

1(a).

$1/\alpha$ .

$$l(w) = \sum_{i=1}^N y^{(i)} \log h(x^{(i)}) + (1 - y^{(i)}) \log (1 - h(x^{(i)}))$$

$$\nabla_w l = \sum_{i=1}^N (y^{(i)} - \sigma(w^T x^{(i)})) x^{(i)}$$

$$\frac{\partial l}{\partial w} = \sum_{i=1}^N -\sigma(w^T x^{(i)}) (1 - \sigma(w^T x^{(i)})) x^{(i)}, \text{ since } \sigma' = \sigma(1 - \sigma)$$

$$= - \sum_{i=1}^N x^{(i)} \sigma(w^T x^{(i)}) (1 - \sigma(w^T x^{(i)})) x^{(i)}$$

So the Hessian matrix can be represented as:

$$H = -X^T \Sigma X, \quad \Sigma = \text{diag}(\sigma^{(n)}(1 - \sigma^{(n)})) = \begin{pmatrix} \sigma(w^T x^{(1)}) (1 - \sigma(w^T x^{(1)})) & 0 & \dots & 0 \\ 0 & \ddots & \ddots & \vdots \\ 0 & \dots & \sigma(w^T x^{(n)}) (1 - \sigma(w^T x^{(n)})) \end{pmatrix}$$

$$\forall z, z^T H z = -z^T X^T \Sigma X z$$

$$= -(Xz)^T \Sigma (Xz) = - \sum_{i=1}^N \sigma^{(i)} (1 - \sigma^{(i)}) (X^{(i)} z)^2$$

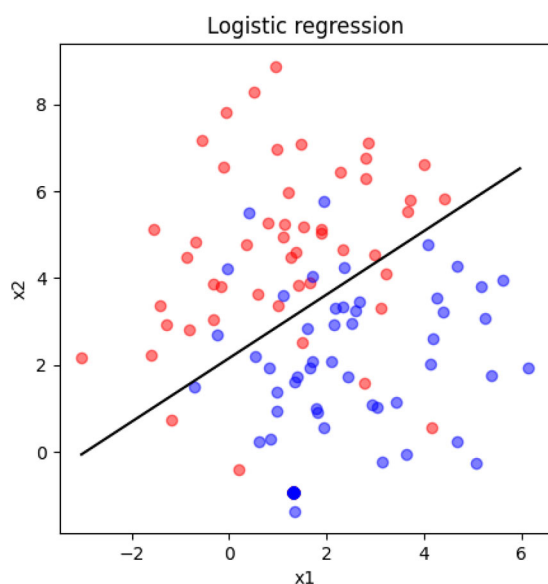
Since  $\sigma(x) = \frac{1}{1 + \exp(-x)} \in (0, 1)$ ,  $z^T H z \leq 0$ , negative semi-definite.

1(b).

The coefficients  $w$  resulting from my fit is [-1.84922892, -0.62814188, 0.85846843]

1(c).

The figure below shows the training data and decision boundary. The red dots represent the true samples, and the blue dots represent the negative samples.



2(a).

$$\begin{aligned} \frac{2}{(a)} \quad \nabla_{w_m} \ell(w) &= \nabla_{w_m} \sum_{i=1}^N \sum_{k=1}^K \log[p(y^{(i)}=k | x^{(i)}, w)]^{I(y^{(i)}=k)} \\ &= \nabla_{w_m} \sum_{i=1}^N I(y^{(i)}=m) \log[p(y^{(i)}=m | x^{(i)}, w)] + \nabla_{w_m} \sum_{i=1}^N \sum_{\substack{k=1 \\ k \neq m}}^K I(y^{(i)}=k) \log[p(y^{(i)}=k | x^{(i)}, w)] \\ &\quad + \nabla_{w_m} \sum_{i=1}^N I(y^{(i)}=K) \log[p(y^{(i)}=K | x^{(i)}, w)] \end{aligned}$$

only  $I(y^{(i)}=m)$  can be non negative.

