

Problem 1: Exponential Family

(1.)

Remark poisson distribution $f(k;\lambda) = P(X=k) = \frac{\lambda^k e^{-\lambda}}{k!}$

Sufficient statistic carry all the information needed to be inference

$$p(x|\theta) = h(x) \exp(\langle \theta, \phi(x) \rangle - A(\theta))$$

Consider

$$\text{Poisson}(x;\lambda) = \frac{\lambda^x}{x!} e^{-\lambda}$$

$$= e^{x \log \lambda} e^{-\lambda} \cdot \frac{1}{x!}$$

$$\text{let } h(x) = \frac{1}{x!}, A(\theta) = \lambda$$

$$\text{then } p(x|\theta) \propto \exp(\theta^t \phi(x))$$

$$= \exp(\log \lambda \cdot x)$$

$$\text{let } \theta^t = \log \lambda \text{ then } \phi(x) = x$$

thus $\phi(x) = [x]$ is the sufficient statistic for poisson dist

(2.)

$$p(x|\lambda) = \frac{\lambda^x}{x!} e^{-\lambda}$$

$$= \frac{1}{x!} \exp(x \log \lambda - \lambda)$$

$$h(x) = \frac{1}{x!}$$

$$\theta = \log \lambda$$

$$\phi(x) = x$$

$$A(\theta) = \lambda = e^\theta$$

by book P214 $p(x|\theta) = h(x) \exp(\langle \theta, \phi(x) \rangle - A(\theta))$

conjugate prior $p(\theta|\tau) = h_c(\theta) \exp(\langle [\begin{smallmatrix} \tau_1 \\ \tau_2 \end{smallmatrix}], [\begin{smallmatrix} \theta \\ -A(\theta) \end{smallmatrix}] \rangle - A_c(\tau))$

canonical form

$$p(x|\alpha, \beta) = h_c(\alpha) \exp(\tau_1 \log \alpha - \tau_2 \lambda - A_c(\tau))$$

$$(*) \quad p(\lambda | \alpha, \beta) \propto h_c(\lambda) \lambda^{\tau_1} e^{-\tau_2 \lambda}$$

note this close to gamma: $f(x; \alpha, \beta) = \frac{\beta^\alpha x^{\alpha-1} e^{-\beta x}}{P(\alpha)}$

let $\exp(-A(\lambda))$ account for $\frac{\beta^\alpha}{P(\alpha)}$

then let $h_c(\lambda) = 1$

$$\left[\begin{matrix} \tau_1 \\ \tau_2 \end{matrix} \right] = \left[\begin{matrix} \alpha-1 \\ \beta \end{matrix} \right]$$

$$\text{then } (*) \propto \lambda^{\tau_1} e^{-\tau_2 \lambda} = \lambda^{\alpha-1} e^{-\beta x}$$

conjugate prior is gamma distribution

Problem 2

consider $L(\theta) = \prod_{i=1}^N P(y_i | x_i; \theta) = \prod_{i=1}^N \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(y_i - ax_i + b)^2}{2\sigma^2}}$

$$\log(L(\theta)) = \ell(\theta) = \log \prod_{i=1}^N \frac{1}{\sqrt{2\pi\sigma^2}} \cdot \log \prod_{i=1}^N e^{-\frac{1}{2\sigma^2}(y_i - ax_i + b)^2}$$

$$= -\frac{1}{2}N \log 2\pi\sigma^2 - \frac{1}{2\sigma^2} \sum_{i=1}^N (y_i - (ax_i + b))^2$$

