

# • Probability

1. probability mass function (PMF)  $P(X=x) = f(x) \leq 1$
2. cumulative distribution function (CDF)  $P(X \leq x) = F(x)$
3. probability density function (PDF)  $f(a \leq x \leq b) = \int_a^b f(x) dx$   

$$f(x) = p^x (1-p)^{1-x}$$
  - Bernoulli Distribution  $E(x) = p$   $\text{Var}(x) = p(1-p)$
  - Geometric Distribution  $f(k) = (1-p)^{k-1} p$   $E(x) = \frac{1}{p}$   $\text{Var}(x) = \frac{1}{p^2} - \frac{1}{p}$
  - Binomial Distribution  $f(k) = \binom{n}{k} p^k (1-p)^{n-k}$   $E(x) = np$   
 $\text{Var}(x) = np(1-p)$

$$\frac{P(X=x)}{P(X=y)} = \frac{f(y)}{f(x)} \quad \text{密度} \uparrow, \text{概率} \downarrow$$

approximate  $\int_a^b h(x) dx$   $f(\lambda) = \frac{1}{b-a} \quad x \in [a, b]$

$$E[(b-a)h(x)] = \int_a^b (b-a)h(x)f(x)dx = \int_a^b h(x)dx$$

•  $x, z$  分布相同.  $\therefore E(g(x)) = E(g(z))$

$$2. E(\sum X_i) = \sum E(X_i)$$

$$\text{Var}(\sum X_i) = \sum \text{Var}(X_i)$$

$$\text{Proof: } Y_i = X_i - E(X_i) \text{ 常数}$$

$$E(Y_i) = E(X_i) - E(X_i) = 0$$

$$\text{Var}(\sum X_i) = E[(\sum (X_i - E(X_i))^2)]$$

$$= E[(\sum Y_i)^2] = \sum_{i,j} E[Y_i Y_j] + E[(Y_i)^2]$$

$$= \sum E[Y_i]^2 = \sum \text{Var}(X_i)$$

Monte Carlo

Uniform sampling  
&  
Law of Large numbers

离散性 点概率为0(但可能发生)

## 4. normal random variables

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

中心极限定理 CLT

$$\bar{x} \sim \text{Normal}(E(x), \frac{\sqrt{\text{Var}(x)}}{\sqrt{n}})$$

$$\sqrt{n}(\bar{x}_n - \mu) \xrightarrow{d} N(0, \sigma^2) \text{ (converges)}$$

## 5. correlation 相关性

$$\rho(x, y) = \frac{\text{Cov}(x, y)}{\sqrt{\text{Var}(x)} \sqrt{\text{Var}(y)}} \in [-1, 1]$$

$$\text{Cov}(x, y) = E[(x - E(x))(y - E(y))]$$

# • Statistic

## 1. likelihood function 似然函数

$$L(\theta | x) = P(x|\theta)$$

↓ 在已知  $\theta$  情况下，发生  $x$  的可能性大小

结果为入，分布函数的参数为  $\theta$  的可能性大小

## 2. 最大似然函数 Maximum Likelihood Estimation (MLE)

$$L(\theta, x) = P(x_1, x_2, \dots, x_n | \theta) = \prod_{i=1}^n P(x_i | \theta) \quad [\text{独立分布}]$$

$$\Rightarrow \log(L) = \sum_{i=1}^n \log(P_i) \quad \hat{\theta} = \underset{\substack{\leftarrow \text{最大参数} \\ \uparrow \text{参数空间}}}{\operatorname{argmax}_{\theta \in \Theta}} L(\theta, x) \quad \hat{\theta} \text{ 是所有可能 } \theta \text{ 中的最大值}$$

[e.g.] Uniform(a, b). find MLE of (a, b)

$$\log(L) = \log\left[\left(\frac{1}{b-a}\right)^n\right] = -n \log(b-a) \quad \hat{a} = \max\{x_i\}, \hat{b} = \min\{x_i\}$$

## 3. 最小二乘估计 Least Squares Estimation LSE

$$(1) \sum_i (\text{Mean} - x_i)^2$$

$$\downarrow I_A(x) = \begin{cases} 1, & x \in A \\ 0, & x \notin A \end{cases}$$

$$[\text{e.g.}] f(x; \theta) = \frac{1}{\theta} \mid \{0 \leq x \leq \theta\}$$

MLE:

$$L(\theta) = \frac{1}{\theta^n} \prod_{i=1}^n I_{\{0 \leq x_i \leq \theta\}} = \begin{cases} \frac{1}{\theta^n}, & \theta \geq \max\{x_i\} \\ 0, & \theta < \max\{x_i\} \end{cases}$$

$$\hat{\theta} = \max\{x_i\}$$

$$\text{LSE: } \sum_{i=1}^n \left( \frac{\theta}{2} - x_i \right)^2 = \sum_{i=1}^n \left( \frac{\theta^2}{4} + x_i^2 - \theta x_i \right)$$

$$\frac{\partial S}{\partial \theta} = - \sum_{i=1}^n x_i + n \frac{\theta}{2} = 0 \Rightarrow \hat{\theta} = \frac{2}{n} \sum_{i=1}^n x_i = 2\bar{x}$$

$$(2) \underset{\alpha, \beta}{\operatorname{Min}} S(\alpha, \beta) = \sum_i (Y_i - \alpha - \beta x_i)^2 \leftarrow X \& Y \text{ is linear}$$

$$\frac{\partial S}{\partial \alpha} = \sum_{i=1}^n 2(Y_i - \alpha - \beta x_i) (-1) = 0$$

$$\Rightarrow \sum_{i=1}^n (Y_i - \beta x_i - \alpha) = 0 \Rightarrow \alpha = \frac{1}{n} \sum_{i=1}^n (Y_i - \beta x_i) = \bar{Y} - \beta \bar{x}$$