$$\begin{aligned} \text{So } \nabla_{w_m} \ell(w) &= \sum_{i=1}^N I(y^{(i)}=m) \nabla_{w_m} \log \frac{\exp(w_m^T \phi(x^{(i)}))}{1 + \sum_{j=1}^{K-1} \exp(w_j^T \phi(x^{(i)}))} \quad \leftarrow \nabla_{w_m} \exp(w_m^T \phi(x^{(i)})) = \phi(x^{(i)}) \exp(w_m^T \phi(x^{(i)})) \\ &= \sum_{i=1}^N I(y^{(i)}=m) \phi(x^{(i)}) \left(1 - \frac{\exp(w_m^T \phi(x^{(i)}))}{1 + \sum_{j=1}^{K-1} \exp(w_j^T \phi(x^{(i)}))}\right) \\ &= \sum_{i=1}^N \phi(x^{(i)}) [I(y^{(i)}=m) - p(y^{(i)}=m | x^{(i)}, w)] \end{aligned}$$

2(b).

The SoftMax regression model implemented using the results in part (a) has an accuracy of 94.0 %

The logistic regression function in sklearn got 92% accuracy.

3(a).

$$\begin{aligned} \frac{3}{(a)} \quad p(y^{(i)}=0 | x^{(i)}) &= p(x^{(i)} | y^{(i)}=0) p(y^{(i)}=0) = \frac{1-\phi}{(2\pi)^{\frac{D}{2}} |\Sigma|^{\frac{1}{2}}} \exp\left(-\frac{1}{2}(x^{(i)}-\mu_0)^T \Sigma^{-1} (x^{(i)}-\mu_0)\right) \\ p(y^{(i)}=1 | x^{(i)}) &= p(x^{(i)} | y^{(i)}=1) p(y^{(i)}=1) = \frac{\phi}{(2\pi)^{\frac{D}{2}} |\Sigma|^{\frac{1}{2}}} \exp\left(-\frac{1}{2}(x^{(i)}-\mu_1)^T \Sigma^{-1} (x^{(i)}-\mu_1)\right) \end{aligned}$$

$$\text{Since } \sigma(\alpha) = \frac{1}{1 + \exp(-\alpha)}, \quad \alpha = \ln\left(\frac{\sigma}{1-\sigma}\right)$$

$$\Rightarrow \alpha = \ln \frac{p(y^{(i)}=1 | x^{(i)})}{p(y^{(i)}=0 | x^{(i)})} = \ln \frac{\exp(-\frac{1}{2}(x^{(i)}-\mu_0)^T \Sigma^{-1} (x^{(i)}-\mu_0))}{\exp(-\frac{1}{2}(x^{(i)}-\mu_1)^T \Sigma^{-1} (x^{(i)}-\mu_1))} + \ln \frac{\phi}{1-\phi}$$

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

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

$$= W^T x + w_0,$$

$$W^T = (\mu_1 - \mu_0)^T \Sigma^{-1}, \quad W_0 = -\frac{1}{2} \mu_1^T \Sigma^{-1} \mu_1 + \frac{1}{2} \mu_0^T \Sigma^{-1} \mu_0 + \ln \frac{\phi}{1-\phi}$$

refine  $x$  to  $M+1$  dim by adding  $x_0=1$ , then the above can be written as

$$W^T x + W_0 x_0 = W^T x_{\text{new}}, \quad x_{\text{new}} \in \mathbb{R}^{(M+1) \times n}$$

$$\Rightarrow p(y=1 | x, \phi, \Sigma, \mu_0, \mu_1) = \frac{1}{1 + \exp(-\alpha)} = \frac{1}{1 + \exp(-W^T x_{\text{new}})}$$

3(b).

$$\begin{aligned} \text{3/cb)} \quad \ell(\mu_0, \mu_1, \Sigma) &= \log \prod_{i=1}^N p(x^{(i)} | y^{(i)}; \mu_0, \mu_1, \Sigma) p(y^{(i)}; \phi) \\ &= \sum_{i=1}^N -\log(\sigma(2\pi)^{\frac{1}{2}}) - \frac{1}{2} (x^{(i)} - \mu_{y^{(i)}})^T \Sigma^{-1} (x^{(i)} - \mu_{y^{(i)}}) + y^{(i)} \log \phi + (1-y^{(i)}) \log(1-\phi) \end{aligned}$$

$$\frac{\partial \ell}{\partial \phi} = 0 \Rightarrow \sum_{i=1}^N \left( \frac{y^{(i)}}{\phi} + \frac{1-y^{(i)}}{\phi-1} \right) = \sum_{i=1}^N \frac{\phi - y^{(i)}}{\phi(\phi-1)} = 0 \Rightarrow \phi = \frac{1}{N} \sum_{i=1}^N y^{(i)} = \frac{1}{N} \sum_{i=1}^N 1\{y^{(i)}=1\}$$

