

Bayesian Classification

Yufei Tao

Department of Computer Science and Engineering
Chinese University of Hong Kong

Review: Classification

Let A_1, \dots, A_d be d **attributes**.

Instance space: $\mathcal{X} = \text{dom}(A_1) \times \text{dom}(A_2) \times \dots \times \text{dom}(A_d)$ where $\text{dom}(A_i)$ represents the set of possible values on A_i .

Label space: $\mathcal{Y} = \{-1, 1\}$ (where -1 and 1 are **class labels**).

Instance-label pair (a.k.a. **object**): a pair (\mathbf{x}, y) in $\mathcal{X} \times \mathcal{Y}$.

- \mathbf{x} is a vector; we use $\mathbf{x}[A_i]$ to represent the vector's value on A_i ($1 \leq i \leq d$).

Denote by \mathcal{D} a probabilistic distribution over $\mathcal{X} \times \mathcal{Y}$.

Review: Classification

Goal: Given an object (\mathbf{x}, y) drawn from \mathcal{D} , we want to predict its label y from its attribute values $\mathbf{x}[A_1], \dots, \mathbf{x}[A_d]$.

Classifier (hypothesis): A function $h: \mathcal{X} \rightarrow \mathcal{Y}$.

Error of h on \mathcal{D} : $err_{\mathcal{D}}(h) = \Pr_{(\mathbf{x}, y) \sim \mathcal{D}}[h(\mathbf{x}) \neq y]$.

namely, if we draw an object (\mathbf{x}, y) according to \mathcal{D} , what is the probability that h mis-predicts the label?

We would like to learn a classifier h with small $err_{\mathcal{D}}(h)$ from a **training set S** where each object is drawn independently from \mathcal{D} .

The Ideal Classifier

Fix a point \mathbf{p} in the instance space. **Think:** given a class label $c \in \mathcal{Y}$, how would you interpret the conditional probability

$$\Pr_{(\mathbf{x}, y) \sim \mathcal{D}}[y = c \mid \mathbf{x} = \mathbf{p}]?$$

Design a classifier h_{opt} as follows:

- $h_{opt}(\mathbf{p}) = -1$ if $\Pr_{(\mathbf{x}, y) \sim \mathcal{D}}[y = -1 \mid \mathbf{x} = \mathbf{p}] \geq 0.5$;
- $h_{opt}(\mathbf{p}) = 1$ otherwise.

This is the best classifier possible.

Its error on \mathcal{D} , namely, $err_{\mathcal{D}}(h_{opt})$, is the **bayesian error**.

We will introduce the **Bayesian method**, which aims to follow the decisions of h_{opt} by approximating the value of $\Pr_{(x,y)\sim\mathcal{D}}[y = c \mid \mathbf{x} = \mathbf{p}]$.

Henceforth, we will abbreviate $\Pr_{(x,y)\sim\mathcal{D}}[y = c \mid \mathbf{x} = \mathbf{p}]$ simply as $\Pr[y = c \mid \mathbf{p}]$.

Example: Suppose that we have the following training set:

| age | education | occupation | loan default |
|------------|------------------|-------------------|---------------------|
| 28 | high school | self-employed | yes |
| 32 | master | programmer | no |
| 33 | undergrad | programmer | yes |
| 37 | undergrad | programmer | no |
| 38 | undergrad | self-employed | yes |
| 45 | master | self-employed | no |
| 48 | high school | programmer | no |
| 50 | master | lawyer | no |
| 52 | master | programmer | no |
| 55 | high school | self-employed | no |

Bayesian classification works most effectively when each attribute has a **small** domain, namely, the attribute has only a small number of possible values. When an attribute has a large domain, we may reduce its domain size through **discretization**.

For example, we may discretize the “age” attribute into a smaller domain: $\{20^+, 30^+, 40^+, 50^+\}$, where “20⁺” corresponds to the interval $[20, 29]$, “30⁺” to $[30, 39]$, and so on. See the next slide for the training set after the conversion.

