

$$1. \quad X \in \mathbb{R}^n \quad X \in \mathbb{R}^{m \times n} \quad \theta \in \mathbb{R}^m \quad Y \in \mathbb{R}^m$$

(a) [10 points] In lecture we saw the average empirical loss for logistic regression:

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

where  $y^{(i)} \in \{0, 1\}$ ,  $h_\theta(x) = g(\theta^T x)$  and  $g(z) = 1/(1 + e^{-z})$ .

Find the Hessian  $H$  of this function, and show that for any vector  $z$ , it holds true that

$$z^T H z \geq 0.$$

$$\begin{aligned} \frac{\partial}{\partial \theta_j} J(\theta) &= \left( y \frac{1}{g(\theta^T x)} - (1-y) \frac{1}{1-g(\theta^T x)} \right) \frac{\partial}{\partial \theta_j} g(\theta^T x) \\ &= \left( y \frac{1}{g(\theta^T x)} - (1-y) \frac{1}{1-g(\theta^T x)} \right) g(\theta^T x) (1-g(\theta^T x)) \frac{\partial}{\partial \theta_j} (\theta^T x) \\ &= [y(1-g(\theta^T x)) - (1-y)g(\theta^T x)] x_j \\ &= [y - yg(\theta^T x) - g(\theta^T x) + yg(\theta^T x)] x_j \\ &= (y - h_\theta(x)) x_j \end{aligned}$$

$$h(x) = \frac{1}{1+e^{-\theta^T x}}$$

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

$$\nabla J(\theta) = -\frac{1}{m} X^T (Y - h_\theta(X))$$

$$-\frac{1}{m} \begin{bmatrix} \sum_{i=1}^m (y^{(i)} - h_\theta(x^{(i)})) x_1^{(i)} \\ \sum_{i=1}^m (y^{(i)} - h_\theta(x^{(i)})) x_2^{(i)} \\ \vdots \\ \sum_{i=1}^m (y^{(i)} - h_\theta(x^{(i)})) x_n^{(i)} \end{bmatrix}$$

$$\frac{\partial^2}{\partial \theta_j \partial \theta_k} J(\theta) = \frac{1}{m} \sum_{i=1}^m [g(\theta^T x^{(i)}) (1-g(\theta^T x^{(i)}))] x_j^{(i)} x_k^{(i)}$$

$$H > \frac{1}{m} X^T h_\theta(X) (1-h_\theta(X)) X \quad \boxed{\frac{1}{m} \begin{bmatrix} \sum_{i=1}^m h_\theta(x^{(i)}) (1-h_\theta(x^{(i)})) x_1^{(i)} x_1^{(i)} & \cdots & \cdots & \cdots & \cdots & \cdots \\ \cdots & \ddots & \cdots & \cdots & \cdots & \cdots \\ \cdots & \cdots & \ddots & \cdots & \cdots & \cdots \\ \cdots & \cdots & \cdots & \sum_{i=1}^m h_\theta(x^{(i)}) (1-h_\theta(x^{(i)})) x_n^{(i)} x_n^{(i)} & \cdots & \cdots \end{bmatrix}}$$

$$H \in \mathbb{R}^{n \times n}$$

$$\mathbb{R}^m$$

$$\begin{aligned}
 Z^T H Z &= \frac{1}{m} \sum_{j=1}^m \sum_{k=1}^n h_{jk} z_j z_k \\
 &= \frac{1}{m} \sum_{i=1}^m \sum_{j=1}^n \sum_{k=1}^n h_\theta(X^{(i)}) (1-h_\theta(X^{(i)})) X_j^{(i)} X_k^{(i)} z_j z_k \\
 &= \frac{1}{m} \sum_{i=1}^m (h_\theta(X^{(i)}) (1-h_\theta(X^{(i)}))) \sum_{j=1}^n \sum_{k=1}^n X_j^{(i)} X_k^{(i)} z_j z_k \\
 &= \frac{1}{m} \sum_{i=1}^m (h_\theta(X^{(i)}) (1-h_\theta(X^{(i)}))) [(X^{(i)})^T Z]^2 \\
 \therefore h_\theta(X^{(i)}) \& 1-h_\theta(X^{(i)}) > 0 \\
 \therefore Z^T H Z &\geq 0
 \end{aligned}$$

