

Bayesian Decision Theory

Bayesian Decision Theory

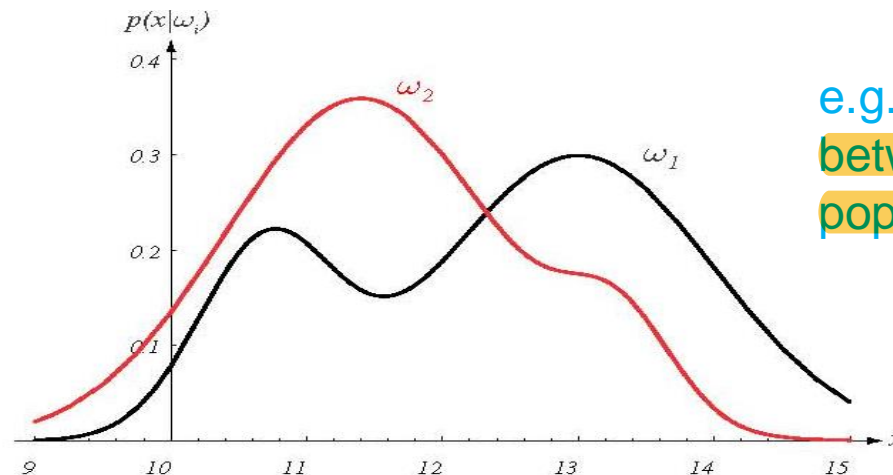
- Design classifiers to recommend **decisions** that minimize some total expected "**risk**".
 - The simplest **risk** is the **classification error** (i.e., costs are equal).
 - Typically, the **risk** includes the **cost** associated with different decisions.

Terminology

- State of nature ω (*random variable*):
 - e.g., ω_1 for sea bass, ω_2 for salmon
- Probabilities $P(\omega_1)$ and $P(\omega_2)$ (*priors*):
 - e.g., prior knowledge of how likely is to get a sea bass or a salmon
- Probability density function $p(x)$ (*evidence*):
 - e.g., how frequently we will measure a pattern with feature value x (e.g., x corresponds to lightness)

Terminology (cont'd)

- Conditional probability density $p(x/\omega_j)$ (*likelihood*) :
 - e.g., how frequently we will measure a pattern with feature value x given that the pattern belongs to class ω_j



e.g., lightness distributions
between salmon/sea-bass
populations

FIGURE 2.1. Hypothetical class-conditional probability density functions show the probability density of measuring a particular feature value x given the pattern is in category ω_i . If x represents the lightness of a fish, the two curves might describe the difference in lightness of populations of two types of fish. Density functions are normalized, and thus the area under each curve is 1.0. From: Richard O. Duda, Peter E. Hart, and David G. Stork, *Pattern Classification*. Copyright © 2001 by John Wiley & Sons,

Terminology (cont'd)

- Conditional probability $P(\omega_j/x)$ (*posterior*) :
 - e.g., the probability that the fish belongs to class ω_j given measurement x .

Decision Rule Using **Prior** Probabilities

Decide ω_1 if $P(\omega_1) > P(\omega_2)$; otherwise decide ω_2

$$P(error) = \begin{cases} P(\omega_1) & \text{if we decide } \omega_2 \\ P(\omega_2) & \text{if we decide } \omega_1 \end{cases}$$

or $P(error) = \min[P(\omega_1), P(\omega_2)]$

- Favours the most likely class.
- This rule will be making the same decision all times.
 - i.e., optimum if no other information is available

Decision Rule Using **Conditional Probabilities**

- Using **Bayes' rule**, the posterior probability of category ω_j given measurement x is given by:

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

where $p(x) = \sum_{j=1}^2 p(x / \omega_j) P(\omega_j)$ (i.e., scale factor – sum of probs = 1)

Decide ω_1 if $P(\omega_1 / x) > P(\omega_2 / x)$; otherwise decide ω_2

or

Decide ω_1 if $p(x/\omega_1)P(\omega_1) > p(x/\omega_2)P(\omega_2)$ otherwise decide ω_2