$$\frac{\partial S}{\partial \beta} = \sum_{i=1}^n 2(Y_i - \alpha - \beta x_i) \cdot (-x_i) = 0$$

$$\Rightarrow \sum_{i=1}^n (Y_i - \beta x_i - \alpha) x_i = 0$$

$$\Rightarrow \sum_{i=1}^n x_i y_i - \beta \sum_{i=1}^n x_i^2 - 2 \sum_{i=1}^n x_i$$

$$= \sum_{i=1}^n x_i y_i - \beta \sum_{i=1}^n x_i^2 - (\bar{y} - \beta \bar{x} \cdot N \bar{x}) = 0$$

$$\Rightarrow \beta \left( \sum_{i=1}^n x_i^2 - N \bar{x}^2 \right) = \sum_{i=1}^n x_i y_i - N \bar{x} \bar{y}$$

$$\beta \left[ \sum_{i=1}^n (x_i - \bar{x})^2 \right] = \sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})$$

$$\Rightarrow \beta = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sum_{i=1}^n (x_i - \bar{x})^2} = \frac{\sum x_i y_i - n \bar{x} \bar{y}}{\sum (x_i)^2 - n \bar{x}^2}$$

•  $y - \alpha - \beta x \sim N(0, \sigma^2)$

$$L(\alpha, \beta) = \frac{1}{(2\pi\sigma^2)^n} \exp \left[ -\frac{1}{2} \frac{\sum (y_i - \beta x_i - \alpha)^2}{\sigma^2} \right] \Rightarrow \text{minimizes } \sum (y_i - \beta x_i - \alpha)^2$$

LSE  
↑

examine

| e<sub>i</sub> independent on x<sub>i</sub> ?

| variance of e<sub>i</sub> independent on x<sub>i</sub> ?

• KL-Divergence (KL 散度 用分布 Q 来估计分布 P 的真实分布的编码损失)

$$D_{KL}(P||Q) = \int P(x) \log \left( \frac{P(x)}{Q(x)} \right) dx$$

$$\theta_{\min KL} = D_{KL}(P(x) || P(x|\theta)) = E_{x \sim P(x)} \log \frac{P(x)}{P(x|\theta)}$$

$$\Rightarrow \theta_{\min KL} = \arg \min_{\theta} E_{x \sim P(x)} [-\log P(x|\theta)] = \arg \max_{\theta} E_{x \sim P(x)} [\log P(x|\theta)] = \hat{\theta}$$

[e.g.] ① Binomial

$$P(X=k) = \binom{n}{k} p^k (1-p)^{n-k}$$

$$l(p, k) = \log \binom{n}{k} + k \log p + (n-k) \log (1-p)$$

$$\hat{p} = \frac{k}{n}$$

② Normal

$$l(u, \sigma^2) = \sum -\frac{1}{2} \log (2\pi\sigma^2) - \frac{(x_i - u)^2}{2\sigma^2}$$

$$\frac{\partial l}{\partial u} = \sum \frac{x_i - u}{\sigma^2} = 0 \Rightarrow u = \frac{1}{n} \sum x_i$$

$$\frac{\partial l}{\partial \sigma^2} = \sum \left[ -\frac{1}{2\sigma^2} + \frac{(x_i - u)^2}{2\sigma^4} \right] = 0$$

$$\sigma^2 = \frac{1}{n} \sum (x_i - u)^2$$

## (4) confidence interval (CI) 置信区间

$$\bar{x} - \frac{b\sigma}{\sqrt{n}} \leq u \leq \bar{x} + \frac{a\sigma}{\sqrt{n}} \Rightarrow \frac{\sqrt{n}(\bar{x}-u)}{\sigma} \sim N(0,1)$$

$$\Rightarrow P\left(\bar{X} - \frac{b\sigma}{\sqrt{n}} \leq u \leq \bar{X} + \frac{a\sigma}{\sqrt{n}}\right) = P\left(-a \leq \frac{\sqrt{n}(\bar{X}-u)}{\sigma} \leq b\right) = \Phi(b) - \Phi(-a) = \Phi(b) + \Phi(a) - 1$$

• Bernoulli:  $\left[ p - 2\frac{\alpha}{n} \sqrt{\frac{p(1-p)}{n}}, p + 2\frac{\alpha}{n} \sqrt{\frac{p(1-p)}{n}} \right] \left( \left[ \bar{x} - \frac{2\frac{\alpha}{n}\sigma}{\sqrt{n}}, \bar{x} + \frac{2\frac{\alpha}{n}\sigma}{\sqrt{n}} \right] \right)$

$$\text{Var}(\hat{p}) = \text{Var}\left(\frac{Y}{n}\right) = \frac{1}{n^2} \text{Var}(Y) = \frac{1}{n^2} \cdot n p(1-p) = \frac{p(1-p)}{n}$$

[e.g.] find a  $1-\alpha$  CI. T?

w.p. (with probability)  $\alpha$ , it is not in T

$$\Rightarrow \forall a, b \quad \phi(b) + \phi(a) - 1 = 1 - \alpha$$

$$\text{let } a=b. \quad 2\phi(a) - 1 = 1 - \alpha$$

$\downarrow \alpha \downarrow$  . length of CI  $\uparrow$

$\downarrow n \uparrow$  .  $\downarrow$

$\downarrow \sigma^2$  .  $\downarrow$

## • optimization

1. feasible set  $\Omega$  (可行集)

$$2. B(x, \epsilon) = \{y : \|y-x\| \leq \epsilon\}$$

$$\forall y \in \Omega \cap B(x, \epsilon), f(x) \leq f(y)$$

$\left| \begin{array}{l} \text{for any } \epsilon > 0, x: \text{global} \\ \text{minimizer} \end{array} \right.$

$\left| \begin{array}{l} \text{for some } \epsilon > 0, x: \text{local} \\ \text{minimizer} \end{array} \right.$

$$3. \text{convex set } x = \theta x_1 + (1-\theta) x_2 \quad \theta \in [0,1]$$

$$\text{convex function } \forall 0 \leq \lambda \leq 1 \quad f(\lambda x + (1-\lambda) y) \leq \lambda f(x) + (1-\lambda) f(y)$$

•  $f(x)$  is a convex function and  $\Omega$  is a convex set.  $\min f(x) \text{ s.t. } x \in \Omega$

Then any local minimum is also a global minimum.

