

Data Mining

INFS 4203/7203

Miao Xu

miao.xu@uq.edu.au

The University of Queensland, 2020 Semester 2

Last week

Construction of a decision tree

For a current node:

- Input:
 - Training set $D = \{(x_1, y_1), \dots, (x_n, y_n)\}$
 - Feature set $A = \{a_1, \dots, a_d\}$
- Process: Function TreeGenerate(D, A)
 1. If **stop criteria** have reached: make the current node a **leaf** node and mark its class, **return**
 2. **Select the optimal** splitting feature a_* from A
 3. For all possible values of a_* , generate a branch (child node) and a subset $D_\nu \subset D$, such that $a_* = \nu$ for all samples in D_ν
 4. For all D_ν
 5. TreeGenerate($D_\nu, A \setminus \{a_*\}$)

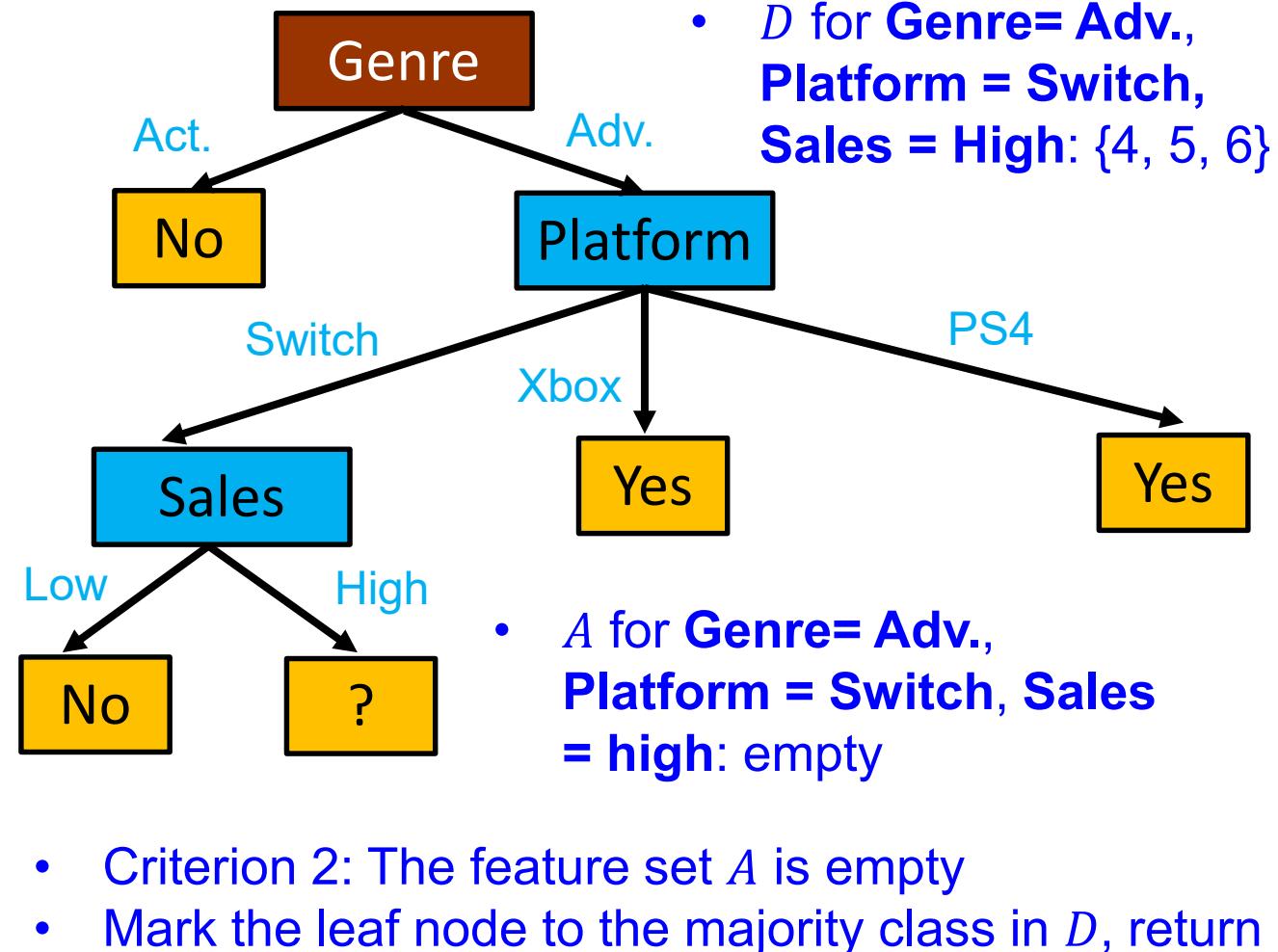
Stop criteria

- A recursive procedure
- It stops when **stop criteria** have reached:
 1. All instances in the current node belong to **the same class C**
 - Mark the leaf node to the class C , **return**
 2. The **feature set A is empty**, or all instances in D have **the same values on A**
 - Mark the leaf node to the **majority** class in D , **return**
 3. D is **empty**
 - Mark the leaf node to the **majority** class of its **parent** node, **return**
- Construct a decision tree from the **root node**

Construction of a decision tree: example

- Assume the features are optimally selected in the order
 - Genre , Platform, Sales

TID	Genre	Platform	Sales	Buy?
1	Action	Xbox	High	No
2	Action	Switch	Low	No
3	Action	PS4	High	No
4	Adventure	Switch	High	Yes
5	Adventure	Switch	High	Yes
6	Adventure	Switch	High	No
7	Adventure	Switch	Low	No
8	Adventure	PS4	High	Yes
9	Adventure	PS4	High	Yes
10	Adventure	PS4	High	No



Three measurements of purity

- Information gain (select the largest, ID3)

$$\text{Ent}(D) = - \sum_{i=1}^m p_i \log_2 p_i$$

$$\text{Gain}(D, a) = \text{Ent}(D) - \sum_{j=1}^V \frac{|D_{v_j}|}{|D|} \text{Ent}(D_{v_j})$$

- Gain ratio (select the largest, or as the heuristics in C4.5)

$$\text{IV}(a) = - \sum_{j=1}^V \frac{|D_{v_j}|}{|D|} \log_2 \frac{|D_{v_j}|}{|D|}$$

$$\text{Gain_ratio}(D, a) = \frac{\text{Gain}(D, a)}{\text{IV}(D, a)}$$

- Gini index (select the lowest, CART)

$$\text{Gini}(D) = 1 - \sum_{i=1}^m p_i^2$$

$$\text{Gini_index}(D, a) = \sum_{j=1}^V \frac{|D_{v_j}|}{|D|} \text{Gini}(D_{v_j})$$

Two pruning approaches

- Pre-pruning (in the construction)
 - Halt tree construction early: do not split a node if this would result in poorer generalization
- Post-pruning (after the construction)
 - Remove branches from a “fully grown” tree, if this would result in better generalization

We use **test data** or **cross validation** to measure the generalization performance.

