# Hands on Machine Learning

## 2019 AI summer program in Asia University

**Yueh-Lin Tsai**
2019 / 07

# About me



- Yueh-Lin Tsai

- Education
  - National Cheng Kung University, M.S., Psychology (2013-2015)
  - National Cheng Kung University, B.S., Psychology (2009-2013)

- Present
  - AI Engineer in Taiwan AI Academy

# What is Artificial Intelligence ?

# Artificial Intelligence

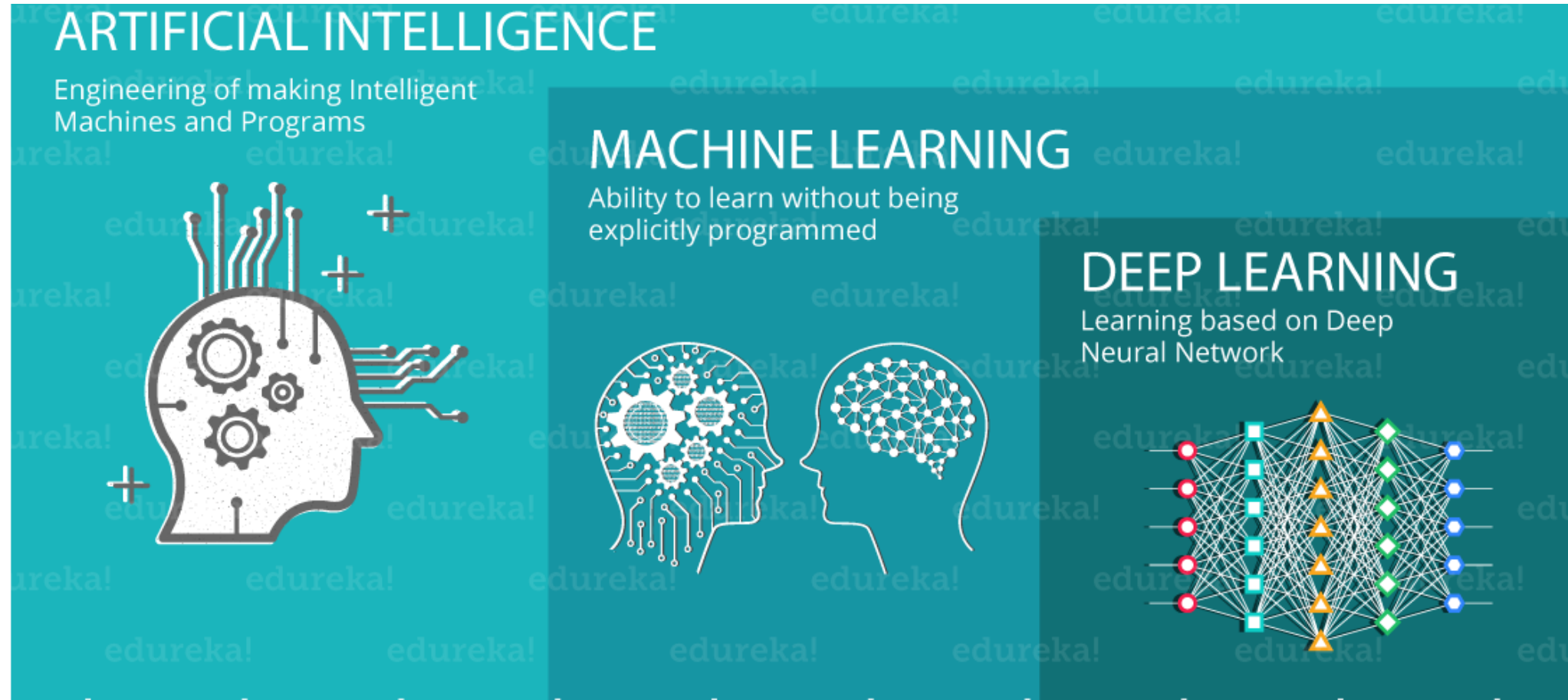- Definition : Intelligence demonstrated by machine

- How ?

