<pre>import matplotlib.image as mpimg img = mpimg.imread('l.jpg') plt.imshow(img) plt.axis('off') plt.show()</pre>
Qu'est-ce que le dataset California Housing ? # Le dataset California Housing est un ensemble de données largement utilisé dans l'apprentissage automatique pour s'entraîner à résoudre des problèmes de régression. Il contient des informations sur les propriétés immobilières en Californie, telles que :#• Données démog population, revenu médian, etc# • Caractéristiques des propriétés: nombre de pièces, âge moyen des maisons, et# • Localisation géographique: latitude, longitu#. • Variable cible: prix médian d maisonsn.
Quel problème d'apprentissage automatique résout-on avec ce dataset ? L'objectif principal est de prédire le prix médian d'une maison en Californie en fonction des autres caractéristiques disponibles dans le dataset. C'est un problème typique de régression. # import librairies and modules import pandas as pd import numpy as np import seaborn as sns import matplotlib.pyplot as plt
<pre>import matplotlin.pyplot as pit from sklearn.model_selection import train_test_split from sklearn.preprocessing import StandardScaler from sklearn.linear_model import LinearRegression from sklearn.metrics import mean_squared_error,mean_absolute_error,r2_score %matplotlib inline # Load the dataset and create a dataframe from this # import the dataset house_data = pd.read_csv('housing.csv')</pre>
Nouse_data Nouse_data Nousing_median_age No
4 -122.25 37.85 52.0 1627.0 280.0 565.0 259.0 3.8462 342200.0 NEAR BAY
20639 -121.24 39.37 16.0 2785.0 616.0 1387.0 530.0 2.3886 89400.0 INLAND 20640 rows × 10 columns Exploring the dataset
Provides a summury of dataframe print (house_data.info()) <class 'pandas.core.frame.dataframe'=""> RangeIndex: 20640 entries, 0 to 20639 Data columns (total 10 columns): # Column Non-Null Count Dtype</class>
1 1 1 1 1 1 1 1 1 1
dtypes: float64 (9), object (1) memory usage: 1.6+ MB None Generates descriptive statistiques for numerical column house_data.describe()
count latitude housing_median_age total_podrous households median_house_value count 20640.00000 20640.00000 20640.00000 20640.00000 20640.00000 20640.00000 20640.00000 mean -119.569704 35.631861 28.639486 2635.763081 537.870553 1425.476744 499.53960 3.870671 206855.816909 std 2.003532 2.135952 12.58558 2181.615252 421.385070 1132.462122 382.329753 1.899822 115395.615874 min -124.350000 32.54000 1.00000 2.000000 787.00000 280.00000 0.499900 14999.000000 25% -121.80000 33.93000 18.00000 1447.75000 296.00000 787.00000 280.00000 2.563400 119600.000000
50% -118.490000 34.260000 29.000000 2127.000000 409.000000 3.534800 179700.000000 75% -118.010000 37.710000 37.00000 3148.00000 647.00000 1725.00000 605.00000 4.743250 264725.000000 max -114.310000 41.950000 52.00000 39320.000000 6445.000000 35682.000000 6082.000000 15.000100 500001.000000 Displays the first few rows of dataframe
3 -122.25 37.85 52.0 1274.0 235.0 558.0 219.0 5.6431 341300.0 NEAR BAY 4 -122.25 37.85 52.0 1627.0 280.0 565.0 259.0 3.8462 342200.0 NEAR BAY Checking missing value in a dataset¶ house_data.isnull().sum()
longitude 0 latitude 0 housing_median_age 0 total_rooms 0 total_bedrooms 207 population 0 households 0 median_income 0 median_house_value 0 ocean_proximity 0
Data cleaning (handling missing values and irrevelant columns) ### Filling missing values in the total_bedrooms column(Median is the mid point value) house_data['total_bedrooms'] = house_data['total_bedrooms'].fillna(house_data['total_bedrooms'].mean()) Mesuring the statistics relationship between a housing data variables
Mesuring the statistics relationship betwen a housing_data variables print(house_data.isnull().sum()) longitude
households 0 median_income 0 median_house_value 0 ocean_proximity 0 dtype: int64 Find outliers
<pre># To calculate the quartiles and IQR Q1 = house_data['median_house_value'].quantile(0.25) Q3 = house_data['median_house_value'].quantile(0.75) IQR = Q3 - Q1 # Define the outliers limits lower_bound = Q1 - 1.5 * IQR upper_bound = Q3 + 1.5 * IQR # identify the outliers values outliers = house_data[(house_data['median_house_value'] < lower_bound) (house_data['median_house_value'] > upper_bound)] print("valeurs aberrantes détéctées: ") print(outliers)</pre>
valeurs aberrantes détéctées: longitude latitude housing_median_age total_rooms \ 89
20426 -118.69
