

语音及语言信息处理国家工程实验室

Pattern Classification (II)

杜俊

jundu@ustc.edu.cn



中国科学技术大学
安徽科大讯飞信息科技
股份有限公司



Review

- Probability & Statistics
 - Bayes' theorem
 - Random variables: discrete vs. continuous
 - Probability distribution: PDF and CDF
 - Statistics: mean, variance, moment
 - Parameter estimation: MLE
- Information Theory
 - Entropy, mutual information, information channel, KL divergence
- Function Optimization
 - Constrained/unconstrained optimization
- Linear Algebra
 - Matrix manipulation



Outline

- Pattern Classification Problems
 - Inference and decision
- Bayesian Decision Theory
 - How to make the optimal decision?
 - **Maximum *a posteriori* (MAP) decision rule**
- **Generative Models**
 - Joint distribution of observation and label sequences
 - Model estimation: MLE, Bayesian learning, discriminative training
- **Discriminative Models**
 - Model the posterior probability directly (discriminant function)
 - Logistic regression, support vector machine, neural network



Bayesian Decision Theory (I)

- Bayesian decision theory is a fundamental statistical approach to all pattern classification problems
- Pattern classification problem is posed in probabilistic terms
 - Observation X is viewed as random variables (vectors,...)
 - Class id C (C_1, C_2, \dots, C_N) is treated as a discrete random variable
 - All info about X and C can be obtained via joint distribution

$$p(X, C) = P(C) \cdot p(X | C)$$

- Bayesian decision theory leads to the optimal classification with $p(X, C)$
 - Optimal \rightarrow guarantee minimum average classification error
 - The minimum classification error is called the Bayes error



Bayesian Decision Theory (II)

- Prior probabilities of each class $P(C)$
 - How likely any pattern from class C before observing any features
 - Prior knowledge from previous experience

$$\sum_{i=1}^N P(C_i) = 1$$

- Class-conditional probability of observed feature $p(X | C)$
 - How the feature X distributes for all patterns belonging to class C
 - If X is continuous, $p(X | C)$ is a PDF

