

COL 341 Minor 2

Chinmay Mittal

TOTAL POINTS

58 / 65

QUESTION 1

1 SVM 13 / 20

+ 0 pts Not Attempted/ Incorrect

✓ + 8 pts *Correct Primal and Lagrange Formation*

- 2 pts Condition of alpha and beta > 0 missing in conditions of Lagrange

+ 6 pts Correct differentiation equations from Lagrange

+ 3 pts Correctly substituting and deriving the final form of Lagrange

✓ + 3 pts *Correctly showing constraints on dual*

+ 0 pts Minor mistakes leading to point adjustment

+ 4 pts Partially correct Primal and Lagrange

+ 4 pts Partial : 2 out of 3 derivatives are correct

+ 2 *Point adjustment*

1 This needs to be derived by substituting with derivatives

QUESTION 2

2 Euclidean Distance 10 / 10

✓ + 10 pts *Correct*

+ 8 pts Minor bugs (may be minor incompleteness)

+ 5 pts Major bugs or incompleteness with partial correctness

+ 0 pts Incorrect/ Not Attempted

QUESTION 3

Decision Tree 10 pts

3.1 (a) 5 / 5

+ 3 pts Partial Correct

✓ + 5 pts *Correct*

+ 0 pts Incorrect/Not Attempted

- 1 pts No Explanation/Justification

3.2 (b) 5 / 5

✓ + 5 pts *Correct*

+ 0 pts Incorrect/Not Attempted

QUESTION 4

AdaBoost 25 pts

4.1 (a) 2 / 2

4a

✓ - 0 pts *Correct*

- 2 pts Incorrect

- 2 pts Unattempted

4.2 (b) 3 / 3

4b

✓ - 0 pts *Correct*

- 3 pts Incorrect/Unattempted.

- 1 pts Wrong Exponent

4.3 (c) 10 / 10

✓ + 10 pts [Click here to replace this description.](#)

+ 0 pts Incorrect/ not attempted

- 2 pts steps are missing

- 1 pts did not explain the normalization term

- 1 pts did not explain the $1/m$ term

+ 3 pts initial steps written

4.4 (d) 10 / 10

+ 0 pts Incorrect

+ 0 pts Not attempted

✓ + 10 pts Correct

Student Name: CHINMAY MITAL

Entry Number: 2020C810336



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

Minor 2

Time: 60 minutes

Maximum Marks: 65

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	2	3(a)	3(b)	4(a)	4(b)	4(c)	4(d)	Total
Max Marks	20	10	5	5	2	3	10	10	65
Earned Marks									

1. [20 marks] Recall the soft-margin SVM covered in the class, where the primal can be written as follows:

$$\begin{aligned} \min_{b, w, \xi} \quad & \frac{1}{2} w^\top w + C \sum_{n=1}^m \xi_n \\ \text{subject to: } & y_n(w^\top x_n + b) \geq 1 - \xi_n \\ & \xi_n \geq 0, \text{ for } n = 1, \dots, m. \end{aligned}$$

We derived the corresponding dual in one of the homework assignments as:

$$\begin{aligned} \min_{\alpha} \quad & \frac{1}{2} \alpha^\top G \alpha - \sum_{n=1}^m \alpha_n \\ \text{subject to: } & y^\top \alpha = 0 \\ & 0 \leq \alpha_n \leq C, \text{ for } n = 1, \dots, m, \end{aligned}$$

for an appropriate matrix G . One can look at the above as including the penalty term $C \|\xi\|_p$, with $p = 1$ in the primal, which led to the constraint $\|\alpha\|_q \leq C$, with $q = \infty$ in the dual. One can prove in general that $\frac{1}{p} + \frac{1}{q} = 1$. Your task in this question is to derive this for $p = \infty$, and $q = 1$.

