

# Lecture 5

## Basic Concepts in Bayesian Decision and Information Theory

Instructor: Shibo Li

shiboli@cs.fsu.edu

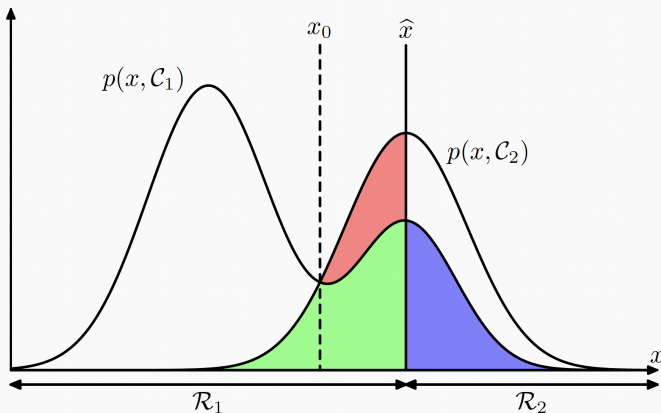


Department of Computer Science  
Florida State University

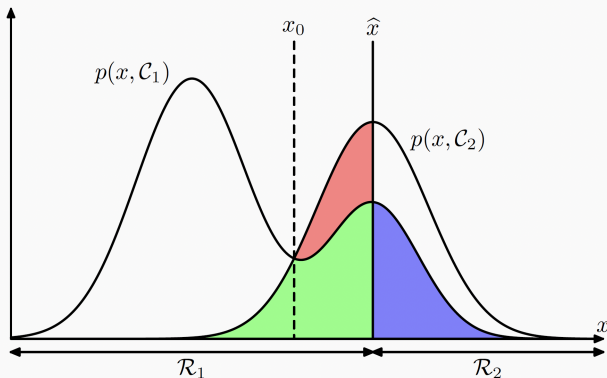
$\mathbf{t}$ : Cancer, Stock price, Weather ...

- Inference step
  - Determine either  $p(\mathbf{t}|\mathbf{x})$  or  $p(\mathbf{x},\mathbf{t})$  (from training data)
- Decision Step
  - For Given  $\mathbf{x}$ , determine optimal  $\mathbf{t}$

- $\mathbf{t} \in \{C_1, \dots, C_K\}$
- Decision regions  $R_k$ : if  $\mathbf{x}$  falls in , predict  $C_k$
- Decision boundaries/surfaces: boundaries between different decision regions




$$\begin{aligned}
 p(\text{mistake}) &= p(\mathbf{x} \in \mathcal{R}_1, \mathcal{C}_2) + p(\mathbf{x} \in \mathcal{R}_2, \mathcal{C}_1) \\
 &= \int_{\mathcal{R}_1} p(\mathbf{x}, \mathcal{C}_2) d\mathbf{x} + \int_{\mathcal{R}_2} p(\mathbf{x}, \mathcal{C}_1) d\mathbf{x}.
 \end{aligned}$$



Question: where shall we set the decision boundary to minimize the misclassification rate? Why?

- In general for  $K$  classes

$$\begin{aligned} p(\text{correct}) &= \sum_{k=1}^K p(\mathbf{x} \in \mathcal{R}_k, \mathcal{C}_k) \\ &= \sum_{k=1}^K \int_{\mathcal{R}_k} p(\mathbf{x}, \mathcal{C}_k) d\mathbf{x} \end{aligned}$$



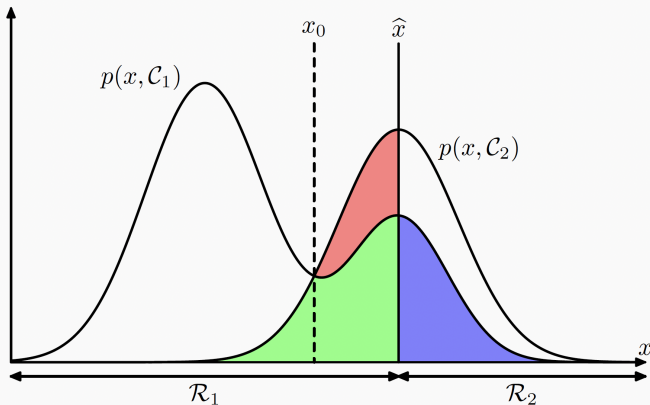
$p(\mathcal{C}_k | \mathbf{x}) p(\mathbf{x})$

How to find regions that maximize the probability of correctness?

