

(1.a)

$$\begin{aligned}
\frac{\partial J}{\partial \theta} &= -\frac{1}{m} \sum_{i=1}^m y^{(i)} \frac{1}{h_{\theta}(x^{(i)})} h_{\theta}(x^{(i)}) (1 - h_{\theta}(x^{(i)})) (-x^{(i)}) \\
&\quad + (1 - y^{(i)}) \frac{1}{1 - h_{\theta}(x^{(i)})} h_{\theta}(x^{(i)}) (1 - h_{\theta}(x^{(i)})) x^{(i)} \\
&= \frac{1}{m} \sum_{i=1}^m (y^{(i)} - h_{\theta}(x^{(i)})) x^{(i)} \\
H = \partial \frac{\partial J}{\partial \theta} &= \frac{1}{m} \sum_{i=1}^m h_{\theta}(x^{(i)}) (1 - h_{\theta}(x^{(i)})) x^{(i)} x^{(i)T}
\end{aligned}$$

To prove  $H$  is positive semidefinite, we show  $z^T H z \geq 0$  for all  $z$

$$\begin{aligned}
z^T H z &= \frac{1}{m} z^T \left( \sum_{i=1}^m h(x^{(i)}) (1 - h(x^{(i)})) x^{(i)} x^{(i)T} \right) z \\
&= \frac{1}{m} \sum_{i=1}^m h(x^{(i)}) (1 - h(x^{(i)})) z^T x^{(i)} x^{(i)T} z \\
&= \frac{1}{m} \sum_{i=1}^m h(x^{(i)}) (1 - h(x^{(i)})) (z^T x^{(i)})^2 \geq 0
\end{aligned}$$

(1.c)

For shorthand, we let  $\mathcal{H} = \{\phi, \Sigma, \mu_0, \mu_1\}$  denote the parameters for the problem. since the given formulae are conditioned on  $y$ , use Bayes rule to get:

$$\begin{aligned}
p(y = 1 \mid x; \phi, \Sigma, \mu_0, \mu_1) &= \frac{p(x \mid y = 1; \phi, \Sigma, \mu_0, \mu_1) p(y = 1; \phi, \Sigma, \mu_0, \mu_1)}{p(x; \phi, \Sigma, \mu_0, \mu_1)} \\
&= \frac{p(x \mid y = 1; \mathcal{H}) p(y = 1; \mathcal{H})}{p(x \mid y = 1; \mathcal{H}) p(y = 1; \mathcal{H}) + p(x \mid y = 0; \mathcal{H}) p(y = 0; \mathcal{H})} \\
&= \frac{\exp\left(-\frac{1}{2}(x - \mu_1)^T \Sigma^{-1} (x - \mu_1)\right) \phi}{\exp\left(-\frac{1}{2}(x - \mu_1)^T \Sigma^{-1} (x - \mu_1)\right) \phi + \exp\left(-\frac{1}{2}(x - \mu_0)^T \Sigma^{-1} (x - \mu_0)\right) (1 - \phi)} \\
&= \frac{1}{1 + \frac{1-\phi}{\phi} \exp\left(-\frac{1}{2}(x - \mu_0)^T \Sigma^0 (x - \mu_0) + \frac{1}{2}(x - \mu_1)^T \Sigma^0 (x - \mu_1)\right)} \\
&= \frac{1}{1 + \exp\left(\log\left(\frac{1-\phi}{\phi}\right) - \frac{1}{2}(x - \mu_0)^T \Sigma^{-1} (x - \mu_0) + \frac{1}{2}(x - \mu_1)^T \Sigma^{-1} (x - \mu_1)\right)}
\end{aligned}$$

Now, we expand and rearrange the difference of quadratic terms in the preceding expression, finding that

