

House Pricing

Hirley Dayan Lourenço da Silva e Marcia Maria Parmigiani Martins

Leitura dos Dados

O dataset utilizado nesse trabalho refere-se a dados de imóveis como total de cômodos, idade, etc, o qual foi dividido entre treino, validação e teste, com o objetivo de criação de um modelo utilizando algoritmo de **Regressão Linear** para predição do preço de imóveis.

O dataset de treino possui **12384** observações e **10** features.

O dataset de validação possui **4128** observações e **10** features.

O dataset de teste possui **4128** observações e **10** features.

Tratamento dos dados

O dataset possui uma feature categórica **ocean_proximity** com 5 níveis. Considerando que a Regressão Linear assume que todas as variáveis independentes são numéricas, iremos utilizar a técnica de **hot encoding** para transformar a feature em numérica, atribuindo valor 1 se o caso se enquadre na determinada categoria. A inclusão da feature categórica possibilitou um resultado melhor para todos os modelos testados.

Além disso, foram removidas observações de features sem anotações (NA) nos dados de treino, validação e teste.

Normalização dos dados

A técnica de normalização aplicada ao dataset é a **Min-Max** onde os dados são dimensionados em um intervalo fixo, normalmente de 0 a 1. A feature target **median_house_value** não foi incluída na normalização.

longitude	latitude	housing_median_age	total_rooms	total_bedrooms	population	households	median_income
0.61	0.16	1.00	0.05	0.07	0.05	0.07	0.13
0.60	0.15	0.88	0.03	0.05	0.02	0.05	0.18
0.20	0.58	0.63	0.00	0.00	0.00	0.00	0.25
0.53	0.30	0.67	0.00	0.00	0.01	0.01	0.10
0.57	0.17	0.63	0.14	0.12	0.06	0.12	0.43

Regressão Linear

O algoritmo de **Regressão Linear** foi utilizado para prever os preços dos imóveis onde a variável que se deseja encontrar é a **median_house_value** que representa o valor do imóvel baseado em suas features.

As medidas de avaliação utilizadas foram a **MAE (Mean Absolute Error)** que calcula a média da diferença absoluta entre os valores preditos e os observados e o Coeficiente de Determinação (R2) que indica quão bem o modelo consegue se ajustar sobre um conjunto de predições e seus valores verdadeiros, podendo variar entre 0 e 1, sendo que 0 indica que o modelo não consegue explicar a variabilidade dos dados e 1 indica que as predições se ajustam perfeitamente aos dados (modelo explica toda a variabilidade das predições) o que indica **overfitting**.

Resultado dos modelos

Avaliação para o modelo **baseline**: MAE **50470.57** e R² **63.43%**.

Avaliação para o modelo **complexo baseado na combinação de features existentes**: MAE **50999.99** e R² **63.43%**.

Avaliação para os modelos **complexos baseado em regressão polinomial**:

Fórmula grau 1: MAE = **50470.57** e R² = **63.43%**

Fórmula grau 2: MAE = **50527.82** e R² = **63.91%**

Fórmula grau 3: MAE = **48052.91** e R² = **66.26%**

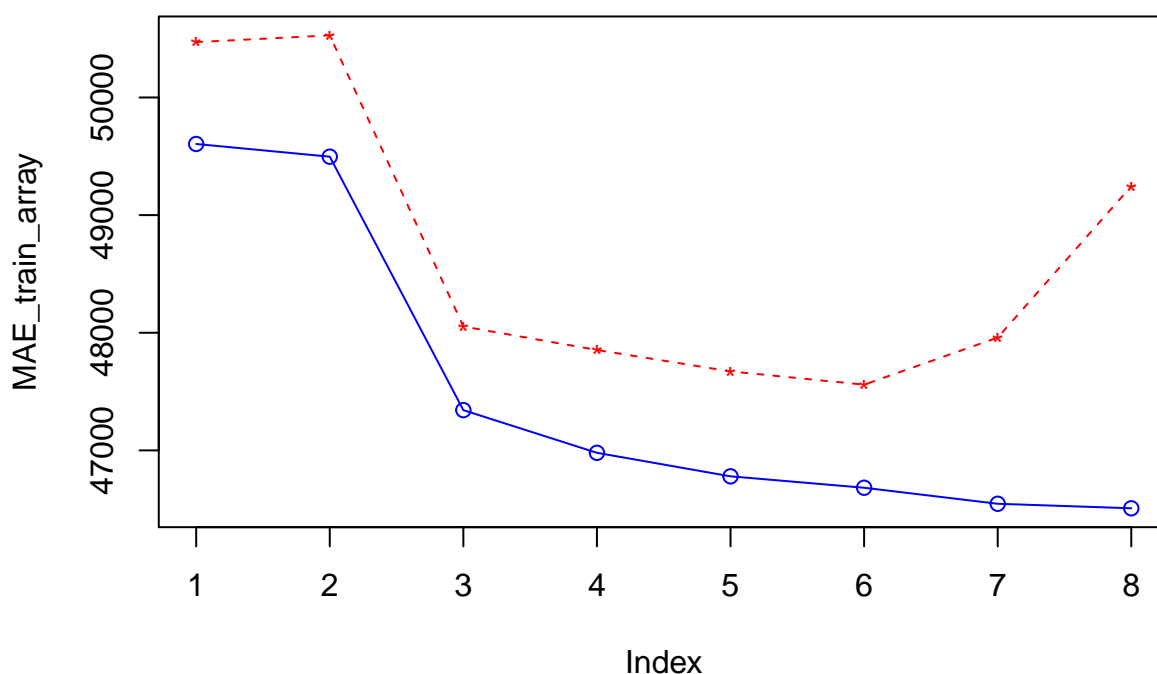
Fórmula grau 4: MAE = **47853.44** e R² = **66.81%**

Fórmula grau 5: MAE = **47671.35** e R² = **67%**

Fórmula grau 6: MAE = **47558.52** e R² = **67.21%**

Fórmula grau 7: MAE = **47958.79** e R² = **61.92%**

Fórmula grau 8: MAE = **49245.76** e R² = **-19.87%**



A partir da função de grau 7 percebe-se que os modelos tem um bom desempenho nos dados de treinamento, porém não generaliza, tornando-se muito especializado, o que caracteriza **overfitting**.

Resultado do melhor modelo no conjunto de teste

O modelo escolhido para uso no conjunto de teste é o polinomial de grau 6, pois na curva do MAE foi o que apresentou menor diferença entre o valor real e o predito durante o treinamento. A escolha da função foi feita de modo a garantir uma convergência entre o treinamento e a validação.

Avaliação para o modelo **polinomial de grau 6**: MAE **46984.26** e R² **68.44%**.