

CS534 Written Assignment 2

Aashish Adhikari

New Link

Date
Page

Q.1) Solution:

The question is analogous to putting n balls in k bins and thus the distribution is multinomial.

pmf for the parameters can be taken from multinomial distribution.

$$\text{Likelihood}(p_a, p_c, p_g, p_t) = p_a^{n_a} \times p_c^{n_c} \times p_g^{n_g} \times p_t^{n_t}$$

$$\text{We have the constraint } p_a + p_c + p_g + p_t = 1$$

In such cases where we have constraints, for optimization, we have Lagrange Multiplier.

We need to find maximum of likelihood (p_a, p_c, p_g, p_t) subject to $p_a + p_c + p_g + p_t = 1$.

$$\Rightarrow F(p_a, p_c, p_g, p_t) = L(p_a, p_c, p_g, p_t) - \lambda (p_a + p_c + p_g + p_t - 1)$$

$$\text{Or, } F(p_a, p_c, p_g, p_t) = \left[p_a^{n_a} \times p_c^{n_c} \times p_g^{n_g} \times p_t^{n_t} \right] - \lambda p_a - \lambda p_c - \lambda p_g - \lambda p_t + \lambda$$

Now, we take partial derivatives wrt p_a, p_c, p_g, p_t & λ and set them to zero to maximize.

~~$\therefore \frac{\partial F}{\partial p_a} =$~~ To make computation easier, let's take log of the likelihood.

$$F'(\lambda) = \log(\lambda) = \log(p_a^{n_a}) + \log(p_c^{n_c}) + \log(p_g^{n_g}) + \log(p_t^{n_t})$$

Let's represent a by 1, c by 2, g by 3, t by 4.

$$\Rightarrow \log(\lambda) = \sum_{i=1}^4 \log p_i^{n_i}$$

$$= \sum_{i=1}^4 n_i \log \frac{p_i}{\sum p_i}$$

Also, $\lambda(p_a + p_c + p_g + p_t - 1)$ becomes $\lambda \left[\left[\sum_{i=1}^4 p_i \right] - 1 \right] = \lambda G$

$$\therefore \frac{\partial [F'(\lambda) - \lambda G]}{\partial p_i} = \frac{n_i \times p_i^{n_i-1}}{(p_i^{n_i})} - \lambda [1]$$

Setting it to 0;

$$n_i \times p_i^{n_i-1-h_i} = \lambda$$

$$\text{or, } \frac{p_i^{h_i}}{p_i} = \lambda$$

$$\Rightarrow p_i = \frac{n_i}{\lambda}$$

~~$$\text{Also, } \frac{\partial [F'(\lambda) - \lambda G]}{\partial \lambda} = 0 - [4\lambda + 1]$$~~

~~$$\text{Setting to 0 } \Rightarrow 4\lambda + 1 = 0$$~~

$$\Rightarrow \lambda = -\frac{1}{4}$$

$$\text{Also, } \frac{\partial [F'(\lambda) - \lambda G]}{\partial \lambda} = 0 - p_a - p_c - p_g - p_t + 1$$

$$\text{Setting to 0 } \Rightarrow 0 = \frac{n_a}{\lambda} + \frac{n_c}{\lambda} + \frac{n_g}{\lambda} + \frac{n_t}{\lambda} - 1$$

$$\text{or, } +\lambda = n_a + n_c + n_g + n_t = n_{\text{total}}$$

$$\text{or, } \lambda = +n_{\text{total}}$$

$$\therefore p_i^o = \frac{n_i}{n_{\text{total}}}$$

maximum likelihood estimators are:

$$\Rightarrow p_a = \frac{n_a}{(n_a + n_c + n_g + n_t)}$$

$$p_c = \frac{n_c}{n_a + n_c + n_g + n_t}$$

$$p_g = \frac{n_g}{n_a + n_c + n_g + n_t}$$

$$p_t = \frac{n_t}{n_a + n_c + n_g + n_t}$$

Q.2. Solution

Naive Bayes Classifier uses

$$P(A|B) = \frac{P(B|A) \times P(A)}{P(B)}$$

$$= \frac{P(B|A) \times P(A)}{\sum_B P(A|B) P(B)}$$

$$\sum_B P(A|B) P(B)$$

Here, we have 3 features A, B, C & output Y.

