

I) Logistic Regression

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

where $h(x) = \sigma(\omega^T x) = \frac{1}{1 + \exp(-\omega^T x)}$

$$1-h(x) = \frac{\exp(-\omega^T x)}{1 + \exp(-\omega^T x)}$$

(a) Hessian of $l(\omega)$ $\nabla_\omega h(\omega) = h(x)(1-h(x))x$
 \Rightarrow First calculate the gradient,

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

$$= \sum_{i=1}^N y^{(i)} \frac{\nabla_\omega h(x^{(i)})}{h(x^{(i)})} + (1-y^{(i)}) \frac{-\nabla_\omega h(x^{(i)})}{1-h(x^{(i)})}$$

$$= \sum_{i=1}^N y^{(i)} \frac{h(x^{(i)}) (1-h(x^{(i)})) x^{(i)}}{h(x^{(i)})} - \frac{(1-y^{(i)}) h(x^{(i)}) (1-h(x^{(i)})) x^{(i)}}{1-h(x^{(i)})}$$

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

$$= \sum_{i=1}^N [y^{(i)} - y^{(i)} h(x^{(i)}) - h(x^{(i)}) + y^{(i)} h(x^{(i)})] x^{(i)}$$

$\nabla_\omega l(\omega) = \sum_{i=1}^N [y^{(i)} - h(x^{(i)})] x^{(i)}$

Now, calculating the Hessian

$$\nabla_w^2 \ell(\omega) = \sum_{i=1}^n [-\nabla_\omega h(x^{(i)})] x^{(i)}$$

$$\boxed{\nabla_w^2 \ell(\omega) = - \sum_{i=1}^n h(x^{(i)}) (1-h(x^{(i)})) x^{(i)} (x^{(i)})^T}$$

(b) show $\nabla^T H z \leq 0$ | let $x^{(i)} \in \mathbb{R}^n$

$$H = - \sum_{i=1}^n h(x^{(i)}) (1-h(x^{(i)})) x^{(i)} (x^{(i)})^T$$

* $h(x^{(i)}) = \frac{1}{1 + \exp(-\omega^T x^{(i)})}$, where $0 < h(x^{(i)}) < 1$

so $h(x^{(i)}) (1-h(x^{(i)})) > 0$

for any vector $z_j \in \mathbb{R}^{n \times 1}$ | $H \in \mathbb{R}^{n \times n}$

$$\begin{aligned} \nabla^T H z &= - \sum_{i=1}^n h(x^{(i)}) (1-h(x^{(i)})) \underbrace{z^T x^{(i)}}_{\substack{1 \times n \\ 1 \times 1}} \underbrace{x^{(i)T} z}_{\substack{n \times 1 \\ 1 \times 1}} \\ &= - \sum_{i=1}^n h(x^{(i)}) (1-h(x^{(i)})) \underbrace{\underbrace{z^T x^{(i)}}_{1 \times n} \underbrace{(z^T x^{(i)})^T}_{n \times 1}}_{1 \times 1} \\ &= - \sum_{i=1}^n h(x^{(i)}) (1-h(x^{(i)})) (\underbrace{z^T x^{(i)}}_{\geq 0})^2 \geq 0 \end{aligned}$$

Hence,

$$\sum_i H_{ii} z_i^2 \leq 0, \text{ hence}$$

H is semi-definite $\Rightarrow l$ is concave
and has no local minima other than
global one.

(c) update rule implied by Newton's
method

update rule:

$$w \leftarrow w - H^{-1} \nabla_w l(w)$$

where

$$H = \left[\sum_{i=1}^n h(x^{(i)}) (1 - h(x^{(i)})) x^{(i)} (x^{(i)})^T \right]$$

