Basics of Machine Learning

Likelihood Ratio Test

Professor: Cheolsoo Park





What is Machine Learning (ML)?

- **Components of Machine Learning (ML)**
 - Computer Program learning experience (E)
 - Class of task (T) corresponding to E 7
 - Performance Measure (P) of Task
 - → Machine Learning Algorithm : Program (or algorithm) improving performance (P) of task (T) using the experience (E)
- ML finds a target function f between data $X = (x_1, x_2, ..., x_3)$ and states $Y = (y_1, y_2, ..., y_n)$ 7
- ML sets a hypothesis function, f', which best approximates target function f with some 7 performance measure
- **Problem Setting** 7
 - Set of possible instance(domain): X
 - Output: Y 7
 - Unknown target function $f: X \to Y$ 7
 - 7

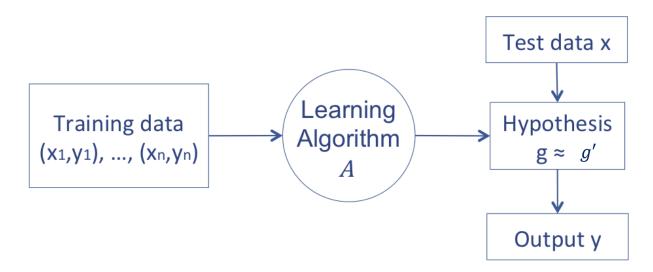
hypothesis function space. Set of hypothesis function space $H \subset \{g | g: X \to Y\}$ Thus limit it into a set to find

- **Input**: Training example $\{\langle x_i, y_i \rangle\}$
 - **Output**: $g' \in H$ that best approximates target function f with some performance measure

It is impossible to check all

What is Machine Learning (ML)?

- Learning using input data (training data)
 - **Supervised learning**: with label $y \{\langle x_i, y_i \rangle\}$
 - Unsupervised learning : without label y
- Return hypothesis, g' among $g \in H$, best approximating target function g'



Probability Theory

- ML— Probabilistic Perspective
 - ML sets a probability density function (PDF), rather than a deterministic function of hypothesis, and find the parameters of the PDF
 - Ex) Assume the data has a Gaussian PDF, and then mean and covariance need to be found
 - **7** function parameter → PDF parameter

Estimation

- In most case, the information of real probability density distribution can't be gotten
- 7 This must be determined by the test data
- Parameter Estimation
 - Estimate parameter about density by MLE (maximum likelihood estimation) method
- Non-parametric Density Estimation
 - **7** k-NNR

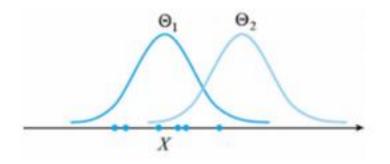
- Maximum Likelihood Estimation (MLE)
 - Estimate random variable's parameter using observation or data Ex) when flipping coin, we can get the p (probability of the front) by counting the number of the front among all trials
 - PDF $f = \{f(\cdot | \theta)\}$ and observation $X = (x_1, x_2, ..., x_n)$ If f is a Gaussian function, θ includes mean μ and covariance Σ . If f is a Bernoulli function, θ includes $0 \le p \le 1$

$$L(\theta; x_1, x_2, ..., x_n) = L(\theta; X) = f(X|\theta) = f(x_1, x_2, ..., x_n|\theta)$$

$$\hat{\theta} = \underset{\theta}{argmax} L(\theta; X) = \underset{\theta}{argmax} f(X|\theta)$$

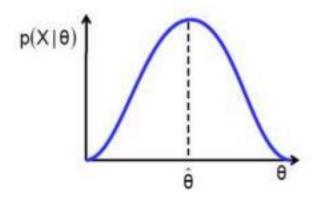
Find the Θ that has maximum likelihood about given X

P(X | Θ1) > P(X | Θ2)



$$\hat{\theta} = \arg \max[p(X \mid \theta)]$$

MAXIMUM LIKELIHOOD



- Estimate parameter θ=[θ1, θ2, ..., θN] using sample data X={x1,x2,...,xN} observed at probability density function
 - The entire sample set is expressed by the joint probability density

