

JAGDEESH BASAVARAJU

ASU ID: 1213004713

## Statistical Machine Learning Assignment #1

1. Given Prior,  $P(\theta)$  is Beta( $\beta_H, \beta_T$ ), Prove that posterior also follows Beta( $\alpha_H + \beta_H, \alpha_T + \beta_T$ ).

Sol<sup>n</sup>:

$$\text{Prior, } P(\theta) = \text{Beta}(\beta_H, \beta_T)$$

$$P(\theta) = \frac{\theta^{\beta_H-1} (1-\theta)^{\beta_T-1}}{B(\beta_H, \beta_T)} \quad \text{--- (1)}$$

$$\text{Posterior } P(\theta|D) \propto P(D|\theta) P(\theta) \quad \text{--- (2)}$$

Consider  $P(D|\theta) P(\theta)$

$$P(D|\theta) P(\theta) = \theta^{\alpha_H} (1-\theta)^{\alpha_T} \cdot \frac{\theta^{\beta_H-1} (1-\theta)^{\beta_T-1}}{B(\beta_H, \beta_T)} \quad [\because \text{from (1)}]$$

$$P(D|\theta) P(\theta) = \frac{1}{B(\beta_H, \beta_T)} \theta^{\alpha_H + \beta_H - 1} (1-\theta)^{\alpha_T + \beta_T - 1} \quad \text{--- (3)}$$

From (1) & (2)

$$P(\theta|D) \propto \frac{1}{C} P(D|\theta) P(\theta)$$

$$P(\theta|D) \propto \theta^{\alpha_H + \beta_H - 1} (1-\theta)^{\alpha_T + \beta_T - 1}$$

$$\therefore P(\theta|D) \sim \text{Beta}(\alpha_H + \beta_H, \alpha_T + \beta_T)$$

$$\text{Mean of posterior} = \text{mean of Beta}(\alpha_H + \beta_H, \alpha_T + \beta_T)$$

$$= \frac{\alpha_H + \beta_H}{\alpha_H + \beta_H + \alpha_T + \beta_T}$$

$$\text{Mode of posterior} = \hat{\theta}_{MAP} = \arg \max P(\theta|D) = \text{mode of beta}$$

$$= \frac{\alpha_H + \beta_H - 1}{\alpha_H + \beta_H + \alpha_T + \beta_T - 2}$$

## 2) Parameter Estimation

i) Prove that  $\hat{\mu}_{MLE}$  is unbiased.

$D = x_1, x_2, \dots, x_N \in R$  and is  $N(\mu, \sigma^2)$

$$\therefore P(D|\mu, \sigma) = \left( \frac{1}{\sigma \sqrt{2\pi}} \right)^N \prod_{i=1}^N \exp \left\{ -\frac{(x_i - \mu)^2}{2\sigma^2} \right\}$$

$$\ln P(D|\mu, \sigma) = -N \ln(\sigma \sqrt{2\pi}) - \sum_{i=1}^N \frac{(x_i - \mu)^2}{2\sigma^2}$$

To find  $\hat{\mu}_{MLE}$ ,  $\frac{d}{d\mu} \ln(P(D|\mu, \sigma)) = 0$

$$\begin{aligned} \frac{d}{d\mu} [\ln P(D|\mu, \sigma)] &= \frac{d}{d\mu} \left[ -N \ln(\sigma \sqrt{2\pi}) - \sum_{i=1}^N \frac{(x_i - \mu)^2}{2\sigma^2} \right] \\ &= -\frac{1}{2\sigma^2} \sum_{i=1}^N 2(x_i - \mu) \cdot (-1) = \frac{1}{\sigma^2} \sum_{i=1}^N (x_i - \mu) = 0 \end{aligned}$$

$$\therefore \sum_{i=1}^N x_i - \sum_{i=1}^N \mu = 0$$

$$\sum_{i=1}^N x_i = NM$$

$$\therefore \hat{\mu}_{MLE} = \frac{1}{N} \sum_{i=1}^N x_i$$

$$\begin{aligned} \text{Expected value of } \hat{\mu}_{MLE} &= E(\hat{\mu}_{MLE}) = \frac{1}{N} \sum_{i=1}^N E(x_i) \\ &= \frac{1}{N} \sum_{i=1}^N \mu = \frac{1}{N} \cdot NM \\ &= \mu \end{aligned}$$

$$\text{Bias} = E(\hat{\mu}_{MLE}) - \mu = \mu - \mu = 0$$

$\therefore \hat{\mu}_{MLE}$  is unbiased.