Proof: let  $\bar{x}$  be a local minimum

$$\Rightarrow \bar{x} \in \Omega. \exists \epsilon > 0 \text{ s.t. } f(\bar{x}) \leq f(x) \text{ for any } x \in \Omega \cap \{x : \|x - \bar{x}\| \leq \epsilon\} \quad ①$$

$$\text{assume } \exists z \in \Omega. f(z) < f(\bar{x})$$

we have  $\lambda \bar{x} + (1-\lambda) z \in \Omega \quad \forall \lambda \in [0,1]$  , then  $f(\lambda \bar{x} + (1-\lambda) z) \leq \lambda f(\bar{x}) + (1-\lambda) f(z)$

$$< \lambda f(\bar{x}) + (1-\lambda) f(\bar{x}) = f(\bar{x}) \quad \text{as } \lambda \rightarrow 1 \quad \lambda \bar{x} + (1-\lambda) z \rightarrow \bar{x}, \text{ which contradicts with } ①$$

• Prove  $H = \{(x, y) : y \leq ax + b\}$  is a convex set

$$(x, y) = \theta(x_1, y_1) + (1-\theta)(x_2, y_2)$$

$$\begin{cases} x = \theta x_1 + (1-\theta)x_2 \\ y = \theta y_1 + (1-\theta)y_2 \end{cases} \quad \theta \in [0, 1]$$

$$\begin{aligned} y_1 \leq ax_1 + b &\Rightarrow \theta y_1 + (1-\theta)y_2 \leq \theta(ax_1 + b) + (1-\theta)(ax_2 + b) \\ y_2 \leq ax_2 + b &= a(\theta x_1 + (1-\theta)x_2) + b \\ &\Rightarrow y \leq ax + b \end{aligned}$$

•  $S_1, S_2$  are convex sets  $\Rightarrow S_1 \cap S_2$  is a convex set

Proof:

$$\text{let } x_1, x_2 \in S_1 \cap S_2$$

$$x = \theta x_1 + (1-\theta)x_2 \quad (\theta \in [0, 1])$$

$S_1, S_2$  are convex sets

$$\Rightarrow x \in S_1, x \in S_2$$

$$\Rightarrow x \in S_1 \cap S_2$$

4. second order condition (SOC) [f is convex]

f is convex iff  $\text{dom}(f)$  is a convex set and  $\begin{cases} f''(x) \geq 0 & \text{for all } x \in \text{dom}(f) \\ \text{concave} & f''(x) \leq 0 \end{cases}$

\* let  $g(\theta) = f(y + \theta e)$  [y 起始点. e 单位向量.  $\theta \in \mathbb{R}$ ]

choose any  $y/e/0$ . prove  $g''(\theta=0) \geq 0$  for any y and e.

[e.g.] prove  $f(x) = \sum_i a_i (x - c_i)^2$  is a convex function in x with  $a_i > 0$

$$g(\theta) = f(y + \theta e) = \sum_i a_i (y_i + \theta e_i - c_i)^2$$

$$g''(\theta) = \sum_i a_i e_i^2 \geq 0 \Rightarrow f \text{ is a convex function } (a_i > 0)$$

[ concave up convex ]

[ concave down concave ]

5. the epigraph of a function 凸數上圖

$$\text{epi } C = \{(x, y) : y \geq f(x), x \in \Omega\}$$

$\text{DC} = \{(x, y) : y \geq f(x), x \in \Omega\}$  is a convex set if  $f(x)$  is a convex function.

new point  $\lambda(x_1, y_1) + (1-\lambda)(x_2, y_2) = (\lambda x_1 + (1-\lambda)x_2, \lambda y_1 + (1-\lambda)y_2)$

$$\Rightarrow \lambda y_1 + (1-\lambda)y_2 \geq \lambda f(x_1) + (1-\lambda)f(x_2) \geq f(\lambda x_1 + (1-\lambda)x_2)$$

$\Rightarrow C$  is a convex set

②  $f(x)$  is a convex function:  $f(C) = \{(x, y) : y \geq f(x), x \in \Omega\}$  is a convex set

$$x_1 \in \Omega, x_2 \in \Omega \quad \text{then } (x_1, f(x_1)) \in C \quad (x_2, f(x_2)) \in C$$

凸集中任兩點連線  $(\lambda x_1 + (1-\lambda)x_2, \lambda f(x_1) + (1-\lambda)f(x_2)) \in C$   
仍屬於凸集

$$\Rightarrow \lambda f(x_1) + (1-\lambda)f(x_2) \geq f(\lambda x_1 + (1-\lambda)x_2)$$

$\Rightarrow f$  is a convex function

6. Suppose  $f$  is differentiable and convex. If  $x^*$  is feasible such that  $\nabla f(x^*) = 0$ , then  $x^*$  is the global minimizer.

