

## Problem Set 1



$$\begin{aligned}
 Q1 \quad (a) \quad \frac{\partial J(\theta)}{\partial \theta_j} &= \frac{\partial}{\partial \theta_j} \left( -\frac{1}{m} \sum_{i=1}^m y^{(i)} \log(h_\theta(x^{(i)})) + (1-y^{(i)}) \log(1-h_\theta(x^{(i)})) \right) \\
 &= -\frac{1}{m} \sum_{i=1}^m (y^{(i)} - h_\theta(x^{(i)})) x_j^{(i)}
 \end{aligned}$$

$$\begin{aligned}
 \frac{\partial^2 J(\theta)}{\partial \theta_i \partial \theta_j} &= -\frac{1}{m} \sum_{k=1}^m \frac{\partial}{\partial \theta_i} (y^{(k)} - g(h_\theta(x^{(k)})) x_j^{(k)}) \\
 &= -\frac{1}{m} \sum_{k=1}^m \frac{\partial}{\partial \theta_i} (y^{(k)} - g(\theta^T x^{(k)})) x_j^{(k)} \\
 &= -\frac{1}{m} \sum_{k=1}^m -g'(\theta^T x^{(k)}) \frac{\partial}{\partial \theta_i} (\theta^T x^{(k)}) x_j^{(k)} \\
 &= -\frac{1}{m} \sum_{k=1}^m -g(\theta^T x^{(k)}) (1-g(\theta^T x^{(k)})) x_i^{(k)} x_j^{(k)} \\
 &= \frac{1}{m} \sum_{k=1}^m g(\theta^T x^{(k)}) (1-g(\theta^T x^{(k)})) x_i^{(k)} x_j^{(k)}
 \end{aligned}$$

$$\begin{aligned}
 (Z^T H Z) &= \sum_i \sum_j z_i \frac{1}{m} \sum_{k=1}^m g(\theta^T x^{(k)}) (1-g(\theta^T x^{(k)})) x_i^{(k)} x_j^{(k)} z_j \\
 &= \frac{1}{m} \sum_{k=1}^m g(\theta^T x^{(k)}) (1-g(\theta^T x^{(k)})) \sum_i \sum_j z_i x_i^{(k)} x_j^{(k)} z_j \\
 &= \frac{1}{m} \sum_{k=1}^m g(\theta^T x^{(k)}) (1-g(\theta^T x^{(k)})) (x^T z)^2 \geq 0
 \end{aligned}$$

$\in \{0, 1\}$

$$x^T H z = 0 \quad \text{iff} \quad x^T z = 0$$

(b) Coding problem.

Requirements  $\rightarrow$  write Hessian  $H = \left[ \frac{\partial^2 J(\theta)}{\partial \theta_i \partial \theta_j} \right]_{\substack{i=1, 2, 3 \\ j=1, 2, 3}}$  in form of matrix multiplication

Define a diagonal matrix.

$$\text{diag} = [\text{diag}]_{i=1,2,\dots,k}$$

$$= \begin{bmatrix} g(\theta^T x^{(1)}) (1 - g(\theta^T x^{(1)})) & & \\ & \ddots & \\ & & g(\theta^T x^{(k)}) (1 - g(\theta^T x^{(k)})) \end{bmatrix}_{k \times k}$$

$$X = \begin{bmatrix} 1 & x_1^{(1)} & x_2^{(1)} \\ 1 & x_1^{(2)} & x_2^{(2)} \\ \vdots & \vdots & \vdots \\ 1 & x_1^{(m)} & x_2^{(m)} \end{bmatrix}^T$$

Thus  $H = X^T \text{diag } X$  (can use matmul)

Similarly,

$$J \in \mathbb{R} \quad \sum_{i=1}^m y^{(i)}$$

$$J = +\frac{1}{m} \cdot (X^T (\text{diag} - y))$$

$$\frac{1}{1+e^{-z}}$$

$$g(\theta^T x)$$

$$\frac{1}{k} \sum_{i=1}^k |\theta - \theta_i| < \epsilon$$

loss at  $i^{\text{th}}$  iteration.

$$\text{loss} = -\frac{1}{m} \left( \sum_{i=1}^m y^{(i)} \log [g(\theta^T x^{(i)})] + (1-y^{(i)}) \log [1-g(\theta^T x^{(i)})] \right)$$

Step 1:  $\theta^0$  is updated to  $\theta^{(1)}$   
Step 2:  $\theta^{(1)}$  is updated to  $\theta^{(2)}$  at  $\rightarrow$  loss  $-\frac{1}{m} \sum_{i=1}^m y^{(i)}$

