

# Algorithmic Learning Theory

## Spring 2017

### Lecture 2

**Instructor:** Farid Alizadeh

**Scribe:** Chien-Ming Huang

**Edit:** Yuan Qu

1/25/2017

1. Review Bayes Theory(Lecture 1)
2. Random Variable and Distribution
  - (a) Random variable
    - i. DRV, discrete random variable
    - ii. CRV, continuous random variable
  - (b) Distribution function
    - i. CDF, cumulative distribution function
    - ii. pdf or pmf, probability density(Mass) function
  - (c) Discrete distribution
    - i. Discrete uniform distribution
    - ii. Beunoulli's distribution
    - iii. Binomial distribution
  - (d) Continuous distribution
    - i. Continuous uniform distribution
    - ii. Normal distribution
3. Multivariate Distributions
  - (a) Random vector
  - (b) Discrete multivariate distribution
  - (c) Binormal distribution
  - (d) Marginal distribution
  - (e) Conditional distribution
4. Bayes Classification

## 1 Review Bayes Theory(Lecture 1)

See notes in "Lecture 1".

## 2 Random Variable and Distribution

### 2.1 Random Variable

#### 2.1.1 Discrete Random Variable(D.R.V)

$x \rightarrow t_1, t_2, \dots, t_n$

#### 2.1.2 Continuous Random Variable(C.R.V)

$x \rightarrow [a, b]$  a range of value.

