

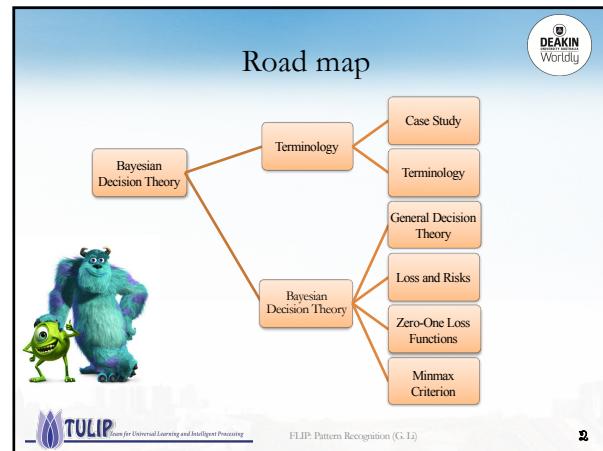
Lecture Notes on
Pattern Recognition

Session 02(A): Bayesian Decision Theory (I)

Gang Li
School of Information Technology
Deakin University, VIC 3125, Australia

DEAKIN
Worldly

TULIP Team for Universal Learning and Intelligent Processing



Terminology

- A Case Study
- Terminology

DEAKIN
Worldly

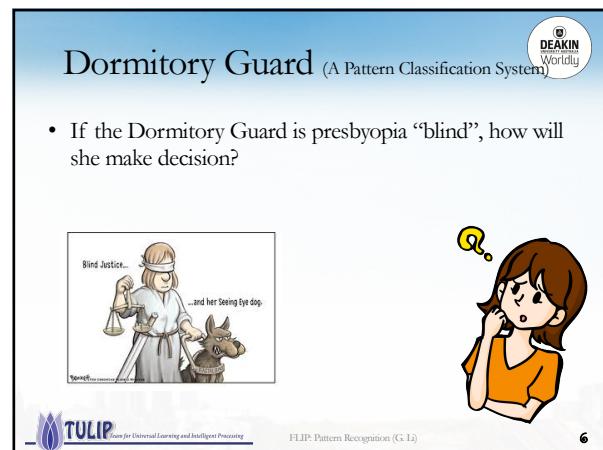
TULIP Team for Universal Learning and Intelligent Processing



Dormitory Guard (A Pattern Classification System)

DEAKIN
Worldly

TULIP Team for Universal Learning and Intelligent Processing



Terminology

- **Random variable:** state of nature ω
 - e.g., ω_1 for girl, ω_2 for boy
- **Priors:** probabilities $P(\omega_1)$ and $P(\omega_2)$
 - e.g., prior knowledge of how likely is to meet a girl or a boy nearby



Decision Rule Using Priors Only

- Decide ω_1 if $P(\omega_1) > P(\omega_2)$; otherwise decide ω_2
- $$P(\text{error}) = \begin{cases} P(\omega_1) & \text{if we decide } \omega_2 \\ P(\omega_2) & \text{if we decide } \omega_1 \end{cases} = \min(p(\omega_1), p(\omega_2))$$
- This rule makes the same decision all times!
 - Unless the prior probability changes
 - Favours the most likely class
 - optimum if no other info is available

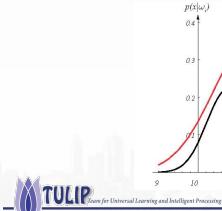
Decision Rule Upon Evidence

- **Evidence:** probability density function $p(x)$
 - how frequently we will measure a pattern with feature value x (e.g., x is **length of the hair**)



Decision Rule Upon Evidence

- **Likelihood:** conditional probability density $p(x | \omega_i)$
 - how frequently we will measure a pattern with feature value x given that the pattern belongs to class ω_i



Decision Rule Upon Evidence

- **Posterior:** conditional probability density $P(\omega_j | x)$
 - the probability that the person belongs to class ω_j given measurement x .



$$\begin{aligned} P(\omega_1 | x) &= \frac{p(x | \omega_1)P(\omega_1)}{p(x)} \\ P(\omega_2 | x) &= \frac{p(x | \omega_2)P(\omega_2)}{p(x)} \end{aligned}$$