Show that

$$(C). p(y=1|x; \phi, \mu_0, \mu_1, \Sigma) = \frac{1}{1 + \exp(-(\phi^T x + \phi_0))} \quad \begin{matrix} \text{since done separately} \\ \text{don't add} \end{matrix}$$

where  $\phi \in \mathbb{R}^n$  &  $\phi_0 \in \mathbb{R}$  are appropriate functions of  $\phi, \Sigma$ ,  $\mu_0$  and  $\mu_1$

$$p(y=1|x; \phi, \mu_0, \mu_1, \Sigma) = \frac{p(x|y=1) p(y=1)}{p(x|y=1) p(y=1) + p(x|y=0) p(y=0)}$$

$$= \frac{\phi}{(2\pi)^{n/2} |\Sigma|^{1/2}} \exp -\frac{1}{2} (x - \mu_1)^T \Sigma^{-1} (x - \mu_1)$$

$$\frac{\phi}{(2\pi)^{n/2} |\Sigma|^{1/2}} \exp -\frac{1}{2} (x - \mu_1)^T \Sigma^{-1} (x - \mu_1) + \frac{1-\phi}{(2\pi)^{n/2} |\Sigma|^{1/2}} \exp -\frac{1}{2} (x - \mu_0)^T \Sigma^{-1} (x - \mu_0)$$

$$= \frac{\phi \exp -\frac{1}{2} (x - \mu_1)^T \Sigma^{-1} (x - \mu_1)}{\{\phi \exp -\frac{1}{2} (x - \mu_1)^T \Sigma^{-1} (x - \mu_1)\} + (1-\phi) \exp -\frac{1}{2} (x - \mu_0)^T \Sigma^{-1} (x - \mu_0)}$$

$$= \frac{1}{1 + \frac{1-\phi}{\phi} \frac{\exp -\frac{1}{2} (x - \mu_0)^T \Sigma^{-1} (x - \mu_0)}{\exp -\frac{1}{2} (x - \mu_1)^T \Sigma^{-1} (x - \mu_1)}}$$

$$= \frac{1}{1 + \exp \left[ \ln \frac{1-\phi}{\phi} \right] \exp \left( -\frac{1}{2} (x - \mu_0)^T \Sigma^{-1} (x - \mu_0) - \exp \left( -\frac{1}{2} (x - \mu_1)^T \Sigma^{-1} (x - \mu_1) \right) \right)}$$

$$= \frac{1}{1 + \exp \left[ -\frac{1}{2} (x - \mu_0)^T \Sigma^{-1} (x - \mu_0) + \frac{1}{2} (x - \mu_1)^T \Sigma^{-1} (x - \mu_1) + \ln \frac{1-\phi}{\phi} \right]}$$

$$= \frac{1}{1 + \exp \left[ \frac{1}{2} \left( -x^T \Sigma^{-1} x + x^T \Sigma^{-1} \mu_0 + \mu_0^T \Sigma^{-1} x - \mu_0^T \Sigma^{-1} \mu_0 + x^T \Sigma^{-1} \mu_1 - \mu_1^T \Sigma^{-1} x + \mu_1^T \Sigma^{-1} \mu_1 \right) + \ln \frac{1-\phi}{\phi} \right]}$$

$$= \frac{1}{1 + \exp \frac{1}{2} \left( 2 \mu_0^T \Sigma^{-1} x - 2 \mu_1^T \Sigma^{-1} x + \mu_1^T \Sigma^{-1} \mu_1 - \mu_0^T \Sigma^{-1} \mu_0 \right) + \ln \frac{1-\phi}{\phi}}$$

$$= \frac{1}{1 + \exp \left( (\mu_0^\top \Sigma^{-1} x - \mu_1^\top \Sigma^{-1} x + \frac{\mu_1^\top \Sigma^{-1} (\mu_1 - \mu_0)^\top \mu_0 + \ln \frac{1-\phi}{\phi}}{2}) \right)}$$

$$\text{Thus } \nabla \Theta_0 = - \left[ (\mu_0 - \mu_1)^\top \Sigma^{-1} \right]^\top = -\Sigma^{-1}(\mu_0 - \mu_1) = \Sigma^{-1}(\mu_1 - \mu_0)$$

$$\Theta_0 = - \left( \frac{1}{2} \left( \mu_1^\top \Sigma^{-1} \mu_1 - \mu_0^\top \Sigma^{-1} \mu_0 \right) + \ln \frac{1-\phi}{\phi} \right).$$

$$\nabla \Theta = \frac{1}{2} \left( \mu_0^\top \Sigma^{-1} \mu_0 - \mu_1^\top \Sigma^{-1} \mu_1 \right) + \ln \frac{1-\phi}{\phi}$$

$$= \frac{1}{2} \left( (\mu_0 + \mu_1)^\top \Sigma^{-1} (\mu_0 - \mu_1) \right) - \ln \frac{1-\phi}{\phi}.$$

