



Support Vector Regression (SVR)

#### Activity #5

5.1 5.2 5.3

Data Preparation Model Training and Cross Validation Model Evaluation

5.1 Data Exploration / Transform / Feature Selection / Time Series windows

**5.2 Model Training and Cross Validation** 

**5.3 Model Evaluation** 

**TOPICS** 

#### **LIBRARIES**

import numpy as np import pandas as pd import matplotlib.pyplot as plt import seaborn as sns • from sklearn.linear model import LinearRegression from sklearn.svm import SVR • from sklearn import preprocessing from sklearn import metrics import pandas datareader.data as web



## 5.1

### **Data Preparation**

# 5.1 (a) Get Data

- 1
- #Read stock data use pandas\_datareader.data from web
- # Get Stock Data
- stk tickers = ['MSFT', 'IBM', 'GOOGL']
- ccy\_tickers = ['DEXJPUS', 'DEXUSUK']
- idx\_tickers = ['SP500', 'DJIA', 'VIXCLS']
- stk\_data = web.DataReader(stk\_tickers, 'yahoo')
- ccy\_data = web.DataReader(ccy\_tickers, 'fred')
- idx\_data = web.DataReader(idx\_tickers, 'fred')

- Select columns
  - Base = stk\_data.loc[:, ('Adj Close', 'MSFT')]
  - X1 = stk\_data.loc[:, ('Adj Close', ('GOOGL', 'IBM'))]
  - X2 = ccy data
  - X3 = idx\_data

1

- #Standardized data (X1, X2, X3) with kept index (date)
  - X1 = standard\_scaler.fit\_transform(X1.values),index = X1.index, columns=X1.columns)

- Calculate ความแตกต่างของค่า ราคา 'Adj Close', 'MSFT') ย้อนหลัง return\_period วัน
  - Y = base. shift(-return\_period)
  - X4\_3DT = base.diff(3\*return\_period).shift(-3\*return\_period)
  - X4 6DT = base.diff(6\*return period).shift(-6\*return period)
  - X4\_12DT = base.diff(12\*return\_period).shift(-12\*return\_period)
  - X4 = pd.concat([X4\_3DT, X4\_6DT, X4\_12DT], axis=1)
  - X4.columns = ['MSFT 3DT', 'MSFT 6DT', 'MSFT 12DT']
  - X4 = pd.DataFrame(standard\_scaler.fit\_transform(X4.values), index = X4.index,columns=X4.columns)

2

- # Forming Dataset
  - X = pd.concat([X1, X2, X3, X4], axis=1)
  - dataset = pd.concat([Y, X], axis=1)

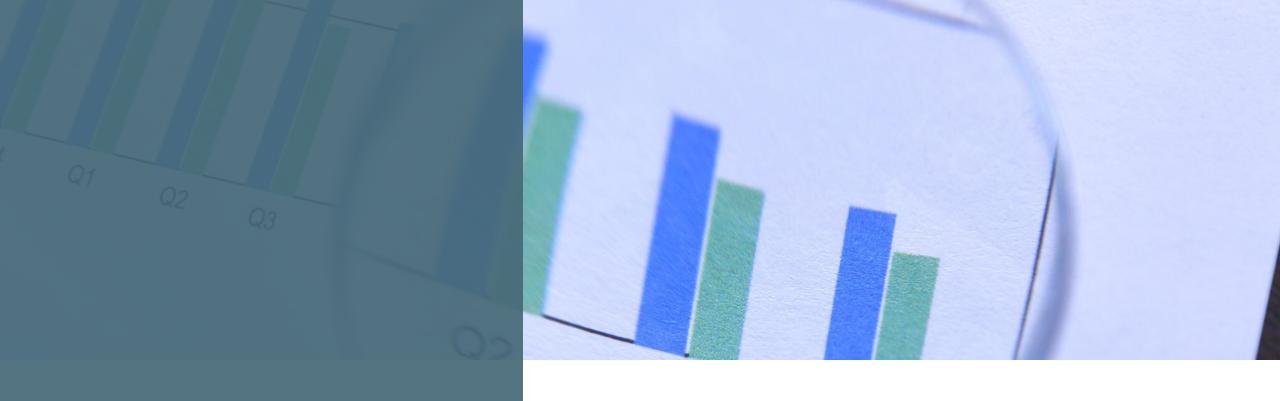
## 5.1 (c) Data Preparation

• # Drop NA

• Dropna()

