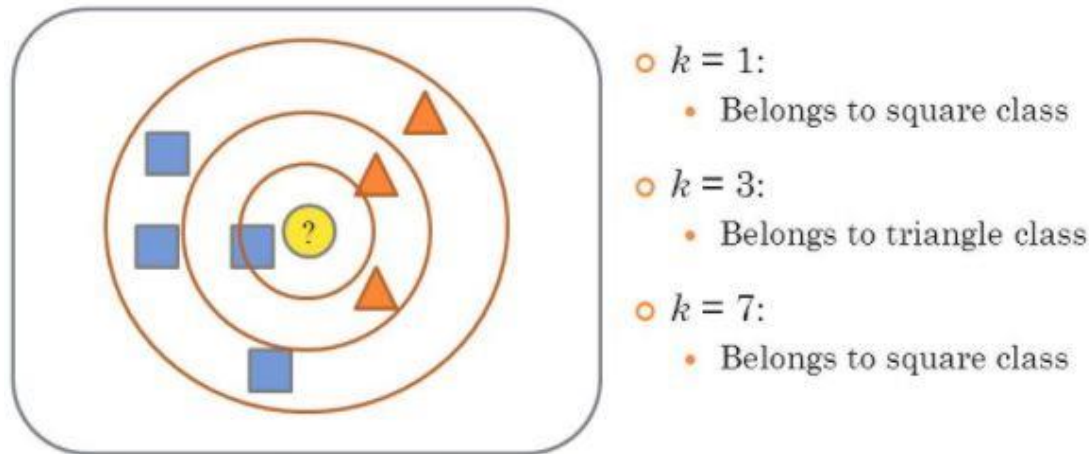# Classification Supervised Learning
# (Part1)

- **K-Nearest Neighbor**

- **Decision Tree**
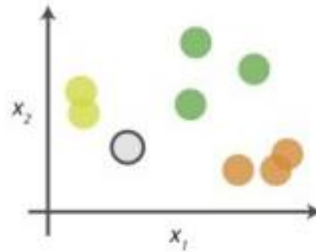  - **CART**
  - **ID3**

# K-Nearest Neighbor

- **A non-parametric method used for classification and regression**
- **Also called kNN**
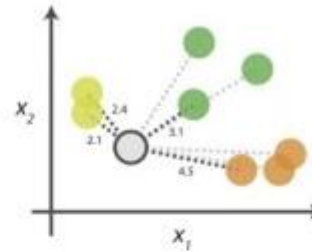  - **"k" mean how many neighbors should be considered to help classification/regression**



- $k = 1$:
  - Belongs to square class
- $k = 3$:
  - Belongs to triangle class
- $k = 7$:
  - Belongs to square class

**kNN intuitive concept**

15-4

- **L1 distance (Manhattan distance)**
- **L2 distance (Euclidean distance)**

L1 (Manhattan) distance

$$d_1(I_1, I_2) = \sum_p \left| I_1^p - I_2^p \right|$$

K = 1

L2 (Euclidean) distance

$$d_2(I_1, I_2) = \sqrt{\sum_p \left( I_1^p - I_2^p \right)^2}$$

K = 1

15-6

Euclidean distance = $\sqrt{(5-1)^2 + (4-1)^2}$ = 5

Manhattan distance = $|5-1| + |4-1|$ = 7

- **K is small**
  - **sensitive to noise points**

- **K is large**
  - **neighborhood may include points from other classes**
  - **smoother boundary**
  - **If too large, machine always predict majority class**

K=1    K=15

- **http://vision.stanford.edu/teaching/cs231n-demos/knn/**

- **1-NN**
  - **Voronoi Diagram**

**Don't use kNN on images**
**(Distance between pixels are meaningless)**



| Original | Boxed | Shifted | Tinted |

All 3 images have same L2 distance to the one on the left!