Decision Rule Using Conditional pdf (cont'd)

$$p(x/\omega_j) \qquad P(\omega_1) = \frac{2}{3} \quad P(\omega_2) = \frac{1}{3} \qquad P(\omega_j/x)$$

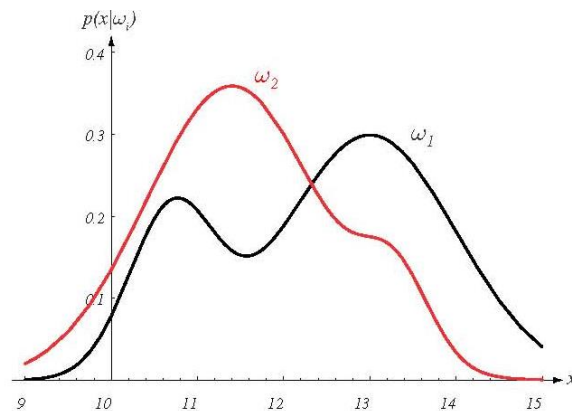


FIGURE 2.1. Hypothetical class-conditional probability density functions show the probability density of measuring a particular feature value x given the pattern is in category ω_i . If x represents the lightness of a fish, the two curves might describe the difference in lightness of populations of two types of fish. Density functions are normalized, and thus the area under each curve is 1.0. From: Richard O. Duda, Peter E. Hart, and David G. Stork, *Pattern Classification*. Copyright © 2001 by John Wiley & Sons,

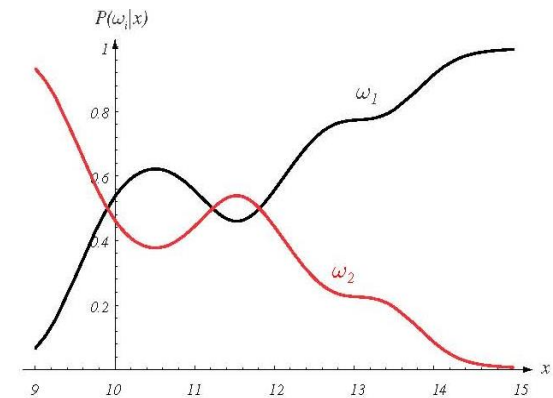


FIGURE 2.2. Posterior probabilities for the particular priors $P(\omega_1) = 2/3$ and $P(\omega_2) = 1/3$ for the class-conditional probability densities shown in Fig. 2.1. Thus in this case, given that a pattern is measured to have feature value $x = 14$, the probability it is in category ω_2 is roughly 0.08, and that it is in ω_1 is 0.92. At every x , the posteriors sum to 1.0. From: Richard O. Duda, Peter E. Hart, and David G. Stork, *Pattern Classification*. Copyright © 2001 by John Wiley & Sons, Inc.

Probability of Error

- The probability of error is defined as:

$$P(error / x) = \begin{cases} P(\omega_1 / x) & \text{if we decide } \omega_2 \\ P(\omega_2 / x) & \text{if we decide } \omega_1 \end{cases}$$

or $P(error/x) = \min[P(\omega_1/x), P(\omega_2/x)]$

- What is the **average probability error**?

$$P(error) = \int_{-\infty}^{\infty} P(error, x) dx = \int_{-\infty}^{\infty} P(error / x) p(x) dx$$

- The Bayes rule is **optimum**, that is, it minimizes the average probability error!

Where do Probabilities Come From?

- There are two competitive answers to this question:
 - (1) **Relative frequency (objective)** approach.
 - Probabilities can only come from experiments.
 - (2) **Bayesian (subjective)** approach.
 - Probabilities may reflect degree of belief and can be based on opinion.

