# Naïve Bayes Classifier

## Ke Chen

# Outline

- Background

- Probability Basics

- Probabilistic Classification

- Naïve Bayes

  - Principle and Algorithms

  - Example: Play Tennis

- Relevant Issues

- Summary

# Background

- There are three methods to establish a classifier

  *a*) Model a classification rule directly

      Examples: k-NN, decision trees, perceptron, SVM

  *b*) Model the probability of class memberships given input data

      Example: perceptron with the cross-entropy cost

  *c*) Make a probabilistic model of data within each class

      Examples: naive Bayes, model based classifiers

- *a*) and *b*) are examples of discriminative classification
- *c*) is an example of generative classification
- *b*) and *c*) are both examples of probabilistic classification

# Probability Basics

- Prior, conditional and joint probability for random variables

  - Prior probability: $P(X)$

  - Conditional probability: $P(X_1 \mid X_2), P(X_2 \mid X_1)$

  - Joint probability: $\mathbf{X} = (X_1, X_2), P(\mathbf{X}) = P(X_1, X_2)$

  - Relationship: $P(X_1, X_2) = P(X_2 \mid X_1)P(X_1) = P(X_1 \mid X_2)P(X_2)$

  - Independence: $P(X_2 \mid X_1) = P(X_2), P(X_1 \mid X_2) = P(X_1), P(X_1, X_2) = P(X_1)P(X_2)$

- Bayesian Rule

$$P(C \mid \mathbf{X}) = \frac{P(\mathbf{X} \mid C)P(C)}{P(\mathbf{X})} \qquad Posterior = \frac{Likelihood \times Prior}{Evidence}$$

# Probability Basics

- <span style="color:red">Quiz</span>: We have two six-sided dice. When they are tolled, it could end up with the following occurance: (<span style="color:red">A</span>) dice 1 lands on side "3", (<span style="color:red">B</span>) dice 2 lands on side "1", and (<span style="color:red">C</span>) Two dice sum to eight. Answer the following questions:

1) $P(A) = ?$

2) $P(B) = ?$

3) $P(C) = ?$

4) $P(A \mid B) = ?$

5) $P(C \mid A) = ?$

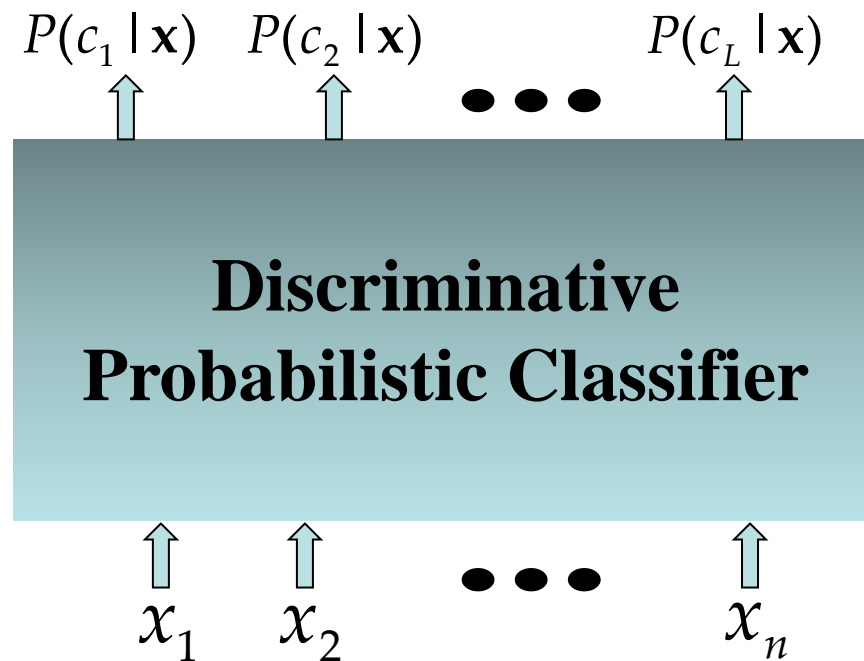6) $P(A, B) = ?$

7) $P(A, C) = ?$

8) Is $P(A, C)$ equal to $P(A) * P(C)$?