- In general for  $K$  classes

$$\begin{aligned} p(\text{correct}) &= \sum_{k=1}^K p(\mathbf{x} \in \mathcal{R}_k, \mathcal{C}_k) \\ &= \sum_{k=1}^K \int_{\mathcal{R}_k} p(\mathbf{x}, \mathcal{C}_k) d\mathbf{x} \\ &\quad \quad \quad \updownarrow \\ &\quad \quad \quad \boxed{p(\mathcal{C}_k|\mathbf{x})p(\mathbf{x})} \end{aligned}$$

Each  $\mathbf{x}$  should be assigned the class having the largest posterior probability  $p(\mathcal{C}_k|\mathbf{x})$



- In practice, mistakes in predicting different classes may lead to different costs

Example: classify medical images as 'cancer' or 'normal'

		Decision	
		cancer	normal
Truth	cancer	0	1000
	normal	1	0

- Define a *cost function*, associate the cost of classifying  $k$  to  $j$  with  $L_{kj}$

$$\mathbb{E}[L] = \sum_k \sum_j \int_{\mathcal{R}_j} L_{kj} p(\mathbf{x}, \mathcal{C}_k) d\mathbf{x}.$$

- We want to find the decision regions  $R_j$  that minimize the expected loss

- Define a *cost function*, associate the cost of classifying  $k$  to  $j$  with  $L_{kj}$

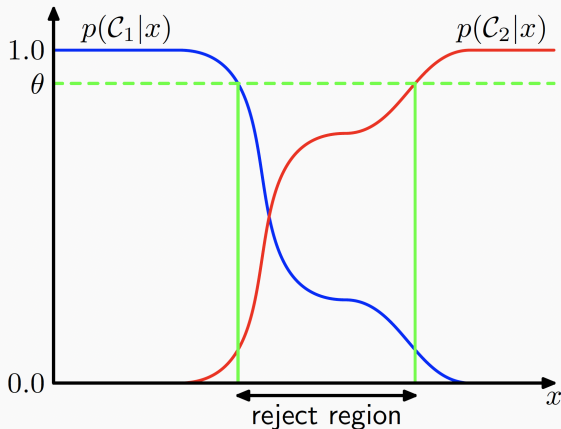
$$\mathbb{E}[L] = \sum_k \sum_j \int_{\mathcal{R}_j} L_{kj} p(\mathbf{x}, \mathcal{C}_k) d\mathbf{x}.$$

- Rule: Assign each  $\mathbf{x}$  to the class for which

$$\sum_k L_{kj} p(\mathcal{C}_k | \mathbf{x})$$

is a minimum

- When the largest posterior probability is still too small



- Inference step
  - Determine  $p(\mathbf{x}, t)$
- Decision step
  - For any given  $\mathbf{x}$ , make optimal prediction  $y(\mathbf{x})$  for  $t$

- Inference step
  - Determine  $p(\mathbf{x}, t)$
- Decision step
  - For any given  $\mathbf{x}$ , make optimal prediction  $y(\mathbf{x})$  for  $t$

Loss function

$$\mathbb{E}[L] = \iint L(t, y(\mathbf{x})) p(\mathbf{x}, t) d\mathbf{x} dt$$

Minimize  $\mathbb{E}[L] = \iint \{y(\mathbf{x}) - t\}^2 p(\mathbf{x}, t) \, d\mathbf{x} \, dt$

Minimize  $\mathbb{E}[L] = \iint \{y(\mathbf{x}) - t\}^2 p(\mathbf{x}, t) \, d\mathbf{x} \, dt$



$$y(\mathbf{x}) = \mathbb{E}[t|\mathbf{x}]$$

- What is the decision? What is the difference between the decision and inference?
- How to find optimal decision regions for classification?
- How to find optimal decisions for continuous variables?

