

CENG 463

Machine Learning

Lecture 02 - Bayes Decision Theory

Bayes Decision Theory

- It is a statistical approach to machine learning.
- Assumptions:
 - Decision problem is probabilistic
 - All relevant probability values are known
- We derive decision rules that are optimal in the sense that they either minimize average probability of **error** or overall **risk**.

Bayes Decision Theory

Example: Credit scoring for bank customers

- Inputs are income and savings: $x = [x_1, x_2]^T$
- There are two classes: low-risk customers ($C=1$) and high-risk customers ($C=0$).
- Probabilities:
 - $P(C=1 | x_1, x_2)$: probability that input belongs to class 1
 - $P(C=0 | x_1, x_2)$: probability that input belongs to class 0
 - $P(C=0) + P(C=1) = 1$
- Prediction: Choose $\begin{cases} C=1 \text{ if } P(C=1 | x_1, x_2) > 0.5 \\ C=0 \text{ otherwise} \end{cases}$

Bayes' Rule (Two-Class Case)

The diagram shows the equation for Bayes' Rule in a two-class case. The equation is $P(C | \mathbf{x}) = \frac{P(C) P(\mathbf{x} | C)}{P(\mathbf{x})}$. Arrows point from labels to parts of the equation: 'posterior' points to $P(C | \mathbf{x})$, 'prior' points to $P(C)$, 'likelihood' points to $P(\mathbf{x} | C)$, and 'evidence' points to $P(\mathbf{x})$.

$$\text{posterior} \quad \text{prior} \quad \text{likelihood}$$
$$P(C | \mathbf{x}) = \frac{P(C) P(\mathbf{x} | C)}{P(\mathbf{x})}$$

evidence

- $P(C=0) + P(C=1) = 1$
- $P(\mathbf{x}) = P(\mathbf{x}|C=0) \cdot P(C=0) + P(\mathbf{x}|C=1) \cdot P(C=1)$
- $P(C=0|\mathbf{x}) + P(C=1|\mathbf{x}) = 1$

Bayes' Rule: Example

We have found that the word 'Rolex' occurs in 250 of 2000 spam messages, and in 5 of 1000 non-spam messages. What is the probability that a new message containing the word 'Rolex' is spam? Assuming it is equally likely that this new message is spam or non-spam.

Lets denote the class spam by S and non-spam by NS.

Our feature (x) is observing the word Rolex.

$$P(S|x) = \frac{P(S) P(x|S)}{P(x)} = \frac{P(S) P(x|S)}{P(x|S)P(S) + P(x|NS)P(NS)}$$

Bayes' Rule: Example

We have found that the word 'Rolex' occurs in 250 of 2000 spam messages, and in 5 of 1000 non-spam messages.

$$P(x|S) = 250/2000 \quad P(x|NS) = 5/1000$$

We assume S and NS are equally likely:

$$P(S | x) = \frac{P(S) P(x|S)}{P(x)} = \frac{0.5 \cdot 0.125}{0.5 \cdot 0.125 + 0.5 \cdot 0.005} = \frac{125}{130}$$

What is the probability of being non-spam?

Note that the evidence $P(x)$ is only necessary for normalization purposes; it does not affect the decision.

Decision

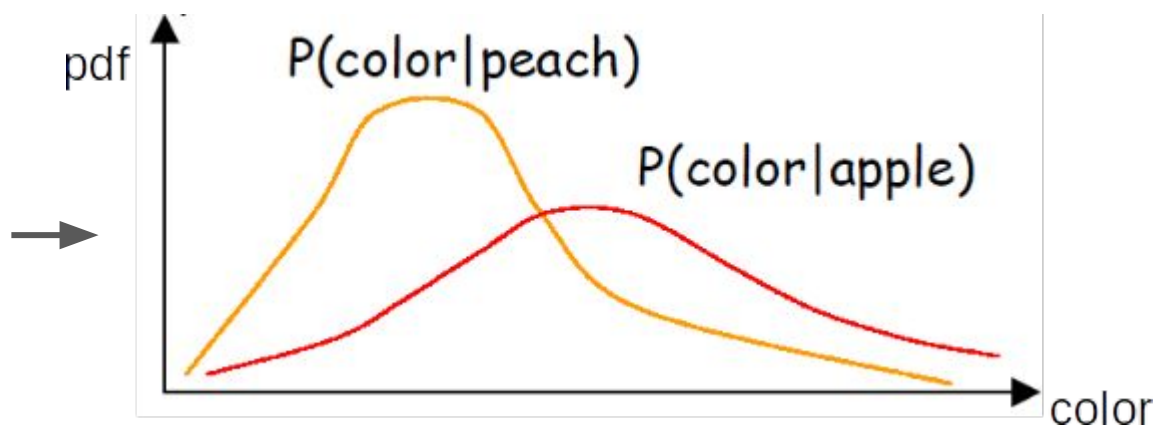
Example: Classification problem of apple and peach by observing their color.