# Probabilistic Classification

- Establishing a probabilistic model for classification
  - **Discriminative model**

$$P(C \mid \mathbf{X}) \quad C = c_1, \cdots, c_L, \ \mathbf{X} = (X_1, \cdots, X_n)$$

$$P(c_1 \mid \mathbf{x}) \quad P(c_2 \mid \mathbf{x}) \quad \cdots \quad P(c_L \mid \mathbf{x})$$

**Discriminative Probabilistic Classifier**

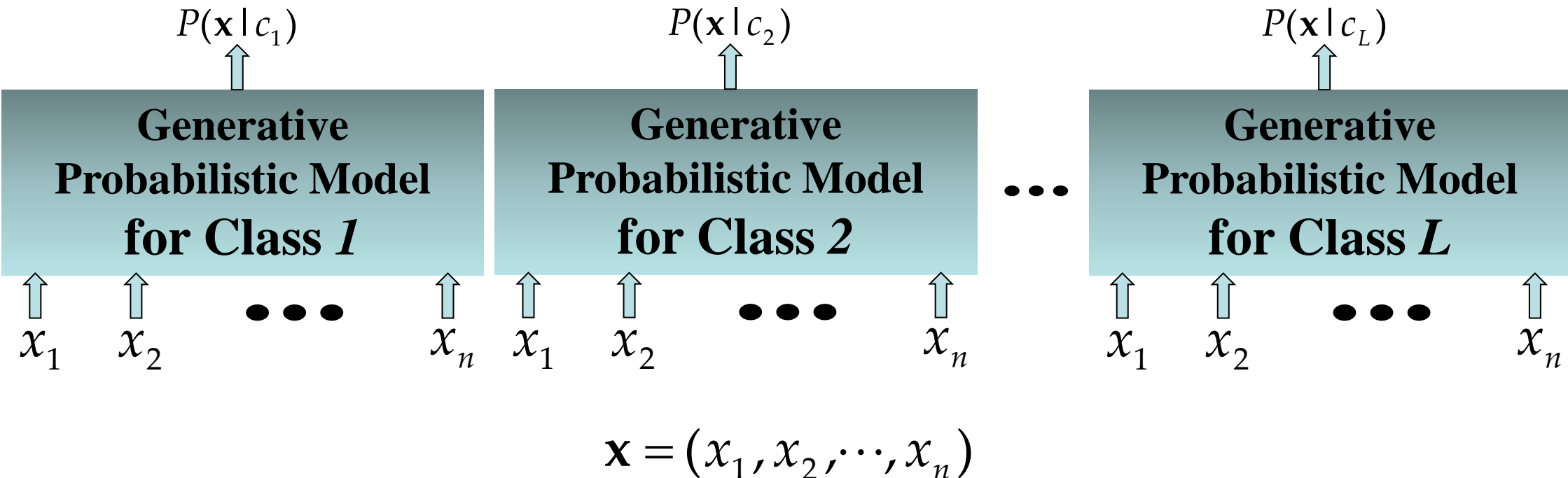$$x_1 \quad x_2 \quad \cdots \quad x_n$$

$$\mathbf{x} = (x_1, x_2, \cdots, x_n)$$

# Probabilistic Classification

- Establishing a probabilistic model for classification (cont.)
  - **Generative model**

$$P(\mathbf{X}|C) \quad C = c_1, \cdots, c_L, \ \mathbf{X} = (X_1, \cdots, X_n)$$