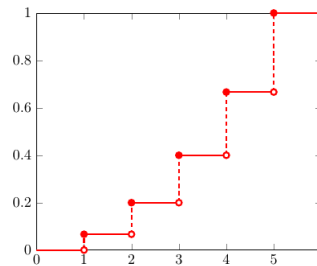
### 2.2 Distribution Function

#### 2.2.1 CDF:

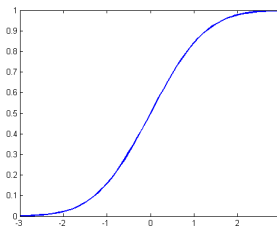
Cumulative Distribution Function

$F_x(t) = P[x \leq t]$ , probability can only increase

1. For Discrete:



2. For Continuous:

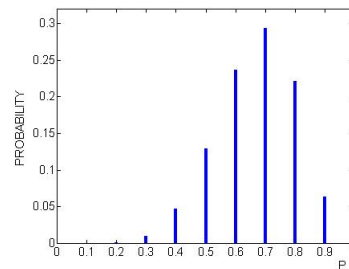


### 2.2.2 pdf or pmf:

Probability Density(Mass) Function

1. For Discrete:

$$x : F_x(t) = P[x = t]$$

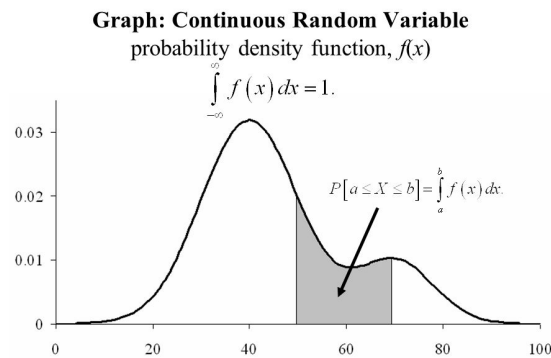


2. For Continuous:

$$f_x(t) = \frac{d}{dt} F_x(t), F_x(t) = \int_{-\infty}^t f_x(t) dt$$

$$i \ f_x(t) \geq 0$$

$$ii \ \int_{-\infty}^{\infty} f_x(t) dt = 1$$

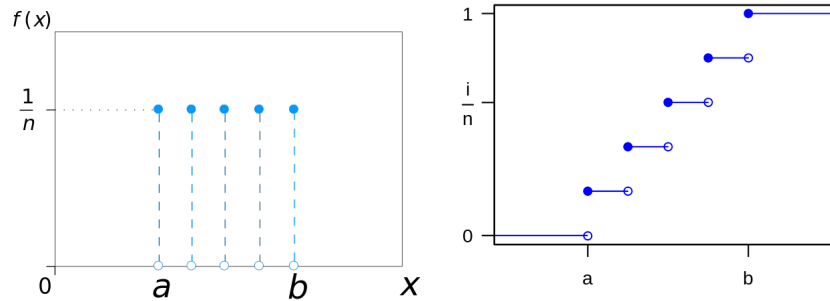


## 2.3 Discrete Distribution

### 2.3.1 Discrete Uniform Distribution

$$x : 1, 2, 3, \dots, k$$

$$\text{pdf} : u_x(t) = \begin{cases} \frac{1}{n}, & \text{if } t = 1, 2, \dots, n \\ 0, & \text{otherwise} \end{cases}$$



### 2.3.2 Bernoulli Distribution

$$\text{pdf: } f_x(t) = \begin{cases} p, & x = 1 \\ 1 - p, & x = 0 \\ 0, & \text{otherwise} \end{cases}$$

$$\text{CDF: } F_x(t) = \begin{cases} 0, & x \leq 0 \\ 1 - p, & 0 \leq x < 1 \\ 1, & x \geq 1 \end{cases}$$

### 2.3.3 Binomial Distribution

numbers of 0's in independent Bernoulli trial with  $P[0] = p$

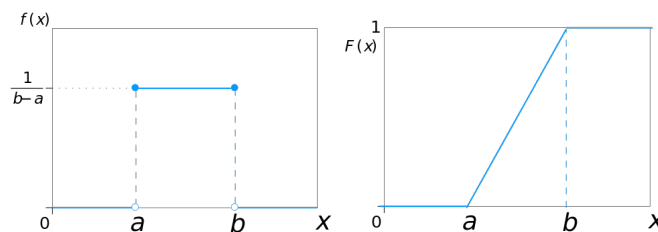
$$\text{pdf: } b(t|p, n) = \binom{n}{t} p^t (1-p)^{n-t}, \quad \binom{n}{t} = \frac{n!}{t!(n-t)!}$$

$$\text{CDF: } B(t|p, n) = \sum_{n=0}^t b(t|p, n)$$

## 2.4 Continuous distribution

### 2.4.1 Continuous uniform distribution

$$u(t|a, b) = \begin{cases} \frac{1}{b-a}, & a \leq t \leq b \\ 0, & \text{otherwise} \end{cases}$$



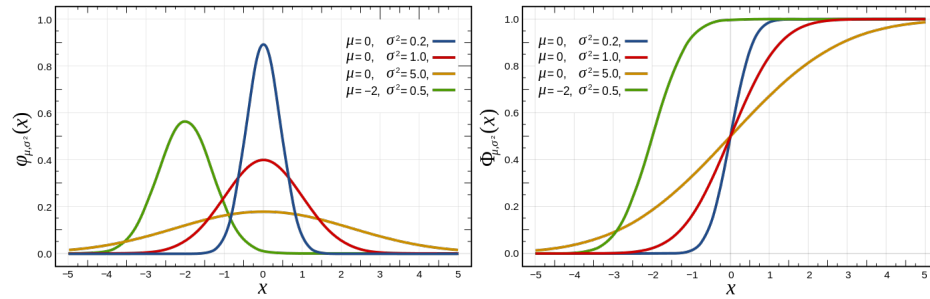
### 2.4.2 Normal distribution

mean =  $\mu$  and std. =  $\sigma$

$$\text{pdf: } f(x|\mu, \sigma^2) = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$

$$\text{CDF: } \frac{1}{2} [1 + \text{erf}(\frac{x-\mu}{\sigma\sqrt{2}})]$$

$$\text{erf}(x) = \frac{1}{\sqrt{\pi}} \int_{-x}^x e^{-t^2} dt$$



## 3 Multivariate Distributions

### 3.1 Random Vector

$X_N = (x_1, x_2, \dots, x_N)$  can be continuous or discrete.

