

# COL 341 Minor 1

Chinmay Mittal

TOTAL POINTS

**65.5 / 70**

QUESTION 1

## Generative - Linear Regression 25 pts

### 1.1 Q1a 8 / 10

✓ + 10 pts Correct

+ 0 pts Incorrect

- 2 Point adjustment

- Had to mention about independence assumption in step  $P(D/w) = \text{Product}(P(x_i, y_i/w))$ .

### 1.2 Q1b 10 / 10

✓ + 10 pts Correct

+ 0 pts Click here to replace this description.

### 1.3 Q1c 5 / 5

✓ + 5 pts Correct

+ 0 pts incorrect/not attempted

- 0.5 pts others

+ 1 pts did not answer the question

- 2 pts Click here to replace this description.

- Note that the answer is partially correct.

You should have proven the hessian to be positive semidefinite.

2) If the hessian matrix is diagonal, then

diagonal elements of the hessian  $\geq 0$  iff it is positive semidefinite i.e., all its eigenvalues are  $\geq 0$ . This holds true because for a diagonal matrix, the diagonal elements are its eigenvalues.

3) However, in this case, the hessian matrix is not diagonal. So, you have to show that it is positive semidefinite in some other way.

QUESTION 2

## GNB - Logistic 20 pts

### 2.1 Q2a 4.5 / 5

✓ + 5 pts Correct

+ 3 pts Partially Correct

+ 0 pts Incorrect

- 0.5 Point adjustment

- The expression for Gaussian is wrong in the denominator.

### 2.2 Q2b 10 / 10

Click here to replace this description.

- 10 pts Incorrect/Unattempted

✓ - 0 pts correct

### 2.3 Q2c 5 / 5

+ 4.5 pts correct

+ 0 pts incorrect/not attempted

+ 3.5 pts partially correct

✓ + 5 pts correct

+ 1 pts [Click here to replace this description.](#)

#### QUESTION 3

##### 3 VC Dimension 8 / 10

+ 10 pts Correct

✓ + 8 pts *Mostly Correct - minor bug*

+ 4 pts Major bugs/ Significant explanation  
missing

+ 2 pts Incorrect - some reasonable attempt

+ 0 pts Not Attempted/ Incorrect

💬 VC-d is infinite.

#### QUESTION 4

##### 4 PLA 15 / 15

✓ + 15 pts Correct

+ 0 pts Incorrect or not attempted

+ 0 pts Correct but incomplete arguments

(partial marks awarded via 'Point Adjustment')

Student Name: CHIN MAY MITTAL

Entry Number: 2020CS10336



Department of Computer Science and Engineering  
Indian Institute of Technology Delhi  
COL341: Fundamentals of Machine Learning

## Minor 1

Time: 60 minutes

Maximum Marks: 70

Number of Questions: 4

**Instructions:** Please attempt all questions. If you feel any question/statement is ambiguous, please write your assumptions clearly, and then answer as per your assumptions.

Question	1(a)	1(b)	1(c)	2(a)	2(b)	2(c)	3	4	Total
Max Marks	10	10	5	5	10	5	10	15	70
Earned Marks									

1. In the class we discussed loss/error function for linear regression in a discriminative style. In this question we will try to derive the expression in a generative style. Consider the dataset  $\mathcal{D} = \{\mathbf{x}^{(i)}, y^{(i)}\}$ . Assume that the value of  $y$  is observed after corruption with a Gaussian noise,  $\epsilon \sim \mathcal{N}(0, \sigma^2)$ . In other words,  $\mathbb{P}(y | \mathbf{x}, \mathbf{w}) = \mathcal{N}(y | \mathbf{w}^\top \mathbf{x}, \sigma^2)$

(a) [10 marks] The maximum likelihood estimate of  $\mathbf{w}$  is given by:  $\mathbf{w}_{\text{MLE}}^* = \arg \max_{\mathbf{w}} \mathbb{P}(\mathcal{D} | \mathbf{w})$ . Show that,  $\mathbf{w}_{\text{MLE}}^* = \arg \min_{\mathbf{w}} \sum_{i=1}^N (y^{(i)} - \mathbf{w}^\top \mathbf{x}^{(i)})^2$