①

$$\begin{aligned} 0 &= \frac{\partial \ell(\theta)}{\partial a} = -\frac{1}{2\sigma^2} \sum_{i=1}^N 2(y_i - ax_i - b) \cdot (-x_i) \\ &= -\frac{1}{\sigma^2} \left[\sum_{i=1}^N (y_i - b)x_i - a \sum_{i=1}^N x_i^2 \right] \end{aligned}$$

thus $\sum_{i=1}^N x_i^2 \cdot a^{ml} + \sum_{i=1}^N x_i \cdot b^{ml} = \sum_{i=1}^N x_i y_i \quad (\star)$

$$0 \stackrel{\text{want}}{=} \frac{\partial \ell(\theta)}{\partial b} = -\frac{1}{2\sigma^2} \sum_{i=1}^N 2(y_i - ax_i - b)$$

thus $\sum_{i=1}^N x_i \cdot a^{ml} + N b^{ml} = \sum_{i=1}^N y_i \quad (\star\star)$

$$N \cdot (\star) - \sum_{i=1}^N x_i \cdot (\star\star) \Rightarrow$$

$$a^{ml} (N \sum_{i=1}^N x_i^2 - (\sum_{i=1}^N x_i)^2) = N \sum_{i=1}^N x_i y_i - \sum_{i=1}^N x_i \cdot \sum_{i=1}^N y_i$$

$$a^{ml} = \frac{N \sum_{i=1}^N x_i y_i - \sum_{i=1}^N x_i \sum_{i=1}^N y_i}{N \sum_{i=1}^N x_i^2 - (\sum_{i=1}^N x_i)^2}$$

by $(\star\star)$

$$\begin{aligned} b^{ml} &= \left(\sum_{i=1}^N y_i - \sum_{i=1}^N x_i \cdot a^{ml} \right) \frac{1}{N} \\ &= \left(\sum_{i=1}^N y_i - \sum_{i=1}^N x_i \cdot \frac{N \sum_{i=1}^N x_i y_i - \sum_{i=1}^N x_i \sum_{i=1}^N y_i}{N \sum_{i=1}^N x_i^2 - (\sum_{i=1}^N x_i)^2} \right) \cdot \frac{1}{N} \end{aligned}$$

$$0 \stackrel{\text{want}}{=} \frac{\partial \ell(\theta)}{\partial G} = -\frac{N}{2} \frac{1}{2\pi G^2} \cdot 4\pi G - 2 \cdot (-\frac{1}{2}) \frac{1}{G^3} \sum_{i=1}^N (y_i - (ax_i + b))^2$$

$$= -\frac{N}{G} + \frac{1}{G^3} \left(\sum_{i=1}^N (y_i - (ax_i + b)) \right)^2$$

thus $(G^{ml})^2 = \frac{1}{N} \left(\sum_{i=1}^N (y_i - (a^{ml}x_i + b^{ml})) \right)^2$

Problem 3

$D = \{(x_i, y_i)\}$ x d dimensional

$$MSE(\theta) = \frac{1}{N} \sum_{i=1}^N (y_i - \theta^t x_i)^2$$

$x_i: dx_1$
 $\theta: dx_1$

① let $f(x; \theta) = \theta^t x$

then want $\hat{\theta} = \arg \min_{\theta} MSE(\theta)$

note $\theta^t x_i = x_i \theta^t$

$$MSE(\theta) = \frac{1}{N} \sum_{i=1}^N (y_i - x_i^t \theta)^2$$

$$= \frac{1}{N} (y - X \theta)^t (y - X \theta)$$

where $X = \begin{pmatrix} x_1^t \\ \vdots \\ x_N^t \end{pmatrix}_{N \times d}$

consider $\frac{dMSE(\theta)}{d\theta} = \frac{1}{N} \frac{d}{d\theta} (y^t y - 2y^t X \theta + \theta^t X^t X \theta)$

$$= \frac{1}{N} (-y^t X + \theta^t X^t X) \stackrel{\text{want}}{=} 0$$

$$y = \begin{pmatrix} y_1 \\ \vdots \\ y_N \end{pmatrix}_{N \times 1}$$

thus $\theta^t (X^t X) = y^t X$

$$\theta^t = y^t X (X^t X)^{-1}$$

$$\theta = (X^t X)^{-1} X^t y$$

$$(X^t X) \theta = X^t y$$

let $A = X^t X$ is $(dxN)(N \times d) \Rightarrow d \times d$

$$X = \begin{pmatrix} x_1^t \\ \vdots \\ x_N^t \end{pmatrix}$$

let $b = X^t y$ is $(dxN)(N \times 1) \Rightarrow d \times 1$

Assume A be full rank and cannot take inverse

since otherwise A is singular then θ may has no solution if $b_j \neq 0$ for $j \in \{1, \dots, d\}$
all infinite solutions

