# Probability & Information Theory

Shan-Hung Wu shwu@cs.nthu.edu.tw

Department of Computer Science, National Tsing Hua University, Taiwan

Machine Learning

Shan-Hung Wu (CS, NTHU)

# Outline

- 1 Random Variables & Probability Distributions
- 2 Multivariate & Derived Random Variables
- **3** Bayes' Rule & Statistics
- 4 Application: Principal Components Analysis
- 5 Technical Details of Random Variables
- 6 Common Probability Distributions
- 7 Common Parametrizing Functions
- 8 Information Theory
- 9 Application: Decision Trees & Random Forest

Shan-Hung Wu (CS, NTHU)

# Outline

#### 1 Random Variables & Probability Distributions

- 2 Multivariate & Derived Random Variables
- **3** Bayes' Rule & Statistics
- Application: Principal Components Analysis
- 5 Technical Details of Random Variables
- 6 Common Probability Distributions
- **7** Common Parametrizing Functions
- 8 Information Theory
- 9 Application: Decision Trees & Random Forest

Shan-Hung Wu (CS, NTHU)

### **Random Variables**

- A *random variable* x is a variable that can take on different values randomly
  - E.g.,  $Pr(x = x_1) = 0.1$ ,  $Pr(x = x_2) = 0.3$ , etc.
  - ${\scriptstyle \bullet}\,$  Technically, x is a function that maps events to a real values
- Must be coupled with a *probability distribution* P that specifies how likely each value is
  - $x \sim P(\theta)$  means "x has distribution P parametrized by  $\theta$ "

- If x is discrete, P(x = x) denotes a *probability mass function*  $P_x(x) = Pr(x = x)$ 
  - E.g., the output of a fair dice has discrete uniform distribution with P(x) = 1/6

- If x is discrete, P(x = x) denotes a *probability mass function*  $P_x(x) = Pr(x = x)$ 
  - E.g., the output of a fair dice has discrete uniform distribution with P(x) = 1/6
- If x is continuous, P(x = x) denotes a *probability density function*  $p_x(x) \ge 0$

- If x is discrete, P(x = x) denotes a *probability mass function*  $P_x(x) = Pr(x = x)$ 
  - E.g., the output of a fair dice has discrete uniform distribution with P(x) = 1/6
- If x is continuous, P(x = x) denotes a *probability density function*  $p_x(x) \ge 0$ 
  - Is  $p_x(x)$  a probability?

- If x is discrete, P(x = x) denotes a *probability mass function*  $P_x(x) = Pr(x = x)$ 
  - E.g., the output of a fair dice has discrete uniform distribution with P(x) = 1/6
- If x is continuous, P(x = x) denotes a *probability density function*  $p_x(x) \ge 0$ 
  - Is  $p_x(x)$  a probability? **No**, it is "rate of increase in probability at x"

$$\Pr(a \le \mathbf{x} \le b) = \int_{[a,b]} p(x) dx$$

Shan-Hung Wu (CS, NTHU)

- If x is discrete, P(x = x) denotes a *probability mass function*  $P_x(x) = Pr(x = x)$ 
  - E.g., the output of a fair dice has discrete uniform distribution with P(x) = 1/6
- $\bullet$  If x is continuous,  $\mathbf{P}(\mathbf{x}=x)$  denotes a probability density function  $p_{\mathbf{x}}(x)\geq 0$ 
  - Is  $p_x(x)$  a probability? **No**, it is "rate of increase in probability at x"

$$\Pr(a \le \mathbf{x} \le b) = \int_{[a,b]} p(x) dx$$

- $p_{\rm x}(x)$  can be greater than 1
- E.g., a continuous uniform distribution within [a,b] has p(x) = 1/b-a if  $x \in [a,b]$ ; 0 otherwise

Shan-Hung Wu (CS, NTHU)

# **Marginal Probability**

- $\bullet\,$  Consider a probability distribution over a set of variables, e.g., P(x,y)
- The probability distribution over the subset of random variables called the *marginal probability* distribution:

$$P(x = x) = \sum_{y} P(x, y)$$
 or  $\int p(x, y) dy$ 

Also called the sum rule of probability

Shan-Hung Wu (CS, NTHU)

# **Conditional Probability**

• Conditional density function:

$$P(x = x | y = y) = \frac{P(x = x, y = y)}{P(y = y)}$$

• Defined only when P(y = y) > 0

# **Conditional Probability**

• Conditional density function:

$$P(x = x | y = y) = \frac{P(x = x, y = y)}{P(y = y)}$$

- Defined only when  $P(\boldsymbol{y}=\boldsymbol{y}) > 0$
- Product rule of probability:

$$\mathbf{P}(\mathbf{x}^{(1)},\cdots,\mathbf{x}^{(n)}) = \mathbf{P}(\mathbf{x}^{(1)})\Pi_{i=2}^{n}\mathbf{P}(\mathbf{x}^{(i)} | \mathbf{x}^{(1)},\cdots,\mathbf{x}^{(i-1)})$$

• E.g., 
$$P(a,b,c) = P(a | b,c)P(b | c)P(c)$$

Shan-Hung Wu (CS, NTHU)

# Independence and Conditional Independence

• We say random variables x is *independent* with y iff

 $P(x \,|\, y) = P(x)$ 

- Implies P(x,y) = P(x)P(y)
- ${\scriptstyle \bullet }$  Denoted by  $x \perp y$

# Independence and Conditional Independence

• We say random variables x is *independent* with y iff

$$P(x \,|\, y) = P(x)$$

- Implies P(x,y) = P(x)P(y)
- ${\scriptstyle \bullet }$  Denoted by  $x \perp y$

We say random variables x is *conditionally independent* with y given z iff

$$P(x \,|\, y, z) = P(x \,|\, z)$$

- Implies P(x, y | z) = P(x | z)P(y | z)
- ${\ {\bullet} \ }$  Denoted by  $x \perp y \, | \, z$

#### Expectation

• The *expectation* (or *expected value* or *mean*) of some function f with respect to x is the "average" value that f takes on:<sup>1</sup>

$$\mathbf{E}_{\mathbf{x}\sim\mathbf{P}}[\mathbf{f}(\mathbf{x})] = \sum_{x} P_{\mathbf{x}}(x)f(x) \text{ or } \int p_{\mathbf{x}}(x)f(x)dx = \mu_{\mathbf{f}(\mathbf{x})}$$

<sup>&</sup>lt;sup>1</sup>The bracket  $[\cdot]$  here is used to distinguish the parentheses inside and has nothing to do with functionals.

Prob. & Info. Theory

#### Expectation

 The *expectation* (or *expected value* or *mean*) of some function f with respect to x is the "average" value that f takes on:<sup>1</sup>

$$\mathbf{E}_{\mathbf{x}\sim\mathbf{P}}[\mathbf{f}(\mathbf{x})] = \sum_{x} P_{\mathbf{x}}(x) f(x) \text{ or } \int p_{\mathbf{x}}(x) f(x) dx = \mu_{\mathbf{f}(\mathbf{x})}$$

• Expectation is linear: E[af(x) + b] = aE[f(x)] + b for deterministic a and b

<sup>&</sup>lt;sup>1</sup>The bracket  $[\cdot]$  here is used to distinguish the parentheses inside and has nothing to do with functionals.

Prob. & Info. Theory

#### Expectation

 The *expectation* (or *expected value* or *mean*) of some function f with respect to x is the "average" value that f takes on:<sup>1</sup>

$$\mathbf{E}_{\mathbf{x}\sim\mathbf{P}}[\mathbf{f}(\mathbf{x})] = \sum_{x} P_{\mathbf{x}}(x) f(x) \text{ or } \int p_{\mathbf{x}}(x) f(x) dx = \mu_{\mathbf{f}(\mathbf{x})}$$

- Expectation is linear: E[af(x) + b] = aE[f(x)] + b for deterministic a and b
- E[E[f(x)]] = E[f(x)], as E[f(x)] is deterministic

<sup>&</sup>lt;sup>1</sup>The bracket  $[\cdot]$  here is used to distinguish the parentheses inside and has nothing to do with functionals.

### **Expectation over Multiple Variables**

• Defined over the join probability distribution, e.g.,

$$\mathbf{E}[\mathbf{f}(\mathbf{x},\mathbf{y})] = \sum_{x,y} P_{\mathbf{x},\mathbf{y}}(x,y) f(x,y) \text{ or } \int_{x,y} p_{\mathbf{x},\mathbf{y}}(x,y) f(x,y) dxdy$$

### **Expectation over Multiple Variables**

Defined over the join probability distribution, e.g.,

$$\mathbf{E}[\mathbf{f}(\mathbf{x},\mathbf{y})] = \sum_{x,y} P_{\mathbf{x},y}(x,y) f(x,y) \text{ or } \int_{x,y} p_{\mathbf{x},y}(x,y) f(x,y) dxdy$$

- E[f(x) | y = y] =  $\int p_{x|y}(x|y)f(x)dx$  is called the *conditional* expectation
- $E[f(x)g(y)] = E[f(x)]E[g(y)] \mbox{ if } x \mbox{ and } y \mbox{ are independent [Proof]}$

Shan-Hung Wu (CS, NTHU)

### Variance

• The *variance* measures how much the values of *f* deviate from its expected value when seeing different values of x:

$$\operatorname{Var}[f(x)] = \operatorname{E}\left[(f(x) - \operatorname{E}[f(x)])^2\right] = \sigma_{f(x)}^2$$

•  $\sigma_{\!f(x)}$  is called the *standard deviation* 

### Variance

• The *variance* measures how much the values of *f* deviate from its expected value when seeing different values of x:

$$\operatorname{Var}[f(x)] = \operatorname{E}\left[(f(x) - \operatorname{E}[f(x)])^2\right] = \sigma_{f(x)}^2$$

•  $\sigma_{\mathrm{f(x)}}$  is called the *standard deviation* 

- $Var[f(x)] = E[f(x)^2] E[f(x)]^2$  [Proof]
- $Var[af(x) + b] = a^2 Var[f(x)]$  for deterministic a and b [Proof]

# Covariance I

• *Covariance* gives some sense of how much two values are *linearly* related to each other

 $Cov[f(x),g(y)] = E\left[(f(x) - E[f(x)])(g(y) - E[g(y)])\right]$ 

- If sign positive, both variables tend to take on high values simultaneously
- If sign negative, one variable tend to take on high value while the other taking on low one

# Covariance I

• *Covariance* gives some sense of how much two values are *linearly* related to each other

 $Cov[f(x),g(y)] = E\left[(f(x) - E[f(x)])(g(y) - E[g(y)])\right]$ 

- If sign positive, both variables tend to take on high values simultaneously
- If sign negative, one variable tend to take on high value while the other taking on low one
- If x and y are independent, then Cov(x, y) = 0 [Proof]
  - The converse is *not* true as X and Y may be related in a nonlinear way
  - E.g., y = sin(x) and  $x \sim Uniform(-2\pi, 2\pi)$

# Covariance II

#### • $\operatorname{Var}(a\mathbf{x} + b\mathbf{y}) = a^2 \operatorname{Var}(\mathbf{x}) + b^2 \operatorname{Var}(\mathbf{y}) + 2ab \operatorname{Cov}(\mathbf{x}, \mathbf{y})$ [Proof]

# Covariance II

- Var(ax + by) = a<sup>2</sup>Var(x) + b<sup>2</sup>Var(y) + 2abCov(x,y) [Proof]
   Var(x + y) = Var(x) + Var(y) if x and y are independent
- Cov(ax+b, cy+d) = acCov(x, y) [Proof]