2) $\therefore P(Y) = 1/2$ for $Y=1 \Rightarrow P(Y=0) = 1/2$ & $P(Y=1) = 1/2$

$$p(Y) = 1/2 \text{ for } Y=0$$

$$\begin{aligned} P(A=0|y=0) &= 2/3 \\ P(B=0|y=0) &= 1/3 \\ P(C=0|y=0) &= 1/3 \end{aligned} \quad \left. \right\} \text{for } y=0$$

$$\begin{aligned} P(A=0|y=1) &= 1/3 \\ P(B=0|y=1) &= 1/3 \\ P(C=0|y=1) &= 2/3 \end{aligned} \quad \left. \right\} \text{for } y=1$$

Since the assumption is that the features are independent, probabilities like $P(A=0, B=0, C=1 | y=0)$ can be calculated by multiplying individual probabilities such as $P(A=1|y=0)$, $P(B=0|y=0)$, and $P(C=1|y=0)$.

Hence, above 8 probabilities are all that we need to estimate.

b. Solution

$$\begin{aligned} &P(y=1 | A=1, B=0, C=0) \\ &= \frac{P(A=1, B=0, C=0)}{P(A=1, B=0, C=0)} P(y=1) \end{aligned}$$

Since they are independent, multiply

$$\begin{aligned} &\cancel{\frac{2/3 \times 1/3 \times 2/3 \times 1/2}{2/3 \times 1/3 \times 2/3 \times 1/2 + 1/3 \times 1/3 \times 1/3 \times 1/2}} \\ &\cancel{2/3 \times 1/3 \times 2/3 \times 1/2 + 1/3 \times 1/3 \times 1/3 \times 1/2} \end{aligned}$$

$$\begin{aligned}
 &= \frac{2/3 \times 1/3 \times 2/3 \times 1/2}{2/3 \times 1/3 \times 2/3 \times 1/2 + 1/3 \times 1/3 \times 1/3 \times 1/2} \\
 &= \frac{4 \times 54}{54 \times 5} \\
 &= \frac{4}{5}
 \end{aligned}$$

C. Solution:

No, it is not valid.

It is because Naive-Bayes classifier assumes conditional independence. However, the question has mentioned that A, B, C are independent and independence does not guarantee conditional independence.

For example, consider picking a card from a deck of 52 cards and rolling a die.

A is an event of choosing an Ace.

S is an event of rolling a 6.

probability of $(A|S)$ = probability of A ← Independent

$$\Rightarrow P(A|S) = \frac{P(S|A) P(S)}{P(A)}$$

$$\text{or, } [P(A)]^2 = [P(S)]^2$$

$$\text{or, } \left(\frac{4}{52}\right)^2 = \left(\frac{1}{6}\right)^2 \text{ which is not true. } \cancel{\text{X}}$$

Hence, invalid.

Q3 Solution:

p is the probability of head for a weighted coin.