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

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

$$\ell(\phi, \mu_0, \mu_1, \Sigma) = \log \prod_{i=1}^m \phi^{\pm \mathbb{I}_{y^{(i)}=1}} \exp \frac{1}{2} (x - \mu_i)^\top \Sigma^{-1} (x - \mu_i)$$

$$p(x^{(i)} | y^{(i)}; \mu_0, \mu_1, \Sigma) p(y^{(i)}; \phi) = \begin{cases} \frac{\phi}{(2\pi)^{n/2} |\Sigma|^{1/2}} \exp \frac{-1}{2} (x - \mu_1)^\top \Sigma^{-1} (x - \mu_1) & \text{when } y^{(i)} = 1 \\ \frac{(1-\phi)}{(2\pi)^{n/2}} \exp \frac{-1}{2} (x - \mu_0)^\top \Sigma^{-1} (x - \mu_0) & \text{when } y^{(i)} = 0 \end{cases}$$

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

$$= \sum_{i=1}^m \frac{1}{p(y^{(i)}; \phi)} \frac{\partial}{\partial \phi} p(y^{(i)}; \phi)$$

$\begin{cases} 1 & \text{when } y^{(i)} = 1 \\ -1 & \text{when } y^{(i)} = 0 \end{cases}$

$$= \frac{m_1}{\phi} - \frac{m_2}{1-\phi}$$

where  $m_1 + m_2 = m$  and  $m_1$  are # of 1's &  $m_2$  are no. of 0's

$$\frac{m_1}{\phi} - \frac{m_2}{1-\phi} = 0 \Rightarrow \frac{m_1}{\phi} = \frac{m_2}{1-\phi} \Rightarrow \frac{1-\phi}{\phi} = \frac{m_2}{m_1}$$

$$\Rightarrow \frac{1}{\phi} - 1 = \frac{m_2}{m_1} \Rightarrow \frac{1}{\phi} = 1 + \frac{m_2}{m_1} = \frac{m_1 + m_2}{m_1}$$

$$\Rightarrow \hat{\phi}_1 = \frac{m_1}{m_1 + m_2} = \frac{m_1}{m} = \frac{1}{m} \sum_{i=1}^n \mathbb{1}\{y^{(i)} = 1\}.$$

$$\frac{\partial L(\phi, \mu_0, \mu_1, \xi)}{\partial \mu_0} = \sum_{i=1}^m \frac{\partial}{\partial \phi} \left[ \log p(x^{(i)} | y^{(i)}; \mu_0, \mu_1, \xi) + \log p(y^{(i)}; \phi) \right]$$

$$= \sum_{i=1}^m \frac{1}{p(x^{(i)} | y^{(i)}, \mu_0, \mu_1, \xi)} \underbrace{\frac{\partial}{\partial \phi} p(x^{(i)} | y^{(i)}; \mu_0, \mu_1, \xi)}_{||}$$

$$= \sum_{i=1}^m \underbrace{\begin{cases} \frac{\sqrt{2\pi}\sigma}{\exp -\frac{1}{2\sigma^2}(x^{(i)} - \mu_0)^2} & y^{(i)} = 1 \\ \frac{\sqrt{2\pi}\sigma}{\exp -\frac{1}{2\sigma^2}(x^{(i)} - \mu_1)^2} & y^{(i)} = 0 \end{cases}}_{||} \underbrace{\begin{cases} 0 & \text{when } y^{(i)} = 1 \\ \frac{+\exp(\mu_0)}{\sqrt{2\pi}\sigma} \exp -\frac{1}{2\sigma^2}(x^{(i)} - \mu_0)^2 & \text{when } y^{(i)} = 0. \end{cases}}_{||}$$

$$= \sum_{i=1}^m \mathbb{1}\{y^{(i)} = 0\} \frac{\sqrt{2\pi}\sigma}{\exp -\frac{1}{2\sigma^2}(x^{(i)} - \mu_0)^2} \times \frac{(x^{(i)} - \mu_0)}{\sqrt{2\pi}\sigma^3} \exp -\frac{1}{2\sigma^2}(x^{(i)} - \mu_0)^2$$

$$= \sum_{i=1}^m \mathbb{1}\{y^{(i)} = 0\} \frac{(x^{(i)} - \mu_0)}{\sigma^2} \quad \text{---}$$

$$\sum_{i=1}^m \mathbb{1}\{y^{(i)} = 0\} \frac{(x^{(i)} - \mu_0)}{\sigma^2} = 0$$

$$\Rightarrow \sum_{i=1}^m \mathbb{1}\{y^{(i)} = 0\} x^{(i)} - \sum_{i=1}^m \mathbb{1}\{y^{(i)} = 0\} \mu_0 = 0.$$

$$\Rightarrow \hat{\mu}_0 = \frac{\sum_{i=1}^m \mathbb{1}\{y^{(i)} = 0\} x^{(i)}}{\sum_{i=1}^m \mathbb{1}\{y^{(i)} = 0\}}$$