(b) [8+2 marks] One problem with maximum likelihood estimation is that it can result in over-fitting when the number of samples are small. In the assignment we suggested you to use ridge regression to ameliorate the problem. In this question we will show that ridge regression is equivalent to MAP estimation of  $\mathbf{w}$ . Assuming Gaussian prior over the weight vector  $\mathbf{w}$ , i.e., assuming  $\mathbb{P}(\mathbf{w}) = \prod_j \mathcal{N}(w_j | 0, \tau^2)$ , show that,  $\mathbf{w}_{\text{MAP}}^* = \arg \min_{\mathbf{w}} \left[ \lambda \|\mathbf{w}\|_2^2 + \sum_{i=1}^N (y^{(i)} - \mathbf{w}^\top \mathbf{x}^{(i)})^2 \right]$ . What is the value of  $\lambda$ ?

(c) [5 marks] Notice that the error function for the ridge regression doesn't have an implicit solution, unlike linear regression. Hence, we need to solve it through gradient descent. Can you prove that the function is convex?

$$\begin{aligned}
 p(\mathcal{D} | \mathbf{w}) &= \prod_{i=1}^d p(d_i | \mathbf{w}) = \prod_{i=1}^d p(y_i | \mathbf{w}, \mathbf{x}_i) \\
 &= \prod_{i=1}^d \mathcal{N}(y_i | \mathbf{w}^\top \mathbf{x}_i, \sigma^2)
 \end{aligned}$$

$$\begin{aligned}
 \mathbf{w}_{\text{MLE}}^* &= \arg \max_{\mathbf{w}} \left( \prod_{i=1}^d \mathcal{N}(y_i | \mathbf{w}^\top \mathbf{x}_i, \sigma^2) \right) \\
 &= \arg \min_{\mathbf{w}} \left( -\log \left( \prod_{i=1}^d \mathcal{N}(y_i | \mathbf{w}^\top \mathbf{x}_i, \sigma^2) \right) \right)
 \end{aligned}$$

$$\begin{aligned}
 \mathbf{w}_{MLE}^* &= \arg\min_{\mathbf{w}} \left( - \sum_{i=1}^d \log(y_i | \mathbf{w}^T \mathbf{x}_i, \sigma^2) \right) \\
 &= \arg\min_{\mathbf{w}} \left( - \sum_{i=1}^d \log \left( \frac{1}{\sqrt{2\pi}\sigma} \exp \left( -\frac{(y_i - \mathbf{w}^T \mathbf{x}_i)^2}{2\sigma^2} \right) \right) \right) \\
 &= \arg\min_{\mathbf{w}} \left( \sum_{i=1}^d \frac{(y_i - \mathbf{w}^T \mathbf{x}_i)^2}{\sigma^2} + \text{constant} \right) \\
 &= \arg\min_{\mathbf{w}} \left( \sum_{i=1}^d (y_i - \mathbf{w}^T \mathbf{x}_i)^2 \right) \quad \begin{array}{l} \text{constant} \\ \text{wrt } \sigma \end{array}
 \end{aligned}$$

(b).  $\mathbf{w}_{MAP}^* = \arg\max_{\mathbf{w}} P(\mathcal{D} | \mathbf{w}) P(\mathbf{w})$  (MAP estimate).

$$\begin{aligned}
 &= \arg\min_{\mathbf{w}} \left[ -\log(P(\mathcal{D} | \mathbf{w}) P(\mathbf{w})) \right] = \arg\min_{\mathbf{w}} \left[ -\log(P(\mathbf{w})) + \log(P(\mathcal{D} | \mathbf{w})) \right] \\
 &= \arg\min_{\mathbf{w}} \left[ -\log \left( \prod_j N(\mathbf{w}_j | 0, \tau^2) \right) + \sum_{i=1}^d \frac{(y_i - \mathbf{w}^T \mathbf{x}_i)^2}{\sigma^2} + K \right] \quad \begin{array}{l} \text{from 1a)} \\ \text{from 1a)} \end{array} \\
 &= \arg\min_{\mathbf{w}} \left[ - \sum_{j=1}^d \log \left( \exp \left( -\frac{(\mathbf{w}_j)^2}{2\tau^2} \right) \right) + \sum_{i=1}^d \frac{(y_i - \mathbf{w}^T \mathbf{x}_i)^2}{\sigma^2} + K \right] \\
 &= \arg\min_{\mathbf{w}} \left[ \sum_{j=1}^d \frac{\mathbf{w}_j^2}{\tau^2} + \sum_{i=1}^d \frac{(y_i - \mathbf{w}^T \mathbf{x}_i)^2}{\sigma^2} \right] \\
 &= \arg\min_{\mathbf{w}} \left[ K \|\mathbf{w}\|^2 + \sum_{i=1}^d (y_i - \mathbf{w}^T \mathbf{x}_i)^2 \right] \quad \begin{array}{l} K = \frac{\sigma^2}{\tau^2} \end{array}
 \end{aligned}$$