7. convex function can be discontinuous at the boundary [但需高于单侧极限]

8.  $f(x) = h(g(x))$   $\left\{ \begin{array}{l} h \text{ convex \& non-decreasing, } g \text{ convex} \Rightarrow f \text{ convex} \\ h \text{ convex \& non-increasing, } g \text{ concave} \Rightarrow f \text{ convex} \\ h \text{ concave \& non-decreasing, } g \text{ concave} \Rightarrow f \text{ concave} \\ h \text{ concave \& non-increasing, } g \text{ convex} \Rightarrow \text{concave} \end{array} \right.$

# 9. exercises.

(1.) ① S is a convex set

<1> choose  $x, y \rightarrow \lambda x + (1-\lambda)y$  is in S?

<2>  $S_1, S_2$  convex  $\rightarrow S_1 \cap S_2$  convex

<3>  $f(x)$  is a convex function  $\left\{ \begin{array}{l} \{x : f(x) \leq y, x \in \Omega\} \text{ for any } y \\ \{x, y : y \geq f(x), x \in \Omega\} \end{array} \right.$  are convex sets

②  $f$  is a convex function

<4> by definition  $f(\lambda x_1 + (1-\lambda)x_2) \leq \lambda f(x_1) + (1-\lambda)f(x_2)$   
 | domain  $\Omega$  is a convex set &  $\frac{d^2 f(x+\theta e)}{d\theta^2} \geq 0 \quad \forall \lambda+0e \in \Omega$

(2) composition

[e.g.]  $g(x) = f(Ax+b) - c$  is convex if  $f$  is convex

$$\begin{aligned} g(\lambda x + (1-\lambda)y) &= f[A(\lambda x + (1-\lambda)y) + b] - c \\ &= f[\lambda(Ax+b) + (1-\lambda)(Ay+b)] - c \\ &\leq \lambda f(Ax+b) + (1-\lambda)f(Ay+b) - c \\ &= \lambda(g(x) + (1-\lambda)g(y)) \end{aligned}$$

(3) elementwise maximum

$f_1, f_2, \dots, f_n$  convex  $\Rightarrow f(x) = \max\{f_1(x), f_2(x), \dots, f_n(x)\}$  convex

proof:

pick any  $x, y \in \text{dom}(f)$ ,  $\lambda \in [0, 1]$ . then

$$\begin{aligned} f(\lambda x + (1-\lambda)y) &= f_j(\lambda x + (1-\lambda)y) \quad (\text{for some } j \in \{1, 2, \dots, n\}) \\ &\leq \lambda f_j(x) + (1-\lambda)f_j(y) \\ &\leq \lambda \max\{f_1(x), \dots, f_n(x)\} + (1-\lambda) \max\{f_1(y), \dots, f_n(y)\} \\ &= \lambda f(x) + (1-\lambda)f(y) \end{aligned}$$

· (4) nonnegative weighted sums

$f_1, f_2, \dots, f_n$  convex,  $w_i \geq 0 \Rightarrow f = \sum_{i=1}^n w_i f_i$  convex

$$\forall x, y, 0 \leq \lambda \leq 1 \quad f_m(\lambda x + (1-\lambda)y) \leq \lambda f_m(x) + (1-\lambda)f_m(y)$$

$$\sum_{i=1}^n f_i(\lambda x + (1-\lambda)y) w_i = f \leq \sum_{i=1}^n w_i f_i \quad [g(x) = \int_A w(y) f_y(x) dy \text{ convex}]$$

(5) show  $S = \{(x, y, z) : z \geq x^2 + y^2\} \subset \mathbb{R}^3$  is convex [同知课点 5]

$$a_1 = (\lambda x_1, \lambda y_1, \lambda z_1), a_2 = ((1-\lambda)x_2, (1-\lambda)y_2, (1-\lambda)z_2)$$

$$\Rightarrow \text{show } [\lambda x_1 + (1-\lambda)x_2]^2 + [\lambda y_1 + (1-\lambda)y_2]^2 \leq \lambda z_1 + (1-\lambda)z_2 \quad ①$$

$$[\lambda x_1 + (1-\lambda)x_2]^2 \leq \lambda x_1^2 + (1-\lambda)x_2^2 \quad (x^2 \text{ is convex})$$

$$① \Rightarrow \lambda x_1^2 + (1-\lambda)x_2^2 + \lambda y_1^2 + (1-\lambda)y_2^2$$

$$= \lambda(x_1^2 + y_1^2) + (1-\lambda)(x_2^2 + y_2^2) \leq \lambda z_1 + (1-\lambda)z_2$$

(6) Show that if  $S_1$  and  $S_2$  are convex sets in  $\mathbb{R}^{m \times n}$ , then so is their partial sum

$$S = \{(x, y_1 + y_2) \mid x \in \mathbb{R}^m, y_1, y_2 \in \mathbb{R}^n, (x, y_1) \in S_1, (x, y_2) \in S_2\}$$

Proof:

$$\text{Consider } (\bar{x}, \bar{y}_1 + \bar{y}_2), (\tilde{x}, \tilde{y}_1 + \tilde{y}_2) \in S.$$

$$\text{for } 0 \leq \theta \leq 1$$

$$\theta(\bar{x}, \bar{y}_1 + \bar{y}_2) + (1-\theta)(\tilde{x}, \tilde{y}_1 + \tilde{y}_2)$$

$$= (\theta \bar{x} + (1-\theta)\tilde{x}, (\theta \bar{y}_1 + (1-\theta)\tilde{y}_1) + (\theta \bar{y}_2 + (1-\theta)\tilde{y}_2))$$

since  $S_1, S_2$  is convex.  $(\bar{x}, \bar{y}_1) \in S_1, (\bar{x}, \bar{y}_2) \in S_2, (\tilde{x}, \tilde{y}_1) \in S_1, (\tilde{x}, \tilde{y}_2) \in S_2$

$$\Rightarrow (\theta \bar{x} + (1-\theta)\tilde{x}, \theta \bar{y}_1 + (1-\theta)\tilde{y}_1) \in S_1, (\theta \bar{x} + (1-\theta)\tilde{x}, \theta \bar{y}_2 + (1-\theta)\tilde{y}_2) \in S_2$$

$$\Rightarrow (\theta \bar{x} + (1-\theta)\tilde{x}, (\theta \bar{y}_1 + (1-\theta)\tilde{y}_1) + (\theta \bar{y}_2 + (1-\theta)\tilde{y}_2)) \in S$$