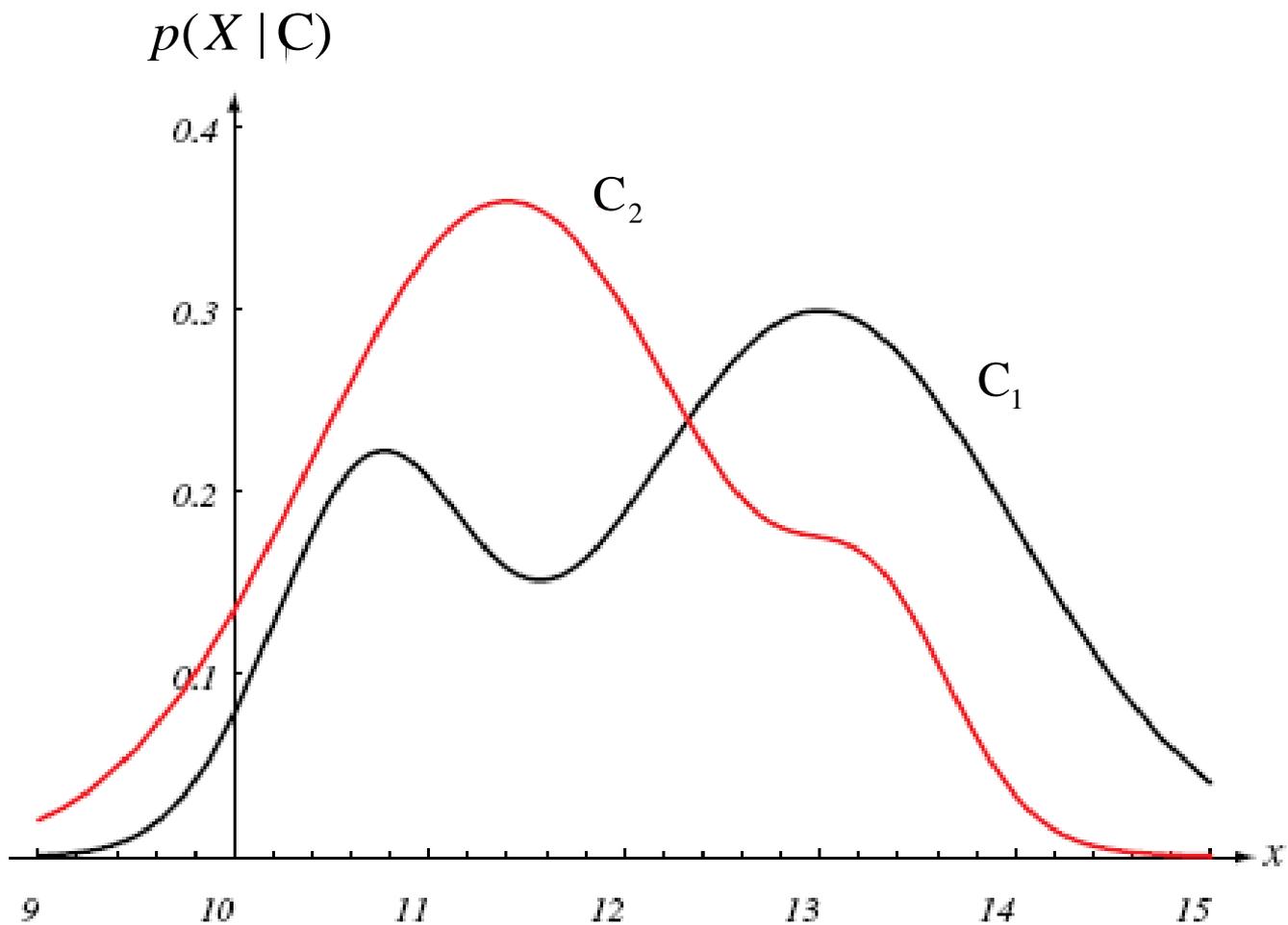
$$\int_X p(X | C_i) \cdot dX = 1$$

- If X is discrete, $p(X | C)$ is a PMF

$$\sum_X p(X | C_i) = 1$$



Examples of Class Conditional Probability



Bayes Decision Rule (I)

- If not observe any feature of an incoming unknown pattern P , classify it based on prior knowledge only
 - Roughly guess it as the class with largest prior probability

$$C_P = \arg \max_C P(C)$$

- If observe some features X of the unknown patter P , we can convert the prior probability $P(C)$ into a posterior probability based on the Bayes' theorem:

$$posterior = \frac{prior \times likelihood}{evidence}$$



Bayes Decision Rule (II)

$$p(C_i | X) = \frac{P(C_i) \cdot p(X | C_i)}{p(X)} = \frac{P(C_i) \cdot p(X | C_i)}{\sum_i P(C_i) \cdot p(X | C_i)}$$

Diagram illustrating the Bayes Decision Rule (II) with labels and arrows:

- $P(C_i)$ is labeled **Prior** (red text) with a blue arrow pointing up to it.
- $p(X | C_i)$ is labeled **Likelihood** (red text) with a blue arrow pointing up to it.
- $p(X)$ is labeled **Evidence** (red text) with a blue arrow pointing down to it.
- $p(C_i | X)$ is labeled **Posterior** (red text) with a blue arrow pointing down to it.

Bayes Decision Rule (III)

- Intuitively, we can classify an unknown pattern into the class with the largest posterior probabilities, resulting in the *maximum a posteriori (MAP) decision rule*, also called *Bayes decision rule*

$$C_p = \arg \max_{C_i} p(C_i | X) = \arg \max_{C_i} P(C_i) \cdot p(X | C_i)$$



The MAP Decision Rule is Optimal (I)

- How well the MAP decision rule behaves??
- Optimality: assume we have complete knowledge $p(X, C)$, the MAP decision rule is optimal to classify patterns, which means it will achieve the lowest average classification error rate.
- Proof of optimality of the MAP rule:
 - Given a pattern P , if its true class id is C_i , but we classify it as C_p , then the classification error is counted as

$$l(C_P | C_i) = \begin{cases} 0 & (C_P = C_i) \\ 1 & (C_P \neq C_i) \end{cases}$$

which is also known as *0-1 loss function*.



