



# Lecture 16: Learning 5

[http://cs.nju.edu.cn/yuy/course\\_ai16.ashx](http://cs.nju.edu.cn/yuy/course_ai16.ashx)



# Previously...



## Learning

Decision tree learning

Neural networks

Why we can learn

Linear models

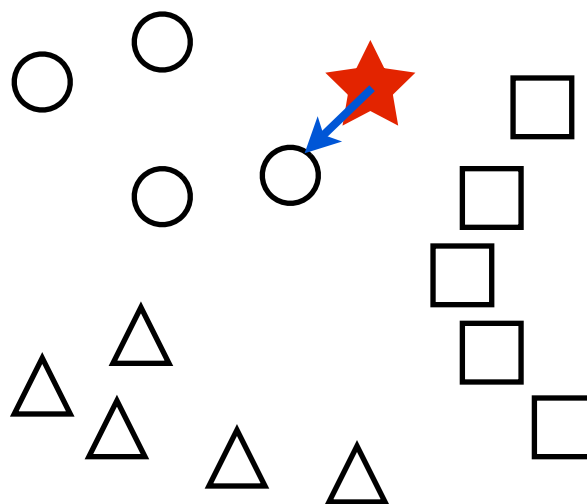


# Nearest Neighbor Classifier

# Nearest neighbor



what looks similar are similar

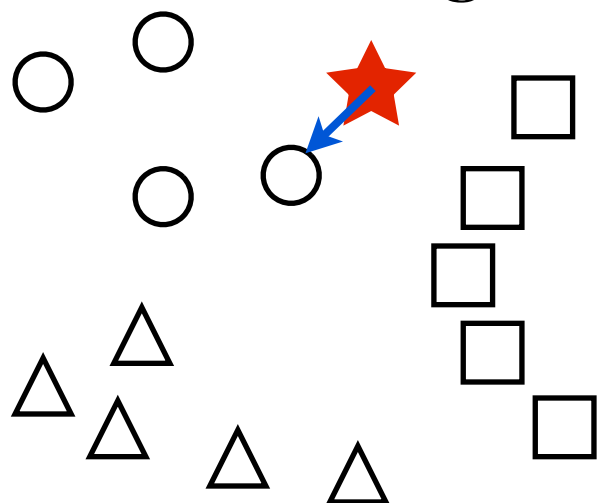


# Nearest neighbor

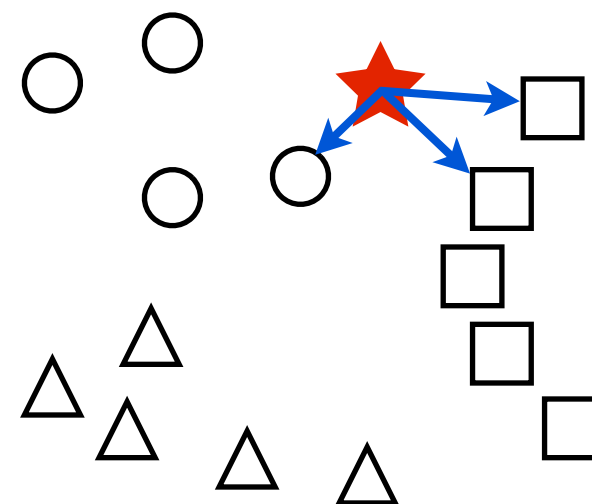


for classification:

1-nearest neighbor:



$k$ -nearest neighbor:



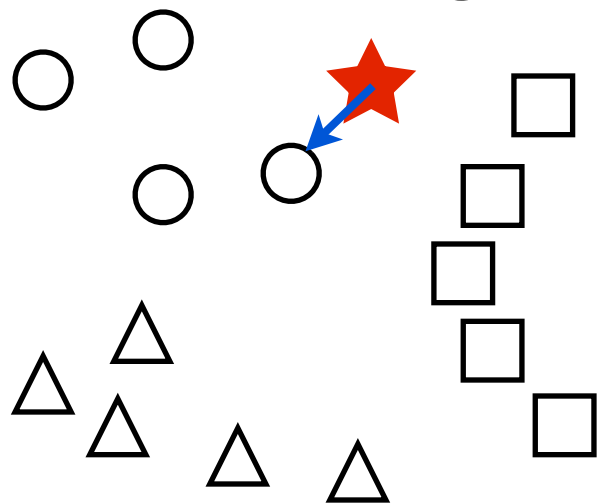
Predict the label as that of the NN  
or the (weighted) majority of the  $k$ -NN

