

①

COL 7/4

Machine Learning
Major Exam

①

Elliptical SVM: 2 Aug 31, 2020

Q.2.

$$h_0(x) = 1 \left\{ \sum_{j=1}^n \frac{x_j^2}{a_j^2} \leq 1 \right\} \quad \Leftrightarrow \text{corresponds to elliptical decision boundary}$$

$$\left[\sum_{j=1}^n x_j^2 / a_j^2 = 1 \right]$$

 $\phi(x_j) =$

let $x_j^2 = t_j$ & $1/a_j^2 = w_j$ ($w_j \geq 0$) — (1)

$\phi(x)$ kernel fn

Then using the transformation defined in (1) we have:

$$\sum_{j=1}^n w_j t_j = 1 \quad \text{as the decision boundary.}$$

Equivalently, we have

$$\sum_{j=1}^n w_j t_j - 1 = 0 \quad \Rightarrow \text{equation of a hyperplane}$$

which can equivalently be written as:

$$\boxed{w'^T t + b = 0}$$

[using the fact that $t_j: -w_j' = w_j * (-b)$]
 $(-b) > 0$

There is an additional scale factor here that we can exploit

$$\Rightarrow 1 \sum w_j' t + b \geq 0$$

$$\Rightarrow h_0(x)$$

$$\& 2 \sum w_j' t + b < 0$$

as -ve class

(wlog we can call above as +ve class)
 [we are flipping the signs of 2 classes]

with additional constraint that

$$\boxed{w' \geq 0}$$

Note: - Since data is a separable using elliptical boundary, to & $\hat{y} \geq 0$ we do not need to explicitly enforce the constraint $-b \geq 0$ (think why).

\Rightarrow converting this into SVM form we get

$$\min_{w, b} \quad \frac{1}{2} w'^T w'$$

$$y^{(i)} [w'^T t^{(i)} + b] \geq 1 \quad \forall i \quad y^{(i)} \in \{1, -1\}$$

\rightarrow Lagrange multipliers α_i

Therefore, Lagrangian ~~multiplier~~ can be written as -

$w' \geq 0$ — Lagrange multipliers: $-\beta$ ($\beta_1 \dots \beta_n$)

$$L(w, b, \alpha, \beta) = \frac{1}{2} w'^T w' + \sum_{i=1}^m \alpha_i [1 - y^{(i)} (w'^T t^{(i)} + b)]$$

Now, to get the dual we will equate gradient wrt primal variables to 0.

$$\nabla_w L(w, b, \alpha, \beta) = w' + \sum_{i=1}^m \alpha_i y^{(i)} t^{(i)} + (-1) \beta = 0$$

Equating this to zero, we get

$$\boxed{w' = \sum_{i=1}^m \alpha_i y^{(i)} t^{(i)} + \beta}$$

Similarly,

$$\nabla_b L(w, b, \alpha, \beta) = \sum_{i=1}^m \alpha_i y^{(i)} (-1) = 0$$

$$\Rightarrow \boxed{\sum_{i=1}^m \alpha_i y^{(i)} = 0}$$

② substituting the value of w' in the Lagrangian, we get:-

$$L(w, b, \alpha, \beta) = \frac{1}{2} \left[\sum_{i=1}^m \alpha_i y^{(i)} t^{(w)} + \beta \right]^T \left[\sum_{i=1}^m \alpha_i y^{(i)} t^{(w)} + \beta \right] + \sum_{i=1}^m \alpha_i - \sum_{i=1}^m \alpha_i y^{(i)} \left[\sum_{i=1}^m \alpha_i y^{(i)} t^{(w)} + \beta \right] t^{(w)} + b - \beta^T \left[\sum_{i=1}^m \alpha_i y^{(i)} t^{(w)} + \beta \right]$$

$$= \frac{1}{2} \left[\sum_{i,j=1}^m \alpha_i \alpha_j y^{(i)} y^{(j)} t^{(w)T} t^{(j)} \right] + \frac{1}{2} [\beta^T \beta]$$

$$+ \cancel{\frac{1}{2} \beta^T \sum_{i=1}^m \alpha_i y^{(i)} t^{(w)}} + \sum_{i=1}^m \alpha_i$$

$$- \sum_{i=1}^m \alpha_i y^{(i)} \left[\sum_{j=1}^m \alpha_j y^{(j)} t^{(j)} + \beta \right]^T t^{(w)} + b$$

$$- \beta^T \sum_{i=1}^m \alpha_i y^{(i)} t^{(w)} + (-1) \beta^T \beta$$

$$= \sum_{i=1}^m \alpha_i - \frac{1}{2} \sum_{i,j=1}^m \alpha_i \alpha_j y^{(i)} y^{(j)} t^{(w)T} t^{(j)} - \frac{1}{2} \beta^T \beta$$

$$- \sum_{i=1}^m [\beta^T \alpha_i y^{(i)} t^{(w)}] - b \sum_{i=1}^m \alpha_i y^{(i)}$$