(c) consider the double derivative along any axis  $\Rightarrow$  Along any axis the function is convex.

$$\begin{aligned}
 \frac{\partial^2 E}{\partial^2 \mathbf{w}_i} &= \frac{\partial}{\partial \mathbf{w}_i} (2\lambda \mathbf{w}_i + 2(y - \mathbf{w}^T \mathbf{x})(-\mathbf{x}_i)) \\
 &= (2\lambda + 2\mathbf{x}_i^2) > 0
 \end{aligned}$$



2. We discussed Naive Bayes classifier in the class. In this question we will discuss Gaussian Naive Bayes, which considers the following data generation model:

- $y$  is Boolean, and governed by a Bernoulli distribution, with parameter  $\theta = \mathbb{P}(y = 1)$
- $\mathbf{x} = \langle x_1, \dots, x_d \rangle$ , where each  $x_i$  is a continuous random variable.
- For each  $x_i$ ,  $\mathbb{P}(x_i | y = y_k)$  is a Gaussian distribution, i.e.  $\mathbb{P}(x_i | y = y_k) = \mathcal{N}(\mu_{ik}, \sigma_i^2)$ . Note here we are assuming the standard deviations  $\sigma_i$  vary from feature to feature, but do not depend on  $y$
- For all  $i$  and  $j \neq i$ ,  $x_i$  and  $x_j$  are conditionally independent given  $y$ .

(a) [5 marks] Show that  $\log \left[ \frac{\mathbb{P}(x_i | y=0)}{\mathbb{P}(x_i | y=1)} \right] = \frac{\mu_{i0} - \mu_{i1}}{\sigma_i^2} x_i + \frac{\mu_{i1}^2 - \mu_{i0}^2}{2\sigma_i^2}$

(b) [10 marks] Show that  $\mathbb{P}(y = 1 | \mathbf{x}) = \frac{1}{1 + \exp\left(\log \frac{1-\theta}{\theta} + \sum_{i=1}^d \log \left[ \frac{\mathbb{P}(x_i | y=0)}{\mathbb{P}(x_i | y=1)} \right]\right)}$

(c) [5 marks] Recall that in the logistic regression we assume that  $\mathbb{P}(y = 1 | \mathbf{x}) = \frac{1}{1 + \exp(w_0 + \sum_{i=1}^d w_i x_i)}$ . Using results from part (a) and (b) above, it is easy to see that  $\mathbb{P}(y = 1 | \mathbf{x})$  can be parameterized into the form used by logistic regression under Gaussian Naive Bayes assumption. The above expression gives an alternate way of estimating the value of the weights  $w_i$  by estimating the probabilities directly. With the above insight, give advantages and disadvantages of using generative and discriminative estimation of the parameters in logistic regression.

$$\begin{aligned} \mathcal{N}(\mu_{ik}, \sigma_i^2) &= \frac{1}{\sqrt{2\pi} \sigma_i} e^{-\frac{(x - \mu_{ik})^2}{2\sigma_i^2}} \\ p[x_i | y=0] &= \frac{1}{\sqrt{2\pi} \sigma_i} e^{-\frac{(x - \mu_{i0})^2}{2\sigma_i^2}} \\ p[x_i | y=1] &= \frac{1}{\sqrt{2\pi} \sigma_i} e^{-\frac{(x - \mu_{i1})^2}{2\sigma_i^2}} \\ \log \left[ \frac{p(x_i | y=0)}{p(x_i | y=1)} \right] &= \log \left( e^{\frac{(x - \mu_{i1})^2}{2\sigma_i^2} - \frac{(x - \mu_{i0})^2}{2\sigma_i^2}} \right) \\ &= \frac{(x - \mu_{i1})^2 - (x - \mu_{i0})^2}{2\sigma_i^2} = \frac{x^2 + \mu_{i1}^2 - 2x\mu_{i1} - (x^2 + \mu_{i0}^2 - 2x\mu_{i0})}{2\sigma_i^2} \\ &= \frac{\mu_{i1}^2 - \mu_{i0}^2}{2\sigma_i^2} + \frac{x(\mu_{i0} - \mu_{i1})}{\sigma_i^2} \end{aligned}$$