- If we assume initial probabilities are not equal, i.e. if $P(\text{apple}) > P(\text{peach})$ and if we do not have a chance to see the fruit, we always predict apple!
- If we see the color of the fruit, we compute posterior (conditional) probabilities:
 $P(\text{apple}|\text{color})=?$, $P(\text{peach}|\text{color})=?$
- We choose the class with **higher conditional probability**.
- How to find these probabilities? **Bayes Rule**

Decision

- Posterior (conditional) probabilities are
 - $P(\text{apple}|\text{color}) = P(\text{color}|\text{apple}) * P(\text{apple}) / P(\text{color})$
 - $P(\text{peach}|\text{color}) = P(\text{color}|\text{peach}) * P(\text{peach}) / P(\text{color})$
- We need likelihoods: $P(\text{color}|\text{apple})$, $P(\text{color}|\text{peach})$.
- A way of describing likelihood is **probability distribution functions (pdf)**:

For this example, assume, color means intensity of light, i.e. how dark or bright is the fruit.



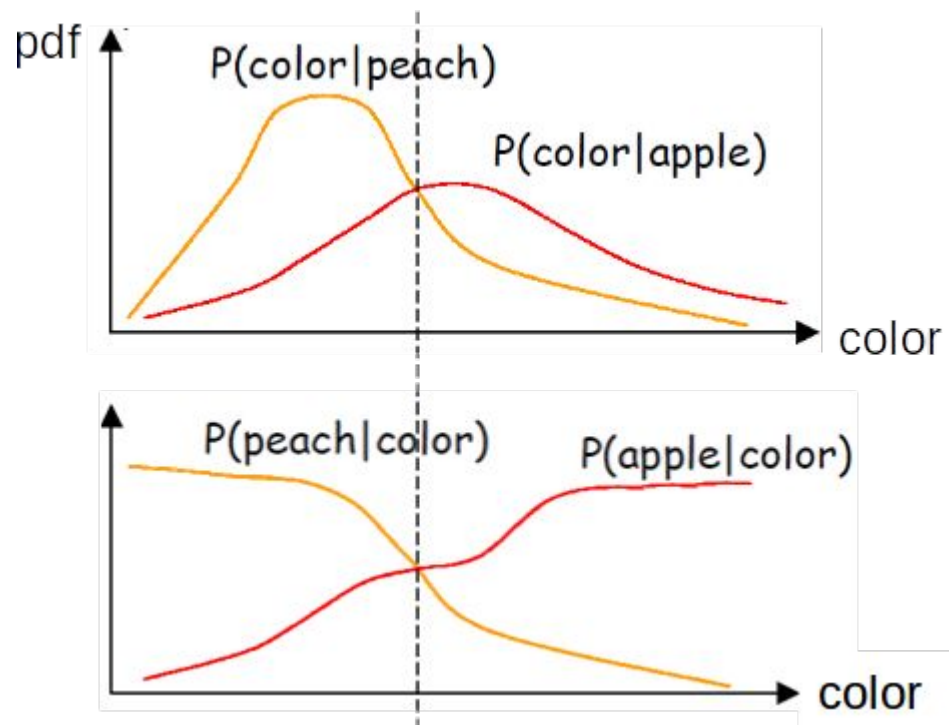
Decision

Using probability distribution functions:

Posterior probabilities after
applying Bayes rule:

Decision:

if $P(\text{apple}|\text{color}) > P(\text{peach}|\text{color})$
then choose apple