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

Thus, we have

$$p(y = 1 \mid x; \mathcal{H}) = \frac{1}{1 + \exp\left(\log \frac{1-\phi}{\phi} + \frac{1}{2}\mu_1^T \Sigma^{-1} \mu_1 - \frac{1}{2}\mu_0^T \Sigma^{-1} \mu_0 + (\mu_0 - \mu_1)^T \Sigma^{-1} x\right)}$$

and setting

$$\theta = \Sigma^{-1} (\mu_1 - \mu_0) \text{ and } \theta_0 = \frac{1}{2} (\mu_0^T \Sigma^{-1} \mu_0 - \mu_1^T \Sigma^{-1} \mu_1) - \log \frac{1-\phi}{\phi}$$

gives that

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

(1.d)

$$\begin{aligned}
\ell(\phi, \mu_0, \mu_1, \Sigma) &= \log \prod_{i=1}^m p(x^{(i)} \mid y^{(i)}; \mu_0, \mu_1, \Sigma) p(y^{(i)}; \phi) \\
&= \sum_{i=1}^m \log p(x^{(i)} \mid y^{(i)}; \mu_0, \mu_1, \Sigma) + \sum_{i=1}^m \log p(y^{(i)}; \phi) \\
&\simeq \sum_{i=1}^m \left[ \frac{1}{2} \log \frac{1}{|\Sigma|} - \frac{1}{2} (x^{(i)} - \mu_{y^{(i)}})^T \Sigma^{-1} (x^{(i)} - \mu_{y^{(i)}}) + y^{(i)} \log \phi + (1 - y^{(i)}) \log(1 - \phi) \right]
\end{aligned}$$

$$\begin{aligned}
\frac{\partial \ell}{\partial \phi} &= \sum_{i=1}^m \left[ \frac{y^{(i)}}{\phi} - \frac{1 - y^{(i)}}{1 - \phi} \right] \\
&= \frac{\sum_{i=1}^m 1 \{y^{(i)} = 1\}}{\phi} - \frac{m - \sum_{i=1}^m 1 \{y^{(i)} = 1\}}{1 - \phi}
\end{aligned}$$

$$\begin{aligned}
\nabla_{\mu_0} \ell &= -\frac{1}{2} \sum_{i: y^{(i)} = -1} \nabla_{\mu_0} (x^{(i)} - \mu_0)^T \Sigma^{-1} (x^{(i)} - \mu_0) \\
&= -\frac{1}{2} \sum_{i: y^{(i)} = -1} \nabla_{\mu_0} \left[ \mu_0^T \Sigma^{-1} \mu_0 - x^{(i)T} \Sigma^{-1} \mu_0 - \mu_0^T \Sigma^{-1} x^{(i)} \right] \\
&= -\frac{1}{2} \sum_{i: y^{(i)} = -1} \nabla_{\mu_0} \text{tr} \left[ \mu_0^T \Sigma^{-1} \mu_0 - x^{(i)T} \Sigma^{-1} \mu_0 - \mu_0^T \Sigma^{-1} x^{(i)} \right] \\
&= -\frac{1}{2} \sum_{i: y^{(i)} = -1} [2\Sigma^{-1} \mu_0 - 2\Sigma^{-1} x^{(i)}]
\end{aligned}$$

For  $\Sigma$ , we find the gradient with respect to  $S = \Sigma^{-1}$  rather than  $\Sigma$  just to simplify the derivation (note that  $|S| = \frac{1}{|\Sigma|}$ ). You should convince yourself that the maximum likelihood estimate  $S_m$  found in this way would correspond to the actual maximum likelihood estimate  $\Sigma_m$  as  $S_m^{-1} = \Sigma_m$

$$\begin{aligned}
\nabla_S \ell &= \sum_{i=1}^m \nabla_S \left[ \frac{1}{2} \log |S| - \frac{1}{2} \underbrace{(x^{(i)} - \mu_{y^{(i)}})^T}_{b_i^T} \underbrace{S(x^{(i)} - \mu_{y^{(i)}})}_{b_i} \right] \\
&= \sum_{i=1}^m \left[ \frac{1}{2|S|} \nabla_S |S| - \frac{1}{2} \nabla_S b_i^T S b_i \right]
\end{aligned}$$

But, we have the following identities:

$$\begin{aligned}
\nabla_S |S| &= |S| (S^{-1})^T \\
\nabla_S b_i^T S b_i &= \nabla_S \text{tr}(b_i^T S b_i) = \nabla_S \text{tr}(S b_i b_i^T) = b_i b_i^T
\end{aligned}$$

In the above, we again used matrix calculus identities, and also the commutativity of the trace operator for square matrices. Putting these into the original equation, we get:

$$\begin{aligned}
\nabla_S \ell &= \sum_{i=1}^m \left[ \frac{1}{2} S^{-1} - \frac{1}{2} b_i b_i^T \right] \\
&= \frac{1}{2} \sum_{i=1}^m [\Sigma - b_i b_i^T]
\end{aligned}$$

(1.g)

The dataset is more concentrated and unseparable.