when A is full rank

only soln $\hat{\theta} = A^{-1} b$

$$\hat{\theta} = A^{-1}b$$

A^{-1} : $d \times d$
 b : $d \times 1$ get (A) need $O(Nd^2)$

inverse matrix A need $O(d^3)$ get b need $O(dN)$
multiplication need $O(d \times 1)$

thus time complexity $O(Nd^2 + d^3)$

Note that if A is full rank we need to assume
 $N > d$ since otherwise let $X_i = \text{ones}(N, d)$
then A will be singular rank(A) not full

thus time complexity is $O(Nd^2)$

thus complexity time is

$$O(Nd^2) + O(d^3)$$

$\approx O(Nd^2)$ when A full rank

Problem 4

$$f(\underline{x}; \theta) = \theta^t \underline{x}$$

Gaussian regression model: $y | \underline{x} \sim N(\theta^t \underline{x}, \sigma^2)$ σ^2 known

independent prior $\theta_j \sim N(\theta_j; 0, s^2)$

$$(1) P(\theta | D) = \frac{P(D|\theta) P(\theta)}{\int P(D|\theta) P(\theta) d\theta} = \frac{P(Y|\underline{x}, \theta) P(\theta)}{\int P(Y|\underline{x}, \theta) P(\theta) d\theta}$$

$$(*) = \prod_{i=1}^N \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{1}{2\sigma^2}(y_i - \theta^t \underline{x}_i)^2} \cdot \prod_{j=1}^d \frac{1}{\sqrt{2\pi s^2}} e^{-\frac{1}{2s^2}(\theta_j)^2}$$

$$P(D) = \int (*) d\theta$$

$$P(\theta | D) = \frac{(*)}{P(D)}$$

$$\log P(\theta | D) = -\frac{N}{2} \log(2\pi\sigma^2) - \frac{1}{2\sigma^2} \sum_{i=1}^N (y_i - \theta^t \underline{x}_i)^2 - \frac{d}{2s^2} \sum_{j=1}^d \theta_j^2 - \log P(D)$$

$$(2) \log P(\theta | D) = -\frac{1}{2\sigma^2} \sum_{i=1}^N (y_i - \theta^t \underline{x}_i)^2 - \frac{1}{2s^2} \sum_{j=1}^d \theta_j^2 + \text{something not depends on } \theta$$

$$\text{Note } \theta^t \underline{x}_i = \underline{x}_i^t \theta$$

$$= -\frac{1}{2\sigma^2} \sum_{i=1}^N (y_i - \underline{x}_i^t \theta)^2 - \frac{1}{2s^2} \sum_{j=1}^d \theta_j^2 + \text{else}$$

\underline{x} and θ

$$= -\frac{1}{2\sigma^2} (\underline{y} - \underline{X}\theta)^t (\underline{y} - \underline{X}\theta) - \frac{1}{2s^2} \theta^t \theta + \text{else}$$

$$\text{where } \underline{y} = \begin{pmatrix} y_1 \\ \vdots \\ y_n \end{pmatrix}, \underline{X} = \begin{pmatrix} \underline{x}_1^t \\ \vdots \\ \underline{x}_n^t \end{pmatrix}_{n \times d}, \theta_{1 \times d}$$

$$\begin{aligned}
 \frac{\partial \ell(\theta)}{\partial \theta} &= -\frac{1}{2G^2} \frac{d}{d\theta} \left(\underline{y^t y} - 2\underline{y^t x\theta} + \underline{\theta^t x^t x\theta} \right) - \frac{1}{2S^2} \frac{d}{d\theta} (\theta^t \theta) \\
 &= -\frac{1}{G^2} (-2\underline{y^t x} + 2\underline{\theta^t x^t x}) - \frac{1}{S^2} 2\theta^t \\
 &= -\frac{1}{G^2} (-\underline{y^t x} + \underline{\theta^t x^t x}) - \frac{\theta^t}{S^2} \\
 &= \frac{\underline{y^t x}}{G^2} - \theta^t \left(\frac{\underline{x^t x}}{G^2} + \frac{1}{S^2} \right)
 \end{aligned}$$