The MAP Decision Rule is Optimal (II)

- The expected (average) classification error

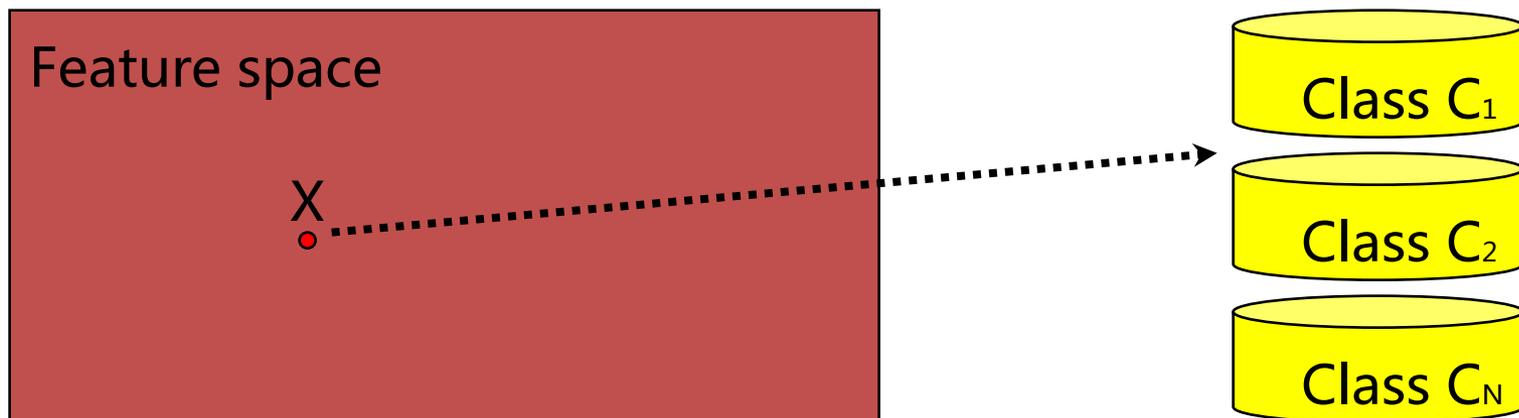
$$R(C_P | X) = \sum_{i=1}^N l(C_P | C_i) \cdot p(C_i | X) = \sum_{C_i \neq C_P} p(C_i | X) = 1 - p(C_P | X)$$

- The optimal classification is to minimize $R(C_P | X)$
 - \rightarrow maximize $p(C_P | X)$
 - \rightarrow the MAP decision rule is optimal

