resultados de cada *loop*, sendo 0 o primeiro e 9 o 10º, e a média final são mostrados na Tabela 5.

Tabela 5 – *Feature Importance* por *loop*

| Variável | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | Média |
|-----------------------|------|------|------|------|------|------|------|------|------|------|-------|
| <i>expected_diff</i> | 0,13 | 0,13 | 0,14 | 0,14 | 0,14 | 0,14 | 0,14 | 0,14 | 0,14 | 0,14 | 0,14 |
| <i>delivered_diff</i> | 0,11 | 0,10 | 0,10 | 0,10 | 0,10 | 0,10 | 0,10 | 0,10 | 0,10 | 0,10 | 0,10 |
| <i>review_diff</i> | 0,09 | 0,09 | 0,09 | 0,09 | 0,09 | 0,09 | 0,09 | 0,09 | 0,09 | 0,09 | 0,09 |
| <i>carrier_diff</i> | 0,09 | 0,09 | 0,09 | 0,08 | 0,09 | 0,08 | 0,08 | 0,08 | 0,08 | 0,08 | 0,08 |
| <i>approve_diff</i> | 0,07 | 0,07 | 0,07 | 0,07 | 0,07 | 0,07 | 0,07 | 0,07 | 0,07 | 0,07 | 0,07 |

Cada variável contém uma importância nesse contexto de avaliação de produtos no comércio eletrônico. O conceito *diff* que acompanha todas as variáveis indica a diferença entre as variáveis apresentadas. Em caso de grande discrepância entre as variáveis analisadas, isso pode a insatisfação do cliente e a perda de confiança na marca.

A variável *expected* representa a expectativa inicial do cliente em relação ao produto ou serviço que está adquirindo. É crucial que as expectativas do cliente sejam atendidas ou superadas durante o processo de compra, para assim garantir a sua satisfação.

A variável *delivered* mostra o que foi efetivamente entregue ao cliente. Garantir que essa entrega corresponda exatamente ao que foi solicitado é fundamental para manter a confiança do cliente e garantir uma experiência positiva de compra.

A variável *review* reflete a diferença entre a revisão esperada do produto e a revisão real. Avaliações positivas são de extrema importância para atração de novos clientes e construção de uma reputação sólida no comércio eletrônico. Uma possível disparidade entre as revisões pode afetar significativamente a credibilidade da empresa.

A variável *carrier* destaca a diferença entre a transportadora esperada e a transportadora real utilizada para entrega do produto. A escolha da transportadora certa afeta diretamente a rapidez, a segurança e a qualidade da entrega, influenciando diretamente a satisfação do cliente.

A variável *approve* indica se o produto entregue foi aprovado pelo cliente. Essa aprovação é um componente fundamental para alcançar a satisfação do consumidor.

Por concluir que essas são as mais relevantes, então repetiu-se o processo em bases separadas por categoria, para que não haja influência do padrão de uma categoria no resultado da outra, apenas com as variáveis tidas como mais relevantes. Após efetuada a análise por categoria, foi considerada a média dos *Feature Importances* médios de cada variável por categoria, assim como a acurácia média, para então ter uma melhor análise dos comportamentos e previsões. O resultado médio da relevância de cada variável e a acurácia para cada categoria é mostrado na Tabela 6, e a acurácia média por categoria.