$D(\alpha, \beta)$

$$\equiv \sum_{i=1}^m \alpha_i - \frac{1}{2} \beta^T \beta - \frac{1}{2} \sum_{i,j=1}^m \alpha_i \alpha_j y^{(i)} y^{(j)} t^{(w)T} t^{(j)} - \beta^T \sum_{i=1}^m \alpha_i y^{(i)} t^{(w)}$$

Dual: -

$$\max_{\alpha, \beta}$$

$$D(\alpha, \beta)$$

$$\alpha, \beta \quad \alpha \geq 0, \beta \geq 0$$

$$\sum_{i=1}^m \alpha_i y_i = 0$$

Q.2. PCA (generalized)

From standard PCA, we know that the PCA components can be obtained by solving

$$(a) \max_{u \perp \phi} u^T \left[\sum_{i=1}^m \frac{\phi(x^{(i)}) \phi(x^{(i)})^T}{m} \right] u \equiv \sum \phi$$

in the transformed space (one component at a time).

This can be done using (finding) / solving eigen the following eigenproblem: value

$$\sum \phi u \phi = \lambda u \phi$$

with $\sum \phi = \sum_{i=1}^m \frac{\phi(x^{(i)}) \phi(x^{(i)})^T}{m}$

substituting this in the above equation, we get

$$\sum_{i=1}^m \frac{\phi(x^{(i)}) \phi(x^{(i)})^T}{m} u \phi = \lambda u \phi$$

(3)

$$\sum_{i=1}^m \phi(x^{(i)}) \frac{[\phi(x^{(i)})^T u_p \phi]}{m} = \lambda u_p \phi$$

$$\Rightarrow u_p \phi = \frac{1}{\lambda m} \left[\sum_{i=1}^m \phi(x^{(i)}) [\phi(x^{(i)})^T u_p \phi] \right]$$

\Rightarrow Desired coefficients are given as:-

$$\boxed{a_p^i = \frac{\phi(x^{(i)})^T u_p \phi}{\lambda m}} \quad \text{s.t.}$$

$$u_p \phi = \sum_{i=1}^m a_p^i \phi(x^{(i)})$$

(b) We have:-

$$u_p \phi = \frac{1}{\lambda m} \sum_{i=1}^m a_p^i \phi(x^{(i)})$$

Pre-multiplying both sides by $\phi(x^{(j)})^T$ for some j ,

we get

$$\phi(x^{(j)})^T u_p \phi = \frac{1}{\lambda m} \phi(x^{(j)})^T \sum_{i=1}^m a_p^i \phi(x^{(i)})$$

$$\Downarrow$$

$$[a_p^j, \lambda m]$$

$$\Rightarrow a_p^j \cancel{\lambda m} = \frac{1}{\lambda m} \sum_{i=1}^m a_p^i \underbrace{\phi(x^{(j)})^T \phi(x^{(i)})}_{K(x^{(j)}, x^{(i)})}$$

$$\boxed{\lambda a_p^j = \frac{1}{m} \sum_{i=1}^m a_p^i K(x^{(j)}, x^{(i)})}$$

This equation
gives an
individual
equation for
 a_p^j

(c) Continuing further we get (using simple matrix algebra)

$$\boxed{\lambda a_p = \frac{1}{m} K^M a_p}$$

[λ is some scalar quantity]

$\Rightarrow a_p$ is eigenvector of the matrix

$$\begin{bmatrix} K^M \\ m \end{bmatrix}$$

& λ corresponding eigenvalue.

$$K^M = \begin{bmatrix} \text{---} & k(x^i, x^j) \end{bmatrix}$$

is the kernel matrix.