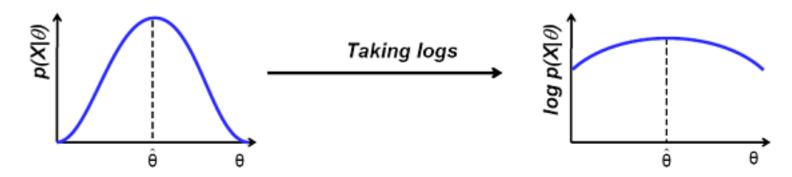
$$p(X \mid \theta) = p(x^{(1)} \mid \theta) p(x^{(2)} \mid \theta) ... p(x^{(N)} \mid \theta) = \prod_{k=1}^{N} p(x^{(k)} \mid \theta)$$

- Θ , having the highest probability, becomes the estimated parameter
- Take the log

$$\hat{\boldsymbol{\theta}} = argmax \Bigg[log \prod_{k=1}^{N} p \Big(\boldsymbol{x}^{(k} \mid \boldsymbol{\theta} \Big) \Bigg] = argmax \Bigg[\sum_{k=1}^{N} log p \Big(\boldsymbol{x}^{(k} \mid \boldsymbol{\theta} \Big) \Bigg]$$

- The log function is convenient to calculate
 - Maximize ∑(sum) is easier than ∏(product)
 - If distribution is similar to gaussian, increase the efficiency

$$\hat{\theta} = \operatorname{argmax}[p(X | \theta)] = \operatorname{argmax}[\log p(X | \theta)]$$



Probability density function p(x) = N(μ, σ).
 If X={x1,x2,...,xN} and σ is fixed. Get the MLE of μ.

$$\begin{split} \theta = \mu \Rightarrow \hat{\theta} = argmax \sum_{k=1}^{N} logp \Big(x^{(k)} \mid \theta \Big) \\ = argmax \sum_{k=1}^{N} log \Bigg(\frac{1}{\sqrt{2\pi}\sigma} exp \bigg(-\frac{1}{2\sigma^2} \big(x^{(k)} - \mu \big)^2 \bigg) \bigg) \\ = argmax \sum_{k=1}^{N} \bigg\{ log \bigg(\frac{1}{\sqrt{2\pi}\sigma} \bigg) - \frac{1}{2\sigma^2} \big(x^{(k)} - \mu \big)^2 \bigg\} \end{split}$$

Maximum point

$$\frac{\partial}{\partial \mu} \sum_{k=1}^{N} \left\{ \log \left(\frac{1}{\sqrt{2\pi\sigma}} \right) - \frac{1}{2\sigma^{2}} (x^{(k} - \mu)^{2}) \right\} = -\frac{1}{2\sigma^{2}} \frac{\partial}{\partial \mu} \sum_{k=1}^{N} (x^{2(k} - 2x^{(k}\mu + \mu^{2}))^{2}) \\
= \frac{1}{2\sigma^{2}} \frac{\partial}{\partial \mu} \sum_{k=1}^{N} (2x^{(k}\mu - \mu^{2})) = \sum_{k=1}^{N} (x^{(k} - \mu)) = \sum_{k=1}^{N} x^{(k} - N\mu) = 0 \Rightarrow \mu = \frac{1}{N} \sum_{k=1}^{N} x^{(k} - N\mu) = 0 \Rightarrow \mu = \frac{1}{N} \sum_{k=1}^{N} x^{(k} - \mu) = 0 \Rightarrow \mu = \frac{1}{N} \sum_{k=1}^{N} x^{(k} - \mu) = 0 \Rightarrow \mu = \frac{1}{N} \sum_{k=1}^{N} x^{(k} - \mu) = 0 \Rightarrow \mu = \frac{1}{N} \sum_{k=1}^{N} x^{(k} - \mu) = 0 \Rightarrow \mu = \frac{1}{N} \sum_{k=1}^{N} x^{(k} - \mu) = 0 \Rightarrow \mu = \frac{1}{N} \sum_{k=1}^{N} x^{(k} - \mu) = 0 \Rightarrow \mu = \frac{1}{N} \sum_{k=1}^{N} x^{(k} - \mu) = 0 \Rightarrow \mu = \frac{1}{N} \sum_{k=1}^{N} x^{(k} - \mu) = 0 \Rightarrow \mu = 0$$

MLE of μ is average of X.

