

1.

$$(a) J(\theta) = -\frac{1}{m} \sum_{i=1}^m y^{(i)} \log(h_{\theta}(x^{(i)})) + (1-y^{(i)}) \log(1-h_{\theta}(x^{(i)}))$$

$$h_{\theta}(x^{(i)}) = g(\theta^T x^{(i)}) . \quad g(z) = g(z)(1-g(z))$$

$$\begin{aligned}\frac{\partial J(\theta)}{\partial \theta_j} &= -\frac{1}{m} \sum_{i=1}^m y^{(i)} \frac{g(\theta^T x^{(i)})(1-g(\theta^T x^{(i)}))}{g(\theta^T x^{(i)})} x_j^{(i)} + (1-y^{(i)}) \frac{-g(\theta^T x^{(i)})(1-g(\theta^T x^{(i)}))}{1-g(\theta^T x^{(i)})} x_j^{(i)} \\ &= -\frac{1}{m} \sum_{i=1}^m [y^{(i)} - g(\theta^T x^{(i)})] x_j^{(i)} \\ \nabla_{\theta} J(\theta) &= \frac{1}{m} \underset{(n,m)}{X^T} \underset{(m,1)}{g(\theta X) - Y}\end{aligned}$$

$$X = \begin{bmatrix} (x^{(1)})^T \\ \vdots \\ (x^{(m)})^T \end{bmatrix}_{m \times n} \quad H_{jk} = \frac{\partial^2 J(\theta)}{\partial \theta_j \partial \theta_k} = \frac{1}{m} \sum_{i=1}^m g(\theta^T x^{(i)}) (1-g(\theta^T x^{(i)})) x_j^{(i)} x_k^{(i)}$$

$$H = \frac{1}{m} \underset{(n,n)}{X^T} \underset{(n,m)}{g(X\theta)} \underset{(m,n)}{(1-g(X\theta))X}$$

element-wise?

$$\begin{aligned}Z^T H Z &= \frac{1}{m} \sum_{k=1}^n \sum_{j=1}^m \sum_{l=1}^m g(\theta^T x^{(i)}) (1-g(\theta^T x^{(i)})) x_j^{(i)} x_k^{(i)} z_j z_k \\ &= \frac{1}{m} \sum_{i=1}^m g(\theta^T x^{(i)}) (1-g(\theta^T x^{(i)})) [x^{(i)T} Z]^2 \geq 0\end{aligned}$$

(b) Show GDA can be written as linear classifier

$$\begin{aligned}p(y=1|x; \phi, \mu_0, \mu_1, \Sigma) &= \frac{p(x|y=1; \mu_0, \mu_1, \Sigma) \cdot p(y=1; \phi)}{p(x|y=1; \mu_0, \mu_1, \Sigma) p(y=1; \phi) + p(x|y=0; \mu_0, \mu_1, \Sigma) p(y=0; \phi)} \\ &= \frac{1}{1 + \frac{p(x|y=0; \mu_0, \mu_1, \Sigma) p(y=0; \phi)}{p(x|y=1; \mu_0, \mu_1, \Sigma) p(y=1; \phi)}} \\ &= \frac{1}{1 + \exp(-\frac{1}{2} (x - \mu_0)^T \Sigma^{-1} (x - \mu_0) + \frac{1}{2} (x - \mu_1)^T \Sigma^{-1} (x - \mu_1)) \frac{1-\phi}{\phi}} \\ &= \frac{1 + \exp(-\frac{1}{2} (\mu_1^T - \mu_0^T) \Sigma^{-1} x + (\frac{1}{2} \mu_0^T \Sigma^{-1} \mu_0 - \frac{1}{2} \mu_1^T \Sigma^{-1} \mu_1 - \log \frac{1-\phi}{\phi})}{1 + \exp(-\frac{1}{2} (\mu_1^T - \mu_0^T) \Sigma^{-1} x + (\frac{1}{2} \mu_0^T \Sigma^{-1} \mu_0 - \frac{1}{2} \mu_1^T \Sigma^{-1} \mu_1 - \log \frac{1-\phi}{\phi})} \\ &\therefore \theta^T = (\mu_1^T - \mu_0^T) \Sigma^{-1}, \quad \theta = \Sigma^{-1}(\mu_1 - \mu_0) \\ \theta_0 &= \frac{1}{2} \mu_0^T \Sigma^{-1} \mu_0 - \frac{1}{2} \mu_1^T \Sigma^{-1} \mu_1 - \log \frac{1-\phi}{\phi} \\ &= \frac{1}{2} (\mu_0 + \mu_1)^T \Sigma^{-1} (\mu_0 - \mu_1) - \log \frac{1-\phi}{\phi}\end{aligned}$$