2 iii) Prove that  $\hat{\sigma}_{MLE}^2$  is biased.

$$\hat{\sigma}_{MLE}^2 = \frac{1}{N} \sum_{i=1}^N (x_i - \hat{\mu}_{MLE})^2$$

Estimate

$$\text{Expected value of } \hat{\sigma}_{MLE}^2 = E(\hat{\sigma}_{MLE}^2)$$

$$= \frac{1}{N} \sum_{i=1}^N E(x_i - \hat{\mu}_{MLE})^2$$

$E(x_i - \hat{\mu}_{MLE})^2$  can't be equal to  $\sigma^2$  as  $\hat{\mu}_{MLE}$  is just an estimate and we don't know true value.

$$\therefore \text{Bias} = E(\hat{\sigma}_{MLE}^2) - \sigma^2$$

$$= \left\{ \frac{1}{N} \sum_{i=1}^N E(x_i - \hat{\mu}_{MLE})^2 \right\} - \sigma^2$$

$$= \left\{ \frac{1}{N} \sum_{i=1}^N E(x_i^2 - 2x_i \hat{\mu}_{MLE} + \hat{\mu}_{MLE}^2) \right\} - \sigma^2$$

Bias  $\neq 0$

Hence  $\hat{\sigma}_{MLE}^2$  is biased.

$$= \left\{ \frac{1}{N} \sum_{i=1}^N (E(x_i^2) - 2E(x_i) \hat{\mu}_{MLE} + E(\hat{\mu}_{MLE}^2)) \right\} - \sigma^2$$

$$= \left\{ \frac{1}{N} \sum_{i=1}^N (E(x_i^2) - 2\hat{\mu}_{MLE} x_{MLE} + \hat{\mu}_{MLE}^2) \right\} - \sigma^2$$

$$= \left\{ \frac{1}{N} \sum_{i=1}^N E(x_i^2) - 2\hat{\mu}_{MLE}^2 + \hat{\mu}_{MLE}^2 \right\} - \sigma^2$$

$$= \left\{ \frac{1}{N} \sum_{i=1}^N E(x_i^2) - \hat{\mu}_{MLE}^2 \right\} - \{E(\hat{x}_i^2) - \mu^2\}$$

Bias  $\neq 0$

Hence  $\hat{\sigma}_{MLE}^2$  is biased

### 3) Naive Bayes Classifier

i) 13 independent parameters

Consider age as  $X_1$ , income as  $X_2$ , student as  $X_3$ , credit rating as  $X_4$  and buy-computer as  $Y$ .

- 1)  $P(Y = \text{yes})$
- 2)  $P(X_1 = \text{youth} | Y = \text{yes})$
- 3)  $P(X_1 = \text{middle-aged} | Y = \text{yes})$
- 4)  $P(X_2 = \text{high} | Y = \text{yes})$
- 5)  $P(X_2 = \text{medium} | Y = \text{yes})$
- 6)  $P(X_3 = \text{yes} | Y = \text{yes})$
- 7)  $P(X_4 = \text{fair} | Y = \text{yes})$
- 8)  $P(X_1 = \text{youth} | Y = \text{no})$
- 9)  $P(X_1 = \text{middle-aged} | Y = \text{no})$
- 10)  $P(X_2 = \text{high} | Y = \text{no})$
- 11)  $P(X_2 = \text{medium} | Y = \text{no})$
- 12)  $P(X_3 = \text{yes} | Y = \text{no})$
- 13)  $P(X_4 = \text{fair} | Y = \text{no})$

Since the class is binary, we need one parameter (1)

Age has 3 distinct values, so need 2 parameters ~~per class~~ (4)

Income also has 3 distinct values, so 2 per class (4)

Student has 2 distinct values, so 1 per class (2)

Credit-rating ~~also~~ has 2 distinct values, so 1 per class (2)

Hence totally 13 independent parameters

3 iii) Estimated value of all these parameters

Age :  $X_1$ , Income :  $X_2$ , Student :  $X_3$ , Credit rating :  $X_4$

Buyer-Computer :  $Y$

$$i) P(Y = \text{yes}) = \frac{9}{14} \therefore P(Y = \text{no}) = \frac{5}{14}$$

$$2) P(X_1 = \text{youth} | Y = \text{yes}) = \frac{2}{9}$$

$$3) P(X_1 = \text{middle-aged} | Y = \text{yes}) = \frac{4}{9} \therefore P(X_1 = \text{senior} | Y = \text{yes}) = \frac{3}{9}$$

$$4) P(X_2 = \text{high} | Y = \text{yes}) = \frac{2}{9}$$

$$5) P(X_2 = \text{medium} | Y = \text{yes}) = \frac{4}{9} \therefore P(X_2 = \text{low} | Y = \text{yes}) = \frac{3}{9}$$

$$6) P(X_3 = \text{yes} | Y = \text{yes}) = \frac{6}{9} \therefore P(X_3 = \text{no} | Y = \text{yes}) = \frac{3}{9}$$