$P(\mathbf{x}|c_1)$ 　　　　　 $P(\mathbf{x}|c_2)$ 　　　　　 $P(\mathbf{x}|c_L)$

| **Generative Probabilistic Model for Class *1*** | **Generative Probabilistic Model for Class *2*** | $\cdots$ | **Generative Probabilistic Model for Class *L*** |
| --- | --- | --- | --- |

$x_1 \quad x_2 \quad \bullet\bullet\bullet \quad x_n$ 　 $x_1 \quad x_2 \quad \bullet\bullet\bullet \quad x_n$ 　 $x_1 \quad x_2 \quad \bullet\bullet\bullet \quad x_n$

$$\mathbf{x} = (x_1, x_2, \cdots, x_n)$$

# Probabilistic Classification

- MAP classification rule
  - **MAP**: **M**aximum **A** **P**osterior
  - Assign $x$ to $c^*$ if

$$P(C = c^* \mid \mathbf{X} = \mathbf{x}) > P(C = c \mid \mathbf{X} = \mathbf{x}) \quad c \neq c^*, \; c = c_1, \cdots, c_L$$

- Generative classification with the MAP rule
  - Apply Bayesian rule to convert them into posterior probabilities

$$P(C = c_i \mid \mathbf{X} = \mathbf{x}) = \frac{P(\mathbf{X} = \mathbf{x} \mid C = c_i)P(C = c_i)}{P(\mathbf{X} = \mathbf{x})}$$
$$\propto P(\mathbf{X} = \mathbf{x} \mid C = c_i)P(C = c_i)$$
$$\text{for } i = 1, 2, \cdots, L$$

  - Then apply the MAP rule

# Naïve Bayes

- Bayes classification

$$P(C \mid \mathbf{X}) \propto P(\mathbf{X} \mid C)P(C) = P(X_1, \cdots, X_n \mid C)P(C)$$

Difficulty: learning the joint probability $P(X_1, \cdots, X_n \mid C)$

- Naïve Bayes classification

  - Assumption that <span style="color:red">all input features are conditionally independent</span>!

$$P(X_1, X_2, \cdots, X_n \mid C) = P(X_1 \mid X_2, \cdots, X_n, C)P(X_2, \cdots, X_n \mid C)$$
$$= P(X_1 \mid C)P(X_2, \cdots, X_n \mid C)$$
$$= P(X_1 \mid C)P(X_2 \mid C) \cdots P(X_n \mid C)$$

  - MAP classification rule: for $\mathbf{x} = (x_1, x_2, \cdots, x_n)$

$$[P(x_1 \mid c^*) \cdots P(x_n \mid c^*)]P(c^*) > [P(x_1 \mid c) \cdots P(x_n \mid c)]P(c), \quad c \neq c^*, c = c_1, \cdots, c_L$$

# Naïve Bayes

- Algorithm: Discrete-Valued Features

  - Learning Phase: Given a training set $\mathbf{S}$ $\mathrm{of}$ $F$ features and $L$ classes,

    For each target value of $c_i$ $(c_i = c_1, \cdots, c_L)$

    $\hat{P}(C = c_i) \leftarrow$ estimate $P(C = c_i)$ with examples in $\mathbf{S}$;

    For every feature value $x_{jk}$ of each feature $X_j$ $(j = 1, \cdots, F; k = 1, \cdots, N_j)$

    $\hat{P}(X_j = x_{jk} \mid C = c_i) \leftarrow$ estimate $P(X_j = x_{jk} \mid C = c_i)$ with examples in $\mathbf{S}$;

    Output: $F * L$ conditional probabilistic (generative) models

  - Test Phase: Given an unknown instance $\mathbf{X}' = (a_1', \cdots, a_n')$

  "Look up tables" to assign the label $c^*$ to $\mathbf{X}'$ if

  $$[\hat{P}(a_1' \mid c^*) \cdots \hat{P}(a_n' \mid c^*)]\hat{P}(c^*) > [\hat{P}(a_1' \mid c) \cdots \hat{P}(a_n' \mid c)]\hat{P}(c), \quad c \neq c^*, c = c_1, \cdots, c_L$$

