**4.1 a)**

Lets say we have 'd' document available

And as the classification problem is binary (Onion OR Economist)

$$k = 2 \quad y \cdot y = 1 \rightarrow \text{Economist}$$
$$y = 2 \rightarrow \text{Onion}$$

The $y^{(i)}$ represent the category (Onion or Economist) of $i$th document in collection.

$x_j^{(i)}$ corresponds to presence or absence (0/1) of particular word in $i$th document from the list of words in vocabulary v.

Naive Bayes algorithm basically finds the most probable class label (Onion or Economist) given $x_i (x) \ldots - x_d) =$

$$\hat{y} = \underset{y}{\arg\max} \; P(Y = y \mid X)$$
$$= \underset{y}{\arg\max} \; P(y \mid x)$$

As $P(y \mid x) = \dfrac{P(x \mid y) \, P(y)}{P(x)}$

$$= \underset{y}{\arg\max} \; \frac{P(x \mid y) \, P(y)}{P(x)}$$

As $P(x)$ is just a scaling variable, we can ignore it.

$$= \underset{y}{\arg\max} \; P(x \mid y) \, P(y)$$

Now as we have V such feature vectors where each feature vector $(x_i)$ is 1 if $i$th word appears in the article.

So,

$$\underset{y}{\arg\max} \; P(x \mid y) \, P(y)$$
$$\left( P(x_1 \mid y) P(x_2 \mid y) \ldots \ldots P(x_v \mid y) \right) P(y)$$

$$= \prod_{i=1}^{v} \left[ P(x_i \mid y) \right] P(y)$$

$$= \prod_{i=1}^{v} \left[ P(x_i \mid y) \right] P(Y = y)$$

<u>4.1</u> b)

For Representation of $P(x|Y=y)$ where $x$ is feature vector and $y$ is associated label.

In Naive Bayes classified :-

→ Prior $P(Y) = k-1$ if we have 'k' classes

→ Liklihood $P(x|Y=y) = \underline{(2^d-1)k}$ for binary features
$$\downarrow$$

Naive Bayes.
(NO assumption)

If we hold the simplicity Assumption for NB classifier. which states that. given the class label $Y$, pair of feature $x_i$ and $x_j$ $(i \neq j)$ are conditionally independent.

$$P(x|Y=y) = \underbrace{P(x_1|Y)\,P(x_2|Y)\,P(x_3|Y) \cdots P(x_d|Y)}_{\text{'d'}}$$

So,
Liklihood $CP(x|Y=y)$ Now has $\underline{dk+k-1}$ parameters.
(With Assumptn)

Its good to make simplicity Assumption because it is a linear function $(dk+k+1)$ unlike a power-ed function that would rise exponentially if d value grows to million of features.

As a Result, if we have large amound of <u>feature vectors</u> in our dataset, it is <u>good to have</u> <u>Simplicity</u> Assumption true, due to computationals <u>efficiency</u>

4.1 C) Let $X_w$ are independent random vocabulary words where $w \in (0 \cdots v)$

$$P(X_w | Y = y) = \theta_{yw} \quad \& \quad P(X_w = 0 | Y = y) = 1 - \theta_{yw}$$

Acc to Burnoli Distribution:

$$p^{x_i}(1-p)^{1-x_i}$$

$$P(X_w | Y = y)^{x_i} \quad P(X_w = 0 | Y = y)^{1-x_i}$$

Likelihood function

$$\mathcal{L}(p) = \prod_{w=1} P(X_w | Y = y)^{x_w} \, P(X_w = 0 | Y = y)^{1-x_w}$$

$$\mathcal{L}(p) = P(x_1 | Y = y)^{x_1} \, P(x = 0 | Y = y)^{x_1} \cdot P(x_2 | Y = y)^{x_2}$$
$$\cdot P(x_2 = 0 | Y = y)^{1-x_2}$$
$$\cdot P(x_w | Y = y)^{x_w}$$
$$\cdots \cdots P(x_w = 0 | Y = y)^{1-x_w}$$

$$\mathcal{L}(p) = P(x_i | Y = y)^{\Sigma x_i} \, P(x_i = 0 | Y = y)^{\Sigma(1-x_i)}$$

$$\mathcal{L}(p) = P(x_i | Y = y)^{\Sigma x_i} \, P(x_i = 0 | Y = y)^{n - \Sigma x_i}$$