$$(d) \ell(\phi, \mu_0, \mu_1, \Sigma) = \log \prod_{i=1}^m p(x^{(i)}, y^{(i)}; \phi, \mu_0, \mu_1, \Sigma)$$

$$= \log \prod_{i=1}^m p(x^{(i)} | y^{(i)}; \mu_0, \mu_1, \Sigma) p(y^{(i)}; \phi)$$

Maximizing ℓ . = maximum likelihood estimate.

$$\ell = \sum_{i=1}^m \log p(x^{(i)} | y^{(i)}; \mu_0, \mu_1, \Sigma) + \sum_{i=1}^m \log p(y^{(i)}; \phi)$$

$$= \sum_{i=1}^m \log \frac{1}{(2\pi)^{\frac{m}{2}} |\Sigma|^{\frac{1}{2}}} \exp(-\frac{1}{2}(x^{(i)} - \mu_{y^{(i)}})^T \Sigma^{-1} (x^{(i)} - \mu_{y^{(i)}})) + \sum_{i=1}^m \log \phi^{\mathbb{I}\{y^{(i)}=1\}} (1-\phi)^{\mathbb{I}\{y^{(i)}=0\}}$$

$$= -\frac{mn}{2} \log 2\pi - \frac{m}{2} \log |\Sigma| - \frac{1}{2} \sum_{i=1}^m (x^{(i)} - \mu_{y^{(i)}})^T \Sigma^{-1} (x^{(i)} - \mu_{y^{(i)}}) + \sum_{i=1}^m \mathbb{I}\{y^{(i)}=1\} \log \phi + \sum_{i=1}^m \mathbb{I}\{y^{(i)}=0\} \log (1-\phi)$$

$$\textcircled{1} \quad \frac{\partial \ell}{\partial \phi} = \frac{1}{\phi} \sum_{i=1}^m \mathbb{I}\{y^{(i)}=1\} + \frac{1}{1-\phi} (-\sum_{i=1}^m \mathbb{I}\{y^{(i)}=0\})$$

$$\text{set } \frac{\partial \ell}{\partial \phi} = 0, \quad \phi = \frac{1}{m} \sum_{i=1}^m \mathbb{I}\{y^{(i)}=1\}$$

$$\textcircled{2} \quad \frac{\partial \ell}{\partial \mu_0} = \frac{\partial \ell}{\partial \mu_{y^{(i)}}} \cdot \frac{\partial \mu_{y^{(i)}}}{\partial \mu_0} = m \sum_{i=1}^m (x^{(i)} - \mu_{y^{(i)}}) \cdot \mathbb{I}\{y^{(i)}=0\}$$

$$\frac{\partial \ell}{\partial \mu_1} = \frac{\partial \ell}{\partial \mu_{y^{(i)}}} \cdot \frac{\partial \mu_{y^{(i)}}}{\partial \mu_1} = m \sum_{i=1}^m (x^{(i)} - \mu_{y^{(i)}}) \cdot \mathbb{I}\{y^{(i)}=1\}.$$