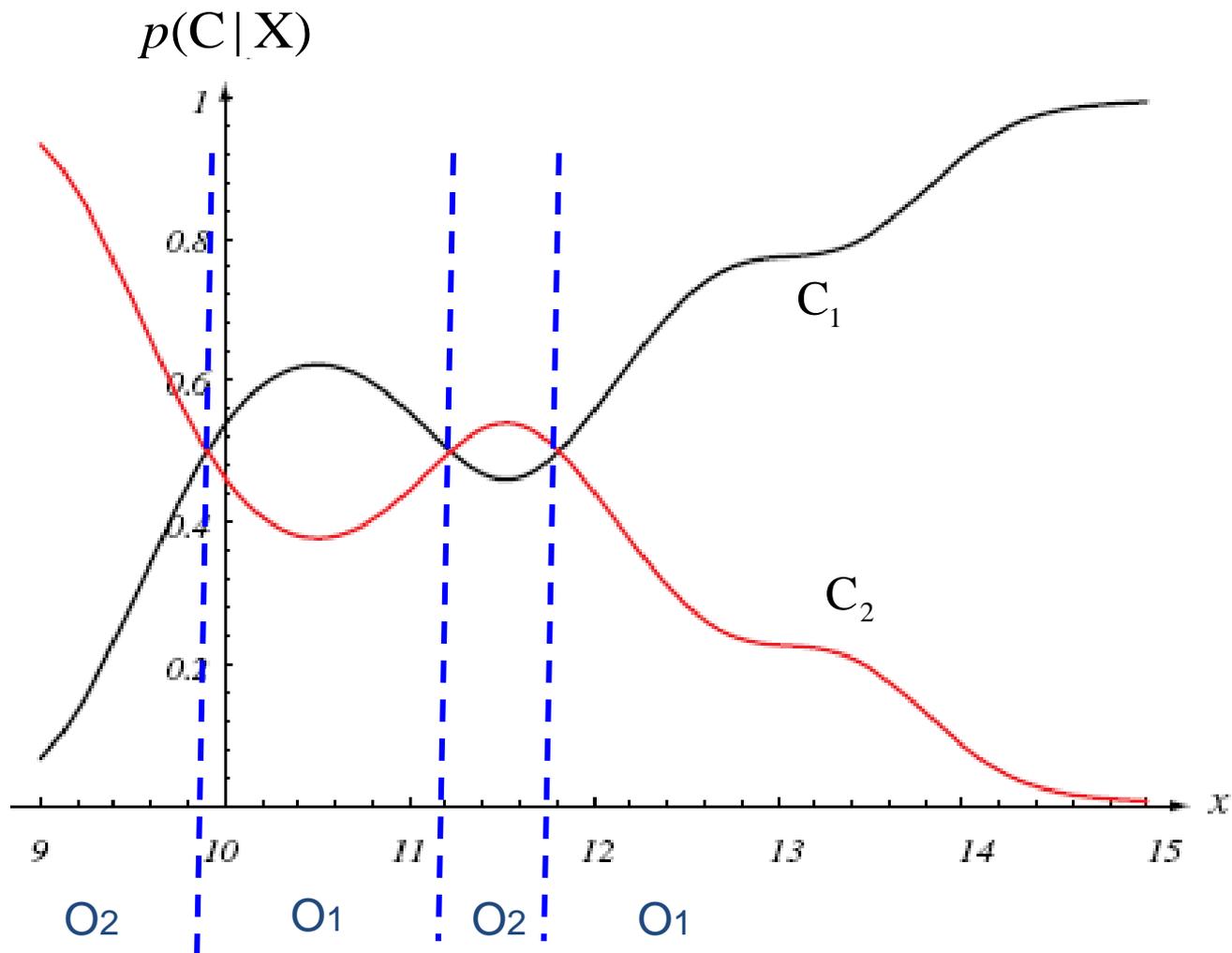


The MAP Decision Rule

- A general decision rule is a mapping function: $X \rightarrow C$
- A decision rule will partition the entire feature space of X into N different regions, O_1, O_2, \dots, O_N . Each region O_i could consist of many contiguous areas.
- If X is located in the region O_i , we classify it as class C_i .
- The MAP decision rule is optimal among all possible decision rules in terms of minimizing average classification errors conditional on that we have complete knowledge about the underlying problem.