③ $y: n \times 1$ $x: n \times d$ $\theta: d \times 1$ $O(x^t x) = O(d \times n \times d) = O(nd^2)$

$$O\left(\frac{\underline{y^t x}}{G^2}\right) = O(nd) + O(d)$$

$$O\left(\frac{\underline{\theta^t x^t x}}{G^2}\right) = O(d^2n) + O(nd) + O(d)$$

$$O(\theta^t / S^2) = O(d)$$

$$O(\text{gradient descent}) = O(MNd^2) \quad \text{proof next page}$$

$$O(\text{normal method}) = O(Nd^4) + O(d^3)$$

normal method has larger time complexity if $d > MN$

else gradient descent will have larger time complexity

double check for gradient method

$$\begin{aligned} & \sum_{i=1}^N (y_i - \theta^T x_i)^2 + \sum_{j=1}^d \theta_j^2 \\ (\star) = & \sum_{i=1}^N \left(y_i - \sum_{j=1}^d \theta_j x_{ij} \right)^2 + \sum_{j=1}^d \theta_j^2 \\ = & \sum_{i=1}^N \left[y_i^2 - 2 y_i \sum_{j=1}^d \theta_j x_{ij} + \sum_{k=1}^d \theta_k x_{ik} \sum_{j=1}^d \theta_j x_{ij} \right] + \underbrace{\sum_{j=1}^d \theta_j^2}_{\sum_{k=1}^d \sum_{j=1}^d \theta_k \theta_j x_{ik} x_{ij}} \end{aligned}$$

$$\frac{\partial C}{\partial \theta_h} = \sum_{i=1}^N \left[-2 y_i \theta_h x_{ih} + \sum_{j=1}^d \theta_j x_{ih} x_{ij} + \sum_{k=1}^d \theta_k x_{ik} x_{ij} \right]$$

$O(1)$ $O(d)$ $O(d)$

$$+ 2 \underbrace{\sum_{j=1}^d \theta_j}_{O(d)}$$

$$\frac{\partial C}{\partial \theta_n} = O(Nd) + O(d) \quad O(d)$$

$$\frac{\partial C}{\partial \theta} = O(Nd^2)$$

Need M steps $O(MNd^2)$

Problem 5

$$MSE(\theta) = \sum_{i=1}^N (y_i - f(x_i; \theta))^2$$

$$D = \{(x_i, y_i)\}$$

$$r(\theta) = \sum_{j=1}^p |\theta_j|$$

$$\text{minimize } MSE(\theta) + \lambda r(\theta)$$

$$\log p(\theta | X, y) = \log p(y | X, \theta) + \log p(\theta) + \text{const}$$

$$\text{new loss function } E(\theta) = \sum_{i=1}^N (y_i - f(x_i; \theta))^2 + \lambda \sum_{j=1}^p |\theta_j|$$

WLOG, say x_i is $d \times 1$

$$\text{let } y = \begin{pmatrix} y_1 \\ \vdots \\ y_n \end{pmatrix}, \quad x = \begin{pmatrix} x_1^t \\ \vdots \\ x_n^t \end{pmatrix}_{n \times d}, \quad \theta_{d \times 1} \quad \theta_{d \times 1}$$

$$\text{then } E(\theta) = (y - x\theta)^t (y - x\theta) + \lambda \sum_{j=1}^p |\theta_j|$$

want $\arg \min_{\theta} E(\theta)$

$$\text{equal to } \arg \max_{\theta} -\frac{1}{2\sigma^2} E(\theta)$$

Consider from ML of Least-Squares Regression of Gaussian in problem 2

$$L(\theta) = \prod_{i=1}^N \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(y_i - x_i^t \theta)^2}{2\sigma^2}}$$

can be consider as likelihood

$$\text{since } l'(\theta) \propto -\frac{1}{2\sigma^2} E(\theta)$$

$$\text{similarly } \prod_{i=1}^p \frac{\lambda \sigma(6)}{2} e^{-\lambda \sigma(6) |\theta_j|}$$

can be considered as prior

called double exponential distribution
(laplace dist.)