Prior and Posterior Probabilities

- $P(A)$ and $P(B)$ are called prior probabilities
- $P(A|B)$, $P(B|A)$ are called posterior probabilities

Example 8.6: Prior versus Posterior Probabilities

- This table shows that the event Y has two outcomes namely A and B , which is dependent on another event X with various outcomes like x_1 , x_2 and x_3 .
- **Case1:** Suppose, we don't have any information of the event A . Then, from the given sample space, we can calculate $P(Y = A) = \frac{5}{10} = 0.5$.
- **Case2:** Now, suppose, we want to calculate $P(X = x_2/Y = A) = \frac{2}{5} = 0.4$.

The later is the conditional or posterior probability, where as the former is the prior probability.

X	Y
x_1	A
x_2	A
x_3	B
x_3	A
x_2	B
x_1	A
x_1	B
x_3	B
x_2	B
x_2	A

Example (objective approach)

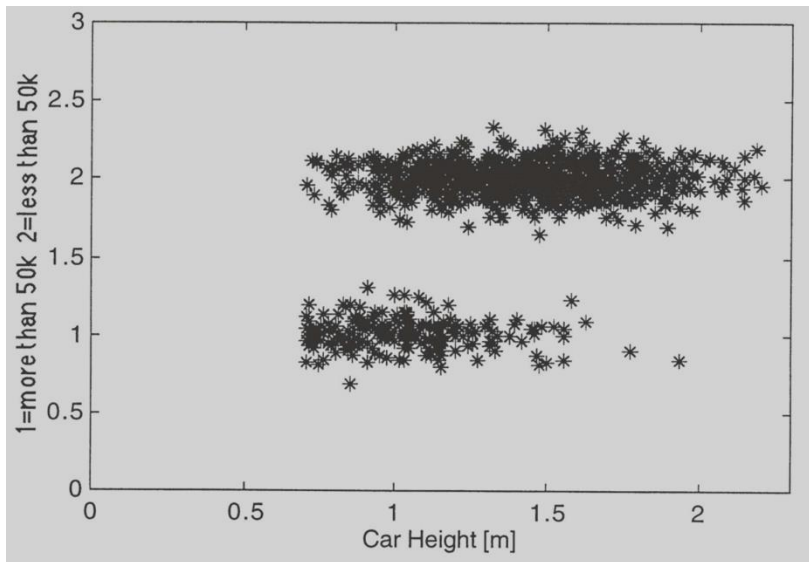
- Classify cars whether they are more or less than \$50K:
 - Classes: C_1 if price > \$50K, C_2 if price ≤ \$50K
 - Features: x , the height of a car
- Use the Bayes' rule to compute the posterior probabilities:

$$P(C_i / x) = \frac{p(x / C_i) P(C_i)}{p(x)}$$

- We need to estimate $p(x/C_1)$, $p(x/C_2)$, $P(C_1)$, $P(C_2)$

Example (cont'd)

- Collect data
 - Ask drivers how much their car was and measure height.
- Determine **prior** probabilities $P(C_1)$, $P(C_2)$
 - e.g., 1209 samples: $\#C_1=221$ $\#C_2=988$



