

Text classification and Naive Bayes (Chapter 13)

Definition 1 (Naive Bayes Classifier)

Naive Bayes (NB) Classifier assumes that the effect of the value of a predictor x on a given class c is class conditional independent. Bayes theorem provides a way of calculating the posterior probability $P(c|x)$ from class prior probability $P(c)$, predictor prior probability $P(x)$ and probability of the predictor given the class $P(x|c)$

$$P(c|x) = \frac{P(x|c)P(c)}{P(x)}$$

and for a vector of predictors $X = (x_1, \dots, x_n)$

$$P(c|X) = \frac{P(x_1|c) \dots P(x_n|c)P(c)}{P(x_1) \dots P(x_n)}.$$

The class with the highest posterior probability is the outcome of prediction.

Exercise 13/1

What is naive about Naive Bayes classifier? Briefly outline its major idea.

Answers can vary. For official definition refer to the Manning book.

Exercise 13/2

Considering the table of observations, use the Naive Bayes classifier to recommend whether to *Play Golf* given a day with *Outlook = Rainy*, *Temperature = Mild*, *Humidity = Normal* and *Windy = True*. Do not deal with the zero-frequency problem.

Outlook	Temperature	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

Table 1: Exercise.

First build the likelihood tables for each predictor

		Play Golf					Play Golf		
		Yes	No				Yes	No	
Outlook	Sunny	3/9	2/5	5/14	Temperature	Hot	2/9	2/5	4/14
	Overcast	4/9	0/5	4/14		Mild	4/9	2/5	6/14
	Rainy	2/9	3/5	5/14		Cool	3/9	1/5	4/14
		9/14	5/14				9/14	5/14	

		Play Golf					Play Golf		
		Yes	No				Yes	No	
Humidity	High	3/9	4/5	7/14	Windy	True	3/9	2/5	5/14
	Normal	6/9	1/5	7/14		False	6/9	3/5	9/14
		9/14	5/14				9/14	5/14	

We see that probability of *Sunny* given *Yes* is $3/9 = 0.33$, probability of *Sunny* is $5/14 = 0.36$ and probability of *Yes* is $9/14 = 0.64$. Then we count the likelihoods of *Yes* and *No*

$$\begin{aligned}
 P(\text{Yes}|\text{Rainy}, \text{Mild}, \text{Normal}, \text{True}) &\propto \\
 &= P(\text{Rainy}|\text{Yes}) \cdot P(\text{Mild}|\text{Yes}) \cdot P(\text{Normal}|\text{Yes}) \cdot P(\text{True}|\text{Yes}) \cdot P(\text{Yes}) \\
 &= \frac{2}{9} \cdot \frac{4}{9} \cdot \frac{6}{9} \cdot \frac{3}{9} \cdot \frac{9}{14} = 0.014109347 \\
 P(\text{No}|\text{Rainy}, \text{Mild}, \text{Normal}, \text{True}) &\propto \\
 &= P(\text{Rainy}|\text{No}) \cdot P(\text{Mild}|\text{No}) \cdot P(\text{Normal}|\text{No}) \cdot P(\text{True}|\text{No}) \cdot P(\text{No}) \\
 &= \frac{3}{5} \cdot \frac{2}{5} \cdot \frac{1}{5} \cdot \frac{3}{5} \cdot \frac{5}{14} = 0.010285714
 \end{aligned} \tag{1}$$

and suggest *Yes*. We can normalize the likelihoods to obtain the % confidence:

$$P(\text{Yes}|\text{Rainy}, \text{Mild}, \text{Normal}, \text{True}) = \frac{0.014109347}{0.014109347 + 0.010285714} = 57.84\%$$

$$P(\text{No}|\text{Rainy}, \text{Mild}, \text{Normal}, \text{True}) = \frac{0.010285714}{0.014109347 + 0.010285714} = 42.16\%$$

Definition 2 (A Linear Classifier)

Our linear classifier finds the hyperplane that bisects and is perpendicular to the connecting line of the closest points from the two classes. The separating (decision) hyperplane is defined in terms of a normal (weight) vector \mathbf{w} and a scalar intercept term b as

$$f(x) = \mathbf{w} \cdot \mathbf{x} + b$$

where \cdot is the dot product of vectors. Finally, the classifier becomes

$$\text{class}(x) = \text{sgn}(f(x)).$$

Exercise 13/3

Draw a sketch explaining the concept of our linear classifier. Include the equation of the separation hyperplane. Is our classifier equivalent to support vector machines (SVM)? What are limitations of our classifier?

Answers can vary. For official definition refer to the Manning book.

Exercise 13/4

Build a linear classifier for the training set $\{([1, 1], -1), ([2, 0], -1), ([2, 3], +1)\}$.