Expertise system

Machine learning

Explicit rule from human

Get knowledge from data

# Modern Artificial Intelligence

# Machine learning

- Extract relations/patterns from data automatically

- Apply those rules to unseen data

$$f(x) = y$$

# Type of machine learning

- Supervised learning
  - Regression
  - Classification

**Problems with answer**

- Unsupervised learning
  - Cluster
  - Dimension reduction

**Problems without answer**

- Reinforcement learning

**Problems with fuzzy metric**

# CLASSICAL MACHINE LEARNING

Data is pre-categorized or numerical

Data is not labeled in any way

## SUPERVISED

## UNSUPERVISED

Predict a category

Predict a number

Divide by similarity

Identify sequences

### CLASSIFICATION

«Divide the socks by color»

### CLUSTERING

«Split up similar clothing into stacks»

Find hidden dependencies

### ASSOCIATION

«Find what clothes I often wear together»

### REGRESSION

«Divide the ties by length»

### DIMENSION REDUCTION (generalization)

«Make the best outfits from the given clothes»

# Machine learning workflow

1. Problem definition

2. Data collection

3. Data exploration / preprocessing
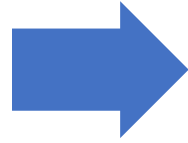
4. Build model

5. Model evaluation

Cheat sheet – scikit learn

# Time for practice

Familiar with colab and python

# Machine Learning with Python

**Collect Data**

- BeautifulSoup
- Lxml
- Requests
- Pandas

**Preprocessing and EDA**

- Numpy
- Pandas
- Scikit-learn
- Matplotlib
- NLTK
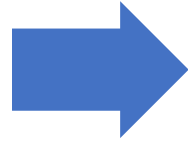
**Analysis and Modeling**

- Statsmodels
- Scikit-learn
- Tensorflow
- Keras
- Pytorch

# Machine Learning with Python

**Collect Data**
- BeautifulSoup
- Lxml
- Requests
- Pandas

**Preprocessing and EDA**
- Numpy
- **Pandas**
- **Scikit-learn**
- Matplotlib
- NLTK

**Analysis and Modeling**
- Statsmodels
- **Scikit-learn**
- Tensorflow
- Keras
- Pytorch

# Build up your ML model in one slide

```python
import pandas as pd
from sklearn import preprocessing, linear_model, model_selection, metrics

data = pd.read_csv('example_data.csv')

data_y = data['target']
data = data.drop('target', axis = 1, inplace = True)

one_hot_data = pd.get_dummies(data)

ss = preprocessing.StandardScaler()
scale_data = ss.fit_transform(data)

train_x, test_x, train_y, test_y = model_selection.train_test_split(data, data_y, test_size = 0.2, random_state = 99)

model = linear_model.LinearRegression() # LogisticRegression()
model.fit(train_x, train_y)

test_prediction = model.predict(test_x)
print('r-square of linear regression : {:.3f}'.format(metrics.r2_score(test_prediction, test_y)))
```

# Exploration and preprocessing

# Data exploration

- Get to know your dataset

    - How many data do I have ?

    - Which column/feature do I have ?

    - Statistics and relations between columns ?

    - Outlier or missing data ?

