

Introduction to ML (CS771), 2024-2025-Sem-I		Total Marks	25
Quiz 3. October 24, 2024		Duration	45 minutes
Name		Roll No.	

**Instructions:**

1.	Clearly write your name (in block letters) and roll number in the provided boxes above.
2.	Write your final answers concisely in the provided space. You may use blue/black pen.
3.	We won't be able to provide clarifications during the quiz. If any aspect of some question appears ambiguous/unclear to you, please state your assumption(s) and answer accordingly.

**Question 1:** Write **T** or **F** for True/False in the box next to each question given below, with a brief (1-2 sentences at most) explanation in the provided space in the box below the question. Marks will be awarded only when the answer (T/F) and explanation both are correct. (3 x 2 = 6 marks)

1.1	To predict the label using a generative classification model, comparing the probabilities $p(y = k \mathbf{x})$ for different values of $k$ is equivalent to comparing the class-conditional probability densities $p(\mathbf{x} y = k)$ for different values of $k$	<b>F</b>
$p(y = k \mathbf{x}) \propto p(y = k)p(\mathbf{x} y = k)$ so it also incorporate the class prior (class marginal) distribution $p(y = k)$		

1.2	A Gaussian prior $p(\mathbf{w}) = \mathcal{N}(\mathbf{w} \mathbf{w}_0, \lambda^{-1}\mathbf{I})$ on the weight vector $\mathbf{w} \in \mathbb{R}^D$ will cause a regularization effect and encourage the entries in $\mathbf{w}$ to take small values.	<b>F</b>
This prior corresponds to a regularizer of the form $\lambda\ \mathbf{w} - \mathbf{w}_0\ ^2$ which will encourage each entry of the vector $\mathbf{w}$ to be close the the corresponding entry in the vector $\mathbf{w}_0$ . Only when $\mathbf{w}_0$ is the zero vector, the statement above would be true but in general it would be false.		

1.3	Even though the MAP estimate is the mode of the posterior distribution, to compute the MAP estimate, it is not necessary to compute the posterior distribution.	<b>T</b>
Recall that the posterior is $p(\theta y) = \frac{p(\theta)p(y \theta)}{p(y)}$ . Because of the denominator (marginal likelihood) is independent of $\theta$ , maximization of the posterior only requires maximization of the numerator $p(\theta)p(y \theta)$ (or $\log p(\theta) + \log p(y \theta)$ ) and we don't need to compute the full posterior for the maximization.		

**Question 2:** Answer the following questions concisely in the space provided below the question.

2.1	Consider the RBF kernel $k(\mathbf{x}_i, \mathbf{x}_j) = \exp(-\gamma\ \mathbf{x}_i - \mathbf{x}_j\ ^2)$ where $\mathbf{x}_i$ and $\mathbf{x}_j$ are $D$ dim inputs. Consider two cases: (1) when bandwidth hyperparameter $\gamma$ is set as very-very large, and (2) when $\gamma$ is set as very-very small. For each of these two cases, answer (with brief justification) whether the resulting kernel function would be practically useful. (4 marks)
<p>(1) When <math>\gamma</math> is very-very large, <math>k(\mathbf{x}_i, \mathbf{x}_j)</math> will be nonzero (will equal 1) only when <math>\mathbf{x}_i</math> and <math>\mathbf{x}_j</math> are nearly identical. For all other pairs of inputs, the kernel will give 0 similarity.</p> <p>(2) When <math>\gamma</math> is very-very small, <math>k(\mathbf{x}_i, \mathbf{x}_j)</math> will be close to 1 for all pairs of inputs, thus treating all pairs of inputs as equally similar to each other.</p> <p>Clearly, neither of these two extreme cases are desirable.</p>	