Deal with continuous values

	0.3	0.4	0.6	0.8	0.9	1.0	1.1	1.2	1.3	1.4
	0.35	0.5	0.7	0.85	0.95	1.05	1.15	1.25	1.35	
Gini index:	0.4444	0.4000	0.3429	0.2667	0.4000	0.4667	0.4190	0.4750	0.4444	

$$\text{Gini_index}(\{1 - 10\}, \text{Genre}) = 0.3429$$

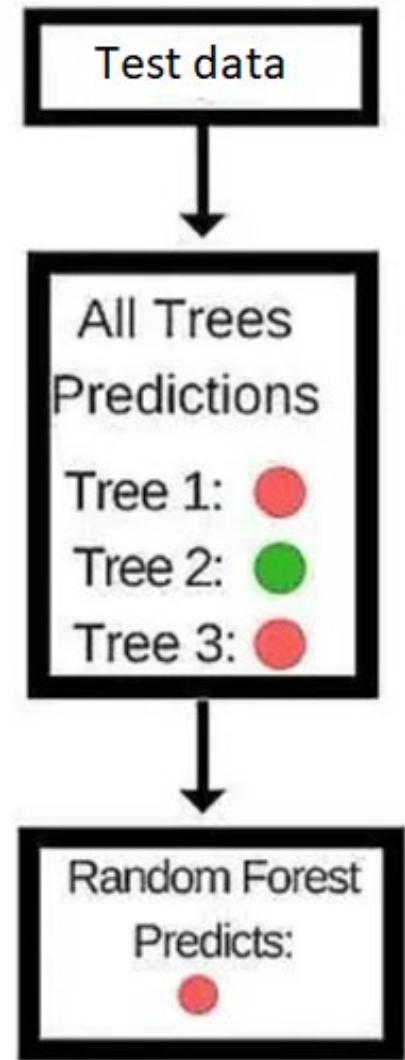
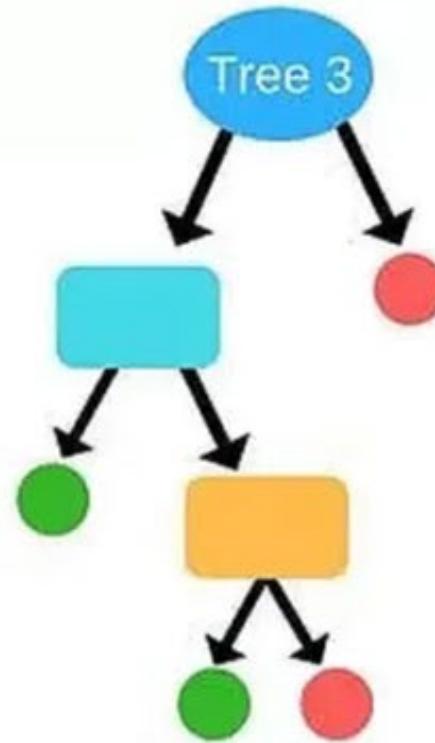
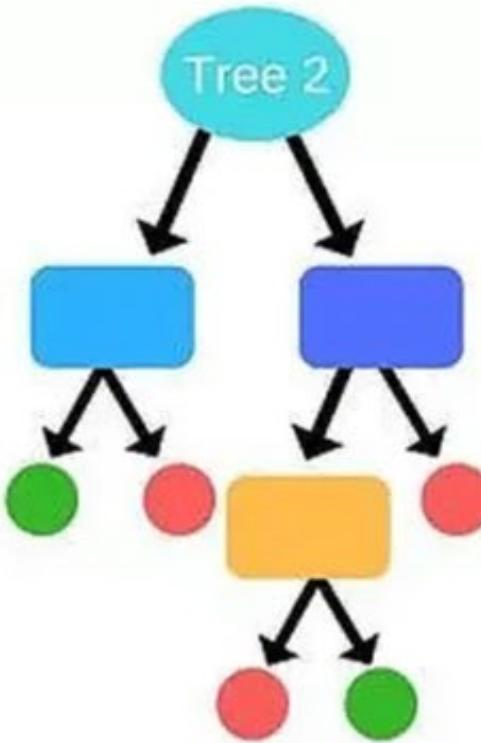
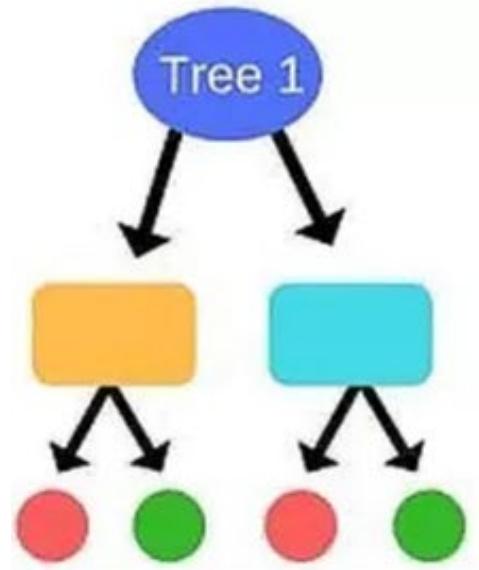
$$\text{Gini_index}(\{1 - 10\}, \text{Platform}) = 0.44$$

- Select “Sales = 0.85” as the first splitting feature and value!

TID	Genre	Platform	Sales	Buy?
1	Action	Xbox	0.8M	No
2	Action	Switch	0.3M	No
3	Action	PS4	0.6M	No
4	Adventure	Switch	0.9M	Yes
5	Adventure	Switch	1.0M	Yes
6	Adventure	Switch	1.1M	No
7	Adventure	Switch	0.4M	No
8	Adventure	PS4	1.2M	Yes
9	Adventure	PS4	1.3M	Yes
10	Adventure	PS4	1.4M	No

Random forest

- Construct several decision trees instead of only one decision tree
- Decision trees are constructed by
 - Determine k , which is a positive integer
 - Randomly sample k features among all d features
 - Construct a decision tree based on the sampled k features of all data
- Combine the output of all the trees together
 - E.g. majority voting
 - Ensemble learning method: predict with multiple classifiers



Lecture 9: Classification 3-Naïve Bayesian and k Nearest Neighbor

Outline

- Naïve Bayes classifier
- k -Nearest Neighbor classifier
- Summary

Classification through conditional probability

- (\mathbf{x}, y) are randomly generated from an unknown distribution
- Conditional probability $p(y|\mathbf{x})$: the likelihood of y given \mathbf{x}
- Example:
 - $p(\text{Yes} | (\text{Adventure, Switch, High}))$
 - $p(\text{No} | (\text{Action, Xbox, Low}))$
- The higher $p(y|\mathbf{x})$ is, the more likely \mathbf{x} is classified as y
- Assume there are m classes in total: C_1, C_2, \dots, C_m
 - Predict \mathbf{x} as C_j , where $p(C_j|\mathbf{x})$ is the maximum value among $\{p(C_1|\mathbf{x}), p(C_2|\mathbf{x}), \dots, p(C_m|\mathbf{x})\}$

How do we estimate $p(y|\mathbf{x})$?