$$\mu_0 = \frac{\sum_{i=1}^m \mathbb{I}\{y^{(i)}=0\} x^{(i)}}{\sum_{i=1}^m \mathbb{I}\{y^{(i)}=0\}}$$

$$\mu_1 = \frac{\sum_{i=1}^m \mathbb{I}\{y^{(i)}=1\} x^{(i)}}{\sum_{i=1}^m \mathbb{I}\{y^{(i)}=1\}}.$$

$$\frac{\partial a^T X b}{\partial X} = ab^T$$

$$\textcircled{3} \quad \frac{\partial \ell}{\partial \Sigma} = -\frac{m}{2} \Sigma^{-1} + \frac{1}{2} \sum_{i=1}^m (x^{(i)} - \mu_{y^{(i)}})(x^{(i)} - \mu_{y^{(i)}})^T \Sigma^{-2}$$

$$\therefore \Sigma = \underbrace{\frac{1}{m} \sum_{i=1}^m (x^{(i)} - \mu_{y^{(i)}})(x^{(i)} - \mu_{y^{(i)}})^T}_{\sim}$$

(ii) take logarithm of x_2 . in Dataset 1

or squared root.

BOX-COX transformation.

2.

(a) Given $t^{(i)}, y^{(i)}$ and $x^{(i)}$ are conditionally independent.

$$P(y^{(i)}=1 | t^{(i)}=1, x^{(i)}) = P(y^{(i)}=1 | t^{(i)})$$

labeled
positive
data
truth
value
being
positive

which says labeled data are randomly selected from the set of positive examples. (All $x^{(i)}$ are positive samples, so $x^{(i)}$ is irrelevant.)

→ 给定数据真实值是1，数据有标记的假象与 $x^{(i)}$ 无关。

$$\text{prove: } P(t^{(i)}=1 | x^{(i)}) = P(y^{(i)}=1 | x^{(i)}) / \alpha.$$

$$\begin{aligned} P(y^{(i)}=1 | x^{(i)}) &= \sum_t P(y^{(i)}=1, t^{(i)}=1 | x^{(i)}) \\ &= P(y^{(i)}=1, t^{(i)}=1 | x^{(i)}) + P(y^{(i)}=1, t^{(i)}=0 | x^{(i)}) \end{aligned}$$

$$\begin{aligned} P(A|B, C)P(C|B) &= P(A, B|C) = P(y^{(i)}=1, t^{(i)}=1 | x^{(i)}) \\ &= P(y^{(i)}=1 | t^{(i)}=1, x^{(i)}) P(t^{(i)}=1 | x^{(i)}) \\ &= P(y^{(i)}=1 | t^{(i)}=1) P(t^{(i)}=1 | x^{(i)}) \\ \therefore \alpha &= P(y^{(i)}=1 | t^{(i)}=1) \end{aligned}$$

$$(b) h(x^{(i)}) \approx P(y^{(i)}=1 | x^{(i)}) = P(y^{(i)}=1 | t^{(i)}=1) P(t^{(i)}=1 | x^{(i)}) \approx P(y^{(i)}=1 | t^{(i)}=1) = \alpha$$

so $h(x^{(i)}) \approx \alpha$ for all $x^{(i)} \in V$

$$\frac{1}{1 + \exp(-\theta^T x)} = \frac{1}{2}$$

What is corrected theta?

$$\frac{2}{2} - 1 = \exp(-\theta^T x)$$

$$\theta^T x + \log\left(\frac{2}{2} - 1\right) = 0$$

$$\therefore \theta' = \theta + \frac{\text{add to } \theta}{\log\left(\frac{2}{2} - 1\right)} [1, 0, 0]$$

3. Poisson Regression.

$$(a) p(y; \lambda) = \frac{e^{-\lambda} \lambda^y}{y!} = \frac{1}{y!} \exp(-\lambda + y \ln \lambda)$$

$$p(y; \eta) = b(y) \exp(\eta^\top T(y) - a(\eta))$$

$$\therefore b(y) = \frac{1}{y!}$$

$$\eta = \ln \lambda,$$

$$a(\eta) = e^\eta$$

$$T(y) = y$$

(b) canonical response function: $g(\eta) = E[T(y); \eta]$

$$g(\eta) = E[y; \eta] = \lambda = e^\eta.$$

(c) In order to construct GLM.