Decision Rule Upon Evidence

- **Bayes Rule:**
 - the posterior probability of category ω_j given measurement x is given by
- $$P(\omega_j | x) = \frac{p(x | \omega_j)P(\omega_j)}{p(x)} = \frac{\text{likelihood} \times \text{prior}}{\text{evidence}}$$
- $$p(x) = \sum_{j=1}^2 p(x | \omega_j)P(\omega_j)$$
- Decide ω_1 if $P(\omega_1 | x) > P(\omega_2 | x)$; otherwise decide ω_2
 - Decide ω_1 if $p(x | \omega_1)P(\omega_1) > p(x | \omega_2)P(\omega_2)$ otherwise decide ω_2

Decision Rule Upon Evidence

- In this case, the probability of error is:
$$P(\text{error} / x) = \begin{cases} P(\omega_1 / x) & \text{if we decide } \omega_2 \\ P(\omega_2 / x) & \text{if we decide } \omega_1 \end{cases}$$
- What is the average probability error?
$$P(\text{error}) = \int_{-\infty}^{\infty} P(\text{error}, x) dx = \int_{-\infty}^{\infty} P(\text{error} / x) p(x) dx$$
- Bayes rule is optimum**
 - it minimizes the average probability error since
$$P(\text{error}/x) = \min[P(\omega_1/x), P(\omega_2/x)]$$

TULIP Team for Universal Learning and Intelligent Processing FLIP: Pattern Recognition (G. Li) 18

Decision Rule Upon Evidence

$$P(\omega_j / x) = \frac{p(x / \omega_j)P(\omega_j)}{p(x)} = \frac{\text{likelihood} \times \text{prior}}{\text{evidence}}$$

- Bayes Rule Special Case (1): Uniform priors**
 - $P(\omega_1) = P(\omega_2) = \dots = P(\omega_c) = 1/c$
 - Then the decision only depends on the **likelihood**
- Decide ω_1 if $p(x/\omega_1) > p(x/\omega_2)$ otherwise decide ω_2 (binary)

TULIP Team for Universal Learning and Intelligent Processing FLIP: Pattern Recognition (G. Li) 19

Decision Rule Upon Evidence

$$P(\omega_j / x) = \frac{p(x / \omega_j)P(\omega_j)}{p(x)} = \frac{\text{likelihood} \times \text{prior}}{\text{evidence}}$$

- Bayes Rule Special Case (2): Uniform likelihood**
 - $P(x | \omega_1) = \dots = P(x | \omega_c)$
 - Then the decision only depends on the **priors**
- Decide ω_1 if $p(\omega_1) > p(\omega_2)$ otherwise decide ω_2 (binary)

TULIP Team for Universal Learning and Intelligent Processing FLIP: Pattern Recognition (G. Li) 15

Where do Probabilities Come From?

- In order to use Bayes Rule, it is necessary to know
 - Priors:** $P(\omega_i)$
 - Likelihood:** $P(x | \omega_i)$
- Two competitive answers to this question:
 - Relative frequency** (objective) approach.
 - Probabilities can only come from experiments.
 - Bayesian** (subjective) approach.
 - Probabilities may reflect degree of belief and can be based on opinion as well as experiments.
- Density Estimation (covered in later sessions)

TULIP Team for Universal Learning and Intelligent Processing FLIP: Pattern Recognition (G. Li) 16

Dormitory Guard

(A Pattern Classification System)

- Suppose in a dorm area, the priors of girls is 2/3, the priors of boys is 1/3.
 - If the observation is on “**with long hair or not**”, what is the chance that a particular long hair person is a “girl”?



TULIP Team for Universal Learning and Intelligent Processing FLIP: Pattern Recognition (G. Li) 17

Dormitory Guard

(A Pattern Classification System)

- Suppose in a dorm area, the priors of girls is 2/3, the priors of boys is 1/3.
 - If the observation is on “**wearing glasses or not**”, what is the chance that a glasses-wearing person is a “girl”?



