

CS 771A: Intro to Machine Learning, IIT Kanpur			Midsem Exam (16 Feb 2020)	
Name	SAMPLE SOLUTION			80 marks
Roll No		Dept.		Page 1 of 6

Instructions:

1. This question paper contains 3 pages (6 sides of paper). Please verify.
2. Write your name, roll number, department in **block letters neatly** with ink **on each page** of this question paper.
3. If you don't write your name and roll number on **all** pages, **pages may get lost** when we unstaple to scan pages
4. Write your final answers neatly **with a blue/black pen**. Pencil marks may get smudged.
5. Don't overwrite/scratch answers especially in MCQ and T/F. We will entertain no requests for leniency.

Q1. Write T or F for True/False (write only in the box on the right hand side) (3x2 = 6 marks)

1	The variance of a real-valued (continuous or discrete) random variable must always be a strictly positive number i.e. it can never be zero or negative.	F
2	Suppose X, Y are two (not necessarily independent) r.v. with support \mathbb{R} . We are told that $\text{Var}[X] \neq 0$ and $\text{Var}[Y] \neq 0$. Then it must be the case that $\text{Var}[X + Y] \neq 0$.	F
3	Consider a doubly differentiable function $f: \mathbb{R} \rightarrow \mathbb{R}$ that always takes values in the interval $[0,1]$ i.e. $f(x) \in [0,1]$ for all $x \in \mathbb{R}$. Such a function can never be convex.	F

Q2. Fill the circle (don't tick) next to all the correct options (many may be correct). (2x3=6 marks)

2.1 Let $X \sim \mathcal{N}(0, \sigma^2)$ be a Gaussian r.v.. Let Y be another r.v. defined as $Y = X$ if $X \in [-1,1]$ else $Y = 2X$ if $X \notin [-1,1]$. Suppose we keep increasing σ . Then as $\sigma \rightarrow \infty$, which of the following is true?

A	The variance $\text{Var}[X]$ goes up as $\sigma \rightarrow \infty$	<input checked="" type="radio"/>
B	The variance $\text{Var}[X]$ goes down as $\sigma \rightarrow \infty$	<input type="radio"/>
C	The PDF $\mathbb{P}[Y X \in [-1,1]]$ starts looking more and more like $\text{UNIF}[-1,1]$ as $\sigma \rightarrow \infty$	<input checked="" type="radio"/>
D	The probability $\mathbb{P}[Y \in [-1,1]]$ approaches 1 as $\sigma \rightarrow \infty$	<input type="radio"/>

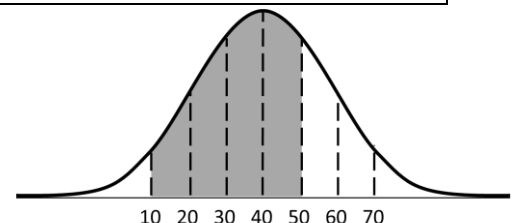
2.2 Suppose X, Y are two independent r.v. Which of the following is always true?

A	$\mathbb{E}[X + Y] = \mathbb{E}X + \mathbb{E}Y$	<input checked="" type="radio"/>
B	$\text{Var}[X + Y] = \text{Var}[X] + \text{Var}[Y]$	<input checked="" type="radio"/>
C	$\mathbb{E}[X Y] = \mathbb{E}X$	<input checked="" type="radio"/>
D	$\text{Var}[X Y] = \text{Var}[X]$	<input checked="" type="radio"/>

Q3 $X \sim \mathcal{N}(\mu, \sigma^2)$ with $\mu = 40$ and $\sigma = 10$.

Approximate $\mathbb{P}[10 \leq X \leq 50]$ using the 68-95-99.7 rule that tells us that the probability that X is within c standard deviations of its mean is approximately 0.68, 0.95 and 0.997 for $c = 1, 2, 3$ respectively. Explain your answer briefly.