# Covariance II

Var(ax+by) = a<sup>2</sup>Var(x) + b<sup>2</sup>Var(y) + 2abCov(x,y) [Proof]
 Var(x+y) = Var(x) + Var(y) if x and y are independent

• 
$$Cov(ax+b, cy+d) = acCov(x, y)$$
 [Proof]

• Cov(ax + by, cw + dv) = acCov(x, w) + adCov(x, v) + bcCov(y, w) + bdCov(y, v) [Proof]

# Outline

1 Random Variables & Probability Distributions

#### 2 Multivariate & Derived Random Variables

- **3** Bayes' Rule & Statistics
- 4 Application: Principal Components Analysis
- 5 Technical Details of Random Variables
- 6 Common Probability Distributions
- 7 Common Parametrizing Functions
- Information Theory
- 9 Application: Decision Trees & Random Forest

Shan-Hung Wu (CS, NTHU)

- A multivariate random variable is denoted by  $\mathbf{x} = [\mathbf{x}_1, \cdots, \mathbf{x}_d]^{ op}$ 
  - Normally,  $x_i$ 's (*attributes* or *variables* or *features*) are dependent with each other
  - $P(\boldsymbol{x})$  is a joint distribution of  $x_1,\cdots,x_d$

- A multivariate random variable is denoted by  $\mathbf{x} = [\mathbf{x}_1, \cdots, \mathbf{x}_d]^{ op}$ 
  - Normally,  $x_i$ 's (*attributes* or *variables* or *features*) are dependent with each other
  - $P(\mathbf{x})$  is a joint distribution of  $x_1, \cdots, x_d$
- The *mean* of **x** is defined as  $\mu_{\mathbf{x}} = \mathrm{E}(\mathbf{x}) = [\mu_{x_1}, \cdots, \mu_{x_d}]^\top$

- A multivariate random variable is denoted by  $\mathbf{x} = [\mathbf{x}_1, \cdots, \mathbf{x}_d]^{ op}$ 
  - Normally, x<sub>i</sub>'s (*attributes* or *variables* or *features*) are dependent with each other
  - $P(\mathbf{x})$  is a joint distribution of  $x_1, \cdots, x_d$
- The *mean* of **x** is defined as  $\mu_{\mathbf{x}} = \mathrm{E}(\mathbf{x}) = [\mu_{x_1}, \cdots, \mu_{x_d}]^\top$
- The *covariance matrix* of **x** is defined as:

$$\Sigma_{\mathbf{x}} = \begin{bmatrix} \sigma_{\mathbf{x}_1}^2 & \sigma_{\mathbf{x}_1,\mathbf{x}_2} & \cdots & \sigma_{\mathbf{x}_1,\mathbf{x}_d} \\ \sigma_{\mathbf{x}_2,\mathbf{x}_1} & \sigma_{\mathbf{x}_2}^2 & \cdots & \sigma_{\mathbf{x}_2,\mathbf{x}_d} \\ \vdots & \vdots & \ddots & \vdots \\ \sigma_{\mathbf{x}_d,\mathbf{x}_1} & \sigma_{\mathbf{x}_d,\mathbf{x}_2} & \cdots & \sigma_{\mathbf{x}_d}^2 \end{bmatrix}$$

• 
$$\sigma_{\mathbf{x}_i,\mathbf{x}_j} = \operatorname{Cov}(\mathbf{x}_i,\mathbf{x}_j) = \operatorname{E}[(\mathbf{x}_i - \mu_{\mathbf{x}_i})(\mathbf{x}_j - \mu_{\mathbf{x}_j})] = \operatorname{E}(\mathbf{x}_i\mathbf{x}_j) - \mu_{\mathbf{x}_i}\mu_{\mathbf{x}_j}$$

Shan-Hung Wu (CS, NTHU)

- A multivariate random variable is denoted by  $\mathbf{x} = [\mathbf{x}_1, \cdots, \mathbf{x}_d]^{ op}$ 
  - Normally,  $x_i$ 's (*attributes* or *variables* or *features*) are dependent with each other
  - $P(\mathbf{x})$  is a joint distribution of  $x_1, \cdots, x_d$
- The *mean* of **x** is defined as  $\mu_{\mathbf{x}} = \mathrm{E}(\mathbf{x}) = [\mu_{x_1}, \cdots, \mu_{x_d}]^\top$
- The *covariance matrix* of **x** is defined as:

$$\Sigma_{\mathbf{x}} = \begin{bmatrix} \sigma_{x_{1}}^{2} & \sigma_{x_{1},x_{2}} & \cdots & \sigma_{x_{1},x_{d}} \\ \sigma_{x_{2},x_{1}} & \sigma_{x_{2}}^{2} & \cdots & \sigma_{x_{2},x_{d}} \\ \vdots & \vdots & \ddots & \vdots \\ \sigma_{x_{d},x_{1}} & \sigma_{x_{d},x_{2}} & \cdots & \sigma_{x_{d}}^{2} \end{bmatrix}$$

• 
$$\sigma_{\mathbf{x}_i,\mathbf{x}_j} = \operatorname{Cov}(\mathbf{x}_i,\mathbf{x}_j) = \operatorname{E}[(\mathbf{x}_i - \boldsymbol{\mu}_{\mathbf{x}_i})(\mathbf{x}_j - \boldsymbol{\mu}_{\mathbf{x}_j})] = \operatorname{E}(\mathbf{x}_i\mathbf{x}_j) - \boldsymbol{\mu}_{\mathbf{x}_i}\boldsymbol{\mu}_{\mathbf{x}_j}$$
  
•  $\Sigma_{\mathbf{x}} = \operatorname{Cov}(\mathbf{x}) = \operatorname{E}\left[(\mathbf{x} - \boldsymbol{\mu}_{\mathbf{x}})(\mathbf{x} - \boldsymbol{\mu}_{\mathbf{x}})^{\top}\right] = \operatorname{E}(\mathbf{x}\mathbf{x}^{\top}) - \boldsymbol{\mu}_{\mathbf{x}}\boldsymbol{\mu}_{\mathbf{x}}^{\top}$ 

Shan-Hung Wu (CS, NTHU)

 ${\ensuremath{\, \bullet \, }}\xspace \Sigma_x$  is always symmetric

- $\Sigma_x$  is always symmetric
- $\boldsymbol{\Sigma}_{\boldsymbol{X}}$  is always positive semidefinite [Homework]

- $\Sigma_x$  is always symmetric
- $\boldsymbol{\Sigma}_{\boldsymbol{X}}$  is always positive semidefinite [Homework]
- $\boldsymbol{\Sigma}_{\boldsymbol{X}}$  is nonsingular iff it is positive definite

- $\Sigma_{\mathbf{x}}$  is always symmetric
- $\Sigma_x$  is always positive semidefinite [Homework]
- $\Sigma_x$  is nonsingular iff it is positive definite
- $\Sigma_x$  is singular implies that x has either:
  - Deterministic/independent/non-linearly dependent attributes causing zero rows, or
  - Redundant attributes causing linear dependency between rows

### **Derived Random Variables**
# Outline

- 1 Random Variables & Probability Distributions
- 2 Multivariate & Derived Random Variables
- **3** Bayes' Rule & Statistics
- 4 Application: Principal Components Analysis
- 5 Technical Details of Random Variables
- 6 Common Probability Distributions
- 7 Common Parametrizing Functions
- 8 Information Theory
- 9 Application: Decision Trees & Random Forest

Shan-Hung Wu (CS, NTHU)

What Does Pr(x = x) Mean?

Shan-Hung Wu (CS, NTHU)

Prob. & Info. Theory

Machine Learning 19 / 78

What Does Pr(x = x) Mean?

Bayesian probability: it's a degree of belief or qualitative levels of certainty

What Does Pr(x = x) Mean?

- **Bayesian probability**: it's a degree of belief or qualitative levels of certainty
- Prequentist probability: if we can draw samples of x, then the proportion of frequency of samples having the value x is equal to Pr(x = x)

### Bayes' Rule

$$P(y | x) = \frac{P(x | y)P(y)}{P(x)} = \frac{P(x | y)P(y)}{\sum_{y} P(x | y = y)P(y = y)}$$

• Bayes' Rule is so important in statistics (and ML as well) such that each term has a name:

$$posterior of y = \frac{(likelihood of y) \times (prior of y)}{evidence}$$

### Bayes' Rule

$$P(y | x) = \frac{P(x | y)P(y)}{P(x)} = \frac{P(x | y)P(y)}{\sum_{y} P(x | y = y)P(y = y)}$$

 Bayes' Rule is so important in statistics (and ML as well) such that each term has a name:

$$\textit{posterior of } y = \frac{(\textit{likelihood of } y) \times (\textit{prior of } y)}{\textit{evidence}}$$

• Why is it so important?

### Bayes' Rule

$$P(y | x) = \frac{P(x | y)P(y)}{P(x)} = \frac{P(x | y)P(y)}{\sum_{y} P(x | y = y)P(y = y)}$$

 Bayes' Rule is so important in statistics (and ML as well) such that each term has a name:

$$\textit{posterior of } y = \frac{(\textit{likelihood of } y) \times (\textit{prior of } y)}{\textit{evidence}}$$

- Why is it so important?
- E.g., a doctor diagnoses you as having a disease by letting x be "symptom" and y be "disease"
  - $\circ~P(x\,|\,y)$  and P(y) may be estimated from sample frequencies more easily

Shan-Hung Wu (CS, NTHU)

### **Point Estimation**

• **Point estimation** is the attempt to estimate some fixed but unknown quantity  $\theta$  of a random variable by using sample data

### **Point Estimation**

- **Point estimation** is the attempt to estimate some fixed but unknown quantity  $\theta$  of a random variable by using sample data
- Let {x<sup>(1)</sup>,...,x<sup>(n)</sup>} be a set of n independent and identically distributed (*i.i.d.*) samples of a random variable x, a *point estimator* or *statistic* is a function of the data:

$$\hat{\theta}_n = g(x^{(1)}, \cdots, x^{(n)})$$

•  $\hat{\theta}_n$  is called the **estimate** of  $\theta$ 

Shan-Hung Wu (CS, NTHU)

• Given  $X = [\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(n)}]^\top \in \mathbb{R}^{n \times d}$  the i.i.d samples, what are the estimates of the mean and covariance of  $\mathbf{x}$ ?

