

Data Valorization: Naïve Bayes Classification

Lionel Fillatre

fillatre@unice.fr

Outline

- Introduction
- Probability Basics
- Naïve Bayes
- Discrete Case
- Continuous Case
- Conclusion



1 Introduction

Background

- There are three methods to establish a classifier
 - a) Model a classification rule directly

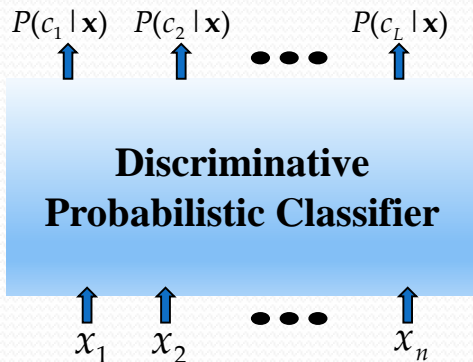
Examples: k-NN, decision trees, perceptron, SVM
 - b) Model the probability of class memberships given input data

Example: logistic regression, perceptron with the cross-entropy cost
 - c) Make a probabilistic model of data within each class

Examples: naive Bayes, hypothesis testing
- *a)* and *b)* are examples of **discriminative** classification
- *c)* is an example of **generative** classification
- *b)* and *c)* are both examples of **probabilistic** classification

Probabilistic Classification

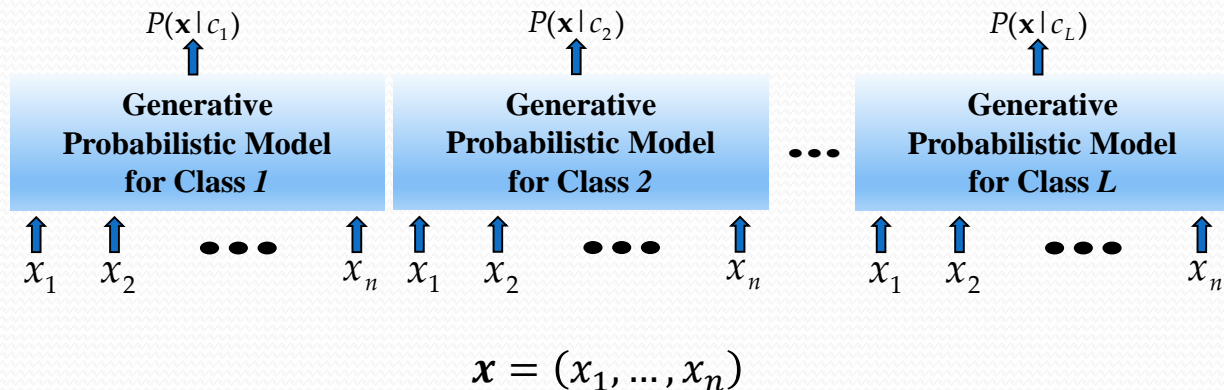
- Establishing a probabilistic model for classification
 - Discriminative model: $P(C|X), C \in \{c_1, \dots, c_L\}, X = (X_1, \dots, X_n)$



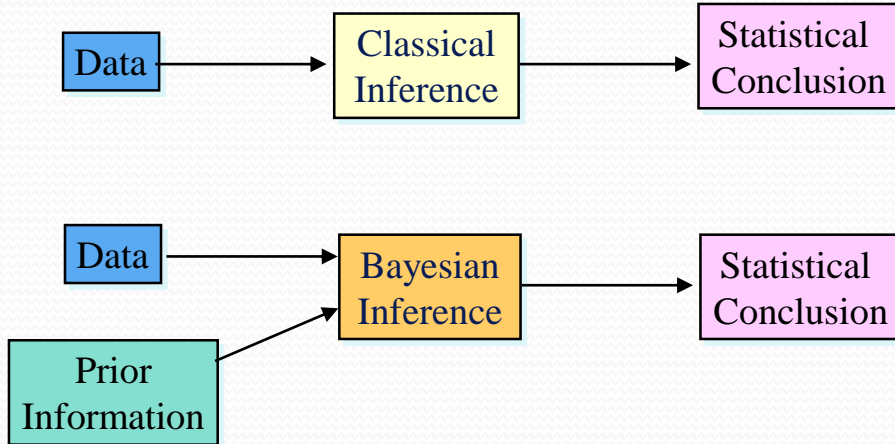
$\mathbf{x} = (x_1, \dots, x_n)$: data, features, etc.