Tabela 6 – Feature Importance médio de cada variável e Acurácia por categoria

| Categoria | expected_ diff | delivered_ diff | review_ diff | carrier_ diff | approve_ diff | Acurácia |
|------------------------|-------------------|--------------------|-----------------|------------------|------------------|----------|
| cama_mesa_banheiro | 22,75 | 21,77 | 19,74 | 18,87 | 16,87 | 78,61 |
| beleza_saude | 28,22 | 19,78 | 19,11 | 17,94 | 14,92 | 83,17 |
| esporte_lazer | 26,32 | 22,44 | 18,32 | 17,90 | 14,99 | 80,45 |
| informatica_acessorios | 22,47 | 20,23 | 19,74 | 19,25 | 18,29 | 78,03 |
| moveis_decoracao | 23,21 | 20,26 | 19,36 | 18,93 | 18,21 | 72,81 |
| utilidades_domesticas | 26,85 | 20,71 | 18,63 | 17,51 | 16,28 | 84,07 |
| relogios_presentes | 24,07 | 23,39 | 19,72 | 17,69 | 15,15 | 68,37 |
| telefonias | 22,30 | 21,29 | 20,42 | 19,42 | 16,54 | 69,43 |
| automotivo | 24,88 | 21,35 | 19,09 | 18,05 | 16,60 | 79,33 |
| brinquedos | 27,59 | 19,99 | 19,38 | 18,26 | 14,75 | 83,00 |

Com resultado unânime entre as categorias tem-se a variável *expected_diff* que é a de maior relevância para a predição do modelo e ilustra que o fator entrega têm um grande peso ao se tratar de compras feitas pela internet, mesmo antes da pandemia da COVID-19.

O desafio das redes varejistas é conseguir acompanhar a evolução deste novo tipo de compra, e consumidores cada vez mais empoderados, que querem um prazo cada vez menor e pontualidade de entrega, tendo a urgência com grande relevância (Shetty et al., 2018).

Para suprir a necessidade de uma entrega rápida e pontual ao cliente, é necessário que cada vez mais as empresas foquem em desenvolver uma boa logística, com processos claros e eficientes.

O presente estudo mostrou que as categorias em que essa necessidade se torna mais presente são de beleza/saúde ($\bar{x}_{FI} = 28,73\%$) e brinquedos ($\bar{x}_{FI} = 27,60\%$), também é importante ressaltar que ambas obtiveram bons resultados de acurácia, reforçando o diagnóstico.

Segundo pesquisa feita pela McKinsey & Company, focada no investimento de alguns países para classificações de bem-estar, é possível observar uma grande importância das áreas de saúde e aparência, se refere ao investimento em produtos de beleza para o Brasil, o que traz um reflexo na necessidade crescente de prazos mais pontuais para os consumidores da categoria (Callaghan, 2021).

Na visão contrária, pode-se concluir que as categorias de informática/acessórios ($\bar{x}_{FI} = 22,47\%$) e móveis/decoração ($\bar{x}_{FI} = 23,21\%$) são as categorias com a melhor distribuição entre as variáveis, mas com acurácias não tão altas.

Entretanto, segundo pesquisa feita pela Opinion Box (2021), consumidores do mercado de móveis e decoração *online* são mais apegados a marca em si, seguidos

por outros fatores como material, durabilidade, preço e, por último, o frete. Diante disso, estas comparações permitem confirmar que os comportamentos podem ser particulares a cada categoria.

CONCLUSÃO

O estudo permitiu atingir o objetivo de identificar os fatores que impactam no *review score* dado pelo consumidor após a experiência via Olist Store, o que ilustra a satisfação com o produto e serviço prestados por meio da técnica de *Machine Learning* do *Random Forest*, e após isso aplicado o *Feature Importance* para ranqueamento da relevância das variáveis perante a predição.

Com a forte relevância das *startups* no cenário econômico brasileiro e mundial, e com a importância da Olist dentro desse cenário, associado a aplicação de técnicas de Inteligência Artificial para a resolução de problemas complexos e trabalhosos para um ser humano, deu-se a presente pesquisa. Por meio da aplicação da técnica do *Random Forest*, modelo de *Machine Learning*, conseguiu-se ensinar a máquina a tentar prever qual *score_class* seria atribuído a cada compra, e com a aplicação do *Feature Importance* descobriu-se quais variáveis tiveram uma maior influência nessa nota.

Após realizar o processo em uma base geral, houve a quebra e repetição do processo dentro de cada categoria, dentre as mais relevantes em volume de compras, permitindo assim a conclusão de que as categorias de beleza/saúde e brinquedos são as mais afetadas pela diferença de dias entre o esperado para a entrega do produto, e a data realmente entregue, reforçando a necessidade que as empresas de *e-commerce* estejam muito alinhadas com os processos de logística.