Bayes' theorem

- Joint probability: $p(x, y)$
- Joint probability can be calculated by:

$$p(x, y) = p(x|y) \times p(y)$$

or

$$p(x, y) = p(y|x) \times p(x)$$

$$p(x|y) \times p(y) = p(y|x) \times p(x)$$

Class-conditional probability/likelihood

- Bayes' theorem:

$$p(y|x) = \frac{p(x|y) \times p(y)}{p(x)}$$

- For different y , $p(x)$ remains the same

– To get the maximum value $p(y|x)$, we evaluate $p(x|y)$ and $p(y)$

How to evaluate the prior $p(y)$

- The Law of Large Numbers (LLN): (informally) the average of the results obtained from a large number of trials should be close to the expected value
- D_y is the set of all samples belong to class y in training set D
- $p(y)$ is estimated by the ratio of examples belonging to class y in the training dataset, i.e.,

$$p(y) = \frac{|D_y|}{|D|}$$

How to evaluate the prior $p(y)$: example

- $y \in \{\text{Yes}, \text{No}\}$
- Yes occurs 4 times. There are totally 10 training instances
$$p(\text{Yes}) = \frac{4}{10} = 0.4$$
- No occurs 6 times. There are totally 10 training instances
$$p(\text{No}) = \frac{6}{10} = 0.6$$

TID	Genre	Platform	Sales	Buy?
1	Action	Xbox	High	No
2	Action	Switch	Low	No
3	Action	PS4	High	No
4	Adventure	Switch	High	Yes
5	Adventure	Switch	High	Yes
6	Adventure	Switch	High	No
7	Adventure	Switch	Low	No
8	Adventure	PS4	High	Yes
9	Adventure	PS4	High	Yes
10	Adventure	PS4	High	No

How to evaluate the likelihood $p(x|y)$

- If we do similar estimation of $p(y)$ to estimate $p(x|y)$
 - The training data may not cover all x
 - E.g. (Adv., Xbox, High)
 - When the number of features becomes large, it is hard to estimate all the $p(x|y)$
- For example, if there are 100 features, each feature has two values, the possible values of x can be 2^{100}

TID	Genre	Platform	Sales	Buy?
1	Action	Xbox	High	No
2	Action	Switch	Low	No
3	Action	PS4	High	No
4	Adventure	Switch	High	Yes
5	Adventure	Switch	High	Yes
6	Adventure	Switch	High	No
7	Adventure	Switch	Low	No
8	Adventure	PS4	High	Yes
9	Adventure	PS4	High	Yes
10	Adventure	PS4	High	No

Any other way to estimate $p(x|y)$?

Naïve Bayes Classifier

- Compute $p(y)$
- Compute $p(\mathbf{x}|y)$ by assuming all features are **independent** from each other **given** the class information y , for all $y \in \{C_1, C_2, \dots, C_m\}$
 - “Attribute conditional independence” assumption
 - The name “Naïve” is given for this strong assumption
- If $\mathbf{x} = [x_1, x_2, \dots, x_d]$, then
$$p(\mathbf{x}|y) = p(x_1|y) \times p(x_2|y) \times \cdots \times p(x_d|y)$$
- Predict \mathbf{x} as C_j , where $p(C_j)p(\mathbf{x}|C_j)$ is the maximum value among $\{p(C_1)p(\mathbf{x}|C_1), p(C_2)p(\mathbf{x}|C_2), \dots, p(C_m)p(\mathbf{x}|C_m)\}$

How do we estimate $p(x_i|y)$?

How to estimate $p(x_i|y)$

- D_y is the set of all samples belong to class y in training set D
- D_{y,x_i} is the set of samples in D_y with the i -th feature equal to x_i
- $p(x_i|y)$ is

$$p(x_i|y) = \frac{|D_{y,x_i}|}{|D_y|}$$

Compute $p(x|y)$ in Naïve Bayes classifier

- For feature Genre
 - Buy? = No: 6 examples
 - Genre = Action: 3 examples
$$p(\text{Genre} = \text{Action}|\text{Buy?} = \text{No}) = \frac{3}{6} = 0.5$$
 - Genre = Adventure: 3 examples
$$p(\text{Genre} = \text{Adventure}|\text{Buy?} = \text{No}) = \frac{3}{6} = 0.5$$
 - Buy? = Yes: 4 examples
 - Genre = Action: 0 examples
$$p(\text{Genre} = \text{Action}|\text{Buy?} = \text{Yes}) = \frac{0}{4} = 0$$
 - Genre = Adventure: 4 examples
$$p(\text{Genre} = \text{Adventure}|\text{Buy?} = \text{Yes}) = \frac{4}{4} = 1$$

TID	Genre	Platform	Sales	Buy?
1	Action	Xbox	High	No
2	Action	Switch	Low	No
3	Action	PS4	High	No
4	Adventure	Switch	High	Yes
5	Adventure	Switch	High	Yes
6	Adventure	Switch	High	No
7	Adventure	Switch	Low	No
8	Adventure	PS4	High	Yes
9	Adventure	PS4	High	Yes
10	Adventure	PS4	High	No

Compute $p(x|y)$ in Naïve Bayes classifier

- For feature Platform

- Buy? = No: 6 examples

- Platform = Xbox: 1 examples

$$p(\text{Platform} = \text{Xbox} | \text{Buy?} = \text{No}) = \frac{1}{6} = 0.2$$

- Platform = Switch: 3 examples

$$p(\text{Platform} = \text{Switch} | \text{Buy?} = \text{No}) = \frac{3}{6} = 0.5$$

- Platform = PS4: 2 examples

$$p(\text{Platform} = \text{PS4} | \text{Buy?} = \text{No}) = \frac{2}{6} = 0.3$$

TID	Genre	Platform	Sales	Buy?
1	Action	Xbox	High	No
2	Action	Switch	Low	No
3	Action	PS4	High	No
4	Adventure	Switch	High	Yes
5	Adventure	Switch	High	Yes
6	Adventure	Switch	High	No
7	Adventure	Switch	Low	No
8	Adventure	PS4	High	Yes
9	Adventure	PS4	High	Yes
10	Adventure	PS4	High	No

Activity: Compute $p(x|y)$ in Naïve Bayes classifier

- For feature Platform
 - Buy? = Yes: ? examples
 - Platform = Xbox: ? examples
 $p(\text{Platform} = \text{Xbox} | \text{Buy?} = \text{Yes}) = ?$
 - Platform = Switch: ? examples
 $p(\text{Platform} = \text{Switch} | \text{Buy?} = \text{Yes}) = ?$
 - Platform = PS4: ? examples
 $p(\text{Platform} = \text{PS4} | \text{Buy?} = \text{Yes}) = ?$

TID	Genre	Platform	Sales	Buy?
1	Action	Xbox	High	No
2	Action	Switch	Low	No
3	Action	PS4	High	No
4	Adventure	Switch	High	Yes
5	Adventure	Switch	High	Yes
6	Adventure	Switch	High	No
7	Adventure	Switch	Low	No
8	Adventure	PS4	High	Yes
9	Adventure	PS4	High	Yes
10	Adventure	PS4	High	No

Compute $p(x|y)$ in Naïve Bayes classifier

- For feature Platform
 - Buy? = Yes: 4 examples
 - Platform = Xbox: 0 examples
$$p(\text{Platform} = \text{Xbox} | \text{Buy?} = \text{Yes}) = \frac{0}{4} = 0$$
 - Platform = Switch: 2 examples
$$p(\text{Platform} = \text{Switch} | \text{Buy?} = \text{Yes}) = \frac{2}{4} = 0.5$$
 - Platform = PS4: 2 examples
$$p(\text{Platform} = \text{PS4} | \text{Buy?} = \text{Yes}) = \frac{2}{4} = 0.5$$

TID	Genre	Platform	Sales	Buy?
1	Action	Xbox	High	No
2	Action	Switch	Low	No
3	Action	PS4	High	No
4	Adventure	Switch	High	Yes
5	Adventure	Switch	High	Yes
6	Adventure	Switch	High	No
7	Adventure	Switch	Low	No
8	Adventure	PS4	High	Yes
9	Adventure	PS4	High	Yes
10	Adventure	PS4	High	No