- # View Statistics
  - Describe()
- # Assign X, Y (drop datetime index)
  - Y = dataset (1<sup>st</sup> column)
  - X = dataset (2<sup>nd</sup> :last column)
- # feature selection (correlation)
  - Calculate correlation between variables for only continuous data columns
    - corr()
  - Reduce Corr() to Lower Matrix
  - Drop columns if | correlation value | > 0.9
- # Train / Test Preparation
- Test size = 0.3 \* len( X )
- Train size = 0.7 \* len(X)
- X\_train, X\_test = X[0:train\_size], X[train\_size:len(X)]
- Y\_train, Y\_test = Y[0:train\_size], Y[train\_size:len(X)]

4



**5.2** 

**Model Training and Cross Validation** 

## 5.2 (a) Model Training and Cross Validation

1

- #Set number of fold / Seed value
  - Num\_fold
  - Seed

2

- # Cross Validation Model
- # set k-fold crossvalidation with shuffle
- kfold = model\_selection.KFold(n\_splits=mum\_fold, shuffle = True, random\_state=seed)

• # Model selection

- Model\_LM = LinearRegression(copy\_X=True, fit\_intercept=True, n\_jobs=1, normalize=False)
- c\_val ลองอย่างน้อย 3 ค่า [0.1, 1, 10, 100]
- svr lin = SVR(kernel='linear', C=c\_val)
- svr\_rbf = SVR(kernel='rbf', C=c\_val, gamma=0.01)
- svr\_poly = SVR(kernel='poly', C=c\_val, degree=2)

## 5.2 (b) Model Training and Cross Validation

- 1
- # Calculate accuracy score for each model
  - score LM = model selection.cross val score(model LM, X train, y train, cv=kfold)
  - score\_lin = model\_selection.cross\_val\_score(svr\_lin, X\_train, y\_train, cv=kfold)
  - score\_rbf = model\_selection.cross\_val\_score(svr\_rbf, X\_train, y\_train, cv=kfold)
  - score\_poly = model\_selection.cross\_val\_score(svr\_poly, X\_train, y\_train, cv=kfold)
- # View score k-fold
  - # Valication score comparison
  - score = pd.DataFrame({'Linear Model':score\_LM,'SVR\_linear':score\_lin, 'SVR\_rbf': score\_rbf, 'SVR\_poly': score\_poly})
  - score\_mean = pd.DataFrame({'AVG Linear Model':[score\_LM.mean()],'AVG SVR\_linear':[score\_lin.mean()], 'AVG SVR\_rbf': [score\_rbf.mean()], 'AVG SVC\_poly': [score\_poly.mean()]})

- 2
- Print(score)
- Print(score\_mean)
- Display( plot ( score ))



5.3

#### **Model Evaluation**

## 5.3 Prediction and Evaluation

1

- # Predict all models (LM, SVR\_linear, SVR\_rbf, SVR Poly)
- LM\_pred = model\_LM.fit(X\_train, y\_train).predict(X\_test)

2

- # Scatter Plot ( X\_test, Predict ) for all model ( LM, SVR\_linear, SVR\_rbf, SVR Poly )
- plt.scatter(X\_test, LM\_pred,c='magenta')
- # Model prediction performance evaluation for all model ( LM, SVR\_linear, SVR\_rbf, SVR Poly )
   MSE
  - LM MSE = metrics.mean squared error(y test, LM pred)
  - R2
    - LM\_r2 = metrics.r2\_score(y\_test, LM\_pred)

ろ

- # Display Prediction MSE, R2 for all models
  - Bar()