Finally we have to ensure that

$$\underline{u_p^T u_p = 1}$$

$$\Rightarrow \sum_{i=1}^m a_p^i \phi(x^i)^T \sum_{j=1}^m a_p^j \phi(x^j) = 1$$

$$\Rightarrow \sum_{i,j=1}^m a_p^i \phi(x^i)^T \phi(x^j) a_p^j = 1$$

$$\Rightarrow \boxed{a_p^T K^M a_p = 1}$$

Putting this in equation

$$a_p = \frac{1}{\lambda m} K^M a_p \quad (\text{Pre-multiply by } a_p^T)$$

$$a_p^T a_p = \frac{1}{\lambda m} a_p^T K^M a_p$$

$$\text{we get } \boxed{a_p^T a_p = \frac{1}{\lambda m} \cdot 1} \quad \square$$

$$\begin{aligned} (d) \quad & \phi(x^i)^T u_p \\ &= \phi(x^i)^T \sum_{j=1}^m a_p^j \phi(x^j) \\ &= \sum_{j=1}^m a_p^j \underbrace{\phi(x^i)^T \phi(x^j)}_{k(x^i, x^j)} \\ & \text{ } a_p^j \text{ can be computed in terms of } K \text{ function} \\ & \Rightarrow \phi(x^i)^T u_p \text{ can be computed in terms of } K \text{ fn.} \end{aligned}$$

Q3. (a) In general. For m points. $LL(\theta) =$ for logistic regression

$$\nabla_0 LL(\theta) = \nabla_0 \sum_{i=1}^m y_i \log \left[\frac{1}{1 + e^{-\theta^T x_i}} \right]$$

$$\begin{aligned} LL(\theta) &= \sum_{i=1}^m y_i \left[\log \frac{1}{1 + e^{-\theta^T x_i}} \right] + \sum_{i=1}^m (1 - y_i) \left[\log \frac{e^{-\theta^T x_i}}{1 + e^{-\theta^T x_i}} \right] \\ &\Rightarrow -\nabla_0 \left[\sum_{i=1}^m y_i \log [1 + e^{-\theta^T x_i}] + \sum_{i=1}^m (1 - y_i) \theta^T x_i \right] \end{aligned}$$

$$= -\nabla_0 \left[\sum_{i=1}^m \log [1 + e^{-\theta^T x_i}] + \sum_{i=1}^m (1 - y_i) \theta^T x_i \right]$$

$$= \sum_{i=1}^m \left[\frac{e^{-\theta^T x_i}}{1 + e^{-\theta^T x_i}} x_i - (1 - y_i) x_i \right]$$

$$\boxed{\nabla_0 LL(\theta) = \sum_{i=1}^m \left[y_i - \frac{1}{1 + e^{-\theta^T x_i}} \right] x_i}$$

$$\nabla_0^2 LL(\theta) = \sum_{i=1}^m \frac{e^{-\theta^T x_i}}{(1 + e^{-\theta^T x_i})^2} x_i x_i^T$$

$$\nabla_0^2 LL(\theta) = \sum_{i=1}^m \frac{1}{(1 + e^{-\theta^T x_i})^2} x_i x_i^T$$

This is the same as the sec. defn.

Since -

$$- \sum_{l=1}^m u^T K^{(l)} x^{(l)} x^{(l)T} u$$

+ve number

$$= - \sum_{l=1}^m K^{(l)} [u^T x^{(l)}]^2 \leq 0$$

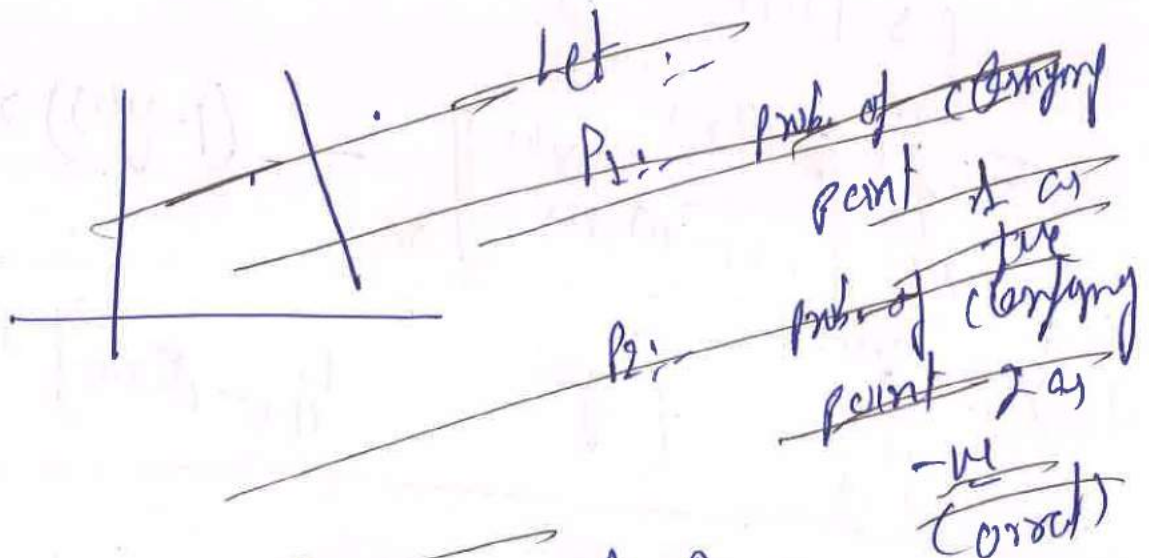
$\Rightarrow L(u)$ is concave.