(7)  $C$  is a convex set  $x_1, x_2, x_3 \in C, \theta_1 + \theta_2 + \theta_3 = 1 (\theta_i > 0)$  show  $\theta_1 x_1 + \theta_2 x_2 + \theta_3 x_3 \in C$

Proof:

$$\theta_1 x_1 + (1-\theta_1) y \Rightarrow \theta_1 x_1 + (\theta_2 + \theta_3) \left( \frac{\theta_2}{\theta_2 + \theta_3} x_2 + \frac{\theta_3}{\theta_2 + \theta_3} x_3 \right)$$

$$\frac{\theta_2}{\theta_2 + \theta_3} + \frac{\theta_3}{\theta_2 + \theta_3} = 1 \Rightarrow y \in C \Rightarrow \theta_1 x_1 + (1-\theta_1) y = \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_3 \in C$$

$$\Rightarrow \theta_i > 0, \sum_i \theta_i = 1 \text{ then } \sum_i \theta_i x_i \in C$$

$$\Rightarrow f(\theta_1 x_1 + \theta_2 x_2 + \theta_3 x_3) = f(\theta_1 x_1 + (1-\theta_1) y) \leq \theta_1 f(x_1) + (1-\theta_1) f(y)$$

$$\leq \theta_1 f(x_1) + \theta_2 f(x_2) + \theta_3 f(x_3) \quad [f(x) \text{ is convex in } x \text{ for the real space}]$$

$$\Rightarrow \theta_i > 0, \sum_i \theta_i = 1 \quad f(\sum_i \theta_i x_i) \leq \sum_i \theta_i f(x_i)$$

$$\Rightarrow X \text{ be the random variable } f(E(X)) \leq E(f(X))$$

(8) Show  $f(x) = \sum_{i=1}^r \alpha_i x_{[i]}$  is a convex function of  $x$ .  $\alpha_1 \geq \alpha_2 \geq \dots \geq \alpha_r \geq 0$

$x_{[i]}$  denotes the  $i$ -th largest component of  $x$ .

$$f(x) = \sum_{i=1}^r \alpha_i x_{[i]} = (\alpha_1 - \alpha_2) x_{[1]} + (\alpha_2 - \alpha_3) (x_{[1]} + x_{[2]}) + \dots + (\alpha_{r-1} - \alpha_r) (x_{[1]} + \dots + x_{[r-1]}) + \alpha_r (x_{[1]} + \dots + x_{[r]})$$

Since  $\alpha_1 \geq \alpha_2 \geq \dots \geq \alpha_r$ ,  $\sum_{i=1}^r f_{[i]}$  is convex

by nonnegative weighted sums, we can conclude  $f(x)$  is convex

### (9) 报童问题 [The NewsVendor's Problem]

c 成本 P 销售价 s 回收价 D 需求 [PDF  $h(x)$ ] Q { } 累积分布  
 $(P > c > s)$ ,  $(Q-D)^+ = \max\{Q-D, 0\}$

$$\int_0^\infty h(x) = 1$$

$$\pi(Q) = p \mathbb{E}[\min(Q, D)] + s \mathbb{E}[(Q-D)^+] - cQ$$

$$= p \int_0^\infty \min(Q, x) h(x) dx + s \mathbb{E}[(Q-D)^+] - cQ$$

$$= p \left( \int_0^Q x h(x) dx + \int_Q^\infty Q h(x) dx \right) + s \int_0^Q (Q-x) h(x) dx - cQ$$

$$\frac{\partial \pi}{\partial Q} = p \left[ Q h(Q) + \int_Q^\infty h(x) dx - Q h(Q) \right] + s \left( 1 - \int_Q^\infty h(x) dx \right) - c$$

$$= (p-s) \int_Q^\infty h(x) dx + s - c$$

$$\frac{\partial \pi}{\partial Q} = (Q h(Q))$$

$$\pi'(Q) = 0 \Rightarrow \int_Q^\infty h(x) dx = \frac{c-s}{p-s}$$

$$\frac{\partial \pi}{\partial Q} = \frac{Q \int_Q^\infty h(x) dx}{\partial Q}$$

$$\Rightarrow \int_0^{Q^*} h(x) dx = 1 - \frac{c-s}{p-s} = \frac{p-c}{p-s}$$

$$\begin{aligned} &= \int_Q^\infty h(x) dx - Q \frac{\int_Q^\infty h(x) dx}{\partial Q} \\ &= \int_Q^\infty h(x) dx - Q [h(Q) - h(\infty)] \\ &= \int_Q^\infty h(x) dx - Q h(Q) \end{aligned}$$

$$X \sim N(\mu, \sigma^2), \int_0^{Q^*} h(x) dx = A$$

转换为均值为0，方差为1的标准正态分布。

$$Z = \frac{Q^* - \mu}{\sigma}, \phi(Z) = A \text{ (查表得 Z)}$$

$$\frac{\partial \pi}{\partial Q} = \frac{Q \int_Q^\infty h(x) dx - \int_Q^\infty x h(x) dx}{\partial Q}$$

$$\begin{aligned} &= Q h(Q) + \int_Q^\infty h(x) dx - Q h(Q) \\ &= \int_Q^\infty h(x) dx = 1 - \int_Q^\infty h(x) dx \end{aligned}$$

$$\Rightarrow Q^* = Z\sigma + \mu$$

$$\pi''(Q) = -(p-s) h(Q) = (s-p) h(Q) \leq 0 \Rightarrow \text{concave function}$$