TULIP Team for Universal Learning and Intelligent Processing FLIP: Pattern Recognition (G. Li) 18

Dormitory Guard

(A Pattern Classification System)

- Suppose in a dorm area, the priors of girls is $2/3$, the priors of boys is $1/3$.
 - If the observation is on “**wearing skirts or not**”, what is the chance that a glasses-wearing person is a “girl”?



TULIP Team for Universal Learning and Intelligent Processing FLIP: Pattern Recognition (G. Li) 19

Dormitory Guard

(A Pattern Classification System)

- Now suppose the mistakes/errors in decisions exist, but with different cost
 - If letting a girl enter a boy’s dorm building, then the guard will be **fined** for **2 days’ salary**
 - If letting a boy enter a girl’s dorm building, then the guard will be **fined** for **4 days’ salary**
 - If rejecting a boy enter a boy’s dorm, or rejecting a girl enter the girl’s dorm, then the guard will be **fined** for **1 day’s salary**



TULIP Team for Universal Learning and Intelligent Processing FLIP: Pattern Recognition (G. Li) 20

Bayesian Decision Theory



- General Decision Theory
- Loss and Risks
- Zero-One Loss Functions
- MinMax Criterion

TULIP Team for Universal Learning and Intelligent Processing FLIP: Pattern Recognition (G. Li) 21

Bayesian Decision Theory

- Fundamental statistical approach that quantifies the trade-offs between various classification decisions using probabilities and the costs associated with such decisions.
 - Design classifiers to recommend actions that **minimize some total expected “risk”**.
 - The simplest risk is the classification error
 - i.e., all costs are equal

TULIP Team for Universal Learning and Intelligent Processing FLIP: Pattern Recognition (G. Li) 22

General Decision Theory

- Use more than one features.
- Allow more than two categories.
- Allow actions other than classifying the input to one of the possible categories (e.g., rejection).
- Employ a more general error function which penalizes some errors more than others.
 - loss function**
 - i.e., associate “costs” with actions

TULIP Team for Universal Learning and Intelligent Processing FLIP: Pattern Recognition (G. Li) 23

General Decision Theory

- Features form a vector
- A finite set of c categories $\omega_1, \omega_2, \dots, \omega_c$
- Bayes rule (i.e., using vector notation):

$$P(\omega_j / \mathbf{x}) = \frac{p(\mathbf{x} / \omega_j)P(\omega_j)}{p(\mathbf{x})} \quad \text{where } p(\mathbf{x}) = \sum_{j=1}^c p(\mathbf{x} / \omega_j)P(\omega_j)$$
- A finite set of l actions $\alpha_1, \alpha_2, \dots, \alpha_l$
- A **loss function** $\lambda(\alpha_i | \omega_j)$
 - the loss or cost incurred for taking action α_i when the classification category is ω_j

TULIP Team for Universal Learning and Intelligent Processing FLIP: Pattern Recognition (G. Li) 24

Expected Loss (or Conditional Risk)

- Suppose we observe x and take action α_i
- If the true classification category is ω_j , we incur the loss $\lambda(\alpha_i | \omega_j)$
- The expected loss (or conditional risk) with taking action α_i :

$$R(a_i / \mathbf{x}) = \sum_{j=1}^c \lambda(a_i / \omega_j) P(\omega_j / \mathbf{x})$$



FLIP: Pattern Recognition (G. Li)

25

Overall Risk

- Overall risk R

$$R = \int R(a(\mathbf{x}) / \mathbf{x}) p(\mathbf{x}) d\mathbf{x}$$

where $a(\mathbf{x})$ is a general decision rule that determines which action $\alpha_1, \alpha_2, \dots, \alpha_c$ to take for every \mathbf{x} .

- To minimize R , we need a decision rule $a(\mathbf{x})$ that chooses the action with the minimum conditional risk $R(a_i | \mathbf{x})$ for every \mathbf{x} .