Compute $p(x|y)$ in Naïve Bayes classifier

- For feature Sales
 - Buy? = No: 6 examples
 - Sales = High: 4 examples
$$p(\text{Sales} = \text{High} | \text{Buy?} = \text{No}) = \frac{4}{6} = 0.6$$
 - Sales = Low : 2 examples
$$p(\text{Sales} = \text{Low} | \text{Buy?} = \text{No}) = \frac{2}{6} = 0.3$$
 - Buy? = Yes: 4 examples
 - Sales = High: 4 examples
$$p(\text{Sales} = \text{High} | \text{Buy?} = \text{Yes}) = \frac{4}{4} = 1$$
 - Sales = Low: 0 examples
$$p(\text{Sales} = \text{Low} | \text{Buy?} = \text{Yes}) = \frac{0}{4} = 0$$

TID	Genre	Platform	Sales	Buy?
1	Action	Xbox	High	No
2	Action	Switch	Low	No
3	Action	PS4	High	No
4	Adventure	Switch	High	Yes
5	Adventure	Switch	High	Yes
6	Adventure	Switch	High	No
7	Adventure	Switch	Low	No
8	Adventure	PS4	High	Yes
9	Adventure	PS4	High	Yes
10	Adventure	PS4	High	No

Prediction with Naïve Bayes classifier

- Predict \mathbf{x} as C_j , where $p(C_j)p(\mathbf{x}|C_j)$ is the maximum value among $\{p(C_1)p(\mathbf{x}|C_1), p(C_2)p(\mathbf{x}|C_2), \dots, p(C_m)p(\mathbf{x}|C_m)\}$
- For **(Adventure, Switch, High)**, we have

$$\begin{aligned} & p(\text{Yes}) \times p(\text{Genre} = \text{Adventure}, \text{Platform} = \text{Switch}, \text{Sales} = \text{High} | \text{Buy?} = \text{Yes}) \\ &= p(\text{Yes}) \times p(\text{Genre} = \text{Adventure} | \text{Buy?} = \text{Yes}) p(\text{Platform} = \text{Switch} | \text{Buy?} = \text{Yes}) p(\text{Sales} = \text{High} | \text{Buy?} = \text{Yes}) \\ &= 0.4 \times 1 \times 0.5 \times 1 = 0.2 \end{aligned}$$

$$\begin{aligned} & p(\text{No}) \times p(\text{Genre} = \text{Adventure}, \text{Platform} = \text{Switch}, \text{Sales} = \text{High} | \text{Buy?} = \text{No}) \\ &= p(\text{No}) \times p(\text{Genre} = \text{Adventure} | \text{Buy?} = \text{No}) p(\text{Platform} = \text{Switch} | \text{Buy?} = \text{No}) p(\text{Sales} = \text{High} | \text{Buy?} = \text{No}) \\ &= 0.6 \times 0.5 \times 0.5 \times 0.6 = 0.09 < 0.2 \quad \text{Prediction: Yes} \end{aligned}$$

- In this way, **(Adventure, Switch, High)** is predicted as **Yes**
- How about **(Action, PS4, High)**?

Activity: Predict by Naïve Bayes classifier

- What is the prediction of (Action, PS4, High) ?

- $y \in \{\text{Yes}, \text{No}\}$
- Yes occurs 4 times. There are totally 10 training instances

$$p(\text{Yes}) = \frac{4}{10} = 0.4$$

- No occurs 6 times. There are totally 10 training instances

$$p(\text{No}) = \frac{6}{10} = 0.6$$

- For feature Genre

- Buy? = No: 6 examples

- Genre = Action: 3 examples

$$p(\text{Genre} = \text{Action} | \text{Buy?} = \text{No}) = \frac{3}{6} = 0.5$$

- Genre = Adventure: 3 examples

$$p(\text{Genre} = \text{Adventure} | \text{Buy?} = \text{No}) = \frac{3}{6} = 0.5$$

- Buy? = Yes: 4 examples

- Genre = Action: 0 examples

$$p(\text{Genre} = \text{Action} | \text{Buy?} = \text{Yes}) = \frac{0}{4} = 0$$

- Genre = Adventure: 4 examples

$$p(\text{Genre} = \text{Adventure} | \text{Buy?} = \text{Yes}) = \frac{4}{4} = 1$$

- For feature Platform

- Buy? = Yes: 4 examples

- Platform = Xbox: 0 examples

$$p(\text{Platform} = \text{Xbox} | \text{Buy?} = \text{Yes}) = \frac{0}{4} = 0$$

- Platform = Switch: 2 examples

$$p(\text{Platform} = \text{Switch} | \text{Buy?} = \text{Yes}) = \frac{2}{4} = 0.5$$

- Platform = PS4: 2 examples

$$p(\text{Platform} = \text{PS4} | \text{Buy?} = \text{Yes}) = \frac{2}{4} = 0.5$$

- For feature Platform

- Buy? = No: 6 examples

- Platform = Xbox: 1 examples

$$p(\text{Platform} = \text{Xbox} | \text{Buy?} = \text{No}) = \frac{1}{6} = 0.2$$

- Platform = Switch: 3 examples

$$p(\text{Platform} = \text{Switch} | \text{Buy?} = \text{No}) = \frac{3}{6} = 0.5$$

- Platform = PS4: 2 examples

$$p(\text{Platform} = \text{PS4} | \text{Buy?} = \text{No}) = \frac{2}{6} = 0.3$$

- For feature Sales

- Buy? = No: 6 examples

- Sales = High: 4 examples

$$p(\text{Sales} = \text{High} | \text{Buy?} = \text{No}) = \frac{4}{6} = 0.6$$

- Sales = Low : 2 examples

$$p(\text{Sales} = \text{Low} | \text{Buy?} = \text{No}) = \frac{2}{6} = 0.3$$

- Buy? = Yes: 4 examples

- Sales = High: 4 examples

$$p(\text{Sales} = \text{High} | \text{Buy?} = \text{Yes}) = \frac{4}{4} = 1$$

- Sales = Low: 0 examples

$$p(\text{Sales} = \text{Low} | \text{Buy?} = \text{Yes}) = \frac{0}{4} = 0$$

Prediction with Naïve Bayes classifier

- For (Action, PS4, High), we have

$$\begin{aligned} & p(\text{Yes}) \times p(\text{Genre} = \text{Action}, \text{Platform} = \text{PS4}, \text{Sales} = \text{High} | \text{Buy?} = \text{Yes}) \\ &= p(\text{Yes}) \times p(\text{Genre} = \text{Action} | \text{Buy?} = \text{Yes}) \ p(\text{Platform} = \text{PS4} | \text{Buy?} = \text{Yes}) \ p(\text{Sales} = \text{High} | \text{Buy?} = \text{Yes}) \\ &= 0.4 \times 0 \times 0.5 \times 1 = 0 \end{aligned}$$

$$\begin{aligned} & p(\text{No}) \times p(\text{Genre} = \text{Action}, \text{Platform} = \text{Xbox}, \text{Sales} = \text{High} | \text{Buy?} = \text{No}) \\ &= p(\text{No}) \times p(\text{Genre} = \text{Action} | \text{Buy?} = \text{No}) \ p(\text{Platform} = \text{PS4} | \text{Buy?} = \text{No}) \ p(\text{Sales} = \text{High} | \text{Buy?} = \text{No}) \\ &= 0.6 \times 0.5 \times 0.3 \times 0.6 = 0.054 \end{aligned}$$