i) $y|x; \theta \sim \text{Exponential Family}(n)$

ii) $h(x) = E[y|x]$

iii) $\eta = \theta^\top x$

$$\begin{aligned} \log p(y^{(i)}|x^{(i)}; \theta) &= \log b(y^{(i)}) \exp(\eta^\top T(y^{(i)}) - a(\eta^{(i)})) \\ &= \log \left(\frac{1}{y^{(i)}!} \right) \exp(\theta^\top x^{(i)})^T y^{(i)} - e^{\theta^\top x^{(i)}} \\ &= -\log y^{(i)!} + (x^{(i)\top} \theta)^{(i)} - e^{\theta^\top x^{(i)}} \end{aligned}$$

$$h(x) = e^{\theta^\top x}$$

$$\frac{\partial h(x)}{\partial \theta} = x^{(i)\top} y^{(i)} - x^{(i)\top} e^{\theta^\top x^{(i)}}$$

$$= x^{(i)\top} (y^{(i)} - h(x^{(i)}))$$

$$\therefore \theta := \theta + \alpha (\eta^{(i)} - h(x^{(i)})) x^{(i)} \quad \text{Stochastic Gradient Ascent Rule}$$

4. Convexity of Generalized Linear Models.

$$\begin{aligned}
 (a) \quad \frac{\partial}{\partial \eta} p(y; \eta) &= \frac{\partial}{\partial \eta} b(y) \exp(\eta y - a(\eta)) \\
 &= b(y) \exp(\eta y - a(\eta)) (y - \frac{\partial}{\partial \eta} a(\eta)) \\
 &= p(y; \eta) (y - \frac{\partial}{\partial \eta} a(\eta)) \\
 &= y p(y; \eta) - p(y; \eta) \frac{\partial}{\partial \eta} a(\eta)
 \end{aligned}$$

$$\begin{aligned}
 E[Y|X; \theta] &= \int y p(y; \eta) dy \\
 &= \int \frac{\partial}{\partial \eta} p(y; \eta) + \frac{\partial}{\partial \eta} a(\eta) \cdot p(y; \eta) dy \\
 &= \frac{\partial}{\partial \eta} \int p(y; \eta) dy + \frac{\partial}{\partial \eta} a(\eta) \cdot \int p(y; \eta) dy \\
 &= \frac{\partial}{\partial \eta} a(\eta)
 \end{aligned}$$

$$\begin{aligned}
 (b) \quad \text{Var}[Y|X; \theta] &= \text{Var}[Y; \eta] \\
 &= \int (\eta - \frac{\partial}{\partial \eta} a(\eta))^2 p(y; \eta) dy \\
 &= \int (\eta^2 - 2\eta \frac{\partial}{\partial \eta} a(\eta) + \frac{\partial^2}{\partial \eta^2} a(\eta)) p(y; \eta) dy \\
 &= \int \eta^2 p(y; \eta) dy - \int 2\eta \frac{\partial}{\partial \eta} a(\eta) p(y; \eta) dy + \int \frac{\partial^2}{\partial \eta^2} a(\eta) p(y; \eta) dy \\
 &= \int \eta^2 p(y; \eta) dy - 2 \frac{\partial}{\partial \eta} a(\eta) \int y p(y; \eta) dy + \frac{\partial^2}{\partial \eta^2} a(\eta) \int p(y; \eta) dy
 \end{aligned}$$

$$\begin{aligned}
 \frac{\partial}{\partial \eta} \int y p(y; \eta) dy &= \int y \frac{\partial}{\partial \eta} p(y; \eta) dy \\
 &= \int \eta^2 p(y; \eta) - y p(y; \eta) \frac{\partial}{\partial \eta} a(\eta) dy
 \end{aligned}$$

$$\therefore \int y^2 p(y; \eta) dy = \frac{\partial^2}{\partial \eta^2} a(\eta)$$

$$\therefore \text{Var}[Y|X; \theta] = \frac{\partial^2}{\partial \eta^2} a(\eta).$$