- $L(u)$ is ~~convex~~ convex

Holds: - for m points

\Rightarrow holds for m=2 points.

(b)



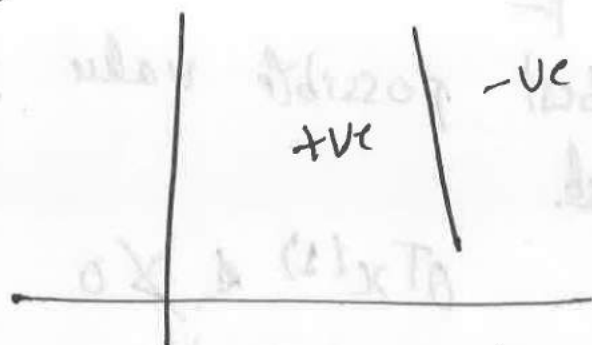
$$L(u) = \log P_1 + \log P_2$$

$$\Rightarrow \text{Dist} = -\alpha$$

$$\text{Dist} = -\eta \epsilon$$

5)

Q.3 (b)



Two points in \mathbb{R}^2 for visualization (only).

First, we will show that optimal decision boundary must classify the 2 points correctly.

$$LL(\theta) = \log \left[\frac{1}{1 + e^{-\theta^T x^{(1)}}} \right] + \log \left[\frac{e^{-\theta^T x^{(2)}}}{1 + e^{-\theta^T x^{(2)}}} \right]$$

\rightarrow likelihood for 1st point likelihood for 2nd point

$$= \log \left[\frac{1}{1 + e^{-\theta^T x^{(1)}}} \right] + \log \left[\frac{1}{1 + e^{\theta^T x^{(2)}}} \right]$$

To maximize this, we would need to maximize both terms

clearly if $\theta^T x^{(1)} > 0$ & $\theta^T x^{(2)} < 0$
 (i.e. decision boundary is between the two points)

we can make

by scaling θ arbitrarily
 the decision surface] $\theta^T x^{(1)} \rightarrow \infty$ & $\theta^T x^{(2)} \rightarrow -\infty$
 (not changing θ)
 Replace θ by $k\theta$
 $k > 0$

as $k \rightarrow \infty$

$$L(0) \rightarrow 1 + \frac{2}{1} = 1 + 1 = 2$$

This is the best possible value that can be achieved.

Note that if $\nabla T(x(1)) \neq 0$

or $\nabla T(x(2)) \neq 0$

i.e. the points are separated by the boundary (with the +ve & -ve on correct sides): -

Then at least one of the terms will not be able to reach 1

i.e. if $\nabla T(x(i)) < 0$ then
maximum value of $\frac{1}{1 + e^{-\nabla T(x(i))}}$

$$\text{is } \frac{1}{2} = 0.5$$

\Rightarrow i.e. $L(0) \leq 1.5$ in this

\Rightarrow Q can not be optimal.

~~(c) For using part (b) optimal class~~

Therefore, for the optimal boundary, it must be somewhere between the +ve & -ve points, specifically the boundary, with equal distance from the +ve & -ve points will be optimal. (one of the)

② Maximum value of $L(\theta)$ in this setting as shown in part (b) above is 2.
But it is only in the limit

$$\left[\frac{1}{1+e^{-\theta^T x^{(1)}k}} + \frac{1}{1+e^{+\theta^T x^{(2)}k}} \right] = L(\theta)$$

$$\underline{k \rightarrow \infty}$$

But for any finite value of k

$$\boxed{L(\theta) \neq 2}$$

But that is not an issue since, during practical implementation, we stop learning when

$$|L(\theta^{(t+1)}) - L(\theta^{(t)})| < \epsilon$$

Alternatively, we can write $\nabla L(\theta) \rightarrow 0$

for some ϵ . Thus will

happen since $L(\theta) \rightarrow 2$ in the limit $\Rightarrow \theta^{(t)}$ & $\theta^{(t+1)}$ as we take gradient steps such that