Probabilistic Classification

- Establishing a probabilistic model for classification
 - Generative model: $P(\mathbf{X}|C), C \in \{c_1, \dots, c_L\}, \mathbf{X} = (X_1, \dots, X_n)$



Bayesian and Classical Statistics



Bayesian statistical analysis incorporates a prior probability distribution and likelihoods of observed data to determine a posterior probability distribution of events.

Hypothesis Testing: Bayes Approach

- In the previous lecture, we have assumed that the statistical model of the samples depends on some **non-random** classes C (or some unknown parameter θ)
- The Bayes approach assumes that the class is **random**
- The distribution $P(C)$ of the class is called the **prior distribution**.
 - It is equivalent to a prior distribution $P(\theta)$ on the parameter (in case there is a parameter defining each class).
- The prior distribution is
 - known
 - or assumed to be known
 - or estimated from a training data set



2 Probability Basics

Probability Basics

- Prior, conditional and joint probability for random variables
 - Prior probability: $P(C)$
 - Conditional probability: $P(X | C), P(C | X)$
 - Joint probability: $Z = (X, C), P(Z) = P(X, C)$
 - Relationship: $P(X, C) = P(X | C)P(C) = P(C | X)P(X)$
 - Independence: $P(X | C) = P(X), P(C | X) = P(C), P(X, C) = P(X)P(C)$
- Bayesian Rule

$$P(C | \mathbf{X}) = \frac{P(\mathbf{X} | C)P(C)}{P(\mathbf{X})}$$

$$Posterior = \frac{Likelihood \times Prior}{Evidence}$$

Probability Basics

- We have two six-sided dice. When they are tolled, it could end up with the following occurrence: (*A*) dice 1 lands on side "3", (*B*) dice 2 lands on side "1", and (*C*) Two dice sum to eight.

1) $P(A) = ?$

2) $P(B) = ?$

3) $P(C) = ?$

4) $P(A | B) = ?$

5) $P(C | A) = ?$

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

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

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



Types of errors

		Decision	
		$Z = 0$	$Z = 1$
Ground Truth	$C = 1$	False Negative (FN)	True Positive (TP)
	$C = 0$	True Negative (TN)	False Positive (FP)

Example: rare disease detection

- A rare disease is affecting 0.1% of the population:
 $P(C = 1) = 0,001$ where $C = 1$ when the patient is ill
- Let us consider a medical test for this rare disease. The decision of the test is denoted Z : it indicates that the person is ill if $Z = 1$, otherwise $Z = 0$
- When administered to an ill person, the test will indicate so with probability 0.92 :
 - It follows that $P(Z = 1|C = 1) = 0,92 \Rightarrow P(Z = 0|C = 1) = 0,08$
 - The event $(Z = 0|C = 1)$ is a false negative (misdetection or type II error)
- When administered to a person who is not ill, the test will erroneously give a positive result with probability 0.04 :
 - It follows that $P(Z = 1|C = 0) = 0,04 \Rightarrow P(Z = 0|C = 0) = 0,96$
 - The event $(Z = 1|C = 0)$ is a false positive (false alarm or type I error).