• (16 points) The NewsVendor's Problem: MLE + Convex Optimization

A newspaper vendor faces the task of determining the optimal quantity of newspapers to stock, given uncertain demand  $D$ . The cost structure for ordering newspapers is as follows: the first twelve copies incur a cost of \$8 per copy, while any additional copies beyond twelve incur a cost of \$5 per copy. Each sold copy yields a revenue of \$10 for the vendor, while each unsold copy has a salvage value of \$2. Remark: for ease of analysis, we allow the vendor to order a non-integer number of copies.

- (a) (4 points) Given that the demand  $D$  follows a uniform distribution with a support of  $[\theta, 2\theta]$ , where  $\theta$  is unknown, the objective is to estimate the maximum likelihood estimator (MLE) of  $\theta$ . Five samples of  $D$  are provided: {12, 15, 13, 20, 14}. Calculate the MLE of  $\theta$ . Explain the result in detail.
- (b) (8 points) Assuming that  $D$  follows a uniform distribution with a support of  $[\theta, 2\theta]$ , where  $\theta$  is the MLE obtained in (a), calculate the vendor's optimal order decision. Explain the result in detail.
- (c) (4 points) For (b), what's the vendor's optimal order decision, assuming that the vendor can only order an integer number of copies? Explain the result in detail.

$$(a) \max \left\{ \frac{x_i}{2} \right\} = 10, \min \left\{ x_i \right\} = 12 \quad L(\theta) = \left( \frac{1}{\theta} \right)^5 \Rightarrow \hat{\theta} = 10$$

(b) order decision  $Q$ ,

$$\textcircled{1} Q \leq 12 \quad C = 8Q \quad \textcircled{2} Q \geq 12 \quad C = 12 \cdot 8 + 5(Q-12) = 5Q + 36$$

$$W(Q) = 10 \min(D, Q) + 2 \max(0, Q-D) - C$$

$$\textcircled{1} Q > 12 \quad W(Q) = 10 \min(D, Q) + 2 \max\{0, Q-D\} - 5Q - 36$$

$$D \sim U(10, 20)$$

$$\begin{aligned} E[W(Q)] &= \frac{1}{10} \left[ \int_{10}^Q (10D + 2(Q-D) - 5Q - 36) dD + \int_Q^{20} (10Q - 5Q - 36) dD \right] \\ &= \frac{1}{10} \left[ \int_{10}^Q (8D - 3Q - 36) dD + \int_Q^{20} (5Q - 36) dD \right] \\ &= \frac{1}{10} (-4Q^2 + 10Q - 760) \quad Q = 16.25 \Rightarrow E[W(Q)] \text{ maximum} \end{aligned}$$

$$\textcircled{2} Q \leq 12 \Rightarrow Q \leq 10 \quad E[W(Q)] = 2Q, \quad Q \in [10, 12] \quad E(W(Q)) = -\frac{1}{5}Q^2 + 10Q - 40$$

$$E(W(Q)) \leq E(W(12))$$

$$(c) Q = 16 \Rightarrow 29.6 \quad Q = 17 \Rightarrow 29.4$$

$$\Rightarrow Q = 16$$

## • Machine Learning

### 1. K-Nearest Neighbors (KNN. K近邻算法) Non-parametric

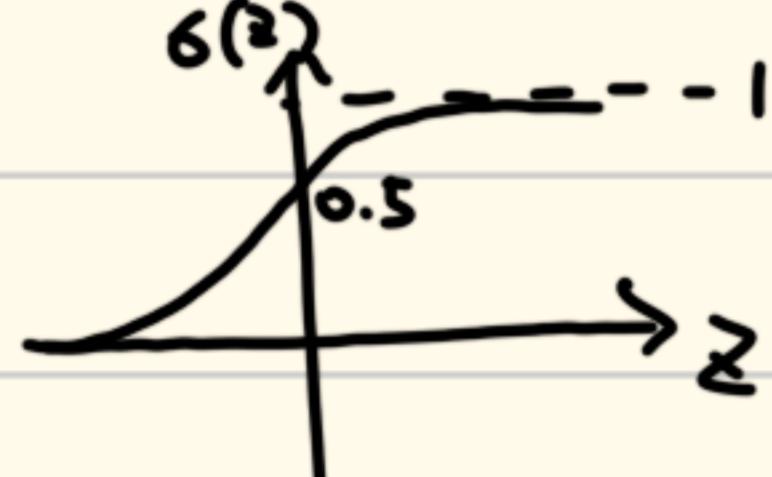
一个新点与其最近的  $K$  个点产生关联

record the entry (disi, labeli) sort by disi. pick the first  $k$  entries.