$$\text{CDF: } F_X(t_1, t_2, \dots, t_n) = P[x_1 \leq t_1, x_2 \leq t_2, \dots, x_n \leq t_n]$$

$$\text{pdf: } \begin{cases} \frac{\partial}{\partial t_1 \partial t_2 \partial t_3 \dots \partial t_n} F(t_1, t_2, \dots, t_n) = f_X(t_1, \dots, t_n), & \text{all continuous} \\ P[x_1 = t_1, x_2 = t_2, \dots, x_n = t_n] = f_X(t_1, \dots, t_n), & \text{all discrete} \end{cases}$$

both are joint distribution R.V.  $x_1, \dots, x_n$

### 3.2 Discrete Multivariate Distribution

$$Y \rightarrow 1, 2, \dots, r; \quad P[Y = r_1] = P_u; \quad \sum P_u = 1$$

repeat  $n$  times,  $x_k$  = number of times  $Y = k$  occurs

$$\underline{x} = (x_1, \dots, x_n)$$

$$\text{pdf: } f_X(x_1, x_2, \dots, x_n) = P[x_1 = t_1, x_2 = t_2, \dots, x_n = t_n] \binom{n}{t_1, t_2, \dots, t_n} p_1^{t_1} p_2^{t_2} \dots p_n^{t_n}$$

$$\binom{n}{t_1, t_2, \dots, t_r} = \frac{n!}{t_1! t_2! \dots t_n!}$$

### 3.3 Binormal distribution

$\underline{x} = (x_1, x_2)$ , both continuous

$$\underline{\mu} = \begin{pmatrix} \mu_1 \\ \mu_2 \end{pmatrix}, \sigma_1^2 \rightarrow x_1, \sigma_2^2 \rightarrow x_2, \sigma_{21} \rightarrow x_1, x_2$$

$$\text{Covariance Matrix: } \Sigma = \begin{pmatrix} \sigma_1^2 & \sigma_{12} \\ \sigma_{12} & \sigma_2^2 \end{pmatrix}$$

$$\phi(t_1, t_2 | \underline{\mu}, \Sigma) = \frac{1}{\sqrt{2\pi \cdot \text{Det}(\Sigma)}} \exp[-(\underline{t} - \underline{\mu})^T \Sigma^{-1} (\underline{t} - \underline{\mu})], \underline{t} = \begin{pmatrix} t_1 \\ t_2 \end{pmatrix}$$

$$\phi(t_1) = \int_{-\infty}^{\infty} \phi(t_1, t_2 | \dots) dt_2$$

Joint pdf:  $f_{(x_1, x_2)}(t_1, t_2)$

In a similar way, for multi-normal distribution,

$$\Sigma = \begin{bmatrix} \sigma_1^2 & \sigma_{12} & \sigma_{13} & \cdots & \sigma_{1n} \\ \sigma_{12} & \sigma_2^2 & \cdots & \cdots & \sigma_{2n} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ \sigma_{n1} & \cdots & \cdots & \cdots & \sigma_n^2 \end{bmatrix}$$

### 3.4 Marginal Distribution

$$\underline{x} = (x_1, x_2, \dots, x_n) \rightarrow f_{x_{\sim}}(t_1, \dots, t_n)$$

To calculate pdf,

$$f_{x_{\sim}}(t_1, t_2, \dots, t_n) = \int_{t_{k+1}, \dots, t_n}^{\infty} f_{x_{\sim}}(t_1, \dots, t_n) dt_{k+1}, \dots, dt_n = \sum_{t_{k+1}} \sum_{t_{k+2}} \cdots \sum_{t_n} f_{x_{\sim}}(t_1, \dots, t_n)$$

Joint pdf:  $f(t_1, t_2)$

### 3.5 Conditional Distribution (2 Vars)

$$\underline{x} = (x_1, x_2) \rightarrow \text{joint: } f(t_1, t_2)$$

Conditional distribution means:  $f_{x_1|x_2}(t_1|t_2 = a)$ ,  $a$  is a given constant,  $x_2$  is given, fixed.

