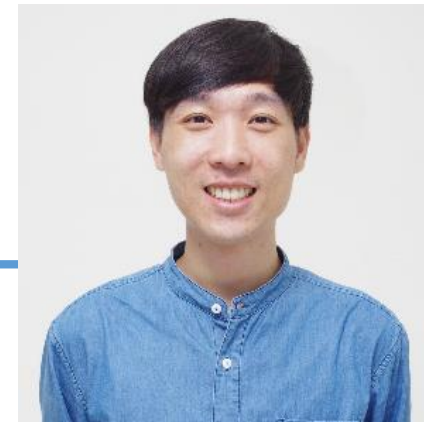# Data Science

Hands on machine learning
with Scikit-learn

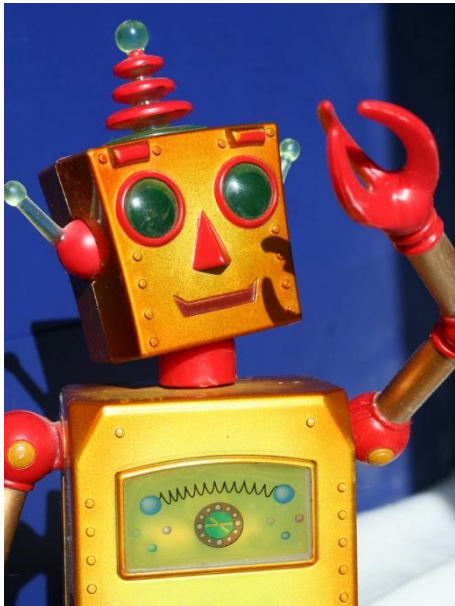**Yueh-Lin Tsai**
**2019 / 07**

# 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 machine learning ?

- Extract relations/patterns from data automatically

- Apply those rules to unseen data

$$f(x) = y$$

# Machine learning workflow

- Problem definition

- Data exploration / preprocessing

- Build model

- Model evaluation

# Type of machine learning

- Supervised learning
  - Regression
  - Classification

- Unsupervised learning
  - Cluster
  - Dimension reduction

- Reinforcement learning

# Scikit - learn

- Package for machine learning in python

- What can sklearn do

  - Preprocess

  - MI models

  - Evaluation metrics



scikit-learn
Machine Learning in Python

- Simple and efficient tools for data mining and data analysis
- Accessible to everybody, and reusable in various contexts
- Built on NumPy, SciPy, and matplotlib
- Open source, commercially usable - BSD license

**Classification**

Identifying to which category an object belongs to.

Applications: Spam detection, Image recognition.
Algorithms: SVM, nearest neighbors, random forest, ... — Examples

**Regression**

Predicting a continuous-valued attribute associated with an object.

Applications: Drug response, Stock prices.
Algorithms: SVR, ridge regression, Lasso, ... — Examples

**Clustering**

Automatic grouping of similar objects into sets.

Applications: Customer segmentation, Grouping experiment outcomes
Algorithms: k-Means, spectral clustering, mean-shift, ... — Examples

**Dimensionality reduction**

Reducing the number of random variables to consider.

Applications: Visualization, Increased efficiency
Algorithms: PCA, feature selection, non-negative matrix factorization. — Examples