2.2	<p>Briefly explain why using kernels with the landmarks approach or the random features approach is faster at test time than using kernels in the standard manner? <b>(3 marks)</b></p> <p>When using landmarks or random features approach, we use the kernel to construct an <math>L</math> dimensional feature representation <math>\psi(\mathbf{x}_n)</math>, and train a linear model on these representations to get a weight vector <math>\mathbf{w}</math> that is <math>L</math> dimensional. Thus, for a test input <math>\mathbf{x}_*</math> the prediction cost for computing <math>\mathbf{w}^\top \psi(\mathbf{x}_*)</math> is also <math>O(L)</math>. In contrast, when using the kernel in the standard manner, this cost of <math>O(N)</math> which can be very high if the number of training inputs is very large</p>
2.3	<p>Given a dataset <math>\mathbf{X}</math> as the <math>N \times D</math> input matrix with <math>N</math> inputs and <math>D</math> features, write down the <math>K</math>-means hard-clustering problem for this dataset in form of an equivalent matrix factorization problem, clearly specifying the meanings of the variables involved in the matrix factorization, their dimensions, and constraints on them, if any. <b>(4 marks)</b></p> <p><math>\{\hat{\mathbf{Z}}, \hat{\boldsymbol{\mu}}\} = \operatorname{argmin}_{\mathbf{Z}, \boldsymbol{\mu}} \ \mathbf{X} - \mathbf{Z}\boldsymbol{\mu}\ ^2</math>. Here <math>\mathbf{Z}</math> is the <math>N \times K</math> matrix with row <math>n</math> (<math>\mathbf{z}_n</math>) being a one-hot vector denoting which cluster the input <math>\mathbf{x}_n</math> belongs to, and <math>\boldsymbol{\mu}</math> denotes the <math>K \times D</math> matrix with row <math>k</math> (<math>\boldsymbol{\mu}_k</math>) denoting the mean of the <math>k^{th}</math> cluster.</p> <p>Constraints on <math>\mathbf{z}_n</math>: Must be a one-hot vector</p> <p>Constraints on <math>\boldsymbol{\mu}_k</math>: None</p>
2.4	<p>Why is it difficult to compute the predictive distribution of a logistic regression model which, by definition, is given by <math>p(y_* = 1 \mathbf{x}_*, \mathbf{X}, \mathbf{y}) = \int p(y_* = 1 \mathbf{w}, \mathbf{x}_*)p(\mathbf{w} \mathbf{X}, \mathbf{y})d\mathbf{w}</math>. Suggest a method to approximate it and clearly show the necessary equations. <b>(3 marks)</b></p> <p>It is difficult because the integral here is not tractable (involves integrating <math>p(y_* = 1 \mathbf{w}, \mathbf{x}_*)</math> which is a sigmoid function over the posterior <math>p(\mathbf{w} \mathbf{X}, \mathbf{y})</math> and even if the latter is Gaussian (like in Laplace approximation), the integral still is intractable. To approximate the integral, one way is to use Monte-Carlo approximation where we draw <math>S</math> i.i.d. samples <math>\mathbf{w}^{(1)}, \mathbf{w}^{(2)}, \dots, \mathbf{w}^{(S)}</math> from the posterior and approximate the predictive distribution as</p> $p(y_* = 1 \mathbf{x}_*, \mathbf{X}, \mathbf{y}) \approx \frac{1}{S} \sum_{s=1}^S p(y_* = 1 \mathbf{w}^{(s)}, \mathbf{x}_*)$
2.5	<p>Show that, for generative classification with uniform class marginal and Gaussian class conditionals <math>\mathcal{N}(\mathbf{x} \boldsymbol{\mu}_k, \boldsymbol{\Sigma})</math>, the posterior probability of input <math>\mathbf{x}</math> belonging to class <math>k</math>, i.e., <math>p(y = k \mathbf{x}) \propto \exp(\mathbf{w}_k^\top \mathbf{x} + b_k)</math>, and write down the expressions for <math>\mathbf{w}_k</math> and <math>b_k</math> <b>(5 marks)</b></p> <p>Since we have a uniform class marginal, the posterior probability will be</p> $p(y = k \mathbf{x}) = \frac{p(\mathbf{x} y = k)}{p(\mathbf{x})} \propto p(\mathbf{x} y = k)$ <p>Since the class conditional <math>p(\mathbf{x} y = k) = \mathcal{N}(\mathbf{x} \boldsymbol{\mu}_k, \boldsymbol{\Sigma})</math>, we have</p> $p(\mathbf{x} y = k) \propto \exp\left(-\frac{1}{2}((\mathbf{x} - \boldsymbol{\mu}_k)\boldsymbol{\Sigma}^{-1}(\mathbf{x} - \boldsymbol{\mu}_k))\right) \propto \exp\left(\boldsymbol{\mu}_k^\top \boldsymbol{\Sigma}^{-1} \mathbf{x} - \frac{1}{2} \boldsymbol{\mu}_k^\top \boldsymbol{\Sigma} \boldsymbol{\mu}_k\right),$ <p>where in the last expression (after the proportionality sign), we are ignoring any terms that are not specific to class <math>k</math>. Thus <math>p(y = k \mathbf{x}) \propto p(\mathbf{x} y = k) \propto \exp\left(\boldsymbol{\mu}_k^\top \boldsymbol{\Sigma}^{-1} \mathbf{x} - \frac{1}{2} \boldsymbol{\mu}_k^\top \boldsymbol{\Sigma} \boldsymbol{\mu}_k\right)</math> which is clearly in the form of <math>\exp(\mathbf{w}_k^\top \mathbf{x} + b_k)</math> where <math>\mathbf{w}_k = (\boldsymbol{\mu}_k^\top \boldsymbol{\Sigma}^{-1})^\top = \boldsymbol{\Sigma}^{-1} \boldsymbol{\mu}_k</math> and <math>b_k = -\frac{1}{2} \boldsymbol{\mu}_k^\top \boldsymbol{\Sigma} \boldsymbol{\mu}_k</math></p> <p>Side note (not required for the answer): Note that the above implies that this generative classification model has a similar form as softmax classification model, although the weight vector is learned using a generative manner, and not using GD as we do in case of softmax classification</p>