

Notation: let $X_p^{(i)} \in \mathbb{R}^d$, where $i = 1, \dots, n$. $p = 1, 2, \dots, d$
 n is the sample size
 d is the dimension of features

$$\bar{X} = \begin{bmatrix} X_p^{(1)} \\ X_p^{(2)} \\ \vdots \\ X_p^{(n)} \end{bmatrix} \text{ is the total sample data}$$

1.a. Let $Z^{(i)} = \theta^T X^{(i)}$ where $\theta \in \mathbb{R}^d$

$$f_1(Z^{(i)}) = -\log\left(\frac{1}{1+e^{-Z^{(i)}}}\right) \quad \textcircled{1}$$

$$f_2(Z^{(i)}) = -\log\left(1 - \frac{1}{1+e^{-Z^{(i)}}}\right) = Z^{(i)} + f_1(Z^{(i)}) \quad \textcircled{2}$$

We can rewrite $J(\theta)$ as

$$J(\theta) = \frac{1}{n} \sum_{i=1}^n [y^{(i)} f_1(Z^{(i)}) + (1-y^{(i)}) f_2(Z^{(i)})] \quad \textcircled{3}$$

$$\text{since } \frac{\partial f_1(Z^{(i)})}{\partial \theta_p} = -\frac{\partial}{\partial \theta_p} \log\left(\frac{1}{1+e^{-Z^{(i)}}}\right) = -\frac{\partial}{\partial Z^{(i)}} \log\left(\frac{1}{1+e^{-Z^{(i)}}}\right) \cdot \frac{\partial Z^{(i)}}{\partial \theta_p}$$

$$= [-1+g(Z^{(i)})] X_p^{(i)}$$

$$\frac{\partial f_2(Z^{(i)})}{\partial \theta_p \partial \theta_q} = \frac{\partial}{\partial \theta_q} ([-1+g(Z^{(i)})] \cdot X_p^{(i)}) = \frac{\partial}{\partial Z^{(i)}} (-1+g(Z^{(i)})) \cdot \frac{\partial Z^{(i)}}{\partial \theta_q}$$

$$= g(Z^{(i)}) [1-g(Z^{(i)})] X_p^{(i)} \cdot X_q^{(i)} \quad \textcircled{4}$$

$$\textcircled{2} \Rightarrow \frac{\partial f_2(Z^{(i)})}{\partial \theta_p} = \frac{\partial}{\partial \theta_p} [Z^{(i)} + f_1(Z^{(i)})] = X_p^{(i)} + \frac{\partial}{\partial \theta_p} f_1(Z^{(i)})$$

$$\frac{\partial f_2(Z^{(i)})}{\partial \theta_p \partial \theta_q} = \frac{\partial}{\partial \theta_q} [X_p^{(i)} + \frac{\partial}{\partial \theta_p} f_1(Z^{(i)})] = \frac{\partial}{\partial \theta_p \partial \theta_q} f_1(Z^{(i)}) \quad \textcircled{5}$$

$$\textcircled{4} \textcircled{5} \Rightarrow H_{pq} = \frac{\partial^2 J(\theta)}{\partial \theta_p \partial \theta_q} = \frac{1}{n} \sum_{i=1}^n [y^{(i)} \frac{\partial}{\partial \theta_p \partial \theta_q} f_1(Z^{(i)}) + (1-y^{(i)}) \frac{\partial}{\partial \theta_p \partial \theta_q} f_2(Z^{(i)})]$$

$$= \frac{1}{n} \sum_{i=1}^n \frac{\partial}{\partial \theta_p \partial \theta_q} f_1(Z^{(i)}) = \frac{1}{n} \sum_{i=1}^n g(Z^{(i)}) [1-g(Z^{(i)})] X_p^{(i)} X_q^{(i)} \quad \textcircled{6}$$

$$\textcircled{6} H = \frac{1}{n} \bar{X}^T \begin{bmatrix} g'(Z^{(1)}) \\ g'(Z^{(2)}) \\ \vdots \\ g'(Z^{(n)}) \end{bmatrix} \bar{X} = \frac{1}{n} \bar{X}^T \operatorname{diag} \{g'(Z^{(i)})\}_{i=1, \dots, n} \bar{X}$$

Therefore, for any vector Z , $Z^T H Z = \frac{1}{n} (\bar{X} Z)^T \operatorname{diag} \{g'(Z^{(i)})\}_{i=1, \dots, n} (\bar{X} Z)$

since $g'(Z^{(i)}) = g(Z^{(i)}) (1-g(Z^{(i)})) \geq 0$, for $i = 1, \dots, n$, we can conclude that

$$Z^T H Z \geq 0$$

$$(I.C) \text{ Define } A = -\frac{1}{2}(X - \mu_1)^T \Sigma^{-1} (X - \mu_1) \quad ①$$

