

Name: Roll No.: Dept.: 

IIT Kanpur  
CS771 Intro to ML  
End-semester Examination  
Date: November 29, 2018

**Instructions:****Total: 120 marks**

1. This question paper contains a total of 10 pages (10 sides of paper). Please verify.
2. Please write your name, roll number, department on **every side of every sheet** of this booklet.
3. You may write your answers using pencil but your handwriting should be bold and prominently visible.
4. **Important:** Please do not give derivations/elaborate steps unless specifically asked for it. Feel free to use standard results (e.g., solution of least squares regression) without deriving them from scratch.
5. The last page of the question paper lists some formulae if you need them.

**Section 1 (True or False:  $12 \times 1 = 12$  marks).** For each of the following simply write **T** or **F** in the box.

1. ☐ The kernel SVM weight vector  $\mathbf{w}$  can be written explicitly as a finite dimensional vector only when using a linear kernel (assuming we aren't using any approximations, such as landmarks).
2. ☐ Learning a single hidden layer neural network with infinite many hidden units is equivalent to learning a kernelized model with an RBF kernel.
3. ☐ Both alternating optimization (ALT-OPT), as well as the expectation maximization (EM) algorithm, are sensitive to initialization.
4. ☐ It is possible to get closed form solutions for all the parameters of a fully supervised generative classification model with Gaussian class-conditionals.
5. ☐ A  $K$ -nearest neighbors classifier that uses Euclidean distances can only learn a linear decision boundary regardless of the value of  $K$ .
6. ☐ A depth-1 decision tree will usually have a higher bias than a depth-5 decision tree (here "bias" is used in the sense of this word as in the bias-variance trade-off).
7. ☐ If the training inputs and test inputs for a classification problem are drawn from the same distribution then the test error of the learned model will be zero.
8. ☐ Iteration  $t + 1$  of Adaboost is faster than iteration  $t$  because iteration  $t + 1$  trains only using the misclassified examples from iteration  $t$ .
9. ☐ MAP estimation for a parameter, when using a Gaussian prior with zero mean and spherical covariance, is equivalent to doing MLE for the parameter.
10. ☐ If the gap between training and test error is large for a model then retraining the model with larger training set may reduce the gap.
11. ☐ Probabilistic PCA with noise variance equal to zero and classic PCA will give the same solution for the projection matrix, assuming mean-centered data for both the methods.
12. ☐ A feedforward neural network's output layer computes a convex combination of the outputs of the last hidden layer nodes.

**Section 2 (MCQ:  $12 \times 2 = 24$  marks).** **Tick-mark** ☒ all the options that you think are correct. **Important:** Marks for a question will be awarded only when all correct options (and only those) are selected.

1. Which of these can be used for regression? ☐ A Decision Tree, ☐ B  $K$ -nearest neighbor, ☐ C Perceptron, ☐ D Feedforward neural network, ☐ E Logistic regression.
2. Which of these can be kernelized? ☐ A  $K$ -nearest neighbors, ☐ B Decision trees, ☐ C  $K$ -means clustering, ☐ D Principal Component Analysis, ☐ E Prototype based classification.
3. Which of these are linear dimensionality reduction methods: ☐ A Probabilistic PCA, ☐ B Standard PCA, ☐ C Fisher Discriminant Analysis, ☐ D Stochastic Neighbor Embedding, ☐ E Locally linear embedding.
4. Which of these objectives are non-differentiable? ☐ A Squared loss with  $\ell_1$  regularizer, ☐ B Hinge loss with  $\ell_2$  reg., ☐ C Hinge loss with  $\ell_1$  reg., ☐ D Huber loss with  $\ell_2$  reg., ☐ E Huber loss with  $\ell_1$  reg.

Name: Roll No.: Dept.: 