$$d = 1.4142 \quad\quad\quad\quad d = 1.4142$$

counter-intuitive results

- Curse of dimensionality

Dimensions = 1
Points = 4

Dimensions = 2
Points = $4^2$

Dimensions = 3
Points = $4^3$

add features

add features

Feature 1

Feature 2

Feature 1

Feature 3

Feature 2

Feature 1

Linear separable in high dimensionality

# Curse of dimensionality

- **Increase dimensionality may obtain perfect classfication**

- **However, extend too many dimensionality(features) lead to overfitting**

- **Example**
  - **KNN**
    - example/supervised learning
- **Practice**
  - **Try to use knn to predict different varieties of wheat**
    - dataset/seeds_dataset.csv
    - practice/supervised learning
  - **More information about the dataset**
    - https://archive.ics.uci.edu/ml/datasets/seeds#

# Decision Tree

- **A decision support tool that uses a tree-like graph of decisions and their possible consequences**

- **Common method in decision tree**
  - **ID3**
  - **CART**



Table Data

target indicator

| | A | B | C | ... | Z |
|---|---|---|---|---|---|
| 1 | | | | | |
| 2 | | | | | |
| 3 | | | | | |
| 4 | | | | | |
| 5 | | | | | |
| 6 | | | | | |
| ... | | | | | |

atom

indicators

Decision Tree

root

node

leaf

15-18

ROOT Node

Branch/ Sub-Tree

**Splitting**

Decision Node

A Decision Node

Terminal Node

Decision Node

Terminal Node

Terminal Node

B

C

Terminal Node

Terminal Node

**Note:-** A is parent node of B and C.

How to define a good split

**TibaMe**

- **Information/Gini gain**
  - **Index to decide how to split each node**
  - **Usually, we choose max information/gini gain as candidate to split**

**CART**

$$Gini\ Gain = Gini(before\ splitting) - E[Gini(after\ splitting)]$$

**ID3**

$$Information\ Gain = Entropy(before\ splitting) - E[Entropy(after\ splitting)]$$

# Decision Tree - CART

- **Classification and Regression Trees(CART) model is a binary tree**
- **Split Based on One Variable**
- **Use Gini impurity to define attribute complexity under each feature**
- **Use Gini gain to split tree**

J classes and each pi is probability of class i

$$\sum_{i=1}^{J} p_i(1 - p_i) = \sum_{i=1}^{J}(p_i - p_i{}^2) = \sum_{i=1}^{J}p_i - \sum_{i=1}^{J}p_i{}^2 = 1 - \sum_{i=1}^{J}p_i{}^2$$

| Class 1 | 0 |
|---|---|
| Class 2 | 6 |

$$p(class\ 1) = \frac{0}{6}, \qquad p(class\ 2) = \frac{6}{6}$$
$$Gini = 1 - (\frac{0}{6})^2 - (\frac{6}{6})^2 = 0$$

| Class 1 | 1 |
|---|---|
| Class 2 | 5 |

$$p(class\ 1) = \frac{1}{6}, \qquad p(class\ 2) = \frac{5}{6}$$
$$Gini = 1 - (\frac{1}{6})^2 - (\frac{5}{6})^2 = 0.278$$

| Class 1 | 2 |
|---|---|
| Class 2 | 4 |

$$p(class\ 1) = \frac{2}{6}, \qquad p(class\ 2) = \frac{4}{6}$$
$$Gini = 1 - (\frac{2}{6})^2 - (\frac{4}{6})^2 = 0.444$$

15-24

Gini **Large** ⟹ **Less** Purity

Gini **Small** ⟹ **More** Purity

before splitting

| Class 1 | 6 |
|---|---|
| Class 2 | 6 |

Gini(before splitting) = 0.5

- - - - - - - - - - - - - - - - - - - - - - - - - - - - - -

after splitting

suppose there are two ways (A or B) to split the data



**A**

| Class 1 | 4 |
|---|---|
| Class 2 | 3 |

| Class 1 | 2 |
|---|---|
| Class 2 | 3 |

Gini = 0.489        Gini = 0.48

$$E[\text{Gini(after splitting)}]$$
$$=\frac{7}{12}*0.489+\frac{5}{12}*0.48=0.4852$$