Example: The training set after discretizing “age”:

| age | education | occupation | loan default |
|------------|------------------|-------------------|---------------------|
| 20+ | high school | self-employed | yes |
| 30+ | master | programmer | no |
| 30+ | undergrad | programmer | yes |
| 30+ | undergrad | programmer | no |
| 30+ | undergrad | self-employed | yes |
| 40+ | master | self-employed | no |
| 40+ | high school | programmer | no |
| 50+ | master | lawyer | no |
| 50+ | master | programmer | no |
| 50+ | high school | self-employed | no |

Bayes' Theorem:

$$\Pr[X | Y] = \frac{\Pr[Y | X] \cdot \Pr[X]}{\Pr[Y]}$$

Given an instance \mathbf{x} , (as in h_{opt}) we predict its label as -1 if and only if

$$\Pr[y = -1 | \mathbf{x}] \geq \Pr[y = 1 | \mathbf{x}].$$

Applying Bayes' theorem, we get:

$$\Pr[y = 1 | \mathbf{x}] = \frac{\Pr[\mathbf{x} | y = 1] \cdot \Pr[y = 1]}{\Pr[\mathbf{x}]}.$$

Similarly:

$$\Pr[y = -1 | \mathbf{x}] = \frac{\Pr[\mathbf{x} | y = -1] \cdot \Pr[y = -1]}{\Pr[\mathbf{x}]}.$$

It suffices to decide which of the following is larger:

- $\Pr[\mathbf{x} | y = 1] \cdot \Pr[y = 1]$, or
- $\Pr[\mathbf{x} | y = -1] \cdot \Pr[y = -1]$.

Bayesian classification **estimates** $\Pr[\mathbf{x} | y = 1] \cdot \Pr[y = 1]$ and $\Pr[\mathbf{x} | y = -1] \cdot \Pr[y = -1]$ using the training set. Next, we will explain only the former, because the estimate of the latter is similar.

The objective, obviously, is to estimate two terms:

- $\Pr[y = 1]$
- $\Pr[\mathbf{x} | y = 1]$

We will discuss each term in turn.

$Pr[y = 1]$

This is the probability for an object drawn from \mathcal{D} to have label 1.

Naturally, we estimate $Pr[y = 1]$ as the percentage of yes objects in the training set S .

Example: In Slide 8, $Pr[y = 1] = 0.3$.

$Pr[\mathbf{x} \mid y = 1]$

This is the probability for a “yes”-object drawn from \mathcal{D} to carry exactly the attribute values $\mathbf{x}[A_1], \dots, \mathbf{x}[A_d]$.

We could estimate $Pr[\mathbf{x} \mid y = 1]$ as the percentage of objects having attribute values $\mathbf{x}[A_1], \dots, \mathbf{x}[A_d]$ among all the yes objects in S . But this is a bad idea because S may have **very few (even none)** such objects, rendering the estimate unreliable (losing statistical significance).

This situation forces us to introduce assumptions which — if satisfied — would allow us to obtain a more reliable estimate of $Pr[\mathbf{x} \mid y = 1]$.

$Pr[x | y = 1]$ (cont.)

Bayesian classification makes an **assumption** here:

$$Pr[x | y = 1] = \prod_{i=1}^d Pr[x[A_i] | y = 1].$$

For each $i \in [1, d]$, we estimate $Pr[x[A_i] | y = 1]$ as the percentage of objects with attribute value $x[A_i]$ among all the yes objects in S .

Example: In Slide 8, $Pr[30+, \text{high-school}, \text{programmer} | y = 1]$ is assumed to be the product of

- $Pr[30+ | y = 1]$, which is estimated as $2/3$
- $Pr[\text{high-school} | y = 1]$, which is estimated as $1/3$
- $Pr[\text{programmer} | y = 1]$, which is estimated as $1/3$.

The product equals $2/27$.

$Pr[\mathbf{x} | y = 1]$ (cont.)

The estimate of $Pr[\mathbf{x}[A_i] | y = 1]$ would be 0 if S does not have any yes-object with attribute value $\mathbf{x}[A_i]$. But that would force our estimate of $Pr[\mathbf{x} | y = 1]$ to be 0. Instead, we replace the 0 estimate with a very small value, for example, 0.000001.