Example (cont.): Conditional Probabilities

Known probabilities:

$$P(C = 1) = 0.001$$

$$P(C = 0) = 0.999$$

$$P(Z = 1|C = 1) = 0.92$$

$$P(Z = 1|C = 0) = 0.04$$

$$P(Z = 0|C = 1) = 0.08$$

$$P(Z = 0|C = 0) = 0.96$$

Conditional probabilities:

$$P(C = 1|Z = 1) = 0.0225$$

$$P(C = 0|Z = 1) = 0.9775$$

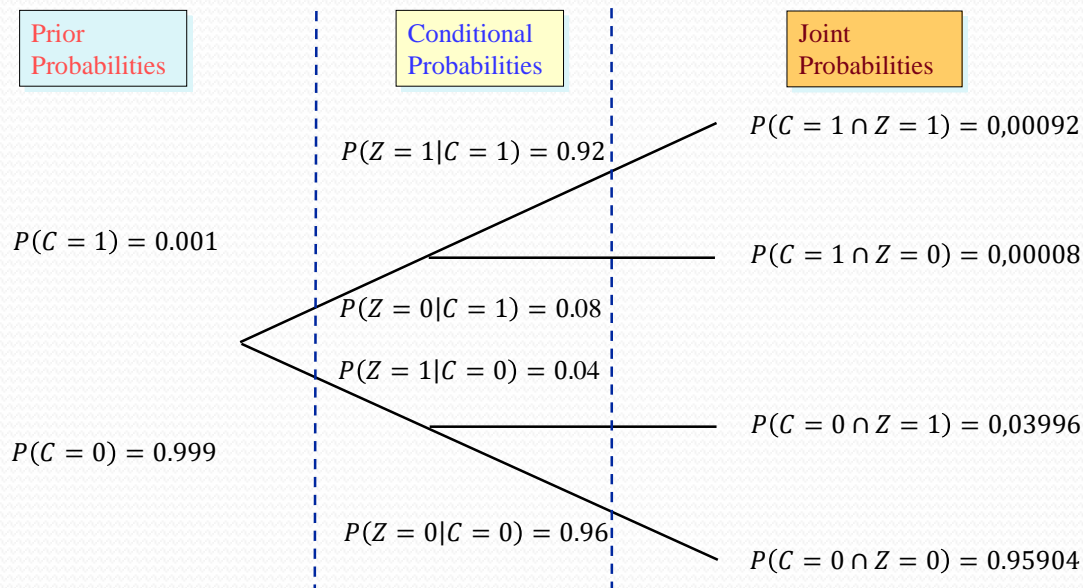
$$P(C = 1|Z = 0) = 0.00008$$

$$P(C = 0|Z = 0) = 0.99992$$

Proof:

$$\begin{aligned} P(C = 1|Z = 1) &= \frac{P(C = 1 \cap Z = 1)}{P(Z = 1)} = \frac{P(C = 1 \cap Z = 1)}{P(C = 0 \cap Z = 1) + P(C = 1 \cap Z = 1)} \\ &= \frac{P(Z = 1|C = 1)P(C = 1)}{P(Z = 1|C = 0)P(C = 0) + P(Z = 1|C = 1)P(C = 1)} = \frac{0.92 \times 0.001}{0.04 \times 0.999 + 0.92 \times 0.001} = 0.0225 \end{aligned}$$

Example (cont.): Error Probability



Global Decision Error

- The global decision error is $P(E) = P(Z \neq C)$

$$\begin{aligned}P(E) &= P(Z = 0 \cap C = 1) + P(Z = 1 \cap C = 0) \\&= P(Z = 0|C = 1)P(C = 1) + P(Z = 1|C = 0)P(C = 0)\end{aligned}$$

- This shows that

$$P(E) = \text{false negative} \times P(C = 1) + \text{false positive} \times P(C = 0)$$

- Optimality:

- **The Bayes test (not the naïve one) minimizes $P(E)$**
- In case of independent features, the naive Bayes test is optimal.

- Example (cont.):

$$P(E) = 0,00008 + 0,03996 = 0,04004$$

- In this example, the two errors are quite different: why? Does it care?
- $P(E) \approx P(Z = 1 \cap C = 0) \approx P(Z = 1|C = 0)$



3 Naïve Bayes

MAP Rule

- MAP classification rule minimizes the Bayes decision error
 - **MAP: M**aximum **A** Posterior
 - Assign \mathbf{x} to c^* if

$$P(C = c^* | \mathbf{X} = \mathbf{x}) > P(C = c | \mathbf{X} = \mathbf{x}), \quad c^* \neq c, \quad c \in \{c_1, \dots, c_L\}$$