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

Similarly by symmetry,  $\hat{\mu}_1 = \frac{\sum_{i=1}^m \mathbb{1}\{y^{(i)} = 0\} x^{(i)}}{\sum_{i=1}^m \mathbb{1}\{y^{(i)} = 0\}}$

$$\frac{\partial \ell(\phi, \mu_0, \mu_1, \sigma^2)}{\partial \sigma} = \sum_{i=1}^m \frac{\partial}{\partial \sigma} \left[ \log p(x^{(i)} | y^{(i)}; \mu_0, \mu_1, \sigma) + \log p(y^{(i)}; \phi) \right]$$

$$= \sum_{i=1}^m \frac{\partial}{\partial \sigma} \left[ \log p(x^{(i)} | y^{(i)}; \mu_0, \mu_1, \sigma) + 0 \right]$$

$$= \sum_{i=1}^m \frac{1}{p(x^{(i)} | y^{(i)}; \mu_0, \mu_1, \sigma)} \frac{\partial}{\partial \sigma} p(x^{(i)} | y^{(i)}; \mu_0, \mu_1, \sigma)$$

$$= \sum_{i=1}^m \frac{(2\pi)^{-\frac{1}{2}} \sigma^{-\frac{1}{2}} e^{-\frac{1}{2} \frac{(x-\mu_0)^2}{\sigma^2}}}{\exp -\frac{1}{2} \frac{(x-\mu_0)^2}{\sigma^2}} \left[ \frac{1}{\sqrt{2\pi}\sigma} - \frac{(x-\mu_0)^2(-2)}{2\sigma^3} \exp -\frac{1}{2} \frac{(x-\mu_0)^2}{\sigma^2} + \frac{-1}{\sqrt{2\pi}\sigma^2} \exp -\frac{1}{2} \frac{(x-\mu_0)^2}{\sigma^2} \right]$$

$$+ \sum_{i=1}^m \left[ \frac{(2\pi)^{-\frac{1}{2}}}{\exp -\frac{1}{2} \frac{(x-\mu_1)^2}{\sigma^2}} + \frac{1}{2} \delta y^{(i)} = 1 \right] \dots$$

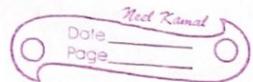
$$= \sum_{i=1}^m \left( \begin{array}{l} \left[ \frac{1}{2} \delta y^{(i)} = 0 \right] \left[ + \frac{(x-\mu_0)^2}{\sigma^3} - \frac{1}{\sigma} \right] \\ + \left[ \frac{1}{2} \delta y^{(i)} = 1 \right] \left[ \frac{(x-\mu_1)^2}{\sigma^3} - \frac{1}{\sigma} \right] \end{array} \right)$$

$$= \left( \frac{(x-\mu_0)^2}{\sigma^3} - \frac{1}{\sigma} \right) m_1 + \left( \frac{(x-\mu_1)^2}{\sigma^3} - \frac{1}{\sigma} \right) m_2$$

$$= \left[ \frac{m_1 (x-\mu_0)^2 + m_2 (x-\mu_1)^2}{\sigma^3} + \frac{m}{\sigma} \right] =$$

$$= \sum_{i=1}^m -\frac{1}{\sigma} + \sum_{i=1}^m \left( \begin{array}{l} \left[ \frac{1}{2} \delta y^{(i)} = 0 \right] \left[ \frac{(x-\mu_0)^2}{\sigma^3} + \right. \\ \left. \left[ \frac{1}{2} \delta y^{(i)} = 1 \right] \left[ \frac{(x-\mu_1)^2}{\sigma^3} \right] \right] \end{array} \right)$$

$$\frac{\partial l}{\partial \sigma}(ab)$$



Equating to 0.

$$\frac{1}{\sigma} \left( m_1 (x^{(i)} - \mu_0)^2 + m_2 (x^{(i)} - \mu_1)^2 - m \right) = 0.$$

$$\Rightarrow \frac{m_1 (x^{(i)} - \mu_0)^2 + m_2 (x^{(i)} - \mu_1)^2}{\sigma^2} = \hat{\sigma}^2$$

$$\Rightarrow \frac{1}{m} \sum_{i=1}^m (x^{(i)} - \mu y^{(i)}) (x^{(i)} - \mu y^{(i)})^T = \hat{\sigma}^2$$