(1.h)

Use logarithmic function to preprocess the data.

(2.a)

Because:

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

then:

$$\begin{aligned} left &= \frac{p(y^{(i)} = 1, t^{(i)} = 1, x^{(i)})}{p(t^{(i)} = 1, x^{(i)})} = \frac{p(y^{(i)} = 1, t^{(i)} = 1, x^{(i)})}{p(t^{(i)} = 1 | x^{(i)}) p(x^{(i)})} \\ &= \frac{p(y^{(i)} = 1, x^{(i)})}{p(t^{(i)} = 1 | x^{(i)}) p(x^{(i)})} = p(y^{(i)} = 1 | t^{(i)} = 1) \end{aligned}$$

So:

$$\begin{aligned} p(t^{(i)} = 1 | x^{(i)}) &= \frac{p(y^{(i)} = 1 | x^{(i)})}{p(y^{(i)} = 1 | t^{(i)} = 1)} = \frac{p(y^{(i)} = 1 | x^{(i)})}{\alpha} \\ \alpha &= p(y^{(i)} = 1 | t^{(i)} = 1) \end{aligned}$$

(2.b)

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

(3.a)

Rewrite the distribution function as:

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

Comparing with the standard form for the exponential family:

$$\begin{aligned} b(y) &= \frac{1}{y!} \\ \eta &= \log \lambda \\ T(y) &= y \\ a(\eta) &= e^\eta \end{aligned}$$

(3.b)

The canonical response function for the GLM model will be:

$$\begin{aligned} g(\eta) &= E[y; \eta] \\ &= \lambda \\ &= e^\eta \end{aligned}$$

(3.c)

The log-likelihood of an example  $(x^{(i)}, y^{(i)})$  is defined as  $\ell(\theta) = \log p(y^{(i)} | x^{(i)}; \theta)$ . To derive the stochastic gradient ascent rule, use the results in part (a) and the standard GLM assumption that  $\eta = \theta^T x$

$$\begin{aligned}\frac{\partial \ell(\theta)}{\partial \theta_j} &= \frac{\partial \log p(y^{(i)} | x^{(i)}; \theta)}{\partial \theta_j} \\&= \frac{\partial \log\left(\frac{1}{y^{(i)}!} \exp(\eta^T y^{(i)} - e^\eta)\right)}{\partial \theta_j} + \frac{\partial \log\left(\frac{1}{y^{(i)}!}\right)}{\partial \theta_j} \\&= \frac{\partial \log\left(\exp\left((\theta^T x^{(i)})^T y^{(i)} - e^{\theta^T x^{(i)}}\right)\right)}{\partial \theta_j} \\&= \frac{\partial \left((\theta^T x^{(i)})^T y^{(i)} - e^{\theta^T x^{(i)}}\right)}{\partial \theta_j} \\&= \frac{\partial \left(\left(\sum_k \theta_k x_k^{(i)}\right) y^{(i)} - e^{\sum_k \theta_k x_k^{(i)}}\right)}{\partial \theta_j} \\&= x_j^{(i)} y^{(i)} - e^{\sum_k \theta_k x_k^{(i)}} x_j^{(i)} \\&= \left(y^{(i)} - e^{\theta^T x^{(i)}}\right) x_j^{(i)}\end{aligned}$$

Thus the stochastic gradient ascent update rule should be:

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

which reduces here to:

$$\theta_j := \theta_j + \alpha \left(y^{(i)} - e^{\theta^T x}\right) x_j^{(i)}$$

(4.a)

$$\begin{aligned}\frac{\partial}{\partial \eta} \int p(y; \eta) dy &= \int \frac{\partial}{\partial \eta} p(y; \eta) dy \\&= \int b(y)(y - a'(\eta)) \exp(\eta y - a(\eta)) dy \\&= \int y b(y) \exp(\eta y - a(\eta)) dy - \int b(y) a'(\eta) \exp(\eta y - a(\eta)) dy \\&= \mathbb{E}_y(y | \eta) - a'(\eta) \int b(y) \exp(\eta y - a(\eta)) dy \\&= \mathbb{E}_y(y | \eta) - a'(\eta) = 0\end{aligned}$$

