

Análise de Dados

Trabalho final

Professor: Paulo Cotta

Entrega: 20/11/2020

Notas: 40 pts

Alunos:

- João Marcelo
- Matheus Reis
- Matheus Sena
- Ygor Oliveira
- Thiago Costa

O conjunto de dados para este projeto se origina do [repositório de Machine Learning da UCI](https://archive.ics.uci.edu/ml/datasets/Housing) (<https://archive.ics.uci.edu/ml/datasets/Housing>). Os dados de imóveis de Boston foram coletados em 1978 e cada uma das 489 entradas representa dados agregados sobre 14 atributos para imóveis de vários subúrbios de Boston.

Neste projeto, você irá avaliar um conjunto de dados coletado dos imóveis dos subúrbios de Boston, Massachusetts. O principal objetivo deste trabalho é realizar a análise e começar a trabalhar com funções e métodos que serão utilizados no dia a dia de vocês como Engenheiros de Dados e/ou Engenheiros de Machine Learning.

In [1]:

```
1 # Verificação se o sklearn está instalado na sua máquina
2 import sklearn
3 print("A versão do scikit-learn é ", sklearn.__version__)
```

A versão do scikit-learn é 0.23.1

O sklearn é um framework que já possui alguns algoritmos de Machine Learning (ML) prontos. Eu recomento que utilizem sempre a ultima versão do framework.

Documentação: [link \(https://scikit-learn.org/stable/\)](https://scikit-learn.org/stable/).

Mediante ao cenário seguinte:

Os dados de imóveis de Boston foram coletados em 1978 e cada uma das 489 entradas representa dados agregados sobre 14 atributos para imóveis de vários subúrbios de Boston. Para o propósito deste projeto, os passos de pré-processamento a seguir foram feitos para esse conjunto de dados:

- 16 observações de dados possuem um valor 'MEDV' de 50.0. Essas observações provavelmente contêm **valores ausentes ou censurados** e foram removidas.
- 1 observação de dados tem um valor 'RM' de 8.78. Essa observação pode ser considerada **valor atípico (outlier)** e foi removida.
- Os atributos 'RM', 'LSTAT', 'PTRATIO', and 'MEDV' são essenciais. O resto dos **atributos irrelevantes** foram excluídos.

- O atributo 'MEDV' foi **escalonado multiplicativamente** para considerar 35 anos de inflação de mercado.

Fica mais tranquilo efetuar o trabalho conhecendo um pouco sobre o conjunto de dados (dataset).

In [2]:

```
1 # Execute a célula de código abaixo para carregar o conjunto dos dados dos imóv
2 # Importar as bibliotecas necessárias para este projeto
3 import numpy as np
4 import pandas as pd
5 from sklearn.model_selection import ShuffleSplit
6
7 # Formatação mais bonita para os notebooks
8 %matplotlib inline
9
10 # Executar o conjunto de dados de imóveis de Boston
11 data = pd.read_csv('housing.csv')
12 prices = data['MEDV']
13 # Dropando a coluna com maior índice de valores ausentes
14 features = data.drop('MEDV', axis = 1)
15
16 data.info()
17
18 # Êxito
19 print("O conjunto de dados de imóveis de Boston tem {} pontos com {} variáveis
```

```
<class 'pandas.core.frame.DataFrame'>
```

```
RangeIndex: 489 entries, 0 to 488
```

```
Data columns (total 4 columns):
```

```
#   Column   Non-Null Count  Dtype
---  -
0    RM      489 non-null      float64
1   LSTAT    489 non-null      float64
2  PTRATIO  489 non-null      float64
3   MEDV    489 non-null      float64
```

```
dtypes: float64(4)
```

```
memory usage: 15.4 KB
```

```
O conjunto de dados de imóveis de Boston tem 489 pontos com 4 variáveis em cada.
```

In [3]:

```
1 data.head()
```

Out[3]:

	RM	LSTAT	PTRATIO	MEDV
0	6.575	4.98	15.3	504000.0
1	6.421	9.14	17.8	453600.0
2	7.185	4.03	17.8	728700.0
3	6.998	2.94	18.7	701400.0
4	7.147	5.33	18.7	760200.0

Explorando os dados

Você aluno deve efetuar uma investigação sobre os dados de imóveis de Boston e fornecerá suas observações. Familiarizar-se com os dados durante o processo de exploração é uma prática fundamental que ajuda você a entender melhor e justificar seus resultados.

Dado que o objetivo principal deste projeto é construir um modelo de trabalho que tem a capacidade de estimar valores dos imóveis, vamos precisar separar os conjuntos de dados em **atributos** e **variável alvo**. O **atributos**, 'RM' , 'LSTAT' e 'PTRATIO' , nos dão informações quantitativas sobre cada ponto de dado. A **variável alvo**, 'MEDV' , será a variável que procuramos estimar. Eles são armazenados em `features` e `prices` , respectivamente.