# Data preprocessing

- Handle missing data

  - Delete data which have missing values (row or column)

  - Missing imputation

- Handle outliers

  - Distribution transformation

  - Replace outliers

# Data preprocessing

- Convert categorical data to numerical data

  - Label encoding

  - One-hot encoding

| Name | Score |
|------|-------|
| Amy | 78 |
| Bob | 90 |
| Chris | 65 |
| Amy | 86 |
| Chris | 67 |

| Name_label | Score |
|------------|-------|
| 1 | 78 |
| 2 | 90 |
| 3 | 65 |
| 1 | 86 |
| 3 | 67 |

| Amy_oh | Bob_oh | Chris_oh | Score |
|--------|--------|----------|-------|
| 1 | 0 | 0 | 78 |
| 0 | 1 | 0 | 90 |
| 0 | 0 | 1 | 65 |
| 1 | 0 | 0 | 86 |
| 0 | 0 | 1 | 67 |

**Label encoding**

**One-hot encoding**

# Data preprocessing

- Normalize data

  - Standard scale

  - Min-max scale

$$x_{standard} = \frac{x - \mu}{\sigma}$$

$$x_{minmax} = \frac{x - Min(x)}{Max(x) - Min(x)}$$

| Name | Score |
|------|------:|
| Amy | 78 |
| Bob | 90 |
| Chris | 65 |
| Amy | 86 |
| Chris | 67 |

| Name | Score |
|------|------:|
| Amy | 0.0719 |
| Bob | 1.1509 |
| Chris | -1.0969 |
| Amy | 0.7912 |
| Chris | -0.9171 |

**Standard scale**

| Name | Score |
|------|------:|
| Amy | 0.48 |
| Bob | 0 |
| Chris | 1 |
| Amy | 0.16 |
| Chris | 0.92 |