$$7) P(X_4 = \text{fair} | Y = \text{yes}) = \frac{6}{9} \therefore P(X_4 = \text{excellent} | Y = \text{yes}) = \frac{3}{9}$$

$$8) P(X_1 = \text{youth} | Y = \text{no}) = \frac{3}{5}$$

$$9) P(X_1 = \text{middle-aged} | Y = \text{no}) = \frac{0}{5} = 0 \therefore P(X_1 = \text{senior} | Y = \text{no}) = \frac{2}{5}$$

$$10) P(X_2 = \text{high} | Y = \text{no}) = \frac{2}{5} \therefore P(X_2 = \text{low} | Y = \text{no}) = \frac{1}{5}$$

$$11) P(X_2 = \text{medium} | Y = \text{no}) = \frac{2}{5} \therefore P(X_3 = \text{no} | Y = \text{no}) = \frac{4}{5}$$

$$12) P(X_3 = \text{yes} | Y = \text{no}) = \frac{1}{5} \therefore P(X_4 = \text{excellent} | Y = \text{no}) = \frac{3}{5}$$

$$13) P(X_4 = \text{fair} | Y = \text{no}) = \frac{2}{5}$$

3 iii)  $X = (\text{youth}, \text{medium}, \text{yes}, \text{fair})$ , find  $P(Y=\text{yes} | X)$

$$P(Y=1 | X_1 \dots X_n) = \frac{P(Y=1) P(X_1 \dots X_n | Y=1)}{P(Y=1) P(X_1 \dots X_n | Y=1) + P(Y=0) P(X_1 \dots X_n | Y=0)}$$

$$= \frac{P(Y=1) P(X_1 | Y=1) \dots P(X_n | Y=1)}{P(Y=1) P(X_1 | Y=1) \dots P(X_n | Y=1) + P(Y=0) P(X_1 | Y=0) \dots P(X_n | Y=0)}$$

$$\therefore P(Y=\text{yes} | X) = P(Y=\text{yes}) P(X_1=\text{youth} | Y=\text{yes}) P(X_2=\text{medium} | Y=\text{yes}) P(X_3=\text{fair} | Y=\text{yes}) \dots P(X_n=\text{fair} | Y=\text{yes})$$

$$[P(Y=\text{yes}) P(X_1=\text{youth} | Y=\text{yes}) P(X_2=\text{medium} | Y=\text{yes}) P(X_3=\text{fair} | Y=\text{yes})]$$

$$+ P(Y=\text{no}) P(X_1=\text{youth} | Y=\text{no}) P(X_2=\text{medium} | Y=\text{no}) P(X_3=\text{fair} | Y=\text{no})]$$

$$\therefore P(Y=\text{yes} | X) = \frac{\left(\frac{9}{14}\right) \left(\frac{2}{9}\right) \left(\frac{4}{9}\right) \left(\frac{6}{9}\right) \left(\frac{6}{9}\right)}{\left[\left(\frac{9}{14}\right) \left(\frac{2}{9}\right) \left(\frac{4}{9}\right) \left(\frac{6}{9}\right) \left(\frac{6}{9}\right) + \left(\frac{5}{14}\right) \left(\frac{3}{5}\right) \left(\frac{2}{5}\right) \left(\frac{1}{5}\right) \left(\frac{2}{5}\right)]}$$

$$= \frac{0.0282}{0.0282 + 0.006857} = \frac{0.0282}{0.035057}$$

$$P(Y=\text{yes} | X) = 0.8044$$

Since it's greater than 0.5, the classifier has predicted  $y=\text{yes}$  for this person.  
i.e buyer-computer = yes

#### 4. Logistic Regression

$$x_1 = (1, 0) \quad x_2 = (0, -1) \quad x_3 = (0, 1) \quad x_4 = (-1, 0)$$

$$y_1 = 1 \quad y_2 = 1 \quad y_3 = 0 \quad y_4 = 0$$

i) Initial weight vector,  $\omega^{(0)} = (0, 0, 0)'$

$$\text{Iteration 1: } \hat{P}(y_j=1/x_j, \omega) = \frac{\exp(\omega_0^{(t)} + \sum_{i=1}^2 \omega_i^{(t)} x_i)}{1 + \exp(\omega_0^{(t)} + \sum_{i=1}^2 \omega_i^{(t)} x_i)}$$

$$\hat{P}(y_1=1/x_1, \omega^{(0)}) = \frac{\exp(0 + 0 \cdot 1 + 0 \cdot 0)}{1 + \exp(0 + 0 \cdot 1 + 0 \cdot 0)} = \frac{1}{2} = 0.5$$