- Given  $X = [\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(n)}]^\top \in \mathbb{R}^{n \times d}$  the i.i.d samples, what are the estimates of the mean and covariance of  $\mathbf{x}$ ?
- A sample mean:

$$\hat{\boldsymbol{\mu}}_{\mathbf{x}} = \frac{1}{n} \sum_{i=1}^{n} \boldsymbol{x}^{(i)}$$

- Given  $X = [\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(n)}]^\top \in \mathbb{R}^{n \times d}$  the i.i.d samples, what are the estimates of the mean and covariance of  $\mathbf{x}$ ?
- A sample mean:

$$\hat{\boldsymbol{\mu}}_{\mathbf{x}} = \frac{1}{n} \sum_{i=1}^{n} \boldsymbol{x}^{(i)}$$

• A sample covariance matrix:

$$\hat{\Sigma}_{\mathbf{x}} = \frac{1}{n} \sum_{i=1}^{n} (\mathbf{x}^{(i)} - \hat{\boldsymbol{\mu}}_{\mathbf{x}}) (\mathbf{x}^{(i)} - \hat{\boldsymbol{\mu}}_{\mathbf{x}})^{\top}$$

• 
$$\hat{\sigma}_{\mathbf{x}_i,\mathbf{x}_j}^2 = \frac{1}{n} \sum_{s=1}^n (x_i^{(s)} - \hat{\mu}_{\mathbf{x}_i}) (x_j^{(s)} - \hat{\mu}_{\mathbf{x}_j})$$

Shan-Hung Wu (CS, NTHU)

- Given  $X = [\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(n)}]^\top \in \mathbb{R}^{n \times d}$  the i.i.d samples, what are the estimates of the mean and covariance of  $\mathbf{x}$ ?
- A sample mean:

$$\hat{\boldsymbol{\mu}}_{\mathbf{x}} = \frac{1}{n} \sum_{i=1}^{n} \boldsymbol{x}^{(i)}$$

• A sample covariance matrix:

$$\hat{\Sigma}_{\mathbf{x}} = \frac{1}{n} \sum_{i=1}^{n} (\mathbf{x}^{(i)} - \hat{\boldsymbol{\mu}}_{\mathbf{x}}) (\mathbf{x}^{(i)} - \hat{\boldsymbol{\mu}}_{\mathbf{x}})^{\top}$$

• 
$$\hat{\sigma}_{\mathbf{x}_i,\mathbf{x}_j}^2 = \frac{1}{n} \sum_{s=1}^n (x_i^{(s)} - \hat{\mu}_{\mathbf{x}_i}) (x_j^{(s)} - \hat{\mu}_{\mathbf{x}_j})$$
  
• If each  $\mathbf{x}^{(i)}$  is centered (by subtracting  $\hat{\mu}_{\mathbf{x}}$  first), then  $\hat{\Sigma}_{\mathbf{x}} = \frac{1}{n} \mathbf{X}^\top \mathbf{X}$ 

Shan-Hung Wu (CS, NTHU)

# Outline

- 1 Random Variables & Probability Distributions
- 2 Multivariate & Derived Random Variables
- **3** Bayes' Rule & Statistics

#### **4** Application: Principal Components Analysis

- 5 Technical Details of Random Variables
- 6 Common Probability Distributions
- 7 Common Parametrizing Functions
- 8 Information Theory
- 9 Application: Decision Trees & Random Forest

Shan-Hung Wu (CS, NTHU)

- Give a collection of data points  $\mathbb{X} = \{ \pmb{x}^{(i)} \}_{i=1}^N$ , where  $\pmb{x}^{(i)} \in \mathbb{R}^D$
- Suppose we want to lossily compress X, i.e., to find a function f such that  $f(\mathbf{x}^{(i)}) = \mathbf{z}^{(i)} \in \mathbb{R}^{K}$ , where K < D
- How to keep the maximum info in X?

- Let  $\boldsymbol{x}^{(i)}$ 's be i.i.d. samples of a random variable  $\mathbf{x}$
- Let f be linear, i.e.,  $f(\mathbf{x}) = \mathbf{W}^{\top}\mathbf{x}$  for some  $\mathbf{W} \in \mathbb{R}^{D \times K}$

- Let  $x^{(i)}$ 's be i.i.d. samples of a random variable x
- Let f be linear, i.e.,  $f(\mathbf{x}) = \mathbf{W}^{\top}\mathbf{x}$  for some  $\mathbf{W} \in \mathbb{R}^{D \times K}$
- **Principal Component Analysis (PCA)** finds K orthonormal vectors  $W = [w^{(1)}, \dots, w^{(K)}]$  such that the transformed variable  $\mathbf{z} = W^{\top}\mathbf{x}$  has the most "spread out" attributes, i.e., each attribute  $z_j = w^{(j)\top}\mathbf{x}$  has the maximum variance  $\operatorname{Var}(z_j)$ 
  - $w^{(1)}, \cdots, w^{(K)}$  are called the *principle components*

- Let  $\mathbf{x}^{(i)}$ 's be i.i.d. samples of a random variable  $\mathbf{x}$
- Let f be linear, i.e.,  $f(\mathbf{x}) = \mathbf{W}^{\top}\mathbf{x}$  for some  $\mathbf{W} \in \mathbb{R}^{D \times K}$
- **Principal Component Analysis (PCA)** finds K orthonormal vectors  $W = [w^{(1)}, \dots, w^{(K)}]$  such that the transformed variable  $\mathbf{z} = W^{\top}\mathbf{x}$  has the most "spread out" attributes, i.e., each attribute  $z_j = w^{(j)\top}\mathbf{x}$  has the maximum variance  $\operatorname{Var}(z_j)$
- \$\mathbf{w}^{(1)}\$,\$\dots\$,\$\mathbf{w}^{(K)}\$ are called the *principle components*Why \$\mathbf{w}^{(1)}\$,\$\dots\$,\$\mathbf{w}^{(K)}\$ need to be orthogonal with each other?

- Let  $\mathbf{x}^{(i)}$ 's be i.i.d. samples of a random variable  $\mathbf{x}$
- Let f be linear, i.e.,  $f(\mathbf{x}) = \mathbf{W}^{\top}\mathbf{x}$  for some  $\mathbf{W} \in \mathbb{R}^{D \times K}$
- **Principal Component Analysis (PCA)** finds K orthonormal vectors  $W = [w^{(1)}, \dots, w^{(K)}]$  such that the transformed variable  $\mathbf{z} = W^{\top}\mathbf{x}$  has the most "spread out" attributes, i.e., each attribute  $z_j = w^{(j)\top}\mathbf{x}$  has the maximum variance  $\operatorname{Var}(z_j)$ 
  - $w^{(1)}, \cdots, w^{(K)}$  are called the *principle components*
- Why  $w^{(1)}, \cdots, w^{(K)}$  need to be orthogonal with each other?
  - Each  $w^{(j)}$  keeps information that cannot be explained by others, so together they preserve the most info

- Let  $x^{(i)}$ 's be i.i.d. samples of a random variable x
- Let f be linear, i.e.,  $f(\mathbf{x}) = \mathbf{W}^{\top}\mathbf{x}$  for some  $\mathbf{W} \in \mathbb{R}^{D \times K}$
- **Principal Component Analysis (PCA)** finds K orthonormal vectors  $W = [w^{(1)}, \dots, w^{(K)}]$  such that the transformed variable  $\mathbf{z} = W^{\top}\mathbf{x}$  has the most "spread out" attributes, i.e., each attribute  $z_j = w^{(j)\top}\mathbf{x}$  has the maximum variance  $\operatorname{Var}(z_j)$ 
  - $w^{(1)}, \dots, w^{(K)}$  are called the *principle components*
- Why  $w^{(1)}, \cdots, w^{(K)}$  need to be orthogonal with each other?
  - Each  $w^{(j)}$  keeps information that cannot be explained by others, so together they preserve the most info
- Why  $\|\boldsymbol{w}^{(j)}\| = 1$  for all j?

- Let  $\mathbf{x}^{(i)}$ 's be i.i.d. samples of a random variable  $\mathbf{x}$
- Let f be linear, i.e.,  $f(\mathbf{x}) = \mathbf{W}^{\top}\mathbf{x}$  for some  $\mathbf{W} \in \mathbb{R}^{D \times K}$
- **Principal Component Analysis (PCA)** finds K orthonormal vectors  $W = [w^{(1)}, \dots, w^{(K)}]$  such that the transformed variable  $\mathbf{z} = W^{\top}\mathbf{x}$  has the most "spread out" attributes, i.e., each attribute  $z_j = w^{(j)\top}\mathbf{x}$  has the maximum variance  $\operatorname{Var}(z_j)$ 
  - $w^{(1)}, \cdots, w^{(K)}$  are called the *principle components*
- Why  $w^{(1)}, \cdots, w^{(K)}$  need to be orthogonal with each other?
  - Each  $w^{(j)}$  keeps information that cannot be explained by others, so together they preserve the most info
- Why  $||\mathbf{w}^{(j)}|| = 1$  for all j?
  - $\bullet\,$  Only directions matter—we don't want to maximize  $Var(z_j)$  by finding a long  $\pmb{w}^{(j)}$

Shan-Hung Wu (CS, NTHU)

• For simplicity, let's consider K = 1 first

• How to evaluate Var(z<sub>1</sub>)?

- For simplicity, let's consider K = 1 first
- How to evaluate Var(z<sub>1</sub>)?
  - Recall that  $z_1 = w^{(1)\top} x$  implies  $\sigma_{z_1}^2 = w^{(1)\top} \Sigma_x w^{(1)}$  [Homework]
  - How to get  $\Sigma_x$ ?

- For simplicity, let's consider K = 1 first
- How to evaluate  $Var(z_1)$ ?
  - Recall that  $z_1 = w^{(1)\top} x$  implies  $\sigma_{z_1}^2 = w^{(1)\top} \Sigma_x w^{(1)}$  [Homework]
  - How to get  $\Sigma_{\mathbf{x}}$ ?
  - An estimate:  $\hat{\Sigma}_{\mathbf{x}} = \frac{1}{N} \mathbf{X}^{\top} \mathbf{X}$  (assuming  $\mathbf{x}^{(i)}$ 's are centered first)

- For simplicity, let's consider K = 1 first
- How to evaluate Var(z<sub>1</sub>)?
  - Recall that  $z_1 = w^{(1)\top} x$  implies  $\sigma_{z_1}^2 = w^{(1)\top} \Sigma_x w^{(1)}$  [Homework]
  - How to get  $\Sigma_x$ ?
  - An estimate:  $\hat{\Sigma}_{\mathbf{x}} = \frac{1}{N} \mathbf{X}^{\top} \mathbf{X}$  (assuming  $\mathbf{x}^{(i)}$ 's are centered first)
- Optimization problem to solve:

$$\arg\max_{\pmb{w}^{(1)}\in\mathbb{R}^D}\pmb{w}^{(1)\top}\pmb{X}^{\top}\pmb{X}\pmb{w}^{(1)}, \text{ subject to } \|\pmb{w}^{(1)}\|=1$$

• For simplicity, let's consider K = 1 first

- How to evaluate Var(z<sub>1</sub>)?
  - Recall that  $z_1 = w^{(1)\top} x$  implies  $\sigma_{z_1}^2 = w^{(1)\top} \Sigma_x w^{(1)}$  [Homework]
  - How to get  $\Sigma_x$ ?
  - An estimate:  $\hat{\Sigma}_{\mathbf{x}} = \frac{1}{N} \mathbf{X}^{\top} \mathbf{X}$  (assuming  $\mathbf{x}^{(i)}$ 's are centered first)
- Optimization problem to solve:

$$\arg\max_{\pmb{w}^{(1)}\in\mathbb{R}^D}\pmb{w}^{(1)\top}\pmb{X}^{\top}\pmb{X}\pmb{w}^{(1)}, \text{ subject to } \|\pmb{w}^{(1)}\|=1$$

•  $X^{ op}X$  is symmetric thus can be eigendecomposed

• For simplicity, let's consider K = 1 first

- How to evaluate Var(z<sub>1</sub>)?
  - Recall that  $z_1 = w^{(1)\top} \mathbf{x}$  implies  $\sigma_{z_1}^2 = w^{(1)\top} \Sigma_{\mathbf{x}} w^{(1)}$  [Homework]
  - How to get  $\Sigma_x$ ?
  - An estimate:  $\hat{\Sigma}_{\mathbf{x}} = \frac{1}{N} \mathbf{X}^{\top} \mathbf{X}$  (assuming  $\mathbf{x}^{(i)}$ 's are centered first)
- Optimization problem to solve:

$$\arg \max_{\boldsymbol{w}^{(1)} \in \mathbb{R}^D} \boldsymbol{w}^{(1)\top} \boldsymbol{X}^\top \boldsymbol{X} \boldsymbol{w}^{(1)}, \text{ subject to } \|\boldsymbol{w}^{(1)}\| = 1$$