Bayes rule

$$Q2b) \quad P(y=1 | x) = \frac{P(y=1) P(x|y=1)}{P(y=0) P(x|y=0) + P(y=1) P(x|y=1)}$$

$$= \frac{1}{1 + \frac{P(y=0) \prod_{i=1}^d P(x_i|y=0)}{P(y=1) \prod_{i=1}^d P(x_i|y=1)}}$$

conditional independence

$$= \frac{1}{1 + \exp \left( \frac{\log \left( \frac{P(y=0)}{P(y=1)} \right) + \sum_{i=1}^d \log \left( \frac{P(x_i|y=0)}{P(x_i|y=1)} \right)}{1} \right)}$$

$$= \frac{1}{1 + \exp \left( \log \left( \frac{P(y=0)}{P(y=1)} \right) + \sum_{i=1}^d \log \left( \frac{P(x_i|y=0)}{P(x_i|y=1)} \right) \right)}$$

$$= \frac{1}{1 + \exp \left( \log \left( \frac{1-\sigma}{\sigma} \right) + \sum_{i=1}^d \log \left( \frac{P(x_i|y=0)}{P(x_i|y=1)} \right) \right)}$$

Q2c). Discriminative Modelling is better for classification with a lot of datapoints. Since find  $w_i$  directly is a much easier and less time consuming process than modelling it indirectly using probabilities.

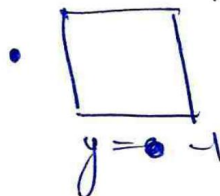
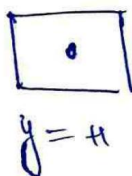
⇒ With large numbers of datapoints, discriminative models are preferred (low bias, high variance) → general  $w_i x_i + w_0$  instead of a specific form as in Naive Bayes.

With lower numbers of points they ~~are~~ tend to overfit and Generative modelling works much better. Because of priors are inserted into the model (high bias, low variance). We also get  $P(x|y)$  for generating data which is not possible through discriminative modelling.



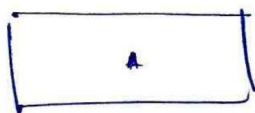
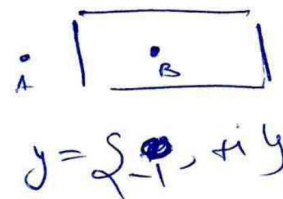
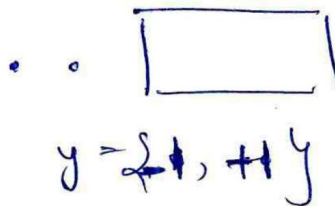
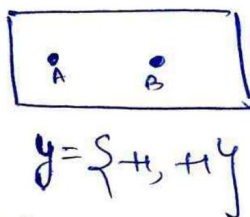
3. [10 marks] Let  $\mathcal{H}$  consists of all hypotheses in two dimensions  $h: \mathbb{R}^2 \rightarrow \{+1, -1\}$  that are positive inside some convex set and negative elsewhere (a set is convex if the line segment connecting any two points in the set lies entirely within the set). Compute the VC dimension of  $\mathcal{H}$ .

for  $N=1$



shattered

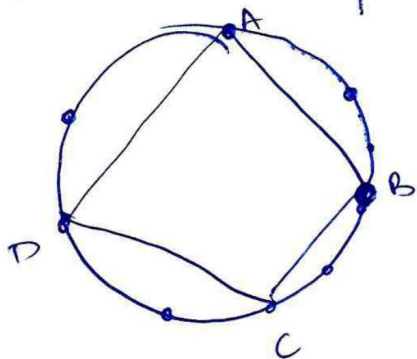
for  $N=2$



$y = \{-, +\}$

shattered

for  $N \geq 3$  consider  $N$  points which lie on the diameter of a circle.