$L(\theta^{(t)})$ is as close to 2 as possible

$$\nabla L(\theta) = \sum_{i=1}^2 \left[y^{(i)} - h_{\theta}(x^{(i)}) \right] x^{(i)}$$

$h_{\theta}(x^{(1)}) \rightarrow 1$ or 0 as the case might be
 $h_{\theta}(x^{(2)}) = \left[\frac{1}{1+e^{-\theta^T x^{(2)}}} \right]$ as $\theta \rightarrow k \rightarrow \infty$
 $h_{\theta}(x^{(2)}) \rightarrow 1$

as $k \rightarrow \infty$ $\log(1) + \log(1) = 0 + 0 = 0$
 $L(0) \rightarrow 1 + \frac{1}{1} = 2$

This is the best possible value that can be achieved.

Note that if $\partial T x(1) \neq 0$

or $\partial T x(2) \neq 0$

i.e. the points are separated by the boundary (with the +ve & -ve on correct sides): -

Then at least one of the terms will not be able to reach 1

i.e. if $\partial T x(i) < 0$ then
 maximum value of $\frac{1}{1 + e^{-\partial T x(i)}}$

is $1/2 = 0.5$

\Rightarrow (for $L(0) \leq 2 \times 0.5$) in this case (other term can be at most 0)

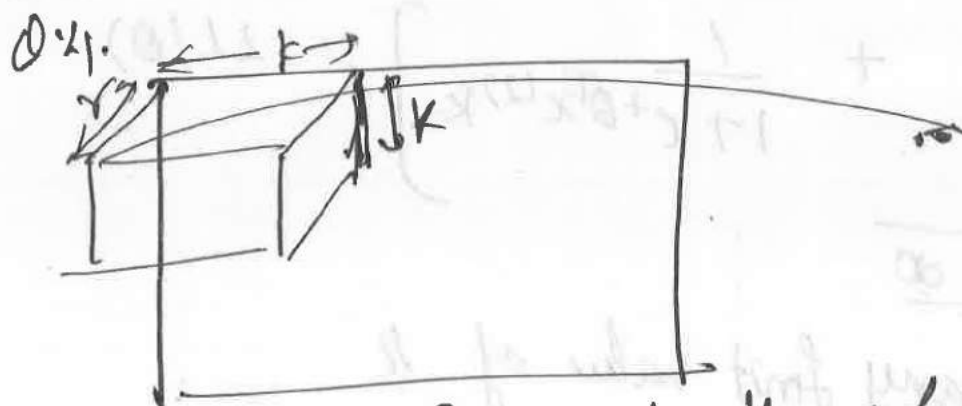
\Rightarrow Q can not be optimal.

~~For using part (b) optimal class~~

Therefore, for the optimal boundary, it must be somewhere between the two +ve points, specifically the boundary, with equal distance from the two -ve points will be optimal. (one of the)

gradient goes to zero similarly for the second term

$$L_{ex} : \{ h_o(x^{(i)}) = \frac{1}{1 + e^{-\theta^T x^{(i)}}} \}$$



let us first write the relevant equations:-

let the kernel be denoted by \underline{K}

i, j vary over width / height of the image

let $x_{i,j,l}^{(t)}$ denote the feature map at (i,j) position & depth l ($1 \leq i, j \leq W$)

$$1 \leq l \leq P$$

P :- kernel index

$$1 \leq P \leq \gamma$$

$x^{(t)}$:- Feature after convolution then (before average pooling)

Feature map before the pooling layer

$$\left[\sum_{l=1}^P K_{u,v,l} \cdot x_{i+u, j+v, l}^{(t)} + b^p \right] = x_{i,j,p}^{(t)}$$

$1 \leq i, j \leq W$
 $1 \leq l \leq \gamma$
 $1 \leq P \leq \gamma$

① Average Pooling:-

$$x_{i,j,p}^{(t+1)} = \frac{\sum_{\alpha,\beta=1}^2 x'_{i+\alpha,j+\beta,p}(t)}{4}$$