$$\hat{P}(y_2=1/x_2, \omega^{(0)}) = \frac{\exp(0 + 0 \cdot 0 + 0 \cdot (-1))}{1 + \exp(0 + 0 \cdot 0 + 0 \cdot (-1))} = \frac{1}{2} = 0.5$$

$$\hat{P}(y_3=1/x_3, \omega^{(0)}) = \frac{\exp(0 + 0 \cdot 0 + 0 \cdot 1)}{1 + \exp(0 + 0 \cdot 0 + 0 \cdot 1)} = \frac{1}{2} = 0.5$$

$$\hat{P}(y_4=1/x_4, \omega^{(0)}) = \frac{\exp(0 + 0 \cdot (-1) + 0 \cdot 0)}{1 + \exp(0 + 0 \cdot (-1) + 0 \cdot 0)} = \frac{1}{2} = 0.5$$

$$\omega_0^{(t+1)} = \omega_0^{(t)} + \eta \sum [y_j - \hat{P}(y_j=1/x_j, \omega^{(t)})]$$

$$\therefore \omega_0^{(1)} = \omega_0^{(0)} + \eta [(1-0.5) + (1-0.5) + (0-0.5) + (0-0.5)]$$

$$= 0 + \eta [0] = 0$$

$$\omega_i^{(t+1)} = \omega_i^{(t)} + \eta \sum_{j=1}^4 [y_j - \hat{P}(y_j=1/x_j, \omega^{(t)})]$$

$$\therefore \omega_1^{(1)} = \omega_1^{(0)} + \eta [(1-0.5) + 0(1-0.5) + 0(0-0.5) + (-1)(0-0.5)]$$

$$= 0 + \eta [0.5 + 0.5] = \eta$$

$$\therefore \omega_2^{(1)} = \omega_2^{(0)} + \eta [0(1-0.5) + (-1)(1-0.5) + 1(0-0.5) + 0(0-0.5)]$$

$$= 0 + \eta [-0.5 - 0.5] = -\eta$$

$$\therefore \omega^{(1)} = (0, \eta, -\eta)'$$

Iteration 2 :-

$$\hat{P}(Y=1 | X^1, \omega^{(1)}) = \frac{\exp(0 + \eta(1) + (-\eta)(0))}{1 + \exp(0 + \eta(1) + (-\eta)(0))} = \frac{\exp(\eta)}{1 + \exp(\eta)}$$

$$\hat{P}(Y=1 | X^2, \omega^{(1)}) = \frac{\exp(0 + \eta(0) + (-\eta)(-1))}{1 + \exp(0 + \eta(0) + (-\eta)(-1))} = \frac{\exp(\eta)}{1 + \exp(\eta)}$$

$$\hat{P}(Y=1 | X^3, \omega^{(1)}) = \frac{\exp(0 + \eta(0) + (-\eta)(1))}{1 + \exp(0 + \eta(0) + (-\eta)(1))} = \frac{\exp(-\eta)}{1 + \exp(-\eta)} = \frac{1}{1 + \exp(\eta)}$$

$$\hat{P}(Y=1 | X^4, \omega^{(1)}) = \frac{\exp(0 + \eta(-1) + (-\eta)(0))}{1 + \exp(0 + \eta(-1) + (-\eta)(0))} = \frac{\exp(-\eta)}{1 + \exp(-\eta)} = \frac{1}{1 + \exp(\eta)}$$

$$\begin{aligned} \omega_0^{(2)} &= \omega_0^{(1)} + \eta \left[ (1 - \frac{\exp(\eta)}{1 + \exp(\eta)}) + (1 - \frac{\exp(\eta)}{1 + \exp(\eta)}) + (0 - \frac{1}{1 + \exp(\eta)}) + (0 - \frac{1}{1 + \exp(\eta)}) \right] \\ &= 0 + \eta \left[ \frac{2}{1 + \exp(\eta)} - \frac{2}{1 + \exp(\eta)} \right] = 0 \end{aligned}$$

$$\begin{aligned} \omega_1^{(2)} &= \omega_1^{(1)} + \eta \left[ 1 \left( \frac{1}{1 + \exp(\eta)} \right) + 0 \left( \frac{1}{1 + \exp(\eta)} \right) + 0 \left( \frac{-1}{1 + \exp(\eta)} \right) + (-1) \left( \frac{-1}{1 + \exp(\eta)} \right) \right] \\ &= \eta + \eta \left[ \frac{2}{1 + \exp(\eta)} \right] = \eta + \frac{2\eta}{1 + \exp(\eta)} \end{aligned}$$

$$\begin{aligned} \omega_2^{(2)} &= \omega_2^{(1)} + \eta \left[ 0 \left( \frac{1}{1 + \exp(\eta)} \right) + (-1) \left( \frac{1}{1 + \exp(\eta)} \right) + 1 \left( \frac{-1}{1 + \exp(\eta)} \right) + 0 \left( \frac{-1}{1 + \exp(\eta)} \right) \right] \\ &= -\eta + \eta \left[ \frac{-2}{1 + \exp(\eta)} \right] = -\eta - \frac{2\eta}{1 + \exp(\eta)} \end{aligned}$$

$$\therefore \omega^{(2)} = \left( 0, \eta + \frac{2\eta}{1 + \exp(\eta)}, -\eta - \frac{2\eta}{1 + \exp(\eta)} \right)$$