$$\text{since } l'(\theta) \propto -\lambda \sigma(6) \sum_{j=1}^p |\theta_j|$$

which is symmetric
a exponential dist by
y axis and rescale sum to 1

consider rescale $e^{-\frac{(y_i - x_i^t \theta)^2}{2\sigma^2}}$

let $\bar{y}_i = y_i/\sqrt{\sigma^2}$, $\bar{\theta} = \theta/\sqrt{\sigma^2}$

then $\Rightarrow e^{-\frac{(\bar{y}_i - x_i^t \bar{\theta})^2}{2}}$ with kernel mean $x_i^t \bar{\theta}$ and variance 1

use same parameters rescaled

$$\lambda^*(\sigma) e^{-\lambda^*(\sigma)|\theta|} \Rightarrow \frac{\lambda^{**}}{2\sigma^2} e^{-\lambda^{**}|\theta|/\sqrt{\sigma^2}}$$

\Downarrow force sum to 1

$$\Rightarrow \frac{\lambda^{**}}{2\sqrt{\sigma^2}} e^{-\lambda^{**}|\bar{\theta}|/\sqrt{\sigma^2}}$$

thus consider prior as $p(\bar{\theta}|\sigma^2) = \pi(\sigma^2) \prod_{j=1}^P \frac{\lambda^{**}}{2\sqrt{\sigma^2}} e^{-\lambda^{**}|\bar{\theta}|/\sqrt{\sigma^2}}$
where $\pi(\sigma^2)$ is hyper prior

check for unimodal

$$p(\bar{\theta}|\mathbf{Y}, \mathbf{X}, \sigma^2) = \prod_{i=1}^N \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(y_i - x_i^t \theta)^2}{2\sigma^2}} \prod_{j=1}^P \frac{\lambda^{**}}{2\sqrt{\sigma^2}} e^{-\lambda^{**}|\bar{\theta}|/\sqrt{\sigma^2}} \cdot \pi(\sigma^2)$$

change of variables let $\bar{\theta} = \theta/\sqrt{\sigma^2}$, $\bar{y} = y/\sqrt{\sigma^2}$

$$\ell(p(\cdot)) = \log(\pi(\sigma^2)) - \frac{n+p-1}{2} \ln(\sigma^2) - \frac{1}{2} \|\bar{y} - \mathbf{X}\bar{\theta}\|_2^2 - \lambda^{**} |\bar{\theta}|$$

$\ell(p(\cdot))$ is convex since \log , $\|\cdot\|$ are convex operator on $\bar{\theta}$ and σ^2

choose $\pi(\sigma^2)$ to be convex also

then $p(\bar{\theta}|\mathbf{Y}, \mathbf{X}, \sigma^2)$ is unimodal

$$\text{also } \ell(\mathbf{CPC}(\cdot)) \propto -\frac{1}{2} \|\bar{\mathbf{y}} - \mathbf{X}\bar{\boldsymbol{\theta}}\|_2^2 - \lambda^{**} |\bar{\boldsymbol{\theta}}|$$

$$\text{choose } \lambda = \frac{1}{2} \lambda^{**}$$

$$\text{then } -\ell(\mathbf{CPC}(\bar{\boldsymbol{\theta}} | \mathbf{Y}, \mathbf{X}, \mathbf{G})) \propto \|\bar{\mathbf{y}} - \mathbf{X}\bar{\boldsymbol{\theta}}\|_2^2 + \lambda |\bar{\boldsymbol{\theta}}|$$

thus lasso gaussian prior is

$$\pi(G^2) \prod_{j=1}^p \frac{\lambda}{2\sqrt{G^2}} e^{-\lambda |\theta_j|/\sqrt{G^2}}$$

is the double exponential distribution (laplace)