**B**

| Class 1 | 1 |
|---|---|
| Class 2 | 4 |

| Class 1 | 5 |
|---|---|
| Class 2 | 2 |

Gini = 0.32        Gini = 0.408

$$E[\text{Gini(after splitting)}]$$
$$=\frac{5}{12}*0.32+\frac{7}{12}*0.408=0.37$$

before splitting

| Class 1 | 6 |
|---------|---|
| Class 2 | 6 |

Gini(before splitting) = 0.5

- - - - - - - - - - - - - - - - - - - - - - - - -

after splitting

Gini Gain on A way
=Gini(before splitting) -E[Gini(after splitting)]
=0.015

Gini Gain on B way
= Gini(before splitting) -E[Gini(after splitting)]
=0.13

**Split on B way is better**

before splitting

| Class 1 | 6 |
|---------|---|
| Class 2 | 6 |

Gini(before splitting) = 0.5

- - - - - - - - - - - - - - - - - - - - - - - - - -

after splitting     suppose there are two ways (A or B) to split the data

A

| Class 1 | 4 |
|---------|---|
| Class 2 | 3 |

| Class 1 | 2 |
|---------|---|
| Class 2 | 3 |

B

| Class 1 | 1 |
|---------|---|
| Class 2 | 4 |

| Class 1 | 5 |
|---------|---|
| Class 2 | 2 |

**Split on B way is better**

**Training Data**

**Model: Decision Tree**

- **There are many different way to deal with continuous attributes when building decision tree**
  - **The most simple way is to split by average of continuous attributes**

# Decision Tree – ID3

- **Iterative Dichotomiser 3(ID3) is a famous algorithm to generate decision tree**
- **Use information gain as index to split each node**
- **Note that ID3 can split multiple branch at each node**



General tree

Binary Tree

- **Entropy**

$$E(S) = \sum_{i=1}^{c} - p_i \log_2 p_i$$

| | |
|---|---|
| Class 1 | 0 |
| Class 2 | 6 |

$$p(class\ 1) = \frac{0}{6}, \qquad p(class\ 2) = \frac{6}{6}$$
$$Entropy = -0 * \log(0) - 1 * \log(1) = 0$$

| | |
|---|---|
| Class 1 | 1 |
| Class 2 | 5 |

$$p(class\ 1) = \frac{1}{6}, \qquad p(class\ 2) = \frac{5}{6}$$
$$Entropy = -\frac{1}{6} * \log\left(\frac{1}{6}\right) - \frac{5}{6} * \log\left(\frac{5}{6}\right) = 0.65$$

| | |
|---|---|
| Class 1 | 2 |
| Class 2 | 4 |

$$p(class\ 1) = \frac{2}{6}, \qquad p(class\ 2) = \frac{4}{6}$$
$$Entropy = -\frac{2}{6} * \log\left(\frac{2}{6}\right) - \frac{4}{6} * \log\left(\frac{4}{6}\right) = 0.91$$

15-33

$$Entropy = -x * \log(x) - (1 - x) * log(1 - x)$$

before splitting

| Class 1 | 6 |
|---------|---|
| Class 2 | 6 |

Entropy(before splitting) = 0.301

- - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - -

after splitting    suppose there are two ways (A or B) to split the data

**A**

| Class 1 | 4 |
|---------|---|
| Class 2 | 3 |

| Class 1 | 2 |
|---------|---|
| Class 2 | 3 |

Entropy = 0.297    Entropy = 0.292

$$E[\text{Gini(after splitting)}]$$
$$=\frac{7}{12} * 0.297 + \frac{5}{12} * 0.292 = 0.294$$

**B**

| Class 1 | 1 |
|---------|---|
| Class 2 | 4 |

| Class 1 | 5 |
|---------|---|
| Class 2 | 2 |

Entropy = 0.217    Entropy = 0.259

$$E[\text{Gini(after splitting)}]$$
$$=\frac{5}{12} * 0.217 + \frac{7}{12} * 0.259 = 0.242$$