$$f_{x_1|x_2}(t_1|t_2) = \frac{f_{x_{\sim}}(t_1, t_2)}{f_{x_2}(t_2)}$$

By definition,  $f_{x_{\sim}}(t_1, t_2)$  is joint distribution,  $f_{x_2}(t_2)$  is marginal distribution. And, we have,

$$f_{x_{\sim}}(t_1, t_2) = f_{x_1|x_2}(t_1|t_2) f_{x_2}(t_2)$$

$$f_{x_2|x_1}(t_2|t_1) = \frac{f_{x_2}(t_1, t_2)}{f_{x_1}(t_1)}$$

$$f_{x_2}(t_1, t_2) = f_{x_2|x_1}(t_2|t_1)f_{x_1}(t_1)$$

So,

$$f_{x_1|x_2}(t_1|t_2) = \frac{f_{x_2|x_1}(t_2|t_1)f_{x_1}(t_1)}{f_{x_2}(t_2)} \rightarrow \text{Bayes Rule}$$

e.g.

$x_2 \backslash x_1$	1	2	3	$f_{x_2}(t)$
0	0.1	0.4	0.2	0.7
1	0.2	0.05	0.05	0.3
$f_{x_1}(t)$	0.3	0.45	0.45	

$$P[x_2 = 0] = P[x_2 = 0|x_1 = 1] + P[x_2 = 0|x_1 = 2] + P[x_2 = 0|x_1 = 3] = 0.7$$

$$P[x_2 = 1] = 1 - P[x_2 = 0] = 0.3$$

## 4 Bayes Classification

$$f_{x_1|x_2}(t_1|t_2) = \frac{f_{x_2|x_1}(t_2|t_1)f_{x_1}(t_1)}{f_{x_2}(t_2)} \rightarrow \text{Bayes Rule}$$

$$\text{For discrete: } f_{x_2}(t_2) = \sum_{t_1} f_{x_2|x_1}(t_2|t_1)f_{x_1}(t_1)$$

$$\text{For continuous: } f_{x_2}(t_2) = \int_{-\infty}^{\infty} f_{x_2|x_1}(t_2|t_1)f_{y_1}(t_1)dt$$

e.g. **Height** Supposed that:  $x_1 \rightarrow \text{Height}$ ,  $x_2 \rightarrow \text{Gender}$ ,  $\begin{cases} 0 & \text{male} \\ 1 & \text{female} \end{cases}$

$$\text{For } x_2 = 0 \rightarrow \text{Height} \sim N(69, 4.5) \Leftrightarrow f_{x_1|x_2}(t_1|t_2 = 0)\phi(t_1|\mu = 69, \sigma = 4.5)$$

$$\text{For } x_2 = 1 \rightarrow \text{Height} \sim N(65, 4.2) \Leftrightarrow f_{x_1|x_2}(t_1|t_2 = 1)\phi(t_1|\mu = 65, \sigma = 4.2)$$

Marginal distribution of height for people:

$$\begin{aligned} f_{x_1}(t_1) &= f_{x_1|x_2}(t_1|t_2 = 0) * f_{x_2}(0) + f_{x_1|x_2}(t_1|t_2 = 1) * f_{x_2}(1) \\ &= \phi(t_1|69, 4.5) \cdot 0.5 + \phi(t_1|65, 4.2) \cdot 0.5 \\ &= \phi(t_1|\frac{69+65}{2}, \sqrt{\frac{4.5^2 + 4.2^2}{2}}) \end{aligned}$$

A person has height 6'7", calculate the probability of each gender.

$$f(x_2 = 0|x_1 = 67) = \frac{f_{x_1|x_2}(67|x_2 = 0)f_{x_2}(0)}{f_{x_1}(67)} = \frac{\phi(67|69, 4.5) \cdot 0.5}{f_{x_1}(67)}$$

$$f(x_2 = 1 | x_1 = 67) = \frac{f_{x_1|x_2}(67|x_2 = 1)f_{x_2}(1)}{f_{x_1}(67)} = \frac{\phi(67|65, 4.2) \cdot 0.5}{f_{x_1}(67)}$$

$$\text{loss}(\hat{f}, \mathbf{x}|f) \rightarrow \text{Minimum } E_{\mathbf{x}} : \text{Loss}(\hat{f}|f)$$

$$\text{Misclassification Rate: } \text{loss}(\hat{f}, \mathbf{x}|f) : \begin{cases} 0, & f_{\mathbf{x}} = \hat{f}(\mathbf{x}) \\ 1, & f_{\mathbf{x}} \neq \hat{f}(\mathbf{x}) \end{cases}, \text{Risk} = E_{\mathbf{x}}[\text{loss}(\hat{f}, \mathbf{x}|f)]$$

$$\text{Probability of Misclassification: } E_{\mathbf{x}} = \begin{cases} \sum_{t_i} t_i f_{\mathbf{x}}(t_i), & \text{for discrete} \\ \int_0^1 t f_{\mathbf{x}}(t) dt, & \text{for continuous} \end{cases}$$