(4 marks)



Denote $p_c \triangleq \mathbb{P}[|X - \mu| \leq c \cdot \sigma]$. Then the desired probability is given by $\frac{(p_3 + p_1)}{2} \approx 0.8385$. Due to the approximate nature of the problem, other acceptable answers include 0.84, 0.839 and 0.838.

Q4 X is a continuous r.v. with support $[0,1]$ and uniform PDF over its support. Define a new discrete r.v. Y with support over $\{0,1,2\}$. Let $0 \leq a \leq b \leq 1$ be two (unknown) constants. Y is defined using

X itself as follows: define $Y = \begin{cases} 0 & X \leq a \\ 1 & a < X < b \\ 2 & X \geq b \end{cases}$ (please pay attention to $<$ vs \leq). First, write down

the PMF for Y . Then write down an expression for $\text{Var}[Y]$. Your PMF and variance expressions must be in terms of the variables a, b and not for specific values of a, b . Then, find a pair of values for a, b for which the r.v. Y has its highest possible variance (give a brief derivation as well). For your chosen value of a, b , write down the value of $\mathbb{E}[X + Y]$ and $\text{Var}[Y]$. **(4+4+4+2+2=16 marks)**

Give PMF for Y here

$$\text{We have } \mathbb{P}[Y = c] = \begin{cases} a & c = 0 \\ b - a & c = 1 \\ 1 - b & c = 2 \end{cases}$$

Give expression for $\text{Var}[Y]$ here

We have $\mathbb{E}Y = 2 - (a + b)$ and $\mathbb{E}[Y^2] = 4 - (a + 3b)$ (giving these expressions is not required to get marks for this part) and thus $\text{Var}[Y] = 3a + b - (a + b)^2$

$$\text{Best value of } a = \left(\frac{1}{2} \right) \quad b = \left(\frac{1}{2} \right)$$

Give derivation for best values of a, b here

Doing the calculations directly would be perfectly fine but cumbersome. To simplify, we notice that variance is a translation invariant statistic i.e. $\text{Var}[Y] = \text{Var}[Y + r]$ for any fixed $r \in \mathbb{R}$. Taking $r = -1$ is advantageous for us since it makes the random variable Y take symmetric values about 0 and gives us $\mathbb{E}[Y - 1] = 1 - (a + b)$ and $\mathbb{E}[(Y - 1)^2] = 1 - (b - a)$. Thus, we have $\text{Var}[Y] = \text{Var}[Y - 1] = \underbrace{1 - (b - a)}_{(A)} - \underbrace{(1 - (a + b))^2}_{(B)}$. Since $b \geq a$, the largest value of (A) is 1 which is achieved when $a = b$. The smallest value of (B) is 0 since it is a squared term and this is achieved when $a + b = 1$. Satisfying both together gives $a = b = \frac{1}{2}$.

$$\mathbb{E}[X + Y] = \left(1.5 \right) \quad \text{Var}[Y] = \left(1 \right)$$

CS 771A: Intro to Machine Learning, IIT Kanpur			Midsem Exam (16 Feb 2020)	
Name	SAMPLE SOLUTION			80 marks
Roll No		Dept.		Page 3 of 6

Q5 Let $a^1, \dots, a^n \in \mathbb{R}$ and $\lambda > 0$ be constants that are given to us. **(4+4+10+6 = 24 marks)**

5.1 Solve the L_2 regularized problem (give brief derivation) $\min_x \frac{\lambda}{2} \cdot x^2 + \frac{1}{2n} \cdot \sum_{i=1}^n (x - a^i)^2$

5.2 We will now solve $\min_x \lambda \cdot |x| + \frac{1}{2n} \cdot \sum_{i=1}^n (x - a^i)^2$ the L_1 regularized problem. We rewrite the optimization problem as shown on the right-hand side. Write down the Lagrangian of this problem by introducing dual variables for the constraints.

$$\begin{aligned} \min_{x,c} \quad & \lambda \cdot c + \frac{1}{2n} \cdot \sum_{i=1}^n (x - a^i)^2 \\ \text{s. t.} \quad & x \leq c \\ & x \geq -c \\ & c \geq 0 \end{aligned}$$

5.3 Using the Lagrangian, create and simplify the dual problem (show brief derivation). If simplified properly, the dual problem should involve only two real valued variables. Try to simplify the dual problem as much as you can (otherwise the next part may be more difficult for you).