(4.b)

$$\begin{aligned}
\frac{\partial}{\partial^2 \eta} \int p(y; \eta) dy &= \int \frac{\partial}{\partial^2 \eta} p(y; \eta) dy \\
&= \int \frac{\partial}{\partial \eta} b(y) (y - a'(\eta)) \exp(\eta y - a(\eta)) dy \\
&= \int \frac{\partial}{\partial \eta} y b(y) \exp(\eta y - a(\eta)) dy - \int \frac{\partial}{\partial \eta} b(y) a'(\eta) \exp(\eta y - a(\eta)) dy \\
&= \mathbb{E}_y[y^2] - a'(\eta) \mathbb{E}[y] - \frac{\partial}{\partial \eta} \int b(y) a'(\eta) \exp(\eta y - a(\eta)) dy \\
&= \mathbb{E}[y^2] - \mathbb{E}[y]^2 - a''(\eta) = 0
\end{aligned}$$

so :  $\text{var}[y|\eta] = a''(\eta)$

(4.c)

$$NLL = -\log b(y) \exp(\eta y - a(\eta)) = -\log b(y) - (\eta y - a(\eta))$$

And  $\eta = \theta^T x$ :

$$\begin{aligned}
\frac{\partial NLL}{\partial \theta} &= \frac{\partial a(\eta)}{\partial \theta} - xy \\
&= \frac{\partial a(\eta)}{\partial \eta} \frac{\partial \eta}{\partial \theta} - xy \\
&= a'(\eta) x - xy
\end{aligned}$$

Then:

$$\begin{aligned}
\frac{\partial NLL}{\partial^2 \theta} &= \frac{\partial a'(\eta)}{\partial \eta} \frac{\partial \eta}{\partial \theta} = x^2 a''(\eta) \\
&= x^2 [\mathbb{E}[y^2] - \mathbb{E}[y]^2] = x^2 \sigma^2 \geq 0
\end{aligned}$$

(5.a)(i)

Let  $W_{ii} = \frac{1}{2} w^{(i)}$ ,  $W_{ij} = 0$  for  $i \neq j$ , let  $\vec{z} = X\theta - \vec{y}$ , i.e.  $z_i = \theta^T x^{(i)} - y^{(i)}$  Then we have:

$$\begin{aligned}
(X\theta - \vec{y})^T W (X\theta - \vec{y}) &= \vec{z}^T W \vec{z} \\
&= \frac{1}{2} \sum_{i=1}^m w^{(i)} z_i^2 \\
&= \frac{1}{2} \sum_{i=1}^m w^{(i)} (\theta^T x^{(i)} - y^{(i)})^2 \\
&= J(\theta)
\end{aligned}$$

(5.a)(ii)

$$\nabla_{\theta} J(\theta) = \nabla_{\theta} (\theta^T X^T W X \theta + \vec{y}^T W \vec{y} - 2 \vec{y}^T W X \theta) = 2 (X^T W X \theta - X^T W \vec{y})$$

so we have  $\nabla_{\theta} J(\theta) = 0$  if and only if

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

These are the normal equations, from which we can get a closed form formula for  $\theta$

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

(5.a)(iii)

$$\begin{aligned} \arg \max_{\theta} \prod_{i=1}^m p(y^{(i)} | x^{(i)}; \theta) &= \arg \max_{\theta} \sum_{i=1}^m \log p(y^{(i)} | x^{(i)}; \theta) \\ &= \arg \max_{\theta} \sum_{i=1}^m \left( \log \frac{1}{\sqrt{2\pi}\sigma^{(i)}} - \frac{(y^{(i)} - \theta^T x^{(i)})^2}{2(\sigma^{(i)})^2} \right) \\ &= \arg \max_{\theta} - \sum_{i=1}^m \frac{(y^{(i)} - \theta^T x^{(i)})^2}{2(\sigma^{(i)})^2} \\ &= \arg \min_{\theta} \frac{1}{2} \sum_{i=1}^m \frac{1}{(\sigma^{(i)})^2} (y^{(i)} - \theta^T x^{(i)})^2 \\ &= \arg \min_{\theta} \frac{1}{2} \sum_{i=1}^m w^{(i)} (y^{(i)} - \theta^T x^{(i)})^2 \end{aligned}$$

where in the last step, we substituted:  $w^{(i)} = \frac{1}{(\sigma^{(i)})^2}$  to get the linear regression form.