where $\pi(G^2)$ is the hyperparameter distribution for unimodal of posterior

Note that if we define MSE as

$$\frac{1}{N} \sum_{i=1}^N (y_i - f_i(x_i; \boldsymbol{\theta}))^2$$

$$\text{then let } \lambda = \frac{\lambda}{N}$$

the lasso gaussian prior will be same as above

Problem 6

$$\textcircled{1} \quad Y_i | \underline{x}_i, \theta \sim \text{Pois}(\lambda(\underline{x})) = \text{Pois}(\theta^t \underline{x}) = \text{Pois}(\underline{x}^t \theta)$$

$$P(Y_i | \lambda) = \frac{\lambda^{y_i} e^{-\lambda}}{y_i!}$$

recall exponential family from problem 1

$$\begin{aligned} P(X | \lambda) &= \frac{\lambda^x}{x!} e^{-\lambda} \\ &= \frac{1}{x!} \exp(x \log \lambda - \lambda) \end{aligned}$$

$$h(x) = \frac{1}{x!}$$

$$\theta = \log \lambda$$

$$\phi(x) = x$$

$$A(\theta) = \lambda = e^\theta$$

$$P(Y | \lambda) = \frac{1}{y!} \exp(y \log \lambda - \lambda)$$

$$h(y) = \frac{1}{y!}$$

$$A(\theta) = \lambda = e^{\theta^t \underline{x}}$$

$$\text{thus } P(Y_i | \underline{x}_i, \theta) = \frac{1}{y_i!} \exp(y_i(\theta^t \underline{x}_i) - e^{\theta^t \underline{x}_i})$$

$$P(D | \theta) = \prod_{i=1}^n \frac{1}{y_i!} e^{y_i(\theta^t \underline{x}_i) - e^{\theta^t \underline{x}_i}}$$

$$\ell(\theta) = \sum_{i=1}^n (y_i(\theta^t \underline{x}_i) - e^{\theta^t \underline{x}_i} - \log(y_i!))$$

$$= \sum_{i=1}^n (y_i(\underline{x}_i^t \theta) - e^{\underline{x}_i^t \theta} - \log(y_i!))$$

$$\theta^t \underline{x}_i : 1 \times d \times d \times 1$$

$$\underline{x}_i^t \theta : 1 \times d \times d \times 1$$

$$\begin{aligned}
 \textcircled{2} \quad \frac{\partial \ell(\theta)}{\partial \theta} &= \sum_{i=1}^n \left(y_i \frac{\partial (\underline{x}_i^t \theta)}{\partial \theta} \right) - \frac{\partial e^{\underline{x}_i^t \theta}}{\partial \theta} - 0 \\
 &= \sum_{i=1}^n y_i \underline{x}_i^t - \underline{x}_i^t e^{\underline{\theta}^t \underline{x}_i} \\
 &= \sum_{i=1}^n \underline{x}_i^t (y_i - e^{\underline{\theta}^t \underline{x}_i})
 \end{aligned}$$

\underline{x}_i	$d \times 1$
\underline{x}_i^t	$1 \times d$
$e^{\underline{x}_i^t \theta}$	$1 \times d \times d = 1$
$y_i \underline{x}_i^t$	$1 \times 1 \times 1 \times d$
$\underline{x}_i^t e^{-\underline{x}_i^t \theta}$	$1 \times d \times$

Problem 7

$$\textcircled{1} \quad f(x_i; \theta) = \frac{1}{1 + e^{-(\theta^T x_i)}}$$

$$CE(\theta) = -\frac{1}{N} \sum_{i=1}^N y_i \underbrace{\log(f(x_i; \theta))}_{\text{part ①}} + (1-y_i) \underbrace{\log(1-f(x_i; \theta))}_{\text{part ②}}$$

$$\frac{d - \log(f(x_i; \theta))}{d\theta} = \frac{d \log(1 + e^{-\theta^T x_i})}{d\theta} = \frac{1}{1 + e^{-\theta^T x_i}} \cdot (-e^{-\theta^T x_i}) \cdot (-x_i)$$