Prediction: **No**

- How to handle the zero probability?
 - If any zero-probability exists, we **have to** use “Laplacian correction” in Naïve Bayes classifier

Laplacian correction

- Assume there are m classes in total: C_1, C_2, \dots, C_m
- Assume for the i -th feature, there are V_i possible values:
 - $\{v_1, v_2, \dots, v_{V_i}\}$
- D_y, D_{y,x_i} defined as before
- After Laplacian correction:

$$\hat{p}(y) = \frac{|D_y| + 1}{|D| + m} \quad \hat{p}(x_i|y) = \frac{|D_{y,x_i}| + 1}{|D_y| + V_i}$$

- It fixes the problem of zero probability. When dataset becomes large, the effect of correction on non-zero probability can be negligible

Laplacian correction: example

- $y \in \{\text{Yes}, \text{No}\}$, 2 classes
- Yes occurs 4 times, there are total 10 training instances

$$\hat{p}(\text{Yes}) = \frac{4 + 1}{10 + 2} = 0.4$$

- No occurs 6 times, there are total 10 training instances

$$\hat{p}(\text{No}) = \frac{6 + 1}{10 + 2} = 0.6$$

TID	Genre	Platform	Sales	Buy?
1	Action	Xbox	High	No
2	Action	Switch	Low	No
3	Action	PS4	High	No
4	Adventure	Switch	High	Yes
5	Adventure	Switch	High	Yes
6	Adventure	Switch	High	No
7	Adventure	Switch	Low	No
8	Adventure	PS4	High	Yes
9	Adventure	PS4	High	Yes
10	Adventure	PS4	High	No

Laplacian correction: example

- For feature Genre, 2 values

 - Buy? = No: 6 examples

 - Genre = Action: 3 examples

$$\hat{p}(\text{Genre} = \text{Action} | \text{Buy?} = \text{No}) = \frac{3 + 1}{6 + 2} = 0.5$$

 - Genre = Adventure: 3 examples

$$\hat{p}(\text{Genre} = \text{Adventure} | \text{Buy?} = \text{No}) = \frac{3 + 1}{6 + 2} = 0.5$$

 - Buy? = Yes: 4 examples

 - Genre = Action: 0 examples

$$\hat{p}(\text{Genre} = \text{Action} | \text{Buy?} = \text{Yes}) = \frac{0 + 1}{4 + 2} = 0.2$$



Non-zero now

 - Genre = Adventure: 4 examples

$$\hat{p}(\text{Genre} = \text{Adventure} | \text{Buy?} = \text{Yes}) = \frac{4 + 1}{4 + 2} = 0.8$$

TID	Genre	Platform	Sales	Buy?
1	Action	Xbox	High	No
2	Action	Switch	Low	No
3	Action	PS4	High	No
4	Adventure	Switch	High	Yes
5	Adventure	Switch	High	Yes
6	Adventure	Switch	High	No
7	Adventure	Switch	Low	No
8	Adventure	PS4	High	Yes
9	Adventure	PS4	High	Yes
10	Adventure	PS4	High	No

Activity: Laplacian correction

- For feature Platform, ? values
 - Buy? = No: ? examples
 - Platform = Xbox: ? examples
 $\hat{p}(\text{Platform} = \text{Xbox} | \text{Buy?} = \text{No}) = ?$
 - Platform = Switch: ? examples
 $\hat{p}(\text{Platform} = \text{Switch} | \text{Buy?} = \text{No}) = ?$
 - Platform = PS4: ? examples
 $\hat{p}(\text{Platform} = \text{PS4} | \text{Buy?} = \text{No}) = ?$

TID	Genre	Platform	Sales	Buy?
1	Action	Xbox	High	No
2	Action	Switch	Low	No
3	Action	PS4	High	No
4	Adventure	Switch	High	Yes
5	Adventure	Switch	High	Yes
6	Adventure	Switch	High	No
7	Adventure	Switch	Low	No
8	Adventure	PS4	High	Yes
9	Adventure	PS4	High	Yes
10	Adventure	PS4	High	No

Laplacian correction: example

- For feature Platform, 3 values
 - Buy? = No: 6 examples
 - Platform = Xbox: 1 examples
$$\hat{p}(\text{Platform} = \text{Xbox} | \text{Buy?} = \text{No}) = \frac{1+1}{6+3} = 0.2$$
 - Platform = Switch: 3 examples
$$\hat{p}(\text{Platform} = \text{Switch} | \text{Buy?} = \text{No}) = \frac{3+1}{6+3} = 0.4$$
 - Platform = PS4: 2 examples
$$\hat{p}(\text{Platform} = \text{PS4} | \text{Buy?} = \text{No}) = \frac{2+1}{6+3} = 0.3$$

TID	Genre	Platform	Sales	Buy?
1	Action	Xbox	High	No
2	Action	Switch	Low	No
3	Action	PS4	High	No
4	Adventure	Switch	High	Yes
5	Adventure	Switch	High	Yes
6	Adventure	Switch	High	No
7	Adventure	Switch	Low	No
8	Adventure	PS4	High	Yes
9	Adventure	PS4	High	Yes
10	Adventure	PS4	High	No

Laplacian correction: example

- For feature Platform, 3 values
 - Buy? = Yes: 4 examples
 - Platform = Xbox: 0 examples
 - Platform = Switch: 2 examples
 - Platform = PS4: 2 examples
$$\hat{p}(\text{Platform} = \text{Xbox} | \text{Buy?} = \text{Yes}) = \frac{0 + 1}{4 + 3} = 0.1$$
 - Non-zero now

TID	Genre	Platform	Sales	Buy?
1	Action	Xbox	High	No
2	Action	Switch	Low	No
3	Action	PS4	High	No
4	Adventure	Switch	High	Yes
5	Adventure	Switch	High	Yes
6	Adventure	Switch	High	No
7	Adventure	Switch	Low	No
8	Adventure	PS4	High	Yes
9	Adventure	PS4	High	Yes
10	Adventure	PS4	High	No

Laplacian correction: example

- For feature Sales, 2 values

 - Buy? = No: 6 examples

 - Sales = High: 4 examples

$$\hat{p}(\text{Sales} = \text{High} | \text{Buy?} = \text{No}) = \frac{4 + 1}{6 + 2} = 0.6$$

 - Sales = Low : 2 examples

$$\hat{p}(\text{Sales} = \text{Low} | \text{Buy?} = \text{No}) = \frac{2 + 1}{6 + 2} = 0.4$$

 - Buy? = Yes: 4 examples

 - Sales = High: 4 examples

$$\hat{p}(\text{Sales} = \text{High} | \text{Buy?} = \text{Yes}) = \frac{4 + 1}{4 + 2} = 0.8$$

 - Sales = Low: 0 examples

$$\hat{p}(\text{Sales} = \text{Low} | \text{Buy?} = \text{Yes}) = \frac{0 + 1}{4 + 2} = 0.2 \quad \text{Non-zero now}$$

TID	Genre	Platform	Sales	Buy?
1	Action	Xbox	High	No
2	Action	Switch	Low	No
3	Action	PS4	High	No
4	Adventure	Switch	High	Yes
5	Adventure	Switch	High	Yes
6	Adventure	Switch	High	No
7	Adventure	Switch	Low	No
8	Adventure	PS4	High	Yes
9	Adventure	PS4	High	Yes
10	Adventure	PS4	High	No