- $X^{ op}X$  is symmetric thus can be eigendecomposed
- By Rayleigh's Quotient, the optimal  $w^{(1)}$  is given by the eigenvector of  $X^{\top}X$  corresponding to the largest eigenvalue

Shan-Hung Wu (CS, NTHU)

• Optimization problem for  $w^{(2)}$ :

```
\arg\max_{\boldsymbol{w}^{(2)}\in\mathbb{R}^D}\boldsymbol{w}^{(2)\top}\boldsymbol{X}^{\top}\boldsymbol{X}\boldsymbol{w}^{(2)}, \text{ subject to } \|\boldsymbol{w}^{(2)}\|=1 \text{ and } \boldsymbol{w}^{(2)\top}\boldsymbol{w}^{(1)}=0
```

• Optimization problem for  $w^{(2)}$ :

 $\arg\max_{\pmb{w}^{(2)}\in\mathbb{R}^D}\pmb{w}^{(2)\top}\pmb{X}^{\top}\pmb{X}\pmb{w}^{(2)}, \text{ subject to } \|\pmb{w}^{(2)}\|=1 \text{ and } \pmb{w}^{(2)\top}\pmb{w}^{(1)}=0$ 

• By Rayleigh's Quotient again,  $w^{(2)}$  is the eigenvector corresponding to the 2-nd largest eigenvalue

• Optimization problem for  $w^{(2)}$ :

 $\arg\max_{\pmb{w}^{(2)}\in\mathbb{R}^D}\pmb{w}^{(2)\top}\pmb{X}^{\top}\pmb{X}\pmb{w}^{(2)}, \text{ subject to } \|\pmb{w}^{(2)}\|=1 \text{ and } \pmb{w}^{(2)\top}\pmb{w}^{(1)}=0$ 

- By Rayleigh's Quotient again,  $w^{(2)}$  is the eigenvector corresponding to the 2-nd largest eigenvalue
- For general case where K > 1, the  $w^{(1)}, \dots, w^{(K)}$  are eigenvectors of  $X^{\top}X$  corresponding to the largest K eigenvalues
  - Proof by induction [Proof]

## Visualization



**Figure:** PCA learns a linear projection that aligns the direction of greatest variance with the axes of the new space. With these new axes, the estimated covariance matrix  $\hat{\Sigma}_{\mathbf{z}} = \mathbf{W}^{\top} \hat{\Sigma}_{\mathbf{x}} \mathbf{W} \in \mathbb{R}^{K \times K}$  is always diagonal.

Shan-Hung Wu (CS, NTHU)

# Outline

- 1 Random Variables & Probability Distributions
- 2 Multivariate & Derived Random Variables
- **3** Bayes' Rule & Statistics
- 4 Application: Principal Components Analysis

#### 5 Technical Details of Random Variables

- 6 Common Probability Distributions
- 7 Common Parametrizing Functions
- Information Theory
- 9 Application: Decision Trees & Random Forest

Shan-Hung Wu (CS, NTHU)

#### Sure and Almost Sure Events

- Given a continuous random variable x, we have Pr(x = x) = 0 for any value x
- Will the event x = x occur?

#### Sure and Almost Sure Events

- Given a continuous random variable x, we have Pr(x = x) = 0 for any value x
- Will the event x = x occur? Yes!
- An event A happens *surely* if always occurs
- An event A happens *almost surely* if Pr(A) = 1 (e.g.,  $Pr(x \neq x) = 1$ )

# Equality of Random Variables I

#### Definition (Equality in Distribution)

Two random variables x and y are *equal in distribution* iff  $Pr(x \le a) = Pr(y \le a)$  for all *a*.

#### Definition (Almost Sure Equality)

Two random variables x and y are *equal almost surely* iff Pr(x = y) = 1.

#### Definition (Equality)

Two random variables x and y are *equal* iff they maps the same events to same values.

# Equality of Random Variables II

• What's the difference between the "equality in distribution" and "almost sure equality?"
## Equality of Random Variables II

- What's the difference between the "equality in distribution" and "almost sure equality?"
- Almost sure equality implies equality in distribution, but converse not true

## Equality of Random Variables II

- What's the difference between the "equality in distribution" and "almost sure equality?"
- Almost sure equality implies equality in distribution, but converse not true
- E.g., let x and y be binary random variables and  $P_x(0) = P_x(1) = P_y(0) = P_y(1) = 0.5$ 
  - They are equal in distribution
  - But  $Pr(x = y) = 0.5 \neq 1$

## Convergence of Random Variables I

#### Definition (Convergence in Distribution)

A sequence of random variables  $\{x^{(1)}, x^{(2)}, \cdots\}$  converges in distribution to x iff  $\lim_{n\to\infty} P(x^{(n)} = x) = P(x = x)$ 

#### Definition (Convergence in Probability)

A sequence of random variables  $\{x^{(1)}, x^{(2)}, \dots\}$  converges in probability to x iff for any  $\varepsilon > 0$ ,  $\lim_{n \to \infty} \Pr(|x^{(n)} - x| < \varepsilon) = 1$ .

#### Definition (Almost Sure Convergence)

A sequence of random variables  $\{x^{(1)}, x^{(2)}, \dots\}$  converges almost surely to x iff  $Pr(\lim_{n\to\infty} x^{(n)} = x) = 1$ .

Shan-Hung Wu (CS, NTHU)

### Convergence of Random Variables II

• What's the difference between the convergence "in probability" and "almost surely?"

### Convergence of Random Variables II

- What's the difference between the convergence "in probability" and "almost surely?"
- Almost sure convergence implies convergence in probability, but converse not true

### Convergence of Random Variables II

- What's the difference between the convergence "in probability" and "almost surely?"
- Almost sure convergence implies convergence in probability, but converse not true
- $\lim_{n\to\infty} \Pr\left(|\mathbf{x}^{(n)} \mathbf{x}| < \varepsilon\right) = 1$  leaves open the possibility that  $|\mathbf{x}^{(n)} \mathbf{x}| > \varepsilon$  happens an infinite number of times
- $Pr\left(\lim_{n\to\infty}x^{(n)}=x\right)=1$  guarantees that  $|x^{(n)}-x|>\epsilon$  almost surely will not occur

## Distribution of Derived Variables I

- Consider a continuous scalar random variable x
- Suppose y = f(x) and  $f^{-1}$  exists, does  $P(y = y) = P(x = f^{-1}(y))$  always hold?

## Distribution of Derived Variables I

- Consider a continuous scalar random variable x
- Suppose y = f(x) and  $f^{-1}$  exists, does  $P(y = y) = P(x = f^{-1}(y))$  always hold? *No*, when x and y are continuous
- 1D example: let  $x \sim \text{Unifrom}(0,1)$  and  $y = 2x \sim \text{Unifrom}(0,2)$ • If  $p_y(y) = p_x(\frac{1}{2}y)$ , then

$$\int_{y=0}^{2} p_{y}(y) dy = \int_{y=0}^{2} p_{x}(\frac{1}{2}y) dy = 2 \neq 1$$

• Violates the unit measure axiom of probability



Shan-Hung Wu (CS, NTHU)

#### **Distribution of Derived Variables II**

• Recall that 
$$Pr(y = y) = p_y(y)dy$$
 and  $Pr(x = x) = p_x(x)dx$ 

• To preserve unit measure, we need to ensure that

$$p_{\mathbf{y}}(f(x))d\mathbf{y} = p_{\mathbf{x}}(x)dx, \forall x$$



We have

$$p_{\mathbf{x}}(x) = p_{\mathbf{y}}(f(x)) \left| \frac{\partial f(x)}{\partial x} \right| \text{ (or } p_{\mathbf{y}}(y) = p_{\mathbf{x}}(f^{-1}(y)) \left| \frac{\partial f^{-1}(y)}{\partial y} \right| \text{)}$$

• Absolute values deal with the case where y = -cx• In 1D example:  $p_y(y) = \frac{1}{2} p_x(\frac{1}{2}y)$  for all y

Shan-Hung Wu (CS, NTHU)

## Distribution of Derived Variables III

• In multivariate case where  $\mathbf{y}=\!f(\mathbf{x})$ , we need to ensure that

$$p_{\mathbf{y}}(\boldsymbol{f}(\boldsymbol{x}))|\det(\Delta)| = p_{\mathbf{x}}(\boldsymbol{x})dx_1dx_2\cdots, \forall \boldsymbol{x}$$

- $dx_1 dx_2 \cdots$ : volume of dx
- |det(Δ)|: volume of dy, where Δ captures the linear relationship between dy and dx



• We have 
$$p_{\mathbf{x}}(\mathbf{x}) = p_{\mathbf{y}}(\mathbf{f}(\mathbf{x})) \left| \frac{1}{dx_1 dx_2 \dots} \det(\Delta) \right| = p_{\mathbf{y}}(\mathbf{f}(\mathbf{x})) \left| \det(\mathbf{J}(\mathbf{f})(\mathbf{x})) \right|$$
  
•  $\mathbf{J}(\mathbf{f})(\mathbf{x})$  is the Jacobian matrix of  $\mathbf{f}$  at input  $\mathbf{x}$ 

• Or equivalently,  $p_{\mathbf{y}}(\mathbf{y}) = p_{\mathbf{x}}(f^{-1}(\mathbf{y})) \left| \det \left( J(f^{-1})(\mathbf{y}) \right) \right|$ 

Shan-Hung Wu (CS, NTHU)

### Distribution of Derived Variables IV

$$\frac{1}{dx_1 dx_2} \det \left( \begin{bmatrix} a & c \\ b & d \end{bmatrix} \right) = \frac{1}{dx_1 dx_2} \det \left( \begin{bmatrix} a & b \\ c & d \end{bmatrix} \right)$$
$$= \det \left( \begin{bmatrix} \frac{a}{dx_1} & \frac{b}{dx_1} \\ \frac{c}{dx_2} & \frac{d}{dx_2} \end{bmatrix} \right) = \det \left( \begin{bmatrix} \frac{\partial y_1}{\partial x_1} & \frac{\partial y_2}{\partial x_1} \\ \frac{\partial y_1}{\partial x_2} & \frac{\partial y_2}{\partial x_2} \end{bmatrix} \right)$$
$$= = \det \left( \begin{bmatrix} \frac{\partial y_1}{\partial x_1} & \frac{\partial y_1}{\partial x_2} \\ \frac{\partial y_2}{\partial x_1} & \frac{\partial y_2}{\partial x_2} \end{bmatrix} \right) = \det (J(f)(\mathbf{x}))$$
$$x_2 \begin{vmatrix} (0, dx_2) & y_2 \\ (dx_1, 0) & y_2 \end{vmatrix}$$
$$(c, d) \qquad (a, b)$$

<u>x</u><sub>1</sub>

*y*<sub>1</sub>

Shan-Hung Wu (CS, NTHU)

# Outline

- 1 Random Variables & Probability Distributions
- 2 Multivariate & Derived Random Variables
- **3** Bayes' Rule & Statistics
- 4 Application: Principal Components Analysis
- 5 Technical Details of Random Variables
- 6 Common Probability Distributions
- 7 Common Parametrizing Functions
- 8 Information Theory
- 9 Application: Decision Trees & Random Forest

Shan-Hung Wu (CS, NTHU)