(c) [5 points] Recall that in GDA we model the joint distribution of  $(x, y)$  by the following equations:

$$\begin{aligned} p(y) &= \begin{cases} \phi & \text{if } y = 1 \\ 1 - \phi & \text{if } y = 0 \end{cases} \\ p(x|y=0) &= \frac{1}{(2\pi)^{n/2}|\Sigma|^{1/2}} \exp\left(-\frac{1}{2}(x - \mu_0)^T \Sigma^{-1}(x - \mu_0)\right) \\ p(x|y=1) &= \frac{1}{(2\pi)^{n/2}|\Sigma|^{1/2}} \exp\left(-\frac{1}{2}(x - \mu_1)^T \Sigma^{-1}(x - \mu_1)\right), \end{aligned}$$

where  $\phi, \mu_0, \mu_1$ , and  $\Sigma$  are the parameters of our model.

Suppose we have already fit  $\phi, \mu_0, \mu_1$ , and  $\Sigma$ , and now want to predict  $y$  given a new point  $x$ . To show that GDA results in a classifier that has a linear decision boundary, show the posterior distribution can be written as

$$p(y=1 | x; \phi, \mu_0, \mu_1, \Sigma) = \frac{1}{1 + \exp(-(\theta^T x + \theta_0))},$$

where  $\theta \in \mathbb{R}^n$  and  $\theta_0 \in \mathbb{R}$  are appropriate functions of  $\phi, \Sigma, \mu_0$ , and  $\mu_1$ .

$$\begin{aligned} P(y=1|x) &= \frac{P(x|y=1) P(y=1)}{P(x)} \\ &= \frac{p(x|y=1) P(y=1)}{p(x|y=1) P(y=1) + p(x|y=0) P(y=0)} \\ &= \frac{\phi \exp(-\frac{1}{2}(x - \mu_1)^T \Sigma^{-1}(x - \mu_1))}{\phi \exp(-\frac{1}{2}(x - \mu_1)^T \Sigma^{-1}(x - \mu_1)) + (1 - \phi) \exp(-\frac{1}{2}(x - \mu_0)^T \Sigma^{-1}(x - \mu_0))} \\ &\quad \checkmark \text{divided by numerator} \end{aligned}$$

$$\frac{1-\phi}{\phi} \exp\left\{\frac{1}{2} \left[ (\underline{x} - \underline{\mu}_1)^T \underline{\Sigma}^{-1} (\underline{x} - \underline{\mu}_1) - (\underline{x} - \underline{\mu}_0)^T \underline{\Sigma}^{-1} (\underline{x} - \underline{\mu}_0) \right]\right\}$$

$$\downarrow \quad \therefore (\underline{\Sigma}^{-1})^T = \underline{\Sigma}^{-1}$$

$$[(\underline{x}^T \underline{\Sigma}^{-1} \underline{x} - \underline{x}^T \underline{\Sigma}^{-1} \underline{\mu}_1 - \underline{\mu}_1^T \underline{\Sigma}^{-1} \underline{x} + \underline{\mu}_1^T \underline{\Sigma}^{-1} \underline{\mu}_1) - (\underline{x}^T \underline{\Sigma}^{-1} \underline{x} - \underline{x}^T \underline{\Sigma}^{-1} \underline{\mu}_0 - \underline{\mu}_0^T \underline{\Sigma}^{-1} \underline{x} + \underline{\mu}_0^T \underline{\Sigma}^{-1} \underline{\mu}_0)]$$

$$= \underline{x}^T \underline{\Sigma}^{-1} (\underline{\mu}_0 - \underline{\mu}_1) + (\underline{\mu}_0^T - \underline{\mu}_1^T) \underline{\Sigma}^{-1} \underline{x} + (\underline{\mu}_1^T \underline{\Sigma}^{-1} \underline{\mu}_1 - \underline{\mu}_0^T \underline{\Sigma}^{-1} \underline{\mu}_0)$$

$$= 2(\underline{\mu}_0 - \underline{\mu}_1)^T \underline{\Sigma}^{-1} \underline{x} - (\underline{\mu}_0 - \underline{\mu}_1)^T \underline{\Sigma}^{-1} (\underline{\mu}_0 + \underline{\mu}_1)$$

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

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

1

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

Let

$$\theta = \Sigma^{-1}(\mu_1 - \mu_0)$$

$$\theta_0 = -\frac{1}{2}(\mu_1 - \mu_0)^T \Sigma^{-1}(\mu_1 + \mu_0) - \ln \frac{1-\phi}{\phi}$$