(c) loss function of one example

$$\begin{aligned}
 \text{NLL} &= -\log p(y; \eta) \\
 l(\theta) &= -\log b(y) \exp(\eta y - a(\eta)) \\
 &= -\log b(y) - \eta y + a(\eta) \\
 &= -\log b(y) - \theta^T x y + a(\theta^T x)
 \end{aligned}$$

$$\frac{\partial l(\theta)}{\partial \theta_j} = -x_j y + \frac{\partial a(\theta^T x)}{\partial \theta_j} \cdot x_j$$

$$\nabla_{\theta} l(\theta) = x \frac{\partial a(\theta^T x)}{\partial \theta} - x y$$

$$H_{jk} = \frac{\partial^2 l(\theta)}{\partial \theta_j \partial \theta_k} = \frac{\partial^2}{\partial \theta_j \partial \theta_k} a(\theta^T x) \cdot x_j x_k$$

$$H = x^T x \frac{\partial^2}{\partial \theta^2} a(\theta^T x)$$

$$\begin{aligned}
 \text{Then, prove PSD} \\
 \forall z \in \mathbb{R}^n, z^T H z &= z^T x^T x \frac{\partial^2}{\partial \theta^2} a(\theta^T x) z \\
 &= z^T x^T x z \frac{\partial^2}{\partial \theta^2} a(\theta^T x) \\
 &= (x^T z)^2 \frac{\partial^2}{\partial \theta^2} a(\theta^T x) \\
 &= (x^T z)^2 \text{Var}[Y|X; \theta] \\
 &\geq 0
 \end{aligned}$$

So H of GLM loss function is PSD.
The NLL loss of GLM is convex.

5. Locally Weighted Linear Regression.

(A) (i)

$$x\theta - y = \begin{bmatrix} x^{(1)}\theta - y^{(1)} \\ x^{(2)}\theta - y^{(2)} \\ \vdots \\ x^{(n)}\theta - y^{(n)} \end{bmatrix} \quad n \times 1$$

$$\pi(\theta) = (x\theta - y)^T w (x\theta - y)$$

$$w = m \times n, \quad w_{ij} = \begin{cases} 0, & i \neq j \\ \frac{1}{z} w^{(i)}, & i = j \end{cases}$$

$$(ii) \quad \frac{\partial}{\partial \theta_j} J(\theta) = \sum_{i=1}^m w^{(i)} (\theta^T x^{(i)} - y^{(i)}) x_j^{(i)}$$

$$\nabla_\theta J(\theta) = 2x^T W (x\theta - y) = 2x^T W x\theta - 2x^T W y$$

$(n, 1) \quad (n, m) (m, m) \quad (m, 1)$

$$\nabla_\theta J(\theta) = 0, \quad \theta = (x^T W x)^{-1} x^T W y$$

$$(iii) \quad l(\theta) = \log P(y | x; \theta)$$

$$= \sum_{i=1}^m \log p(y^{(i)} | x^{(i)}; \theta)$$

$$= \sum_{i=1}^m \log \frac{1}{\sqrt{2\pi\sigma^2}} - \frac{(y^{(i)} - \theta^T x^{(i)})^2}{2(\sigma^2)^2}$$

$$J(\theta) = \frac{1}{2} \sum_{i=1}^m w^{(i)} (\theta^T x^{(i)} - y^{(i)})^2$$

Maximizing $l(\theta)$ equals to minimizing $J(\theta)$

$$w^{(i)} = -\frac{1}{(\sigma^2)^2}$$