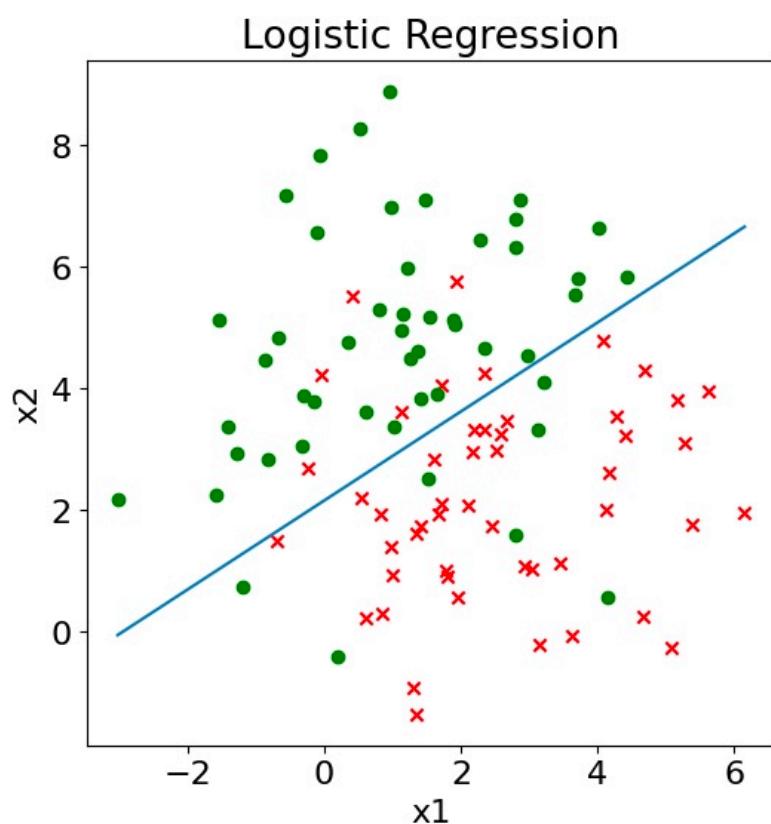
$$\nabla_w l(w) = \sum_{i=1}^n [y^{(i)} - h(x^{(i)})] x^{(i)}$$

(d) Autograder

(e)

$$\omega = [-1.84922892, -0.62814188, 0.85846843]$$
$$\text{Difference} = 1.6653345 \times 10^{-15}$$

(f)



② Softmax Regression via Gradient Ascent.

(a)

$$l(\omega) = \log(L(\omega)) = \sum_{i=1}^n \sum_{k=1}^K \log \left(\left[p(y^{(i)} = k | x^{(i)}, \omega) \right] \right)^{\prod \{y^{(i)} = k\}}$$

Case 1:

considering $y^{(i)} = k = m$,

$$\begin{aligned} \nabla_{w_m} l(\omega) &= \sum_{i=1}^n \prod \{y^{(i)} = m\} \nabla_{w_m} \log(p(y^{(i)} = m | x^{(i)}, \omega)) \\ &= \sum_{i=1}^n \prod \{y^{(i)} = m\} \nabla_{w_m} \left(\log \left(\frac{\exp(\omega_m^\top \phi(x^{(i)}))}{1 + \sum_{j=1}^{K-1} \exp(\omega_j^\top \phi(x^{(i)}))} \right) \right) \\ &= \sum_{i=1}^n \prod \{y^{(i)} = m\} \frac{\cancel{1 + \sum_{j=1}^{K-1} \exp(\omega_j^\top \phi(x^{(i)}))}}{\cancel{\exp(\omega_m^\top \phi(x^{(i)}))}} \times \left(\frac{\cancel{\phi(x^{(i)}) \exp(\omega_m^\top \phi(x^{(i)}))} \left(1 + \sum_{j=1}^{K-1} \exp(\omega_j^\top \phi(x^{(i)})) - \exp(\omega_m^\top \phi(x^{(i)})) \right) \cancel{\exp(\omega_m^\top \phi(x^{(i)}))}}{\cancel{\left(1 + \sum_{j=1}^{K-1} \exp(\omega_j^\top \phi(x^{(i)})) \right)}} \right) \\ &= \sum_{i=1}^n \prod \{y^{(i)} = m\} \phi(x^{(i)}) \left(\frac{1 + \sum_{j=1}^{K-1} \exp(\omega_j^\top \phi(x^{(i)})) - \exp(\omega_m^\top \phi(x^{(i)}))}{1 + \sum_{j=1}^{K-1} \exp(\omega_j^\top \phi(x^{(i)}))} \right) \\ &= \sum_{i=1}^n \phi(x^{(i)}) \left(\prod \{y^{(i)} = m\} - \frac{\prod \{y^{(i)} = m\} \exp(\omega_m^\top \phi(x^{(i)}))}{1 + \sum_{j=1}^{K-1} \exp(\omega_j^\top \phi(x^{(i)}))} \right) \\ \nabla_{w_m} l(\omega) &= \sum_{i=1}^n \phi(x^{(i)}) \left(\prod \{y^{(i)} = m\} - \frac{\exp(\omega_m^\top \phi(x^{(i)}))}{1 + \sum_{j=1}^{K-1} \exp(\omega_j^\top \phi(x^{(i)}))} \right) \end{aligned}$$