- (d) [7 points] For this part of the problem only, you may assume  $n$  (the dimension of  $x$ ) is 1, so that  $\Sigma = [\sigma^2]$  is just a real number, and likewise the determinant of  $\Sigma$  is given by  $|\Sigma| = \sigma^2$ . Given the dataset, we claim that the maximum likelihood estimates of the parameters are given by

$$\begin{aligned}\phi &= \frac{1}{m} \sum_{i=1}^m \mathbb{1}\{y^{(i)} = 1\} \\ \mu_0 &= \frac{\sum_{i=1}^m \mathbb{1}\{y^{(i)} = 0\}x^{(i)}}{\sum_{i=1}^m \mathbb{1}\{y^{(i)} = 0\}} \\ \mu_1 &= \frac{\sum_{i=1}^m \mathbb{1}\{y^{(i)} = 1\}x^{(i)}}{\sum_{i=1}^m \mathbb{1}\{y^{(i)} = 1\}} \\ \Sigma &= \frac{1}{m} \sum_{i=1}^m (x^{(i)} - \mu_{y^{(i)}})(x^{(i)} - \mu_{y^{(i)}})^T\end{aligned}$$

The log-likelihood of the data is

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

By maximizing  $\ell$  with respect to the four parameters, prove that the maximum likelihood estimates of  $\phi$ ,  $\mu_0$ ,  $\mu_1$ , and  $\Sigma$  are indeed as given in the formulas above. (You may assume that there is at least one positive and one negative example, so that the denominators in the definitions of  $\mu_0$  and  $\mu_1$  above are non-zero.)

$$\begin{aligned}& p(x^{(i)} | y^{(i)}; \mu_0, \mu_1, \Sigma) \\ &= \frac{1}{(2\pi)^{\frac{1}{2}} |\Sigma|^{\frac{1}{2}}} \exp\left(-\frac{1}{2} (X^{(i)} - \mu_{y^{(i)}})^T \Sigma^{-1} (X^{(i)} - \mu_{y^{(i)}})\right) \\ & P(y^{(i)}; \phi) = \phi^{\mathbb{1}\{y^{(i)} = 1\}} (1-\phi)^{1-\mathbb{1}\{y^{(i)} = 1\}} \\ & \ell = \sum_{i=1}^m \log p(x^{(i)} | y^{(i)}; \mu_0, \mu_1, \Sigma) + \sum_{i=1}^m \log P(y^{(i)}; \phi) \\ & \frac{\partial}{\partial \phi} \ell = \frac{1}{\phi} \sum_{i=1}^m \mathbb{1}\{y^{(i)} = 1\} - \frac{1}{1-\phi} \sum_{i=1}^m (1-\mathbb{1}\{y^{(i)} = 1\}) \\ & \frac{\partial}{\partial \mu_0} \ell = \left[ \sum_{i=1}^m \left( -\frac{1}{2} (X^{(i)} - \mu_{y^{(i)}})^T \Sigma^{-1} (X^{(i)} - \mu_{y^{(i)}}) \right) \right]' \\ &= \left[ \sum_{i=1}^m \left( -\frac{1}{2} (X^{(i)})^2 + \mu_{y^{(i)}}^2 - 2X^{(i)}\mu_{y^{(i)}} \right) \right]' \\ &= \sum_{i=1}^m (\mathbb{1}\{y=0\}\mu_0 - X^{(i)}) \mathbb{1}\{y=0\} \\ & \frac{\partial}{\partial \mu_1} \ell = \sum_{i=1}^m (\mathbb{1}\{y=1\}\mu_1 - X^{(i)}) \mathbb{1}\{y=1\}\end{aligned}$$

$$\frac{\partial}{\partial \Sigma} \ell = \left( -\frac{1}{2} \sum_{i=1}^m \log |\Sigma| - \frac{1}{2} \sum_{i=1}^m (X - M_{y^{(i)}})^T \Sigma^{-1} (X - M_{y^{(i)}}) \right)^T$$

$$= -\frac{m}{2} \bar{\Sigma}^{-1} + \frac{1}{2} \bar{\Sigma}^{-1} \sum_{i=1}^m (X - M_{y^{(i)}}) (X - M_{y^{(i)}})^T \bar{\Sigma}^{-1}$$