$$B = -\frac{1}{2}(X - \mu_0)^T \Sigma^{-1} (X - \mu_0) \quad ②$$

$$C = \frac{1}{(2\pi)^{d/2} |\Sigma|^{1/2}} \quad ③$$

from Bayes Rule, we can get

$$\begin{aligned} p(y=1|X) &= \frac{p(X|y=1)p(y=1)}{p(X|y=1)p(y=1) + p(X|y=0)p(y=0)} = \frac{Ce^A\phi}{Ce^A\phi + Ce^B(1-\phi)} \\ &= \frac{1}{1 + e^{-(A-B+\ln\phi-\ln(1-\phi))}} \end{aligned} \quad ④$$

$$\begin{aligned} \text{from } ①② \Rightarrow A - B &= -\frac{1}{2}(X - \mu_1)^T \Sigma^{-1} (X - \mu_1) + \frac{1}{2}(X - \mu_0)^T \Sigma^{-1} (X - \mu_0) \\ &= (\mu_1 - \mu_0)^T \Sigma^{-1} X + \frac{1}{2}\mu_0^T \Sigma^{-1} \mu_0 - \frac{1}{2}\mu_1^T \Sigma^{-1} \mu_1 \end{aligned} \quad ⑤$$

$$\text{from } ⑤ \Rightarrow A - B + \ln\phi - \ln(1-\phi) = (\mu_1 - \mu_0)^T \Sigma^{-1} X + \frac{1}{2}\mu_0^T \Sigma^{-1} \mu_0 - \frac{1}{2}\mu_1^T \Sigma^{-1} \mu_1 + \ln\phi - \ln(1-\phi)$$

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

we can get

$$A - B + \ln\phi = \theta^T X + \theta_0$$

so we can rewrite ④ as

$$p(y=1|X) = \frac{1}{1 + \exp(-(\theta^T X + \theta_0))}$$

1.d) The log likelihood of the training dataset is given by

$$\begin{aligned}
 L(\phi, \mu_0, \mu_1, \Sigma) &= \log \prod_{i=1}^n p(x^{(i)}, y^{(i)}; \phi, \mu_0, \mu_1, \Sigma) \\
 &= \sum_{i=1}^n \log p(x^{(i)} | y^{(i)}; \mu_0, \mu_1, \Sigma) + \sum_{i=1}^n \log p(y^{(i)}; \phi) \\
 &= \sum_{i=1}^n \log \frac{1}{(2\pi)^{d/2}} + \frac{1}{2} \sum_{i=1}^n \log \frac{1}{|\Sigma|} - \frac{1}{2} \sum_{i=1}^n \mathbb{1}\{y^{(i)}=0\} (x^{(i)} - \mu_0)^T \Sigma^{-1} (x^{(i)} - \mu_0) \\
 &\quad - \frac{1}{2} \sum_{i=1}^n \mathbb{1}\{y^{(i)}=1\} (x^{(i)} - \mu_1)^T \Sigma^{-1} (x^{(i)} - \mu_1)
 \end{aligned}$$

Now, the likelihood is maximized by setting the derivative with respect to each of the parameters to zero :

$$\text{For } \phi : \frac{\partial L}{\partial \phi} = \frac{1}{\phi} \sum_{i=1}^n \mathbb{1}\{y^{(i)}=1\} - \frac{1}{1-\phi} \sum_{i=1}^n \mathbb{1}\{y^{(i)}=0\}$$

$$\text{setting } \frac{\partial L}{\partial \phi} = 0 \text{ and solving for } \phi \Rightarrow \phi = \frac{1}{n} \sum_{i=1}^n \mathbb{1}\{y^{(i)}=1\}$$

$$\begin{aligned}
 \text{For } \mu_0 : \frac{\partial L}{\partial \mu_0} &= -\frac{1}{2} \sum_{i=1}^n \mathbb{1}\{y^{(i)}=0\} \frac{\partial}{\partial \mu_0} (x^{(i)} - \mu_0)^T \Sigma^{-1} (x^{(i)} - \mu_0) \\
 &= -\frac{1}{2} \sum_{i=1}^n \mathbb{1}\{y^{(i)}=0\} \frac{\partial}{\partial \mu_0} (x^{(i)\top} \Sigma^{-1} x^{(i)}) - x^{(i)\top} \Sigma^{-1} \mu_0 - \mu_0^\top \Sigma^{-1} x^{(i)} + \mu_0^\top \Sigma^{-1} \mu_0 \\
 &= \frac{1}{2} \sum_{i=1}^n \mathbb{1}\{y^{(i)}=0\} \left[\frac{\partial}{\partial \mu_0} (x^{(i)\top} \Sigma^{-1} \mu_0) + \frac{\partial}{\partial \mu_0} (\mu_0^\top \Sigma^{-1} x^{(i)}) - \frac{\partial}{\partial \mu_0} (\mu_0^\top \Sigma^{-1} \mu_0) \right] \\
 &= \frac{1}{2} \sum_{i=1}^n \mathbb{1}\{y^{(i)}=0\} \left[(\Sigma^{-1})^\top x^{(i)} + \Sigma^{-1} x^{(i)} - (\Sigma^{-1} + (\Sigma^{-1})^\top) \mu_0 \right]
 \end{aligned}$$

since Σ^{-1} is symmetric, $(\Sigma^{-1})^\top = \Sigma^{-1}$, Setting $\frac{\partial L}{\partial \mu_0} = 0$ and solving for $\mu_0 \Rightarrow$