**Model selection**

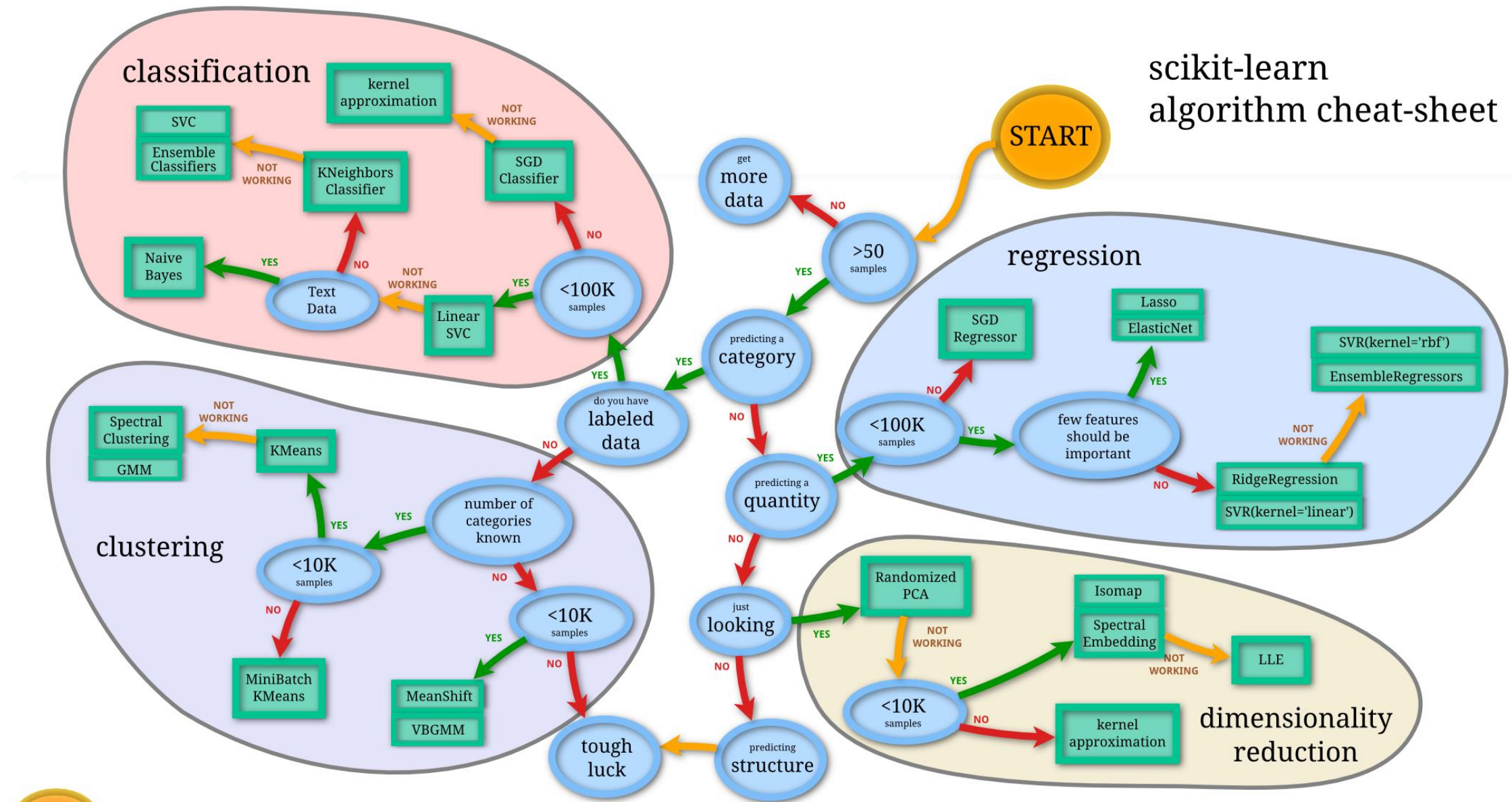Comparing, validating and choosing parameters and models.

Goal: Improved accuracy via parameter tuning
Modules: grid search, cross validation, metrics. — Examples

**Preprocessing**

Feature extraction and normalization.

Application: Transforming input data such as text for use with machine learning algorithms.
Modules: preprocessing, feature extraction. — Examples

# scikit-learn algorithm cheat-sheet

**classification**

- kernel approximation
- SVC
- Ensemble Classifiers
- KNeighbors Classifier
- SGD Classifier
- Naive Bayes
- Linear SVC
- Text Data
- <100K samples

**regression**

- SGD Regressor
- Lasso / ElasticNet
- SVR(kernel='rbf') / EnsembleRegressors
- few features should be important
- <100K samples
- RidgeRegression / SVR(kernel='linear')

**clustering**

- Spectral Clustering / GMM
- KMeans
- <10K samples
- number of categories known
- MiniBatch KMeans
- MeanShift / VBGMM
- <10K samples

**dimensionality reduction**

- Randomized PCA
- Isomap / Spectral Embedding
- LLE
- <10K samples
- kernel approximation