Além disso, observa-se que a pandemia ampliou a expansão do mercado de compras *online*, em que o consumidor está cada vez mais exigente quanto a prazo de entrega, e o cumprimento desses, sendo este um item que merece atenção dos empreendedores. Por isso, os lojistas que utilizam o Olist como canal intermediário para venda de seus produtos estão fortalecendo a presença *online* e aumentando as vendas. Dentre os principais benefícios para este segmento de empresas destacam-se: alta visibilidade, tráfego qualificado, tecnologia robusta e baixos investimentos.

Ao explorar os dados por meio de técnicas de *Machine Learning*, é possível identificar padrões e tendências nas relações de consumo dentro do *SaaS enabled marketplace*, e compreender os motivos relacionados a esses comportamentos. Ao estabelecer conexões mais profundas entre os dados observados e a teoria discutida nesse trabalho, é possível extrair *insights* significativos que fornecem uma compreensão mais holística do fenômeno observado. Com isso, essa abordagem fortalece a análise dos resultados da pesquisa, e também contribui para uma contínua evolução do Olist, fornecendo uma base sólida para futuras estratégias e decisões.

Desde a sua fundação, o Olist vem amadurecendo consideravelmente enquanto um unicórnio e vem apresentando um crescimento significativo. A empresa vem trilhando um caminho de inovação e expansão, consolidando-se como uma das principais plataformas de comércio eletrônico, não só no Brasil, mas no mundo.

Ao se tornar um unicórnio, a Olist demonstrou uma capacidade de atrair investimentos significativos e a habilidade de criar valor e impacto no mercado. Esses são resultados de uma visão estratégica aliada à qualidade dos serviços oferecidos.

Além disso, o Olist vem ampliando sua base de clientes, e diversificando os produtos e serviços oferecidos. Com isso, a empresa vem investindo em tecnologia, logística e experiência do cliente com o objetivo de melhorar sua oferta e manter sua competitividade em um mercado que se mostra cada vez mais dinâmico.

A base de dados utilizada para análise dos pontos levantados foi disponibilizada gratuitamente no Kaggle, entre os anos de 2016 a 2018. Portanto, entende-se que esta é uma limitação temporal e que futuras bases de dados atualizadas devem ser analisadas até como uma forma comparativa, devido ao crescimento da empresa no mercado em vários setores.

Por fim, concluiu-se que as empresas de *e-commerce*, possuem vantagens competitivas ao explorar e investir em tecnologia, para entender melhor seu mercado e padrões de consumo, assim como buscar oportunidades e melhorar a eficiência, desde que haja uma boa coleta e armazenamento de dados.

Como limitação e sugestão para futuras pesquisas, tem-se a possibilidade de olhar para a classificação do *review_score* e, por exemplo, entender se a distância entre comprador e vendedor pode ter uma influência positiva no *expected_diff*, consequentemente melhorando a satisfação com a integração de comércios regionais.

REFERÊNCIAS

AHMED, U.; MUMTAZ, R.; ANWAR, H.; SHAH, A. A.; IRFAN, R.; GARCIA-NIETO, J. Efficient Water Quality Prediction Using Supervised Machine Learning. **Water**, v. 11, n. 11, p. 2210, 2019.

ABSTARTUPS. Modelo de Negócio para Startups: É melhor criar um SaaS ou uma Plataforma Marketplace? **Associação Brasileira de Startups**. Disponível em: <https://abstartups.com.br/modelo-de-negocio-para-startups-e-melhor-criar-um-saas-ou-uma-plataforma-marketplace/>. Acesso em: 03 mar. 2022.

AVUÇLU, E.; ELEN, A. Evaluation of train and test performance of machine learning algorithms and Parkinson diagnosis with statistical measurements. **Medical & Biological Engineering & Computing**, v. 58, n. 11, p. 2775–2788, 2020.