As we continue the number of iterations like this, we get a final  $\omega$ , i.e  $\omega^{(\infty)} = (0, \infty, -\infty)'$

Hence  $\omega^{(\infty)} = (0, \infty, -\infty)'$  when  $\omega^{(0)} = (0, 0, 0)'$

4 ii) Initial weight vector,  $w^{(0)} = (0, 1, 0)^T$

Iteration 1:

$$\hat{P}(Y=1 | X^1, w^{(0)}) = \frac{\exp(0 + 1(1) + 0(0))}{1 + \exp(0 + 1(1) + 0(0))} = \frac{\exp(1)}{1 + \exp(1)} = \frac{e}{1+e}$$

$$\hat{P}(Y=1 | X^2, w^{(0)}) = \frac{\exp(0 + 1(0) + 0(-1))}{1 + \exp(0 + 1(0) + 0(-1))} = \frac{1}{2} = 0.5$$

$$\hat{P}(Y=1 | X^3, w^{(0)}) = \frac{\exp(0 + 1(0) + 0(1))}{1 + \exp(0 + 1(0) + 0(-1))} = \frac{1}{2} = 0.5$$

$$\hat{P}(Y=1 | X^4, w^{(0)}) = \frac{\exp(0 + 1(-1) + 0(0))}{1 + \exp(0 + 1(-1) + 0(0))} = \frac{e^{-1}}{1 + e^{-1}} = \frac{1}{1+e}$$

$$\therefore w_0^{(1)} = w_0^{(0)} + \eta \left[ (1 - \frac{e}{1+e}) + (1 - 0.5) + (0 - 0.5) + (0 - \frac{1}{1+e}) \right] \\ = 0 + \eta \left[ (\frac{1}{1+e}) + (-0.5) + (-0.5) + (\frac{-1}{1+e}) \right] = 0 + \eta(0) = 0$$

$$\therefore w_1^{(1)} = w_1^{(0)} + \eta \left[ 1(\frac{1}{1+e}) + 0(0.5) + 0(-0.5) + (-1)(\frac{-1}{1+e}) \right] \\ = 1 + \eta \left[ \frac{2}{1+e} \right] = 1 + \frac{2\eta}{1+e}$$

$$\therefore w_2^{(1)} = w_2^{(0)} + \eta \left[ 0(\frac{1}{1+e}) + (-1)(0.5) + 1(-0.5) + 0(\frac{-1}{1+e}) \right] \\ = 0 + \eta \left[ -0.5 - 0.5 \right] = -\eta$$

$$\therefore w^{(1)} = (0, 1 + \frac{2\eta}{1+e}, -\eta)$$

Iteration 2:

$$\hat{P}(Y=1 | X^1, w^{(1)}) = \frac{\exp(0 + (1 + \frac{2\eta}{1+e})(1) + (-\eta)(0))}{1 + \exp(0 + (1 + \frac{2\eta}{1+e})(1) + (-\eta)(0))} \\ = \frac{\exp(1 + \frac{2\eta}{1+e})}{1 + \exp(1 + \frac{2\eta}{1+e})}$$

$$\hat{P}(Y^2=1 | X^2, w^{(1)}) = \frac{\exp(0 + (1 + \frac{2n}{1+e})(0) + (-n)(-1))}{1 + \exp(0 + (1 + \frac{2n}{1+e})(0) + (-n)(-1))} \\ = \frac{\exp(n)}{1 + \exp(-n)}$$

$$\hat{P}(Y^3=1 | X^3, w^{(1)}) = \frac{\exp(0 + (1 + \frac{2n}{1+e})(0) + (-n)(1))}{1 + \exp(0 + (1 + \frac{2n}{1+e})(0) + (-n)(1))} \\ = \frac{\exp(-n)}{1 + \exp(-n)} = \frac{1}{1 + \exp(n)}$$

$$\hat{P}(Y^4=1 | X^4, w^{(1)}) = \frac{\exp(0 + (1 + \frac{2n}{1+e})(-1) + (-n)(0))}{1 + \exp(0 + (1 + \frac{2n}{1+e})(-1) + (-n)(0))} \\ = \frac{\exp(-(1 + \frac{2n}{1+e}))}{1 + \exp(-(1 + \frac{2n}{1+e}))} = \frac{1}{1 + \exp(1 + \frac{2n}{1+e})}$$