It has a Beta distribution

$$p \sim \text{Beta}(2, 2)$$

$$\text{pdf } f(p; 2, 2) = \frac{1}{B(\alpha=2, \beta=2)} p^{2-1} (1-p)^{2-1}$$

$$= \frac{1}{B(2, 2)} p \times (1-p)$$

Posterior Distribution of p is given by

$$\frac{1}{B(n_1+2, n_0+2)} p^{n_1+1} \times (1-p)^{n_0+1}$$

where n_1 is the no of heads
 n_0 is the no of tails

Case 1 : $n_1 + n_0 = 5$

$$n_1 = 2$$

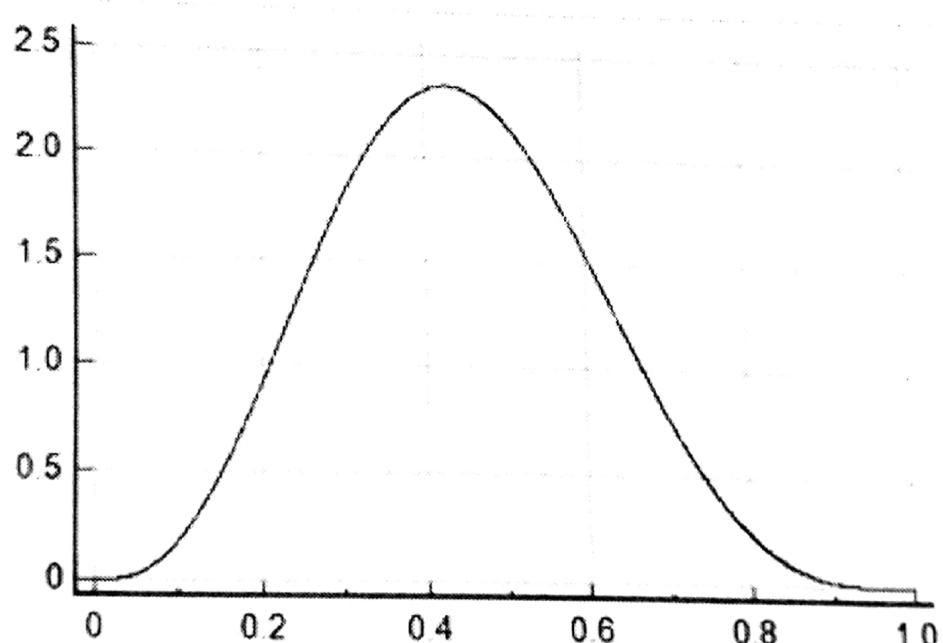
$$\Rightarrow \frac{1}{B(4, 5)} p^3 \times (1-p)^4$$

Case 2: $n_1 + n_0 = 50$

$$n_1 = 20$$

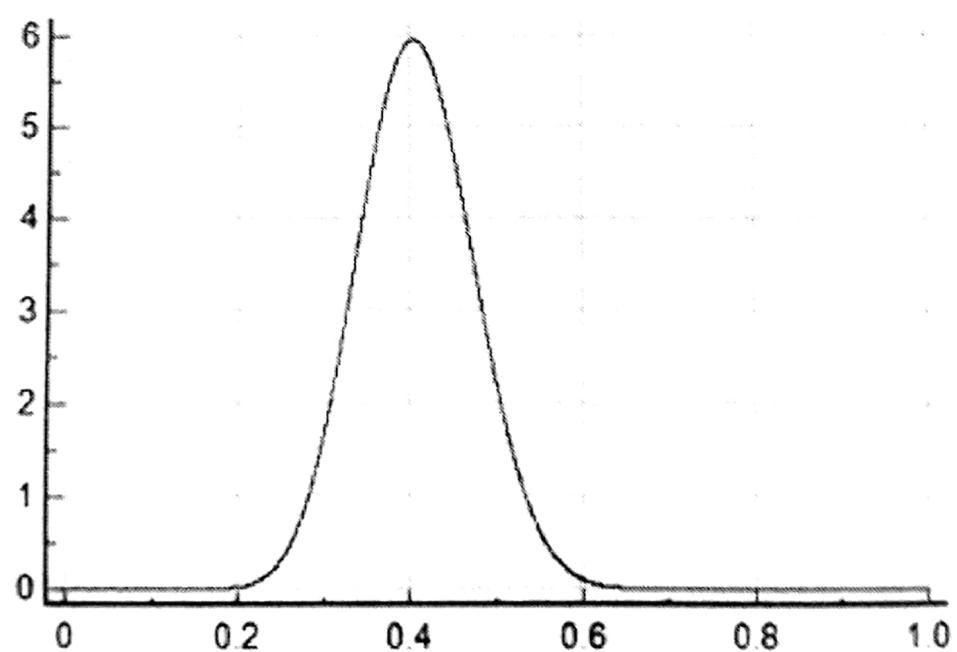
$$\Rightarrow \frac{1}{B(22, 32)} p^{21} (1-p)^{31}$$