# Nearest neighbor

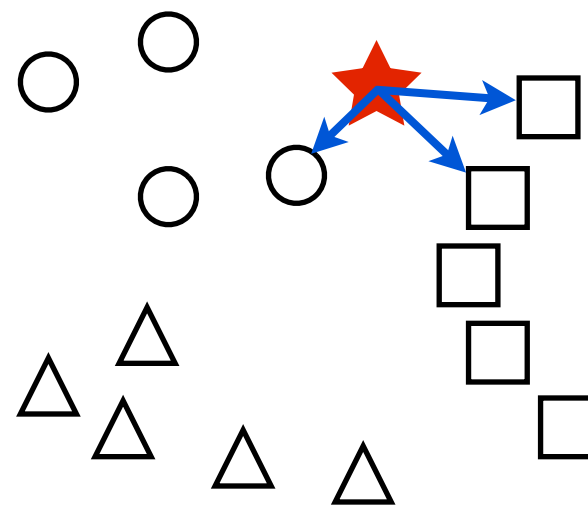


for regression:

1-nearest neighbor:



$k$ -nearest neighbor:

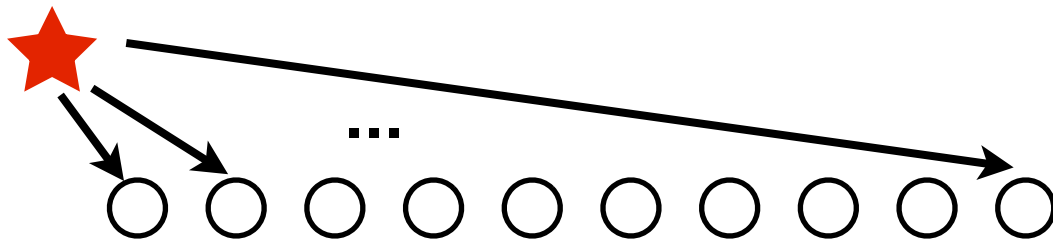


Predict the label as that of the NN  
or the (weighted) *average* of the  $k$ -NN

# Search for the nearest neighbor



Linear search



$n$  times of distance calculations

$$O(dn \ln k)$$

$d$  is the dimension,  $n$  is the number of samples

# Nearest neighbor classifier



- ▶ as classifier, asymptotically less than 2 times of the optimal Bayes error
- ▶ naturally handle multi-class
- ▶ no training time
- ▶ nonlinear decision boundary
  
- ▶ slow testing speed for a large training data set
- ▶ have to store the training data
- ▶ sensitive to similarity function

nonparametric method





# Naive Bayes Classifier

# Bayes rule



classification using posterior probability

for binary classification

$$f(\mathbf{x}) = \begin{cases} +1, & P(y = +1 | \mathbf{x}) > P(y = -1 | \mathbf{x}) \\ -1, & P(y = +1 | \mathbf{x}) < P(y = -1 | \mathbf{x}) \\ \text{random,} & \textit{otherwise} \end{cases}$$

in general

$$f(\mathbf{x}) = \arg \max_y P(y | \mathbf{x})$$

# Bayes rule



classification using posterior probability

for binary classification

$$f(\mathbf{x}) = \begin{cases} +1, & P(y = +1 | \mathbf{x}) > P(y = -1 | \mathbf{x}) \\ -1, & P(y = +1 | \mathbf{x}) < P(y = -1 | \mathbf{x}) \\ \text{random,} & \textit{otherwise} \end{cases}$$

in general

$$\begin{aligned} f(\mathbf{x}) &= \arg \max_y P(y | \mathbf{x}) \\ &= \arg \max_y P(\mathbf{x} | y)P(y)/P(\mathbf{x}) \\ &= \arg \max_y P(\mathbf{x} | y)P(y) \end{aligned}$$

how the  
probabilities be  
estimated

# Naive Bayes

$$f(x) = \arg \max_y P(\mathbf{x} | y)P(y)$$

estimation the a priori by frequency:

$$P(y) \leftarrow \tilde{P}(y) = \frac{1}{m} \sum_i I(y_i = y)$$



# Consider a very simple case



color ←



→ taste ?

id	color	taste
1	red	sweet
2	red	sweet
3	half-red	not-sweet
4	not-red	not-sweet
5	not-red	not-sweet
6	half-red	not-sweet
7	red	sweet
8	not-red	not-sweet
9	not-red	not-sweet
10	half-red	not-sweet
11	red	sweet
12	half-red	not-sweet
13	not-red	not-sweet

$$P(\text{red} \mid \text{sweet}) = 1$$

$$P(\text{half-red} \mid \text{sweet}) = 0$$

$$P(\text{not-red} \mid \text{sweet}) = 0$$

$$P(\text{sweet}) = 4/13$$

$$P(\text{red} \mid \text{not-sweet}) = 0$$

$$P(\text{half-red} \mid \text{not-sweet}) = 4/9$$

$$P(\text{not-red} \mid \text{not-sweet}) = 5/9$$

$$P(\text{not-sweet}) = 9/13$$

# Consider a very simple case



id	color	taste
1	red	sweet
2	red	sweet
3	half-red	not-sweet
4	not-red	not-sweet
5	not-red	not-sweet
6	half-red	not-sweet
7	red	sweet
8	not-red	not-sweet
9	not-red	not-sweet
10	half-red	not-sweet
11	red	sweet
12	half-red	not-sweet
13	not-red	not-sweet

what the  $f'$  would be?