FLIP: Pattern Recognition (G. Li)

26

Bayes Risk

- Bayes decision rule** minimizes R by:
 - Computing $R(\alpha_i | \mathbf{x})$ for every α_i given an input \mathbf{x}
 - Choosing the action α_i with the minimum $R(\alpha_i | \mathbf{x})$
- The resulting minimum overall risk is called **Bayes risk** and is the best (i.e., optimum) performance that can be achieved.

$$R^* = \min R$$



FLIP: Pattern Recognition (G. Li)

27

Two-category classification

- Define
 - α_1 : decide ω_1
 - α_2 : decide ω_2
 - $\lambda_{ij} = \lambda(\alpha_i | \omega_j)$
- The conditional risks are:

$$R(a_i / \mathbf{x}) = \sum_{j=1}^c \lambda(a_i / \omega_j) P(\omega_j / \mathbf{x}) \quad (c=2)$$



$$R(a_1 / \mathbf{x}) = \lambda_{11} P(\omega_1 / \mathbf{x}) + \lambda_{12} P(\omega_2 / \mathbf{x})$$

$$R(a_2 / \mathbf{x}) = \lambda_{21} P(\omega_1 / \mathbf{x}) + \lambda_{22} P(\omega_2 / \mathbf{x})$$

FLIP: Pattern Recognition (G. Li)

28

Two-category classification

- Minimum risk decision rule:
- Decide ω_1 if $R(a_1 / \mathbf{x}) < R(a_2 / \mathbf{x})$; otherwise decide ω_2
- Decide ω_1 if $(\lambda_{21} - \lambda_{11})P(\omega_1 / \mathbf{x}) > (\lambda_{12} - \lambda_{22})P(\omega_2 / \mathbf{x})$; otherwise decide ω_2
- In general, the cost for error action (decision) is larger than the cost for right action (decision)

$$\lambda_{21} - \lambda_{11} > 0 \quad \lambda_{12} - \lambda_{22} > 0$$



FLIP: Pattern Recognition (G. Li)

29

Zero-One Loss Function

- Assign the same loss to all errors:

$$\lambda(\alpha_i, \omega_j) = \begin{cases} 0 & i = j \\ 1 & i \neq j \end{cases} \quad i, j = 1, \dots, c$$

- The conditional risk corresponding to this loss function:

$$R(\alpha_i / \mathbf{x}) = \sum_{j=1}^c \lambda(\alpha_i / \omega_j) P(\omega_j / \mathbf{x}) = \sum_{j \neq i} P(\omega_j / \mathbf{x}) = 1 - P(\omega_i / \mathbf{x})$$

"The risk corresponding to the 0-1 loss function is the average probability of error"



FLIP: Pattern Recognition (G. Li)

30

Zero-One Loss Function

- The overall risk in this case is the average probability error! Minimizing the risk requires maximizing the posterior probability $P(\omega_i | x)$

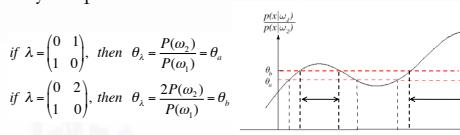
$$R(\alpha_i | x) = 1 - P(\omega_i | x)$$

- For **minimum error rate**
 - Decide ω_i , if $P(\omega_i | x) > P(\omega_j | x), \forall j \neq i$

 TULIP Team for Universal Learning and Intelligent Processing FLIP: Pattern Recognition (G. Li) 31

Decision Boundaries and Decision Regions

- Likelihood Ratio**
Let $\frac{\lambda_{12} - \lambda_{22}}{\lambda_{21} - \lambda_{11}} \cdot \frac{P(\omega_2)}{P(\omega_1)} = \theta_\lambda$ then decide ω_1 if: $\frac{P(x | \omega_1)}{P(x | \omega_2)} > \theta_\lambda$
- If λ is the 0-1 loss function then the threshold involves only the priors:



 TULIP Team for Universal Learning and Intelligent Processing FLIP: Pattern Recognition (G. Li) 32

Dormitory Guard

(A Pattern Classification System)

- Now suppose the mistakes/errors in decisions exist, but with different cost
 - If letting a girl enter a boy's dorm building, then the guard will be **fined for 2 days' salary**
 - If a boy enter a girl's dorm building, then the guard will be **fined for 1 day's salary**
 - If rejecting a boy enter a boy's dorm, or rejecting a girl enter the girl's dorm, then the guard will be **fined for 1 day's salary**