START → >50 samples → get more data (NO) / predicting a category (YES) → do you have labeled data → predicting a quantity → <100K samples → just looking → predicting structure → tough luck

Source

# 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)))
```

# Supervised learning

# Machine learning workflow

- Problem definition

- Data exploration / preprocessing

- Build model

- Model evaluation

# 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

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

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

- Normalize data

  - Standard scale

  - Min-max scale

| 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

# Model selection

- Linear model

  - Linear Regression

  - Logistic Regression

  - Ridge regression

  - Lasso regression

$$y = f(a_0 + a_1 x_1 + a_2 x_2 + \ldots)$$

**The aim of linear model is to find a line which minimize the distance (error) between data and the line**

# Linear model

- Linear regression

# Linear model

$$f(x) = \frac{1}{1 + e^{-(x)}}$$

- Logistic regression

# Model selection

- Tree-based model

  - Decision Tree

  - Random Forest

**The aim of tree-based model is to find good cut-off points to split data repetitively.**



Decision Tree for *PlayTennis*

# 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

- Other machine learning models

  - Support Vector Machine

  - K-Nearest Neighbor

  - Naïve-bayes

  - Neural network

# Model evaluation

- Regression problem

  - R - square, $R^2$

  - Mean Squared Error, MSE

  - Mean Absolute Error, MAE

# Model evaluation

- R-square

  - Range from 0 to 1

  - More is better

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

# Model evaluation

- Mean Square Error

- Mean Absolute Error

  - Less is better

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

The square of the difference between actual and predicted

Divide by the total number of data points

Predicted output value

Actual output value

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

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

# Model evaluation



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

# Model evaluation

- Precision & recall

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

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

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

# Other topics – supervised learning

- Feature engineering and feature selection

- Overfitting and generalizations

- Parameter selection