BANACHEWICZ, K.; MASSARON, L. **The Kaggle Book: Data analytics and machine learning for competitive data science**. Packt Publishing; 1ª Edição, 2022.

BARRETT, P.; HUNTER, J.; MILLER, J.T; HSU, J-C; GREENFIELD, P. Matplotlib – A Portable Python Plotting Package. In: ASTRONOMICAL DATA ANALYSIS SOFTWARE AND SYSTEMS, 91, São Francisco. Proceedings... São Francisco, 2005.

BIRKETT, A. **What is Customer Satisfaction Score (CSAT)?** Disponível em: <https://blog.hubspot.com/service/customer-satisfaction-score>. Acesso em: 03 mar. 2022.

BREIMAN, L. Random Forests. **Machine Learning**, v. 45, n. 1, p. 5–32, 2001.

CASTLE, N. **What is Semi-Supervised Learning? Oracle AI & Data Science Blog.** Disponível em: <https://blogs.oracle.com/ai-and-datascience/post/what-is-semi-supervised-learning>. Acesso em: 03 mar. 2022.

CALLAGHAN, S.; LOSCH, M.; PIONE, A.; TEICHLNER, W. Sentir-se bem: O futuro do mercado de bem-estar de \$ 1,5 trilhão | **McKinsey**, maio 17. Disponível em: <https://www.mckinsey.com/industries/consumer-packaged-goods/our-insights/feeling-good-the-future-of-the-1-5-trillion-wellness-market/pt-BR>. Acesso em: 03 mar. 2022.

COSTI, G. Aprendizagem Não Supervisionada. **Lambda3**. Disponível em: <https://lambda3.com.br>. Acesso em: 03 mar. 2022.

CHAN, D. Y.; VASARHELYI, M. A. Innovation and practice of continuous auditing. **International Journal of Accounting Information Systems**, v. 12, n. 2, p. 152–160, 2011.

ESTRELLA, C. Pandas: Combinando data frames com merge() e concat(). **Data Hackers**, 2020. Disponível em: <https://medium.com/data-hackers/pandas-combinando-data-frames-com-merge-e-concat-10e7d07ca5ec>. Acesso em: 03 mar. 2022.

EUROMONITOR. Understanding Global Marketplace Trends. **Market Research Report**, 2018. Disponível em: <https://www.euromonitor.com/understanding-global-marketplace-trends/report>. Acesso em: 03 mar. 2022.

FACEII, K.; LORENA, A. C.; GAMA, J.; CARVALHO, A. C. P. DE L. F. DE. **Inteligência Artificial: uma abordagem de aprendizado de máquina**. Rio de Janeiro: LTC, 2011.

FIGUEIREDO, V. Seus primeiros passos como Data Scientist: Introdução ao Pandas! **Data Hackers**. Disponível em: <https://medium.com/data-hackers/uma-introdu%C3%A7%C3%A3o-simples-ao-pandas-1e15eea37fa1>. Acesso em: 03 mar. 2022.

GUNAWAN, T. S.; ASHRAF, A.; RIZA, B. S.; HARYANTO, E. V.; ROSNELLY, R.; KARTIWI, M.; JANIN, Z. Development of video-based emotion recognition using deep learning with Google Colab. **TELKOMNIKA (Telecommunication Computing Electronics and Control)**, v. 18, n. 5, p. 2463–2471, 2020.

HASNAIN, M. PASHA, M. F.; GHANI, I.; IMRAN, M.; ALZHRANI, M. Y.; BUDIARTO, R. Evaluating Trust Prediction and Confusion Matrix Measures for Web Services Ranking. **IEEE Access**, v. 8, p. 90847–90861, 2020.

HUTTER, F.; KOTTHOFF, L.; VANSCHOREN, J. **Automated Machine Learning: Methods, Systems, Challenges**. The Springer Series on Challenges in Machine Learning, 2019.

JORDAN, M. I.; MITCHELL, T. M. Machine learning: Trends, perspectives, and prospects. **Science**, v. 349, n. 6245, p. 255–260, 2015.