Second derivative tests to show that indeed the critical pts are pts of maxima.

$$\frac{\partial^2 l(\phi, \mu_0, \mu_1, \varepsilon)}{\partial \phi^2} = -\frac{m_1}{\phi^2} - \frac{m_2}{(1-\phi)^2} < 0.$$

$$\begin{aligned} \frac{\partial^2 l(\phi, \mu_0, \mu_1, \varepsilon)}{\partial \mu_0^2} &= \frac{\partial}{\partial \mu_0} \sum_{i=1}^m \left[ \frac{1}{\sigma^2} y^{(i)} = 0 \right] \frac{(x^{(i)} - \mu_0)}{\sigma^2} \\ &= \sum_{i=1}^m \frac{1}{\sigma^2} y^{(i)} \left( -\frac{1}{\sigma^2} \right) < 0. \end{aligned}$$

$\frac{\partial^2 l}{\partial \mu_1^2}$  similarly by symmetry

$$\frac{\partial^2 l}{\partial \mu_1^2} < 0$$

$$\begin{aligned} \frac{\partial^2 l}{\partial \sigma^2} &= \frac{\partial}{\partial \sigma} \left( \frac{1}{\sigma^3} (m_1 (x - \mu_0)^2 + m_2 (x - \mu_1)^2) - \frac{m}{\sigma} \right) \\ &= -\frac{3}{\sigma^4} (m_1 (x - \mu_0)^2 + m_2 (x - \mu_1)^2) + \frac{m}{\sigma^2} \\ &= \frac{1}{\sigma^2} \left( m - \frac{3(m_1 (x - \mu_0)^2 + m_2 (x - \mu_1)^2)}{\sigma^2} \right) \end{aligned}$$

$$\begin{bmatrix} \mu_{y^{(1)}} & \mu_{y^{(2)}} & \mu_{y^{(3)}} \\ \vdots & \vdots & \vdots \\ \mu_{y^{(k)}} & \mu_{y^{(k)}} & \mu_{y^{(k)}} \end{bmatrix}$$

Note → In python of numpy.

Let  $m = np.array([1, 2, 3])$  then  $m.T$  is same as  
m. Both are  $(3,)$  ~~vectors~~. lists.

Thus to use transpose. first convert to ~~a~~ m, n form. If  $3 \times 1$  needed then use  $(:, 1)$

$m = np.array([1, 2, 3], ndims=2) \cdot T$

(f) run script .../linearclass/qf.py.  
 logreg theta estimate  $\rightarrow$  (1)  $[-2.4078467 \quad 1.0343113 \quad 0.24463071]$   
 (2)  $[-2.0967324 \quad 0.8663494 \quad 0.828187]$   
 gda theta estimate  $\rightarrow$  (1)  $[-2.10437509 \quad 1.59905448 \quad 0.01081916]$   
 $\rightarrow$  (2)  $[-2.11390101 \quad 0.85879252 \quad 0.84927843]$

(g) In ~~dataset~~ 1, logistic regression produces better results than GDA because the distribution is not gaussian while in 2 it is much nearer to gaussian distribution.

(h) Transform is box-cox transformation. Python scipy.stats has a boxcox function. But it can't handle but it requires positive values as inputs. I tried subtracting minimum value & then add a small value to ~~zero~~ non-zero entries but it did not work. Giving some errors.

(b) Box-Cox transformation.

Q2 (a)  $p(t^{(i)}=1 | x^{(i)}) = p(t^{(i)}=1, x^{(i)})$ .

$$= p(t^{(i)}=1, x^{(i)}, y^{(i)}=1) + p(t^{(i)}=1, x^{(i)}, y^{(i)}=0)$$

$$= p(t^{(i)}=1 | x^{(i)}, y^{(i)}=1) p(y^{(i)}=1 | x^{(i)}) + p(y^{(i)}=0 | x^{(i)})$$

$$= p(y^{(i)}=1 | x^{(i)}) + p(y^{(i)}=0 | t^{(i)}=1) p(t^{(i)}=1 | x^{(i)})$$

$$\Rightarrow p(t^{(i)}=1 | x^{(i)}) [1 - p(y^{(i)}=0 | t^{(i)}=1)] = p(y^{(i)}=1 | x^{(i)}).$$

$$\Rightarrow p(t^{(i)}=1 | x^{(i)}) = \frac{p(y^{(i)}=1 | x^{(i)})}{1 - p(y^{(i)}=0 | t^{(i)}=1)}$$

$$\Rightarrow p(t^{(i)}=1 | x^{(i)}) = \frac{p(y^{(i)}=1 | x^{(i)})}{(p(y^{(i)}=1 | t^{(i)}=1))^\alpha}$$