# Unsupervised learning

# Machine learning workflow

- Problem definition

- Data exploration / preprocessing

- Build model

- Model evaluation

# Problem definition and model selection

- clustering

  - K-means clustering

  - Hierarchical clustering

  - DBSCAN

**The aim of clustering is to categorize data base on their similarity**

# K-means

0a. Données d'entrée
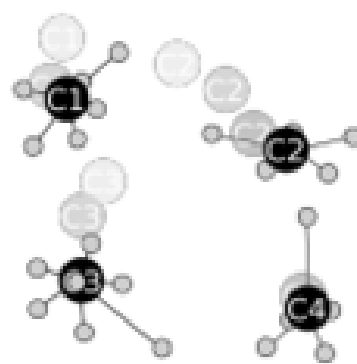
0b. intialisation
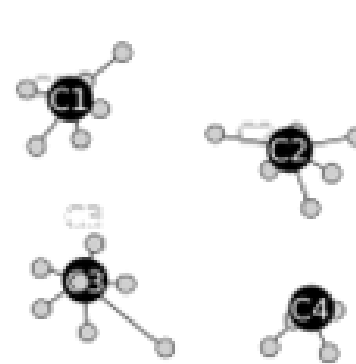
1a. assignation

1b. calcul des points moyens

2a. assignation

2b. calcul des points moyens

3a. assignation

3b. calcul des points moyens

4a. assignation
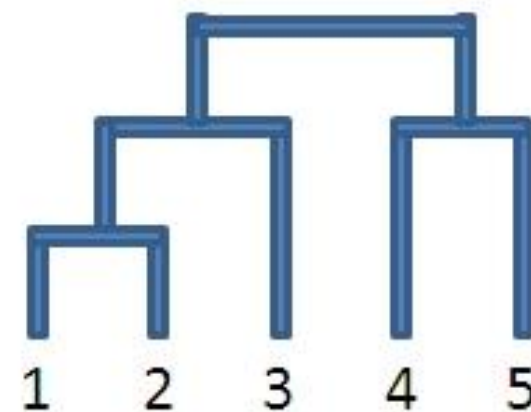clusters stables (fin)

# Hierarchical clustering

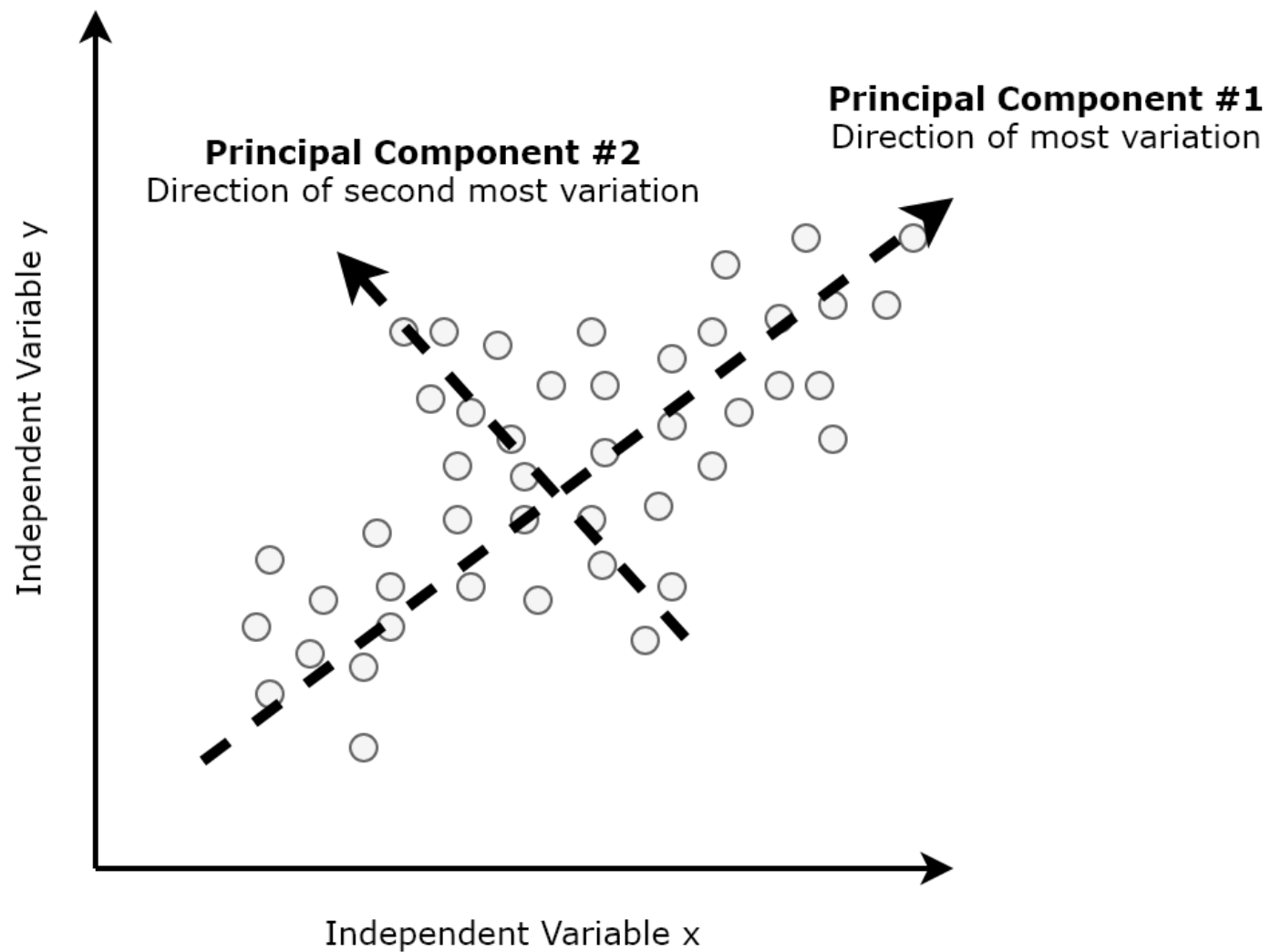階層式

# DBSCAN

$MinPts = 4$

# Problem definition and model selection

- dimension reduction

  - Principal component analysis, PCA

**The aim of dimension reduction is to describe high-dimension data with fewer dimension**

from sklearn.decomposition import PCA

# PCA

# Model evaluation

- Cluster

  - Homogeneity scores

- Dimension reduction

  - Information loss ratio

# Supplementary

# Feature engineering / feature selection

- Feature generation

  - From domain knowledge

  - By exploring data