KANG, K.; RYU, H. Predicting types of occupational accidents at construction sites in Korea using random forest model. **Safety Science**, v. 120, p. 226–236, 2019.

KUMAR, A.; SYED, A. A.; PANDEY, A. Adoption of online resources to improve the marketing performance of SMES. **Asia-Pacific Journal of Health Management**, v. 16, n. 3, 2021.

LIMA, J. M.; CÔRREA, R. O.; das CHAGAS, D. A.; OLIVEIRA, T. de S.; de CARVALHO, G. D. G. Empreendedorismo como aporte para o empoderamento econômico feminino. **Revista Tecnologia e Sociedade**, v.17, n. 48, 2021.

LOOSVELT, L.; PETERS, J.; SKRIVER, H.; DE BAETS, B.; VERHOEST, N. Impact of reducing polarimetric SAR input on the uncertainty of crop classifications based on the random forests' algorithm. **IEEE Transactions on Geoscience and Remote Sensing**, v. 50, n. 10, p. 4185–4200, 2012.

MCCOY, J. T.; AURET, L. Machine learning applications in minerals processing: A review. **Minerals Engineering**, v. 132, p. 95–109, 2019.

MOHRI, M.; ROSTAMIZADEH, A.; TALWALKAR, A. (2018). **Foundations of Machine Learning** (2° ed). MIT Press.

MONARD, M. C.; BARANAUSKAS, J. A. Reviewing some machine learning concepts and methods. **Technical Report 102, Instituto de Ciências Matemáticas e de Computação**, 2000.

NIWATE, T. S. (2021). Impact on revenue generation of Olist ecommerce company on the basis of various product parameters. **Electronic Theses, Projects, and Dissertations**. 1315. Disponível em: <https://scholarworks.lib.csusb.edu/etd/1315>. Acesso em: 03/03/2022.

PAÇO, A. do.; SHIEL, C.; ALVES, H. A New Model for Testing Green Consumer Behaviour. **Journal of Cleaner Production**, 207, 2018.

PEDREGOSA, F., VAROQUAUX, G., GRAMFORT, A.; MICHEL, V., THIRION, B.; GRISEL, O., BLONDEL, M., PRETTENHOFER, P., WEISS, R., DUBOURG, V., VANDERPLAS, J., PASSOS, A., COURNAPEAU, D., BRUCHER, M., PERROT, M., & DUCHESNAY, É. (2011). Scikit-learn: Machine Learning in Python. **The Journal of Machine Learning Research**, v.12, p. 2825–2830.

PINOCHET, L. H. C., ALVES, F. R. R., LOPES, E. L., HERRERO, E., & BRELAZ, G. (2021). 'From cloud to the board' - identification of triggers adoption in Brazilian

companies. **International Journal of Business Information Systems**, v. 38, p. 343-366.

POBEE, F. (2021). Modeling the Factors that Influence Ghanaian Entrepreneurs to Adopt e-Commerce. **International Journal of Innovation and Technology Management**, v. 18 n. 6, 2150029-18, p. 1-18.

RATH, S. B.; BASU, P.; MANDAL, P.; PAUL, S. (2021). Financing models for an online seller with performance risk in an E-commerce marketplace. **Transportation Research Part E: Logistics and Transportation Review**, v. 155, 102468.

RODRIGUES, G. A. B.; PLENS, M.; PIANTINO, N. P. De A. (2020). Gestão em vendas online: Estudo de caso de empresa calçadista com modelo de negócio em marketplace. **Revista Empreenda UniToledo Gestão, Tecnologia e Gastronomia**, v. 4, n. 1.

SAMUEL, N. (2019). **The Python Math Library. Stack Abuse**. Disponível em: <<https://stackabuse.com/the-python-math-library/>>. Acesso em: 03 mar. 2022.

SANTIAGO JR., L. (2019). Entendendo a biblioteca NumPy. **Ensina.AI**. Disponível em: <<https://medium.com/ensina-ai/entendendo-a-biblioteca-numpy-4858fde63355>>. Acesso em: 03 mar. 2022.