5.4 Propose a coordinate ascent/maximization (or descent/minimization if you are writing your dual as a minimization problem) method to solve your simplified dual. Use cyclic coordinate selection and random initialization for simplicity. Give precise expressions in your pseudocode (not vague statements) on how you would process a chosen coordinate taking care of constraints etc.

Give solution to 5.1 here

Applying first order optimality (since this is an unconstrained problem with a differentiable objective) gives us $\lambda \cdot x + \frac{1}{n} \sum_{i=1}^n (x - a^i) = 0$ which is uniquely solved at $x = \bar{a}/(\lambda + 1)$ where $\bar{a} = \frac{1}{n} \sum_{i=1}^n a^i$.

Give solution to 5.2 here

We introduce three new dual variables α, β, γ corresponding to the three constraints to get

$$\mathcal{L}(x, c, \alpha, \beta, \gamma) = \lambda \cdot c + \frac{1}{2n} \cdot \sum_{i=1}^n (x - a^i)^2 + \alpha(x - c) - \beta(x + c) - \gamma c$$

Give solution to 5.3 here

The (unsimplified) dual problem is $\max_{\alpha \geq 0, \beta \geq 0, \gamma \geq 0} \min_{x, c} \mathcal{L}(x, c, \alpha, \beta, \gamma)$. Since the inner optimization problem is unconstrained and differentiable, applying first order optimality w.r.t c gives us

$$\frac{\partial \mathcal{L}}{\partial c} = 0 \Rightarrow \gamma = \lambda - (\alpha + \beta)$$

Since $\gamma \geq 0$, the above tells us that $\alpha + \beta \leq \lambda$. Similarly applying first order optimality w.r.t x

$$\frac{\partial \mathcal{L}}{\partial x} = 0 \Rightarrow x - \bar{a} + \alpha - \beta = 0$$

which tells us that $x = \bar{a} + (\beta - \alpha)$. Using these to simplify the dual gives us

$$\begin{aligned} \min_{\alpha, \beta} \quad & \frac{1}{2} \cdot (\bar{a} + (\beta - \alpha))^2 \\ \text{s.t.} \quad & \alpha + \beta \leq \lambda \\ & \alpha \geq 0 \\ & \beta \geq 0 \end{aligned}$$

We note that we neglected a term $\frac{1}{2n} \cdot \sum_{i=1}^n (a^i)^2$ in the objective above since it is a constant and does not affect the final solution at all.

Give solution to 5.4 here

Suppose we fix α and want to update β to its optimal value, then had it not been for the constraints, the optimal value would have been $\alpha - \bar{a}$. However, just as in the case of the CSVM dual, the constraints can be taken care of in this case by simply projecting this value onto the interval $[0, \lambda - \alpha]$. Similarly, if we fix β and want to update α to its optimal value, that value is simply the unconstrained optimum value $\beta + \bar{a}$ projected to the interval $[0, \lambda - \beta]$. The final coordinate descent algo is given below

1. Initialize $\alpha = \text{UNIF}[0, \lambda], \beta = \lambda - \alpha$ (make sure that $\alpha + \beta \leq \lambda$ is satisfied initially)
2. For $t = 1, 2, \dots$
 - a. If t is odd: # update α
 - i. $\alpha \leftarrow \beta + \bar{a}$
 - ii. If $\alpha < 0$, set $\alpha = 0$
 - iii. If $\alpha > \lambda - \beta$, set $\alpha = \lambda - \beta$
 - b. If t is even: # update β
 - i. $\beta \leftarrow \alpha - \bar{a}$
 - ii. If $\beta < 0$, set $\beta = 0$
 - iii. If $\beta > \lambda - \alpha$, set $\beta = \lambda - \alpha$

Coordinate descent algorithms can sometimes get stuck as well, especially if the optimization problem is not very well behaved. However, in general, they offer some of the fastest solvers, especially when applied to dual problems.