- Let us start with discrete random variables

- How to represent the information contained in the random variables?

$$h(\mathbf{x}) \geq 0$$

$$h(\mathbf{x}, \mathbf{y}) = h(\mathbf{x}) + h(\mathbf{y}) \quad \mathbf{x}, \mathbf{y} \text{ are independent}$$

$$p(\mathbf{x}, \mathbf{y}) = p(\mathbf{x})p(\mathbf{y})$$



$$h(\mathbf{x}) = -\log(p(\mathbf{x}))$$

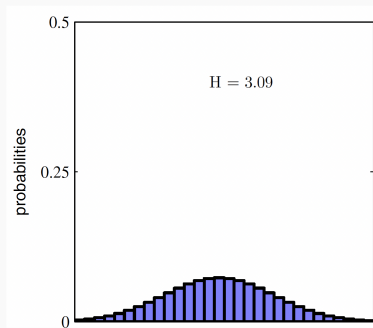
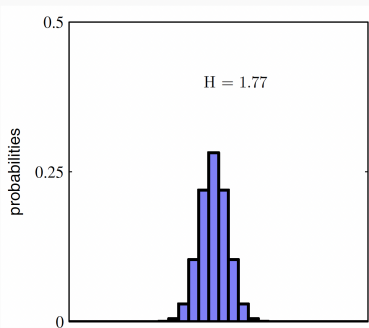
- The average amount of information needed to transmit

$$H(\mathbf{x}) = - \sum_{\mathbf{x}} p(\mathbf{x}) \log(p(\mathbf{x}))$$


$x$	a	b	c	d	e	f	g	h
$p(x)$	$\frac{1}{2}$	$\frac{1}{4}$	$\frac{1}{8}$	$\frac{1}{16}$	$\frac{1}{64}$	$\frac{1}{64}$	$\frac{1}{64}$	$\frac{1}{64}$


$$\begin{aligned}
 H[x] &= -\frac{1}{2} \log_2 \frac{1}{2} - \frac{1}{4} \log_2 \frac{1}{4} - \frac{1}{8} \log_2 \frac{1}{8} - \frac{1}{16} \log_2 \frac{1}{16} - \frac{4}{64} \log_2 \frac{1}{64} \\
 &= 2 \text{ bits}
 \end{aligned}$$

Entropy is also the average code length



- Consider a discrete R.V. with  $M$  possible status. We want to find the distribution has the the maximum entropy  $H[p] = - \sum_i p(x_i) \ln p(x_i)$ .


$$\tilde{H} = - \sum_i p(x_i) \ln p(x_i) + \lambda \left( \sum_i p(x_i) - 1 \right)$$

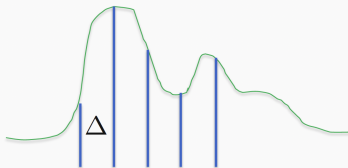

$$p(x_i) = 1/M \quad \text{uniform distribution}$$

- Entropy is naturally defined on discrete random variables.
- But how about continuous variables?

- Let us divide  $x$  into bins of  $\Delta$

*Mean-value theorem*

$$\int_{i\Delta}^{(i+1)\Delta} p(x) dx = p(x_i)\Delta$$



*Entropy on discretized probability*

$$H_{\Delta} = - \sum_i p(x_i)\Delta \ln(p(x_i)\Delta) = - \sum_i p(x_i)\Delta \ln p(x_i) - \ln \Delta$$

$$\sum_i p(x_i)\Delta = 1$$

$$H_{\Delta} = - \sum_i p(x_i) \Delta \ln (p(x_i) \Delta) = - \sum_i p(x_i) \Delta \ln p(x_i) - \ln \Delta$$

Goes to infinity  
Throw out it

$$\lim_{\Delta \rightarrow 0} \left\{ \sum_i p(x_i) \Delta \ln p(x_i) \right\} = \int p(x) \ln p(x) dx$$

$$H[\mathbf{x}] = - \int p(\mathbf{x}) \ln p(\mathbf{x}) d\mathbf{x}$$

- The term that is thrown out reflects that to specify a continuous variable very precisely requires many many bits
- Note: differential entropy can be negative!