---

$$\begin{cases} \frac{\partial}{\partial \phi} \ell = 0 \\ \frac{\partial}{\partial M_0} \ell = 0 \\ \frac{\partial}{\partial M_1} \ell = 0 \\ \frac{\partial}{\partial \Sigma} \ell = 0 \end{cases} \Rightarrow$$

$$\begin{aligned} \frac{1}{\phi} \sum_{i=1}^m \mathbb{I}\{y^{(i)}=1\} - \frac{1}{1-\phi} \sum_{i=1}^m (1 - \mathbb{I}\{y^{(i)}=1\}) &= 0 \\ (1-\phi) \sum_{i=1}^m \mathbb{I}\{y^{(i)}=1\} &= \phi \left( m - \sum_{i=1}^m \mathbb{I}\{y^{(i)}=1\} \right) \end{aligned}$$

$$\boxed{\phi = \frac{1}{m} \sum_{i=1}^m \mathbb{I}\{y^{(i)}=1\}}$$

$$\sum_{i=1}^m \mathbb{I}\{y=0\} M_0 - X^{(i)} \mathbb{I}\{y=0\} = 0$$

$$M_0 \sum_{i=1}^m \mathbb{I}\{y=0\} - \sum_{i=1}^m X^{(i)} \mathbb{I}\{y=0\} = 0$$

$$\boxed{M_0 = -\frac{\sum_{i=1}^m X^{(i)} \mathbb{I}\{y=0\}}{\sum_{i=1}^m \mathbb{I}\{y=0\}}} \quad M_1 = -\frac{\sum_{i=1}^m X^{(i)} \mathbb{I}\{y=1\}}{\sum_{i=1}^m \mathbb{I}\{y=1\}}$$

$$-\frac{m}{2} \cancel{\bar{\Sigma}^{-1}} + \frac{1}{2} \cancel{\bar{\Sigma}^{-1}} \sum_{i=1}^m (X - M_{y^{(i)}}) (X - M_{y^{(i)}})^T \bar{\Sigma}^{-1} = 0$$

$$\boxed{\frac{1}{m} \sum_{i=1}^m (X - M_{y^{(i)}}) (X - M_{y^{(i)}})^T = \bar{\Sigma}}$$

(h)

The dataset where GDA performed worse is far away from Gaussian distribution. GDA has stronger assumptions than Logistic Regression's. It will perform better if dataset is Gaussian distributed.

2.

- (a) [5 points] Suppose that each  $y^{(i)}$  and  $x^{(i)}$  are conditionally independent given  $t^{(i)}$ :

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

Note this is equivalent to saying that labeled examples were selected uniformly at random from the set of positive examples. Prove that the probability of an example being labeled differs by a constant factor from the probability of an example being positive. That is, show that  $p(t^{(i)} = 1 | x^{(i)}) = p(y^{(i)} = 1 | x^{(i)})/\alpha$  for some  $\alpha \in \mathbb{R}$ .

$$\begin{aligned} p(y^{(i)} = 1 | x^{(i)}, t^{(i)} = 1) &= \frac{p(t^{(i)} = 1 | y^{(i)} = 1, x^{(i)}) \cdot p(y^{(i)} = 1 | x^{(i)})}{p(t^{(i)} = 1 | x^{(i)})} \\ p(y^{(i)} = 1 | t^{(i)} = 1) \cdot p(t^{(i)} = 1 | x^{(i)}) &\stackrel{?}{=} p(y^{(i)} = 1 | x^{(i)}) \\ p(t^{(i)} = 1 | x^{(i)}) &= p(y^{(i)} = 1 | x^{(i)}) / \lambda \quad (1) \\ \text{where } \lambda &= p(y^{(i)} = 1 | t^{(i)} = 1) \end{aligned}$$

- (b) [5 points] Suppose we want to estimate  $\alpha$  using a trained classifier  $h$  and a held-out validation set  $V$ . Let  $V_+$  be the set of labeled (and hence positive) examples in  $V$ , given by  $V_+ = \{x^{(i)} \in V | y^{(i)} = 1\}$ . Assuming that  $h(x^{(i)}) \approx p(y^{(i)} = 1 | x^{(i)})$  for all examples  $x^{(i)}$ , show that

$$h(x^{(i)}) \approx \alpha \quad \text{for all } x^{(i)} \in V_+.$$

You may assume that  $p(t^{(i)} = 1 | x^{(i)}) \approx 1$  when  $x^{(i)} \in V_+$ .