(b) ~~REASON~~  $V^+ = \{x^{(i)} \in V | y^{(i)}=1\}$ .

$$h(x^{(i)}) \approx p(y^{(i)}=1 | x^{(i)})$$

$$\Rightarrow p(y^{(i)}=1 | x^{(i)}) \approx h(x^{(i)})$$

$$\Rightarrow p(t^{(i)}=1 | x^{(i)}) \approx h(x^{(i)})$$

$$\text{for all } x^{(i)} \in V^+ \Rightarrow y^{(i)}=1 \Rightarrow t^{(i)}=1 \Rightarrow p(t^{(i)}=1 | x^{(i)})=1$$

$$\Rightarrow \alpha \approx h(x^{(i)}) \quad \forall x^{(i)} \in V^+$$

(c) ✓

(d) ✓

(e) ✓

Q3 (a)  $p(y; \lambda) = \frac{e^{-\lambda} \lambda^y}{y!}$

Exponential family

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

$$\frac{e^{-\lambda} \lambda^y}{y!} = \frac{e^{-\lambda} e^{y \log \lambda}}{y!}$$

$$= \frac{e^{-\lambda} e^{y \log \lambda}}{y!}$$

$$= \frac{y!}{y!} \exp(-\lambda + y \log \lambda)$$

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

$$\eta = \log \lambda$$

$$\Rightarrow \eta^n = \lambda$$

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

$$\eta = \log \lambda$$

$$T(y) = y$$

$$\alpha(\eta) = \lambda = e^\eta$$

(b)  $y | x; \theta \sim \text{Poisson}(\lambda)$

$$\lambda = h(x) = E[y|x; \theta]$$

$$= \lambda$$

$$= e^\eta$$

$$= e^{\theta T x}$$

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

$$\begin{aligned}
 c) \log p(y^{(i)} | x^{(i)}; \theta) &= \log \left[ \frac{1}{y^{(i)}!} \exp(n^T y^{(i)} - e^n) \right] \\
 &= \log \frac{1}{y^{(i)}!} + [n^T y^{(i)} - e^n] \\
 &= \log \frac{1}{y^{(i)}!} + [(O^T x^{(i)})^T y^{(i)} - e^{O^T x^{(i)}}]
 \end{aligned}$$

$$\begin{aligned}
 &\frac{\partial}{\partial \theta_j} \left[ \log \left( \frac{1}{y^{(i)}!} \right) + [(O^T x^{(i)})^T y^{(i)} - e^{O^T x^{(i)}}] \right] \\
 &= \frac{\partial}{\partial \theta_j} \left[ e^{(O_0 + O_1 x_1^{(i)} + \dots + O_d x_d^{(i)})^T y^{(i)}} - e^{(O_0 + O_1 x_1^{(i)} + \dots + O_d x_d^{(i)})} \right] \\
 &= x_j^{(i)} e^{-O^T x^{(i)}} y^{(i)} - x_j^{(i)} e^{O^T x^{(i)}} \\
 &= x_j^{(i)} y^{(i)} - x_j^{(i)} e^{O^T x^{(i)}} \\
 &= [y^{(i)} - e^{O^T x^{(i)}}] x_j^{(i)}
 \end{aligned}$$

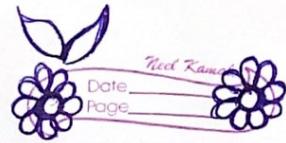
Thus stochastic gradient ascent rule becomes

$$\theta_j := \theta_j + \alpha \frac{\partial l(\theta)}{\partial \theta_j}$$

$\phi$  reduces to

$$\theta_j := \theta_j + \alpha (y^{(i)} - e^{O^T x^{(i)}}) g x_j^{(i)}$$

(d) ✓



## Q4 Convexity of GLM

$$p(y; \eta) = b(y) \exp(\eta y - \alpha(\eta))$$

$$\begin{aligned}
 (a) \quad \frac{\partial}{\partial \eta} p(y; \eta) &= \frac{\partial}{\partial \eta} [b(y) \exp(\eta y - \alpha(\eta))] \\
 &= b(y) \left[ y - \frac{\partial \alpha}{\partial \eta}(\eta) \right] \cancel{\exp(\eta y - \alpha(\eta))} \\
 &= p(y; \eta) \left[ y - \frac{\partial}{\partial \eta} \alpha(\eta) \right].
 \end{aligned}$$

$$\int \frac{\partial}{\partial \eta} p(y; \eta) dy = \int p(y; \eta) \left[ y - \frac{\partial}{\partial \eta} \alpha(\eta) \right] dy.$$