choose the label with the largest frequency

### 2. logistic regression parametric

$$P_{w,b}(C_1|x) = \sigma(z) = \frac{1}{1+e^{-z}} \quad (z = w^T x + b. P_{w,b}(C_1|x) \geq 0.5 \rightarrow \text{output } C_1, \text{ otherwise } C_2)$$



$$\begin{aligned} & y \in \{0, 1\} \quad \max_{\theta, b} L(\theta, b) = \log \prod_{k=1}^m P(y^k | x^k, \theta, b) = \sum_{k=1}^m \log P(y^k | x^k, \theta, b) \\ & (x_k, y_k), k \in \{1, 2, \dots, m\} \end{aligned}$$

$$\begin{aligned} \log P(y^k | x^k, \theta, b) &= \log \left[ \sigma\left(\sum_i \theta_i x_i^k + b\right)^{y^k} \cdot \left[1 - \sigma\left(\sum_i \theta_i x_i^k + b\right)\right]^{1-y^k} \right] \\ &= y^k \log \sigma\left(\sum_i \theta_i x_i^k + b\right) + (1-y^k) \log \left[1 - \sigma\left(\sum_i \theta_i x_i^k + b\right)\right] \\ &= (y^k - 1) \left(\sum_i \theta_i x_i^k + b\right) - \boxed{\log [1 + \exp(-\sum_i \theta_i x_i^k - b)]} \quad \begin{aligned} & y^k \log (1+e^{-z}) + (1-y^k) \log \frac{e^{-z}}{1+e^{-z}} \\ & = -y^k \log (1+e^{-z}) + (1-y^k) [-z - \log(1+e^{-z})] \\ & = (y^k - 1)z - \log (1+e^{-z}) \end{aligned} \end{aligned}$$

for  $\boxed{\square}$ .  $\nabla L(\theta, b)$  as a direction

$$h(t) = \log [1 + \exp(-\sum_i \theta_i x_i^k - b + t(-\sum_i e_i x_i^k - c))]$$

$$= \log [1 + \exp(A + tB)] = \log [1 + e^{z(t)}] \quad (z(t) = A + tB)$$

$$h''(t) = \frac{B^2}{[1 + e^{z(t)}]^2} e^{z(t)} \geq 0$$

$\Rightarrow h(t)$  is convex in  $(\theta, b)$

$\Rightarrow L(\theta, b)$  is concave in  $(\theta, b)$

☆优化算法因凸性而  
可被是收敛

### 3. Gradient Descent 梯度下降

$$x^{(t+1)} = x^{(t)} - \alpha \nabla f(x^{(t)})$$

$$[\min_x f(x)]$$

( $\alpha$  学习率)

$\max_x f(x)$  需要变号

→ 沿增长方向相反的方向偏移

stopping criteria by  $\epsilon$

$$|x^{(t+1)} - x^{(t)}| \leq \epsilon \text{ or } |f'(x^{(t)})| \leq \epsilon$$

## 4. Stochastic gradient descent [随机梯度下降]

for logical regression

$$\left\{ \begin{array}{l} \theta^{(t+1)} \leftarrow \theta^{(t)} + \alpha^{(t)} \frac{1}{|B|} \sum_{k \in B} \left\{ (y^{k-1}) x^k + \frac{\exp(-\theta^{(t)} x^k - b^{(t)})}{1 + \exp(-\theta^{(t)} x^k - b^{(t)})} x^k \right\} \\ b^{(t+1)} \leftarrow b^{(t)} + \alpha^{(t)} \frac{1}{|B|} \sum_{k \in B} \left\{ (y^{k-1}) + \frac{\exp(-\theta^{(t)} x^k - b^{(t)})}{1 + \exp(-\theta^{(t)} x^k - b^{(t)})} \right\} \end{array} \right.$$

### (10 points) Gradient Descent Method

Let's maximize the function  $f(x) = -x^2/2 + x$  through the gradient descent method. We fix the learning rate in the algorithm as a constant,  $\alpha$ . Assume that the initial value of  $x$  is  $x^{(0)} = 2$  and  $x^{(t)}$  is the value of  $x$  after the  $t$ -th iteration of gradient descent. The stopping rule is  $|x^{(t+1)} - x^{(t)}| \leq \epsilon$ , where  $\epsilon$  is a non-negative constant.

- (a) Let  $\epsilon = 1/1000$ . In each of the following situations will the final output be the optimal solution? Explain the reasons.
- (1 point)  $\alpha = 0$ .
  - (1 point)  $\alpha = 1$ .
  - (1 point)  $\alpha = 2$ .
  - (2 points)  $\alpha = 1/2$ .

- (b) (5 points) Let  $\epsilon = 0$ . For what range of  $\alpha$  will  $\lim_{t \rightarrow \infty} x^{(t)}$  converge to the optimal solution? Explain the reasons.

$$\nabla f(x) = -x + 1 \quad | \text{maximize}$$

$$(a) x^{(t+1)} \leftarrow x^{(t)} + \alpha \nabla f(x^{(t)}) = x^{(t)} + \alpha (-x^{(t)} + 1)$$

$$x^{(t)} - 1 = (1-\alpha)(x^{(t-1)} - 1) \Rightarrow x^{(t)} = (1-\alpha)^t + 1 \quad \text{the answer should be } 1$$

①  $\alpha=0 \quad x^{(t)}=2 \quad \text{No}$

②  $\alpha=1 \quad x^{(t)}=1 \quad \text{Yes}$

③  $\alpha=2 \quad x^{(t)} = -t + 1 \quad \text{No}$

④  $\alpha=\frac{1}{2} \quad x^{(t)} = (\frac{1}{2})^t + 1 \rightarrow 1 \quad \text{Yes}$

$$(b) (1-\alpha)^t + 1 \rightarrow 1 \Rightarrow (1-\alpha)^t \rightarrow 0 \quad \alpha \in (0, 1]$$

## 5. k-Means (unsupervised machine learning 未监督式)

① initialize  $k$  cluster centers  $\{c^1, c^2, \dots, c^k\}$  randomly

② 分配每个样本到最近的 cluster assignment

$$\pi(i) = \arg \min_{j=1 \dots k} \|x^i - c^j\|^2$$

③ Center adjustment  $c^j = \frac{1}{|\{i : \pi(i)=j\}|} \sum_{i \in \pi(j)} x^i$