Probability density function p(x) = N(μ, σ).
 If X={x1,x2,...,xN} and σ is fixed. Get the MLE of μ and σ.

$$\hat{\boldsymbol{\theta}} = \begin{bmatrix} \boldsymbol{\theta}_1 = \boldsymbol{\mu} \\ \boldsymbol{\theta}_2 = \boldsymbol{\sigma}^2 \end{bmatrix} \Rightarrow \nabla_{\boldsymbol{\theta}} = \begin{bmatrix} \frac{\partial}{\partial \boldsymbol{\theta}_1} \sum_{k=1}^{N} logp(\boldsymbol{x}^{(k} \mid \boldsymbol{\theta})) \\ \frac{\partial}{\partial \boldsymbol{\theta}_2} \sum_{k=1}^{N} logp(\boldsymbol{x}^{(k} \mid \boldsymbol{\theta})) \end{bmatrix} = \sum_{k=1}^{N} \begin{bmatrix} \frac{1}{\boldsymbol{\theta}_2} (\boldsymbol{x}^{(k} - \boldsymbol{\theta}_1)) \\ -\frac{1}{2\boldsymbol{\theta}_2} + \frac{(\boldsymbol{x}^{(k} - \boldsymbol{\theta}_1)^2}{2\boldsymbol{\theta}_2^2} \end{bmatrix} = 0$$

• Simplify them with respect to θ_1 and θ_2

$$\hat{\theta}_1 = \frac{1}{N} \sum_{k=1}^{N} x^{(k)}; \ \hat{\theta}_2 = \frac{1}{N} \sum_{k=1}^{N} (x^{(k)} - \hat{\theta}_1)^2$$

- Disadvantage of MLE
 - Too sensitive to the observation

EX) when we only get the front of coin during flipping coin trials, then p=1 by MLE

Solution

Maximum a Posteriori Estimation (MAP)

Maximum a Posteriori Estimation (MAP)

- $\mathsf{MLE}: \hat{\theta} = \underset{\theta}{argmax} f(X|Y,\theta)$
- $\mathsf{MAP}: \hat{\theta} = \underset{\theta}{argmax} f(Y|X, \theta)$
- While MLE produces parameters based on the current observation, MAP uses general knowledge about the data as well as the observation (prior knowledge)
- EX) Flipping coin example
- → MAP uses Bayes' Theorem

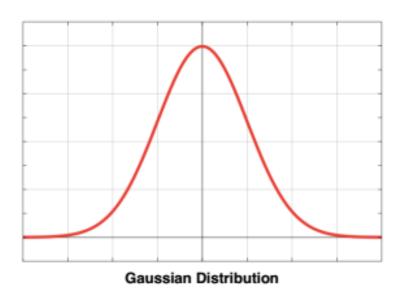
$$\hat{\theta} = \underset{\theta}{\operatorname{argmax}} f(Y|X,\theta) = \underset{\theta}{\operatorname{argmax}} \frac{f(X|Y,\theta)f(Y)}{f(X|\theta)}$$

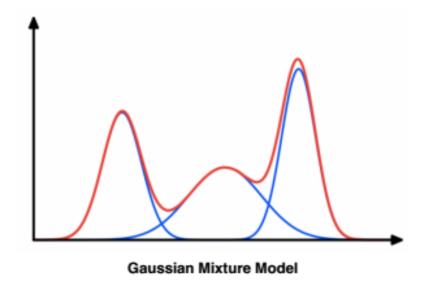
$$\Rightarrow \hat{\theta} = \underset{\theta}{\operatorname{argmax}} f(X|Y,\theta)f(Y)$$

