

PROJETO DE CIÊNCIA DE DADOS: Perfil de Clientes (Clusterização)

Introdução

Esse projeto é uma iniciativa pessoal voltada para o aprendizado.

Durante o projeto um dataset sobre Perfil de Clientes é explorado, para passar por etapas de Análise, Transformação, Clusterização e entrega de análises.

Link para o dataset: [Marketing Campaign \(kaggle.com\)](https://www.kaggle.com/marketing-campaign)

Tecnologias utilizadas:

Python 3 - [Welcome to Python.org](https://www.python.org/)

Anaconda - [Anaconda | The World's Most Popular Data Science Platform](https://www.anaconda.com/)

Jupyter Notebook - [Project Jupyter | Home](https://projectjupyter.org/)

Power Bi - [Power BI - Visualização de Dados | Microsoft Power Platform](https://powerbi.microsoft.com/pt-br/)

Análise Exploratória e Limpeza e Transformação de Dados

Durante a fase de análise exploratória, eu busquei compreender o conjunto de dados e encontrar qualquer tipo de problema que possa causar complicações nas próximas fases do projeto.

O primeiro problema encontrado foi a presença de um outlier nos valores de renda, identificado ao utilizar o método "describe()" no dataframe. Para visualizar melhor o problema, utilizei 2 gráficos: Boxplot e Dispersão

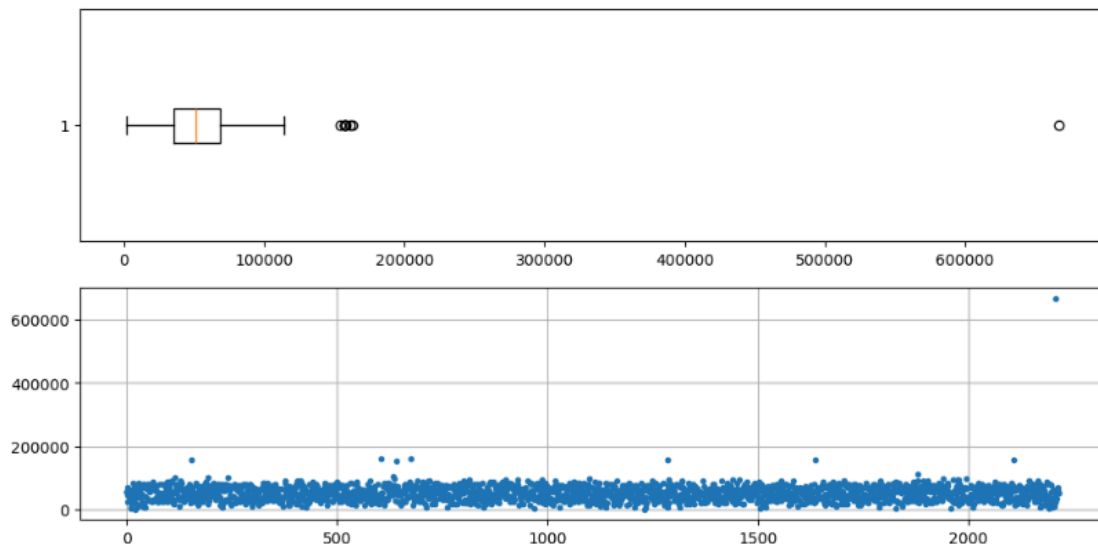


Figura 1: Detecção de Outlier

Após visualizar esse valor separadamente, cheguei a conclusão que esse valor pode ter sido um erro de digitação, uma vez que estando tudo dentro dos conformes, o valor da renda seja muito específico: “666666”. Como se trata de apenas um valor, mas que estava interferindo consideravelmente nas medidas devido ao seu alto valor, optei por remover essa observação do conjunto de dados.

Também é possível visualizar alguns outros outliers no boxplot, entretanto, eles parecem ser valores reais e não apresentam implicações significativas no conjunto de dados, portanto, não vejo necessidade de tratá-los.