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

similarly, for μ_1 :

$$\frac{\partial L}{\partial \mu_1} = \frac{1}{2} \sum_{i=1}^n \mathbb{1}\{y^{(i)}=1\} \left[(\Sigma^{-1})^\top x^{(i)} + \Sigma^{-1} x^{(i)} - (\Sigma^{-1} + (\Sigma^{-1})^\top) \mu_1 \right]$$

since Σ^{-1} is symmetric, $(\Sigma^{-1})^\top = \Sigma^{-1}$, Setting $\frac{\partial L}{\partial \mu_1} = 0$ and solving for $\mu_1 \Rightarrow$

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

1.d. for Σ :

$$\frac{\partial L}{\partial \Sigma} = \frac{1}{2} \Sigma^{-1} - \frac{n}{2} \sum_{i=1}^n \left[\log\left(\frac{1}{\Sigma}\right) - \frac{1}{2} \sum_{i=1}^n (\mathbf{x}^{(i)} - \boldsymbol{\mu}_{y^{(i)}})^T \Sigma^{-1} (\mathbf{x}^{(i)} - \boldsymbol{\mu}_{y^{(i)}}) \right]$$

with $\frac{\partial}{\partial \Sigma} \log\left(\frac{1}{\Sigma}\right) = -(\Sigma^{-1})^T = -\Sigma^{-1}$ (since Σ is symmetric)

$$\frac{\partial}{\partial \Sigma} \Sigma^{-1} = -\Sigma^{-1} \Sigma^{-1} \quad (\text{derived by } \frac{\partial \Sigma}{\partial \Sigma} = \frac{\partial}{\partial \Sigma} (\Sigma^{-1} \Sigma) = 0)$$

we can further get

$$\frac{\partial L}{\partial \Sigma} = -\frac{1}{2} n \Sigma^{-1} + \frac{1}{2} \sum_{i=1}^n (\mathbf{x}^{(i)} - \boldsymbol{\mu}_{y^{(i)}})(\mathbf{x}^{(i)} - \boldsymbol{\mu}_{y^{(i)}})^T (\Sigma^{-1})(\Sigma^{-1})$$

Setting $\frac{\partial L}{\partial \Sigma} = 0$ and solving for $\Sigma \Rightarrow$

$$\Sigma = \frac{1}{n} \sum_{i=1}^n (\mathbf{x}^{(i)} - \boldsymbol{\mu}_{y^{(i)}})(\mathbf{x}^{(i)} - \boldsymbol{\mu}_{y^{(i)}})^T$$

1.f observations: GDA seems to be affected more towards the outliers in the validation dataset 1. This indicates the GDA may have a higher bias compared to logistic regression.

1.g GDA seems to perform worse than logistic regression for Dataset 1;
GDA makes more specific assumptions about the dataset than logistic regression, and GDA will perform well if those assumptions are true, so the difference might be caused by the difference of distribution between Dataset 1 and Dataset 2.

1.h The z-transformation given by

$$z^{(i)} = \frac{\mathbf{x}^{(i)} - \boldsymbol{\mu}}{s}, \quad i=1, \dots, n$$

where $\boldsymbol{\mu}$ is the mean of the dataset and s is the standard deviation of the dataset might help improve the GDA's performance with Dataset 1.

2.a) since the exponential family is represented by

$$p(y|x; \eta) = b(y) e^{\eta^T T(y) - a(\eta)}$$

and since the poisson distribution can be re-written as

$$\begin{aligned} p(y; \lambda) &= \frac{e^{-\lambda} \lambda^y}{y!} \\ &= \frac{e^{-\lambda} e^{\lambda y}}{y!} = \frac{1}{y!} \exp(y \ln \lambda - \lambda) \end{aligned}$$

with $b(y) = \frac{1}{y!}$, $\eta = \ln \lambda$, $T(y) = y$, and $a(\eta) = \lambda$, we can conclude that poisson distribution is in the exponential family

2.b) Canonical response function is given by

$$g(\eta) = E[T(y); \eta] = E[y; \eta] = \lambda = \exp(\eta)$$

2.c)

The log likelihood of an example $(x^{(i)}, y^{(i)})$ is defined by

$L(\theta) = -\log p(y^{(i)}|x^{(i)}; \theta)$ and from 2.a, we can further derive

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

$$\text{so } \frac{\partial L(\theta)}{\partial \theta} = y^{(i)\top} x^{(i)} - \exp(\theta^T x^{(i)}) \cdot x^{(i)}$$

$$\Rightarrow \theta^{(i)} := \theta^{(i)} + \alpha \frac{\partial L(\theta)}{\partial \theta}$$

$$= \theta^{(i)} + \alpha [y^{(i)\top} x^{(i)} - \exp(\theta^T x^{(i)})] x^{(i)}$$

$$= \theta^{(i)} + \alpha [y^{(i)} - \exp(\theta^T x^{(i)})] x^{(i)} \quad \text{where } i=1, \dots, n$$