We first take the closest two points from the respective classes: $[1, 1]$ and $[2, 3]$. We have $\mathbf{w} = a \cdot ([1, 1] - [2, 3]) = [a, 2a]$. Now we calculate a and b

$$a + 2a + b = -1$$

$$2a + 6a + b = 1$$

for the points $[1, 1]$ and $[2, 3]$, respectively. The solution is

$$a = \frac{2}{5} \quad b = \frac{-11}{5}$$

building the weight vector

$$\mathbf{w} = \left[\frac{2}{5}, \frac{4}{5} \right]$$

and the final classifier becomes

$$\text{class}(x) = \text{sgn} \left(\frac{2}{5}x_1 + \frac{4}{5}x_2 - \frac{11}{5} \right).$$

Exercise 13/5

Explain the concept of classification based on neural networks. Draw a sketch and comment on all components.

Answers can vary. For official definition refer to the Manning book.

Exercise 13/6

What is the difference between supervised and unsupervised learning? Give examples.

Answers can vary. For official definition refer to the Manning book.

Flat clustering (Chapter 16)

Algorithm 1 K-means($\{\vec{x}_1, \dots, \vec{x}_N\}, K, \text{stopping criterion}$)

```

1:  $(\vec{s}_1, \dots, \vec{s}_K) \leftarrow \text{SelectRandomSeeds}(\{\vec{x}_1, \dots, \vec{x}_N\}, K)$ 
2: for  $k \leftarrow 1$  to  $K$  do
3:    $\vec{\mu}_k \leftarrow \vec{s}_k$ 
4: end for
5: repeat
6:   for  $k \leftarrow 1$  to  $K$  do
7:      $\omega_k \leftarrow \{\}$ 
8:   end for
9:   for  $n \leftarrow 1$  to  $N$  do
10:     $j \leftarrow \text{argmin}_{j'} |\vec{\mu}_{j'} - \vec{x}_n|$ 
11:     $\omega_j \leftarrow \omega_j \cup \{\vec{x}_n\}$  ▷ reassigning vectors
12:   end for
13:   for  $k \leftarrow 1$  to  $K$  do
14:     $\vec{\mu}_k \leftarrow \frac{1}{|\omega_k|} \sum_{\vec{x} \in \omega_k} \vec{x}$  ▷ recomputing centroids
15:   end for
16: until a stopping criterion has been met
17: return  $\{\vec{\mu}_1, \dots, \vec{\mu}_K\}$ 

```

Exercise 16/1

Use the K -means algorithm with Euclidean distance to cluster the following $N = 8$ examples into $K = 3$ clusters: $A_1 = (2, 10)$, $A_2 = (2, 5)$, $A_3 = (8, 4)$, $A_4 = (5, 8)$, $A_5 = (7, 5)$, $A_6 = (6, 4)$, $A_7 = (1, 2)$, $A_8 = (4, 9)$. Suppose that the initial seeds (centers of each cluster) are A_1 , A_4 and A_7 . Run the K -means algorithm for 3 epochs. After each epoch, draw a 10×10 space with all the 8 points and show the clusters with the new centroids.

$d(A, B)$ denotes the Euclidean distance between $A = (a_1, a_2)$ and $B = (b_1, b_2)$. It is calculated as $d(A, B) = \sqrt{(a_1 - b_1)^2 + (a_2 - b_2)^2}$.

Take seeds $\vec{s}_1 = A_1 = (2, 10)$, $\vec{s}_2 = A_4 = (5, 8)$, $\vec{s}_3 = A_7 = (1, 2)$.

By 1 we count the alignment for epoch 1: $A_1 \in \omega_1$, $A_2 \in \omega_3$, $A_3 \in \omega_2$, $A_4 \in \omega_2$, $A_5 \in \omega_2$, $A_6 \in \omega_2$, $A_7 \in \omega_3$, $A_8 \in \omega_2$; and we get the clusters: $\omega_1 = \{A_1\}$, $\omega_2 = \{A_3, A_4, A_5, A_6, A_8\}$, $\omega_3 = \{A_2, A_7\}$.

Centroids of the clusters: $\vec{\mu}_1 = (2, 10)$, $\vec{\mu}_2 = ((8+5+7+6+4)/5, (4+8+5+4+9)/5) = (6, 6)$, $\vec{\mu}_3 = ((2+1)/2, (5+2)/2) = (1.5, 3.5)$.