$$\therefore w_0^{(2)} = w_0^{(1)} + \eta \left[ \left( 1 - \frac{\exp(1 + \frac{2n}{1+e})}{1 + \exp(1 + \frac{2n}{1+e})} \right) + \left( 1 - \frac{\exp(n)}{1 + \exp(n)} \right) \right. \\ \left. + \left( 0 - \frac{1}{1 + \exp(n)} \right) + \left( 0 - \frac{1}{1 + \exp(1 + \frac{2n}{1+e})} \right) \right] \\ = 0 + \eta \left[ \frac{1}{1 + \exp(1 + \frac{2n}{1+e})} + \frac{1}{1 + \exp(n)} - \frac{1}{1 + \exp(n)} \right. \\ \left. - \frac{1}{1 + \exp(1 + \frac{2n}{1+e})} \right] \\ = 0 + \eta [0] = 0$$

$$\begin{aligned}
 \therefore w_1^{(2)} &= w_1^{(1)} + \eta \left[ 1 \cdot \left( \frac{1}{1 + \exp\left(\frac{1+2n}{1+e}\right)} \right) + 0() + 0(0) + (-1) \left( \frac{-1}{1 + \exp\left(\frac{1+2n}{1+e}\right)} \right) \right] \\
 &= \left( 1 + \frac{2n}{1+e} \right) + \eta \left[ \frac{2}{1 + \exp\left(\frac{1+2n}{1+e}\right)} \right] \\
 &= 1 + \frac{2n}{1+e} + \frac{2n}{1 + \exp\left(\frac{1+2n}{1+e}\right)} \\
 w_2^{(2)} &= w_2^{(1)} + \eta \left[ 0() + (-1) \left( \frac{1}{1 + \exp(n)} \right) + 1 \left( \frac{-1}{1 + \exp(n)} \right) + 0() \right] \\
 &= -n + \eta \left[ \frac{-2}{1 + \exp(n)} \right] = -n - \frac{2n}{1 + \exp(n)} \\
 &= - \left[ n + \frac{2n}{1 + \exp(n)} \right] \\
 \therefore w^{(2)} &= \left( 0, 1 + \frac{2n}{1+e} + \frac{2n}{1 + \exp\left(\frac{1+2n}{1+e}\right)}, - \left[ n + \frac{2n}{1 + \exp(n)} \right] \right)^T
 \end{aligned}$$

As we continue the number of iterations, the pattern continues and we get a final  $w^{(*)}$

$$\text{i.e } w^{(*)} = (0, \infty, -\infty)^T$$

Hence the final vector,  $w^{(*)}$  will be the same for the two different initial values.

5) Naive Bayes Classifier and Logistic Regression

i) Gaussian Naive Bayes & Logistic Regression

- Number of independent parameters in GNB?

Sol:-  $4d + 1$ ,  $d$  being the number of features

i.e  $P(Y=1)$

$$d \begin{cases} \mu_{1,0} \\ \mu_{2,0} \\ \mu_{3,0} \\ \vdots \\ \mu_{d,0} \end{cases} \quad d \begin{cases} \mu_{1,1} \\ \mu_{2,1} \\ \mu_{3,1} \\ \vdots \\ \mu_{d,1} \end{cases} \quad d \begin{cases} \sigma_{1,0} \\ \sigma_{2,0} \\ \sigma_{3,0} \\ \vdots \\ \sigma_{d,0} \end{cases} \quad d \begin{cases} \sigma_{1,1} \\ \sigma_{2,1} \\ \sigma_{3,1} \\ \vdots \\ \sigma_{d,1} \end{cases}$$

Hence  $4d + 1$  independent parameters

- No, 'w' cannot be translated into the parameters of an equivalent GNB classifier without any extra assumption.

Extra assumption: Variance is the same across all the classes.

i.e  $\sigma_{ik} = \sigma_i$  ( $K$  no. of classes)

With this assumption, we can be translated.