Laplacian correction: prediction

- For (Adventure, Switch, High), we have

$$\begin{aligned}\hat{p}(\text{Yes}) \times \hat{p}(\text{Genre} = \text{Adventure}, \text{Platform} = \text{Switch}, \text{Sales} = \text{High} | \text{Buy?} = \text{Yes}) \\ = \hat{p}(\text{Yes}) \times \hat{p}(\text{Genre} = \text{Adventure} | \text{Buy?} = \text{Yes}) \hat{p}(\text{Platform} = \text{Switch} | \text{Buy?} = \text{Yes}) \hat{p}(\text{Sales} = \text{High} | \text{Buy?} = \text{Yes}) \\ = 0.4 \times 0.8 \times 0.4 \times 0.8 = 0.1024\end{aligned}$$

$$\begin{aligned}\hat{p}(\text{No}) \times \hat{p}(\text{Genre} = \text{Adventure}, \text{Platform} = \text{Switch}, \text{Sales} = \text{High} | \text{Buy?} = \text{No}) \\ = \hat{p}(\text{No}) \times \hat{p}(\text{Genre} = \text{Adventure} | \text{Buy?} = \text{No}) \hat{p}(\text{Platform} = \text{Switch} | \text{Buy?} = \text{No}) \hat{p}(\text{Sales} = \text{High} | \text{Buy?} = \text{No}) \\ = 0.6 \times 0.5 \times 0.4 \times 0.6 = 0.072 < 0.1024\end{aligned}$$

Prediction: Yes

- For (Action, PS4, High), we have

$$\begin{aligned}\hat{p}(\text{Yes}) \times \hat{p}(\text{Genre} = \text{Action}, \text{Platform} = \text{PS4}, \text{Sales} = \text{High} | \text{Buy?} = \text{Yes}) \\ = \hat{p}(\text{Yes}) \times \hat{p}(\text{Genre} = \text{Action} | \text{Buy?} = \text{Yes}) \hat{p}(\text{Platform} = \text{PS4} | \text{Buy?} = \text{Yes}) \hat{p}(\text{Sales} = \text{High} | \text{Buy?} = \text{Yes}) \\ = 0.4 \times 0.2 \times 0.4 \times 0.8 = 0.0256\end{aligned}$$

$$\begin{aligned}\hat{p}(\text{No}) \times \hat{p}(\text{Genre} = \text{Action}, \text{Platform} = \text{PS4}, \text{Sales} = \text{High} | \text{Buy?} = \text{No}) \\ = \hat{p}(\text{No}) \times \hat{p}(\text{Genre} = \text{Action} | \text{Buy?} = \text{No}) \hat{p}(\text{Platform} = \text{PS4} | \text{Buy?} = \text{No}) \hat{p}(\text{Sales} = \text{High} | \text{Buy?} = \text{No}) \\ = 0.6 \times 0.5 \times 0.3 \times 0.6 = 0.054 > 0.0256\end{aligned}$$

Prediction: No

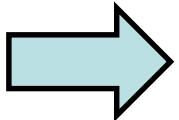
Continuous values

Training data

TID	Genre	Platform	Sales	Buy?
1	Action	Xbox	High	No
2	Action	Switch	Low	No
3	Action	PS4	High	No
4	Adventure	Switch	High	Yes
5	Adventure	Switch	High	Yes
6	Adventure	Switch	High	No
7	Adventure	Switch	Low	No
8	Adventure	PS4	High	Yes
9	Adventure	PS4	High	Yes
10	Adventure	PS4	High	No

Training data

TID	Genre	Platform	Sales	Buy?
1	Action	Xbox	0.8M	No
2	Action	Switch	0.3M	No
3	Action	PS4	0.6M	No
4	Adventure	Switch	0.9M	Yes
5	Adventure	Switch	1.0M	Yes
6	Adventure	Switch	1.1M	No
7	Adventure	Switch	0.4M	No
8	Adventure	PS4	1.2M	Yes
9	Adventure	PS4	1.3M	Yes
10	Adventure	PS4	1.4M	No



- Probability density estimation
 - Assume the values of the feature follow an univariate gaussian distribution
 - Use data of class y to estimate parameters of distribution (e.g., mean μ and variance σ^2)
 - Use the probability distribution to estimate $p(x_i|y)$

$$x_i \sim N(\mu, \sigma^2): p(x_i|y) = \frac{1}{\sqrt{2\pi}\sigma} \exp \left\{ -\frac{(x_i-\mu)^2}{2\sigma^2} \right\}$$

Continuous value: example

- Data of class No:

$$-\mu: \frac{0.8+0.3+0.6+1.1+0.4+1.4}{6} = 0.77$$

$$-\sigma^2: \frac{[(0.8-0.77)^2 + (0.3-0.8)^2 + (0.6-0.8)^2 + (1.1-0.8)^2 + (0.4-0.8)^2 + (1.4-0.8)^2]}{5} = 0.18$$

$$-\sigma: 0.42$$

$$-p(x_i|No) = \frac{1}{\sqrt{2\pi} \times 0.42} \exp \left\{ -\frac{(x_i - 0.77)^2}{2 \times 0.18} \right\}$$

– For example

$$p(0.5M|No) = \frac{1}{\sqrt{2\pi} \times 0.42} \exp \left\{ -\frac{(0.5 - 0.77)^2}{2 \times 0.18} \right\} = 0.78$$

Sales	Buy?
0.8M	No
0.3M	No
0.6M	No
1.1M	No
0.4M	No
1.4M	No

Activity: continuous value

- What is the following value?

$$p(0.5M|Yes) = ?$$

Sales	Buy?
0.9M	Yes
1.0M	Yes
1.2M	Yes
1.3M	Yes

$$p(x_i|y) = \frac{1}{\sqrt{2\pi}\sigma} \exp \left\{ -\frac{(x_i - \mu)^2}{2\sigma^2} \right\}$$

Continuous value: example

- Data of class Yes:

- $\mu: (0.9+1.0+1.2+1.3)/4 = 1.1$

- $\sigma^2: \frac{[(0.9-1.1)^2 + (1.0-1.1)^2 + (1.2-1.1)^2 + (1.3-1.1)^2]}{3} = 0.0333$

- $\sigma: 0.1826$

- $p(x_i|\text{Yes}) = \frac{1}{\sqrt{2\pi} \times 0.1826} \exp \left\{ -\frac{(x_i-1.1)^2}{2 \times 0.0333} \right\}$

- For example

$$p(0.5M|\text{Yes}) = \frac{1}{\sqrt{2\pi} \times 0.1826} \exp \left\{ -\frac{(0.5-1.1)^2}{2 \times 0.0333} \right\} = 0.0099$$

Sales	Buy?
0.9M	Yes
1.0M	Yes
1.2M	Yes
1.3M	Yes

Summarize of NB classifier

- Training:
 - Calculate all $p(y)$
 - Calculate all possible $p(x_i|y)$ (assume continuous values follows univariate gaussian distribution)
- Prediction:
 - If $\mathbf{x} = [x_1, x_2, \dots, x_d]$, then calculate for all possible y
$$p(y) \times p(x_1|y) \times p(x_2|y) \times \cdots \times p(x_d|y) = p(y) \times p(\mathbf{x}|y)$$
 - Predict \mathbf{x} as C_j , where $p(C_j)p(\mathbf{x}|C_j)$ is the **maximum** value among $\{p(C_1)p(\mathbf{x}|C_1), p(C_2)p(\mathbf{x}|C_2), \dots, p(C_m)p(\mathbf{x}|C_m)\}$

