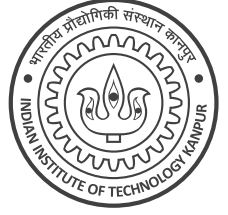


CS 771A: Intro to Machine Learning, IIT Kanpur			Quiz II (20 Mar 2024)	
Name	MELBO			20 marks Page 1 of 2
Roll No	240007	Dept.	AWSM	

Instructions:

1. This question paper contains 1 page (2 sides of paper). Please verify.
2. Write your name, roll number, department above in **block letters neatly with ink**.
3. Write your final answers neatly **with a blue/black pen**. Pencil marks may get smudged.
4. Don't overwrite/scratch answers especially in MCQ – such cases will get straight 0 marks.
5. Do not rush to fill in answers. You have enough time to solve this quiz.



(Noise to Regularize) The underlying principle behind the deep learning technique *dropout* is that adding noise to data can prevent models from overfitting. Let us derive this fact formally.

Q1. Let $\epsilon \in \{-1, +1\}^D$ be a D -dim Rademacher vector with coordinates chosen i.i.d. $\epsilon_j = 1$ or -1 uniformly randomly. Find the following (no derivation) Note: $j, k \in [D], j \neq k$. **(6 x 1 = 6 marks)**

$$\mathbb{E}[\epsilon_j + \epsilon_k] = 0$$

$$\text{Var}[\epsilon_j + \epsilon_k] = 2$$

$$\mathbb{E}[\epsilon_j \epsilon_k] = 0$$

$$\text{Var}[\epsilon_j \epsilon_k] = 1$$

$$\mathbb{E}[\epsilon_j / \epsilon_k] = 0$$

$$\text{Var}[\epsilon_j / \epsilon_k] = 1$$

Q2. Let $y, \lambda \in \mathbb{R}$ and $\mathbf{w}, \mathbf{x} \in \mathbb{R}^D$ be constants and $\epsilon \in \{-1, +1\}^D$ be a Rademacher vector sampled independently of $y, \lambda, \mathbf{w}, \mathbf{x}$. Obtain a simplified expression (expectation is over the choice of ϵ only). Give brief derivation. Your expression should not contain any ϵ_j terms. **(2 + 4 = 6 marks)**

Write final expression in the box

$$\mathbb{E}_{\epsilon} \left[(y - \mathbf{w}^T (\mathbf{x} + \lambda \cdot \epsilon))^2 \right] = (y - \mathbf{w}^T \mathbf{x})^2 + \lambda^2 \cdot \|\mathbf{w}\|_2^2 \text{ or else } (y - \mathbf{w}^T \mathbf{x})^2 + \lambda^2 \cdot \mathbf{w}^T \mathbf{w}$$

Brief derivation for simplification

Expanding the expression gives us $\mathbb{E}_{\epsilon} [y^2 + (\mathbf{w}^T \mathbf{x} + \lambda \cdot \mathbf{w}^T \epsilon)^2 - 2y(\mathbf{w}^T (\mathbf{x} + \lambda \cdot \epsilon))]$

Using linearity of expectation yields $y^2 + (\mathbf{w}^T \mathbf{x})^2 - 2y \cdot \mathbf{w}^T \mathbf{x} + \mathbb{E}_{\epsilon} [(\lambda \cdot \mathbf{w}^T \epsilon)^2 - 2y\lambda \cdot \mathbf{w}^T \epsilon]$

Using the fact that $\mathbb{E}_{\epsilon} [\epsilon] = \mathbf{0}$ gives us $y^2 + (\mathbf{w}^T \mathbf{x})^2 - 2y \cdot \mathbf{w}^T \mathbf{x} + \lambda^2 \cdot \mathbb{E}_{\epsilon} [(\mathbf{w}^T \epsilon)^2]$

Expanding the last term gives us $\mathbb{E}_{\epsilon} [(\mathbf{w}^T \epsilon)^2] = \mathbb{E}_{\epsilon} \left[\sum_{d \in [D]} w_d^2 \epsilon_d^2 + \sum_{\substack{d \neq d' \\ d, d' \in [D]}} w_d w_{d'} \epsilon_d \epsilon_{d'} \right]$

Using results from Q1 simplifies this to $\mathbb{E}_{\epsilon} [(\mathbf{w}^T \epsilon)^2] = \sum_{d \in [D]} w_d^2 = \mathbf{w}^T \mathbf{w} = \|\mathbf{w}\|_2^2$

Completing the squares then yields the final expression.