- Feature selection

  - Based on correlation

  - Based on coefficients in lasso / ridge regression model

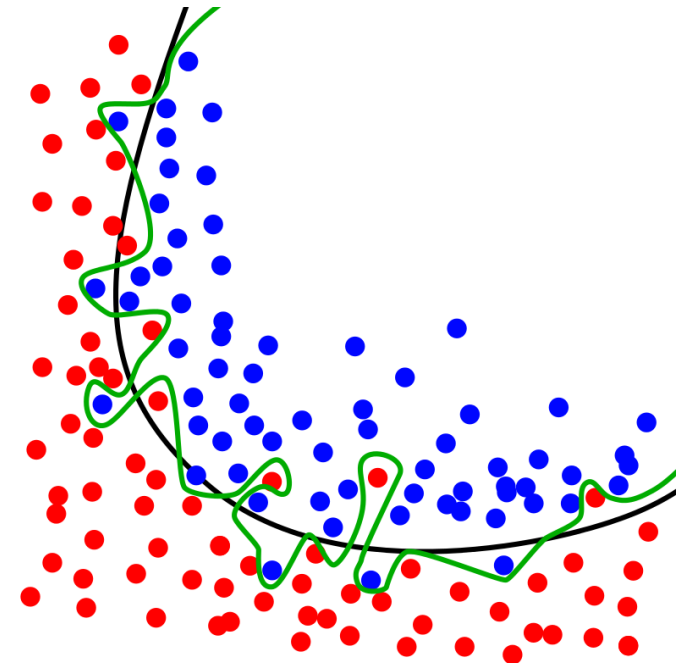  - Based on feature importance index in tree-based models

# Other ML models

- MI models that usually showed on ml competitions

  - Bagging ： [Random Forest](#)

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

  - Neural Network

  - Stacking model

# Overfitting and generalization

We can always get 100% accuracy on training set with a powerful model, but …

- How to detect overfitting ?

- How to resolve / avoid overfitting ?

  - Decrease the power of present model
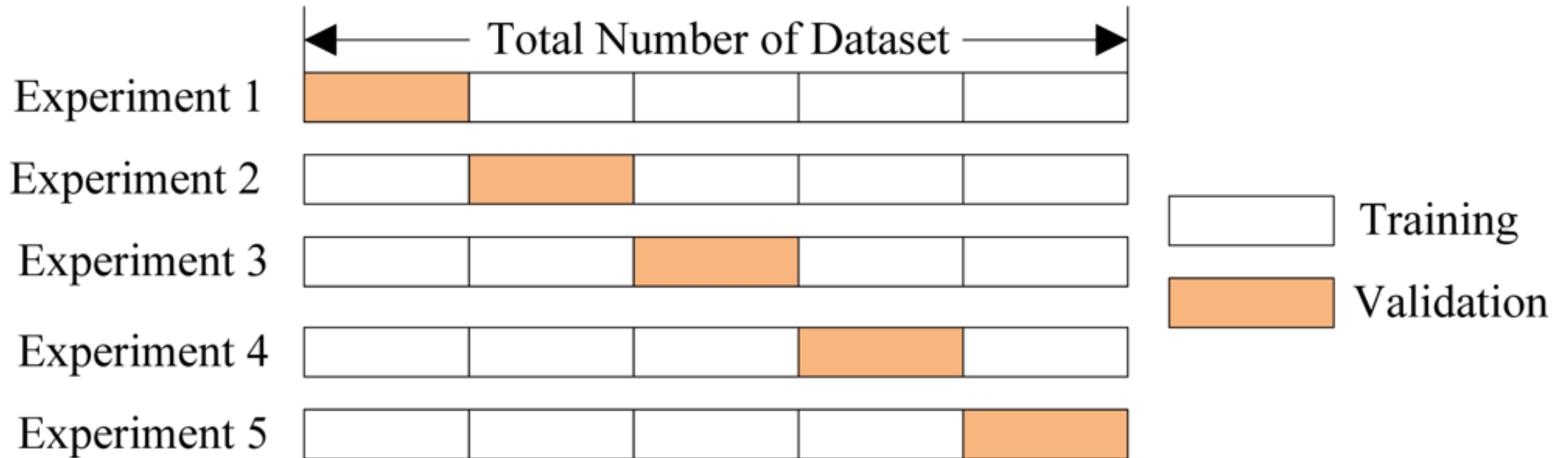
  - Ensemble methods
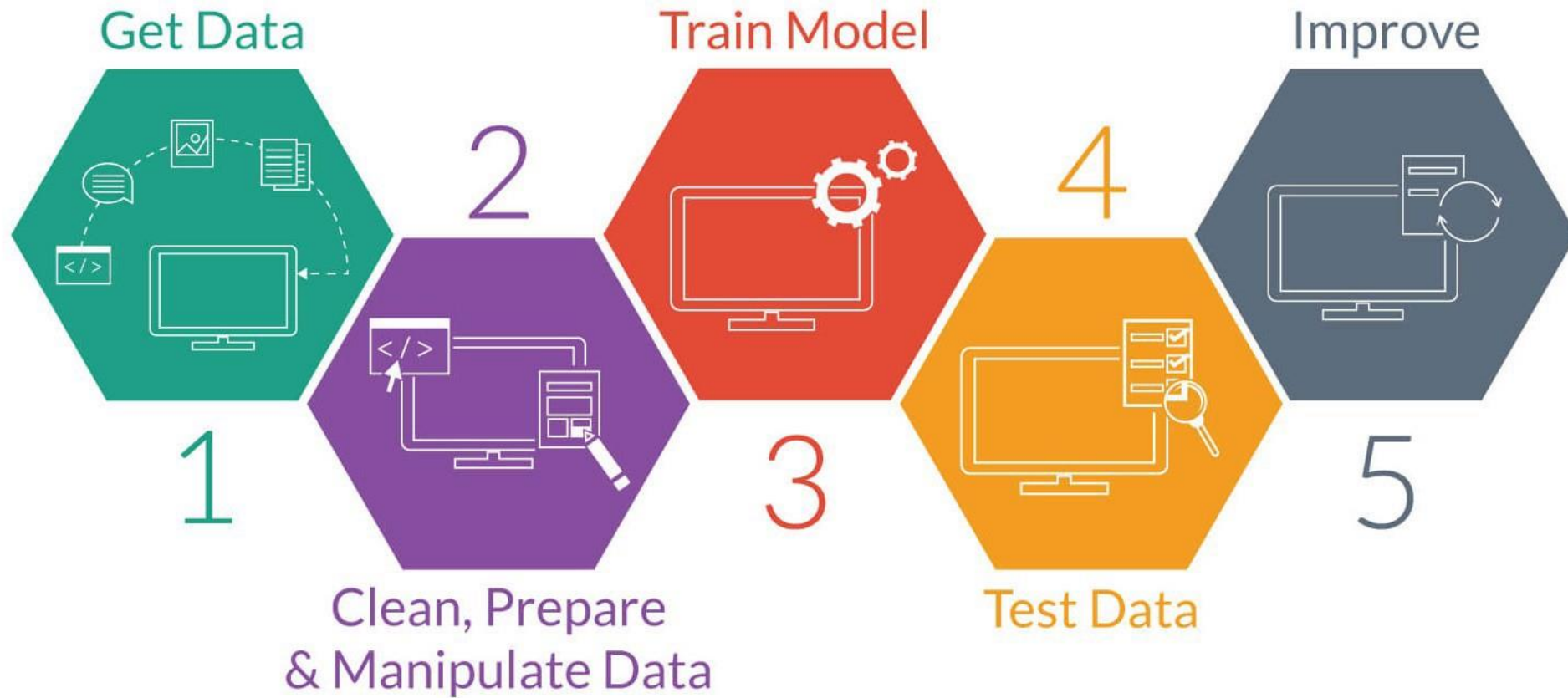
# Parameter selection

- Select a proper set of parameters would affect the power of model hence make a better prediction

- How to decide which parameter set should be used ?

# Parameter selection

## Cross Validation

# Summary of Machine learning



Get Data

1

Clean, Prepare
& Manipulate Data

2

Train Model

3

4

Test Data

Improve

5

# How to build a great model ?

- Before model fitting

  - Preprocessing
  - Data exploration
  - Data engineering

- After preliminary model

  - Parameter tuning
  - Cross validation