To maximize likelihood $\dfrac{\partial \mathcal{L}(p)}{\partial p} = 0$

Log Transformation:

$$\log \mathcal{L}(p) = \Sigma x_i \log P(x_i | Y = y) + (-1)(n - \Sigma x_i) \log P(x_i = 0 | Y = y)$$

$$\frac{\partial \log(\mathcal{L}(p))}{\partial p} = \frac{\Sigma x_i}{P(x_i | Y = y)} + \frac{n - \Sigma x_i}{P(x_i = 0 | Y = y)} = 0$$

$$\Sigma x_i \, (P(x_i = 0 | Y = y)) = n - \Sigma x_i \, (P(x_i | Y = y))$$
$$\Sigma x_i \, (P(x_i = 0 | Y = y)) = n P(x_i | Y = y) - \Sigma x_i \, P(x_i | Y = y)$$
$$\Sigma x_i \, (P(x_i = 0 | Y = y) + \Sigma x_i \, P(x_i | Y = y) = n P(x_i | Y = y)$$
$$\Sigma x_i (1) = n P(x_i | Y = y) \quad (\text{as } 1 - \theta_{yw} + \theta_{yw} = 1)$$
$$\sum_{i=0}^{M} \frac{\Sigma x_i}{n} = \hat{P}(x_i | Y = y)$$

$P(Y=1) = P$     $P(Y=2) = 1-P$

Burnolli Representation :-

$$(P)^{x_1} (1-P)^{x_2}$$

Likehood function:

$$\mathcal{L} = P(Y=1)^{x_1} \; P(Y=0)^{x_2}$$

$$P^{x_1} (1-P)^{x_2}$$

Taking derivative of $\mathcal{L}$ to maximize

$$\frac{\partial \mathcal{L}(\theta)}{\partial \theta} = 0$$

But we Need to log transform first

$$\log(\mathcal{L}) = x_1 \log(P) + x_2 (\log(1-P) \times (-1)$$

$$\frac{d \log(\theta)}{d(\theta)} = \frac{x_1}{P} - \frac{x_2}{1-P} = 0$$

$$\frac{x_1}{P} = \frac{x_2}{1-P}$$

$$x_1 - x_1 P = x_2 P$$

$$x_1 = x_2 P + x_1 P$$

$$x_1 = P(x_1 + x_2)$$

$$\boxed{P = \frac{x_1}{x_1 + x_2}}$$

## 4.2

4 features : $(x_1, x_2, x_3, x_4)$

3 labels : $(+1, 0, -1)$

a) For prior probability

$$P(Y = +1) = \frac{2}{7}$$

$$P(Y = 0) = \frac{2}{7}$$

$$P(Y = -1) = \frac{3}{7}$$

Because

$\hat{P}(y = y_k)$

$\left\{ \dfrac{\text{No. of } \hat{y} = y}{\text{Total values in Data.}} \right\}$

b)  $P((x_1 = 1) \mid Y = +1) = \frac{1}{2}$

$P(x_2 \mid Y = +1) = 1$

$P(x_3 \mid Y = +1) = \frac{1}{2}$

$P(x_4 \mid Y = +1) = \frac{1}{2}$


$P(x_1 = 1 \mid Y = 0) = \frac{1}{2}$

$P(x_2 = 1 \mid Y = 0) = \frac{1}{2}$

$P(x_3 = 1 \mid Y = 0) = 1$

$P(x_4 = 1 \mid Y = 0) = 1$


$P(x_1 = 1 \mid Y = -1) = \frac{1}{3}$

$P(x_2 = 1 \mid Y = -1) = \frac{1}{3}$

$P(x_3 = 1 \mid Y = -1) = \frac{1}{3}$

$P(x_4 = 1 \mid Y = -1) = \frac{2}{3}$

| $x_i$ $y_i \rightarrow$ | $y=+1$ | $y=0$ | $y=-1$ |
|---|---|---|---|
| $x_1=1$ | 1/2 | 1/2 | 1/3 |
| $x_2=1$ | 1 | 1/2 | 1/3 |
| $x_3=1$ | 1/2 | 1 | 1/3 |
| $x_4=1$ | 1/2 | 1 | 2/3 |

c) The given date point with feature $X_1, X_2, X_3, X_4$

$$\left( \prod_{i=1}^{V} P(x|y=y) \right) P(Y=y)$$

If $y=0$ → $P(X|y=0) \times P(x_2|y=0) \times P(x_3|y=0) \times P(x_4|y=0) \ P(y=0)$

→ $\frac{1}{2} \times \frac{1}{2} \times \times 1 \times 1 \times 2/7$

→ $1/14$

If $y=+1$ → $P(x_1|y=+1) \times P(x_2|y=+1) \times P(x_3|y=+1) \times P(x_4|y=+1) P(y=1)$

→ $\frac{1}{2} \times \frac{1}{2} \times 1 \times \frac{1}{2} \times 2/7$

→ $1/28$

If $y=-1$ → $P(x_1|y=-1) \times P(x_2|y=-1) \times P(x_3|y=-1) \times P(x_4|y=-1) \ P(y=-1)$

→ $\frac{1}{3} + \frac{1}{3} \times \frac{1}{3} \times \frac{2}{3} \times 3/7$

→ $2/189$

As.

Objective function states to classify in the class which has maximum probability

From above calculation we find that

$$\frac{1}{14} > \frac{1}{28} > \frac{2}{189}$$

So, The data point should be classified into
y=0 class (Max probability)