#### **Random Experiments**

- The value of a random variable x can be think of as the outcome of an random experiment
- $\bullet \ \ \, \text{Helps us define } P(x)$

## Bernoulli Distribution (Discrete)

• Let  $x \in \{0,1\}$  be the outcome of tossing a coin, we have:

Bernoulli(x = x; 
$$\rho$$
) =   

$$\begin{cases}
\rho, & \text{if } x = 1 \\
1 - \rho, & \text{otherwise}
\end{cases} \text{ or } \rho^{x}(1 - \rho)^{1 - x}$$

Properties: [Proof]

• 
$$E(x) = \rho$$
  
•  $Var(x) = \rho(1 - \rho)$ 

## Categorical Distribution (Discrete)

• Let  $x \in \{1, \dots, k\}$  be the outcome of rolling a *k*-sided dice, we have:

Categorical(x = x; 
$$\rho$$
) =  $\prod_{i=1}^{k} \rho_i^{1(x;x=i)}$ , where  $\mathbf{1}^{\top} \rho = 1$ 

## Categorical Distribution (Discrete)

• Let  $x \in \{1, \dots, k\}$  be the outcome of rolling a *k*-sided dice, we have:

Categorical(x = x; 
$$\rho$$
) =  $\prod_{i=1}^{k} \rho_i^{1(x;x=i)}$ , where  $\mathbf{1}^{\top} \rho = 1$ 

• An extension of the Bernoulli distribution for k states

Shan-Hung Wu (CS, NTHU)

Prob. & Info. Theory

Machine Learning 42 / 78

## Multinomial Distribution (Discrete)

• Let  $\mathbf{x} \in \mathbb{R}^k$  be a random vector where  $\mathbf{x}_i$  the number of the outcome *i* after rolling a *k*-sided dice *n* times:

Multinomial(
$$\mathbf{x} = \mathbf{x}; n, \rho$$
) =  $\frac{n!}{x_1! \cdots x_k!} \prod_{i=1}^k \rho_i^{x_i}$ , where  $\mathbf{1}^\top \rho = 1$  and  $\mathbf{1}^\top \mathbf{x} = n$ 

## Multinomial Distribution (Discrete)

• Let  $\mathbf{x} \in \mathbb{R}^k$  be a random vector where  $\mathbf{x}_i$  the number of the outcome *i* after rolling a *k*-sided dice *n* times:

Multinomial(
$$\mathbf{x} = \mathbf{x}; n, \rho$$
) =  $\frac{n!}{x_1! \cdots x_k!} \prod_{i=1}^k \rho_i^{x_i}$ , where  $\mathbf{1}^\top \rho = 1$  and  $\mathbf{1}^\top \mathbf{x} = n$ 

Properties: [Proof]

• 
$$E(\mathbf{x}) = n\rho$$
  
•  $Var(\mathbf{x}) = n \left( diag(\rho) - \rho \rho^{\top} \right)$   
(i.e.,  $Var(\mathbf{x}_i) = n\rho_i(1 - \rho_i)$  and  $Var(\mathbf{x}_i, \mathbf{x}_j) = -n\rho_i\rho_j$ )

Shan-Hung Wu (CS, NTHU)

#### Theorem (Central Limit Theorem)

The sum x of many independent random variables is approximately normally/Gaussian distributed:

$$\mathcal{N}(\mathbf{x}=x;\boldsymbol{\mu},\boldsymbol{\sigma}^2) = \sqrt{\frac{1}{2\pi\sigma^2}} \exp\left(-\frac{1}{2\sigma^2}(x-\boldsymbol{\mu})^2\right).$$

#### Theorem (Central Limit Theorem)

The sum x of many independent random variables is approximately normally/Gaussian distributed:

$$\mathcal{N}(\mathbf{x}=x;\boldsymbol{\mu},\boldsymbol{\sigma}^2) = \sqrt{\frac{1}{2\pi\sigma^2}} \exp\left(-\frac{1}{2\sigma^2}(x-\boldsymbol{\mu})^2\right).$$

Holds regardless of the original distributions of individual variables

#### Theorem (Central Limit Theorem)

The sum x of many independent random variables is approximately normally/Gaussian distributed:

$$\mathcal{N}(\mathbf{x}=x;\boldsymbol{\mu},\boldsymbol{\sigma}^2) = \sqrt{\frac{1}{2\pi\sigma^2}} \exp\left(-\frac{1}{2\sigma^2}(x-\boldsymbol{\mu})^2\right).$$

• Holds regardless of the original distributions of individual variables •  $\mu_x=\mu$  and  $\sigma_x^2=\sigma^2$ 

#### Theorem (Central Limit Theorem)

The sum x of many independent random variables is approximately normally/Gaussian distributed:

$$\mathcal{N}(\mathbf{x}=x;\boldsymbol{\mu},\boldsymbol{\sigma}^2) = \sqrt{\frac{1}{2\pi\sigma^2}} \exp\left(-\frac{1}{2\sigma^2}(x-\boldsymbol{\mu})^2\right).$$

- Holds regardless of the original distributions of individual variables
- $\mu_{\mathrm{x}} = \mu$  and  $\sigma_{\mathrm{x}}^2 = \sigma^2$
- To avoid inverting  $\sigma^2$ , we can parametrize the distribution using the *precision*  $\beta$ :

$$\mathcal{N}(\mathbf{x}=x;\boldsymbol{\mu},\boldsymbol{\beta}^{-1}) = \sqrt{\frac{\beta}{2\pi}} \exp\left(-\frac{\beta}{2}(x-\boldsymbol{\mu})^2\right)$$

Shan-Hung Wu (CS, NTHU)

### **Confidence Intervals**



**Figure:** Graph of  $\mathcal{N}(\mu, \sigma^2)$ .

Shan-Hung Wu (CS, NTHU)

Prob. & Info. Theory

Machine Learning 45 / 78

### **Confidence Intervals**



**Figure:** Graph of  $\mathcal{N}(\mu, \sigma^2)$ .

• We say the interval  $[\mu - 2\sigma, \mu + 2\sigma]$  has about the 95% confidence

Shan-Hung Wu (CS, NTHU)

Shan-Hung Wu (CS, NTHU)

Prob. & Info. Theory

Machine Learning 46 / 78

1 It can model noise in data (e.g., Gaussian white noise)

• Can be considered to be the accumulation of a large number of small independent latent factors affecting data collection process

1) It can model noise in data (e.g., Gaussian white noise)

- Can be considered to be the accumulation of a large number of small independent latent factors affecting data collection process
- <sup>2</sup> Out of all possible probability distributions (over real numbers) with the same variance, it encodes the maximum amount of uncertainty
  - Assuming  $P(y \,|\, x) \sim \mathscr{N}$  , we insert the least amount of prior knowledge into a model

1 It can model noise in data (e.g., Gaussian white noise)

- Can be considered to be the accumulation of a large number of small independent latent factors affecting data collection process
- <sup>2</sup> Out of all possible probability distributions (over real numbers) with the same variance, it encodes the maximum amount of uncertainty
  - Assuming  $P(y \,|\, x) \sim \mathscr{N}$  , we insert the least amount of prior knowledge into a model
- ③ Convenient for many analytical manipulations
  - Closed under affine transformation, summation, marginalization, conditioning, etc.
  - Many of the integrals involving Gaussian distributions that arise in practice have simple closed form solutions

#### Properties

 Closed under affine transformation: if x ~ N(μ, σ<sup>2</sup>), then ax+b ~ N(aμ+b, a<sup>2</sup>σ<sup>2</sup>) for any deterministic a, b ∈ ℝ, a ≠ 0 [Proof]
 z = x-μ/σ ~ N(0,1) the *z*-normalization or standardization of x

#### Properties

- Closed under affine transformation: if x ~ N(μ, σ<sup>2</sup>), then ax + b ~ N(aμ + b, a<sup>2</sup>σ<sup>2</sup>) for any deterministic a, b ∈ ℝ, a ≠ 0 [Proof]
   z = x-μ/σ ~ N(0,1) the *z*-normalization or standardization of x
- Closed under summation: if  $x^{(1)} \sim \mathcal{N}(\mu^{(1)}, \sigma^{2(1)})$  is independent with  $x^{(2)} \sim \mathcal{N}(\mu^{(2)}, \sigma^{2(2)})$ , then  $x^{(1)} + x^{(2)} \sim \mathcal{N}(\mu^{(1)} + \mu^{(2)}, \sigma^{2(1)} + \sigma^{2(2)})$  [Homework:  $p_{x^{(1)}+x^{(2)}}(x) = \int p_{x^{(1)}}(x-y)p_{x^{(2)}}(y)dy$  the convolution]

#### Properties

- Closed under affine transformation: if x ~ (μ, σ<sup>2</sup>), then ax + b ~ (aμ + b, a<sup>2</sup>σ<sup>2</sup>) for any deterministic a, b ∈ ℝ, a ≠ 0 [Proof]
   z = x-μ/σ ~ (0,1) the *z*-normalization or standardization of x
- Closed under summation: if  $\mathbf{x}^{(1)} \sim \mathcal{N}(\boldsymbol{\mu}^{(1)}, \sigma^{2(1)})$  is independent with  $\mathbf{x}^{(2)} \sim \mathcal{N}(\boldsymbol{\mu}^{(2)}, \sigma^{2(2)})$ , then  $\mathbf{x}^{(1)} + \mathbf{x}^{(2)} \sim \mathcal{N}(\boldsymbol{\mu}^{(1)} + \boldsymbol{\mu}^{(2)}, \sigma^{2(1)} + \sigma^{2(2)})$ [Homework:  $p_{\mathbf{x}^{(1)} + \mathbf{x}^{(2)}}(x) = \int p_{\mathbf{x}^{(1)}}(x - y) p_{\mathbf{x}^{(2)}}(y) dy$  the convolution] • *Not* true if  $\mathbf{x}^{(1)}$  and  $\mathbf{x}^{(2)}$  are dependent

• When **x** is sum of many random vectors:

$$\mathcal{N}(\mathbf{x} = \mathbf{x}; \boldsymbol{\mu}, \boldsymbol{\Sigma}) = \sqrt{\frac{1}{(2\pi)^d \det(\boldsymbol{\Sigma})}} \exp\left[-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^\top \boldsymbol{\Sigma}^{-1}(\mathbf{x} - \boldsymbol{\mu})\right]$$

• 
$$\mu_{\mathbf{x}} = \mu$$
 and  $\Sigma_{\mathbf{x}} = \Sigma$  (must be nonsingular)

• When **x** is sum of many random vectors:

$$\mathcal{N}(\mathbf{x} = \mathbf{x}; \boldsymbol{\mu}, \boldsymbol{\Sigma}) = \sqrt{\frac{1}{(2\pi)^d \det(\boldsymbol{\Sigma})}} \exp\left[-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^\top \boldsymbol{\Sigma}^{-1}(\mathbf{x} - \boldsymbol{\mu})\right]$$

•  $\mu_{\mathbf{x}} = \mu$  and  $\Sigma_{\mathbf{x}} = \Sigma$  (must be nonsingular)

- If  $\mathbf{x} \sim \mathscr{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$ , then each attribute  $x_i$  is univariate normal
  - Converse not true

• When x is sum of many random vectors:

$$\mathcal{N}(\mathbf{x} = \mathbf{x}; \boldsymbol{\mu}, \boldsymbol{\Sigma}) = \sqrt{\frac{1}{(2\pi)^d \det(\boldsymbol{\Sigma})}} \exp\left[-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^\top \boldsymbol{\Sigma}^{-1}(\mathbf{x} - \boldsymbol{\mu})\right]$$