Example



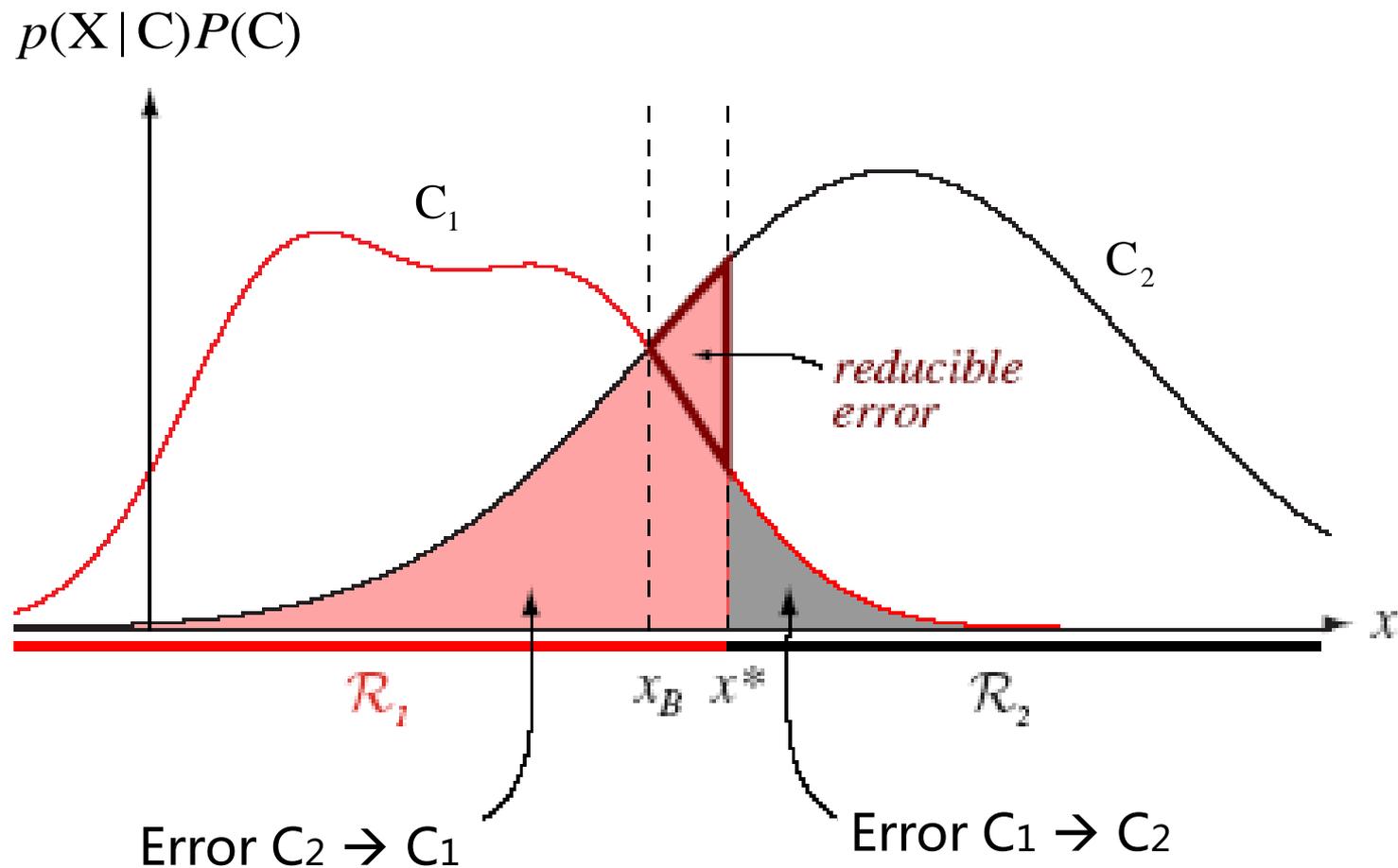
Classification Error Probability

- Assume N-class problem, any a decision rule partitions the feature space into N regions, O_1, O_2, \dots, O_N .
- $\Pr(X \in O_i, C_j)$ denotes the probability of the observation X with true class id C_j in the region O_i .
- The overall classification error probability of the decision rule is:

$$\begin{aligned}\Pr(\text{error}) &= 1 - \Pr(\text{correct}) = 1 - \sum_{i=1}^N \Pr(X \in O_i, C_i) \\ &= 1 - \sum_{i=1}^N \Pr(X \in O_i | C_i) \cdot P(C_i) \\ &= 1 - \sum_{i=1}^N \int_{O_i} p(X | C_i) \cdot P(C_i) dX\end{aligned}$$



Example



Bayes Error

- Bayes error: error probability of the Bayes (MAP) decision rule.
- Since Bayes decision rule guarantees the minimum error, the Bayes error is the lower bound of all possible error probabilities.
- It is difficult to calculate the Bayes error, even for the very simple cases because of discontinuous nature of the decision regions in the integral, especially in high dimensions.
- Some approximation methods to estimate an upper bound.
 - Chernoff bound
 - Bhattacharyya bound
- Evaluate on an independent test set.



Example: X is Discrete (I)

- A simple case (Binomial model): 2-class (C_1, C_2), feature vector is d -dimensional vector, whose components are binary-valued and conditionally independent.

$$X = (x_1, x_2, \dots, x_d)^t \quad x_i = 0, 1 \quad (1 \leq i \leq d)$$

$$p_i = \Pr(x_i = 1 | C_1) \quad q_i = \Pr(x_i = 1 | C_2)$$

$$p(X | C_1) = \prod_{i=1}^d p_i^{x_i} (1 - p_i)^{1-x_i}$$

$$p(X | C_2) = \prod_{i=1}^d q_i^{x_i} (1 - q_i)^{1-x_i}$$



Example: X is Discrete (II)

- The MAP decision rule:

classify to C_1 if $P(C_1) \cdot p(X | C_1) \geq P(C_2) \cdot p(X | C_2)$, otherwise C_2