# Example

- Example: Play Tennis

*PlayTennis*: training examples

| Day | Outlook | Temperature | Humidity | Wind | PlayTennis |
|-----|---------|-------------|----------|------|------------|
| D1 | Sunny | Hot | High | Weak | No |
| D2 | Sunny | Hot | High | Strong | No |
| D3 | Overcast | Hot | High | Weak | Yes |
| D4 | Rain | Mild | High | Weak | Yes |
| D5 | Rain | Cool | Normal | Weak | Yes |
| D6 | Rain | Cool | Normal | Strong | No |
| D7 | Overcast | Cool | Normal | Strong | Yes |
| D8 | Sunny | Mild | High | Weak | No |
| D9 | Sunny | Cool | Normal | Weak | Yes |
| D10 | Rain | Mild | Normal | Weak | Yes |
| D11 | Sunny | Mild | Normal | Strong | Yes |
| D12 | Overcast | Mild | High | Strong | Yes |
| D13 | Overcast | Hot | Normal | Weak | Yes |
| D14 | Rain | Mild | High | Strong | No |

# Example

- Learning Phase

| Outlook | Play=*Yes* | Play=*No* |
|---------|------------|-----------|
| *Sunny* | 2/9 | 3/5 |
| *Overcast* | 4/9 | 0/5 |
| *Rain* | 3/9 | 2/5 |

| Temperature | Play=*Yes* | Play=*No* |
|-------------|------------|-----------|
| *Hot* | 2/9 | 2/5 |
| *Mild* | 4/9 | 2/5 |
| *Cool* | 3/9 | 1/5 |

| Humidity | Play=*Yes* | Play=*No* |
|----------|------------|-----------|
| *High* | 3/9 | 4/5 |
| *Normal* | 6/9 | 1/5 |

| Wind | Play=*Yes* | Play=*No* |
|------|------------|-----------|
| *Strong* | 3/9 | 3/5 |
| *Weak* | 6/9 | 2/5 |

$$P(\text{Play}=Yes) = 9/14 \qquad P(\text{Play}=No) = 5/14$$

# Example

- Test Phase
  - Given a new instance, predict its label

    **x′**=(Outlook=*Sunny,* Temperature=*Cool,* Humidity=*High,* Wind=*Strong*)
  - Look up tables achieved in the learning phrase

    P(Outlook=*Sunny*|Play=*Yes*) = 2/9

    P(Temperature=*Cool*|Play=*Yes*) = 3/9

    P(Huminity=*High*|Play=*Yes*) = 3/9

    P(Wind=*Strong*|Play=*Yes*) = 3/9

    P(Play=*Yes*) = 9/14

    P(Outlook=S*unny*|Play=*No*) = 3/5

    P(Temperature=*Cool*|Play==*No*) = 1/5

    P(Huminity=*High*|Play=*No*) = 4/5

    P(Wind=*Strong*|Play=*No*) = 3/5

    P(Play=*No*) = 5/14

  - Decision making with the MAP rule

    P(*Yes*|**x′**) ≈ [P(*Sunny*|Y*es*)P(*Cool*|*Yes*)P(*High*|*Yes*)P(*Strong*|*Yes*)]P(Play=*Yes*) = 0.0053

    P(*No*|**x′**) ≈ [P(*Sunny*|N*o*) P(*Cool*|N*o*)P(*High*|N*o*)P(*Strong*|N*o*)]P(Play=*No*) = 0.0206

    Given the fact P(*Yes*|**x′**) < P(*No*|**x′**), we label **x′** to be "*No*".