$$f(x) = \arg \max_y P(x | y)P(y)$$

# Consider a very simple case



id	color	taste
1	red	sweet
2	red	sweet
3	half-red	not-sweet
4	not-red	not-sweet
5	not-red	not-sweet
6	half-red	not-sweet
7	red	sweet
8	not-red	not-sweet
9	not-red	not-sweet
10	half-red	not-sweet
11	red	sweet
12	half-red	not-sweet
13	not-red	not-sweet

what the  $f'$  would be?

$$f(x) = \arg \max_y P(x | y)P(y)$$

$$P(\text{red} | \text{sweet})P(\text{sweet}) = 4/13$$

$$P(\text{red} | \text{not-sweet})P(\text{not-sweet}) = 0$$

# Consider a very simple case



id	color	taste
1	red	sweet
2	red	sweet
3	half-red	not-sweet
4	not-red	not-sweet
5	not-red	not-sweet
6	half-red	not-sweet
7	red	sweet
8	not-red	not-sweet
9	not-red	not-sweet
10	half-red	not-sweet
11	red	sweet
12	half-red	not-sweet
13	not-red	not-sweet

what the  $f'$  would be?

$$f(x) = \arg \max_y P(x | y)P(y)$$

$$P(\text{red} | \text{sweet})P(\text{sweet}) = 4/13$$

$$P(\text{red} | \text{not-sweet})P(\text{not-sweet}) = 0$$

$$P(\text{half-red} | \text{sweet})P(\text{sweet}) = 0$$

$$P(\text{half-red} | \text{not-sweet})P(\text{not-sweet}) = \frac{4}{9} \times \frac{9}{13} = \frac{4}{13}$$



# Consider a very simple case



id	color	taste
1	red	sweet
2	red	sweet
3	half-red	not-sweet
4	not-red	not-sweet
5	not-red	not-sweet
6	half-red	not-sweet
7	red	sweet
8	not-red	not-sweet
9	not-red	not-sweet
10	half-red	not-sweet
11	red	sweet
12	half-red	not-sweet
13	not-red	not-sweet

what the  $f'$  would be?

$$f(x) = \arg \max_y P(x | y)P(y)$$

$$P(\text{red} | \text{sweet})P(\text{sweet}) = 4/13$$

$$P(\text{red} | \text{not-sweet})P(\text{not-sweet}) = 0$$

$$P(\text{half-red} | \text{sweet})P(\text{sweet}) = 0$$

$$P(\text{half-red} | \text{not-sweet})P(\text{not-sweet}) = \frac{4}{9} \times \frac{9}{13} = \frac{4}{13}$$

*perfect  
but not realistic*



# Naive Bayes

$$f(x) = \arg \max_y P(\mathbf{x} | y)P(y)$$

estimation the a priori by frequency:

$$P(y) \leftarrow \tilde{P}(y) = \frac{1}{m} \sum_i I(y_i = y)$$

assume features are conditional independence given the class (**naive assumption**):

$$\begin{aligned} P(\mathbf{x} | y) &= P(x_1, x_2, \dots, x_n | y) \\ &= P(x_1 | y) \cdot P(x_2 | y) \cdot \dots \cdot P(x_n | y) \end{aligned}$$

decision function:

$$f(x) = \arg \max_y \tilde{P}(y) \prod_i \tilde{P}(x_i | y)$$

# Naive Bayes



color={0,1,2,3} weight={0,1,2,3,4}

color	weight	sweet?
3	4	yes
2	3	yes
0	3	no
3	2	no
1	4	no

$$P(y = \text{yes}) = 2/5$$

$$P(y = \text{no}) = 3/5$$

$$P(\text{color} = 3 \mid y = \text{yes}) = 1/2$$

...

# Naive Bayes



color={0,1,2,3} weight={0,1,2,3,4}

color	weight	sweet?
3	4	yes
2	3	yes
0	3	no
3	2	no
1	4	no

$$P(y = \text{yes}) = 2/5$$

$$P(y = \text{no}) = 3/5$$

$$P(\text{color} = 3 \mid y = \text{yes}) = 1/2$$

...

$f(y \mid \text{color} = 3, \text{weight} = 3) \rightarrow$

# Naive Bayes



color={0,1,2,3} weight={0,1,2,3,4}

color	weight	sweet?
3	4	yes
2	3	yes
0	3	no
3	2	no
1	4	no

$$P(y = \text{yes}) = 2/5$$

$$P(y = \text{no}) = 3/5$$

$$P(\text{color} = 3 \mid y = \text{yes}) = 1/2$$

...