In [4]:

```
1  # TODO: Preço mínimo dos dados
2  minimum_price = np.amin(prices)
3
4  # TODO: Preço máximo dos dados
5  maximum_price = np.amax(prices)
6
7  # TODO: Preço médio dos dados
8  mean_price = np.mean(prices)
9
10 # TODO: Preço mediano dos dados
11 median_price = np.median(prices)
12
13 # TODO: Desvio padrão do preço dos dados
14 std_price = np.std(prices)
15
16 # Mostrar as estatísticas calculadas
17 print("Estatísticas para os dados dos imóveis de Boston:\n")
18 print("Preço mínimo: ${:,.2f}".format(minimum_price))
19 print("Preço máximo: ${:,.2f}".format(maximum_price))
20 print("Preço médio: ${:,.2f}".format(mean_price))
21 print("Preço mediano: ${:,.2f}".format(median_price))
22 print("Desvio padrão dos preços: ${:,.2f}".format(std_price))
```

Estatísticas para os dados dos imóveis de Boston:

Preço mínimo: \$105,000.00
Preço máximo: \$1,024,800.00
Preço médio: \$454,342.94
Preço mediano: \$438,900.00
Desvio padrão dos preços: \$165,171.13

In [5]:

```
1 # Trouxe os mesmo resultados apresentados e inclusive os Quartis
2 data.describe()
```

Out[5]:

	RM	LSTAT	PTRATIO	MEDV
count	489.000000	489.000000	489.000000	4.890000e+02
mean	6.240288	12.939632	18.516564	4.543429e+05
std	0.643650	7.081990	2.111268	1.653403e+05
min	3.561000	1.980000	12.600000	1.050000e+05
25%	5.880000	7.370000	17.400000	3.507000e+05
50%	6.185000	11.690000	19.100000	4.389000e+05
75%	6.575000	17.120000	20.200000	5.187000e+05
max	8.398000	37.970000	22.000000	1.024800e+06

Questão 1

Após apresentação dos dados, em formato não supervisionado, desenvolva utilizando K-means e ou GMM um classificador e apresente em formato de Data Visualization os clusters.

Imports

In [21]:

```
1 # imports
2 import numpy as np
3 import pandas as pd
4 import matplotlib.pyplot as plt
5 from sklearn.cluster import KMeans
6 from sklearn.model_selection import train_test_split
7 from sklearn.metrics import accuracy_score, classification_report, confusion_ma
8
9 %matplotlib inline
```

Dataset

In [7]:

```
1 data.head()
```

Out[7]:

	RM	LSTAT	PTRATIO	MEDV
0	6.575	4.98	15.3	504000.0
1	6.421	9.14	17.8	453600.0
2	7.185	4.03	17.8	728700.0
3	6.998	2.94	18.7	701400.0
4	7.147	5.33	18.7	760200.0

In [8]:

```
1 data.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 489 entries, 0 to 488
Data columns (total 4 columns):
#   Column      Non-Null Count  Dtype
---  -
0    RM          489 non-null    float64
1    LSTAT        489 non-null    float64
2    PTRATIO      489 non-null    float64
3    MEDV         489 non-null    float64
dtypes: float64(4)
memory usage: 15.4 KB
```

Colunas

<https://www.cs.upc.edu/~belanche/Docencia/mineria/Practiques/Boston.dat>
[\(https://www.cs.upc.edu/~belanche/Docencia/mineria/Practiques/Boston.dat\)](https://www.cs.upc.edu/~belanche/Docencia/mineria/Practiques/Boston.dat)

- RM : average number of rooms per dwelling
- LSTAT : % lower status of the population
- PTRATIO : pupil-teacher ratio by town
- MEDV : Median value of owner-occupied homes in \$1000's

Train Test Split

In [9]:

```
1 x_train, x_test = train_test_split(data, test_size=0.3, random_state=101)
```

Modelo

In [10]:

```
1 kmeans = KMeans(4)
2 kmeans.fit(x_train)
```

Out[10]:

KMeans(n_clusters=4)