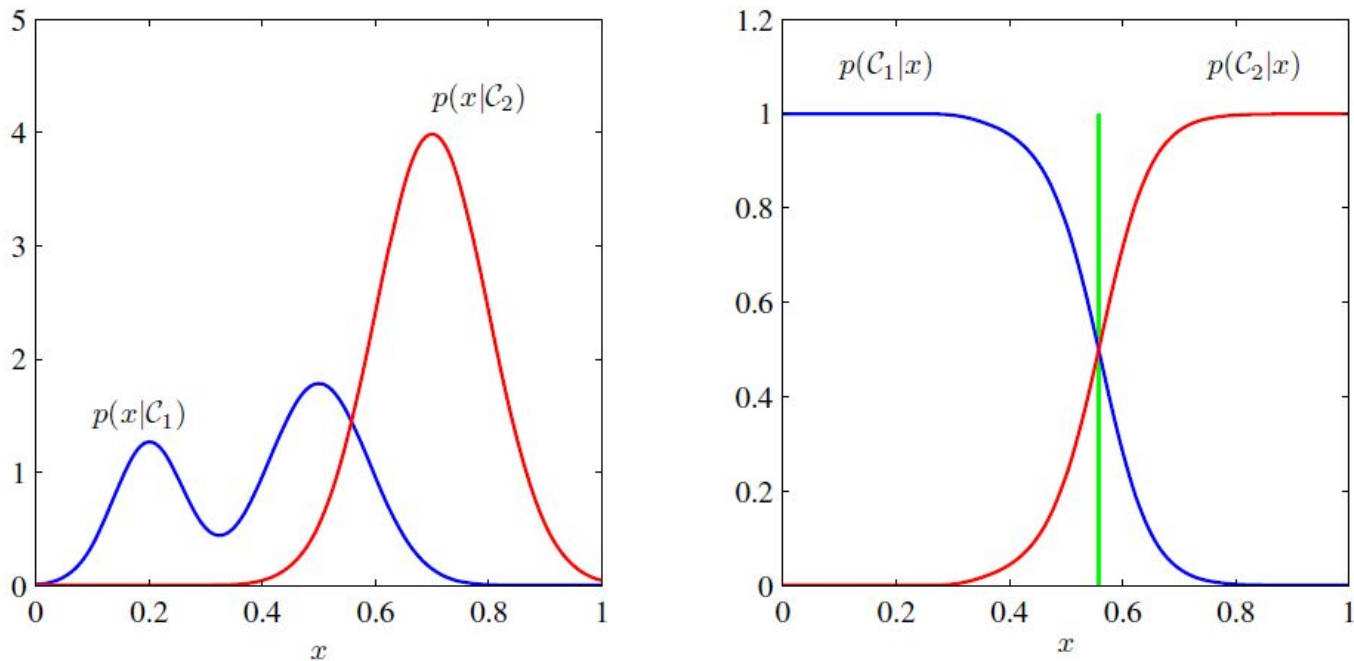


Intersections are co-located
only when prior probabilities are
equal, $P(\text{apple})=P(\text{peach})$.

Two-Class Example

- Assume our color is red and we have the following likelihoods and priors:
 - $P(\text{red}|\text{apple})=0.5$, $P(\text{red}|\text{peach})=0.2$
 - $P(\text{apple})=0.4$, $P(\text{peach})=0.6$
- Given a red fruit, would you say it is an apple or a peach?
- Posteriors:
 - $P(\text{apple}|\text{red})=0.5*0.4=0.2$
 - $P(\text{peach}|\text{red})=0.2*0.6=0.12$
- You would say '**apple**'!
- Normalized posteriors:
 - $P(\text{apple}|\text{red})= 0.2/0.32=0.625$
 - $P(\text{peach}|\text{red})=0.12/0.32=0.375$ } sum = 1

Another Visual Two-Class Example



Probability distributions (left plot) and their posterior probabilities (right plot). Note that the left plot affects the posterior probabilities only when $0.35 < x < 0.8$.

The vertical green line in the right plot shows the decision boundary.