SANTOS, H. G. dos.; NASCIMENTO, C. F. do.; IZBICKI, R.; DUARTE, Y. A. O.; CHIAVEGATTO FILHO, A. D. P. (2019). Machine learning para análises preditivas em saúde: Exemplo de aplicação para predizer óbito em idosos de São Paulo, Brasil. **Cadernos de Saúde Pública**, v. 35 p. 7.

SCUTARIU, A-L.; SUSU, S.; HUIDUMAC-PETRESCU, A-E.; GOGONEA, R-M. (2022). A Cluster Analysis Concerning the Behavior of Enterprises with E-Commerce Activity in the Context of the COVID-19 Pandemic. **Journal Theoretical and Applied Electronic Commerce Research**, v. 17 n. 1, p. 47-68.

SHEHATA, G.; MONTASH, M. (2020). Driving the internet and e-business technologies to generate a competitive advantage in emerging markets Evidence from Egypt. **Information Technology & People**, v. 33 n. 2, p. 389–423.

SHETTY, A. S.; JEEVANANDA, S. (2018). How to win back the disgruntled consumer? The omni-channel way. **Journal of Business & Retail Management Research (JBRMR)**, v. 12 n. 4, p. 200-207.

SILVA, V. M. da. (2021). Estudos de futuro e foresight para ciência, tecnologia e inovação: tendências do uso de big data e machine learning. Tese de Doutorado. **Universidade Estadual de Campinas**. Disponível em: <http://repositorio.unicamp.br/Acervo/Detalhe/1164751> Acesso em: 03/03/2022.

SPEISER, J. L.; DURKALSKI, V. L.; LEE, W. M. (2015). Random forest classification of etiologies for and orphan disease, **Statistic in Medicine**, v. 34 n. 5, p. 887-899.

SPEISER, J. L.; MILLER, M. E.; TOOZE, J.; IP, E. A comparison of random forest variable selection methods for classification prediction modeling. **Expert Systems with Applications**, v. 134, p. 93-101, 2019.

TEMBUSAI, Z. R.; MAWENGGANG, H.; ZARLIS, M. K-Nearest Neighbor with K-Fold Cross Validation and Analytic Hierarchy Process on Data Classification. **International Journal of Advances in Data and Information Systems**, v. 2, n. 1, p. 1-8, 2021.

TSUNODA, D. F.; da CONCEIÇÃO MOREIRA, P. S.; GUIMARÃES, A. J. R. Machine Learning e revisão sistemática de literatura automatizada: uma revisão sistemática. **Revista Tecnologia e Sociedade**, v. 16, n. 45, 2020.

WANG, J.; WU, J.; SUN, S.; WANG, S. The relationship between attribute performance and customer satisfaction: An interpretable machine learning approach. **Data Science and Management**, 2024.

WASKOM, M. Seaborn: Statistical data visualization. **Journal of Open-Source Software**, v. 6, n. 60, p. 3021, 2021.

ZAGHLOUL, M.; BARAKAT, S.; REZK, A. Predicting E-commerce customer satisfaction: Traditional machine learning vs. deep learning approaches. **Journal of Retailing and Consumer Services**, v. 79, 2024.

Recebido: 09/05/2023

Aprovado: 20/05/2024

DOI: 10.3895/rts.v20n60.16920

Como citar:

LEITE, Gabriela Amaral de Alencar; PINOCHET, Luis Hernan Contreras; PARDIM, Vanessa Itacaramby et al. Uso de Machine Learning para entendimento das relações de consumo no modelo de negócio (SaaS enabled Marketplace) do Olist. *Tecnol. Soc., Curitiba*, v. 20, n. 60, p.328-350, abr./jun., 2024. Disponível em:

<https://periodicos.utfpr.edu.br/rts/article/view/16920>

Acesso em: XXX.

Correspondência:

Direito autoral: Este artigo está licenciado sob os termos da Licença Creative Commons-Atribuição 4.0 Internacional.