- Generative classification with the MAP rule
 - Apply Bayesian rule to calculate posterior probabilities for all c_i

$$P(C = c_i | \mathbf{X} = \mathbf{x}) = \frac{P(\mathbf{X} = \mathbf{x} | C = c_i)P(C = c_i)}{P(\mathbf{X} = \mathbf{x})} \propto P(\mathbf{X} = \mathbf{x} | C = c_i)P(C = c_i)$$

- Then apply the MAP rule

Naïve Bayes

- Difficulty: learning the joint probability $P(\mathbf{X}_1, \dots, \mathbf{X}_n | C)$

$$P(C = c_i | \mathbf{X} = \mathbf{x}) \propto P(\mathbf{X} = \mathbf{x} | C = c_i) P(C = c_i) = P(\mathbf{X}_1 = \mathbf{x}_1, \dots, \mathbf{X}_n = \mathbf{x}_n | C = c_i) P(C = c_i)$$

- Naïve Bayes classification
 - **Assumption:** all input features are **conditionally independent!**

$$\begin{aligned} P(\mathbf{X}_1, \dots, \mathbf{X}_n | C) &= P(\mathbf{X}_1 | \mathbf{X}_2, \dots, \mathbf{X}_n, C) P(\mathbf{X}_2, \dots, \mathbf{X}_n | C) \\ &= P(\mathbf{X}_1 | C) P(\mathbf{X}_2, \dots, \mathbf{X}_n | C) \\ &= \dots \\ &= P(\mathbf{X}_1 | C) \dots P(\mathbf{X}_n | C) \end{aligned}$$

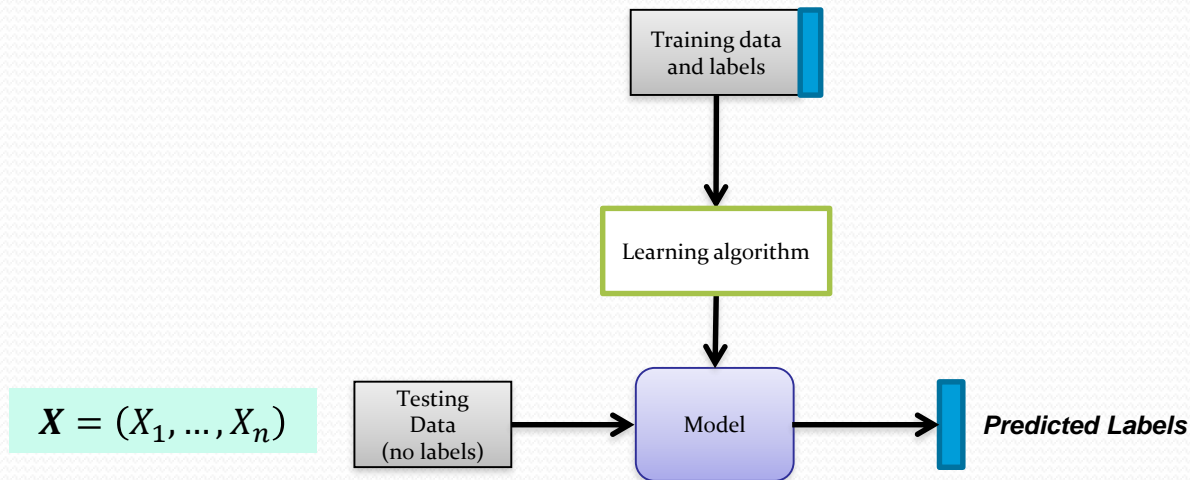
- MAP classification rule

$$c^* = \operatorname{argmax}_{c \in \{c_1, \dots, c_L\}} P(\mathbf{X}_1 = \mathbf{x}_1 | C = c) \dots P(\mathbf{X}_n = \mathbf{x}_n | C = c) P(C = c)$$

- In case of tie-breaking (several maxima), choose randomly the class

The Supervised Learning Pipeline

$$S = \{(\mathbf{X}_1, y_1), \dots, (\mathbf{X}_N, y_N)\}, y_i \in \{c_1, \dots, c_L\}$$