$$\Rightarrow \int \frac{\partial}{\partial \eta} p(y; \eta) dy = \int y p(y; \eta) dy - \int \frac{\partial}{\partial \eta} \alpha(\eta) p(y; \eta) dy.$$

$$\Rightarrow \int \frac{\partial}{\partial \eta} \int p(y; \eta) dy = E[y; \eta] - \frac{\partial}{\partial \eta} \alpha(\eta) \int p(y; \eta) dy$$

$$\Rightarrow E[y; \eta] = \frac{\partial}{\partial \eta} \alpha(\eta)$$

$$(b) \quad \text{Show that } \text{Var}(Y|x; \theta) = \frac{\partial^2}{\partial \eta^2} \alpha(\eta)$$

$$\frac{\partial}{\partial \eta} p(y; \eta) = p(y; \eta) \left[ y - \frac{\partial}{\partial \eta} \alpha(\eta) \right].$$

$$\frac{\partial^2}{\partial \eta^2} p(y; \eta) = \left[ \frac{\partial}{\partial \eta} p(y; \eta) \right] \left[ y - \frac{\partial}{\partial \eta} \alpha(\eta) \right] + p(y; \eta) \left[ -\frac{\partial^2}{\partial \eta^2} \alpha(\eta) \right]$$

$$\begin{aligned}
 \frac{\partial^2}{\partial \eta^2} p(y; \eta) &= p(y; \eta) \left[ y - \frac{\partial}{\partial \eta} \alpha(\eta) \right]^2 + \\
 &\quad p(y; \eta) \left[ -\frac{\partial^2}{\partial \eta^2} \alpha(\eta) \right]
 \end{aligned}$$

$$\frac{\partial^2}{\partial \eta^2} p(y; \eta) = y^2 p(y; \eta) + \left[ \frac{\partial \alpha(\eta)}{\partial \eta} \right]^2 p(y; \eta) \\ - 2y \frac{\partial}{\partial \eta} \alpha(\eta) p(y; \eta) + p(y; \eta) \left[ - \frac{\partial^2 \alpha(\eta)}{\partial \eta^2} \right]$$

$$\int \frac{\partial^2}{\partial \eta^2} p(y; \eta) dy = \int y^2 p(y; \eta) dy + \left[ \frac{\partial \alpha(\eta)}{\partial \eta} \right]^2 \int p(y; \eta) dy \\ - \frac{\partial \alpha(\eta)}{\partial \eta} \int 2y p(y; \eta) dy - \frac{\partial^2 \alpha(\eta)}{\partial \eta^2} \int p(y; \eta) dy.$$

$$\Rightarrow \int \frac{\partial^2}{\partial \eta^2} p(y; \eta) dy = \int y^2 p(y; \eta) dy + \left[ \frac{\partial \alpha(\eta)}{\partial \eta} \right]^2 - \\ \frac{\partial \alpha(\eta)}{\partial \eta} 2 \frac{\partial \alpha(\eta)}{\partial \eta} - \frac{\partial^2 \alpha(\eta)}{\partial \eta^2}$$

$$\Rightarrow 0 = \int y^2 p(y; \eta) dy - \left[ \frac{\partial \alpha(\eta)}{\partial \eta} \right]^2 - \frac{\partial^2 \alpha(\eta)}{\partial \eta^2}$$

$$\Rightarrow \left( \frac{\partial \alpha(\eta)}{\partial \eta} \right)^2 =$$

$$\Rightarrow \frac{\partial^2 \alpha(\eta)}{\partial \eta^2} = \int y^2 p(y; \eta) dy - [E(y; \eta)]^2$$

(c)  $L(\theta) = -\log - \text{log-likelihood.}$

$$= -\log [b(y^{(1)}) \exp(y^{(1)} - \alpha(\eta)) \dots b(y^{(n)}) \exp(y^{(n)} - \alpha(\eta))]$$

$$= -\log [b(y^{(1)}) \dots b(y^{(n)}) \exp(y^{(1)} - \alpha(\eta)) \dots \exp(y^{(n)} - \alpha(\eta))].$$