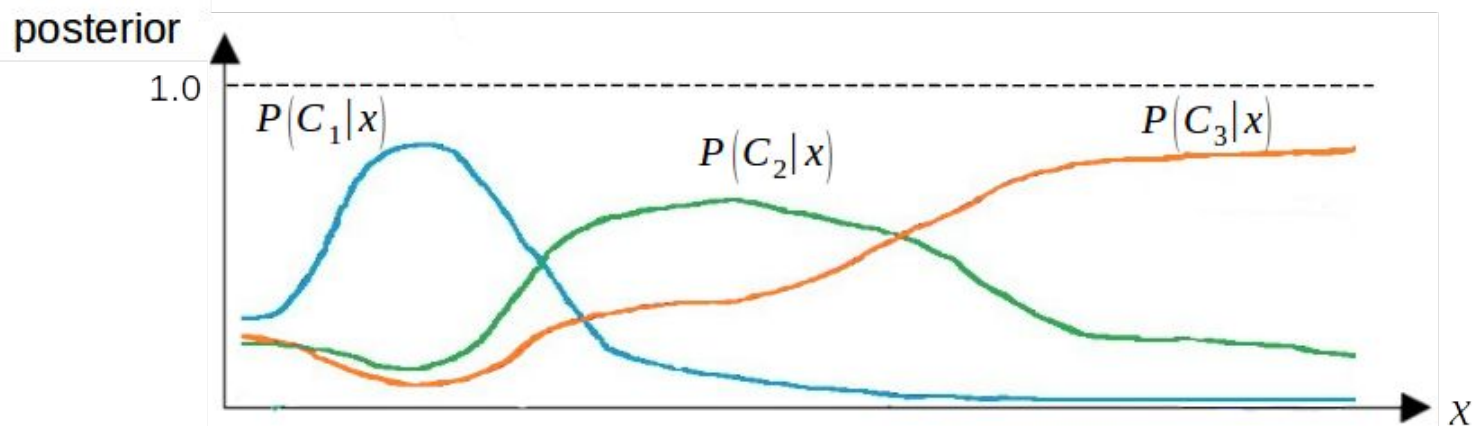
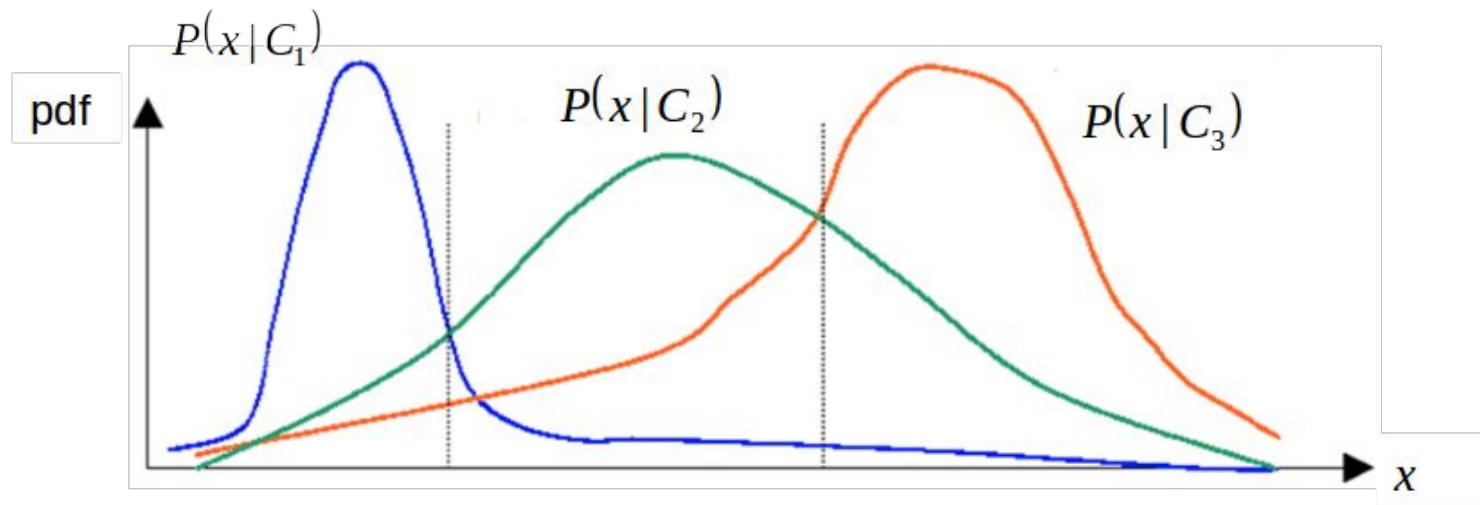
Bayes' Rule (Multi-Class Case)

$$\begin{aligned} P(C_i | \mathbf{x}) &= \frac{P(\mathbf{x} | C_i)P(C_i)}{P(\mathbf{x})} \\ &= \frac{P(\mathbf{x} | C_i)P(C_i)}{\sum_{k=1}^K P(\mathbf{x} | C_k)P(C_k)} \end{aligned}$$