# Naïve Bayes

- ## Algorithm: Continuous-valued Features

  - Numberless values taken by a continuous-valued feature

  - Conditional probability often modeled with the normal distribution

$$\hat{P}(X_j \mid C = c_i) = \frac{1}{\sqrt{2\pi}\sigma_{ji}} \exp\left( - \frac{(X_j - \mu_{ji})^2}{2\sigma_{ji}^2} \right)$$

$\mu_{ji} :$ mean (avearage) of feature values $X_j$ of examples for which $C = c_i$

$\sigma_{ji} :$ standard deviation of feature values $X_j$ of examples for which $C = c_i$

  - Learning Phase: for $\mathbf{X} = (X_1, \cdots, X_n)$, $C = c_1, \cdots, c_L$

    Output: $n \times L$ normal distributions and $P(C = c_i)$ $i = 1, \cdots, L$

  - Test Phase: Given an unknown instance $\mathbf{X}' = (a_1', \cdots, a_n')$

    - Instead of looking-up tables, calculate conditional probabilities with all the normal distributions achieved in the learning phrase

    - Apply the MAP rule to make a decision

# Naïve Bayes

- Example: Continuous-valued Features
  - Temperature is naturally of continuous value.

    **Yes**: 25.2, 19.3, 18.5, 21.7, 20.1, 24.3, 22.8, 23.1, 19.8

    **No**: 27.3, 30.1, 17.4, 29.5, 15.1

  - Estimate mean and variance for each class

  $$\mu = \frac{1}{N}\sum_{n=1}^{N} x_n, \quad \sigma^2 = \frac{1}{N}\sum_{n=1}^{N}(x_n - \mu)^2$$

  $\mu_{Yes} = 21.64, \; \sigma_{Yes} = 2.35$

  $\mu_{No} = 23.88, \; \sigma_{No} = 7.09$

  - **Learning Phase**: output two Gaussian models for P(temp|C)

  $$\hat{P}(x \mid Yes) = \frac{1}{2.35\sqrt{2\pi}}\exp\left(-\frac{(x-21.64)^2}{2\times 2.35^2}\right) = \frac{1}{2.35\sqrt{2\pi}}\exp\left(-\frac{(x-21.64)^2}{11.09}\right)$$

  $$\hat{P}(x \mid No) = \frac{1}{7.09\sqrt{2\pi}}\exp\left(-\frac{(x-23.88)^2}{2\times 7.09^2}\right) = \frac{1}{7.09\sqrt{2\pi}}\exp\left(-\frac{(x-23.88)^2}{50.25}\right)$$

# Relevant Issues

- Violation of Independence Assumption

    - For many real world tasks, $P(X_1,\cdots,X_n \mid C) \neq P(X_1 \mid C)\cdots P(X_n \mid C)$

    - Nevertheless, naïve Bayes works surprisingly well anyway!

- Zero conditional probability Problem

    - If no example contains the attribute value $X_j = a_{jk}$, $\hat{P}(X_j = a_{jk} \mid C = c_i) = 0$

    - In this circumstance, $\hat{P}(x_1 \mid c_i)\cdots\hat{P}(a_{jk} \mid c_i)\cdots\hat{P}(x_n \mid c_i) = 0$ during test

    - For a remedy, conditional probabilities estimated with

$$\hat{P}(X_j = a_{jk} \mid C = c_i) = \frac{n_c + mp}{n + m}$$

$n_c$ : number of training examples for which $X_j = a_{jk}$ and $C = c_i$

$n$ : number of training examples for which $C = c_i$

$p$ : prior estimate (usually, $p = 1/t$ for $t$ possible values of $X_j$)

$m$ : weight to prior (number of "virtual" examples, $m \geq 1$)

# Summary

- Naïve Bayes: the conditional independence assumption

  - Training is very easy and fast; just requiring considering each attribute in each class separately

  - Test is straightforward; just looking up tables or calculating conditional probabilities with estimated distributions

- A popular generative model

  - Performance competitive to most of state-of-the-art classifiers even in presence of violating independence assumption

  - Many successful applications, e.g., spam mail filtering

  - A good candidate of a base learner in ensemble learning

  - Apart from classification, naïve Bayes can do more...