$$= \frac{x_i e^{-\theta^T x_i}}{1 + e^{-\theta^T x_i}} = x_i (1 - f(x_i; \theta))$$

$$\frac{d - \log(1 - f(x_i; \theta))}{d\theta} = \frac{d \log\left(\frac{e^{-\theta^T x_i}}{1 + e^{-\theta^T x_i}}\right)^{-1}}{d\theta} = \frac{d \log(1 + e^{-\theta^T x_i}) - \log(e^{-\theta^T x_i})}{d\theta}$$

$$= x_i (1 - f(x_i; \theta)) + x_i \quad \text{first part same as above}$$

thus $\frac{dCE(\theta)}{d\theta} = \frac{1}{N} \sum_{i=1}^N y_i x_i (1 - f(x_i; \theta)) + (1 - y_i) [x_i (1 - f(x_i; \theta)) + x_i]$

(2)

Theorem 2.6.1: If the function f has a second derivative which is non-negative (positive) everywhere, then the function is convex (strictly convex).

WTS hessian matrix is SPD (semi-positive definite)

hessian for part ①

$$\begin{aligned}\frac{\partial \underline{x}_i(1-f(\underline{x};\theta))}{\partial \theta} &= \frac{\partial -\underline{x}_i f(\underline{x};\theta)}{\partial \theta} \\ &= \frac{\partial \frac{-\underline{x}_i}{1+e^{-\theta^T \underline{x}}}}{\partial \theta} \\ &= \frac{\theta + \underline{x}_i^2 e^{-\theta^T \underline{x}_i}}{(1+e^{-\theta^T \underline{x}_i})^2} = -f(\underline{x}_i;\theta)(1-f(\underline{x}_i;\theta))\underline{x}_i \underline{x}_i^t\end{aligned}$$

let $\underline{z} \in \mathbb{R}_{d \times 1}$

$$\text{then } C(\underline{x}): \underline{z}^T \underbrace{f(\underline{x};\theta)(1-f(\underline{x};\theta)) \underline{x}_i \underline{x}_i^T}_{\text{since } f(\underline{x};\theta) \in (0,1]} \underline{z}$$

$$\text{then } f(\underline{x}_i;\theta)(1-f(\underline{x}_i;\theta)) \geq 0$$

$$\text{thus } C(\underline{x}) = \text{nonnegative number} \cdot (\underline{x}_i^T \underline{z})^2 \geq 0 \quad \text{SPD}$$

hessian for part ②

$$\frac{\partial \underline{x}_i(1-f(\underline{x};\theta)) + \underline{x}_i}{\partial \theta} \stackrel{\text{same as above}}{=} f(\underline{x}_i;\theta)(1-f(\underline{x}_i;\theta)) \underline{x}_i \underline{x}_i^t$$

Same as $C(\underline{x})$ part ② hessian is SPD

$$\text{Since } 0 \leq y^i \in \{0,1\} \quad 0 \leq 1-y^i \in \{0,1\}$$

$$CE(\theta) = -\frac{1}{N} \sum_{i=1}^N y_i \log(f(x_i, \theta)) + (1-y_i) \log(1-f(x_i, \theta))$$

Since if f convex func g convex func $\lambda_1 \geq 0$ $\lambda_2 \geq 0$

$$\text{then } (\lambda_1 f + \lambda_2 g)(dx + (1-d)y)$$

$$\stackrel{\text{convex}}{=} \lambda_1 f(dx + (1-d)y) + \lambda_2 g(dx + (1-d)y)$$

$$\leq \lambda_1 f dx + \lambda_1 f(1-d)y + \lambda_2 g dx + \lambda_2 g(1-d)y$$

$$= \alpha(\lambda_1 f + \lambda_2 g)x + (1-\alpha)(\lambda_1 f + \lambda_2 g)y$$

$\lambda_1 f + \lambda_2 g$ also convex

thus since $-\log f(x_i, \theta)$ and $-\log(1-f(x_i, \theta))$ convex
 $y_i, 1-y_i \geq 0$

thus $CE(\theta)$ is convex