Explanation:

From Gaussian Naive Bayes,

$$P(Y=1|X) = \frac{P(Y=1) P(X|Y=1)}{P(Y=1) P(X|Y=1) + P(Y=0) P(X|Y=0)}$$

$$P(Y=1|X) = \frac{1}{1 + \exp\left(\ln\left(\frac{P(Y=0)}{P(Y=1)}\right) + \sum_i \ln \frac{P(X_i|Y=0)}{P(X_i|Y=1)}\right)}$$

Assuming  $P(Y=1) = \theta$ ,  $P(Y=0) = 1-\theta$  &  $\sigma_{ik} = \sigma_i$ ,  
i.e variance same across classes

we get,

$$P(Y=1|X) = \frac{1}{1 + \exp\left(\ln\frac{1-\theta}{\theta} + \sum_i \left(\frac{\mu_{io}-\mu_{ii}}{\sigma_i^2}\right) X_i + \left(\frac{\mu_{ii}-\mu_{io}}{2\sigma_i^2}\right)\right)}$$

This is of the form

$$P(Y=1|X) = \frac{1}{1 + \exp(w_0 + \sum_i w_i X_i)}$$

which is the form of Logistic Regression. Hence  
Hence parameters of logistic regression  $w$ , can  
be translated into equivalent Gaussian Naive  
Bayes classifier.

## **Question 5.2: Implementation of Gaussian Naïve Bayes and Logistic Regression**

### i) Pseudocode

Load X(Samples), y(class) from the dataset

For fraction in fractions (0.01, 0.02, 0.05, 0.1, 0.625, 1)

    Divide data into 3 folds making one of the folds test and other 2 folds training set

    We will end up with 3 different combinations of train, test data

    In each such combination, randomly pick a fraction of train data

    Repeat the above step 5 times so that we will have 5 sets of data for each combination

        Train the model for each of the set and learn the parameters

        Predict the class for the test data from the learnt parameters and calculate the accuracy

        Calculate the mean of accuracies which will be the accuracy for the fraction

Plot accuracy vs the fraction size curve

### **Gaussian Naïve Bayes**

Function **Train**(X\_train, y\_train):

$P(y=1)$  = number of positive samples / total number of samples

Find  $\text{Mean}_{i,1}$  for training data of class  $y=1$

Find  $\text{Mean}_{i,0}$  for training data of class  $y=0$

Find  $\text{Variance}_{i,1}$  for training data of class  $y=1$

Find  $\text{Variance}_{i,0}$  for training data of class  $y=0$

Learnt\_parameters =  $P(y=1)$ ,  $\text{Mean}_{i,1}$ ,  $\text{Mean}_{i,0}$ ,  $\text{Variance}_{i,1}$ ,  $\text{Variance}_{i,0}$

Return Learnt\_parameters

Function **Predict**(X\_test, y\_test, Learnt\_parameters):

$P(X/Y=1) = (1/\sqrt{2\pi \text{Variance}_{i,1}}) \exp(-((X_{\text{test}} - \text{Mean}_{i,1})^2)/(2 \text{Variance}_{i,1}))$

$H_{\text{pos}} = P(Y=1) * \prod_i P(X_i/Y=1)$

$P(X/Y=0) = (1/\sqrt{2\pi \text{Variance}_{i,0}}) \exp(-((X_{\text{test}} - \text{Mean}_{i,0})^2)/(2 \text{Variance}_{i,0}))$

$H_{\text{zero}} = P(Y=0) * \prod_i P(X_i/Y=0)$

If  $H_{\text{pos}} > H_{\text{zero}}$  then

```

y_pred = 1.0
else
    y_pred = 0.0
Accuracy = (# y_pred == y_test) / (# y_pred)
Return Accuracy