Graph



Function: PDFBeta(x .4 ,5)

Graph

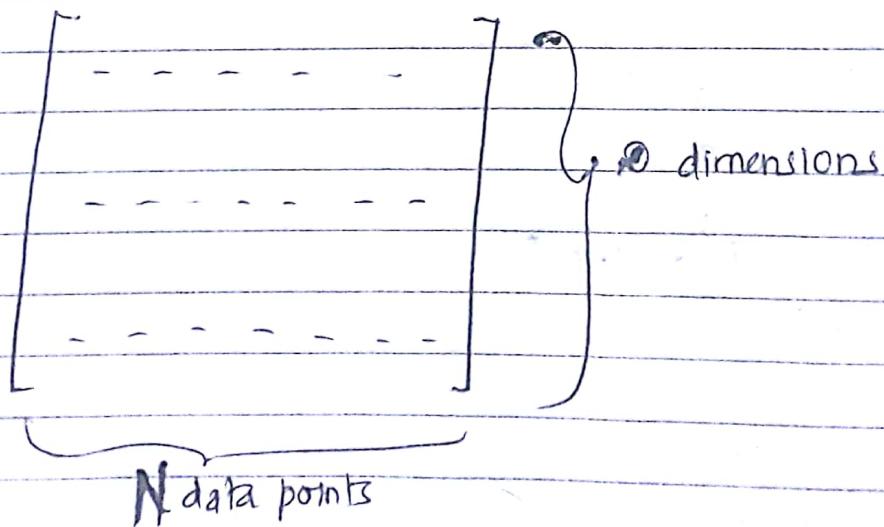


Function: PDFBeta(x .22 ,32)

We can see that as the number of coin tosses from this coin increases, the probability gets nearer and nearer to 0.4. Thus as the number of coin tosses tends to a large number, the probability of getting a head ≈ 0.4 according to the distribution plot i.e., it reaches to toward the real value.

Q.4. Solution

a.) We have N data points in d dimensions.



we map each point \vec{x}_n to $(d+1)$ dimensional point $\langle \vec{x}_n, e_n \rangle$ where e_n is a N -dimensional vector of all zeros but one 1 at the n^{th} position.

After adding the dimensions

New Link

Date _____

Page _____

individual points	$\xrightarrow{\quad \quad}$	$\begin{bmatrix} 1 & x_1 & x_2 & x_3 & \dots & x_d & 1 & 0.0 & \dots & 0 \\ 1 & x_1 & x_2 & x_3 & \dots & x_d & 0 & 1 & 0 & \dots & 0 \\ 1 & & & & & & & & & \\ 1 & & & & & & & & & \\ 1 & & & & & & & & & \\ 1 & x_1 & x_2 & x_3 & \dots & x_d & 0 & 0 & 0 & \dots & 1 \end{bmatrix}$	$\xleftarrow{\quad \quad}$	
			$d+n+1$ dimensions	\longrightarrow

The perceptron loss is given as $L_p(g(w, x), y) = \max \begin{cases} 0 \\ -y w^T x \end{cases}$

so, when we increase the dimension by n , the new dimension is $(d+n)$.

Let points x_0, x_1, \dots, x_n are in a ~~$d+1$~~ ^(d+n) dimensional space as shown in the matrix above. These points suppose come under 2 different sets i.e., are linearly separable. For this, there must exist $(d+n+1)$ real numbers $w_0, w_1, \dots, w_n, w_{n+1}$ such that each point x_i in A satisfies $\sum_{i=0}^n w_i x_i \geq k$ and

each point x_i in B satisfies $\sum_{i=0}^n w_i x_i < k$.

Consider we are adding 5 new dimensions.

$$X \begin{bmatrix} \dots & \dots & 1 & 0 & 0 \\ \dots & \dots & 0 & 1 & 0 \\ \dots & \dots & 0 & 0 & 1 \end{bmatrix}$$

the new weights added are w_1, w_2, w_3 .

We need to ~~find the~~ make sure that the $\sum w_i x_i$ is greater or less than k.