 $f(\theta)$ is a prior knowledge or prior assumption

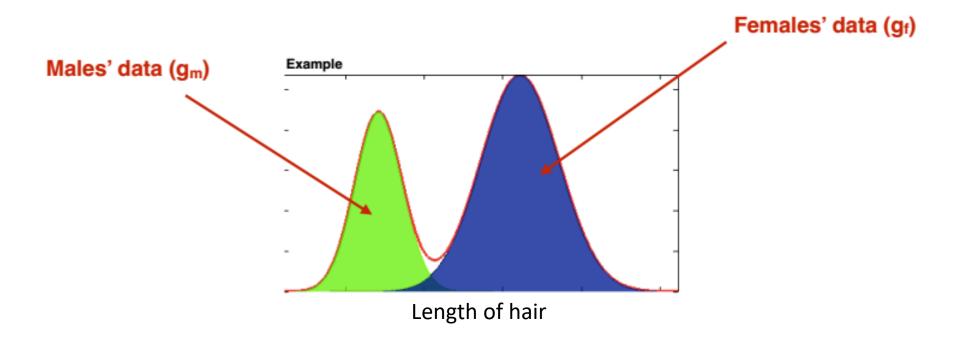
Maximum a Posteriori Estimation (MAP)

- Prior knowledge could be prejudice
- The results are depending on the prior, and thus it is very important to set the prior probability function carefully



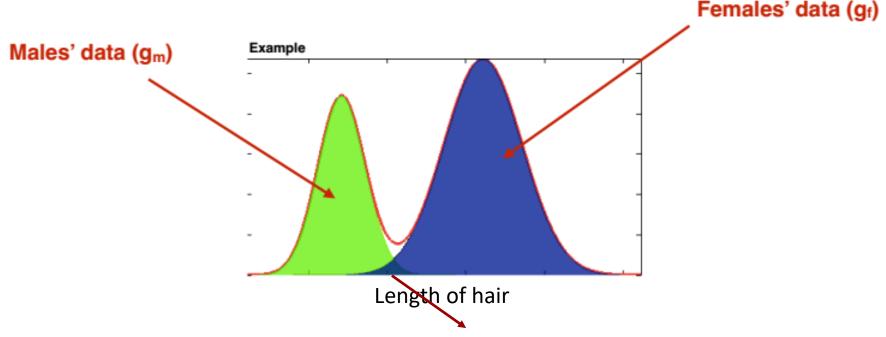


* PDF of hair length of the students in our class



https://www.youtube.com/watch?v=U8mjbpqsOTI

PDF of hair length of the students in our class



How can we decide whether this is for male or female?

• Assume we are going to classify an object into classes based on the evidence provided by a measurement (or feature vector) \mathbf{x}

- Choose a class most 'probable' given the observed feature vector x
- In formal

Evaluate the posterior probability of each class $P(w_i|x)$ and choose a class with the largest $P(w_i|x)$

- 2 class classification example
- The decision rule

if
$$P(w_1|x) > P(w_2|x)$$
, choose w_1 else, choose w_2

Or

$$P(w_1|x) \underset{w_2}{\overset{w_1}{\gtrless}} P(w_2|x)$$

Using Bayes theorem

$$\frac{P(x|w_1)P(w_1)}{P(x)} \underset{w_2}{\gtrless} \frac{P(x|w_2)P(w_2)}{P(x)} \longrightarrow \Lambda(x) = \frac{P(x|w_1)}{P(x|w_2)} \underset{w_2}{\gtrless} \frac{P(w_2)}{P(w_1)}$$