**Min-max scale**

# Data preprocessing

- Data splitting
  - Training set
  - Validation set
  - Testing set

- Cross validation

# Time for practice

Data exploration and preprocessing

# Build model

# Model selection

- Model comparison

  - Linear model

    - Focus on global information

    - Have data hypothesis

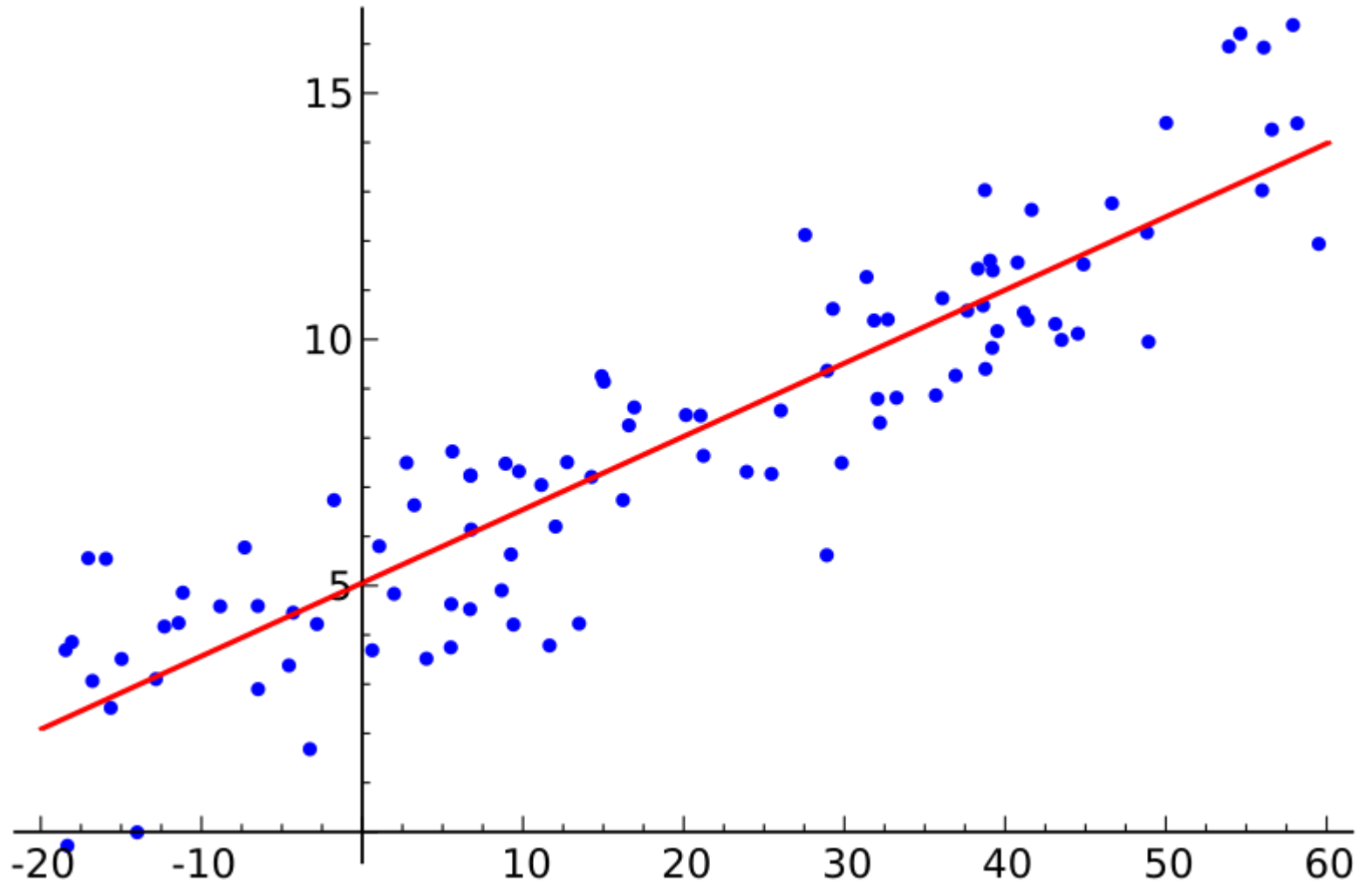    **Encoding matters**

    **Normalization is needed**

  - Tree-based model

    - Clear rules provided by model

    - Focus on local information
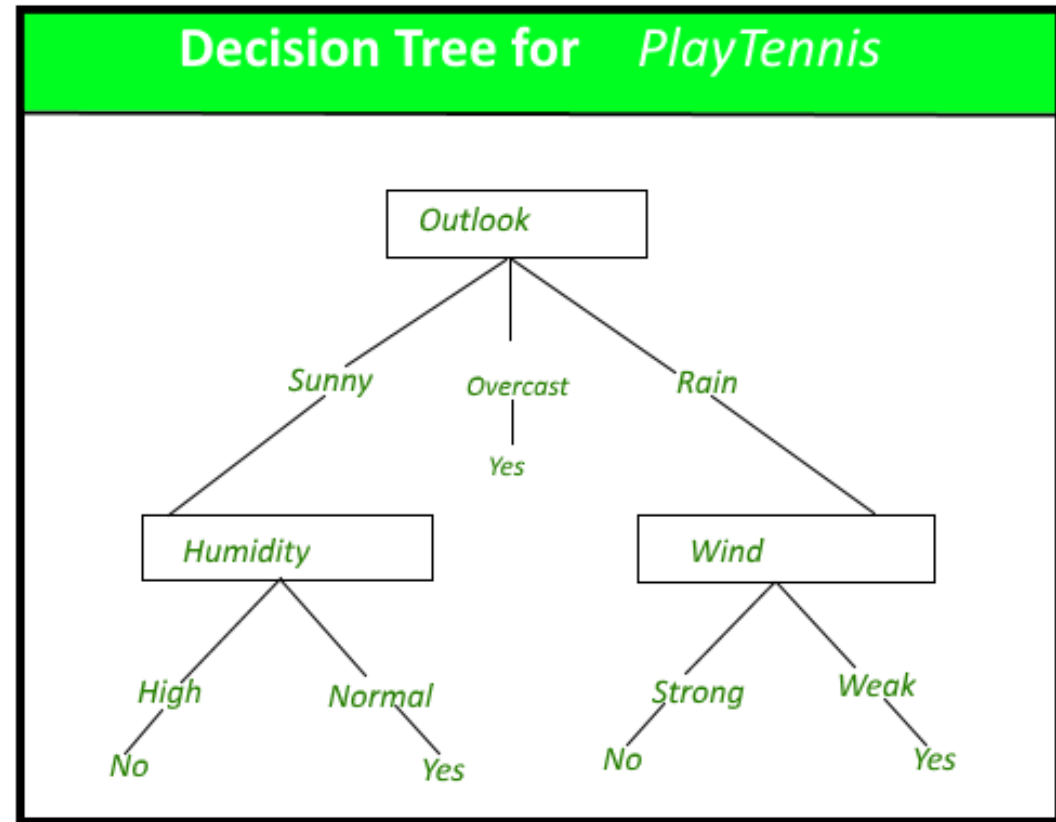
    **No need to normalize data**