from (1),  $\lambda \approx p(y^{(i)} = 1 | x^{(i)}) = h(x^{(i)}) \text{ for all } x^{(i)} \in V_+$

$$\begin{aligned} (E) \quad \frac{1}{1 + e^{-\theta^T x}} &= 0.5 \lambda \quad \theta_1 X_1 + \theta_2 X_2 + \theta_0 = \left(n \left(\frac{\lambda}{2 - \lambda}\right)\right) \\ \frac{1}{1 + e^{-\theta^T x}} &= \frac{\lambda}{2} \quad \theta_2 X_2 = -(\theta_1 X_1 + \left(n \left(\frac{\lambda}{2} - 1\right)\right) + \theta_0) \\ \theta^T x &= \left(n \left(\frac{\lambda}{2} - 1\right)\right) \quad X_2 = -\left(\frac{\theta_1}{\theta_2} X_1 + \frac{\left(n \left(\frac{\lambda}{2} - 1\right) + \theta_0\right)}{\theta_2}\right) \\ \frac{\theta_0}{\theta_2} \cdot C &= \frac{\left(n \left(\frac{\lambda}{2} - 1\right) + \theta_0\right)}{\theta_2} \\ C &= \frac{1}{\theta_0} \left(n \left(\frac{\lambda}{2} - 1\right) + 1\right) \end{aligned}$$

3.

- (a) [5 points] Consider the Poisson distribution parameterized by  $\lambda$ :

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

Show that the Poisson distribution is in the exponential family, and clearly state the values for  $b(y)$ ,  $\eta$ ,  $T(y)$ , and  $a(\eta)$ .

$$\begin{aligned} p(y; \lambda) &= \frac{e^{-\lambda} \lambda^y}{y!} = \frac{1}{y!} e^{-\lambda} \cdot e^{y \ln \lambda} \\ &= \frac{1}{y!} \exp(y \ln \lambda - \lambda) \\ \text{where } &\left\{ \begin{array}{l} b(y) = \frac{1}{y!} \\ T(y) = y \\ \eta = \ln \lambda \\ \lambda(\eta) = \lambda = e^\eta \end{array} \right. \end{aligned}$$

- (b) [3 points] Consider performing regression using a GLM model with a Poisson response variable. What is the canonical response function for the family? (You may use the fact that a Poisson random variable with parameter  $\lambda$  has mean  $\lambda$ .)

$$\begin{aligned} g(\eta) &= E(Y^{(t)}) | \eta = \frac{\partial}{\partial \eta} \lambda(\eta) = e^\eta \\ h_\theta(x) &= e^{\theta^T x} \end{aligned}$$

- (c) [7 points] For a training set  $\{(x^{(i)}, y^{(i)}); i = 1, \dots, m\}$ , let the log-likelihood of an example be  $\log p(y^{(i)}|x^{(i)}; \theta)$ . By taking the derivative of the log-likelihood with respect to  $\theta_j$ , derive the stochastic gradient ascent update rule for learning using a GLM model with Poisson responses  $y$  and the canonical response function.

$$\text{let } \eta = \theta^T X$$

$$\therefore p(y^{(i)}|X^{(i)}; \theta) = \frac{1}{y^{(i)}} \exp(y^{(i)} \theta^T X^{(i)} - e^{\theta^T X^{(i)}})$$

$$\log p(y^{(i)}|X^{(i)}; \theta) = -\log y^{(i)} + (y^{(i)} \theta^T X^{(i)} - e^{\theta^T X^{(i)}})$$

$$\frac{\partial}{\partial \theta_j} \log p(y^{(i)}|X^{(i)}; \theta) = y^{(i)} X_j^{(i)} - X_j^{(i)} e^{\theta^T X^{(i)}}$$

$$\theta_j := \theta_j + 2(X_j^{(i)} y^{(i)} - X_j^{(i)} e^{\theta^T X^{(i)}})$$

4.

- (a) [5 points] Derive an expression for the mean of the distribution. Show that  $\mathbb{E}[Y | X; \theta]$  can be represented as the gradient of the log-partition function  $\lambda$  with respect to the natural parameter  $\eta$ .