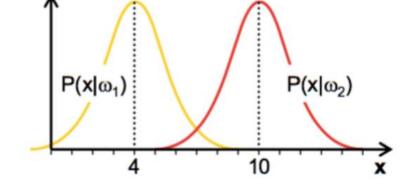
Likelihood ratio test

Example

Derive a classification rule using the likelihood ratio test with assuming the priors of two classes are equal

$$P(x|w_1) = \frac{1}{\sqrt{2\pi}} exp\{-\frac{1}{2}(x-4)^2\}$$

$$P(x|w_2) = \frac{1}{\sqrt{2\pi}} exp\{-\frac{1}{2}(x-10)^2\}$$

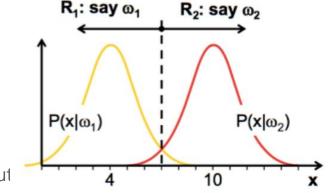


$$\Lambda(x) = \frac{\frac{1}{\sqrt{2\pi}} \exp\{-\frac{1}{2}(x-4)^2\}}{\frac{1}{\sqrt{2\pi}} \exp\{-\frac{1}{2}(x-10)^2\}} \underset{w_2}{\overset{w_1}{\geq}} \frac{1}{1}$$

$$(x-4)^2 - (x-10)^2 \underset{w_2}{\gtrless} 0 \implies x \underset{w_2}{\lessgtr} 7$$

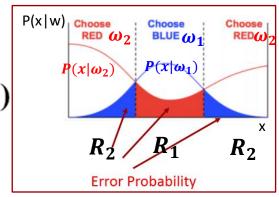






Decided by the borders

$$\begin{aligned} P(\text{error}) &= P(x \in R_2, \omega_1) + P(x \in R_1, \omega_2) \\ &= P(x \in R_2 \mid \omega_1) P(\omega_1) + P(x \in R_1 \mid \omega_2) P(\omega_2) \\ &= \int p(x \mid \omega_1) P(\omega_1) dx + \int p(x \mid \omega_2) P(\omega_2) dx \end{aligned}$$



$$P(error) = \sum_{i=1}^{C} P(error \mid \omega_i) P(\omega_i)$$

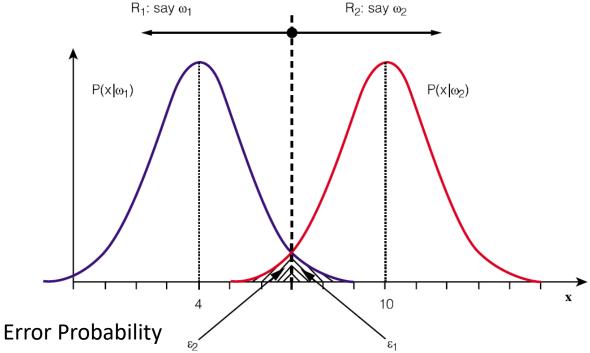
, where

P(error
$$|\omega_i|$$
) = P(choose $\omega_j |\omega_i|$) = $\int_{R_i} P(x |\omega_i) dx$

Error Probability

$$P[error] = P[\omega_1] \underbrace{\int_{R_2} P(x \mid \omega_1) dx + P[\omega_2] \int_{R_1} p(x \mid \omega_2) dx}_{\epsilon_1}$$

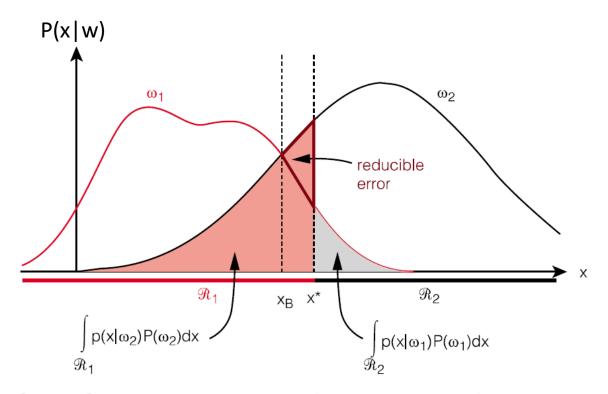
P(error) = $(\varepsilon_1 + \varepsilon_2)/2$, given $P(w_1) = P(w_2) = \frac{1}{2}$



nputing & Machine Learning (BCML) Lab

Error Probability

■ The optimal border = minimum error probability



[그림 5-3] 오류확률에 의한 결정 경계의 결정 [Duda, Hart, Stock, 2001]