Equivalently, we have the decision function :

$$g(X) = \sum_{i=1}^d \left[x_i \ln \frac{p_i}{q_i} + (1 - x_i) \ln \frac{1 - p_i}{1 - q_i} \right] + \ln \frac{P(C_1)}{P(C_2)} = \sum_{i=1}^d \lambda_i x_i + \lambda_0$$

$$\lambda_i = \ln \frac{p_i(1 - q_i)}{q_i(1 - p_i)} \quad \lambda_0 = \sum_{i=1}^d \ln \frac{1 - p_i}{1 - q_i} + \ln \frac{P(C_1)}{P(C_2)}$$

If $g(X) \geq 0$, classify to C_1 , otherwise C_2 .



Example: X is Continuous

- Gaussian model: 2-class (C_1, C_2), the feature vector is a scalar which is real-valued

$$P(x | C_1) = N(x; \mu_1, \sigma_1^2) = \frac{1}{\sqrt{2\pi\sigma_1^2}} e^{-(x-\mu_1)^2/2\sigma_1^2}$$

$$P(x | C_2) = N(x; \mu_2, \sigma_2^2) = \frac{1}{\sqrt{2\pi\sigma_2^2}} e^{-(x-\mu_2)^2/2\sigma_2^2}$$

- The MAP decision rule:

classify to C_1 if $P(C_1) \cdot p(x | C_1) \geq P(C_2) \cdot p(x | C_2)$, otherwise C_2

.



Missing Features/Data (I)

- If we know the full probability structure of a problem, we can construct the optimal Bayes decision rule.
- In some practical situations, for some patterns, we can't observe the full feature vector described in the probability structure. Only partial information of the feature vector is observed, but some components are missing.
- How to classify such corrupted inputs to obtain minimum average error?
- Let the full feature vector $X=[X_g, X_b]$, X_g represents the observed or good features, X_b represents the missing or bad ones.
- In this case, the optimal decision rule is constructed as follows:

$$C_p = \arg \max_{C_i} p(C_i | X_g)$$



Missing Features/Data (II)

$$p(C_i | X_g) = \frac{p(C_i, X_g)}{p(X_g)} = \frac{\int p(C_i, X_g, X_b) dX_b}{p(X_g)}$$

$$(1) = \frac{\int p(C_i | X_g, X_b) \cdot p(X_g, X_b) dX_b}{p(X_g)} = \frac{\int p(C_i | X) \cdot p(X) dX_b}{\int p(X) dX_b}$$

$$(2) = \frac{\int P(C_i) \cdot p(X_g, X_b | C_i) dX_b}{\sum_{C_i} \int P(C_i) \cdot p(X_g, X_b | C_i) dX_b} = \frac{\int P(C_i) \cdot p(X | C_i) dX_b}{\sum_{C_i} \int P(C_i) \cdot p(X | C_i) dX_b}$$



Practical Issue

- The optimal Bayes decision rule is not feasible in practice.
 - In any practical problem, we can not have a complete knowledge about the problem.
 - E.g., the class-conditional probability are always unavailable and extremely hard to estimate.
- However, possible to collect a set of sample data for each class in question.
 - The sample data are always far from enough to estimate a reliable PDF by using sample data themselves ONLY.
- Question: How to build a reasonable classifier based on a limited set of sample data, instead of the true PDF?