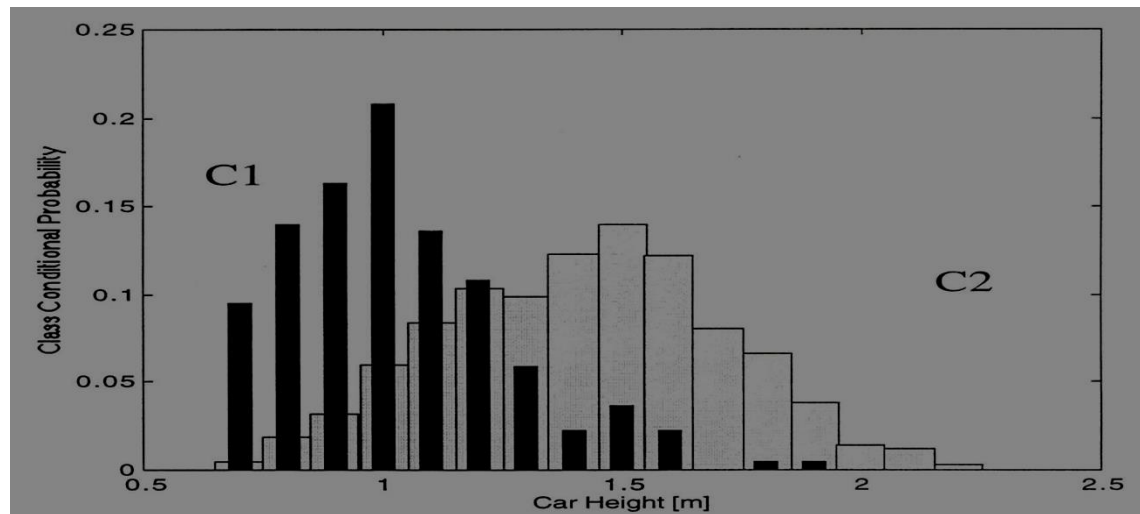
$$P(C_1) = \frac{221}{1209} = 0.183$$

$$P(C_2) = \frac{988}{1209} = 0.817$$

Example (cont'd)

- Determine **class conditional probabilities** (*likelihood*)
 - Discretize car height into bins and use normalized histogram

$$p(x / C_i)$$



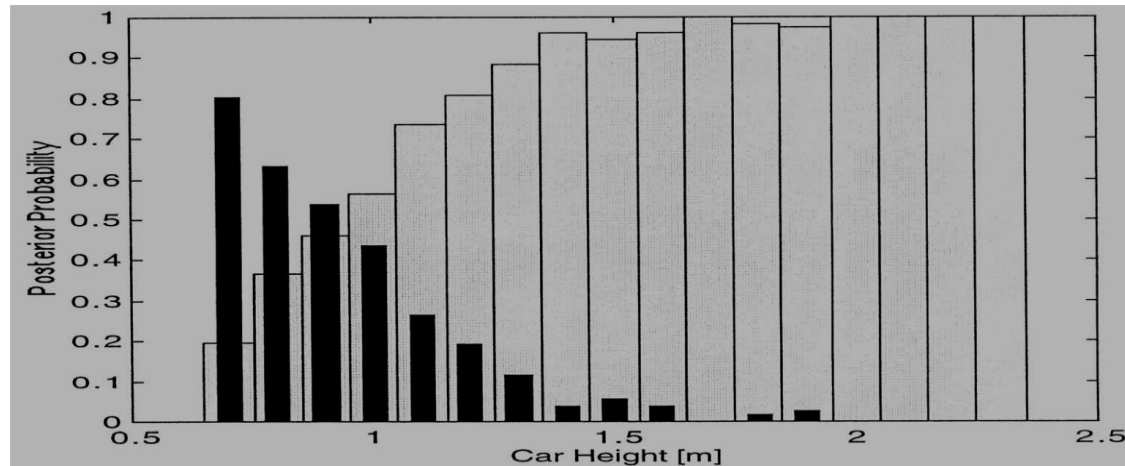
Example (cont'd)

- Calculate the **posterior** probability for each bin:

$$P(C_1 / x = 1.0) = \frac{p(x = 1.0 / C_1) P(C_1)}{p(x = 1.0 / C_1) P(C_1) + p(x = 1.0 / C_2) P(C_2)} =$$

$$= \frac{0.2081 * 0.183}{0.2081 * 0.183 + 0.0597 * 0.817} = 0.438$$

$P(C_i / x)$



A More General 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 (i.e., “**risk**” function) by associating a “**cost**” (“**loss**” function) with each error (i.e., wrong action).