Loss (or Cost) function

- In the previous slides, the error was 0 or 1 (misclassification number) when the classification was correct or wrong.
- However, the number of error can be a summation of squared difference between target data and estimated data. (This is one example)
 - → Cost function or Lost function
- For example, let's think about buying a wrong plane ticket to find a person in the world
- While true value is constant, the estimated value can be changed depending on the updated parameters

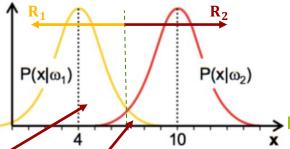
$$L(\theta, \hat{\theta})$$
: Loss function

 $(\theta : \text{true parameter}, \hat{\theta} : \text{estimated parameter})$

- So far we have assumed that the penalty of misclassifying a class w_1 as class w2 is the same as the opposite way $(w_2 \rightarrow w_1)$
- In general, this is not the case:
 - For example, misclassifying a cancer sufferer as a healthy patient is a much more serious problem than the other way around
- 7 This concept can be also formalised in terms of a cost (or lost) function C_{ij}
 - C_{ij} represents the cost of choosing (wrong) class w_i when class w_j is the true class
- We define the Bayes Risk as the expected value of the cost

$$\Re = E[C] = \sum_{i=1}^{2} \sum_{j=1}^{2} C_{ij} \cdot P(choose \ w_i \ and \ x \in w_j) = \sum_{i=1}^{2} \sum_{j=1}^{2} C_{ij} \cdot P(x \in R_i | w_j) \cdot P(w_j)$$

$$\begin{split} R = E[C] = \sum_{i=1}^2 \sum_{j=1}^2 C_{ij} \cdot P(choose \ \omega_i \ , \ x \in \omega_j) = \sum_{i=1}^2 \sum_{j=1}^2 C_{ij} \cdot P(x \in R_i \mid \omega_j) \cdot P(\omega_j) \\ \text{, where } P(x \in R_i \mid \omega_j) = \int\limits_{R_i} P(x \mid \omega_j) dx \\ R = \int\limits_{R_1} [C_{11}P(\omega_1)P(x \mid \omega_1) + C_{12}P(\omega_2)P(x \mid \omega_2)] dx \\ \int\limits_{R_1} [C_{21}P(\omega_1)P(x \mid \omega_1) + C_{22}P(\omega_2)P(x \mid \omega_2)] dx \end{split}$$



Decision rule to reduce the Bayes risk

$$\int_{R_1} P(x \mid \omega_i) dx + \int_{R_2} P(x \mid \omega_i) dx = \int_{R_1 \cup R_2} P(x \mid \omega_i) dx = 1$$

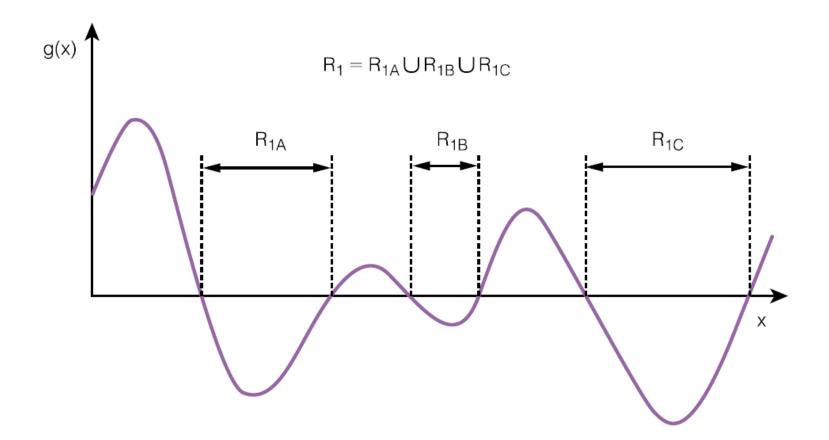
$$\begin{split} R = & C_{11}P(\omega_{1})\int\limits_{R_{1}}P(x\,|\,\omega_{1})dx \\ + & C_{21}P(\omega_{1})\int\limits_{R_{2}}P(x\,|\,\omega_{1})dx \\ + & C_{21}P(\omega_{1})\int\limits_{R_{2}}P(x\,|\,\omega_{1})dx \\ + & C_{21}P(\omega_{1})\int\limits_{R_{1}}P(x\,|\,\omega_{1})dx \\ - & C_{21}P(\omega_{1})\int\limits_{R_{1}}P(x\,|\,\omega_{1})dx \\ - & C_{21}P(\omega_{1})\int\limits_{R_{1}}P(x\,|\,\omega_{1})dx \\ - & C_{22}P(\omega_{2})\int\limits_{R_{1}}P(x\,|\,\omega_{2})dx \\ - & C_{22}P(\omega_{2})\int\limits_{R_{1}}P(x\,|\,\omega_{2})dx \\ \end{split}$$