459 NEAR BAY 489 NEAR BAY 493 NEAR BAY 20422 <1H OCEAN 20426 <1H OCEAN 20427 <1H OCEAN 20436 <1H OCEAN 20433 <1H OCEAN
<pre>[1071 rows x 10 columns] # to delete the ouliers value house_data_cleaned = house_data[(house_data['median_house_value'] >= lower_bound) & (house_data['median_house_value'] <= upper_bound)] print(" données sans valeurs aberrantes : ") print(house_data_cleaned) données sans valeurs aberrantes : longitude latitude housing_median_age total_rooms total_bedrooms \</pre>
0 -122.23
population households median_income definition and defi
20637 1007.0 433.0 1.7000 92300.0 20638 741.0 349.0 1.8672 84700.0 20639 1387.0 530.0 2.3886 89400.0
20635 INLAND 20636 INLAND 20637 INLAND 20638 INLAND 20639 INLAND [19569 rows x 10 columns] Dropping a ocean_proximity column
house_data_cleaned =house_data_cleaned.drop(columns=['ocean_proximity']) house_data_cleaned.corr() longitude latitude housing_median_age total_rooms total_bedrooms population households median_income median_house_value longitude 1.00000 -0.924031 -0.101691 0.044735 0.070191 0.101638 0.056346 -0.010808 -0.047342
latitude -0.924031 1.00000 0.006104 -0.033798 -0.068286 -0.114359 -0.073547 -0.076249 -0.149100 housing_median_age -0.101691 0.006104 1.000000 -0.372715 -0.326596 -0.294893 -0.310343 -0.197628 0.061480 total_rooms 0.044735 -0.033798 -0.372715 1.000000 0.931110 0.859642 0.921278 0.226557 0.147526 total_bedrooms 0.070191 -0.068286 -0.326596 0.931110 1.000000 0.875347 0.973927 0.024602 0.078245 population 0.101638 -0.114359 -0.294893 0.859642 0.875347 1.000000 0.999218 0.045413 0.017764 households 0.056346 -0.073547 -0.310343 0.973927 0.999218 1.000000 0.049055 0.099352
median_income
Nouse_data_cleaned['median_house_value'] Nouse_data_cleaned.head(3)
<pre>X_train, X_test, y_train, y_test= train_test_split(X,y, test_size=0.2, random_state = 42) print(X_train.shape, X_test.shape) (15655, 8) (3914, 8) Standardize the features</pre>
Scaler = StandardScaler() X_train_scaled = Scaler.fit_transform(X_train) X_test_scaled = Scaler.transform(X_test) initialise and train the regression model
model = LinearRegression() model.fit(X_train_scaled, y_train) ▼ LinearRegression LinearRegression() Making predictions
y_pred= model.predict (X_test_scaled) y_pred array([74747.38050415, 128726.2701118 , 261189.13900352,, 169719.22186104, 228841.10418274, 133428.93214741]) Evaluate the model: using Mean Squared Error (MSE), Mean Absolute Error (MAE) and R2 (determination coefficient)
<pre>MAE = mean_absolute_error(y_test,y_pred) MSE = mean_squared_error(y_test,y_pred) r2 = r2_score(y_test,y_pred) print("Mean Absolute Error:", MAE) print("Mean Squared Error:", MSE) print("R- Squared:", r2)</pre> Mean Absolute Error: 44585.379349951014
Mean Squared Error: 3585071186.9282236 R- Squared: 0.6123494081149208 predicting a small sample of data Sample_data = X_test.head() Sample_data_scaled = Scaler.transform(Sample_data) predictions = model.predict(Sample_data_scaled)
print(" Prédictions: ", predictions) print("Actual values: ", y_test.head().values) Prédictions: [74747.38050415 128726.2701118 261189.13900352 171205.84624638 199909.7129558] Actual values: [55200. 97800. 245300. 162100. 163100.] Create a scatter plot
<pre>plt.figure(figsize= (12,8)) sns.scatterplot(x=y_test, y=y_pred, alpha = 0.5) plt.plot([min(y_test), max(y_test)], [min(y_test), max(y_test)], color ='red', linestyle='') plt.xlabel('Actual House Prices') plt.ylabel('Predicted House Prices') plt.title("Actual Vs Predicted House Prices ") plt.show()</pre> Actual Vs Predicted House Prices
Actual Vs Predicted House Prices 400000 -
200000 - O - O - O - O - O - O - O - O -
T
-400000 - -600000 - 0 100000 200000 300000 400000 500000 Actual House Prices
Creating a bar plot sample_size = 20 indices = np.arange(sample_size) actual_prices = y_test[:sample_size].values predicted_prices = y_pred[:sample_size]
<pre>predicted_prices = y_pred[:sample_size] Sample_df = pd.DataFrame({ 'Actual' :actual_prices, 'predicted' : predicted_prices }) plt.figure(figsize=(14,7)) Sample_df.plot (kind='bar', figsize=(14,7)) plt.xlabel('Sample Index') plt.ylabel('House price') plt.title('Actual vs Predicted house for sample data') plt.show()</pre>
Figure size 1400x700 with 0 Axes> Actual vs Predicted house for sample data 400000 -
300000 - 300000 -
100000 -
Sample Index
Save the housing model import joblib joblib.dump(model,'california_house_prices_scaler.pkl') model_scaler = {

Out[152]: ['model_scaler.pkl']