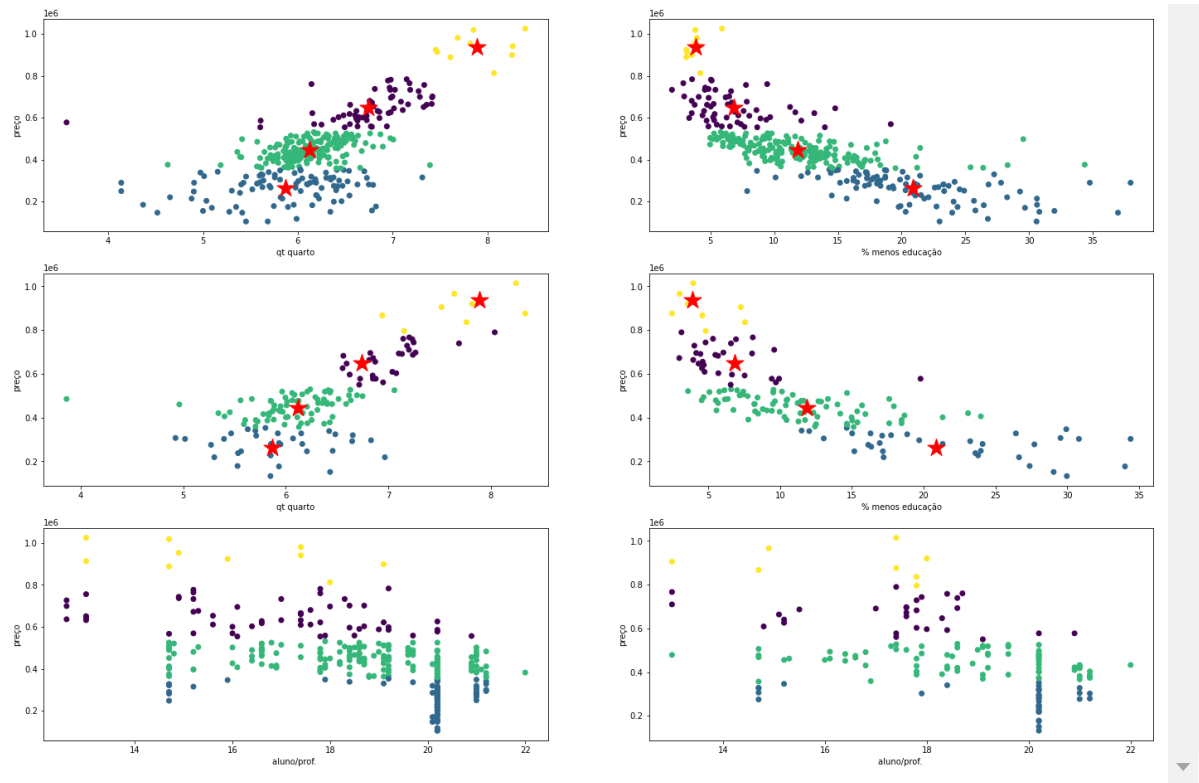
In [11]:

```
1 pred_train = kmeans.predict(x_train)
2 pred_test = kmeans.predict(x_test)
```

Visualização do cluster

In [42]:

```
1 fig, axs = plt.subplots(3, 2, figsize=(24, 16))
2
3 # quantidade de quartos X preço
4 axs[0][0].scatter(x_train['RM'], x_train['MEDV'], c=pred_train)
5 axs[0][0].set_xlabel('qt quarto')
6 axs[0][0].set_ylabel('preço')
7 axs[1][0].scatter(x_test['RM'], x_test['MEDV'], c=pred_test)
8 axs[1][0].set_xlabel('qt quarto')
9 axs[1][0].set_ylabel('preço')
10
11 # educação X preço
12 axs[0][1].scatter(x_train['LSTAT'], x_train['MEDV'], c=pred_train)
13 axs[0][1].set_xlabel('% menos educação')
14 axs[0][1].set_ylabel('preço')
15
16 axs[1][1].scatter(x_test['LSTAT'], x_test['MEDV'], c=pred_test)
17 axs[1][1].set_xlabel('% menos educação')
18 axs[1][1].set_ylabel('preço')
19
20 # alunos/educação X preço
21 axs[2][0].scatter(x_train['PTRATIO'], x_train['MEDV'], c=pred_train)
22 axs[2][0].set_xlabel('aluno/prof. ')
23 axs[2][0].set_ylabel('preço')
24
25 axs[2][1].scatter(x_test['PTRATIO'], x_test['MEDV'], c=pred_test)
26 axs[2][1].set_xlabel('aluno/prof. ')
27 axs[2][1].set_ylabel('preço')
28
29 # centróides
30 for centroid in kmeans.cluster_centers_:
31     axs[0][0].scatter(centroid[0], centroid[3], marker='*', color='red', s=500)
32     axs[1][0].scatter(centroid[0], centroid[3], marker='*', color='red', s=500)
33     axs[0][1].scatter(centroid[1], centroid[3], marker='*', color='red', s=500)
34     axs[1][1].scatter(centroid[1], centroid[3], marker='*', color='red', s=500)
35
36 plt.show()
```



Conclusão

- Casas com maior quantidade de quartos, tendem a ser mais caras.
- Casas tender a ser mais caras em bairros onde há o maior número de pessoas alfabetizadas.
- A relação de aluno/professor não afeta o preço das casas

In []:

```
1
```