- Given a continuous variable  $x$  with mean  $\mu$  and variance  $\sigma^2$ , which distribution has the largest entropy?

$$\int_{-\infty}^{\infty} p(x) dx = 1$$

$$\int_{-\infty}^{\infty} xp(x) dx = \mu$$

$$\int_{-\infty}^{\infty} (x - \mu)^2 p(x) dx = \sigma^2.$$

$$\begin{aligned} \max \quad & - \int_{-\infty}^{\infty} p(x) \ln p(x) \, dx + \lambda_1 \left( \int_{-\infty}^{\infty} p(x) \, dx - 1 \right) \\ & + \lambda_2 \left( \int_{-\infty}^{\infty} xp(x) \, dx - \mu \right) + \lambda_3 \left( \int_{-\infty}^{\infty} (x - \mu)^2 p(x) \, dx - \sigma^2 \right) \end{aligned}$$



$$p(x) = \frac{1}{(2\pi\sigma^2)^{1/2}} \exp \left\{ -\frac{(x - \mu)^2}{2\sigma^2} \right\} \quad \text{Gaussian distribution!}$$

- Given  $\mathbf{x}$ , how much information is left for  $\mathbf{y}$

$$H[\mathbf{y}|\mathbf{x}] = - \iint p(\mathbf{y}, \mathbf{x}) \ln p(\mathbf{y}|\mathbf{x}) \, d\mathbf{y} \, d\mathbf{x}$$

$$H[\mathbf{x}, \mathbf{y}] = H[\mathbf{y}|\mathbf{x}] + H[\mathbf{x}] \quad \text{Prove it by yourself}$$

- Also called relative entropy

$$\begin{aligned}\text{KL}(p\|q) &= - \int p(\mathbf{x}) \ln q(\mathbf{x}) \, d\mathbf{x} - \left( - \int p(\mathbf{x}) \ln p(\mathbf{x}) \, d\mathbf{x} \right) \\ &= - \int p(\mathbf{x}) \ln \left\{ \frac{q(\mathbf{x})}{p(\mathbf{x})} \right\} \, d\mathbf{x}.\end{aligned}$$

If we use  $q$  to transmit information for  $p$ , how much extra information do we need

- KL divergence is widely used to measure the difference between two distributions

$$\text{KL}(p\|q) \geq 0 \quad =0 \text{ iff } p = q$$

Prove it with convexity  
And Jensen's inequality

- However, it is not symmetric!

$$\text{KL}(p\|q) \neq \text{KL}(q\|p)$$

- KL divergence plays the key role in approximate inference
- All the deterministic approximate methods aim to minimize the KL divergence between the true and approximate posteriors (or in the reversed direction)
- In general, we have alpha divergence
- We will discuss these in detail later

How many information do the two random variables share?

$$\begin{aligned} I[\mathbf{x}, \mathbf{y}] &\equiv \text{KL}(p(\mathbf{x}, \mathbf{y}) \| p(\mathbf{x})p(\mathbf{y})) \\ &= - \iint p(\mathbf{x}, \mathbf{y}) \ln \left( \frac{p(\mathbf{x})p(\mathbf{y})}{p(\mathbf{x}, \mathbf{y})} \right) d\mathbf{x} d\mathbf{y} \end{aligned}$$



$$I[\mathbf{x}, \mathbf{y}] = H[\mathbf{x}] - H[\mathbf{x}|\mathbf{y}] = H[\mathbf{y}] - H[\mathbf{y}|\mathbf{x}]$$

Prove it by  
yourself

- Definition of entropy
- How is differential entropy is derived
- Entropy is an indicator for uncertainty
- KL divergence and properties (especially asymmetric)
- Mutual information