•  $\mu_{\mathbf{x}} = \mu$  and  $\Sigma_{\mathbf{x}} = \Sigma$  (must be nonsingular)

• If  $\mathbf{x} \sim \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$ , then each attribute  $x_i$  is univariate normal

- Converse not true
- However, if  $x_1, \dots, x_d$  are Gaussian and independent with each other, then  $\mathbf{x} \sim \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$ , where  $\boldsymbol{\mu} = [\boldsymbol{\mu}_1, \dots, \boldsymbol{\mu}_d]^\top$  and  $\boldsymbol{\Sigma} = \operatorname{diag}(\sigma_1^2, \dots, \sigma_d^2)$

Shan-Hung Wu (CS, NTHU)

• When x is sum of many random vectors:

$$\mathcal{N}(\mathbf{x} = \mathbf{x}; \boldsymbol{\mu}, \boldsymbol{\Sigma}) = \sqrt{\frac{1}{(2\pi)^d \det(\boldsymbol{\Sigma})}} \exp\left[-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^\top \boldsymbol{\Sigma}^{-1}(\mathbf{x} - \boldsymbol{\mu})\right]$$

•  $\mu_{\mathbf{x}} = \mu$  and  $\Sigma_{\mathbf{x}} = \Sigma$  (must be nonsingular)

• If  $\mathbf{x} \sim \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$ , then each attribute  $x_i$  is univariate normal

- Converse not true
- However, if  $x_1, \dots, x_d$  are Gaussian and independent with each other, then  $\mathbf{x} \sim \mathscr{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$ , where  $\boldsymbol{\mu} = [\boldsymbol{\mu}_1, \dots, \boldsymbol{\mu}_d]^\top$  and  $\boldsymbol{\Sigma} = \text{diag}(\sigma_1^2, \dots, \sigma_d^2)$
- What does the graph of  $\mathscr{N}(\mu, \Sigma)$  look like?

Shan-Hung Wu (CS, NTHU)

#### **Bivariate Example I**

• Consider the *Mahalanobis distance* first

$$\mathcal{N}(\boldsymbol{\mu},\boldsymbol{\Sigma}) = \sqrt{\frac{1}{(2\pi)^d \det(\boldsymbol{\Sigma})}} \exp\left[-\frac{1}{2}(\boldsymbol{x}-\boldsymbol{\mu})^\top \boldsymbol{\Sigma}^{-1}(\boldsymbol{x}-\boldsymbol{\mu})\right]$$
### Bivariate Example I

• Consider the *Mahalanobis distance* first

$$\mathcal{N}(\boldsymbol{\mu},\boldsymbol{\Sigma}) = \sqrt{\frac{1}{(2\pi)^d \det(\boldsymbol{\Sigma})}} \exp\left[-\frac{1}{2}(\boldsymbol{x}-\boldsymbol{\mu})^\top \boldsymbol{\Sigma}^{-1}(\boldsymbol{x}-\boldsymbol{\mu})\right]$$



Shan-Hung Wu (CS, NTHU)





Cov(x1,x2)<0

• The level sets closer to the center  $\mu_x$  are lower

### Bivariate Example I

• Consider the *Mahalanobis distance* first

$$\mathcal{N}(\boldsymbol{\mu},\boldsymbol{\Sigma}) = \sqrt{\frac{1}{(2\pi)^d \det(\boldsymbol{\Sigma})}} \exp\left[-\frac{1}{2}(\boldsymbol{x}-\boldsymbol{\mu})^\top \boldsymbol{\Sigma}^{-1}(\boldsymbol{x}-\boldsymbol{\mu})\right]$$



Shan-Hung Wu (CS, NTHU)



 $Cov(x_1,x_2)=0, Var(x_2)>Var(x_2)$ 

Cov(x1,x2)<0



- The level sets closer to the center  $\mu_{\mathbf{x}}$  are lower
- Increasing  $Cov[x_1,x_2]$  stretches the level sets along the  $45^\circ$  axis
- Decreasing  $Cov[x_1,x_2]$  stretches the level sets along the  $-45^\circ$  axis

#### **Bivariate Example II**

• The hight of  $\mathscr{N}(\mu, \Sigma) = \sqrt{\frac{1}{(2\pi)^d \det(\Sigma)}} \exp\left[-\frac{1}{2}(\mathbf{x} - \mu)^\top \Sigma^{-1}(\mathbf{x} - \mu)\right]$  in its graph is inversely proportional to the Mahalanobis distance



### Bivariate Example II

• The hight of  $\mathscr{N}(\mu, \Sigma) = \sqrt{\frac{1}{(2\pi)^d \det(\Sigma)}} \exp\left[-\frac{1}{2}(\mathbf{x} - \mu)^\top \Sigma^{-1}(\mathbf{x} - \mu)\right]$  in its graph is inversely proportional to the Mahalanobis distance



#### • A multivariate Gaussian distribution is *isotropic* iff $\Sigma = \sigma I$

Shan-Hung Wu (CS, NTHU)

#### Properties

- Closed under affine transformation: if  $\mathbf{x} \sim \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$ , then  $\mathbf{w}^{\top} \mathbf{x} \sim \mathcal{N}(\mathbf{w}^{\top} \boldsymbol{\mu}, \mathbf{w}^{\top} \boldsymbol{\Sigma} \mathbf{w})$  for any deterministic  $\mathbf{w} \in \mathbb{R}^d$ 
  - More generally, given  $W \in \mathbb{R}^{d \times k}$ , k < d, we have  $W^{\top} \mathbf{x} \sim \mathcal{N}(W^{\top} \mu, W^{\top} \Sigma W)$  that is k-variate normal
  - I.e., the projection of  $\mathbf{x}$  onto a k-dimensional subspace is still normal

#### Properties

#### Properties

- Closed under affine transformation: if x ~ N(μ,Σ), then w<sup>T</sup>x ~ N(w<sup>T</sup>μ, w<sup>T</sup>Σw) for any deterministic w ∈ ℝ<sup>d</sup> More generally, given W ∈ ℝ<sup>d×k</sup>, k < d, we have W<sup>T</sup>x ~ N(W<sup>T</sup>μ, W<sup>T</sup>ΣW) that is k-variate normal
  I.e., the projection of x onto a k-dimensional subspace is still normal

  Consider x = 

  x<sub>1</sub>
  N(μ = 

  μ<sub>1</sub>
  μ<sub>2</sub>
  S<sub>2,1</sub>
  Σ<sub>2,2</sub>

  Closed under marginalization: x<sub>1</sub> ~ N(μ<sub>1</sub>,Σ<sub>1,1</sub>) [Proof: P(x<sub>1</sub>) = 

  x<sub>2</sub>
  µ,Σ)dx<sub>2</sub>
- Closed under conditioning:  $(\mathbf{x}_1 | \mathbf{x}_2) \sim \mathscr{N}(\boldsymbol{\mu}_1 + \boldsymbol{\Sigma}_{1,2}\boldsymbol{\Sigma}_{2,2}^{-1}(\mathbf{x}_2 - \boldsymbol{\mu}_2), \boldsymbol{\Sigma}_{1,1} - \boldsymbol{\Sigma}_{1,2}\boldsymbol{\Sigma}_{2,2}^{-1}\boldsymbol{\Sigma}_{2,1}) \text{ [Proof]}$

Shan-Hung Wu (CS, NTHU)

### **Exponential Distribution (Continuous)**

 $\bullet\,$  In deep learning, we often want to have a probability distribution with a sharp point at x=0

### **Exponential Distribution (Continuous)**

- In deep learning, we often want to have a probability distribution with a sharp point at  $\mathbf{x} = \mathbf{0}$
- To accomplish this, we can use the *exponential distribution*:

Exponential 
$$(x = x; \lambda) = \lambda 1(x; x \ge 0) \exp(-\lambda x)$$



### Laplace Distribution (Continuous)

 Laplace distribution can be think of as a "two-sided" exponential distribution centered at μ:

Laplace 
$$(\mathbf{x} = x; \boldsymbol{\mu}, b) = \frac{1}{2b} \exp\left(-\frac{|\mathbf{x} - \boldsymbol{\mu}|}{b}\right)$$



Shan-Hung Wu (CS, NTHU)

# Dirac Distribution (Continuous)

 $\bullet\,$  In some cases, we wish to specify that all of the mass in a probability distribution clusters around a single data point  $\mu\,$ 

# Dirac Distribution (Continuous)

- In some cases, we wish to specify that all of the mass in a probability distribution clusters around a single data point  $\mu$
- This can be accomplished by using the *Dirac distribution*:

 $\operatorname{Dirac}(\mathbf{x}=\mathbf{x};\boldsymbol{\mu})=\boldsymbol{\delta}(\mathbf{x}-\boldsymbol{\mu}),$ 

where  $\delta(\cdot)$  is the Dirac delta function that

- f 1 Is zero-valued everywhere except at input f 0
- Integrals to 1

# **Empirical Distribution (Continuous)**

- Given a dataset  $\mathbb{X} = \{ \pmb{x}^{(i)} \}_{i=1}^N$  where  $\pmb{x}^{(i)}$ 's are i.i.d. samples of  $\pmb{\mathrm{x}}$
- What is the distribution  $P(\theta)$  that maximizes the likelihood  $P(\theta|\mathbb{X})$  of  $\mathbb{X}?$

# **Empirical Distribution (Continuous)**

- Given a dataset  $\mathbb{X} = \{ \pmb{x}^{(i)} \}_{i=1}^N$  where  $\pmb{x}^{(i)}$ 's are i.i.d. samples of  $\pmb{\mathrm{x}}$
- What is the distribution  $P(\theta)$  that maximizes the likelihood  $P(\theta|X)$  of X?
- If **x** is discrete, the distribution simply reflects the empirical frequency of values:

$$\text{Empirical}(\mathbf{x} = \mathbf{x}; \mathbb{X}) = \frac{1}{N} \sum_{i=1}^{N} \mathbb{1}(\mathbf{x}; \mathbf{x} = \mathbf{x}^{(i)})$$

# **Empirical Distribution (Continuous)**

- Given a dataset  $\mathbb{X} = \{m{x}^{(i)}\}_{i=1}^N$  where  $m{x}^{(i)}$ 's are i.i.d. samples of  $m{x}$
- What is the distribution  $P(\theta)$  that maximizes the likelihood  $P(\theta|X)$  of X?
- If **x** is discrete, the distribution simply reflects the empirical frequency of values:

Empirical
$$(\mathbf{x} = \mathbf{x}; \mathbb{X}) = \frac{1}{N} \sum_{i=1}^{N} \mathbb{1}(\mathbf{x}; \mathbf{x} = \mathbf{x}^{(i)})$$

• If **x** is continuous, we have the *empirical distribution*:

Empirical
$$(\mathbf{x} = \mathbf{x}; \mathbb{X}) = \frac{1}{N} \sum_{i=1}^{N} \delta(\mathbf{x} - \mathbf{x}^{(i)})$$

Shan-Hung Wu (CS, NTHU)

• We may define a probability distribution by combining other simpler probability distributions  $\{\mathbf{P}^{(i)}(\boldsymbol{\theta}^{(i)})\}_i$ 

- We may define a probability distribution by combining other simpler probability distributions  $\{\mathbf{P}^{(i)}(\boldsymbol{\theta}^{(i)})\}_i$
- E.g., the *mixture model*:

Mixture 
$$(\mathbf{x} = \mathbf{x}; \boldsymbol{\rho}, \{\boldsymbol{\theta}^{(i)}\}_i) = \sum_i \mathbf{P}^{(i)}(\mathbf{x} = \mathbf{x} | \mathbf{c} = i; \boldsymbol{\theta}^{(i)})$$
Categorical  $(\mathbf{c} = i; \boldsymbol{\rho})$ 

- We may define a probability distribution by combining other simpler probability distributions  $\{\mathbf{P}^{(i)}(\boldsymbol{\theta}^{(i)})\}_i$
- E.g., the *mixture model*:

Mixture
$$(\mathbf{x} = \mathbf{x}; \boldsymbol{\rho}, \{\boldsymbol{\theta}^{(i)}\}_i) = \sum_i \mathbf{P}^{(i)}(\mathbf{x} = \mathbf{x} | \mathbf{c} = i; \boldsymbol{\theta}^{(i)})$$
Categorical $(\mathbf{c} = i; \boldsymbol{\rho})$ 

• The empirical distribution is a mixture distribution (where  $ho_i = 1/N$ )

- We may define a probability distribution by combining other simpler probability distributions {P<sup>(i)</sup>(θ<sup>(i)</sup>)}<sub>i</sub>
- E.g., the *mixture model*:

Mixture
$$(\mathbf{x} = \mathbf{x}; \boldsymbol{\rho}, \{\boldsymbol{\theta}^{(i)}\}_i) = \sum_i \mathbf{P}^{(i)}(\mathbf{x} = \mathbf{x} | \mathbf{c} = i; \boldsymbol{\theta}^{(i)})$$
Categorical $(\mathbf{c} = i; \boldsymbol{\rho})$ 

The empirical distribution is a mixture distribution (where ρ<sub>i</sub> = 1/N)
 The component identity variable c is a *latent variable*

• Whose values are not observed

Shan-Hung Wu (CS, NTHU)

#### Gaussian Mixture Model

• A mixture model is called the *Gaussian mixture model* iff  $P^{(i)}(\mathbf{x} = \mathbf{x} | \mathbf{c} = i; \theta^{(i)}) = \mathcal{N}^{(i)}(\mathbf{x} = \mathbf{x} | \mathbf{c} = i; \mu^{(i)}, \Sigma^{(i)}), \forall i$ 

#### Gaussian Mixture Model

• A mixture model is called the *Gaussian mixture model* iff  $P^{(i)}(\mathbf{x} = \mathbf{x} | \mathbf{c} = i; \theta^{(i)}) = \mathcal{N}^{(i)}(\mathbf{x} = \mathbf{x} | \mathbf{c} = i; \mu^{(i)}, \Sigma^{(i)}), \forall i$ 

• Variants:  $\Sigma^{(i)} = \Sigma$  or  $\Sigma^{(i)} = \text{diag}(\sigma)$  or  $\Sigma^{(i)} = \sigma I$ 

### Gaussian Mixture Model

• A mixture model is called the *Gaussian mixture model* iff  $P^{(i)}(\mathbf{x} = \mathbf{x} | \mathbf{c} = i; \boldsymbol{\theta}^{(i)}) = \mathcal{N}^{(i)}(\mathbf{x} = \mathbf{x} | \mathbf{c} = i; \boldsymbol{\mu}^{(i)}, \boldsymbol{\Sigma}^{(i)}), \forall i$ 

• Variants:  $\Sigma^{(i)} = \Sigma$  or  $\Sigma^{(i)} = \text{diag}(\sigma)$  or  $\Sigma^{(i)} = \sigma I$ 

• Any smooth density can be approximated by a Gaussian mixture model with enough components



# Outline

- 1 Random Variables & Probability Distributions
- 2 Multivariate & Derived Random Variables
- **3** Bayes' Rule & Statistics
- Application: Principal Components Analysis
- 5 Technical Details of Random Variables
- 6 Common Probability Distributions
- **7** Common Parametrizing Functions
- 8 Information Theory
- 9 Application: Decision Trees & Random Forest

Shan-Hung Wu (CS, NTHU)

#### **Parametrizing Functions**

- A probability distribution  $P(\theta)$  is parametrized by  $\theta$
- $\, \bullet \,$  In ML,  $\, \theta \,$  may be the output value of a deterministic function
  - Called *parametrizing function*

### **Logistic Function**

• The *logistic function* (a special case of *sigmoid functions*) is defined as:

$$\sigma(x) = \frac{\exp(x)}{\exp(x) + 1} = \frac{1}{1 + \exp(-x)}$$



### **Logistic Function**

• The *logistic function* (a special case of *sigmoid functions*) is defined as:



• Always takes on values between (0,1)

Shan-Hung Wu (CS, NTHU)

### **Logistic Function**

• The *logistic function* (a special case of *sigmoid functions*) is defined as:



- Always takes on values between (0,1)
- $\circ$  Commonly used to produce the ho parameter of Bernoulli distribution

Shan-Hung Wu (CS, NTHU)

### **Softplus Function**

• The *softplus function* :

$$\zeta(x) = \log(1 + \exp(x))$$



Shan-Hung Wu (CS, NTHU)

Prob. & Info. Theory

Machine Learning 61 / 78

### **Softplus Function**

• The *softplus function* :

$$\zeta(x) = \log(1 + \exp(x))$$



• A "softened" version of  $x^+ = \max(0, x)$ 

Shan-Hung Wu (CS, NTHU)

### **Softplus Function**

• The *softplus function* :

$$\zeta(x) = \log(1 + \exp(x))$$



• A "softened" version of  $x^+ = \max(0, x)$ 

- Range:  $(0,\infty)$
- ${\, \bullet \, }$  Useful for producing the  $\beta$  or  $\sigma$  parameter of Gaussian distribution

Shan-Hung Wu (CS, NTHU)

### Properties [Homework]

• 
$$1 - \sigma(x) = \sigma(-x)$$
  
•  $\log \sigma(x) = -\zeta(-x)$   
•  $\frac{d}{dx}\sigma(x) = \sigma(x)(1 - \sigma(x))$   
•  $\frac{d}{dx}\zeta(x) = \sigma(x)$   
•  $\forall x \in (0,1), \sigma^{-1}(x) = \log\left(\frac{x}{1-x}\right)$   
•  $\forall x > 0, \zeta^{-1}(x) = \log(\exp(x) - 1)$   
•  $\zeta(x) = \int_{-\infty}^{x} \sigma(y) dy$   
•  $\zeta(x) - \zeta(-x) = x$   
•  $\zeta(-x)$  is the softened  $x^{-} = \max(0, -x)$   
•  $x = x^{+} - x^{-}$ 

Shan-Hung Wu (CS, NTHU)

# Outline

- 1 Random Variables & Probability Distributions
- 2 Multivariate & Derived Random Variables
- **3** Bayes' Rule & Statistics
- 4 Application: Principal Components Analysis
- 5 Technical Details of Random Variables
- 6 Common Probability Distributions
- 7 Common Parametrizing Functions

#### 8 Information Theory

9 Application: Decision Trees & Random Forest

Shan-Hung Wu (CS, NTHU)

### What's Information Theory

 Probability theory allows us to make uncertain statements and reason in the presence of uncertainty

### What's Information Theory

- Probability theory allows us to make uncertain statements and reason in the presence of uncertainty
- Information theory allows us to *quantify* the amount of uncertainty

### **Self-Information**

• Given a random variable x, how much information you receive when seeing an event x = x?

### Self-Information

- Given a random variable x, how much information you receive when seeing an event x = x?
- 1 Likely events should have low information
  - E.g., we are less surprised when tossing a biased coins
## Self-Information

- Given a random variable x, how much information you receive when seeing an event x = x?
- 1 Likely events should have low information
  - E.g., we are less surprised when tossing a biased coins
- 2 Independent events should have additive information
  - E.g, "two heads" should have twice as much info as "one head"

## Self-Information

- Given a random variable x, how much information you receive when seeing an event x = x?
- Likely events should have low information
  - $\, \bullet \,$  E.g., we are less surprised when tossing a biased coins
- 2 Independent events should have additive information
  - E.g, "two heads" should have twice as much info as "one head"
  - The *self-information*:

$$I(x = x) = -\log P(x = x)$$

Shan-Hung Wu (CS, NTHU)

## Self-Information

- Given a random variable x, how much information you receive when seeing an event x = x?
- Likely events should have low information
  - E.g., we are less surprised when tossing a biased coins
- 2 Independent events should have additive information
  - E.g, "two heads" should have twice as much info as "one head"
  - The *self-information*:

$$I(x = x) = -\log P(x = x)$$

- Called bit if base-2 logarithm is used
- Called *nat* if base-*e*

Shan-Hung Wu (CS, NTHU)

• Self-information deals with a particular outcome

- Self-information deals with a particular outcome
- We can quantify the amount of uncertainty in an entire probability distribution using the *entropy*:

$$H(\mathbf{x} \sim \mathbf{P}) = \mathbf{E}_{\mathbf{x} \sim \mathbf{P}}[\mathbf{I}(\mathbf{x})] = -\sum_{x} P(x) \log P(x) \text{ or } -\int p(x) \log p(x) dx$$

• Let 
$$0\log 0 = \lim_{x\to 0} x\log x = 0$$

0

- Self-information deals with a particular outcome
- We can quantify the amount of uncertainty in an entire probability distribution using the *entropy*:

$$H(\mathbf{x} \sim \mathbf{P}) = E_{\mathbf{x} \sim \mathbf{P}}[\mathbf{I}(\mathbf{x})] = -\sum_{x} P(x) \log P(x) \text{ or } -\int p(x) \log p(x) dx$$

- Let  $0\log 0 = \lim_{x\to 0} x\log x = 0$
- Called **Shannon entropy** when x is discrete; **differential entropy** when x is continuous

- Self-information deals with a particular outcome
- We can quantify the amount of uncertainty in an entire probability distribution using the *entropy*:

$$H(\mathbf{x} \sim \mathbf{P}) = E_{\mathbf{x} \sim \mathbf{P}}[\mathbf{I}(\mathbf{x})] = -\sum_{x} P(x) \log P(x) \text{ or } -\int p(x) \log p(x) dx$$

- Let  $0\log 0 = \lim_{x\to 0} x\log x = 0$
- Called *Shannon entropy* when x is discrete; *differential entropy* when x is continuous



**Figure:** Shannon entropy H(x) over Bernoulli distributions with different  $\rho$ .

Shan-Hung Wu (CS, NTHU)

### Average Code Length

• Shannon entropy gives a lower bound on the number of "bits" needed on average to encode values drawn from a distribution P

### Average Code Length

- Shannon entropy gives a lower bound on the number of "bits" needed on average to encode values drawn from a distribution P
- $\,$  o Consider a random variable x  $\sim$  Uniform having 8 equally likely states
  - To send a value x to receiver, we would encode it into 3 bits
  - Shannon entropy:  $H(x \sim Uniform) = -8 \times \frac{1}{8} \log_2 \frac{1}{8} = 3$

### Average Code Length

- Shannon entropy gives a lower bound on the number of "bits" needed on average to encode values drawn from a distribution P
- $\,$  o Consider a random variable x  $\sim$  Uniform having 8 equally likely states
  - To send a value x to receiver, we would encode it into 3 bits
  - Shannon entropy:  $H(x \sim Uniform) = -8 \times \frac{1}{8} \log_2 \frac{1}{8} = 3$
- If the probabilities of the 8 states are  $(\frac{1}{2}, \frac{1}{4}, \frac{1}{8}, \frac{1}{16}, \frac{1}{64}, \frac{1}{64}, \frac{1}{64}, \frac{1}{64})$  instead
  - H(x) = 2
  - The encoding 0, 10, 110, 1110, 111100, 111101, 111110, 111111 gives the average code length 2

# Kullback-Leibler (KL) Divergence

• How many extra "bits" needed in average to transmit a value drawn from distribution P when we use a code that was designed for another distribution Q?