The primal optimization problem becomes $\min_{b, w, \xi} \frac{1}{2} w^\top w + C \|\xi\|_\infty$

$$\text{s.t. } y_n(w^\top x_n + b) \geq 1 - \xi_n$$

$$\xi_n \geq 0$$

$$\|\xi\|_\infty = \max_i \xi_i$$

We can form an equivalent optimization problem as follows.

$$\min_{b, w, \xi} \frac{1}{2} w^T w + C \sum \xi^*$$

subject to $\xi^* \geq \xi_n \quad \forall n$

$$y_n (w^T x_n + b) \geq 1 - \xi_n \quad \forall n$$

$$\xi_n \geq 0 \quad \forall n$$

This encodes the constraint
 $\max_n \xi_n = \xi^*$

ξ, β_n, c_n
 for 3 types of constraints

The Lagrangian of this equivalent primal optimization is as follows

$$L(b, w, \xi, \gamma_n, \beta_n, c_n) = \frac{1}{2} w^T w + C \xi^* + \sum \gamma_n (1 - \xi_n - y_n (w^T x_n + b)) + \sum \beta_n (-\xi_n) + \sum c_n (\xi_n - \xi^*)$$

$$\text{s.t. } \gamma_n \geq 0, \beta_n \geq 0, c_n \geq 0$$

from KKT conditions

$$\frac{\partial L}{\partial \xi_n} = 0 \Rightarrow -\gamma_n - \beta_n + c_n = 0 \quad \forall n$$

$$\gamma_n + \beta_n = c_n$$

$$\gamma_n \geq 0, \beta_n \geq 0, c_n \geq 0$$

$$\frac{\partial L}{\partial \xi^*} = 0 \Rightarrow C - \sum c_n = 0 \Rightarrow \sum c_n = C$$

Substituting this in the Lagrangian we get (β_n and c_n in terms of γ_n)

$$L(b, w, \xi, \gamma_n) = \frac{1}{2} w^T w - \sum \gamma_n y_n (w^T x_n + b) + \xi^* (C - \sum c_n) + \sum \gamma_n (c_n - \gamma_n - \beta_n)$$

Lagrangian is same as

hard margin SVM subject

to different constraints.

The form of the optimization problem comes out

exactly the same $\min_{\alpha} \frac{1}{2} \alpha^T G \alpha - \sum_{n=1}^m \gamma_n$ subject to $y^T \alpha = 0$

The new constraints are

$$C = \sum c_n \Rightarrow C = \sum \gamma_n + \sum \beta_n$$

$$\Rightarrow \|\alpha\|_1 \leq C$$

$$\Rightarrow \text{When } q=1 \quad p=1$$

Since $\beta_n \geq 0$

$$\text{we get } 0 \leq \sum \gamma_n \leq C \quad \gamma_n \geq 0$$

2. [10 marks] Kernels are general idea which can be used to work with infinite dimensional feature space in variety of ML techniques. We saw it extensively in SVMs, and also mentioned that it could be used for k-NN classifier as well. However, to make it work in k-NN we need to compute the distances in the D -dimensional feature space, where D can be ∞ as well. Suppose we have a kernel $K(\cdot, \cdot)$, such that there is an implicit high dimensional feature map $\phi: \mathbb{R}^d \rightarrow \mathbb{R}^D$ that satisfies $\forall x, z \in \mathbb{R}^d, K(x, z) = \phi(x) \cdot \phi(z)$, where $\phi(x) \cdot \phi(z) = \sum_{i=1}^D \phi(x)_i \phi(z)_i$ is the dot product in the D -dimensional space. Show, how to calculate the Euclidean distance in the D -dimensional space:

$$\|\phi(x) - \phi(z)\| = \sqrt{\sum_{i=1}^D (\phi(x)_i - \phi(z)_i)^2},$$

without explicitly calculating the values in the D -dimensional vectors.