$$\nabla_{w_m} l(\omega) = \sum_{i=1}^n \phi(x^{(i)}) \left[\mathbb{I}\{y^{(i)}=m\} - p(y^{(i)}=m|x^{(i)}, \omega) \right]$$

case 2: considering $y^{(i)} = k \neq m$,

$$\begin{aligned} \nabla_{w_k} l(\omega) &= \nabla_{w_k} \left[\sum_{i=1}^n \mathbb{I}(y^{(i)}=k) \log \left(p(y^{(i)}=k|x^{(i)}, \omega) \right) \right] \\ &= \nabla_{w_k} \left[\sum_{i=1}^n \mathbb{I}(y^{(i)}=k) \times \log \left(\frac{1}{1 + \sum_{j=1}^{K-1} \exp(\omega_j^\top \phi(x))} \right) \right] \end{aligned}$$

$$\nabla_{w_k} l(\omega) = \textcircled{O}$$

no w_k term
so, $\rightarrow 0$

Hence,

$$\nabla_{w_m} l(\omega) = \sum_{i=1}^n \phi(x^{(i)}) \left[\mathbb{I}\{y^{(i)}=m\} - p(y^{(i)}=m|x^{(i)}, \omega) \right]$$

for $m = 1, \dots, K-1$

(b) Auto grader

(c) Accuracy on the test data
is 94 %

(3)

Gaussian Discriminate Analysis

a

$$P(y=1 | x) = \frac{P(x | y=1) P(y=1)}{P(x)}$$

$$= \frac{P(x | y=1) P(y=1)}{P(x | y=0) P(y=0) + P(x | y=1) P(y=1)}$$

$$= \frac{\frac{1}{(2\pi)^{N/2} |\Sigma|^{\frac{1}{2}}} \exp\left(-\frac{1}{2}(x^{(i)} - \mu_1)^T \Sigma^{-1} (x^{(i)} - \mu_1)\right) \phi}{\frac{1}{(2\pi)^{N/2} |\Sigma|^{\frac{1}{2}}} \exp\left(-\frac{1}{2}(x^{(i)} - \mu_0)^T \Sigma^{-1} (x^{(i)} - \mu_0)\right) (1-\phi) + \frac{1}{(2\pi)^{N/2} |\Sigma|^{\frac{1}{2}}} \exp\left(-\frac{1}{2}(x^{(i)} - \mu_1)^T \Sigma^{-1} (x^{(i)} - \mu_1)\right) (\phi)}$$

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

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

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

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

$$= \frac{1}{1 + \exp\left(\log\left(\frac{1-\phi}{\phi}\right) + \mathbf{x}^{(i)\top} \Sigma^{-1} (\boldsymbol{\mu}_0 - \boldsymbol{\mu}_1) - \frac{1}{2} \boldsymbol{\mu}_0^\top \Sigma \boldsymbol{\mu}_0 + \frac{1}{2} \boldsymbol{\mu}_1^\top \Sigma \boldsymbol{\mu}_1\right)}$$

$$P(y=1 | x; \phi, \Sigma, \boldsymbol{\mu}_0, \boldsymbol{\mu}_1) = \frac{1}{1 + \exp(-w^\top \hat{\mathbf{x}})}$$

where

$$\hat{\mathbf{X}} = \begin{bmatrix} & \\ & \\ & \mathbf{x}^{(i)\top} \\ & \end{bmatrix}_{(M+1 \times 1)} \quad \text{and}$$

$$\omega = \begin{bmatrix} \log\left(\frac{1-\phi}{\phi}\right) + \frac{1}{2} \boldsymbol{\mu}_0^\top \Sigma \boldsymbol{\mu}_0 - \frac{1}{2} \boldsymbol{\mu}_1^\top \Sigma \boldsymbol{\mu}_1 \\ - \Sigma^{-1} (\boldsymbol{\mu}_0 - \boldsymbol{\mu}_1) \end{bmatrix}$$