$$\frac{\partial \ell}{\partial \mu_0} = 0 \Rightarrow \sum_{i=1}^N \Sigma^{-1} (x^{(i)} - \mu_0) 1\{y^{(i)}=0\} = 0 \Rightarrow \mu_0 = \frac{\sum_{i=1}^N 1\{y^{(i)}=0\} x^{(i)}}{\sum_{i=1}^N 1\{y^{(i)}=0\}}$$

$$\frac{\partial \ell}{\partial \mu_1} = 0 \Rightarrow \sum_{i=1}^N \Sigma^{-1} (x^{(i)} - \mu_1) 1\{y^{(i)}=1\} = 0 \Rightarrow \mu_1 = \frac{\sum_{i=1}^N 1\{y^{(i)}=1\} x^{(i)}}{\sum_{i=1}^N 1\{y^{(i)}=1\}}$$

$$\begin{aligned} \frac{\partial \ell}{\partial \Sigma} = 0 &\Rightarrow \text{let } S = \Sigma^{-1}, \quad -\frac{1}{2} \sum_{i=1}^N \mathbb{E}[-\log |S| + (x^{(i)} - \mu_{y^{(i)}})^T S (x^{(i)} - \mu_{y^{(i)}})] = 0 \\ &\Rightarrow -\frac{1}{2} \sum_{i=1}^N (-S^{-1} + (x^{(i)} - \mu_{y^{(i)}})(x^{(i)} - \mu_{y^{(i)}})^T) \\ &= \sum_{i=1}^N \frac{1}{2} S = -\frac{1}{2} \sum_{i=1}^N (x^{(i)} - \mu_{y^{(i)}})(x^{(i)} - \mu_{y^{(i)}})^T = 0 \\ &\Rightarrow \frac{1}{2} N S = \frac{1}{2} \sum_{i=1}^N (x^{(i)} - \mu_{y^{(i)}})(x^{(i)} - \mu_{y^{(i)}})^T \\ &\Rightarrow \Sigma = \frac{1}{N} \sum_{i=1}^N (x^{(i)} - \mu_{y^{(i)}})(x^{(i)} - \mu_{y^{(i)}})^T \end{aligned}$$

3(c).

$$\text{3/c).} \quad \phi \text{ remains the same as (b): } \phi = \frac{1}{N} \sum_{i=1}^N 1\{y^{(i)}=1\}.$$

$$\begin{aligned} \frac{d\ell}{d\mu_0} &= \sum_{i=1}^N -\frac{1}{2} \frac{d}{d\mu_0} [(x^{(i)} - \mu_0)^T \Sigma^{-1} (x^{(i)} - \mu_0)] 1\{y^{(i)}=0\} \\ &= \sum_{i=1}^N -\Sigma^{-1} ((x^{(i)} - \mu_0) 1\{y^{(i)}=0\}) \end{aligned}$$

$$\begin{aligned} \frac{d\ell}{d\mu_1} &= \sum_{i=1}^N -\frac{1}{2} \frac{d}{d\mu_1} [(x^{(i)} - \mu_1)^T \Sigma^{-1} (x^{(i)} - \mu_1)] 1\{y^{(i)}=1\} \\ &= \sum_{i=1}^N -\Sigma^{-1} ((x^{(i)} - \mu_1) 1\{y^{(i)}=1\}) \end{aligned}$$

$$\Rightarrow \mu_0 = \frac{\sum_{i=1}^N 1\{y^{(i)}=0\} x^{(i)}}{\sum_{i=1}^N 1\{y^{(i)}=0\}}, \quad \mu_1 = \frac{\sum_{i=1}^N 1\{y^{(i)}=1\} x^{(i)}}{\sum_{i=1}^N 1\{y^{(i)}=1\}}$$

**4(a).**

The Naive Bayes model has a test error rate of 1.625 %

**4(b).**

Using the model trained in (a), the top 5 tags that are most indicative of spam categories are "httpaddr", "spam", "unsubscribe", "ebay", "valet"

**4(c).**

Below is a plot of the test error with respect to the training set size. It can be seen that the training set containing 1400 texts gives the smallest test error of 1.6250%