4 Discrete Case

$\mathbf{X}_j \in \{a_{j1}, \dots, a_{jK_j}\}$ with K_j possible values for all j

- Given a training set $S = \{(\mathbf{X}_1, y_1), \dots, (\mathbf{X}_N, y_N)\}$, $y_i \in \{c_1, \dots, c_L\}$
- Learning phase**
 - For each c_i , calculate an estimate $\hat{P}(C = c_i)$ of $P(C = c_i)$ from S
 - For every feature value a_{jk} of each feature \mathbf{X}_j , calculate an estimate $\hat{P}(\mathbf{X}_j = a_{jk} | C = c_i)$ of $P(\mathbf{X}_j = a_{jk} | C = c_i)$ from S
 - Finally, we get $n \times L$ conditional probabilistic (generative) models and L estimates $\hat{P}(C = c_i)$
- Test phase**
 - Given an unknown instance $\mathbf{x}' = (x'_1, \dots, x'_n)$
 - Assign the label c^* to \mathbf{x}' such that

$$c^* = \operatorname{argmax}_{c \in \{c_1, \dots, c_L\}} \hat{P}(\mathbf{X}_1 = x'_1 | C = c) \cdots \hat{P}(\mathbf{X}_n = x'_n | C = c) \hat{P}(C = c)$$

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

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

Example: Test Phase

- Given a new instance, predict its label

$\mathbf{x}' = (\text{Outlook}=\textit{Sunny}, \text{Temperature}=\textit{Cool}, \text{Humidity}=\textit{High}, \text{Wind}=\textit{Strong})$

- Look up tables achieved in the learning phase

$$P(\text{Outlook}=\textit{Sunny} \mid \text{Play}=\textit{Yes}) = 2/9$$

$$P(\text{Temperature}=\textit{Cool} \mid \text{Play}=\textit{Yes}) = 3/9$$

$$P(\text{Humidity}=\textit{High} \mid \text{Play}=\textit{Yes}) = 3/9$$

$$P(\text{Wind}=\textit{Strong} \mid \text{Play}=\textit{Yes}) = 3/9$$

$$P(\text{Play}=\textit{Yes}) = 9/14$$

$$P(\text{Outlook}=\textit{Sunny} \mid \text{Play}=\textit{No}) = 3/5$$

$$P(\text{Temperature}=\textit{Cool} \mid \text{Play}=\textit{No}) = 1/5$$

$$P(\text{Humidity}=\textit{High} \mid \text{Play}=\textit{No}) = 4/5$$

$$P(\text{Wind}=\textit{Strong} \mid \text{Play}=\textit{No}) = 3/5$$

$$P(\text{Play}=\textit{No}) = 5/14$$

- Decision making with the MAP rule

$$P(\text{Yes} \mid \mathbf{x}') \approx [P(\textit{Sunny} \mid \textit{Yes})P(\textit{Cool} \mid \textit{Yes})P(\textit{High} \mid \textit{Yes})P(\textit{Strong} \mid \textit{Yes})]P(\text{Play}=\textit{Yes}) = 0.0053$$

$$P(\text{No} \mid \mathbf{x}') \approx [P(\textit{Sunny} \mid \textit{No})P(\textit{Cool} \mid \textit{No})P(\textit{High} \mid \textit{No})P(\textit{Strong} \mid \textit{No})]P(\text{Play}=\textit{No}) = 0.0206$$

Given the fact $P(\text{Yes} \mid \mathbf{x}') < P(\text{No} \mid \mathbf{x}')$, we label \mathbf{x}' to be “No”.

Relevant Issues

- **Independence Assumption**

- For many real world tasks, $P(\mathbf{X}_1, \dots, \mathbf{X}_n | C) \neq P(\mathbf{X}_1 | C) \cdots P(\mathbf{X}_n | C)$
- Nevertheless, naïve Bayes works surprisingly well anyway!