(b)

$$l(\phi, \boldsymbol{\mu}_0, \boldsymbol{\mu}_1, \Sigma) = \log \prod_{i=1}^N p(x^{(i)}, y^{(i)}; \phi, \boldsymbol{\mu}_0, \boldsymbol{\mu}_1, \Sigma)$$

$$= \log \prod_{i=1}^N p(x^{(i)} | y^{(i)}; \phi, \boldsymbol{\mu}_0, \boldsymbol{\mu}_1, \Sigma) p(y^{(i)}; \phi)$$

$$= \sum_{i=1}^N \log(p(x^{(i)} | y^{(i)}; \phi, \boldsymbol{\mu}_0, \boldsymbol{\mu}_1, \Sigma)) + \log(p(y^{(i)}; \phi))$$

$$= \sum_{i=1}^N \log \left(\frac{1}{(2\pi)^{M/2} |\Sigma|^{\frac{1}{2}}} \exp \left(-\frac{1}{2} (x^{(i)} - \mu_{y(i)})^T \Sigma^{-1} (x^{(i)} - \mu_{y(i)}) \right) \right) + \log \left(\frac{\phi^{y(i)}}{(1-\phi)^{1-y(i)}} \right)$$

$$= \sum_{i=1}^N -\log \left((2\pi)^{M/2} |\Sigma|^{\frac{1}{2}} \right) + \left(-\frac{1}{2} (x^{(i)} - \mu_{y(i)})^T \Sigma^{-1} (x^{(i)} - \mu_{y(i)}) \right) + y^{(i)} \log(\phi) + (1-y^{(i)}) \log(1-\phi)$$

$$\ell(\phi, \mu_0, \mu_1, \Sigma) = \sum_{i=1}^N -\frac{M}{2} \log(2\pi) - \frac{1}{2} \log(|\Sigma|) - \frac{1}{2} (x^{(i)} - \mu_{y(i)})^T \Sigma^{-1} (x^{(i)} - \mu_{y(i)}) + y^{(i)} \log(\phi) + (1-y^{(i)}) \log(1-\phi)$$

$$\nabla_\phi \ell(\phi, \mu_0, \mu_1, \Sigma) = \sum_{i=1}^N \left(\frac{y^{(i)}}{\phi} + \frac{(1-y^{(i)})}{1-\phi} x - 1 \right)$$

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

$$0 = \sum_{i=1}^N \frac{\mathbb{I}\{y^{(i)}=1\}}{\phi} + \frac{\mathbb{I}\{y^{(i)}=0\}}{1-\phi} - \frac{N}{1-\phi}$$

$$\Rightarrow 0 = \sum_{i=1}^N (\mathbb{I}\{y^{(i)}=1\}(1-\phi) + \mathbb{I}\{y^{(i)}=0\}\phi) - N\phi$$

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

Hence proved

$$\nabla_{\mu_0} \ell(\phi, \mu_0, \mu_1, \Sigma) = \sum_{i=1}^N \nabla_{\mu_0} \left(\frac{1}{2} (x^{(i)} - \mu_{y(i)})^T \Sigma^{-1} (x^{(i)} - \mu_{y(i)}) \right)$$

$$= \sum_{i=1}^N -\frac{1}{2} \nabla_{\mu_0} \left[x^{(i)T} \Sigma^{-1} x^{(i)} - 2 \mu_{y(i)}^T \Sigma^{-1} x^{(i)} + \mu_{y(i)}^T \Sigma^{-1} \mu_{y(i)} \right]$$

$$= \binom{1}{2} \sum_{i=1}^2 \left[-2 \sum_{j=1}^1 x^{(i)} \mathbb{I}\{y^{(i)} = 0\} + 2 \sum_{j=1}^1 x_0 \mathbb{I}\{y^{(i)} = 0\} \right]$$

$$0 = \sum_{i=1}^2 \mathbb{I}\{y^{(i)} = 0\} \sum_{j=1}^1 x^{(i)} - \sum_{i=1}^2 \mathbb{I}\{y^{(i)} = 0\} \sum_{j=1}^1 x_0$$