$$R = \begin{bmatrix} C_{21}P(\omega_{1}) & + & C_{22}P(\omega_{2}) \\ + & (C_{12} - C_{22})P(\omega_{2}) \int_{R_{1}} P(x \mid \omega_{2}) dx \\ - & (C_{21} - C_{11})P(\omega_{1}) \int_{R_{1}} P(x \mid \omega_{1}) dx \end{bmatrix}$$

earning (BCML) Lab

$$\begin{split} R_{1} &= \underset{R_{1}}{arg \, min} \left\{ \int\limits_{R_{1}} [(C_{12} - C_{22}) P[\omega_{2}] P(x \, | \, \omega_{2}) - (C_{21} - C_{11}) P[\omega_{1}] P(x \, | \, \omega_{1})] dx \right\} \\ &= \underset{R_{1}}{arg \, min} \left\{ \int\limits_{R_{1}} g(x) dx \right\} \end{split}$$

Find the decision area based on the Bayes risk



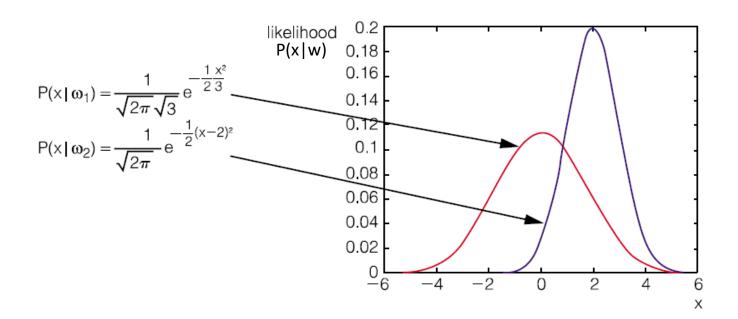
$$R_{1} = \arg\min \left\{ \int_{R_{1}} [(C_{12} - C_{22})P[\omega_{2}]P(x \mid \omega_{2}) - (C_{21} - C_{11})P[\omega_{1}]P(x \mid \omega_{1})]dx \right\}$$

$$= \arg\min \left\{ \int_{R_{1}} g(x)dx \right\}$$

$$(C_{12}-C_{22})P[\omega_{2}]P(x|\omega_{2})$$
 $< (C_{21}-C_{11})P[\omega_{1}]P(x|\omega_{1})$

$$\frac{P(x \mid \omega_{1})}{P(x \mid \omega_{2})} > \frac{(C_{12} - C_{22})P[\omega_{2}]}{(C_{21} - C_{11})P[\omega_{1}]}$$

Ex) Given the likelihood functions below for two classes, their prior probabilities, $P(w_1) = P(w_2) = 0.5$, and costs $C_{11} = C_{22} = 0$, $C_{12} = 1$, $C_{21} = 3^{\frac{1}{2}}$, find the likelihood ratio based on the Bayes risk and decision boundary to minimize the error probability