- **Zero conditional probability Problem**

- If no example contains the feature value $\mathbf{X}_j = a_{jk}$, $\hat{P}(\mathbf{X}_j = a_{jk} | C = c_i) = 0$
- In this circumstance, during test, if $\mathbf{x}'_j = a_{jk}$, then

$$\hat{P}(\mathbf{X}_1 = \mathbf{x}'_1 | C = c_i) \cdots \hat{P}(\mathbf{X}_j = a_{jk} | C = c_i) \cdots \hat{P}(\mathbf{X}_n = \mathbf{x}'_n | C = c_i) = 0$$

- Lidstone's smoothing: for a remedy, conditional probabilities can be re-estimated with

$$\hat{P}(\mathbf{X}_j = a_{jk} | C = c_i) = \frac{N_{jki} + \lambda}{N_i + \lambda K_j}$$

- N_{jki} : number of training examples for which $\mathbf{X}_j = a_{jk}$ and $C = c_i$
- N_i : number of training examples such that $C = c_i$
- K_j : number of possible values of \mathbf{X}_j
- $\lambda > 0$: smoothing parameter
 - $\lambda = 0$: no-smoothing, $\lambda = 1$: add-one smoothing and $\lambda = \infty$: uniform probability $1/K_j$



5 Continuous Case

$X_j \in \mathbb{R}$ for all j

- Given a training set $S = \{(\mathbf{X}_1, y_1), \dots, (\mathbf{X}_N, y_N)\}$, $y_i \in \{c_1, \dots, c_L\}$
- Learning phase**
 - For each c_i , calculate an estimate $\hat{P}(C = c_i)$ of $P(C = c_i)$ from S
 - Conditional probability of each feature \mathbf{X}_j often modeled with the normal distribution

$$\hat{P}(\mathbf{X}_j = \mathbf{x}_j | C = c_i) = \frac{1}{\sqrt{2\pi\sigma_{ji}^2}} \exp\left(-\frac{(\mathbf{x}_j - \mu_{ji})^2}{2\sigma_{ji}^2}\right)$$

- μ_{ji} : mean of feature \mathbf{X}_j of examples for which $C = c_i$
 - σ_{ji} : standard deviation of feature \mathbf{X}_j of examples for which $C = c_i$
 - For every feature \mathbf{X}_j , calculate some estimates $\hat{\mu}_{ji}$ of μ_{ji} and $\hat{\sigma}_{ji}^2$ of σ_{ji}^2 from S
 - Finally, we get $n \times L$ conditional probabilistic (generative) models and L estimates $\hat{P}(C = c_i)$
- Test phase**
 - Given an unknown instance $\mathbf{x}' = (\mathbf{x}'_1, \dots, \mathbf{x}'_n)$
 - Calculate $\hat{P}(\mathbf{X}_j = \mathbf{x}'_j | C = c_i)$ for all j and all i
 - Assign the label c^* to \mathbf{x}' such that $c^* = \underset{c \in \{c_1, \dots, c_L\}}{\operatorname{argmax}} \hat{P}(\mathbf{X}_1 = \mathbf{x}'_1 | C = c) \cdots \hat{P}(\mathbf{X}_n = \mathbf{x}'_n | C = c) \hat{P}(C = c)$

Example: Monitoring a System Temperature

- Two classes:
 - Yes: the system works well
 - No: the system has a failure
- Temperature is naturally a continuous value X :

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:

$$\hat{\mu} = \frac{1}{N} \sum_{i=1}^N x_i, \hat{\sigma}^2 = \frac{1}{N-1} \sum_{i=1}^N (x_i - \hat{\mu})^2 \text{ gives } \begin{cases} \hat{\mu}_{\text{Yes}} = 21.64, & \hat{\sigma}_{\text{Yes}} = 2.35 \\ \hat{\mu}_{\text{No}} = 23.88, & \hat{\sigma}_{\text{No}} = 7.09 \end{cases}$$

- Learning Phase:** output two Gaussian models for $P(\text{temperature}|C)$

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

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



6 Conclusion

Conclusion

- Naïve Bayes is based on the conditional independence assumption between attributes
 - 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
- The Bayes test is generally more difficult to build when taking care of the dependence between the features.
- 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 target test for many learning based classifiers.