Then:- we have: \downarrow
(Gradient wrt ~~pool~~ feature map going into pooling layer)

$$\textcircled{1} \quad \frac{\partial L}{\partial x'_{i,j,p}(t)} = \sum_{\alpha,\beta=1}^2 \frac{\partial L}{\partial x_{i+\alpha,j+\beta,p}^{(t+1)}} \cdot \left[\frac{\partial x_{i+\alpha,j+\beta,p}^{(t+1)}}{\partial x'_{i,j,p}(t)} \right]$$

$$= \frac{1}{4} \sum_{\alpha,\beta=1}^2 \frac{\partial L}{\partial x_{i+\alpha,j+\beta,p}^{(t+1)}} \quad \text{From } \textcircled{2} \quad \Downarrow \quad \frac{1}{4} \quad \text{--- (1)}$$

~~the~~ Gradient going into ~~on the layer before~~ part the feature map convolution ~~operation~~ layer.

$$\textcircled{2} \quad \frac{\partial L}{\partial x_{i,j,p}^{(t)}} = \sum_{u,v=1}^K \sum_{\alpha,\beta=1}^2 \left[\frac{\partial L}{\partial x_{i+\alpha,j+\beta,p}^{(t)}} \right] \left[\frac{\partial x_{i+\alpha,j+\beta,p}^{(t)}}{\partial x_{i,j,p}^{(t)}} \right]$$

\Downarrow from Eq ①

~~from~~ Eq ①

When ~~t=t~~ t=t we know that

$$\left[\frac{\partial L}{\partial x^{(t)}} \right] \text{ is given to us}$$

$$\boxed{\frac{\partial x^{(t)}_{i,j,p}}{\partial x^{(t)}_{i+u,j+v,e}} = \cancel{\frac{\partial x^{(t)}_{i,j,p}}{\partial K_{u,v,e}}} \cdot K_{u,v,e}^p}$$

Finally,

$$\textcircled{3} \quad \frac{\partial L}{\partial K_{u,v,e}} = \sum_{t=1}^d \left[\sum_{i,j=1}^W \left[\frac{\partial L}{\partial x^{(t)}_{i,j,p}} \right] \cdot \left[\frac{\partial x^{(t)}_{i,j,p}}{\partial K_{u,v,e}} \right] \right] \quad \text{From } \textcircled{1}$$

$$\boxed{\frac{\partial x^{(t)}_{i,j,p}}{\partial K_{u,v,e}} = \cancel{\frac{\partial x^{(t)}_{i,j,p}}{\partial x^{(t)}_{i+u,j+v,e}}} \cdot x^{(t)}_{i+u,j+v,e}}$$

$$\textcircled{4} \quad \frac{\partial L}{\partial b^p} = \sum_{t=1}^d \left[\sum_{i,j=1}^W \frac{\partial L}{\partial x^{(t)}_{i,j,p}} \cdot \frac{\partial x^{(t)}_{i,j,p}}{\partial b^p} \right] \quad \text{From } \textcircled{1}$$

$$\boxed{\frac{\partial x^{(t)}_{i,j,p}}{\partial b^p} = 1}$$

Need to double check the computation

⑧

0.5 - Generative Process -

$$z^{(i)} \sim \text{Multinomial}(\Phi) \quad z^{(i)} \in \{1, \dots, K\}$$

$$y^{(i)} \sim \text{Bernoulli}(p) \quad \therefore \text{does not depend on } z^{(i)} \text{ (as given in the question)}$$

$$\forall i: - x_j^{(i)} | y_j^{(i)}, z^{(i)} = l \quad l \in \{1, \dots, K\} \quad t \in \{1, 0\}$$

$$\equiv \text{Bernoulli}(\theta_{j|y=y^{(i)}, z=z^{(i)}})$$

$$\equiv \text{Bernoulli}(\theta_{j|y=t, z=l})$$

Parameters of the Model:

$$\Phi: - (\Phi_1, \dots, \Phi_K) \quad \therefore \text{latent mixture parameters} \quad K \text{ parameters (one redundant)}$$

$$\theta \Rightarrow \sum \Phi_l = 1 \quad \phi: - \text{prior class probability. (one parameter)}$$

$$\left\{ \theta_{j|y=t, z=l} \right\}_{j=1, t=0, l=1}^{n, 1, K} \quad n \times 2 \times K = 2nK \text{ parameters}$$