**Hint:** Start with observing that  $\frac{\partial}{\partial \eta} \int p(y; \eta) dy = \int \frac{\partial}{\partial \eta} p(y; \eta) dy$ .

To prove:

$$\mathbb{E}(Y | X; \theta) = \mathbb{E}(y; \theta^T x) = \int y p(y; \eta) dy$$

$$p(y; \eta) = b(y) \exp[\eta^T y - \lambda(\eta)]$$

$$\frac{\partial}{\partial \eta} p(y; \eta) = \left( y - \frac{\partial}{\partial \eta} \lambda(\eta) \right) b(y) \exp(\eta^T y - \lambda(\eta))$$

$$\frac{\partial}{\partial \eta} \int p(y; \eta) dy = y \int p(y; \eta) dy - \frac{\partial}{\partial \eta} \lambda(\eta) \int p(y; \eta) dy \quad (1)$$

$$\frac{\partial}{\partial \eta} \int p(y; \eta) dy = \underbrace{\int y p(y; \eta) dy}_{\uparrow E(y; \theta^T x)} - \frac{\partial}{\partial \eta} \lambda(\eta) \int p(y; \eta) dy \quad \uparrow \int p(y; \eta) dy = 1$$

$$0 = \mathbb{E}(y; \theta^T x) - \frac{\partial}{\partial \eta} \lambda(\eta)$$

$$\mathbb{E}(y; \theta^T x) = \frac{\partial}{\partial \eta} \lambda(\eta)$$

- (b) [5 points] Next, derive an expression for the variance of the distribution. In particular, show that  $\text{Var}(Y | X; \theta)$  can be expressed as the derivative of the mean w.r.t  $\eta$  (i.e., the second derivative of the log-partition function  $a(\eta)$  w.r.t the natural parameter  $\eta$ .)

$$\text{Var}(Y | X; \theta) = \text{Var}(Y; \theta^T X) = \text{Var}(Y; \eta)$$

$$= E(Y^2; \eta) - (E(Y; \eta))^2$$

$$\frac{\partial}{\partial \eta} p(y; \eta) = y p(y; \eta) - \frac{\partial}{\partial \eta} a(\eta) p(y; \eta)$$

$$\frac{\partial^2}{\partial \eta^2} p(y; \eta) = y [y \cdot p(y; \eta) - \frac{\partial}{\partial \eta} a(\eta) p(y; \eta)]$$

$$- \frac{\partial^2}{\partial \eta^2} a(\eta) p(y; \eta) - \frac{\partial}{\partial \eta} a(\eta) \cdot \frac{\partial}{\partial \eta} p(y; \eta)$$

$$\underbrace{E(y)}_{\text{red}} \quad \underbrace{E(y)}_{\text{red}}$$

$$0 = \int y^2 p(y; \eta) dy - \frac{\partial}{\partial \eta} a(\eta) \int y p(y; \eta) dy - \frac{\partial^2}{\partial \eta^2} a(\eta)$$

$$\frac{\partial^2}{\partial \eta^2} a(\eta) = E(Y^2; \eta) - (E(Y; \eta))^2$$

- (c) [5 points] Finally, write out the loss function  $\ell(\theta)$ , the NLL of the distribution, as a function of  $\theta$ . Then, calculate the Hessian of the loss w.r.t  $\theta$ , and show that it is always PSD. This concludes the proof that NLL loss of GLM is convex.

**Hint:** Use the chain rule of calculus along with the results of the previous parts to simplify your derivations.

For a single training example:

$$\ell(\eta) = -\log p(y^{(i)}; \eta) = \log p(y^{(i)} | X^{(i)}; \theta)$$

$$\ell(\eta) = -\log b(y^{(i)}) - (\eta \cdot y^{(i)} - \lambda(\eta))$$

$$\frac{\partial}{\partial \eta} \ell(\eta) = \frac{\partial}{\partial \eta} \lambda(\eta) - y^{(i)} = E(Y|\eta) - y^{(i)}$$

$$\frac{\partial^2}{\partial \eta^2} \ell(\eta) = \frac{\partial^2}{\partial \eta^2} \lambda(\eta) = \text{Var}(Y|\eta)$$