Bio Computing & Machine Learning (BCML) Lab

Minimise Bayes Risk

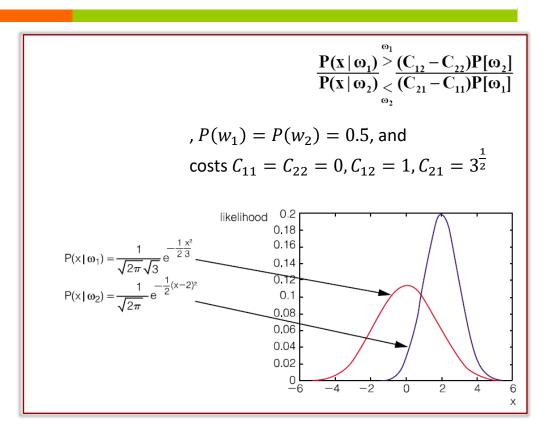
Likelihood Ratio

$$\Lambda(x) = \frac{\frac{1}{\sqrt{2\pi}\sqrt{3}} e^{-\frac{1}{2}\frac{x^2}{3}} \omega_1}{\frac{1}{\sqrt{2\pi}} e^{-\frac{1}{2}(x-2)^2}} > \frac{1}{\omega_2}$$

$$\frac{e^{-\frac{1}{2}\frac{x^2}{3}}}{e^{-\frac{1}{2}(x-2)^2}} > 1$$

$$-\frac{1}{2}\frac{x^2}{3}+\frac{1}{2}(x-2)^2 > 0$$

$$2x^{2}-12x+12 > 0 \Rightarrow x = 4.73, 1.27$$

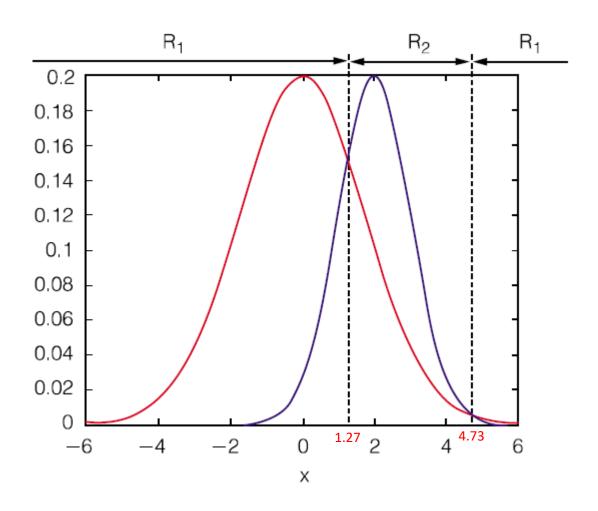


$$\therefore \omega_1 < 1.27, \omega_1 > 4.73$$

 $1.27 < \omega_2 < 4.73$

Bio Computing & Machine Learning (BCML) Lab

Minimise Bayes Risk



Variations of LRT Decision Rule

- Bayes' Criterion
 - → LRT decision rule with reducing the Bayes' risk

$$\Lambda(\mathbf{x}) = \frac{P(\mathbf{x} \mid \omega_1)}{P(\mathbf{x} \mid \omega_2)} > \frac{(C_{12} - C_{22})P[\omega_2]}{(C_{21} - C_{11})P[\omega_1]}$$

- MAP Criterion
 - Symmetric or zero-one cost functions make Bayes' criterion a ratio of posteriori functions → called maximum a posteriori (MAP) criterion

$$C_{ij} = \begin{cases} 0 & i = j \\ 1 & i \neq j \end{cases} \Rightarrow \Lambda(x) = \frac{P(x \mid \omega_1)}{P(x \mid \omega_2)} \stackrel{\omega_1}{>} \frac{P[\omega_2]}{P[\omega_1]} \Leftrightarrow \frac{P(\omega_1 \mid x)}{P(\omega_2 \mid x)} \stackrel{\omega_1}{>} 1$$

Variations of LRT Decision Rule

- ML Criterion
 - Equal prior probabilities and zero-one cost functions make Bayes' criterion a ratio of likelihood functions → called maximum likelihood (ML) criterion

$$C_{ij} = \begin{cases} 0 & i = j \\ 1 & i \neq j \end{cases} \Rightarrow \Lambda(x) = \frac{P(x \mid \omega_1)}{P(x \mid \omega_2)} > 1$$

$$P(\omega_i) = \frac{1}{C} \forall i$$