# Kullback-Leibler (KL) Divergence

- How many extra "bits" needed in average to transmit a value drawn from distribution P when we use a code that was designed for another distribution Q?
- *Kullback-Leibler (KL) Divergence* or (*relative entropy*) from distribution Q to P:

$$D_{KL}(P \| Q) = E_{x \sim P} \left[ log \frac{P(x)}{Q(x)} \right] = -E_{x \sim P} \left[ log Q(x) \right] - H(x \sim P)$$

 $\bullet~$  The term  $-E_{x\sim P}[log\,Q(x)]$  is called the cross~entropy

# Kullback-Leibler (KL) Divergence

- How many extra "bits" needed in average to transmit a value drawn from distribution P when we use a code that was designed for another distribution Q?
- *Kullback-Leibler (KL) Divergence* or (*relative entropy*) from distribution Q to P:

$$D_{KL}(P \| Q) = E_{x \sim P} \left[ log \frac{P(x)}{Q(x)} \right] = -E_{x \sim P} \left[ log Q(x) \right] - H(x \sim P)$$

• The term  $-E_{x\sim P}[\log Q(x)]$  is called the *cross entropy* • If P and Q are independent, we can solve

$$\arg\min_{Q} D_{KL}(P||Q)$$

by

$$\arg\min_{Q} - E_{x \sim P} [\log Q(x)]$$

Shan-Hung Wu (CS, NTHU)

#### Properties

- $\bullet \ D_{KL}(P\|Q) \geq 0 \text{, } \forall P,Q$
- $\bullet \ D_{KL}(P\|Q) = 0$  iff P and Q are equal almost surely
- $\bullet~$  KL divergence is asymmetric, i.e.,  $D_{KL}(P\|Q) \neq D_{KL}(Q\|P)$



Figure: KL divergence for two normal distributions.

Shan-Hung Wu (CS, NTHU)

### Minimizer of KL Divergence

• Given P, we want to find  $Q^*$  that minimizes the KL divergence •  $Q^{*(from)} = \arg \min_{O} D_{KL}(P||Q)$  or  $Q^{*(to)} = \arg \min_{O} D_{KL}(Q||P)$ ?

# Minimizer of KL Divergence

• Given P, we want to find  $Q^*$  that minimizes the KL divergence •  $Q^{*(from)} = arg min_Q D_{KL}(P \| Q)$  or  $Q^{*(to)} = arg min_Q D_{KL}(Q \| P)$ ? •  $Q^{*(from)}$  places high probability where P has high probability •  $Q^{*(to)}$  places low probability where P has low probability



Figure: Approximating a mixture P of two Gaussians using a single Gaussian Q.

Shan-Hung Wu (CS, NTHU)

# Outline

- 1 Random Variables & Probability Distributions
- 2 Multivariate & Derived Random Variables
- **3** Bayes' Rule & Statistics
- Application: Principal Components Analysis
- 5 Technical Details of Random Variables
- 6 Common Probability Distributions
- 7 Common Parametrizing Functions
- Information Theory

#### 9 Application: Decision Trees & Random Forest

Shan-Hung Wu (CS, NTHU)

### **Decision Trees**

- Given a supervised dataset  $\mathbb{X} = \{(\mathbf{x}^{(i)}, y^{(i)})\}_{i=1}^N$
- Can we find out a tree-like function f (i.e, a set of rules) such that  $f(\pmb{x}^{(i)}) = y^{(i)}$ ?



- Start from root which corresponds to all data points  $\{(\pmb{x}^{(i)}, y^{(i)}): \mathsf{Rules} = \pmb{\emptyset})\}$
- Recursively split leaf nodes until data corresponding to children are "pure" in labels

- Start from root which corresponds to all data points  $\{(\pmb{x}^{(i)}, y^{(i)}): \mathsf{Rules} = \pmb{\emptyset})\}$
- Recursively split leaf nodes until data corresponding to children are "pure" in labels
- How to split?

- Start from root which corresponds to all data points  $\{(\mathbf{x}^{(i)}, y^{(i)}) : \mathsf{Rules} = \mathbf{0})\}$
- Recursively split leaf nodes until data corresponding to children are "pure" in labels
- How to split? Find a cutting point (j, v) among all unseen attributes such that after partitioning the corresponding data points  $\mathbb{X}^{\text{parent}} = \{(\mathbf{x}^{(i)}, y^{(i)} : \text{Rules})\}$  into two groups



$$\begin{split} \mathbb{X}^{\mathsf{left}} &= \{ (\pmb{x}^{(i)}, y^{(i)}) : \mathsf{Rules} \cup \{ x_j^{(i)} < v \} \}, \text{ and} \\ \mathbb{X}^{\mathsf{right}} &= \{ (\pmb{x}^{(i)}, y^{(i)}) : \mathsf{Rules} \cup \{ x_j^{(i)} \ge v \} \}, \end{split}$$

the "impurity" of labels drops the most

- Start from root which corresponds to all data points  $\{(\mathbf{x}^{(i)}, y^{(i)}) : \text{Rules} = \mathbf{0})\}$
- Recursively split leaf nodes until data corresponding to children are "pure" in labels
- How to split? Find a cutting point (j, v) among all unseen attributes such that after partitioning the corresponding data points  $\mathbb{X}^{\text{parent}} = \{(\mathbf{x}^{(i)}, y^{(i)} : \text{Rules})\}$  into two groups

I C.

(.) (.)



$$\mathbb{X}^{\mathsf{lert}} = \{ (\boldsymbol{x}^{(i)}, y^{(i)}) : \mathsf{Rules} \cup \{ x_j^{(i)} < v \} \}, \text{ and}$$
$$\mathbb{X}^{\mathsf{right}} = \{ (\boldsymbol{x}^{(i)}, y^{(i)}) : \mathsf{Rules} \cup \{ x_j^{(i)} \ge v \} \},$$
the "impurity" of labels drops the most, i.e., solve
$$\arg \max_{j, v} \left( \mathsf{Impurity}(\mathbb{X}^{\mathsf{parent}}) - \mathsf{Impurity}(\mathbb{X}^{\mathsf{left}}, \mathbb{X}^{\mathsf{right}}) \right)$$

(:)

### Impurity Measure

$$\arg \max_{j,v} \left( \text{Impurity}(\mathbb{X}^{\mathsf{parent}}) - \text{Impurity}(\mathbb{X}^{\mathsf{left}}, \mathbb{X}^{\mathsf{right}}) \right)$$

• What's  $Impurity(\cdot)?$ 

### Impurity Measure

$$\arg\max_{j,v} \left( \mathrm{Impurity}(\mathbb{X}^{\mathsf{parent}}) - \mathrm{Impurity}(\mathbb{X}^{\mathsf{left}}, \mathbb{X}^{\mathsf{right}}) \right)$$

- What's  $Impurity(\cdot)$ ?
- Entropy is a common choice:

$$\begin{split} \text{Impurity}(\mathbb{X}^{\mathsf{parent}}) &= \mathbf{H}[y \sim \text{Empirical}(\mathbb{X}^{\mathsf{parent}})]\\ \\ \text{Impurity}(\mathbb{X}^{\mathsf{left}}, \mathbb{X}^{\mathsf{right}}) &= \sum_{i = \mathsf{left}, \mathsf{right}} \frac{|\mathbb{X}^{(i)}|}{|\mathbb{X}^{\mathsf{parent}}|} \mathbf{H}[y \sim \text{Empirical}(\mathbb{X}^{(i)})] \end{split}$$

Shan-Hung Wu (CS, NTHU)

### Impurity Measure

$$\arg\max_{j,v} \left( \mathrm{Impurity}(\mathbb{X}^{\mathsf{parent}}) - \mathrm{Impurity}(\mathbb{X}^{\mathsf{left}}, \mathbb{X}^{\mathsf{right}}) \right)$$

- What's  $Impurity(\cdot)$ ?
- Entropy is a common choice:

Impurity(
$$\mathbb{X}^{\mathsf{parent}}$$
) = H[ $y \sim \mathsf{Empirical}(\mathbb{X}^{\mathsf{parent}})$ ]

$$Impurity(\mathbb{X}^{\mathsf{left}}, \mathbb{X}^{\mathsf{right}}) = \sum_{i = \mathsf{left}, \mathsf{right}} \frac{|\mathbb{X}^{(i)}|}{|\mathbb{X}^{\mathsf{parent}}|} \mathbf{H}[y \sim \mathsf{Empirical}(\mathbb{X}^{(i)})]$$

 $\bullet~$  In this case,  $Impurity(\mathbb{X}^{parent})-Impurity(\mathbb{X}^{left},\mathbb{X}^{right})$  is called the information~gain

Shan-Hung Wu (CS, NTHU)

• A decision tree can be very deep

- A decision tree can be very deep
- Deeper nodes give more specific rules
  - Backed by less training data
  - May not be applicable to testing data
- How to ensure the generalizability of a decision tree?
  - I.e., to have high prediction accuracy on testing data

- A decision tree can be very deep
- Deeper nodes give more specific rules
  - Backed by less training data
  - May not be applicable to testing data
- How to ensure the *generalizability* of a decision tree?
  - I.e., to have high prediction accuracy on testing data
- Pruning (e.g., limit the depth of the tree)

- A decision tree can be very deep
- Deeper nodes give more specific rules
  - Backed by less training data
  - May not be applicable to testing data
- How to ensure the *generalizability* of a decision tree?
  - I.e., to have high prediction accuracy on testing data
- **1** Pruning (e.g., limit the depth of the tree)
- 2 Random forest: an ensemble of many (deep) trees

Shan-Hung Wu (CS, NTHU)

- 2 Grow a decision tree from the bootstrap samples. At each node:
  - **Randomly select** K features without replacement
  - 2 Find the best cutting point (j, v) and split the node

- 2 Grow a decision tree from the bootstrap samples. At each node:
  - **Randomly select** K features without replacement
  - 2 Find the best cutting point (j, v) and split the node
- 3 Repeat the steps 1 and 2 for T times to get T trees

- 2 Grow a decision tree from the bootstrap samples. At each node:
  - **Randomly select** K features without replacement
  - 2 Find the best cutting point (j, v) and split the node
- 3 Repeat the steps 1 and 2 for T times to get T trees
- Aggregate the predictions made by different trees via the *majority* vote

Randomly pick *M* samples from the training set with replacement
Called the *bootstrap* samples

- 2 Grow a decision tree from the bootstrap samples. At each node:
  - **Randomly select** K features without replacement
  - **2** Find the best cutting point (j, v) and split the node
- 3 Repeat the steps 1 and 2 for T times to get T trees
- Aggregate the predictions made by different trees via the *majority* vote
  - Each tree is trained slightly differently because of Step 1 and 2(a)
  - Provides different "perspectives" when voting

Shan-Hung Wu (CS, NTHU)

### **Decision Boundaries**



#### Decision Trees vs. Random Forests

• Cons of random forests:

Less interpretable model

Shan-Hung Wu (CS, NTHU)
## Decision Trees vs. Random Forests

- Cons of random forests:
  - Less interpretable model
- Pros:
  - Less sensitive to the depth of trees
    - The majority voting can "absorb" the noise from individual trees
  - Can be parallelized
    - Each tree can grow independently