Log-likelihood

$$\sum_{i=1}^n \log [P(x^i, y^i; \theta)]$$

$$= \sum_{i=1}^n \log [P(x^i | y^i; \theta) P(y^i; \phi)]$$

$$= \sum_{i=1}^n \log \left[\sum_{z^i} [P(x^i | y^i, z^i; \theta)] \cdot P(z^i; \Phi) \right]$$

$$= \sum_{i=1}^n \left[\log \left[\sum_{z^i} [P(x^i | y^i, z^i; \theta)] \right] + \log P(y^i; \phi) \right]$$

[since y^i is independent of z^i , $P(y^i; \phi)$ can be expressed in terms of parameters Φ]

$$= \sum_{i=1}^n \log \left[\sum_{z^i} \left[\prod_{j=1}^n P(x_j^i | y^i, z^i; \theta) \right] P(z^i; \Phi) + \log P(y^i; \phi) \right]$$

can be expressed in terms of parameters of the model

↓
can be expressed in terms of parameters

(skipping the parameters in place of corresponding substitutions)

$$= \sum_{i=1}^n \log \left[\sum_{z^i} \{ \prod_{j=1}^n P(x_j^i | y^i, z^i; \theta) \} P(z^i; \Phi) + \log P(y^i; \phi) \right]$$

$\sum_{j=1}^n y_j^i = 1$
 $\sum_{j=1}^n x_j^i = 1$
 $\sum_{j=1}^n (1 - x_j^i) = 1$

$$\begin{aligned}
 & + 1\{y^{(i)}=0\} \left[1\{x_j^{(i)}=1\} \theta_j | y=0, z^{(i)}=0 \right. \\
 & \quad \left. + 1\{x_j^{(i)}=0\} (1-\theta_j) | y=0, z^{(i)}=0 \right] \\
 & + \sum_{i=1}^n y^{(i)} \phi + (-y^{(i)}) (1-\phi)
 \end{aligned}$$

(b) EM based estimation

E-step: \rightarrow

$$p(z^{(i)} | x^{(i)}, y^{(i)}; \theta)$$

$$= \frac{p(z^{(i)}, x^{(i)}, y^{(i)}; \theta)}{p(x^{(i)}, y^{(i)}; \theta)}$$

$$= \frac{p(x^{(i)} | y^{(i)}, z^{(i)}; \theta) p(y^{(i)}; \phi) p(z^{(i)}; \Phi)}{\sum_{z^{(i)}} p(x^{(i)}, y^{(i)} | z^{(i)}; \theta) p(z^{(i)}; \phi)}$$

$$= \frac{\left[\prod_{j=1}^n p(x_j^{(i)} | y^{(i)}, z^{(i)}; \theta) \right] \cdot p(y^{(i)}; \phi) p(z^{(i)}; \Phi)}{\sum_{z^{(i)}} \left[p(x^{(i)} | y^{(i)}, z^{(i)}; \theta) \right] p(y^{(i)}; \phi) p(z^{(i)}; \Phi)}$$

We know how to compute Φ for given set of parameters.

M-step - Given $p(z^{(i)} = 0)$, we want to estimate the parameters of the naive Bayes model in the next step.

→ independent of $z^{(i)}$

$$\phi = \frac{\sum_{i=1}^m 1 \{y^{(i)} = t\}}{m}$$

constant through out

$$\Phi = \frac{\sum_{l=1}^m P(z^{(i)} = l)}{m} \quad \left[\text{softer version of } \frac{\sum_{i=1}^m 1 \{z^{(i)} = l\}}{m} \right]$$

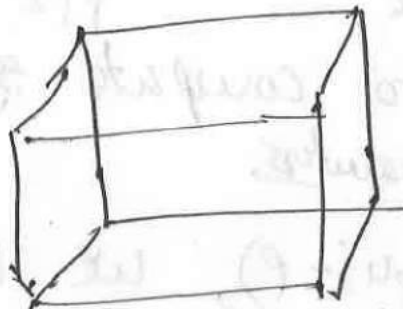
M-step

$$Q_j | y=t, z^{(i)}=l = \frac{\sum_{i=1}^m P(z^{(i)}=l) 1 \{y^{(i)}=t\}}{\sum_{i=1}^m 1 \{y^{(i)}=t\} P(z^{(i)}=l)}$$

softer version of naive Bayes parameter estimation
[can do Laplace smoothing] by

adding 1 to the numerator
& 2 to the denominator

Q8: VC H_{11}^3 :-



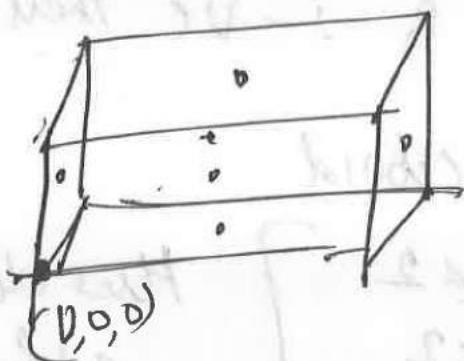
claim:-
VC-dimension of
class H_{11}^3 (axis
11 cuboids)
= 6