$f(y \mid \text{color} = 3, \text{weight} = 3) \rightarrow$

$$P(\text{color} = 3 \mid y = \text{yes})P(\text{weight} = 3 \mid y = \text{yes})P(y = \text{yes}) = 0.5 \times 0.5 \times 0.4 = 0.1$$

$$P(\text{color} = 3 \mid y = \text{no})P(\text{weight} = 3 \mid y = \text{no})P(y = \text{no}) = 0.33 \times 0.33 \times 0.6 = 0.06$$

# Naive Bayes



color={0,1,2,3} weight={0,1,2,3,4}

color	weight	sweet?
3	4	yes
2	3	yes
0	3	no
3	2	no
1	4	no

$$P(y = \text{yes}) = 2/5$$

$$P(y = \text{no}) = 3/5$$

$$P(\text{color} = 3 \mid y = \text{yes}) = 1/2$$

...

$$f(y \mid \text{color} = 3, \text{weight} = 3) \rightarrow$$

$$P(\text{color} = 3 \mid y = \text{yes})P(\text{weight} = 3 \mid y = \text{yes})P(y = \text{yes}) = 0.5 \times 0.5 \times 0.4 = 0.1$$

$$P(\text{color} = 3 \mid y = \text{no})P(\text{weight} = 3 \mid y = \text{no})P(y = \text{no}) = 0.33 \times 0.33 \times 0.6 = 0.06$$

$$f(y \mid \text{color} = 0, \text{weight} = 1) \rightarrow$$

# Naive Bayes



color={0,1,2,3} weight={0,1,2,3,4}

color	weight	sweet?
3	4	yes
2	3	yes
0	3	no
3	2	no
1	4	no

$$P(y = \text{yes}) = 2/5$$

$$P(y = \text{no}) = 3/5$$

$$P(\text{color} = 3 \mid y = \text{yes}) = 1/2$$

...

$$f(y \mid \text{color} = 3, \text{weight} = 3) \rightarrow$$

$$P(\text{color} = 3 \mid y = \text{yes})P(\text{weight} = 3 \mid y = \text{yes})P(y = \text{yes}) = 0.5 \times 0.5 \times 0.4 = 0.1$$

$$P(\text{color} = 3 \mid y = \text{no})P(\text{weight} = 3 \mid y = \text{no})P(y = \text{no}) = 0.33 \times 0.33 \times 0.6 = 0.06$$

$$f(y \mid \text{color} = 0, \text{weight} = 1) \rightarrow$$

$$P(\text{color} = 0 \mid y = \text{yes})P(\text{weight} = 1 \mid y = \text{yes})P(y = \text{yes}) = 0$$

$$P(\text{color} = 0 \mid y = \text{no})P(\text{weight} = 1 \mid y = \text{no})P(y = \text{no}) = 0$$

# Naive Bayes



color={0,1,2,3} weight={0,1,2,3,4}

color	weight	sweet?
3	4	yes
2	3	yes
0	3	no
3	2	no
1	4	no

+

color	sweet?
0	yes
1	yes
2	yes
3	yes

**smoothed (Laplacian correction) probabilities:**

$$P(\text{color} = 0 \mid y = \text{yes}) = (0 + 1) / (2 + 4)$$

$$P(y = \text{yes}) = (2 + 1) / (5 + 2)$$

for counting frequency,  
assume every event  
has happened once.

$$f(y \mid \text{color} = 0, \text{weight} = 1) \rightarrow$$

$$P(\text{color} = 0 \mid y = \text{yes})P(\text{weight} = 1 \mid y = \text{yes})P(y = \text{yes}) = \frac{1}{6} \times \frac{1}{7} \times \frac{3}{7} = 0.01$$

$$P(\text{color} = 0 \mid y = \text{no})P(\text{weight} = 1 \mid y = \text{no})P(y = \text{no}) = \frac{2}{7} \times \frac{1}{8} \times \frac{4}{7} = 0.02$$



# Naive Bayes



advantages:

very fast:

scan the data once, just count:  $O(mn)$

store class-conditional probabilities:  $O(n)$

test an instance:  $O(cn)$  ( $c$  the number of classes)

good accuracy in many cases

parameter free

output a probability

naturally handle multi-class

disadvantages:

# Naive Bayes



advantages:

very fast:

scan the data once, just count:  $O(mn)$

store class-conditional probabilities:  $O(n)$

test an instance:  $O(cn)$  ( $c$  the number of classes)

good accuracy in many cases

parameter free

output a probability

naturally handle multi-class

disadvantages:

the strong assumption may harm the accuracy

does not handle numerical features naturally