$$\begin{aligned} \|\phi(x) - \phi(z)\| &= \sqrt{\sum_{i=1}^D (\phi(x)_i^2 + \phi(z)_i^2 - 2\phi(x)_i \phi(z)_i)} \\ &= \sqrt{\sum_{i=1}^D \phi(x)_i \cdot \phi(x)_i + \sum_{i=1}^D \phi(z)_i \cdot \phi(z)_i - 2 \sum_{i=1}^D \phi(x)_i \phi(z)_i} \\ &= \sqrt{K(x, x) + K(z, z) - 2K(x, z)} \end{aligned}$$

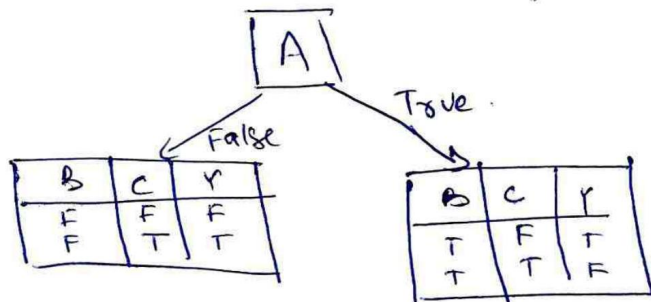
MITTAL

3. Consider the following dataset with $A, B, C \in \{T, F\}$ as the input features, and $Y \in \{T, F\}$ as the output label:

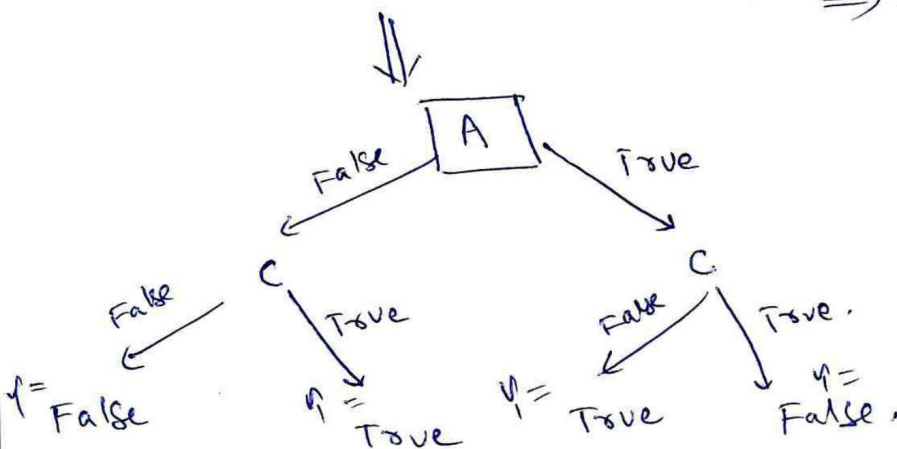
A	B	C	Y
F	F	F	F
F	F	T	T
T	T	F	T
T	T	T	F

- (a) [5 marks] Using the dataset above, we want to build a decision tree which classifies Y as T/F . Draw the tree that would be learned by the greedy algorithm with zero training error. You do not need to show any computation but should describe the justification of each decision in English.
- (b) [5 marks] Is this tree optimal (i.e. does it get zero training error with minimal depth)? Explain in less than two sentences. If it is not optimal, draw the optimal tree as well.

At the first node, the choices are. $A^- [F+]$ $B^- [FT]$ $C^- [FT]$
 $A^+ [TF]$ $B^+ [TF]$ $C^+ [TF]$
 All nodes have same information gain \Rightarrow choose any node, I choose A



For both the nodes
 all data items have same
 value of B
 \Rightarrow split has to be w.r.t C
 at each node



Attributes A and B are symmetric (all data points have same value)

\Rightarrow choosing B first will not give a better tree

\checkmark if we choose C first, then at the second stage we can choose any amongst attribute A/B to perfectly classify all nodes \Rightarrow depth of tree = 2

Hence the tree we created is optimal (no better tree exists)
 \hookrightarrow Hence amongst all possibilities no tree is better than ours \Rightarrow our tree is optimal.

$$f(x) = \sum_{t=1}^T \alpha_t h_t(x)$$

COL341: Minor 2

Student Name: CHINMAY

Entry Number: 2020CS10336

MITTAL

4. We discussed AdaBoost algorithm in the class to learn a classifier $f(x) = \text{sign} \left(\sum_{t=1}^T \alpha_t h_t(x) \right)$.