\*直至中心点稳定

[12 points] Suppose we would like to use the K-means algorithm and L2-norm distance (Euclidian distance) to cluster the 8 data points given in Figure 3 below into  $K = 3$  clusters. The L2-norm distance between points  $\mathbf{x} = (x_1, x_2)$  and  $\mathbf{y} = (y_1, y_2)$  is  $d(\mathbf{x}, \mathbf{y}) = \sqrt{(x_1 - y_1)^2 + (x_2 - y_2)^2}$ . The coordinates of the data points are:

$$\begin{array}{llll} x^1 = (2, 8) & x^2 = (2, 5) & x^3 = (1, 2) & x^4 = (5, 8) \\ x^5 = (7, 3) & x^6 = (6, 4) & x^7 = (8, 4) & x^8 = (4, 7) \end{array}$$

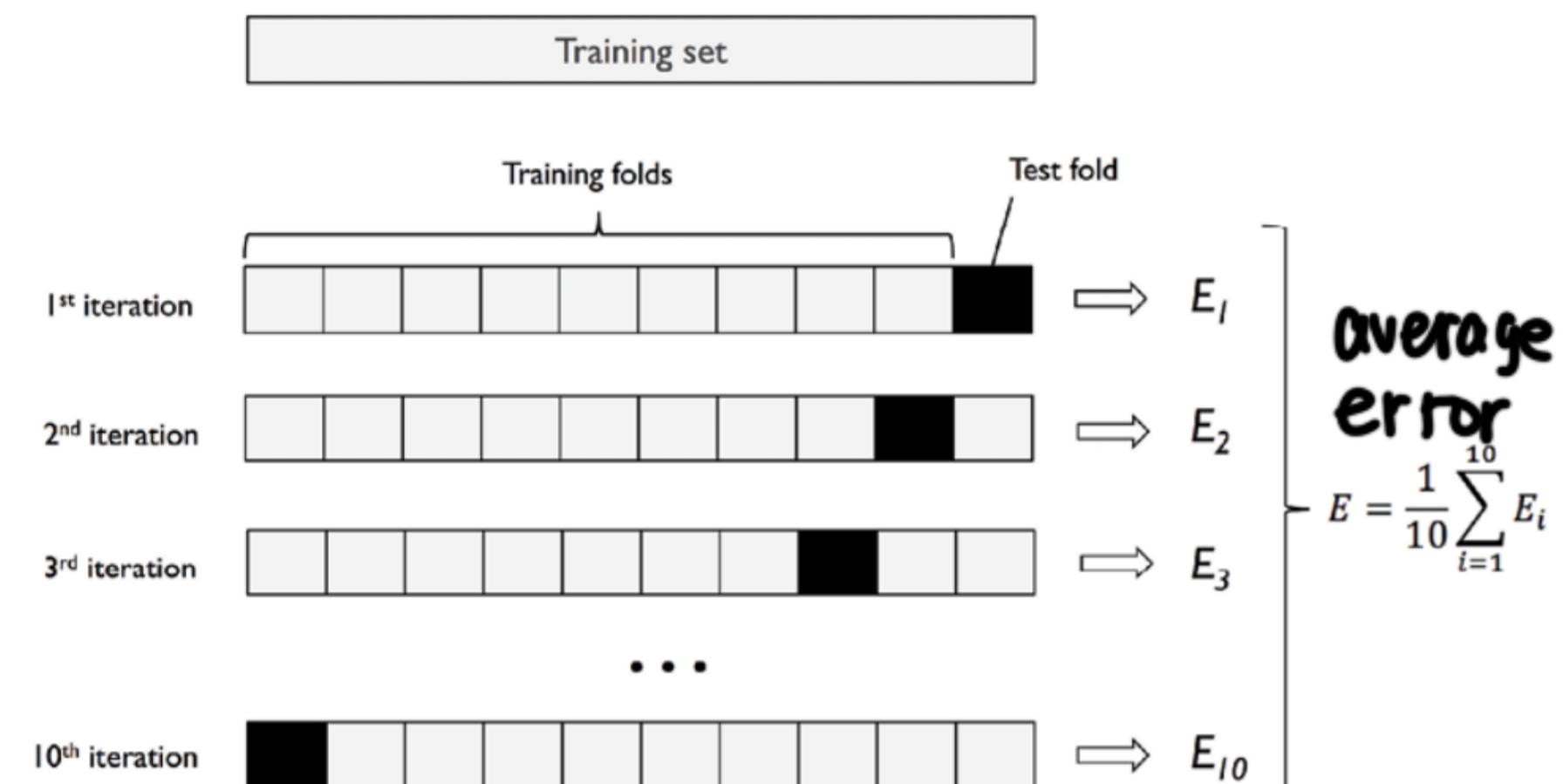
choose  $x_1, x_3, x_5$  as center

1:	$x_1, x_2, x_4, x_8$	(3.25, 7)
2:	$x_3$	(1, 2)
3:	$x_5, x_6, x_7$	(7, 3.67)

iteration → stop after the centers don't change

## 6. Cross-validation 反复验证

k-fold  $(k-1)$  for training, 1 for validation



Model 1 is better iff.

$$\frac{1}{k} \sum_i Err_h(\theta; \cdot) < \frac{1}{k} \sum_i Err_g(\gamma; \cdot)$$

1  $h(\theta, x)$       2  $g(\gamma, x)$

for supervised learning  
 (sol.) CROSS-validation for

Small summary



Algorithm	parametric	supervised	underfitting	overfitting
KNN	non parametric	supervised	large $k$ (oversmooth decision boundaries)	small $k$ (noise-sensitive & complex boundaries)
k-Means	non parametric	unsupervised	small $k$ (fail to capture clusters)	large $k$ (meaningless tiny clusters)
Logic Regression	parametric	supervised	poor features & overly simple model	too many features & weak regression