**IIT Kanpur**  
**CS771 Intro to ML**  
**End-semester Examination**  
 Date: November 29, 2018

5. Which of these can only learn linear decision boundaries? ☐ A SVM with quadratic kernel, ☐ B Decision tree classifier, ☐ C Prototype based classification with Euclidean distances, ☐ D Single hidden layer neural net with ReLU activations, ☐ E Logistic regression with score being linear combination of the features.
6. Which of these learning problems/sub-problems require constrained optimization? ☐ A Solving for the mixing proportion weights in a mixture model, ☐ B Learning the standard Perceptron, ☐ C Learning PPCA, ☐ D Learning the kernel SVM, ☐ E Value-iteration based policy learning in reinforcement learning.
7. Which of the following are true? ☐ A When the regularization hyperparam. tends to infinity, regularization becomes ineffective, ☐ B  $\ell_1$  norm is non-convex, ☐ C  $\ell_2$  norm is convex, ☐ D,  $\ell_1$  norm promotes non-negativity, ☐ E Using a Laplace prior is equivalent to using an  $\ell_2$  regularizer.
8. The output of a matrix factorization model learned using only user-item ratings matrix can be used to: ☐ A Find other users similar to a given user, ☐ B Find other items similar to a given item, ☐ C Recommend existing items to new users, ☐ D Learn clusters of items, ☐ E Learn clusters of users.
9. How can we turn a linear classifier into a nonlinear one? ☐ A First project the inputs to a low-dim space using PCA, ☐ B First project the inputs to a low-dim space using Fisher Discriminant Analysis, ☐ C Use it as a base learner in Adaboost, ☐ D Use scores of  $K > 1$  such classifiers to get  $K$  new features and learn another linear classifier on those features. ☐ E Cluster inputs and learn a linear classifier for each cluster.
10. Posterior can be computed in closed form for: ☐ A Linear regression with Gaussian likelihood, zero mean Gaussian prior, and fixed hyperparams, ☐ B Linear regression with Gaussian likelihood, non-zero mean Gaussian prior, and fixed hyperparams ☐ C Logistic regression with Gaussian prior, ☐ D Bernoulli coin-toss model with Beta prior on coin's bias, ☐ E Gaussian mean estimation with Gaussian prior on mean.
11. Which of the following are true about  $KNN$ : ☐ A Very fast at test time, ☐ B Tends to underfit as  $K$  increases, ☐ C Have zero error on training data, ☐ D Equivalent to prototype based classification for  $K = 1$ , ☐ E Training them is computationally very expensive.
12. Which of the following is true about support vector machines (SVM)? ☐ A They are faster than decision trees at test time, ☐ B Multiclass SVMs are equivalent to softmax regression, ☐ C For linear SVM, every training example is a support vector, ☐ D Maximizing the SVM margin is equivalent to maximizing the  $\ell_2$  norm of SVM weight vector, ☐ E Increasing margin leads to more number of misclassified training examples.

**Section 3 (Short Answer:  $8 \times 4 = 32$  marks).** Write your answers precisely and concisely in the provided box.

1. Consider a generative model for binary classification. Suppose each input has 2 features, where the first feature takes one of 5 possible values and the second feature is binary. With naive Bayes assumption, how many parameters would we need to learn for this generative classification model. Justify your answer.

Name: Roll No.: Dept.: 

IIT Kanpur  
**CS771 Intro to ML**  
**End-semester Examination**  
*Date:* November 29, 2018

2. You are given  $N$  inputs  $\{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N\}$ . Suppose, for each  $\mathbf{x}_n \in \mathbb{R}^D$ , you want to obtain a  $K$  dimensional and *non-sparse* feature vector  $\mathbf{z}_n$ , where the sum of the  $K$  features is one. Briefly describe how you would compute such feature vectors  $\{\mathbf{z}_1, \mathbf{z}_2, \dots, \mathbf{z}_N\}$ , using a  $K$ -means clustering algorithm on this data?

3. Consider a linear model with a regularizer  $R(\mathbf{w}) = \|\mathbf{w}\|^2 + \sum_{d=1}^D \sum_{d'=d+1}^D (w_d - w_{d'})^2$ . What will be the effect of such a regularizer on  $\mathbf{w}$  when minimizing the objective  $\sum_{n=1}^N \ell(y_n, \mathbf{w}^\top \mathbf{x}_n) + \lambda R(\mathbf{w})$  w.r.t.  $\mathbf{w}$ ?

4. In at most 1-3 sentences (preferably only words, no equations!), describe how additional unlabeled data can be utilized within an algorithm for learning the parameters of a generative classification model.

5. Can we compute the squared  $\ell_2$  norm  $\|\mathbf{w}\|^2$  of the kernel ridge regression weights  $\mathbf{w} = \sum_{n=1}^N \alpha_n \phi(\mathbf{x}_n)$ , assuming  $\phi$  to be the feature mapping of an RBF kernel? If yes, show how it can be done. If no, clearly state why it can't be done. Also answer the same question if we want to compute the  $\ell_1$  norm of  $\mathbf{w}$ .

6. Assume a model  $f$  applied to a two-class data. Suppose the PDF of scores  $s \in (-\infty, \infty)$  of  $f$  on positive examples is  $p_1(s)$ , whereas the PDF of  $f$ 's scores on negative examples is  $p_0(s)$  (assume positive examples to be 1 and negative examples to be 0). Suppose  $F_0$  denotes the CDF of  $p_0$  and  $F_1$  denotes the CDF of  $p_1$ . What's the false negative rate (FNR) and the true negative rate (TNR) of  $f$ ?

Name: Roll No.: Dept.: 

IIT Kanpur  
CS771 Intro to ML  
End-semester Examination  
Date: November 29, 2018

7. Suppose you have two coins  $c_1$  and  $c_2$  with biases  $\pi_1 \in (0, 1)$  and  $\pi_2 \in (0, 1)$ , respectively. You have another coin  $c_3$  with bias  $\mu \in (0, 1)$ . You do two coin tosses as follows: First, you toss coin  $c_3$ . If it shows heads, you toss coin  $c_1$ ; otherwise you toss coin  $c_2$ . Denote the outcome of the second toss (i.e., when you toss  $c_1$  or  $c_2$ ) as  $x \in \{0, 1\}$ . What is the marginal *distribution*  $p(x)$ ? Clearly write down its expression.