Beyond NB: other ways to estimate $p(x|y)$

- Maximum likelihood estimation (MLE)
- Semi-naïve Bayes classifier
- Bayesian network (or called belief network)
- Expectation–maximization algorithm (EM algorithm) to consider latent variables

Outline

- Naïve Bayes classifier
- k -Nearest Neighbor classifier
- Summary

Basic idea

- Find from historical record those as similar as possible to the new record



Niffler



Nearest Neighbor Classifier

Requires three input:

- The set of **training data**
- **Distance metric** to compute distances between data points
- The value of **k** , the number of nearest neighbours

Nearest Neighbor Classifier: Prediction

For a test example

- Compute its distance to all training samples
- Identify k nearest neighbours
- Use labels of nearest neighbours to determine the label of test example (e.g., by taking majority vote—the majority class among the nearest neighbours)

k -Nearest Neighbour (k -NN for short) classifier is a lazy learning approach (or called lazy learner), because the training data is just memorized, without any model built in the training phase.

Nearest Neighbor Classifier: example

Customer	Age	Income (K)	Games played	Buy?
John	35	35	3	Yes
Rachel	22	50	2	No
Ruth	63	200	1	No
Tom	59	170	1	No
Neil	25	40	4	Yes
David	37	50	2	?

Euclidean Distance from David
$\sqrt{(35 - 37)^2 + (35 - 50)^2 + (3 - 2)^2} = 15.16$
$\sqrt{(22 - 37)^2 + (50 - 50)^2 + (2 - 2)^2} = 15$
$\sqrt{(63 - 37)^2 + (200 - 50)^2 + (1 - 2)^2} = 152.23$
$\sqrt{(59 - 37)^2 + (170 - 50)^2 + (1 - 2)^2} = 122$
$\sqrt{(25 - 37)^2 + (40 - 50)^2 + (4 - 2)^2} = 15.74$

If k = 1, the nearest neighbour is Rachel, and we predict it as “No”
If k = 3, the nearest neighbours are Rachel, John, and Neil, we predict it as “Yes”
If k = 5, the nearest neighbours are Rachel, John, Neil, Tom and Ruth. We predict it as “No”

Data normalization

- Attribute values may have to be scaled to prevent distance measures from being dominated by one of the attributes
- Example:
 - Age of a person may vary from 0 to 120
 - Income of a person may vary from \$10K to \$1M
 - Number of games played of a person may vary from 0 to 1000

Normalization: example of standardization

Customer	Age	Income (K)	Games played	Buy?
John	35	35	3	Yes
Rachel	22	50	2	No
Ruth	63	200	1	No
Tom	59	170	1	No
Neil	25	40	4	Yes
David	37	50	2	?

Standardization:

$$x' = \frac{x - \mu}{\sigma}$$

1. Calculate the mean of each feature of the training data
2. Calculate the variance of each feature of the training data
3. Normalize the training data feature by feature
4. When test data come, normalize the test data in the same way the training data is normalized, and do prediction

Normalization: example of standardization

Customer	Age	Income (K)	Games played	Buy?	Customer	Age	Income (K)	Games played	Buy?
John	35	35	3	Yes	John	-0.30	-0.81	0.61	Yes
Rachel	22	50	2	No	Rachel	-0.98	-0.62	-0.15	No
Ruth	63	200	1	No	Ruth	1.16	1.27	-0.92	No
Tom	59	170	1	No	Tom	0.95	0.89	-0.92	No
Neil	25	40	4	Yes	Neil	-0.83	-0.74	1.38	Yes
David	37	50	2	?	David	-0.20	-0.62	-0.15	?

Normalization: example of standardization

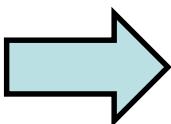
Customer	Age	Income (K)	Games played	Buy?
John	-0.30	-0.81	0.61	Yes
Rachel	-0.98	-0.62	-0.15	No
Ruth	1.16	1.27	-0.92	No
Tom	0.95	0.89	-0.92	No
Neil	-0.83	-0.74	1.38	Yes
David	-0.20	-0.62	-0.15	?

Euclidean Distance from David
0.7897
0.7800
2.4525
2.0483
1.6590

If $k = 1$, the nearest neighbour is Rachel, and we predict it as “No”
If $k = 3$, the nearest neighbours are Rachel, John, and Neil, we predict it as “Yes”
If $k = 5$, the nearest neighbours are Rachel, John, Neil, Tom and Ruth. We predict it as “No”

Non-numerical features

TID	Genre	Platform	Sales	Buy?
1	Action	Xbox	0.8M	No
2	Action	Switch	0.3M	No
3	Action	PS4	0.6M	No
4	Adventure	Switch	0.9M	Yes
5	Adventure	Switch	1.0M	Yes
6	Adventure	Switch	1.1M	No
7	Adventure	Switch	0.4M	No
8	Adventure	PS4	1.2M	Yes
9	Adventure	PS4	1.3M	Yes
10	Adventure	PS4	1.4M	No



TID	Genre	Xbox	Switch	PS4	Sales	Buy?
1	1	1	0	0	0.8M	No
2	1	0	1	0	0.3M	No
3	1	0	0	1	0.6M	No
4	0	0	1	0	0.9M	Yes
5	0	0	1	0	1.0M	Yes
6	0	0	1	0	1.1M	No
7	0	0	1	0	0.4M	No
8	0	0	0	1	1.2M	Yes
9	0	0	0	1	1.3M	Yes
10	0	0	0	1	1.4M	No

- Followed by the k -NN method.
- Normalization could also be applied if necessary.

Summarize k -NN

- Transfer non-numerical features into numerical
- Normalize the training data if necessary
- Specify k and distance metric, and store all training data
- For a test example, do normalization if necessary; then calculate its distance to all training instances
- Take all labels of the k -nearest neighbours in the training data, and then make prediction by majority voting (or other voting methods)

Determine k

Choosing the value of k :

- If k is too small, the prediction can easily be “wrong” if there is noisy point (be sensitive to noisy points)
- If k is too large, neighbourhood may include too many points from other classes

If $k = 1$, the nearest neighbour is Rachel, and we predict it as “No”

If $k = 3$, the nearest neighbours are Rachel, John, and Neil, we predict it as “Yes”

If $k = 5$, the nearest neighbours are Rachel, John, Neil, Tom and Ruth. We predict is as “No”

It is preferred to choose odd numbers for k . Why?