In each time step t we compute the error, $\text{err}_t = \frac{\sum_{i=1}^m w_t^i \mathbb{I}(h_t(x^i) \neq y^i)}{\sum_{i=1}^m w_t^i}$, where m denotes the number of samples. We use w to estimate $\alpha_t = \frac{1}{2} \log \left(\frac{1 - \text{err}_t}{\text{err}_t} \right)$. We then reweigh the samples according to the expression: $w_{t+1}^i = w_t^i \exp(-\alpha_t y^i h_t(x^i))$.

(a) [2 marks] Show that $y^i h_t(x^i) = \mathbb{I}(h_t(x^i) \neq y^i)$

(b) [3 marks] Show that $\exp(-\alpha_t y^i h_t(x^i)) \propto \exp(2\alpha_t \mathbb{I}(h_t(x^i) \neq y^i))$

(c) [10 marks] Assuming $Z_t = \sum_{i=1}^m w_t^i$, and we initialize weight of each sample as uniform in $t = 1$, show that $w_{t+1}^i = \frac{\exp(-y^i \sum_{s=1}^t \alpha_s h_s(x^i))}{\prod_{s=1}^t Z_s}$.

(d) [10 marks] Assuming that the AdaBoost minimizes the loss $\mathcal{L}(f) = \frac{1}{m} \sum_{i=1}^m \exp(-y^i f(x^i))$, show that the loss at time step t , $\mathcal{L}(f_t) = \prod_{s=1}^t Z_s$.

(a) $\frac{1}{2} (1 - h_t(x^i) y^i)$ When the two are not equal $h_t(x^i) \neq y^i$
 $\Rightarrow h_t(x^i) = 1, y^i = -1$ or $h_t(x^i) = -1, y^i = 1 \Rightarrow h_t(x^i) y^i = -1$
 $\Rightarrow \frac{1}{2} (1 - h_t(x^i) y^i) = \frac{1}{2} (1 - (-1)) = 1 = \mathbb{I}(h_t(x^i) \neq y^i)$

When the two are equal $h_t(x^i) = y^i = 1$ or $h_t(x^i) = y^i = -1$
 $\Rightarrow \frac{1}{2} (1 - h_t(x^i) y^i) = \frac{1}{2} (1 - 1) = 0 = \mathbb{I}(h_t(x^i) \neq y^i)$

(b) $\exp \left(\frac{1}{2} (1 - h_t(x^i) y^i) \right) = \exp(\mathbb{I}(h_t(x^i) \neq y^i))$
 $\Rightarrow \exp(1 - h_t(x^i) y^i) = \exp(2 \mathbb{I}(h_t(x^i) \neq y^i))$
 $\Rightarrow \exp(-h_t(x^i) y^i) = \exp(-1) \cdot \exp(2 \mathbb{I}(h_t(x^i) \neq y^i))$
 $\Rightarrow \exp(-\alpha_t h_t(x^i) y^i) = \exp(-\alpha_t) \cdot \exp(2 \alpha_t \mathbb{I}(h_t(x^i) \neq y^i))$
 $\propto \exp(2 \alpha_t \mathbb{I}(h_t(x^i) \neq y^i))$

(c) $w_2^i = \frac{w_1^i}{\alpha} \exp(-\alpha_1 y^i h_1(x^i)) / Z_1$ initial weight: $1/m$ In my proof I explicitly normalize weights at each step.