Example: In Slide 8, $Pr[\text{lawyer} | y = 1]$ is estimated as 0.000001.

Think: At the beginning, we said that Bayesian classification works better on small domains. Why?

The effectiveness of Bayesian classification relies on the accuracy of the assumption:

$$\Pr[\mathbf{x} \mid y = 1] = \prod_{i=1}^d \Pr[x[A_i] \mid y = 1].$$

This assumption is called the **conditional independence** assumption. When this assumption is seriously violated, the accuracy of the method drops significantly.

The approach we have discussed so far is known as **naive Bayes classification**.

The approach can be integrated with alternative (less stringent) conditional independence assumption. Consider the evaluation of

$$\Pr[30+, \text{undergrad}, \text{programmer} \mid y = -1]$$

in the context of Slide 8. Suppose that “age” and “education” are independent after fixing “occupation” and the class label. Then:

$$\begin{aligned} & \Pr[30+, \text{undergrad}, \text{programmer} \mid y = -1] \\ = & \Pr[30+, \text{undergrad} \mid \text{programmer}, y = -1] \cdot \\ & \Pr[\text{programmer} \mid y = -1] \\ = & \Pr[30+ \mid \text{programmer}, y = -1] \\ & \cdot \Pr[\text{undergrad} \mid \text{programmer}, y = -1] \\ & \cdot \Pr[\text{programmer} \mid y = -1] \\ = & \frac{2}{4} \cdot \frac{1}{4} \cdot \frac{4}{7} = 1/14. \end{aligned}$$

Next, we will provide an alternative way to describe the Bayes method (using naive Bayes as an example). Our description will clarify what is actually the set \mathcal{H} of classifiers to be learned from. This allows you to apply the generalization theorem (discussed in the previous lecture) to bound the generalization error of the classifier obtained.

Recall that we have attributes A_1, \dots, A_d .

We assume that each A_i ($i \in [1, d]$) has a finite domain $dom(A_i)$.

For each A_i , we introduce $2|dom(A_i)|$ parameters. Specifically, for each value $a \in dom(A_i)$, there are two parameters:

- $p_i(a | -1)$, which is our estimate of $\Pr[\mathbf{x}[A_i] = a | y = -1]$;
- $p_i(a | 1)$, which is our estimate of $\Pr[\mathbf{x}[A_i] = a | y = 1]$.

Furthermore, we also introduce:

- $p(-1)$, which is our estimate of $\Pr[y = -1]$;
- $p(1)$, which is our estimate of $\Pr[y = 1]$.

In total, we have $2 + 2 \sum_{i=1}^d |dom(A_i)|$ parameters.

Once the values of the $2 + 2 \sum_{i=1}^d |\text{dom}(A_i)|$ parameters have been fixed, the conditional independence assumption (of naive Bayes) gives rise to the following classifier $h(\mathbf{x})$:

- $h(\mathbf{x}) = -1$ if

$$p(-1) \cdot \prod_{i=1}^d p_i(\mathbf{x}[A_i] \mid -1) \geq p(1) \cdot \prod_{i=1}^d p_i(\mathbf{x}[A_i] \mid 1)$$

- $h(\mathbf{x}) = 1$ otherwise.

The set \mathcal{H} contains all such classifiers.

Remark: The Bayes method we explained earlier gives an efficient way for choosing a reasonably good classifier $h \in \mathcal{H}$.

The rest of the slides will not be tested.

Next, we introduce the **Bayesian network** which is a popular way to describe sophisticated conditional independence assumptions.

Let us review some concepts on acyclic directed graphs (DAG):

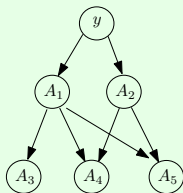
- A DAG G is a directed graph with no cycles.
- A node in G with 0 in-degree is a **root**. Note that G may have multiple roots.
- If a node u has an edge to another node v , then u is a **parent** of v . Note that a node can have multiple parents.

We will use $\text{parents}(v)$ to represent the set of parents of a node v .