Because we do majority voting—no ties using odd numbers for binary classification

Determine k by cross-validation: example of 3-CV

- Randomly partition the data into three folders:

	TID	Genre	Xbox	Switch	PS4	Sales	Buy?
Folder 1	1	1	1	0	0	0.8M	No
	2	1	0	1	0	0.3M	No
	3	1	0	0	1	0.6M	No
Folder 2	4	0	0	1	0	0.9M	Yes
	5	0	0	1	0	1.0M	Yes
	6	0	0	1	0	1.1M	No
Folder 3	7	0	0	1	0	0.4M	No
	8	0	0	0	1	1.2M	Yes
	9	0	0	0	1	1.3M	Yes
	10	0	0	0	1	1.4M	No

- Select an evaluation metric, e.g. accuracy

Activity: k -NN prediction

- Using folder 1 and folder 2 as training data, and folder 3 as test data
- Euclidean distances from training data and test data

TID	1	2	3	4	5	6
7	1.78	1.00	1.74	0.50	0.60	0.70
8	1.78	1.95	1.17	1.45	1.43	1.42
9	1.80	2.00	1.22	1.47	1.45	1.43
10	1.83	2.05	1.28	1.50	1.47	1.45

- Prediction

TID	$k = 1$	$k = 3$	$k = 5$	Ground Truth
7	4 Yes	4,5,6 Yes	4,5,6,2,3 No	No
8	?	?	?	Yes
9				Yes
10				No

TID	Buy?
1	No
2	No
3	No
4	Yes
5	Yes
6	No

Trial 1

- Using folder 1 and folder 2 as training data, and folder 3 as test data
- Euclidean distances from training data and test data

TID	1	2	3	4	5	6
7	1.78	1.00	1.74	0.50	0.60	0.70
8	1.78	1.95	1.17	1.45	1.43	1.42
9	1.80	2.00	1.22	1.47	1.45	1.43
10	1.83	2.05	1.28	1.50	1.47	1.45

- Prediction

TID	k = 1	k = 3	k = 5	Ground Truth
7	4 Yes	4,5,6 Yes	4,5,6,2,3 No	No
8	3 No	3,6,5 No	3,6,5,4,1 No	Yes
9	3 No	3,6,5 No	3,6,5,4,1 No	Yes
10	3 No	3,6,5 No	3,6,5,4,1 No	No

TID	Buy?
1	No
2	No
3	No
4	Yes
5	Yes
6	No

Accuracy for k = 1: 25%
Accuracy for k = 3: 25%
Accuracy for k = 5: 50%

Trial 2

- Using folder 1 and folder 3 as training data, and folder 2 as test data
- Euclidean distances from training data and test data

TID	1	2	3	7	8	9	10
4	1.73	1.17	1.76	0.50	1.45	1.47	1.50
5	1.74	1.22	1.78	0.60	1.43	1.45	1.47
6	1.76	1.28	1.80	0.70	1.42	1.43	1.45

- Prediction

TID	k = 1	k = 3	k = 5	Ground Truth
4	7 No	7,2,8 No	7,2,8,9,10 No	Yes
5	7 No	7,2,8 No	7,2,8,9,10 No	Yes
6	7 No	7,2,8 No	7,2,8,9,10 No	No

Accuracy for k = 1: 33%

Accuracy for k = 3: 33%

Accuracy for k = 5: 33%

Trial 3

- Using folder 2 and folder 3 as training data, and folder 1 as test data
- Euclidean distances from training data and test data

TID	4	5	6	7	8	9	10
1	1.73	1.74	1.76	1.78	1.78	1.80	1.83
2	1.17	1.22	1.28	1.00	1.95	2.00	2.05
3	1.76	1.78	1.80	1.74	1.17	1.22	1.28

- Prediction

TID	k = 1	k = 3	k = 5	Ground Truth
1	4 Yes	4,5,6 Yes	4,5,6,7,8 Yes	No
2	7 No	7,4,5 Yes	7,4,5,6,8 Yes	No
3	8 Yes	8,9,10 Yes	8,9,10,7,4 Yes	No

Accuracy for k = 1: 33%

Accuracy for k = 3: 0%

Accuracy for k = 5: 0%

- Summarize the results of 3-CV

	K = 1	K = 3	K = 5
Trial 1	25%	25%	50%
Trial 2	33%	33%	33%
Trial 3	33%	0%	0%
Average	30%	19%	28%

- We select $k = 1$ for it has the highest average accuracy
- If data is enough, we usually 10-CV or even 30-CV; to save time, we may use 3-CV or 5-CV; however, it still depends on the data

Evaluation by cross-validation: example of 3-CV

- Similar methods can be used to evaluate the classifier if we do not have any test data
- E.g., we use 3-CV to evaluate 1-NN classifier on the games-buying dataset, and the average accuracy of all 3 trials is 30%

Advanced k-NN method (1)

- Avoid having to compute distance to all objects in the training set
 - Multi-dimensional access methods (k-d trees)
 - Fast approximate similarity search
 - Locality Sensitive Hashing (LSH)
- Condensing
 - Determine a smaller set of objects that give the same performance
- Editing
 - Remove objects to improve classification accuracy

Advanced k-NN method (2)

- Assign a weight to each individual attribute
- Assign a weight to each training instance

Outline

- Naïve Bayes classifier
- k -Nearest Neighbor classifier
- Summary

Three classification models

- Decision Tree
- Naïve Bayes classifier
- k -Nearest Neighbor classifier

Advantages of decision tree

- Relatively inexpensive to construct
- Relatively fast at classifying unknown records
- Easy to interpret for small-sized trees
- Robust to noise (especially when methods to avoid overfitting are employed)
- Can handle redundant or irrelevant attributes

Disadvantages of decision tree

- Due to the greedy nature of splitting criterion, interacting attributes (that can distinguish between classes together but not individually) may be passed over in favour of other attributed that are less discriminating.
- Each decision boundary involves only a single attribute

Advantages of Naïve Bayes

- Efficient and easy to construct
- The building process is parallelized
- Relatively fast at classifying unknown records
- Robust to isolated noise points
- Robust to irrelevant attributes

Disadvantages of Naïve Bayes

- The oversimplified assumption of features are conditionally independent with each other
 - Cannot handle redundant features
- Based on numerical computation, no rules generated – less interpretable

Advantages of k-NN

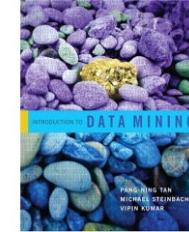
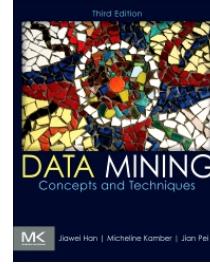
- Very simple implementation
- Few assumption on the distribution of the data
- Classifier can be updated online at very little cost as new instances with unknown classes are presented

Disadvantages of k-NN

- Expensive testing of each instance
- Sensitiveness to noisy or irrelevant attributes
- Sensitiveness to imbalanced datasets, where most entities belong to one or a few classes

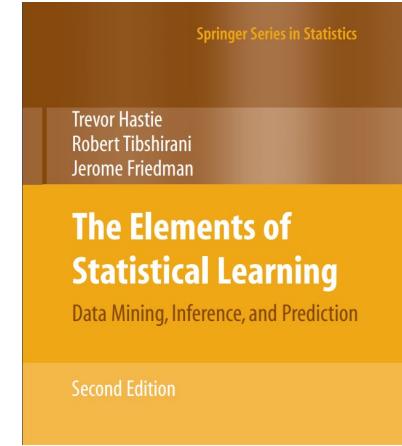
Recommended reading

- [\[Han et al., 2012\]](#)
 - Chapter 8.3
- [\[Tan et al., 2005\]](#)
 - Chapter 5.2, 5.3
- [\[Aggarwal, 2015\]](#)
 - Chapter 10.5, 10.8
- All other references (links) within the slides



Other classification models*

- Logistic regression
 - [Chapter 4.4](#)
- Support vector machine
 - [An introduction to Support Vector Machine](#) by John Shawe-Taylor
- Neural network
 - [Deep Learning](#) by Ian Goodfellow, Yoshua Bengio and Aaron Courville
- ...



Next two weeks

- Mining Text, Web, Stream, Time-series and Graph data
- Final-term revision