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

choose C_i if $P(C_i | \mathbf{x}) = \max_k P(C_k | \mathbf{x})$

Bayes' Rule (Multi-Class Case)



Losses and Risks

- **Action**, α_i : assigning an input to C_i
- **Loss**, λ_{ik} : Loss of α_i when the actual class is C_k
- Expected risk of choosing α is the sum of losses over all classes:

$$R(\alpha_i | \mathbf{x}) = \sum_{k=1}^K \lambda_{ik} P(C_k | \mathbf{x})$$

- We choose the class with minimum expected risk:

$$\text{choose } \alpha_i \text{ if } R(\alpha_i | \mathbf{x}) = \min_k R(\alpha_k | \mathbf{x})$$

- It is meaningful to define losses as: $\lambda_{ik} = \begin{cases} 0 & \text{if } i = k \\ 1 & \text{if } i \neq k \end{cases}$

Losses and Risks

- When we define losses as:

$$\lambda_{ik} = \begin{cases} 0 & \text{if } i = k \\ 1 & \text{if } i \neq k \end{cases}$$

- Then risk becomes:

$$R(\alpha_i | \mathbf{x}) = \sum_{k=1}^K \lambda_{ik} P(C_k | \mathbf{x}) = \sum_{k \neq i} P(C_k | \mathbf{x}) = 1 - P(C_i | \mathbf{x})$$

- For minimum risk, choose the most probable class.

Losses and Risks: Reject

- In some applications, the cost of choosing a wrong class is very high. So, we may think about a third option: **rejecting** to classify the sample.
- Assume we have a loss of rejection, such that $0 < \lambda < 1$

$$\lambda_{ik} = \begin{cases} 0 & \text{if } i = k \\ \lambda & \text{if } i = K + 1, \quad 0 < \lambda < 1 \\ 1 & \text{otherwise} \end{cases}$$

Risk of rejection: $R(\alpha_{K+1} | \mathbf{x}) = \sum_{k=1}^K \lambda P(C_k | \mathbf{x}) = \lambda \sum_{k=1}^K P(C_k | \mathbf{x}) = \lambda$

Risk of choosing C_i : $R(\alpha_i | \mathbf{x}) = \sum_{k \neq i} P(C_k | \mathbf{x}) = 1 - P(C_i | \mathbf{x})$

Losses and Risks: Reject

- Since we have

$$R(\alpha_{K+1} | \mathbf{x}) = \lambda$$

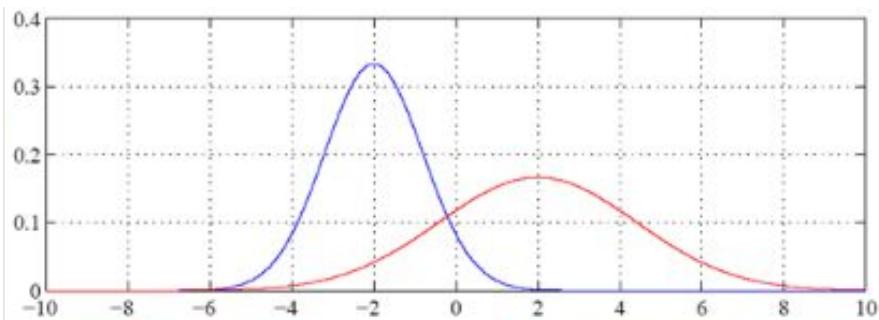
and

$$R(\alpha_i | \mathbf{x}) = 1 - P(C_i | \mathbf{x})$$

- If probability of belonging to class C_i is low or risk is high, we should choose to reject.
- Final decision criterion:

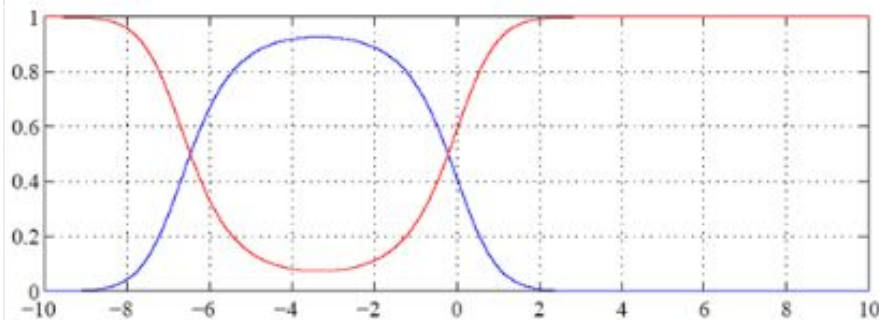
choose C_i if $P(C_i | \mathbf{x}) > P(C_k | \mathbf{x}) \quad \forall k \neq i$ and $P(C_i | \mathbf{x}) > 1 - \lambda$
reject otherwise