 TULIP Team for Universal Learning and Intelligent Processing FLIP: Pattern Recognition (G. Li) 33

Minmax criterion

- When prior probabilities are not known exactly, we need to design classifiers that perform well over a range of prior probabilities.
- R = $\int R(\alpha(x) | x) p(x) dx$
- Consider the case of two categories:

$$\begin{aligned} R &= \int_{\mathcal{X}_1} [\lambda_{11} P(\omega_1 | x) + \lambda_{12} P(\omega_2 | x)] p(x) dx + \\ &\quad \int_{\mathcal{X}_2} [\lambda_{21} P(\omega_1 | x) + \lambda_{22} P(\omega_2 | x)] p(x) dx \\ &= \int_{\mathcal{X}_1} [\lambda_{11} p(x | \omega_1) P(\omega_1) + \lambda_{12} p(x | \omega_2) P(\omega_2)] dx + \\ &\quad \int_{\mathcal{X}_2} [\lambda_{21} p(x | \omega_1) P(\omega_1) + \lambda_{22} p(x | \omega_2) P(\omega_2)] dx \end{aligned}$$

 TULIP Team for Universal Learning and Intelligent Processing FLIP: Pattern Recognition (G. Li) 34

Minmax criterion

- Since $P(\omega_1) + P(\omega_2) = 1$ and $\int_{\mathcal{X}_1} p(x | \omega_1) dx + \int_{\mathcal{X}_2} p(x | \omega_1) dx = 1$
- We rewrite the overall risk as:

$$R(P(\omega_1)) = \lambda_{22} + (\lambda_{12} - \lambda_{22}) \int_{R_1} p(x | \omega_2) dx +$$

$$P(\omega_1)[(\lambda_{11} - \lambda_{22}) - (\lambda_{21} - \lambda_{11}) \int_{R_2} p(x | \omega_1) dx - (\lambda_{12} - \lambda_{22}) \int_{R_2} p(x | \omega_2) dx]$$

- Assuming that the decision boundary is fixed, the overall risk is a linear function of $P(\omega_1)$

 TULIP Team for Universal Learning and Intelligent Processing FLIP: Pattern Recognition (G. Li) 35

Minmax criterion

- Need to find a decision boundary that makes the coefficient of $P(\omega_1)$ zero
 - overall risk would be independent of priors!
- Minmax solution:

$$R_{mm} = \lambda_{22} + (\lambda_{12} - \lambda_{22}) \int_{\mathcal{X}_1} p(x | \omega_2) dx$$

$$R_{mm} = \lambda_{11} + (\lambda_{21} - \lambda_{11}) \int_{\mathcal{X}_2} p(x | \omega_1) dx$$

 TULIP Team for Universal Learning and Intelligent Processing FLIP: Pattern Recognition (G. Li) 36

Minmax criterion

• How to find the minimax solution?
– Search for the **prior** that **maximizes the Bayes risk**

• R_{mn} becomes equal to the worst Bayes risk

Here, we should design our decision boundary for $P(\omega_1)=0.6$

37

Seminar S02

- **Topic**
 - Try to solve a daily life problem using Bayesian Decision Theory, e.g.:
 - How the dorm guard will make decision on different features?
 - Will I get an A in this unit, if I speak 3 times in seminars?
- **Requirements**
 - Prepare a **15 minutes** talk on your chosen topic
 - Make **ppt** to assist your talk
 - Prepare **at least 3 questions** to ask the audience after your talk
 - Get ready to **take questions** from the audience
- **Hints**
 - Do statistics to get what you need for Bayes Formula
 - If real numbers are hard to get, make reasonable assumptions
 - Try cost-sensitive problems to make your talk more **challenging and interesting**

TULIP Team for Universal Learning and Intelligent Processing FLIP: Pattern Recognition (G. Li)

38

Questions?

TULIP Team for Universal Learning and Intelligent Processing FLIP: Pattern Recognition (G. Li)

39