Bayes Classification Rule for Binary: choose  $k$  that  $P[y = k|\mathbf{x}]$ ,

$$k = \underset{k}{\operatorname{argmax}} \frac{P[\mathbf{x}|k]P[k]}{P[\mathbf{x}]} \propto P[\mathbf{x}|k]P[k]$$

For cost matrix  $C$ , assumed that

$C_{ij}$  = the cost of classification if by wrong classified in  $j$

$$\begin{aligned} E_{\mathbf{x}}(\text{loss}(j = \hat{f}(\mathbf{x})|\mathbf{i})) &= \sum_{i=1}^k P[\mathbf{i}|\mathbf{x}]C_{ij} = \sum_{i=1}^k \frac{f(\mathbf{x}|\mathbf{i})P[\mathbf{i}]C_{ij}}{f_{\mathbf{x}}} \cdot C_{ij} \\ &\propto \sum_{i=1}^k f(\mathbf{x}|\mathbf{i})P[\mathbf{i}]C_{ij} \end{aligned}$$

$$\text{Modified Bayes Rule: } \mathbf{c} = \underset{j}{\operatorname{argmin}} \sum_{i=1}^k f(\mathbf{x}|\mathbf{i})P_r[\mathbf{i}]C_{ij}$$

**e.g. Coins** In a box,  $\frac{1}{4}$  of coins are fake,  $\frac{3}{4}$  of coins are real.

$$\text{For fake: } P[\text{head}] = \frac{1}{3}, P[\text{tail}] = \frac{2}{3}$$

$$\text{For real: } P[\text{head}] = \frac{1}{2}, P[\text{tail}] = \frac{1}{2}$$

Take a random coin selected,  $n = 20$  times,  $t = 7$  heads, what is  $P[\text{real}]$  and  $P[\text{false}]$ ?

$x_1$  = number of heads in  $n = 20$  trials.

$$x_2 = \begin{cases} 0, & \text{fake} \\ 1, & \text{real} \end{cases}$$

$$\text{So, } f_{x_2}(0) = \frac{1}{4}, f_{x_2}(1) = \frac{3}{4}$$

$$f_{x_2|x_1}(t_2 = 0 | x_1 = 7, n = 20) = \frac{f_{x_1|x_2}(7|\text{fake}, n = 20)f_{x_2}(0)}{f_{x_1}(7)} = \frac{\binom{20}{7}(\frac{1}{3})^7(\frac{2}{3})^{13} * 0.25}{f_{x_1}(7)} = \frac{0.45}{f_{x_1}(7)}$$



$$f_{x_2|x_1}(t_2 = 1|x_1 = 7, n = 20) = \frac{f_{x_1|x_2}(7|\text{real}, n = 20)f_{x_2}(1)}{f_{x_1}(7)} = \frac{\binom{20}{7}(\frac{1}{2})^7(\frac{1}{2})^{13} * 0.75}{f_{x_1}(7)} = \frac{0.55}{f_{x_1}(7)}$$

$$f_{x_1}(7) = f_{x_1|x_2}(7|\text{fake})f_{x_2}(0) + f_{x_1|x_2}(7|\text{real})f_{x_2}(1)$$

Cost matrix for real and fake:  $C = \begin{pmatrix} C_{11} & C_{12} \\ C_{21} & C_{22} \end{pmatrix}$

For real:  $P[x_2 = 1|x_1 = 7]C_{11} + P[x_2 = 0|x_1 = 7]C_{12}$

For fake:  $P[x_2 = 1|x_1 = 7]C_{21} + P[x_2 = 0|x_1 = 7]C_{22}$