CS 771A: Intro to Machine Learning, IIT Kanpur			Midsem Exam (16 Feb 2020)	
Name	SAMPLE SOLUTION			80 marks
Roll No		Dept.		Page 5 of 6

Q6 Let $M \gg 1$ be some large constant. For $a \in [0.5, M - 0.5]$, define an *almost-uniform* (AU) distribution $\tilde{\mathcal{U}}$ with PDF described below. We are given data $x^1, \dots, x^n \in [0, M]$. We wish to fit a mixture of two AU distributions $\tilde{\mathcal{U}}(x; a^1)$ and $\tilde{\mathcal{U}}(x; a^2)$ for $a^1, a^2 \in [0.5, M - 0.5]$ to the data.

$$\mathbb{P}[x | a] = \tilde{\mathcal{U}}(x; a) \triangleq \begin{cases} 0 & x \notin [0, M] \\ \frac{1}{2M} & x \in \left[0, a - \frac{1}{2}\right) \\ \frac{1}{2M} + \frac{1}{2} & x \in \left[a - \frac{1}{2}, a + \frac{1}{2}\right] \\ \frac{1}{2M} & x \in \left(a + \frac{1}{2}, M\right] \end{cases}$$

Let $z^i \in \{1, 2\}$ be latent variables denoting which distribution generated which data point. Use the data likelihood function $\mathbb{P}[x^i | z^i, a^1, a^2] = \tilde{\mathcal{U}}(x^i; a^{(z^i)})$. Expressions in your answers may contain unspecified normalization constants. Give only brief derivations. **(8+10+6=24 marks)**

6.1 Assuming $\mathbb{P}[z^i = c | a^1, a^2] = \mathbb{P}[z^i = c] = 0.5$ for $c \in \{1, 2\}$ (i.e. uniform prior on z^i), derive an expression for $\mathbb{P}[z^i = 1 | x^i, a^1, a^2]$. Using this, show how to calculate the MAP estimate for the latent variables i.e. $\operatorname{argmax}_{c \in \{1, 2\}} \mathbb{P}[z^i = c | x^i, a^1, a^2]$. Break any ties in any way you like.

6.2 Show how to calculate the MAP estimate for the model i.e. $\operatorname{argmax}_{a^1, a^2 \in [0.5, M-0.5]} \mathbb{P}[a^1, a^2 | \mathbf{X}, \mathbf{z}]$ (note

the shorthand $\mathbf{X} = [x^1, \dots, x^n] \in [0, M]^n$, $\mathbf{z} = [z^1, \dots, z^n] \in \{1, 2\}^n$). You are allowed to use `k-nn` and `r-nn` as library function calls in your solution. For any $x \in [0, M]$, $k \in \mathbb{N}$, $r \in \mathbb{R}_+$, `k-nn`(x, k) returns the k nearest neighbors of x in the dataset \mathbf{X} whereas `r-nn`(x, r) returns all points in the dataset \mathbf{X} which are a distance r or less from x . Both `k-nn` and `r-nn` use the Euclidean distance.

6.3 Using the above, give pseudocode (as we do in lecture slides with sufficient algo details and not necessarily Python code) for a k-means style alternating optimization algorithm (and not “soft” k-means or EM) for estimating the model $a^1, a^2 \in [0.5, M - 0.5]$. Use random initialization for sake of simplicity. Give precise update expressions in pseudocode and not just vague statements.