$$P(x | C_i)$$



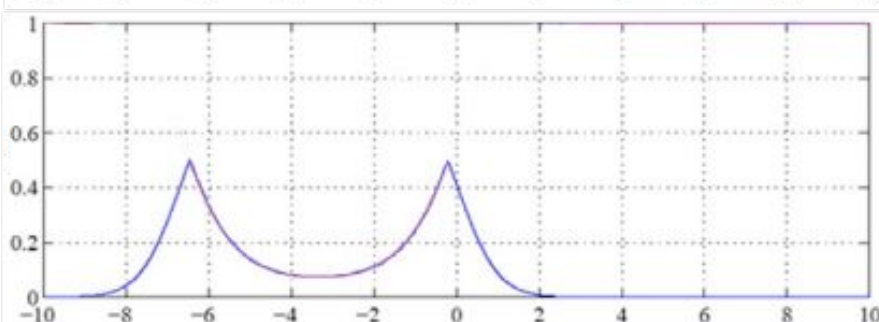
Likelihoods

$$P(C_i | x)$$



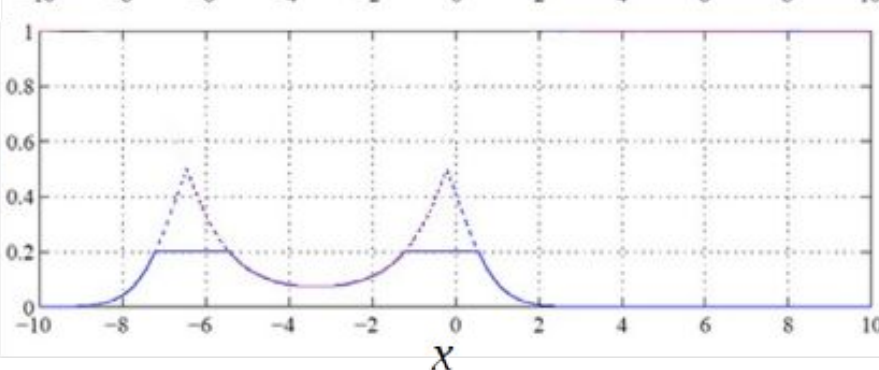
Posteriors

$$R(\alpha_i | x)$$



Risk

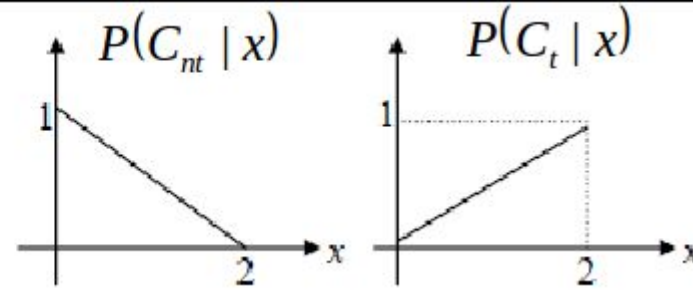
$$R(\alpha_{K+1} | x) = \lambda$$



*Risk with rejection,
rejection loss (λ)=0.2*

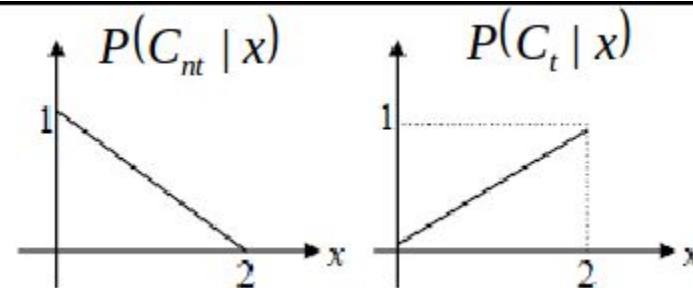
Losses and Risks: Example

- Assume some sensor data, x , is used to detect targets. There are 2 classes: target and no-target.
- If target, system responds, if no-target, system stays idle. A third option is rejection of making a decision.
- The following table is constructed with past experience:

| LOSS | Target | No-target |  | |
|--------------------|----------------|----------------|--------------------------------------------------------------------------------------|--|
| <i>Respond</i> | 0 | 2ε | | |
| <i>Stay Idle</i> | 4ε | 0 | | |
| <i>No-decision</i> | ε | ε | | |

- If your sensor data $x=1.5$, would you respond? Would you stay idle? Would you reject to make a decision?