# Model selection

- Linear model

  - Linear regression

  - Logistic regression

# Model selection
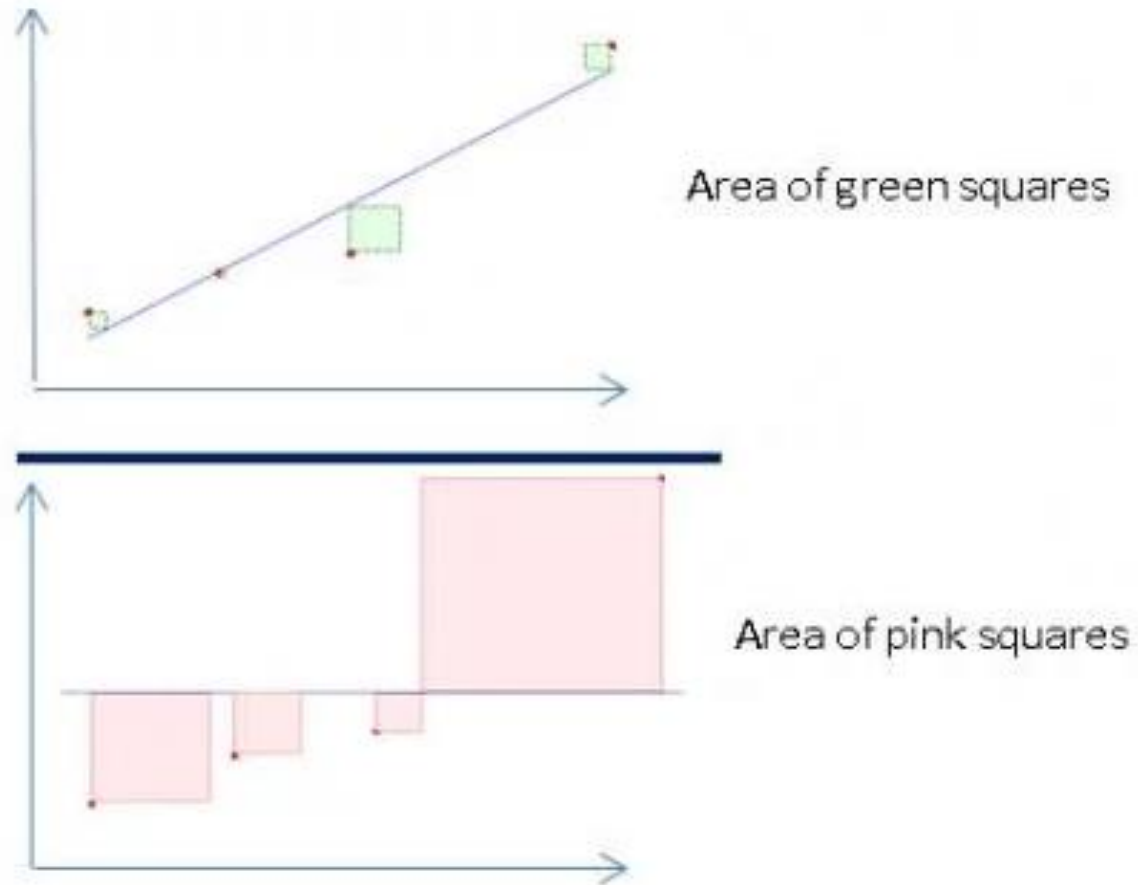
- Tree based model

  - Decision tree



**Decision Tree for** *PlayTennis*

Outlook

Sunny    Overcast    Rain

Yes

Humidity              Wind

High    Normal        Strong    Weak

No        Yes          No        Yes

# Model evaluation

# Model evaluation

- Regression problem

  - R – square, $R^2$

  - Mean Squared Error, MSE

  - Mean Absolute Error, MAE

# Model evaluation



$$R^2 = 1 - \frac{\text{Area of green squares}}{\text{Area of pink squares}}$$

# Model evaluation

$$MSE = \frac{1}{n} \Sigma \left( y - \widehat{y} \right)^2$$

The square of the difference between actual and predicted

Divide by the total number of data points

Predicted output value

Actual output value

$$MAE = \frac{1}{n} \Sigma \left| y - \widehat{y} \right|$$

Sum of

The absolute value of the residual

# Model evaluation

• Classification problem

  • Confusion matrix

  • Accuracy

  • Precision, recall

  • Area under curve (AUC), f1 score

# Model evaluation



$$\text{Accuracy} = \frac{TP + TN}{TP + FP + FN + TN}$$
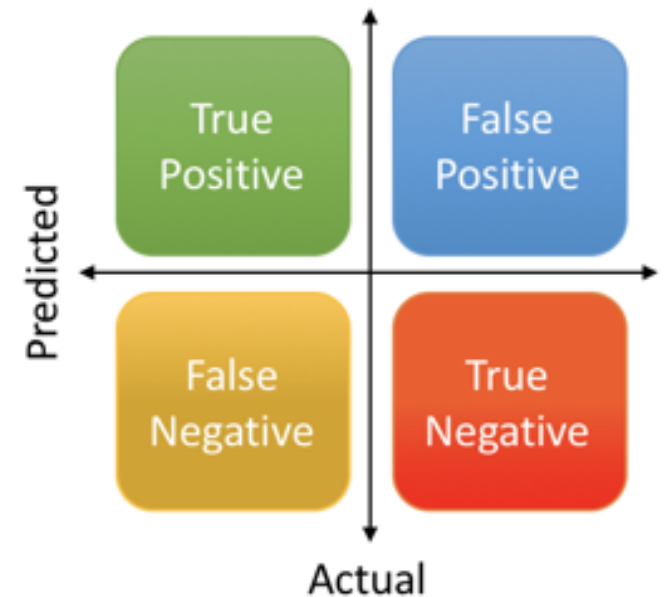
# Model evaluation

# Model evaluation

- Precision & recall

$$\text{Precision} = \frac{\text{True Positive}}{\text{Actual Results}} \quad \text{or} \quad \frac{\text{True Positive}}{\text{True Positive} + \text{False Positive}}$$

$$\text{Recall} = \frac{\text{True Positive}}{\text{Predicted Results}} \quad \text{or} \quad \frac{\text{True Positive}}{\text{True Positive} + \text{False Negative}}$$

$$\text{Accuracy} = \frac{\text{True Positive} + \text{True Negative}}{\text{Total}}$$

# Time for practice

Build your first machine learning model with python

# Advanced topics about ML

# Feature engineering / Feature selection

- Generate new feature

  - Domain knowhow

  - Data exploration

- Feature selection

  - Correlation

  - Lasso, Ridge regression

  - Index of feature importance

# Other ML models

- Other machine learning models

    - Support Vector Machine

    - K-Nearest Neighbor

    - Naïve-bayes

    - Neural network

# Other ML models

- Ml models that usually showed on ml competitions

  - Bagging  : [Random Forest](#)

  - Boosting  :  [XGBoost](#), [LightGBM](#), [CatBoost](#)

  - Neural Network

  - Stacking model