Terminology

- Features form a vector $\mathbf{x} \in R^d$
- A finite set of c categories $\omega_1, \omega_2, \dots, \omega_c$
- Bayes rule (i.e., using vector notation):

$$P(\omega_j / \mathbf{x}) = \frac{p(\mathbf{x} / \omega_j)P(\omega_j)}{p(\mathbf{x})} \quad \text{where} \quad p(\mathbf{x}) = \sum_{j=1}^c p(\mathbf{x} / \omega_j)P(\omega_j)$$

- A finite set of l actions $\alpha_1, \alpha_2, \dots, \alpha_l$
- A loss function $\lambda(\alpha_i / \omega_j)$
 - the cost associated with taking action α_i when the correct classification category is ω_j .
 - For example, in case of binary classification if we take an action of classifying an input feature vector into class 1 when it should have been in class 2, we incur the loss $\lambda(\alpha_1 / \omega_2)$.
 - If $i = j$, then we get a smaller value of the loss as compared to the alternative cases because it corresponds to a correct decision.

Conditional Risk (or Expected Loss)

- Suppose we observe \mathbf{x} and take action α_i
- Suppose that the cost associated with taking action α_i with ω_j being the correct category is $\lambda(\alpha_i / \omega_j)$
- The **conditional risk** (or **expected loss**) with taking action α_i is:

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

Overall Risk

- Suppose $a(\mathbf{x})$ is a general decision rule that determines which action $\alpha_1, \alpha_2, \dots, \alpha_l$ to take for every \mathbf{x} ; then the overall risk is defined as:

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

- The optimum decision rule is the *Bayes rule*

Overall Risk (cont'd)

- The *Bayes decision rule* minimizes R by:
 - (i) Computing $R(\alpha_i/\mathbf{x})$ for every α_i given an \mathbf{x}
 - (ii) 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:

$$R^* = \min R$$

Example: Two-category classification

- Define
 - α_1 : decide ω_1 (c=2)
 - α_2 : decide ω_2
 - $\lambda_{ij} = \lambda(\alpha_i / \omega_j)$
- The conditional risks are:

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



$$\begin{aligned} R(a_1 / \mathbf{x}) &= \lambda_{11} P(\omega_1 / \mathbf{x}) + \lambda_{12} P(\omega_2 / \mathbf{x}) \\ R(a_2 / \mathbf{x}) &= \lambda_{21} P(\omega_1 / \mathbf{x}) + \lambda_{22} P(\omega_2 / \mathbf{x}) \end{aligned}$$

Example: Two-category classification (cont'd)

- Minimum risk decision rule:

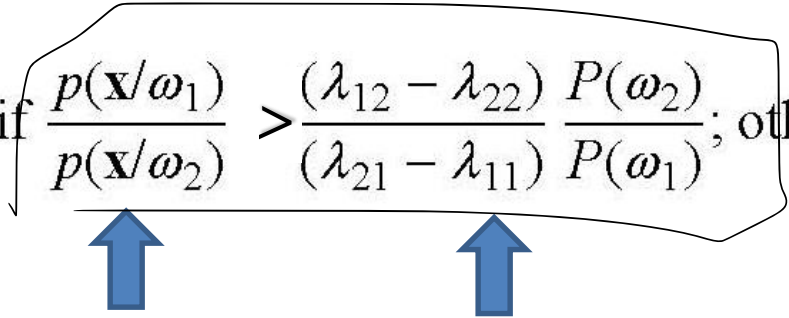
Decide ω_1 if $R(a_1/\mathbf{x}) < R(a_2/\mathbf{x})$; otherwise decide ω_2

or

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

or (i.e., using likelihood ratio)

Decide ω_1 if $\frac{p(\mathbf{x}/\omega_1)}{p(\mathbf{x}/\omega_2)} > \frac{(\lambda_{12} - \lambda_{22})}{(\lambda_{21} - \lambda_{11})} \frac{P(\omega_2)}{P(\omega_1)}$; otherwise decide ω_2