O próximo problema identificado foi a presença de alguns valores nulos, detectados ao utilizar o método `isna().sum()`. Tratei esses valores substituindo-os pela média da respectiva coluna.

Por fim, eu removi clientes que tenham idade superior a 100 anos, pois eles representam um padrão de compra mais complicado para a análise, já que muitas vezes essas compras nem são realizadas por eles, mas sim por seu filhos e cuidadores, entretanto isso representa uma parcela muito pequena no conjunto, portanto não provoca mudanças significativas.

O que é Clusterização?

A clusterização é uma técnica de aprendizado de máquina não supervisionada, que visa dividir um conjunto de objetos em grupos semelhantes. Essa técnica de agrupamento é utilizada para lidar com grandes volumes de dados, comumente usada para encontrar padrões e criar sistemas de recomendação.

Clusterização

Para realizar a clusterização, é necessário selecionar as colunas numéricas relevantes para o agrupamento. Neste projeto, escolhi 14 colunas que fornecem informações sobre o cliente e seus padrões de compra:

'Income', 'Age', 'Kidhome', 'Teenhome', 'MntWines', 'MntFruits', 'MntMeatProducts', 'MntFishProducts', 'MntSweetProducts', 'MntGoldProds', 'NumDealsPurchases', 'NumWebPurchases', 'NumCatalogPurchases', 'NumStorePurchases'

Após a seleção das colunas, é necessário normalizar os dados, ou seja, ajustar todos os valores para uma mesma escala, o que é uma exigência para a entrada de dados no sklearn.

O próximo passo é determinar o número de clusters. Existem várias técnicas para definir o número ideal de clusters para um conjunto de dados. Neste projeto, optei por utilizar o método do cotovelo, que envolve a plotagem da curva de variância em função do número de clusters. Normalmente, o ponto onde a curva começa a se achatar indica o número ideal de clusters para o conjunto de dados.

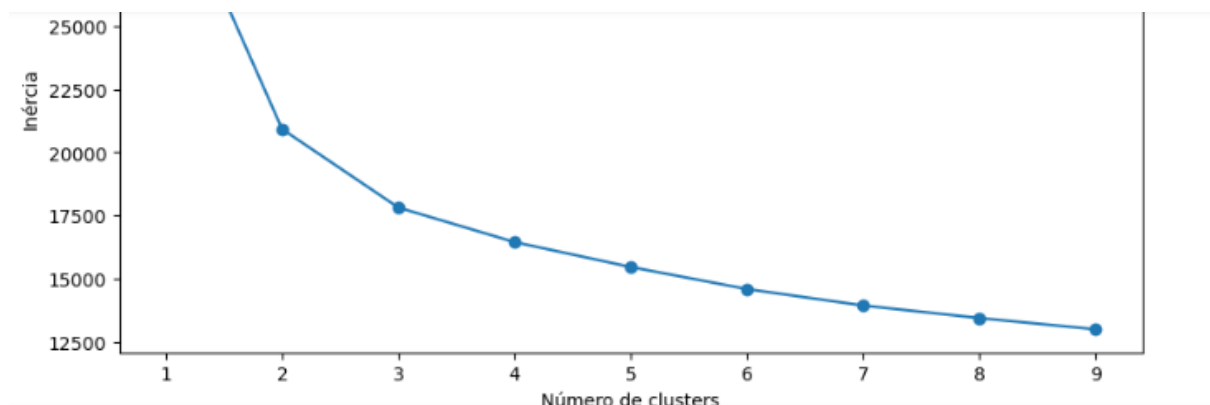


Figura 2 - Método do Cotovelo

Com base na forma da curva no gráfico, é concluído que 3 é um número interessante de clusters para este conjunto de dados.

Com tudo pronto, o próximo passo é treinar o algoritmo para concluir a etapa de clusterização.