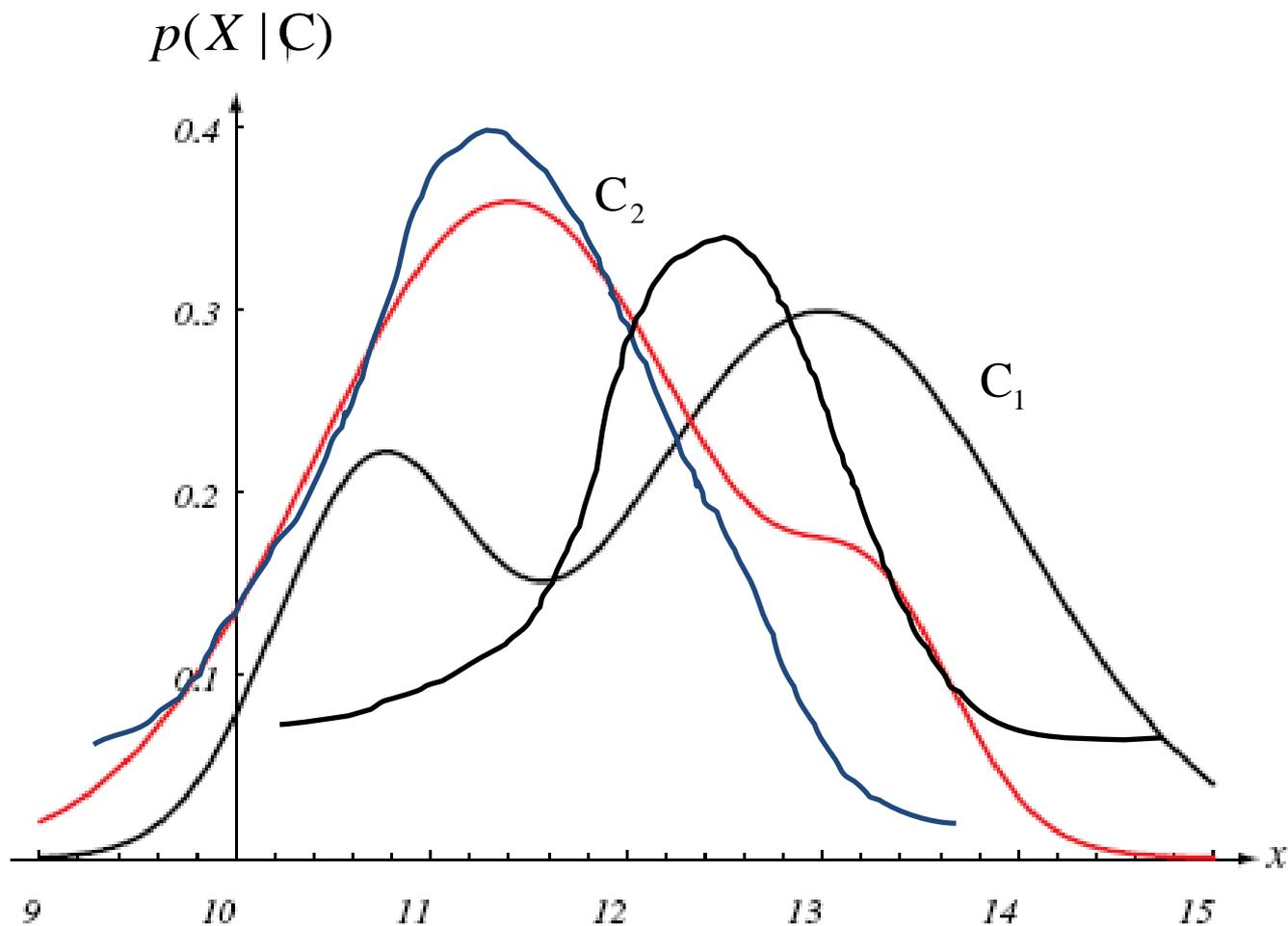


Statistical Data Modeling

- For any real problem, the true PDFs are always unknown
- Statistical data modeling: based on the available sample data set, choose a proper statistical model to fit into the available data set.
 - Data modeling stage: once the statistical model is selected, its function form becomes known except a set of model parameters associated with the model are unknown to us.
 - Learning (training) stage: the unknown parameters can be estimated by fitting the model into the data set based on certain estimation criterion.
 - Decision (test) stage: the estimated PDFs are plugged into the optimal Bayes decision rule in place of the real PDFs , so called **plug-in MAP decision rule**
 - Not optimal but performs reasonably well in practice



Data Modeling Example



Plug-in MAP Decision Rule

- Once the statistical models are estimated, they are treated as if they were true distributions of the data, and plug into the form of the optimal Bayes (MAP) decision rule in place of the unknown true PDFs.
- The plug-in MAP decision rule:

$$C_P = \arg \max_{C_i} p(C_i | X) = \arg \max_{C_i} P(C_i) \cdot p(X | C_i)$$

$$\approx \arg \max_{C_i} \bar{P}_{\Gamma_i}(C_i) \cdot \bar{p}_{\Lambda_i}(X | C_i)$$