Recall:- $VC(H_{11}^2) = 4$

(10) First, we will show that $VC(H_1^3) = 6$

Consider a cuboid of s length, width & height = 2 with bottom left inner corner at $(0,0,0)$

$$0 \leq x \leq 2, 0 \leq y \leq 2, 0 \leq z \leq 2$$



Now, let us place six points at the center of six faces of the cuboid i.e.,

at

$$(\frac{1}{2}, \frac{1}{2}, 0)$$

$$(0, \frac{1}{2}, 0)$$

$$(0, 0, \frac{1}{2})$$

$$(\frac{1}{2}, \frac{1}{2}, 0)$$

$$(\frac{1}{2}, 0, \frac{1}{2})$$

$$(\frac{1}{2}, \frac{1}{2}, \frac{1}{2})$$

$$P_1 (1, 1, 0)$$

$$P_2 (1, 1, 2)$$

top & bottom face

$$P_3 (1, 0, 1)$$

$$P_4 (1, 2, 1)$$

front & back face

$$P_5 (0, 1, 1)$$

$$P_6 (2, 1, 1)$$

left & right face

Then these six points can be easily shattered by moving the plane

at which the point lies by $\pm \epsilon$ ($\epsilon > 0$)
 to give the concerned point a label of
 $+ve$ or $-ve$ (or vice-versa) where ϵ is
 very small

e.g. Consider $P_1: - (1, 1, 0)$

the ~~is~~ if $P_1: +ve$ then

In the cuboid,

$$-\epsilon \leq x \leq 2$$

$$0 \leq y \leq 2$$

$$0 \leq z \leq 2$$

Similarly if $P_1: -ve$ then

In the cuboid

$$\epsilon \leq x \leq 2$$

$$0 \leq y \leq 2$$

$$0 \leq z \leq 2$$

This does not
 affect the
 label of other
 points

\Rightarrow All 96 labelings can be achieved
 by doing $\pm \epsilon$ movement of 6
 faces independently (corresponding
 to label of each point).

⑪

Next, we will show that $VC(H_{11}^3) \geq 6$

i.e. no set of 7 points can be shattered.

Assume Contrary:-

$\exists P_1, \dots, P_7$ s.t. H_{11}^3 can shatter

then we will show that 3 labeling of these 7 points that can not be achieved

let p_x^{\max} :- point with max x value

p_x^{\min} :- point with min x value

similarly define

p_y^{\max}, p_y^{\min} & p_z^{\max}, p_z^{\min}

\exists let b :- remaining point (at least one such point must be there which is not covered above)

Assign b :- VC label
Add other points the label

clearly:-

let $h \in H_{11}^3$ shd achieve this labeling.

$$p_x^{\min}[x] \leq p[x] \leq p_x^{\max}[x]$$

$$p_y^{\min}[y] \leq p[y] \leq p_y^{\max}[y]$$

$$p_z^{\min}[z] \leq p[z] \leq p_z^{\max}[z]$$

Since $p_x^{\min}, p_x^{\max}, p_y^{\min}, p_y^{\max}, p_z^{\min}, p_z^{\max}$ are assigned the label, any point ~~without~~ lying they will be inside the cuboid formed by p

$$p_x^{\min} \leq p' \leq p_x^{\max} \wedge p_y^{\min} \leq p' \leq p_y^{\max}$$

$$\wedge p_z^{\min} \leq p' \leq p_z^{\max} \text{ must also}$$

be inside. $\Rightarrow p$ must be

inside $h \in H_1^3$

$\Rightarrow p$ must have a true label contradiction.

⑥⑦ In general: $H_1^n := \lfloor VC\text{-Dim}(H_1^n) = 2n \rfloor$

A very similar argument follows

$2n$ points can be shattered by considering 2 points on two faces \parallel to each of n ~~axes~~ axes. & then making the ϵ argument as earlier

$2n+1$ points can not be shattered by: assigning the label to ϵ those points which max/min value along any dimension & assigning -ve label to remaining point.