```

### **Logistic Regression**

```

Function train(X_train, y_train, learning_rate)
W_0=0, W=(0, 0, 0, 0)
For 500 iterations
    Z = w_0 +  $\sum_i w_i X_{\text{train}_i}$ 
    P(Y=1/X_train, w) = exp(z)/(1+exp(z)) = 1/(1+exp(-z)) = sigmoid(-z) = a
    w_0 = w_0 + learning_rate * ( $\sum_j (y_{\text{train}_j} - P(Y=1/X_{\text{train}_j}, w))$ )
    w = w + learning_rate * ( $\sum_i X^i (y_{\text{train}_i} - P(Y=1/X_{\text{train}_i}, w))$ )
return w, w_0

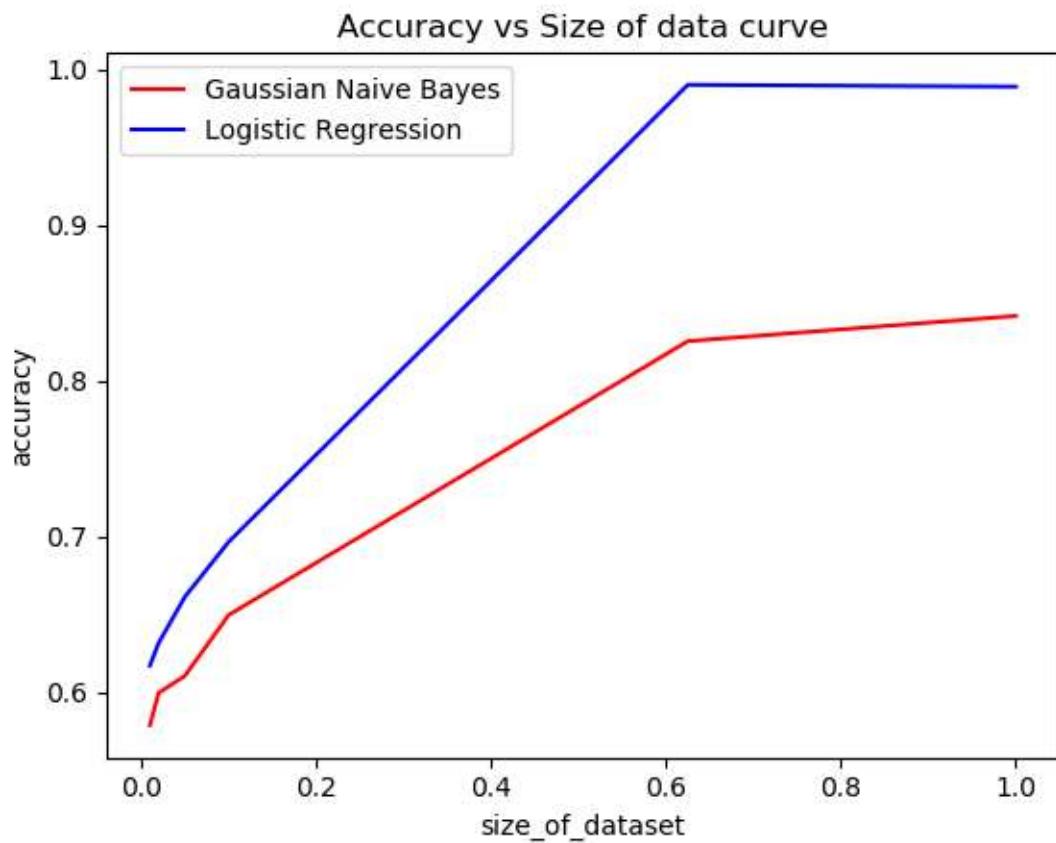
```

```

Function predict(X_test, y_test, w, w_0)
Z = w_0 +  $\sum_i w_i X_{\text{test}_i}$ 
P(Y=1/X_test, w) = exp(z)/(1+exp(z)) = 1/(1+exp(-z)) = sigmoid(-z) = a
If P(Y=1/X_test, w) > 0.5 then
    y_pred = 1.0
else
    y_pred = 0.0
Accuracy = (# y_pred == y_test) / (# y_pred)
Return Accuracy

```

ii) Learning curve



iii) Generate examples and compare their mean and variance with the training data

#### Fold 1

**Accuracy** = 0.838074398249

Mean of training set: [-1.92548409 -1.22379259 2.43710485 -1.18223808]

Mean of generated set: [-1.87141318 -1.2079421 2.60292242 -1.20923457]

**Percentage of Mean Deviated:** [ 2.80817191 1.2951936 6.37043817 2.23252709]

Variance of training set: [ 3.70027822 30.50494584 28.38578698 4.1709176 ]

Variance of generated set: [ 3.44684103 28.88546879 29.63427956 3.83852848]

**Percentage of Variance Deviated:** [ 6.84913883 5.30889996 4.21300131 7.96920828]

#### Fold 2

**Accuracy** = 0.851203501094

Mean of training set: [-1.83526396 -0.80990072 1.90860182 -1.33436393]

Mean of generated set: [-1.8923333 -1.10565946 2.0417084 -1.37504932]

**Percentage of Mean Deviated:** [ 3.01581866 26.74953262 6.51937279 2.9588308 ]

Variance of training set: [ 3.51025097 29.10009686 27.13543712 4.59173298 ]

Variance of generated set: [ 3.31589324 31.0324415 29.20502293 4.53426649 ]

**Percentage of Variance Deviated:** [ 5.53686119 6.22685341 7.08640364 1.25152073 ]

#### Fold 3

**Accuracy** = 0.835886214442

Mean of training set: [-1.84221312 -0.93518316 2.0838102 -1.22831574]

Mean of generated set: [-1.95093491 -1.15531175 1.98556775 -1.19001075]

**Percentage of Mean Deviated:** [ 5.5728046 19.05360942 4.71455866 3.11849693 ]

Variance of training set: [ 3.38007627 27.77767508 27.22473826 4.08723726 ]

Variance of generated set: [ 3.35737334 31.72728121 26.58498813 3.94825671 ]

**Percentage of Variance Deviated:** [ 0.67166924 12.44861199 2.34988531 3.40035454 ]

As it can be observed, the mean and variance of the training samples is matching with the mean and variance of the generated samples. This is because both are generated from the same underlying gaussian distribution,  $N(\text{mean}, \text{variance})$ .