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

Hence Proved

$$\nabla_{\mu_1} \ell(\phi, \mu_0, \mu_1, \Sigma) = \sum_{i=1}^N \nabla_{\mu_1} \left(-\frac{1}{2} (x^{(i)} - \mu_{y^{(i)}})^T \sum_{j=1}^1 (x^{(j)} - \mu_{y^{(j)}}) \right)$$

$$= \sum_{i=1}^N \frac{1}{2} \nabla_{\mu_1} \left[x^{(i)T} \sum_{j=1}^1 x^{(j)} - 2 \mu_{y^{(i)}}^T \sum_{j=1}^1 x^{(j)} + \mu_{y^{(i)}}^T \sum_{j=1}^1 \mu_{y^{(j)}} \right]$$

$$= \binom{1}{2} \sum_{i=1}^2 \left[-2 \sum_{j=1}^1 x^{(i)} \mathbb{I}\{y^{(i)} = 1\} + 2 \sum_{j=1}^1 x_i \mathbb{I}\{y^{(i)} = 1\} \right]$$

$$0 = \sum_{i=1}^2 \mathbb{I}\{y^{(i)} = 1\} \sum_{j=1}^1 x^{(i)} - \sum_{i=1}^2 \mathbb{I}\{y^{(i)} = 1\} \sum_{j=1}^1 x_i$$

$$\mu_1 = \frac{\sum_{i=1}^2 \mathbb{I}\{y^{(i)} = 1\} x^{(i)}}{\sum_{i=1}^2 \mathbb{I}\{y^{(i)} = 1\}}$$

Hence Proved

(G) For $M=1$, $\Sigma = [\sigma^2]$, $|\Sigma| = \sigma^2$

$$l(\phi, \mu_0, \Sigma) = \sum_{i=1}^n -\frac{M}{2} \log(2\pi) - \frac{1}{2} \log(|\Sigma|) - \frac{1}{2} (x^{(i)} - \mu_{y^{(i)}})^T \Sigma^{-1} (x^{(i)} - \mu_{y^{(i)}}) + y^{(i)} \log(\phi) + (1-y^{(i)}) \log(1-\phi)$$

$$\nabla_{\Sigma} l(\phi, \mu_0, \Sigma) = \sum_{i=1}^n \left(\frac{-1}{2} \log(|\Sigma|) - \frac{1}{2} (x^{(i)} - \mu_{y^{(i)}})^T \Sigma^{-1} (x^{(i)} - \mu_{y^{(i)}}) \right)$$

$$0 = \sum_{i=1}^n -\frac{1}{2} \frac{1}{|\Sigma|} - \frac{1}{2} \frac{1}{\Sigma^2} (x^{(i)} - \mu_{y^{(i)}})(x^{(i)} - \mu_{y^{(i)}})^T$$

$$0 = \sum_{i=1}^n -\frac{1}{|\Sigma|} + \frac{1}{\Sigma^2} \sum_{i=1}^n (x^{(i)} - \mu_{y^{(i)}})^T (x^{(i)} - \mu_{y^{(i)}})$$

$$\frac{2}{|\Sigma|} = \sum_{i=1}^n (x^{(i)} - \mu_{y^{(i)}})^T (x^{(i)} - \mu_{y^{(i)}})$$

$$\boxed{\Sigma = \frac{1}{n} \sum_{i=1}^n (x^{(i)} - \mu_{y^{(i)}})^T (x^{(i)} - \mu_{y^{(i)}})}$$

Hence Proved

y) Naive Bayes for classifying SPAM

(a) Objective is to compute, $P(u|x, \gamma)$

For MAP, we want posterior distribution

$$P(u|x, \gamma)$$

using Bayes formula, $P(u|x, \gamma) = \frac{P(x, \gamma|u) P(u)}{P(x, \gamma)}$

①

$P(x, \gamma|u)$ is a likelihood,

$$P(x, \gamma|u) = \prod_{i=1}^K \prod_{j=1}^M (\phi_i)^{N_i} (u_j^i)^{N_j^i} \quad \text{--- (ii)}$$