We define a **Bayesian network** as an acyclic directed graph (DAG) G satisfying:

- 1 G has $d + 1$ nodes, including a node for the class label and a node for each attribute;
- 2 G has a single root node, which must be the class label;
- 3 if attribute u has no path to any of the attributes v_1, \dots, v_x (where $x \geq 1$ can be any integer), then u and (v_1, \dots, v_x) are independent conditioned on $parents(u)$.

Example: The following is a Bayesian network with $d = 5$.



- A_1 and A_2 independent conditioned on y ;
- A_4 and A_5 independent conditioned on A_1, A_2 ;
- A_3 and A_5 independent conditioned on A_1 ;
- A_3 and (y, A_5) independent conditioned on A_1 .

Theorem 1: Given the conditional independence assumptions described by a Bayesian network G , we have

$$\Pr[A_1, \dots, A_d \mid y] = \prod_{i=1}^d \Pr[A_i \mid \text{parents}(A_i)].$$

Before proving the theorem, let us first see an example.

Example: Given the Bayesian network on the previous slide, we have:

$$\Pr[A_1, A_2, \dots, A_5 \mid y] = \\ \Pr[A_1 \mid y] \cdot \Pr[A_2 \mid y] \cdot \Pr[A_3 \mid A_1] \cdot \Pr[A_4 \mid A_1, A_2] \cdot \Pr[A_5 \mid A_1, A_2].$$

We will now proceed to prove the theorem. The following facts about conditional independence will be useful:

Lemma 1: If A and B are independent conditioned on C , then:

- $\Pr[A, B | C] = \Pr[A | C] \cdot \Pr[B | C]$;
- $\Pr[A | C, B] = \Pr[A | C]$.

Proof: The first bullet is the definition of conditional independence, whereas the second bullet holds because

$$\begin{aligned}\Pr[A | C, B] &= \frac{\Pr[A, B | C]}{\Pr[B | C]} \\ &= \frac{\Pr[A | C] \Pr[B | C]}{\Pr[B | C]} = \Pr[A | C].\end{aligned}$$

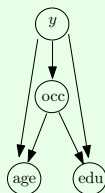
□

Proof of Theorem 1: Without loss of generality, suppose that y, A_1, \dots, A_d is a topological order of G .

$$\begin{aligned} \Pr[A_1, \dots, A_d | y] &= \Pr[A_2, \dots, A_d | y, A_1] \cdot \Pr[A_1 | y] \\ &= \Pr[A_3, \dots, A_d | y, A_1, A_2] \cdot \Pr[A_2 | y, A_1] \cdot \Pr[A_1 | y] \\ &\dots \\ &= \prod_{i=1}^d \Pr[A_i | y, A_1, \dots, A_{i-1}] \\ \text{(by Lemma 1)} &= \prod_{i=1}^d \Pr[A_i | \text{parents}(A_i)] \end{aligned}$$

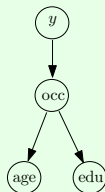
where the last equality used the conditional-independence properties implied by G and the fact that $\text{parents}(A_i) \subseteq \{y, A_1, \dots, A_{i-1}\}$. \square

Example: Consider the training set on Slide 8. If we are given the Bayesian network



then $\Pr[30+, \text{undergrad}, \text{programmer} \mid y = -1]$ is calculated as shown on Slide 17.

Example (cont.): If the Bayesian network is



then $\Pr[30+, \text{undergrad}, \text{programmer} \mid y = -1]$

$$\begin{aligned} &= \Pr[30+, \text{undergrad} \mid \text{programmer}, y = -1] \cdot \\ &\quad \Pr[\text{programmer} \mid y = -1] \\ &= \Pr[30+ \mid \text{programmer}] \cdot \Pr[\text{undergrad} \mid \text{programmer}] \\ &\quad \cdot \Pr[\text{programmer} \mid y = -1] \\ &= \frac{3}{5} \cdot \frac{2}{5} \cdot \frac{4}{7} = 24/175. \end{aligned}$$

Think: What is the set \mathcal{H} of classifiers to be learned from if we are given the Bayesian network on the previous slide?