Q3. We have N datapoints $(\mathbf{x}^n, y^n) \in \mathbb{R}^D \times \mathbb{R}, n \in [N], \lambda \in \mathbb{R}, \mathbf{w} \in \mathbb{R}^D$ all of which can be treated as constants. We also sample N Rademacher vectors $\epsilon^n \in \{-1, +1\}^D, n \in [N]$ i.i.d. of each other as well as independent of the datapoints and λ, \mathbf{w} . Expectation is over the choice of $\{\epsilon^n, n \in [N]\}$ only. Write down a simplified expression for the following (no derivation needed). **(2 marks)**

$$\mathbb{E}_{\{\epsilon^n\}} \left[\sum_{n \in [N]} (y^n - \mathbf{w}^\top (\mathbf{x}^n + \lambda \cdot \epsilon^n))^2 \right] = \left(\sum_{n \in [N]} (y^n - \mathbf{w}^\top \mathbf{x}^n)^2 \right) + N\lambda^2 \cdot \|\mathbf{w}\|_2^2$$

or else $\sum_{n \in [N]} ((y^n - \mathbf{w}^\top \mathbf{x}^n)^2 + \lambda^2 \cdot \|\mathbf{w}\|_2^2)$.

Note: this is exactly the objective function that ridge regression tries to minimize ☺☺

Q4 (IPL Intrigue). Melbo is a big IPL fan and is trying to analyse the performance of MI vs CSK on various kinds of pitches. Let M be the event that MI won a MI-vs-CSK match and C be the event that CSK won a MI-vs-CSK match. There are 3 kinds of pitches $F = \text{flat}, G = \text{green}, D = \text{dusty}$. A total of 24 matches were played between MI and CSK, $1/4^{\text{th}}$ of which were on green pitches and $1/3^{\text{rd}}$ on flat pitches. MI won 6 of the matches played on flat pitches. Both MI and CSK won equal number of matches played on green pitches i.e., $\mathbb{P}[M | G] = \mathbb{P}[C | G]$. Also, both flat and dusty pitches have been equally favourable for MI in that $\mathbb{P}[F | M] = \mathbb{P}[D | M]$. Find out the following quantities as fractions or decimals (no derivations needed). Hint: either use Bayes rule or fill-up a 2×3 matrix showing which team won how many matches on what kind of pitch. **(6 x 1 = 6 marks)**

$$\mathbb{P}[F | M] = \frac{2}{5}$$

$$\mathbb{P}[F | C] = \frac{2}{9}$$

$$\mathbb{P}[G | M] = \frac{1}{5}$$

$$\mathbb{P}[G | C] = \frac{1}{3}$$

$$\mathbb{P}[D | M] = \frac{2}{5}$$

$$\mathbb{P}[D | C] = \frac{4}{9}$$

$$\mathbb{P}[F] = \frac{1}{3}, \mathbb{P}[G] = \frac{1}{4}, \mathbb{P}[D] = \frac{5}{12} \text{ which gives us } \mathbb{P}[M | F] = \frac{\mathbb{P}[M \cap F]}{\mathbb{P}[F]} = \frac{\frac{6}{24}}{\frac{1}{3}} = \frac{3}{4} \text{ and } \mathbb{P}[C | F] = \frac{1}{4}.$$

$$\text{Since } \mathbb{P}[M | G] + \mathbb{P}[C | G] = 1, \text{ we get } \mathbb{P}[M | G] = \mathbb{P}[C | G] = \frac{1}{2}. \text{ Since } \mathbb{P}[F | M] = \mathbb{P}[D | M], \text{ we get } \mathbb{P}[M | D] = \frac{\mathbb{P}[D | M] \cdot \mathbb{P}[M]}{\mathbb{P}[D]} = \frac{\mathbb{P}[F | M] \cdot \mathbb{P}[M]}{\mathbb{P}[D]} = \mathbb{P}[M | F] \cdot \frac{\mathbb{P}[F]}{\mathbb{P}[D]} = \frac{3}{5} \text{ and } \mathbb{P}[C | D] = \frac{2}{5}.$$

The law of total probability then allows us to calculate the following:

$$\mathbb{P}[M] = \mathbb{P}[M | F] \cdot \mathbb{P}[F] + \mathbb{P}[M | G] \cdot \mathbb{P}[G] + \mathbb{P}[M | D] \cdot \mathbb{P}[D] = \frac{5}{8} \text{ and } \mathbb{P}[C] = 1 - \mathbb{P}[M] = \frac{3}{8}.$$

$$\text{Applying the Bayes rule then tells us that } \mathbb{P}[F | M] = \frac{\mathbb{P}[M | F] \cdot \mathbb{P}[F]}{\mathbb{P}[M]} = \frac{2}{5} = \mathbb{P}[D | M] \text{ as promised and } \mathbb{P}[G | M] = 1 - \mathbb{P}[F | M] - \mathbb{P}[D | M] = \frac{1}{5}.$$

$$\text{Applying the Bayes rule again tells us } \mathbb{P}[F | C] = \frac{\mathbb{P}[C | F] \cdot \mathbb{P}[F]}{\mathbb{P}[C]} = \frac{2}{9}, \mathbb{P}[G | C] = \frac{\mathbb{P}[C | G] \cdot \mathbb{P}[G]}{\mathbb{P}[C]} = \frac{1}{3} \text{ and } \mathbb{P}[D | C] = 1 - \mathbb{P}[F | C] - \mathbb{P}[G | C] = \frac{4}{9}.$$