15-35

before splitting

| Class 1 | 6 |
|---------|---|
| Class 2 | 6 |

Entropy(before splitting) = 0.5

- - - - - - - - - - - - - - - - - - - - - - - -

after splitting

Information Gain on A way
=Entropy(before splitting) -E[Entropy(after splitting)]
=0.007

Information Gain on B way
= Entropy (before splitting) -E[Entropy(after splitting)]
=0.069

## Split on B way is better

15-36

before splitting

| Class 1 | 6 |
|---------|---|
| Class 2 | 6 |

Gini(before splitting) = 0.5

- - - - - - - - - - - - - - - - - - - - - - - - - -

after splitting    suppose there are two ways (A or B) to split the data

A

| Class 1 | 4 |
|---------|---|
| Class 2 | 3 |

| Class 1 | 2 |
|---------|---|
| Class 2 | 3 |

B

| Class 1 | 1 |
|---------|---|
| Class 2 | 4 |

| Class 1 | 5 |
|---------|---|
| Class 2 | 2 |

**Split on B way is better**

**Predict if playing golf or not**

| Outlook | Temp | Humidity | Windy | Play Golf |
|---------|------|----------|-------|-----------|
| Rainy | Hot | High | False | No |
| Rainy | Hot | High | True | No |
| Overcast | Hot | High | False | Yes |
| Sunny | Mild | High | False | Yes |
| Sunny | Cool | Normal | False | Yes |
| Sunny | Cool | Normal | True | No |
| Overcast | Cool | Normal | True | Yes |
| Rainy | Mild | High | False | No |
| Rainy | Cool | Normal | False | Yes |
| Sunny | Mild | Normal | False | Yes |
| Rainy | Mild | Normal | True | Yes |
| Overcast | Mild | High | True | Yes |
| Overcast | Hot | Normal | False | Yes |
| Sunny | Mild | High | True | No |

| Outlook | Temp | Humidity | Windy | Play Golf |
|---------|------|----------|-------|-----------|
| Rainy | Hot | High | False | No |
| Rainy | Hot | High | True | No |
| Overcast | Hot | High | False | Yes |
| Sunny | Mild | High | False | Yes |
| Sunny | Cool | Normal | False | Yes |
| Sunny | Cool | Normal | True | No |
| Overcast | Cool | Normal | True | Yes |
| Rainy | Mild | High | False | No |
| Rainy | Cool | Normal | False | Yes |
| Sunny | Mild | Normal | False | Yes |
| Rainy | Mild | Normal | True | Yes |
| Overcast | Mild | High | True | Yes |
| Overcast | Hot | Normal | False | Yes |
| Sunny | Mild | High | True | No |

$$Entropy(before\ split) = -\frac{5}{14} * \log\left(\frac{5}{14}\right) - \frac{9}{14} * \log\left(\frac{9}{14}\right) = 0.94$$

calculate entropy if splitting on **outlook** column

| | | Play Golf | | |
|---|---|---|---|---|
| | | Yes | No | |
| Outlook | Sunny | 3 | 2 | 5 |
| | Overcast | 4 | 0 | 4 |
| | Rainy | 2 | 3 | 5 |
| | | | | 14 |

$$E[Entropy(after\ splitting)]$$
$$= P(sunny) * E(3,2) + P(overcast) * E(4,0) + P(rainy) * E(2,3)$$
$$= \left(\frac{5}{14}\right) * 0.971 + \left(\frac{4}{14}\right) * 0 + \left(\frac{5}{14}\right) * 0.971 = 0.693$$

|  |  | Play Golf | |
|---|---|---|---|
|  |  | Yes | No |
| Outlook | Sunny | 3 | 2 |
|  | Overcast | 4 | 0 |
|  | Rainy | 2 | 3 |
| Information Gain = 0.247 | | | |