=

$$\nabla_{\theta}(\theta^T x) = \left[ \frac{\partial}{\partial \theta_j} \sum_{i=1}^{d+1} \theta_i x_{0,i} \right]_{j=1, 2, \dots, d+1}$$

$$= \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_{d+1} \end{bmatrix} = x.$$

$$\nabla_{\theta}(x^T \theta) = \left[ \frac{\partial}{\partial \theta_j} \nabla_{\theta}(\theta^T x) \right] = x.$$

$$\begin{aligned} l(\theta) &= -\log p(y; \eta) \\ &= -\log [b(y) \exp(\eta^T y - a(\eta))] \\ &= a(\eta) - y^T \eta - \log b(y). \end{aligned}$$

$$\begin{aligned} \nabla_{\theta} l(\theta) &= \left( \frac{\partial}{\partial \eta} a(\eta) \right) \nabla_{\theta} \eta - y \nabla_{\theta} \eta \\ &= \frac{\partial}{\partial \eta} a(\eta) x - y x. \end{aligned}$$

$$\begin{aligned} \nabla_{\theta}^2 l(\theta) &= \left[ \frac{\partial}{\partial \eta} \left( \frac{\partial}{\partial \eta} a(\eta) x - y x \right) \right] \\ &= \cancel{x} \cancel{\frac{\partial}{\partial \eta} \frac{\partial}{\partial \eta} a(\eta)} + \cancel{\frac{\partial}{\partial \eta} a(\eta)} \cancel{\frac{\partial}{\partial \eta} x} = x \frac{\partial^2 a(\eta)}{\partial \eta^2} \nabla_{\theta} \eta \\ &= x \frac{\partial^2 a(\eta)}{\partial \eta^2} \nabla_{\theta} \eta \cancel{x} \cancel{\frac{\partial}{\partial \eta} a(\eta)} \cancel{\nabla_{\theta} \eta} = \frac{\partial^2 a(\eta)}{\partial \eta^2} x x^T \\ &= x \cancel{\frac{\partial}{\partial \eta} a(\eta)} \cancel{\nabla_{\theta} \eta} v_{\text{ax}}(y; \eta) = v_{\text{ax}}(y; \eta) x x^T \end{aligned}$$

Another way

$$l(\theta) = - \left[ \log \prod_{i=1}^n b(y^{(i)}) + \sum_{i=1}^n (\theta^T x^{(i)}) y^{(i)} - \alpha(\theta^T x^{(i)}) \right]$$

$$\frac{\partial l(\theta)}{\partial \theta_j} = - \sum_{i=1}^n \left[ x_j^{(i)} y^{(i)} - \frac{\partial \alpha}{\partial \theta_j}(x_j^{(i)}) \right]$$

$$\frac{\partial^2 l(\theta)}{\partial \theta_k \partial \theta_j} = \sum_{i=1}^n x_j^{(i)} \frac{\partial^2 \alpha}{\partial \theta_j^2} x_k^{(i)} = \text{Var}(y|x,y) \sum_{i=1}^n x_j^{(i)} x_k^{(i)}$$

$$\Rightarrow H(\theta) = \text{Var}(y|x,y) \underset{n \times n}{\underset{\text{symmetric}}{\underline{x^T x}}} \quad \text{here } x = [x_1^{(i)}] = [a_{ij}]$$

Since for  $x \neq 0$   $\text{Var}(y|x,y) > 0$ .

$$\bullet x^T x x^T x = \cancel{(x^T x)^T x^T x} = (x^T x)^2 \geq 0$$

Thus  $H(\theta)$  is Positive Semi definite. Equals zero when  $x^T x = 0$  or  $\text{Var}(y|x,y) = 0$ . When  $\text{Var}(y|x,y) \neq 0$ .

then  $x^T x \neq 0$ . Thus it might happen that the matrix is positive even when  $x \neq 0$ .

Thus the function is convex (not strictly it can have a flat region)

Q5 in new version of the tutorial is on linear regression..

Q5 (a)  $J(\theta) = \frac{1}{2} \sum_{i=1}^N (\theta^T \hat{x}^{(i)} - y^{(i)})^2$

$$\nabla_{\theta} J(\theta) = \frac{1}{2} \sum_{i=1}^N 2(\theta^T \hat{x}^{(i)} - y^{(i)}) \hat{x}^{(i)}$$

Update rule:

$$\theta_j := \theta_j - \lambda \sum_{i=1}^N (\theta^T \hat{x}^{(i)} - y^{(i)}) (\hat{x}^{(i)})_j$$

#

(b) Normal eqn.

$$\nabla_{\theta} J(\theta) = \hat{x}^T \hat{x} \theta - \hat{x}^T \hat{y} = 0$$

$$\Rightarrow \hat{x}^T \hat{x} \theta = \hat{x}^T \hat{y}$$

(c) As  $N$  increases fit higher degree polynomials tend to fit data better, high ~~very~~ high degree polynomials can cause overfitting. be numerically unstable

(d) In the presence of  $\sin(\alpha)$  feature, the models seem to fit the data better/more robustly. However the numerical instability with high degree polynomials remain.

(e) When dataset is small, higher degree polynomials tend to pass through all the points, but qualitatively seem like a poor fit. Numerical instability remains.