Thus, setting learning ~~example~~ values of new weights such that each of them lies on one side of the linear classifier guarantees the solution.

Notice that for each example,

If we set all values upto dimension (1 to d+1) as 1 and a large positive or negative value for every other position in our new vector of weight, then the result is either we get $\hat{g}(x, w)$ is +ve or -ve.

For example

$$[w_0 \ w_1 \ w_2 \ w_3] \quad \begin{bmatrix} 1 \\ x_1 \\ 1 \\ 0 \end{bmatrix} \quad \begin{array}{l} \text{already existing input} \\ \text{added according to the question} \end{array}$$

$$w_3 \times 0 = 0$$

only w_2 is taken to be significant among the added dimensions.

so, setting w_2 to a large negative value can allow us to classify as a negative class.

Hence such a projection forces the data to be linearly separable.

New Link

Date
Page

- b) This mapping does not allow us to generalize, because for each different length of ϵ_n , we need to address the weight uniquely. However, forcing the data to be linearly separable is possible independent of the size of ϵ_n .

Q. B. Solution

$$\text{Given } K(x_i, x_j) = (x_i \cdot x_j + 1)^3$$

Let $x_i = a$ & $x_j = b$.

$$\therefore (a \cdot b + 1)^3 \\ = a^3 + b^3 + 3a^2b + 3ab^2$$

Assume we have p features for each training example

$$\begin{aligned} & \sum_{x=1}^p (a_x)^3 + \sum_{x=1}^p (b_x)^3 + 3 \sum_{x=1}^p (a_x)^2 b_x + 3 \sum_{x=1}^p a_x (b_x)^2 \\ &= \sum_{x=1}^p \sum_{y=1}^p \sum_{z=1}^p a_x a_y a_z + \sum_{x=1}^p \sum_{y=1}^p \sum_{z=1}^p b_x b_y b_z \\ &+ 3 \sum_{x=1}^p \sum_{y=1}^p a_x a_y b_x + 3 \sum_{x=1}^p \sum_{y=1}^p a_x b_x b_y \\ &= \sum_{x=1}^p \sum_{y=1}^p \sum_{z=1}^p a_x a_y a_z + \sum_{x=1}^p \sum_{y=1}^p \sum_{z=1}^p b_x b_y b_z + \\ & 3 \sum_{x=1}^p \sum_{y=1}^p a_x^2 b_x + 3 \sum_{x=1}^p \sum_{y=x+1}^p a_x a_y b_x + \\ & 3 \sum_{x=1}^p \sum_{y=1}^p a_x b_x^2 + 3 \sum_{x=1}^p \sum_{y=x+1}^p a_x b_x b_y \end{aligned}$$

i.e. the kernel corresponding ϕ functions are :

$$\{ \sqrt{3} a_x^3, b_x^3, \sqrt{3} a_x^2, \sqrt{3} b_x, \sqrt{3} a_x a_y, \sqrt{3} a_x b_x, \\ \sqrt{3} a_x, \sqrt{3} b_x b_y \}$$