Com o algoritmo treinado, é possível ver os resultados por meio de um gráfico de distribuição de valores pela idade e renda:

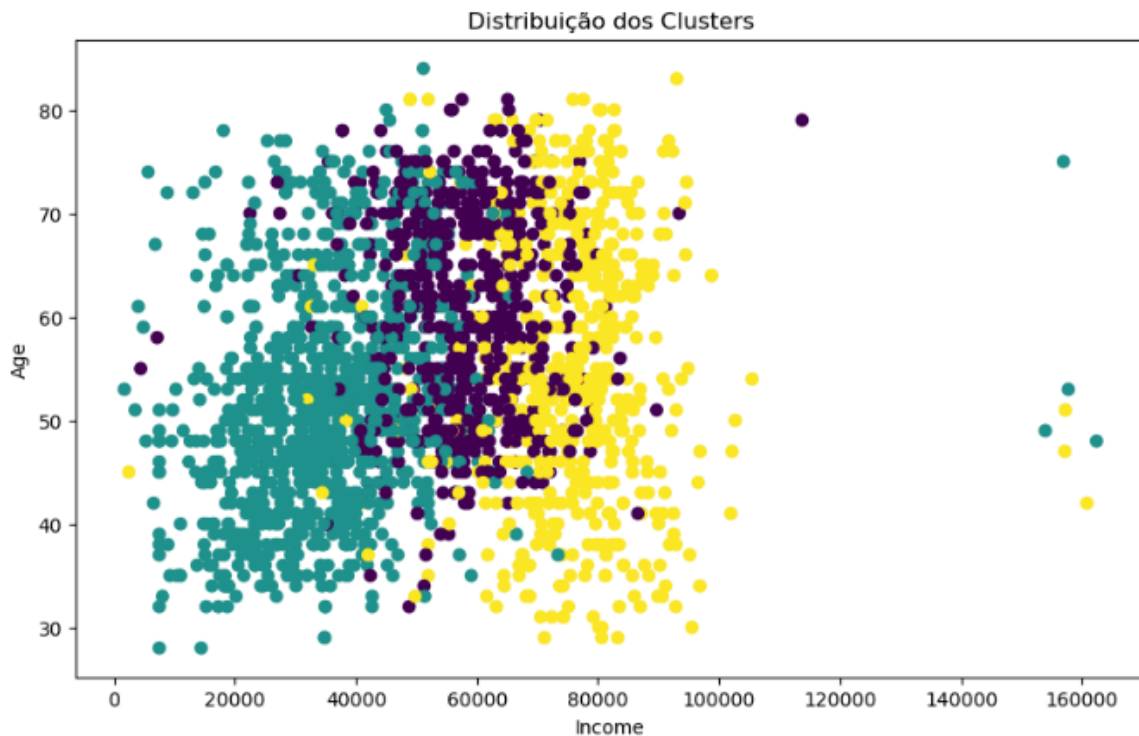


Figura 3: Distribuição dos Clusters

Análises e Insights (Power BI)

Como último passo do projeto, levar os dados preparados para o Power BI é uma excelente forma de montar análises e extrair valor dos dados. O conjunto de dados é focado no perfil dos clientes, permitindo a compreensão de padrões entre os grupos.

Para aproveitar ao máximo os dados disponíveis, dividi o dashboard em quatro partes:

Relação de Idade - Esta seção tem como objetivo apresentar todas as informações relevantes relacionadas à idade do cliente, como a distribuição por faixas etárias, nível de educação, renda, principais produtos adquiridos e canais de compra utilizados.

Relação Familiar - Esta seção fornece informações sobre a situação familiar do cliente, como estado civil e número de filhos, e como esses fatores se relacionam com outras características de cada cliente.

Aceitação das campanhas - A empresa realizou diversas campanhas. Esta seção indica os resultados dessas campanhas, o número de clientes que responderam a elas e os perfis desses clientes.

Distribuição de Compras de Produtos - Mostra a distribuição de compra dos produtos, baseada na renda dos clientes e do grupo, com a possibilidade de filtrar por idade.