After epoch 2 the clusters are $\omega_1 = \{A_1, A_8\}$, $\omega_2 = \{A_3, A_4, A_5, A_6\}$, $\omega_3 = \{A_2, A_7\}$ with centroids $\vec{\mu}_1 = (3, 9.5)$, $\vec{\mu}_2 = (6.5, 5.25)$ and $\vec{\mu}_3 = (1.5, 3.5)$. And finally after epoch 3, the clusters are $\omega_1 = \{A_1, A_4, A_8\}$, $\omega_2 = \{A_3, A_5, A_6\}$, $\omega_3 = \{A_2, A_7\}$ with centroids $\vec{\mu}_1 = (3.66, 9)$, $\vec{\mu}_2 = (7, 4.33)$ and $\vec{\mu}_3 = (1.5, 3.5)$.

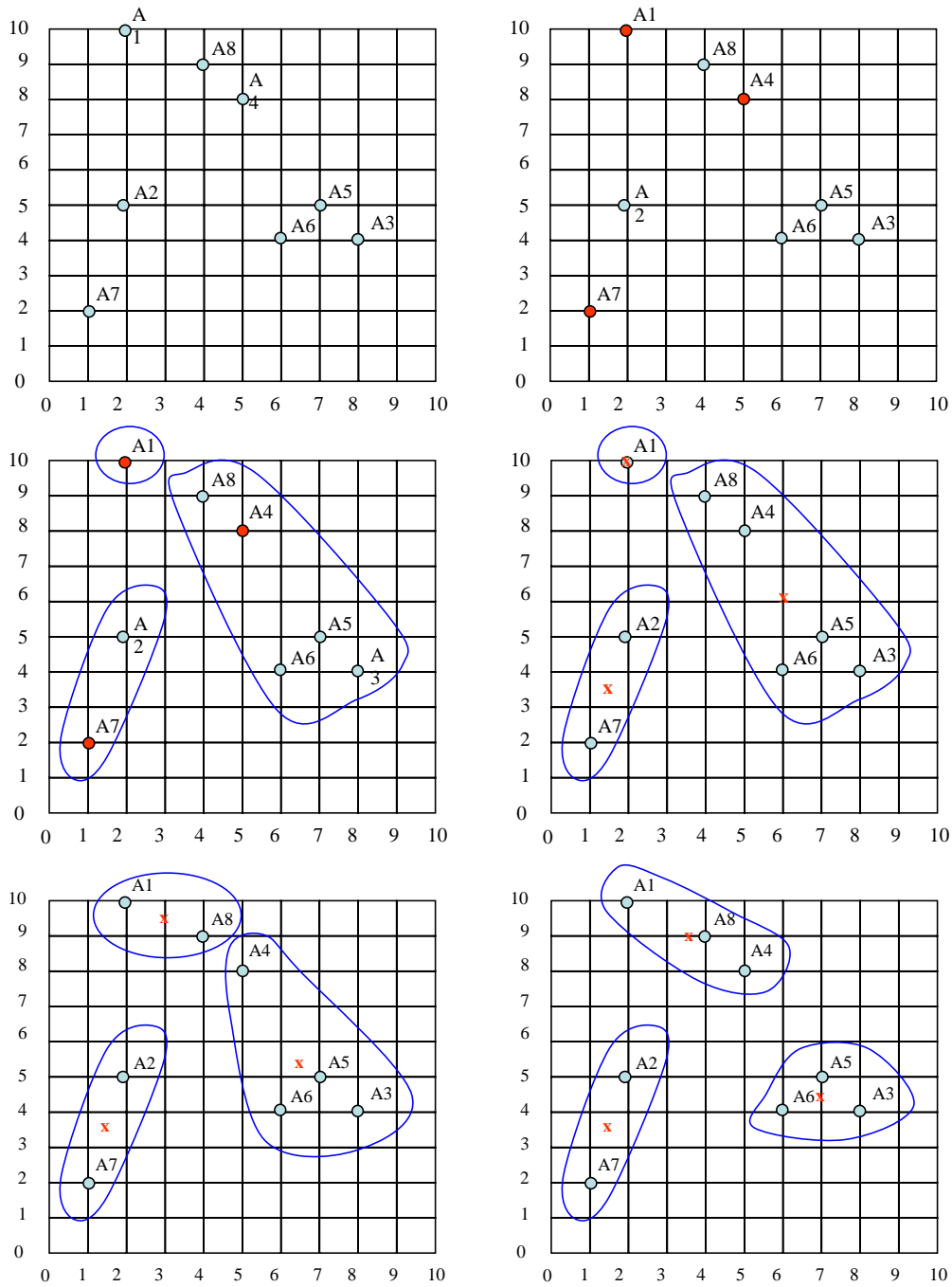


Figure 1: Visualization of K -means clustering algorithm.

Exercise 16/2

What makes a good clustering? Give some clustering evaluation metrics.

Answers can vary. For official definition refer to the Manning book.