The diagram shows the decision rule with a box around the inequality. Two blue arrows point upwards from the labels 'likelihood ratio' and 'threshold' to the terms $\frac{p(\mathbf{x}/\omega_1)}{p(\mathbf{x}/\omega_2)}$ and $\frac{(\lambda_{12} - \lambda_{22})}{(\lambda_{21} - \lambda_{11})} \frac{P(\omega_2)}{P(\omega_1)}$ respectively.

likelihood ratio threshold

Special Case: Zero-One Loss Function

- Assign the same loss to all errors:

$$\lambda(a_i/\omega_j) = \begin{cases} 0 & i = j \\ 1 & i \neq j \end{cases}$$

- The conditional risk corresponding to this loss function:

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

Special Case: Zero-One Loss Function (cont'd)

- The decision rule becomes:

Decide ω_1 if $R(a_1/\mathbf{x}) < R(a_2/\mathbf{x})$; otherwise decide ω_2

or Decide ω_1 if $1 - P(\omega_1/\mathbf{x}) < 1 - P(\omega_2/\mathbf{x})$; otherwise decide ω_2

or Decide ω_1 if $P(\omega_1/\mathbf{x}) > P(\omega_2/\mathbf{x})$; otherwise decide ω_2

- In this case, the **overall risk** is the **average probability error!**
-

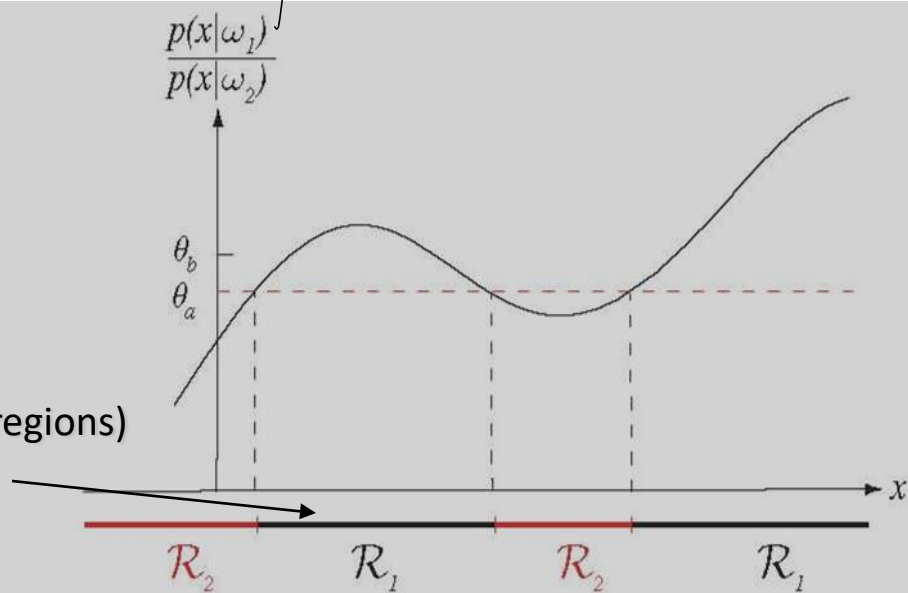
Example

Assuming **general** loss:

Decide ω_1 if $\frac{p(\mathbf{x}/\omega_1)}{p(\mathbf{x}/\omega_2)} > \frac{(\lambda_{12} - \lambda_{22})}{(\lambda_{21} - \lambda_{11})} \frac{P(\omega_2)}{P(\omega_1)}$; otherwise decide ω_2

Assuming **zero-one** loss:

Decide ω_1 if $p(x/\omega_1)/p(x/\omega_2) > P(\omega_2)/P(\omega_1)$ otherwise **decide ω_2**



$$\theta_a = P(\omega_2) / P(\omega_1)$$

$$\theta_b = \frac{P(\omega_2)(\lambda_{12} - \lambda_{22})}{P(\omega_1)(\lambda_{21} - \lambda_{11})}$$

assume: $\lambda_{12} > \lambda_{21}$