Give solution to 6.1 here

Since we have uniform priors i.e. $\mathbb{P}[z^i = c | a^1, a^2] = 0.5$ for all $c \in \{1, 2\}$, we get

$$\operatorname{argmax}_{c \in \{1, 2\}} \mathbb{P}[z^i = c | x^i, a^1, a^2] = \operatorname{argmax}_{c \in \{1, 2\}} \mathbb{P}[x^i | z^i = c, a^1, a^2]$$

The above can be simplified as follows by breaking ties lexicographically

1. If $|x^i - a^1| \leq 0.5$
 - a. Set $z^i = 1$
2. Else
 - a. Set $z^i = 2$

Give solution to 6.2 here

Assume a uniform prior over a^1, a^2 i.e. $\mathbb{P}[a^1, a^2 | \mathbf{z}] = \mathbb{P}[a^1, a^2] = \text{UNIF}[0.5, M - 0.5]$ (any other valid assumption will also get marks) so that we have $\arg \max_{a^1, a^2 \in [0.5, M-0.5]} \mathbb{P}[a^1, a^2 | \mathbf{X}, \mathbf{z}] =$

$\arg \max_{a^1, a^2 \in [0.5, M-0.5]} \mathbb{P}[\mathbf{X} | a^1, a^2, \mathbf{z}]$. Denote $N_1 \triangleq |i: z^i = 1 \wedge |x^i - a^1| \leq 0.5|$. Similarly define N_2 .

Note that $\mathbb{P}[\mathbf{X} | a^1, a^2, \mathbf{z}] \propto \left(\frac{1}{2} + \frac{1}{2M}\right)^{N_1+N_2} \left(\frac{1}{2M}\right)^{n-(N_1+N_2)}$. Since we are told that $M \gg 1$, to maximize $\mathbb{P}[\mathbf{X} | a^1, a^2, \mathbf{z}]$ w.r.t a^1 , we must try to maximize N_1 i.e. find a value of a^1 that contains as many points with $z^i = 1$ as possible within a 0.5 radius ball around it. The following algorithm does this for a^1 . A similar algorithm works to obtain the best value of a^2 .

1. Let $S_1 \triangleq \{i: z^i = 1\}$
2. Let $C_1 = \{x^i + 0.5: i \in S_1\} \cup \{x^i - 0.5: i \in S_1\}$ *# $x^i \pm 0.5$ are candidates for a^1*
3. Let $N_1 = 0, a^1 = 0.5$ *# Initialize to valid values*
4. For all values $p \in C_1$:
 - a. If $p \notin [0.5, M - 0.5]$: *# Only valid candidates considered*
 - i. Continue
 - b. Let $N \triangleq \text{r-nn}(p, 0.5) \cap S_1$ *# All points in S_1 within a 0.5 radius*
 - c. If $N > N_1$: *# Found a better candidate - yay!*
 - i. $N_1 = N, a^1 = p$
5. Return a^1

Give solution to 6.3 here

The following is simply the above two steps carried out alternately.

1. Initialize $a^1, a^2 \sim \text{UNIF}[0.5, M - 0.5]$ *# Can use k-means++ as well*
2. Repeat till converged to till timeout:
 - a. For all $i \in [n]$:
 - i. If $|x^i - a^1| \leq 0.5$, set $z^i = 1$ else set $z^i = 2$
 - b. Define S_c, C_c for $c \in \{1, 2\}$ as shown above
 - c. Initialize $N_c = 0, a^c = 0.5$ for $c \in \{1, 2\}$
 - d. For $c \in \{1, 2\}$:
 - i. For all values $p \in C_c$:
 1. If $p \in [0.5, M - 0.5]$:
 - a. Let $N \triangleq \text{r-nn}(p, 0.5) \cap S_c$
 - b. If $N > N_c$: $N_c = N, a^c = p$
6. Return $\{a^c, c \in \{1, 2\}\}$