Losses and Risks: Example

| LOSS | Target | No-target |  |
|--------------------|--------|-----------|-------------------------------------------------------------------------------------|
| <i>Respond</i> | 0 | 2ε | |
| <i>Stay Idle</i> | 4ε | 0 | |
| <i>No-decision</i> | ε | ε | |

$$R(\alpha_t | x) = \sum_{k \neq t} \lambda_{t,k} P(C_k | x) = \lambda_{t,nt} P(C_{nt} | x) \rightarrow R(\alpha_t | 1.5) = 2\varepsilon \cdot 0.25 = 0.5\varepsilon$$

$$R(\alpha_{nt} | x) = \sum_{k \neq nt} \lambda_{nt,k} P(C_k | x) = \lambda_{nt,t} P(C_t | x) \rightarrow R(\alpha_{nt} | 1.5) = 4\varepsilon \cdot 0.75 = 3\varepsilon$$

$$R(\alpha_{K+1} | x) = \sum_{k=1}^K \lambda P(C_k | x) = \lambda (P(C_t | x) + P(C_{nt} | x)) \rightarrow R(\alpha_{K+1} | 1.5) = \varepsilon \cdot 1 = \varepsilon$$

Discriminant Functions

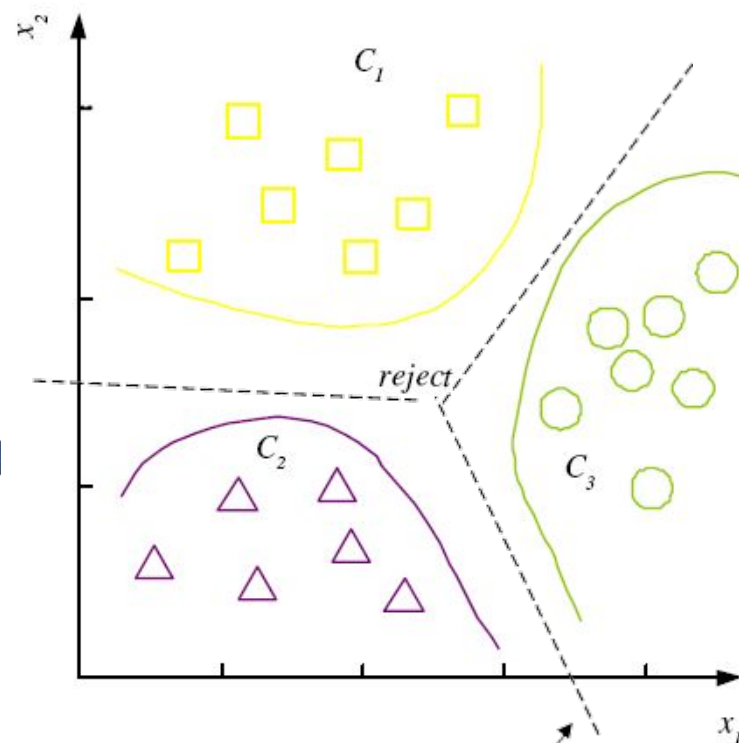
- Classification is done by defining discriminant function $g(\mathbf{x})$.

choose C_i if $g_i(\mathbf{x}) = \max_k g_k(\mathbf{x})$

$$g_i(\mathbf{x}) = \begin{cases} -R(\alpha_i | \mathbf{x}) & \leftarrow \text{minimum risk} \\ P(C_i | \mathbf{x}) & \leftarrow \text{maximum posterior} \\ p(\mathbf{x} | C_i)P(C_i) & \leftarrow \text{unnormalized posterior} \end{cases}$$

- K decision regions R_1, \dots, R_K

$$R_i = \{\mathbf{x} | g_i(\mathbf{x}) = \max_k g_k(\mathbf{x})\}$$



If rejection is
not an option

Summary

We have learned about:

- Bayes' Rule
- How to Make a Decision using Bayes' Rule?
- Risk
- Risk with Rejection Option
- What are Discriminant Functions?