$w_3^i = \frac{w_2^i}{\alpha} \exp(-\alpha_2 y^i h_2(x^i)) / Z_2$ normalization

\vdots

$w_{t+1}^i = \frac{w_t^i}{\alpha} \exp(-\alpha_t y^i h_t(x^i)) / Z_t$

Multiplying all equations

$$w_{t+1}^i = \frac{1}{m} \exp\left(-y^i \left(\sum_{s=1}^T \alpha_s h_s(x^i)\right)\right) \frac{1}{\prod_{s=1}^T Z_s}$$

We had normalized in the computation

(d) $w_{t+1}^i = \frac{1}{m} \exp\left(-y^i f_t(x^i)\right) \frac{1}{\prod_{s=1}^T Z_s}$ since $f_t(x^i) = \sum_{s=1}^T \alpha_s h_s(x^i)$

$$\sum_{i=1}^m w_{t+1}^i = \frac{1}{m} \frac{1}{\prod_{s=1}^T Z_s} \sum_{i=1}^m \exp(-y^i f_t(x^i))$$

$= 1$

$$\Rightarrow 1 = \frac{1}{\prod_{s=1}^T Z_s} \frac{1}{m} \sum_{i=1}^m \exp(-y^i f_t(x^i))$$

Since weights are normalized

$$\Rightarrow \prod_{s=1}^T Z_s = L(f_t)$$

ROUAK

COL341: Minor 2

Student Name: ~~2020010336~~

Entry Number: 2020010336

CHINMAY
MITTAL

$$\frac{1 \cdot 0}{1 - 1} = 0$$

$$\sqrt{\sum_{i=1}^D (\phi(x)_i - \phi(z)_i)^2}$$

$$= \sqrt{\sum_{i=1}^D (\phi(x)_i^2 + \phi(z)_i^2 - 2\phi(x)_i\phi(z)_i)}$$

$$= \sqrt{\sum_{i=1}^D \phi(x)_i^2 + \phi(z)_i^2 - 2\phi(x)_i\phi(z)_i}$$

$$= \sqrt{K(x, x) + K(z, z) - 2K(x, z)}$$

$$e^{\frac{1}{2}(1 - h(x)(y))}$$

$$e^{\frac{1}{2}(1 - h(x)(y))} = e^{\frac{1}{2} - \frac{1}{2}h(x)(y)}$$

$$\frac{1}{2} \exp\left(\frac{1}{2}h(x)(y)\right)$$

$$\frac{1}{2} \exp\left(\frac{1}{2}h(x)(y)\right) = \frac{1}{2} \exp\left(\frac{1}{2}h(x)(y)\right)$$

$$\frac{1}{2} \exp\left(\frac{1}{2}h(x)(y)\right) = \frac{1}{2} \exp\left(\frac{1}{2}h(x)(y)\right)$$

$$(d) \quad w_{t+n}^i = \frac{\exp(-j^i t_t(x^i))}{\sum_{s=1}^m \pi^s 2s}$$

$$\sum w_{t+n}^i = \frac{L(t_t)}{\pi^t 2s}$$

$$\min_{b, w, \epsilon} \quad \frac{1}{2} w^T w + c \sum \xi_i$$

$$\xi_i \geq \epsilon_i \geq 0$$

$$L \rightarrow \frac{1}{2} w^T w + c \sum \xi_i + \sum \alpha_n (1 - \xi_n - y_n (w^T x_n + b)) \\ + \sum \beta_n (-\xi_n) + \sum c_n (\xi_n - \xi)$$

$$c - \alpha_n - \beta_n = 0$$

$$c = \sum c_n$$

$$-\alpha_n - \beta_n + c_n = 0$$

$$L \rightarrow \frac{1}{2} w^T w - \sum \alpha_n y_n (w^T x_n + b) \quad \boxed{c_n = \alpha_n + \beta_n}$$

$$\alpha_n \geq 0$$

$$\beta_n \geq 0$$

$$c_n \geq 0$$

$$\sum \alpha_n + \sum \beta_n = c$$