**Max Gain**

|  |  | Play Golf | |
|---|---|---|---|
|  |  | Yes | No |
| Temp | Hot | 2 | 2 |
|  | Mild | 4 | 2 |
|  | Cool | 3 | 1 |
| Information  Gain = 0.029 | | | |

|  |  | Play Golf | |
|---|---|---|---|
|  |  | Yes | No |
| Humidity | High | 3 | 4 |
|  | Normal | 6 | 1 |
| Information  Gain = 0.152 | | | |

|  |  | Play Golf | |
|---|---|---|---|
|  |  | Yes | No |
| Windy | False | 6 | 2 |
|  | True | 3 | 3 |
| Information  Gain = 0.048 | | | |

**After first splitting, decision tree look like the following**

| Outlook | Temp | Humidity | Windy | Play Golf |
|---------|------|----------|-------|-----------|
| Sunny | Mild | High | FALSE | Yes |
| Sunny | Cool | Normal | FALSE | Yes |
| Sunny | Cool | Normal | TRUE | No |
| Sunny | Mild | Normal | FALSE | Yes |
| Sunny | Mild | High | TRUE | No |

| Outlook | Temp | Humidity | Windy | Play Golf |
|---------|------|----------|-------|-----------|
| Overcast | Hot | High | FALSE | Yes |
| Overcast | Cool | Normal | TRUE | Yes |
| Overcast | Mild | High | TRUE | Yes |
| Overcast | Hot | Normal | FALSE | Yes |

| Outlook | Temp | Humidity | Windy | Play Golf |
|---------|------|----------|-------|-----------|
| Rainy | Hot | High | FALSE | No |
| Rainy | Hot | High | TRUE | No |
| Rainy | Mild | High | FALSE | No |
| Rainy | Cool | Normal | FALSE | Yes |
| Rainy | Mild | Normal | TRUE | Yes |

Outlook — Sunny / Overcast / Rainy

No need to further split overcast because all of target are "Yes"

| Temp. | Humidity | Windy | Play Golf |
|-------|----------|-------|-----------|
| Hot | High | FALSE | Yes |
| Cool | Normal | TRUE | Yes |
| Mild | High | TRUE | Yes |
| Hot | Normal | FALSE | Yes |



15-43

Continue split nodes on same method

| Temp. | Humidity | Windy | Play Golf |
|-------|----------|-------|-----------|
| Mild | High | FALSE | Yes |
| Cool | Normal | FALSE | Yes |
| Mild | Normal | FALSE | Yes |
| Cool | Normal | TRUE | No |
| Mild | High | TRUE | No |

R₁: **IF** (Outlook=Sunny) **AND** (Windy=FALSE) **THEN** Play=Yes

R₂: **IF** (Outlook=Sunny) **AND** (Windy=TRUE) **THEN** Play=No

R₃: **IF** (Outlook=Overcast) **THEN** Play=Yes

R₄: **IF** (Outlook=Rainy) **AND** (Humidity=High) **THEN** Play=No

R₅: **IF** (Outlook=Rain) **AND** (Humidity=Normal) **THEN** Play=Yes



http://www.saedsayad.com/decision_tree.htm

- **Calculate target Entropy**
- **Find the information gain on each attribute**
- **Split tree on an attribute which information gain is max**
- **Repeat**

- **Pruning is a technique that reduces the size of decision trees**
  - **Reduce model complexity and overfitting**



BEFORE PRUNING          AFTER PRUNING

- **Pre-pruning**
  - **Stop the algorithm before it becomes a fully-grown tree**
    - Stop if all instances belong to the same class
    - Stop if number of instances is less than some user-specified threshold
    - Stop if expanding the current node does not improve impurity measures
    - ……

- **Post-pruning**
  - **Grow decision tree to its entirety  and trim the nodes of the decision tree in a bottom-up**
  - **If generalization error improves after trimming, replace sub-tree by a leaf node**

- **Example**
  - **Decision Tree (CART)**
    - example/supervised learning
- **Practice**
  - **Try to use decision tree to predict if abalone is old or young**
    - dataset/abalone.csv
    - practice/supervised learning
    - we assume age > 8 is old and other is young
  - **More information about the dataset**
    - https://archive.ics.uci.edu/ml/datasets/abalone