~~we~~

We show that we can create all dichotomies using  $\mathcal{H}$ .

in the dichotomy where all points are -ve consider the empty convex set.

In the dichotomy where ~~one~~ <sup>one</sup> point is

+ve consider the convex set to be the point itself

~~if~~ When 2 points are +ve consider the convex set to be the line joining the two points.

for  $\geq 3$  points positive consider the convex polygon joining those points, all other points will be -ve.

This shows that for any  $N$  all dichotomies can be generated  $\Rightarrow d_{VC} = N$

4. [15 marks] In the class we discussed Perceptron Learning Algorithm (PLA) as follows:

---

**Algorithm 1** Perceptron Learning Algorithm

---

$w(1) \leftarrow 0, t \leftarrow 1$

**while** any misclassified example left **do**

    Denote the current weight vector as  $w(t)$

    Pick any misclassified sample  $(x_*, y_*) : \text{sign}(w(t)^\top x_*) \neq y_*$

$w(t+1) \leftarrow w(t) + y_* x_*$

$t \leftarrow (t+1)$

**end while**

---

Prove that, if the data is linearly separable, then PLA will be able to find one such separator.

Since the data is linearly separable consider that some optimal  $w^*$  exists which correctly classifies all points.

i.e.  $\text{sign}(w^* x) = y$  for all examples.

Consider the angle  $\theta(t)$  between  $w(t)$  and  $w^* \rightarrow \theta(t)$

$$\cos(\theta(t)) = \frac{(w^*)^\top w(t)}{\|w^*\| \|w(t)\|}$$

Consider the numerator at iteration  $t+1$

$$\begin{aligned} (w^*)^\top [w(t+1)] &= (w^*)^\top [w(t) + y_* x_*] \\ &= (w^*)^\top w(t) + \underbrace{y_* (w^*)^\top x_*}_{>0} \end{aligned}$$

$\Rightarrow$  Numerator ~~increases~~ increases with every iteration.

Since  $\text{sign}(w^* x_*) = y_*$



Consider the denominator,

$$\|w^*\| \|w(t+1)\| = \|w^*\| \left( \|w(t) + y^* x^*\| \right)$$

The angle b/w  $w(t)$  &  $y^* x^*$  is obtuse since  $y^* x^*$  is misclassified by  $w(t)$

$$\Rightarrow \|w(t) + y^* x^*\| < \|w(t)\|$$

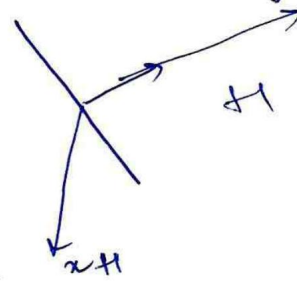
~~At every iteration~~ The denominator decreases with every iteration.

$\Rightarrow \cos(\theta(t))$  increases with every iteration.

Since  $\cos(\theta(t))$  is bounded by 1 it cannot keep increasing in every iteration.  $\Rightarrow$  The Algorithm must terminate.

$\Rightarrow$  All examples are correctly classified by  $w(t)$  at that point.

$\Rightarrow$  The Algorithm will find a classifier which separates the classes.



---

$$N(u_i, \sigma_i^2) =$$

$$A \cdot e^{-\frac{(x-u_i)^2}{2\sigma_i^2}}$$

$$\frac{1}{\sqrt{2\pi\sigma^2}}$$

$$A e^{-\frac{(x-u_i)^2}{2\sigma^2}}$$

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

(ROUGH WORK)

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

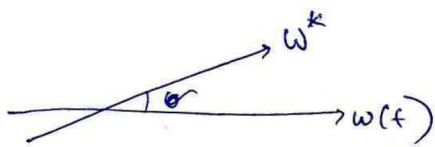
$$-\log(P(y=1|x)) = \log\left(1 + \frac{P(x|y=0)P(y=0)}{P(x|y=1)P(y=1)}\right)$$

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

$$P(w)P(x|w)$$

$$-\log P(w)$$





$$\cos \theta = \frac{w^{*T} (w(t) + y_* x_*)}{|w^*| |w(t)|}$$

$$w^{*T} w(t) + \underbrace{w^{*T} y_* x_*}_{\text{tr.}}$$

$$|w(t+1)|$$

$$2(y - w x)(-x)$$

