Activity #4

Linear Regression



4.1 Data Exploration / Transform / Feature Selection

4.2 PCA: Feature Dimensional Reduction

4.3 Linear Regression

TOPICS

LIBRARIES

import numpy as np import pandas as pd import matplotlib.pyplot as plt • import seaborn as sns from sklearn import preprocessing from sklearn import preprocessing from sklearn.decomposition import PCA from sklearn.model selection import train test split from sklearn.linear_model import LinearRegression from sklearn.metrics import r2_score



4.1

Data Exploration / Transform /
Feature Selection

4.1 (a) Data exploration

Read .csv file

read_csv("CarPrice.csv")

- View Data Array Shape
 - # Variables
 - # Samples
 - # Statistics -> describe()
- Remove
 - 'car_ID','CarName'
- View Variable info
 - Info()
 - Data Type / # non null
- Fill NA
 - fillna()

4.1 (b) Data Transform and Feature Selection

• Standardized Data for continuous data columns for only continuous data columns

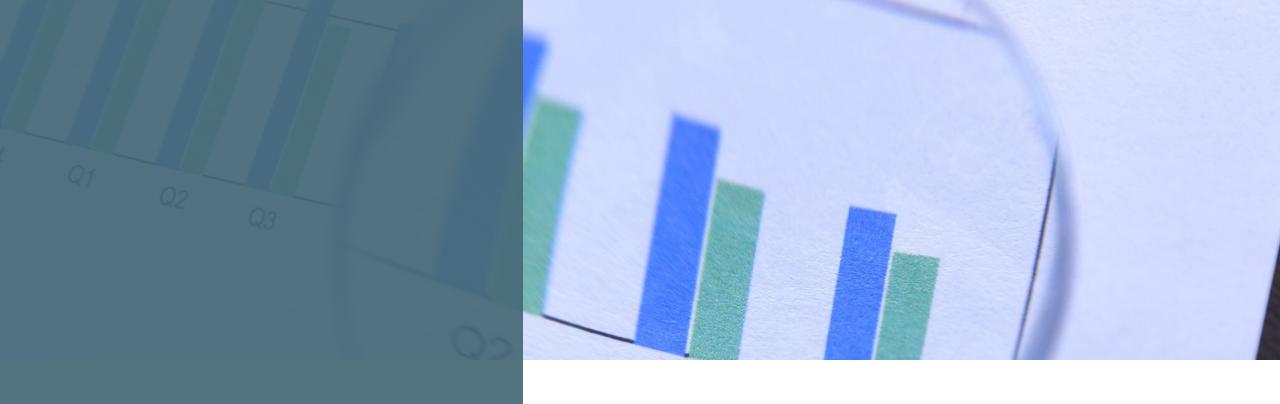
Calculate correlation between variables for only continuous data columns

• corr()

• Reduce Corr() to Lower Matrix

• Drop columns if correlation value > 0.86

- OneHotEncode for categorical columns (try from Pandas)
 - pd.get dummies(data, columns = categorical features, drop first=True)



4.2

PCA: Feature Dimensional Reduction

4.2 PCA Dimensional Reduction

• # PCA n_components (ทดลองเปลี่ยนค่า n_components อย่างน้อย 3 ค่า เพื่อเลือกค่าดีที่สุด)

- # PCA all variables
- pca = PCA()
- X_pca = pca.fit_transform(X_standard)

2

- Visualize Explained Variance Ratio (% eigenvalues)
 - plt.bar() ค่าของ pca.explained_variance_ratio_

• pca2 = PCA(n_components)

• X pca 2 = pca2.fit transform(X standard)



4.3

Linear Regression

4.3 Linear Regression

- # Shuffle Split (Train / Test Split)
- Rseed กำหนด ค่าใดก็ได้
- x_train_set, x_test, y_train_set, y_test = train_test_split(X, Y, test_size = 0.3, random_state = Rseed)
- # Shuffle Split (Train / Validation Split)
- x_train, x_validate, y_train, y_validate = train_test_split(x_train_set, y_train_set, test_size = 0.3, random_state = Rseed)
- # Perform Linear Regression -> All variables
- Ir = LinearRegression()
- # Train
- lr.fit(x_train, y_train)
- # Validate
- y pred lr = lr.predict(x validate)
- # Test
- y_test_pred_lr = lr.predict(x_test)
- # Measure Accuracy Validation and Test
- r2_score(y_pred_lr, y_validate)
- r2_score(y_test_pred_lr, y_test)
- lr.score(x_validate, y_validate)
- lr.score(x_test, y_test)

เปรียบเทียบประสิทธิภาพ Linear Regression