New Link

Date _____

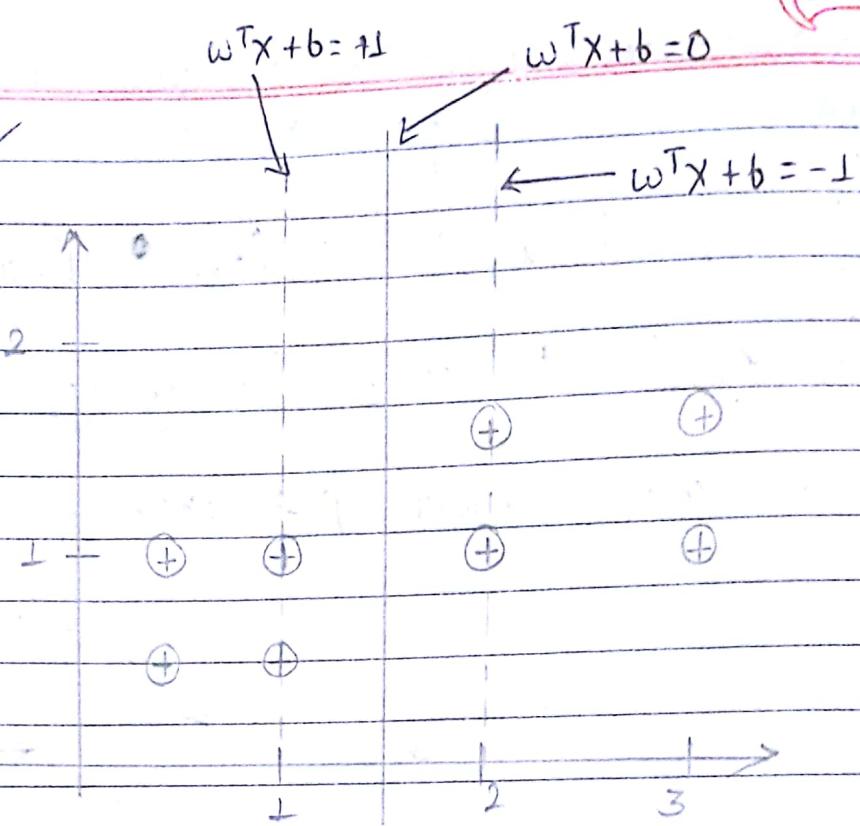
Page _____

Changing indices $x_i y_{j,k}$ to 1, 2, 3 and
converting in terms of x_i & x_j .
 $a = x_i, b = x_j$

$$\{ x_{i,j}^3, x_{j,i}^3, \sqrt{3}(x_{i,j})^2, \sqrt{3}(x_{j,i})^2, \sqrt{3} x_{i,j} x_{j,i}, \\ \sqrt{3} (x_{j,i})^2, \sqrt{3} x_{i,i}, \sqrt{3} x_{j,i} x_{j,i} \}$$

which is our required ϕ function.

Q 6 Solution



- $w^T x + b = 0$ is the decision boundary that lies in the middle of -1 & 1 on X axis.
- ⇒ point (1.5, 0) lies on the decision boundary.
- ⇒ points (1, 0.5) & (1, 1) lie on $w^T x + b = 0.5$
- ⇒ points (2, 1) & (2, 1.5) lie on $w^T x + b = 1$

$$w^T x + b = 0$$

we have the point (1.5, 0) on this boundary

(b) Taking weight vector $w = \begin{pmatrix} w \\ 0 \end{pmatrix}$ & $w^T = [w \ 0]$

$$\begin{bmatrix} w & 0 \end{bmatrix} \begin{bmatrix} 1.5 \\ 0 \end{bmatrix} + b = 0$$

$$\text{gives } b = -1.5w.$$

For $w^T x + b = -1$, there are 2 points $(2, 1)$ & $(2, 1.5)$
 for $w^T x + b = 1$, we have $(1, 0.5)$ & $(1, 1)$.

$$\left[\begin{matrix} w & 0 \end{matrix} \right] \left[\begin{matrix} 1 \\ 0.5 \end{matrix} \right] + b = 1 \quad \left| \quad \left[\begin{matrix} w & 0 \end{matrix} \right] \left[\begin{matrix} 1 \\ 1 \end{matrix} \right] + b = 1 \right.$$

$$w + b = 1 \qquad \qquad \qquad w + b = 1$$

I

plot values for $w^T x + b = 1$

$$\left[\begin{matrix} w & 0 \end{matrix} \right] \left[\begin{matrix} 2 \\ 1 \end{matrix} \right] + b = -1 \quad \left| \quad \left[\begin{matrix} w & 0 \end{matrix} \right] \left[\begin{matrix} 2 \\ 1.5 \end{matrix} \right] + b = -1 \right.$$

$$2w + b = -1 \qquad \qquad \qquad 2w + b = -1$$

II

Solve I & II

$$\begin{array}{r} 2w + b = -1 \\ -w + b = -1 \\ \hline w = -2 \end{array}$$

$$\text{put values in } II \Rightarrow b = -1 - 2(-2) \\ = 3.$$

$$\therefore b = 3, w = -2 \neq$$