$$\begin{aligned} \frac{\partial^2}{\partial \theta_j \partial \theta_k} \ell(\theta) &= \frac{\partial}{\partial \eta} \ell(\eta) \cdot 0 - \frac{\partial^2}{\partial \eta^2} \ell(\eta) \frac{\partial}{\partial \theta_j}(\eta) \frac{\partial}{\partial \theta_k}(\eta) \\ &= \text{Var}(Y|\eta) X_j X_k \end{aligned}$$

$$\begin{aligned} Z^T H Z &= \text{Var}(Y|\eta) \cdot \sum_j \sum_k Z_j X_j X_k Z_k \\ &= \text{Var}(Y|\eta) (X^T Z)^2 / n \end{aligned}$$

$$5.(a) \quad J(\theta) = \frac{1}{2} \sum_{i=1}^m w^{(i)} (\theta^T X^{(i)} - y^{(i)})^2$$

i. [2 points] Show that  $J(\theta)$  can also be written

$$J(\theta) = (X\theta - y)^T W (X\theta - y)$$

for an appropriate matrix  $W$ , and where  $X$  and  $y$  are as defined in class. Clearly specify the value of each element of the matrix  $W$ .

$$W \in \mathbb{R}^{m \times m}$$

$$W = \text{diag}(\frac{1}{2} w^{(i)}) \quad \text{for } i \in \{0, 1, 2, \dots, m\}$$

ii. [4 points] If all the  $w^{(i)}$ 's equal 1, then we saw in class that the normal equation is

$$X^T X \theta = X^T y,$$

and that the value of  $\theta$  that minimizes  $J(\theta)$  is given by  $(X^T X)^{-1} X^T y$ . By finding the derivative  $\nabla_{\theta} J(\theta)$  and setting that to zero, generalize the normal equation to this weighted setting, and give the new value of  $\theta$  that minimizes  $J(\theta)$  in closed form as a function of  $X$ ,  $W$  and  $y$ .

$$\begin{aligned} J_{\theta}(\theta) &= (X\theta - y)^T W (X\theta - y) & W^T = W \\ &= (\theta^T X^T - y^T) W (X\theta - y) \\ &= (\theta^T X^T W X \theta - \theta^T X^T W y - y^T W X \theta + y^T W y) \\ &= (\theta^T X^T W X \theta - 2\theta^T X^T W y) \\ \nabla_{\theta} J(\theta) &= 2X^T W X \theta - 2X^T W y = 0 \\ \Rightarrow X^T W X \theta &= X^T W y \\ \theta &= (X^T W X)^{-1} X^T W y \end{aligned}$$

$$\begin{aligned} \text{or } \nabla_{\theta} J(\theta) &= \nabla_{\theta} (X\theta - y)(W(X\theta - y) + W^T(X\theta - y)) \\ &= 2X^T W (X\theta - y) = 2X^T W X - 2X^T W y \\ \Rightarrow \theta &= (X^T W X)^{-1} X^T W y \end{aligned}$$

- iii. [4 points] Suppose we have a dataset  $\{(x^{(i)}, y^{(i)}); i = 1 \dots, m\}$  of  $m$  independent examples, but we model the  $y^{(i)}$ 's as drawn from conditional distributions with different levels of variance  $(\sigma^{(i)})^2$ . Specifically, assume the model

$$p(y^{(i)}|x^{(i)}; \theta) = \frac{1}{\sqrt{2\pi\sigma^{(i)}}} \exp\left(-\frac{(y^{(i)} - \theta^T x^{(i)})^2}{2(\sigma^{(i)})^2}\right)$$

That is, each  $y^{(i)}$  is drawn from a Gaussian distribution with mean  $\theta^T x^{(i)}$  and variance  $(\sigma^{(i)})^2$  (where the  $\sigma^{(i)}$ 's are fixed, known, constants). Show that finding the maximum likelihood estimate of  $\theta$  reduces to solving a weighted linear regression problem. State clearly what the  $w^{(i)}$ 's are in terms of the  $\sigma^{(i)}$ 's.

$$\begin{aligned} l(\theta) &= \log p(y^{(i)}|x^{(i)}; \theta) = \log \frac{1}{\sqrt{2\pi\sigma^{(i)}}} - \frac{(y^{(i)} - \theta^T x^{(i)})^2}{2(\sigma^{(i)})^2} \\ \frac{\partial}{\partial \theta_j} l(\theta) &= \frac{1}{(\sigma^{(i)})^2} (y^{(i)} - \theta^T x^{(i)}) x_j^{(i)} \end{aligned}$$