$P(u)$ is a prior,

$$P(u) = \prod_{i=1}^K \prod_{j=1}^M (u_j^i)^\alpha \quad \text{--- (iii)}$$

Taking log of eq ① and substituting ② and ③

$$\begin{aligned} \log(P(u|x, \gamma)) &= \sum_{i=1}^K \sum_{j=1}^M N_i \log(\phi_i) + N_j^i \log(u_j^i) \\ &\quad + \sum_{i=1}^K \sum_{j=1}^M \alpha \log(u_j^i) - \log(P(x, \gamma)) \\ &\quad - \log(z) \end{aligned}$$

$$\log(P(u|X, Y)) = \sum_{i=1}^K \sum_{j=1}^M N_i \log(\phi_i) + (N_j^i + \alpha) \log(m_j^i)$$

$\rightarrow \log(P(X, Y))$
 $\rightarrow \log(z)$

→ This is similar to log-likelihood derived in the lecture notes.

This can be further simplified by removing dependent terms (as done in the lecture),

$$\log(P(u|X, Y)) = \sum_{i=1}^K \sum_{j=1}^M N_i \log(\phi_i) + \sum_{i=1}^K \left[\sum_{j=1}^{M-1} (N_j^i + \alpha) \log(m_j^i) + (N_M^i + \alpha) \log\left(1 - \sum_{j=1}^{M-1} m_j^i\right) \right]$$

$\rightarrow \log(P(X, Y)) - \log(z)$

Taking partial derivative w.r.t m_j^i , and equating to zero.

$$0 = \frac{N_j^i + \alpha}{m_j^i} - \frac{N_M^i + \alpha}{1 - \sum_{j=1}^{M-1} m_j^i}$$

$$\Rightarrow \frac{N_j^i + \alpha}{m_j^i} = \text{constant } \forall j$$

$$\Rightarrow m_j^i = \frac{N_j^i + \alpha}{\sum_{j=1}^M N_j^i + \alpha M}$$

Hence Proved

(b)

(i) Auto grader

(ii) Top 5 most indicative tokens are: ['valet' 'ebai' 'unsubscrib'
'spam' 'httpaddr'].

(iii)

Accuracy for 50 mail data : 96.125%.

Accuracy for 100 mail data : 97.375%.

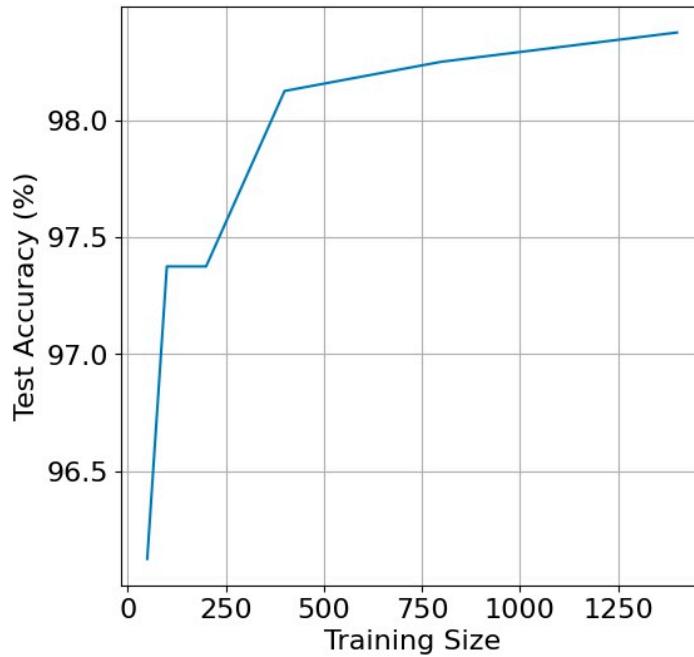
Accuracy for 200 mail data : 97.375%.

Accuracy for 400 mail data : 98.125%.

Accuracy for 800 mail data : 98.25%.

Accuracy for 1400 mail data : 98.375%.

(iv)



(v)

1400 Training set size gives the best accuracy.