8. Taking the example of a single hidden layer feedforward neural network, show that it is necessary to have nonlinearities in the hidden nodes, in the absence of which the network would reduce to a linear model.

**Section 4** (5 problems:  $5 \times 8 = 40$  marks). Write your answers precisely and concisely in the provided box.

1. Derive and write down the SGD update (minibatch size = 1) for the linear regression model with squared loss  $\sum_{n=1}^N (y_n - \mathbf{w}^\top \mathbf{x}_n)^2$ . Using the SGD update equation, formally show that each SGD update does the *right* thing, i.e., it *improves* the model's prediction on the current example  $(\mathbf{x}_n, y_n)$ .

Name: Roll No.: Dept.: 

IIT Kanpur  
**CS771 Intro to ML**  
**End-semester Examination**  
*Date:* November 29, 2018

2. Consider the prototype based classification model, given training data  $\{(\mathbf{x}_n, y_n)\}_{n=1}^N$ , where input  $\mathbf{x}_n \in \mathbb{R}^D$  and label  $y_n \in \{-1, +1\}$ . Suppose we have mapped the inputs to a new feature space  $\phi$  that has an associated kernel function  $k(.,.)$ . Show that the prediction for a new test input  $\mathbf{x}_*$  can be written in form of  $y_* = \text{sign}[f(\mathbf{x}_*)]$  and clearly write down the expression for  $f(\mathbf{x}_*)$ . The expression for  $f(\mathbf{x}_*)$  must be only in terms of the kernel function  $k$ , and must not contain the feature mapping  $\phi$  in it.

3. Consider  $K$ -means clustering where we are trying to learn  $K$  means  $\mu_1, \dots, \mu_K$ , given  $N$  observations  $\{\mathbf{x}_1, \dots, \mathbf{x}_N\}$ , with each  $\mathbf{x}_n \in \mathbb{R}^D$ . Suppose we have some *a priori* information that the  $K$  means are “close” to known vectors  $\mu_1^*, \dots, \mu_K^*$ , respectively. Propose a suitable prior for each mean  $\mu_k$  that makes use of this information. For any iteration of  $K$ -means, given the current observation-to-cluster assignments  $\{z_1, \dots, z_N\}$ , and your proposed prior distribution, derive the update equation for each mean.

Name: Roll No.: Dept.: 

IIT Kanpur  
**CS771 Intro to ML**  
**End-semester Examination**  
*Date:* November 29, 2018

4. Consider learning a linear regression model by minimizing the squared loss function  $\sum_{n=1}^N (y_n - \mathbf{w}^\top \mathbf{x}_n)^2$ . Suppose we decide to mask out or “drop” each feature  $x_{nd}$  of each input  $\mathbf{x}_n \in \mathbb{R}^D$ , independently, with probability  $1 - p$  (equivalently, retaining the feature with probability  $p$ ). Masking or dropping out basically means that we will set the feature  $x_{nd}$  to 0 with probability  $1 - p$ . Essentially, it would be equivalent to replacing each input  $\mathbf{x}_n$  by  $\tilde{\mathbf{x}}_n = \mathbf{x}_n \circ \mathbf{m}_n$ , where  $\circ$  denotes elementwise product and  $\mathbf{m}_n$  denotes the  $D \times 1$  binary mask vector with  $m_{nd} \sim \text{Bernoulli}(p)$  ( $m_{nd} = 1$  means the feature  $x_{nd}$  was retained;  $m_{nd} = 0$  means the feature  $x_{nd}$  was masked/zeroed).

Let us now define a new loss function using these masked inputs as follows:  $\sum_{n=1}^N (y_n - \mathbf{w}^\top \tilde{\mathbf{x}}_n)^2$ . Show that minimizing the *expected* value of this new loss function (where the expectation is used since the mask vectors  $\mathbf{m}_n$  are random) is equivalent to minimizing a **regularized** loss function. Clearly write down the expression of this regularized loss function. (PS: You did something like this in Practice Set 1).

5. Consider the full (not truncated) singular value decomposition (SVD) of an  $N \times D$  matrix  $\mathbf{X}$ . Denote it as  $\mathbf{X} = \mathbf{U}\mathbf{\Lambda}\mathbf{V}^\top$ . Show that the left singular vectors  $\mathbf{U} = [\mathbf{u}_1, \dots, \mathbf{u}_N]$  are also the eigenvectors of  $\mathbf{X}\mathbf{X}^\top$ . Also show that the right singular vectors  $\mathbf{V} = [\mathbf{v}_1, \dots, \mathbf{v}_D]$  are also the eigenvectors of  $\mathbf{X}^\top \mathbf{X}$ .

Name: Roll No.: Dept.: 

IIT Kanpur  
**CS771 Intro to ML**  
**End-semester Examination**  
*Date: November 29, 2018*

**Section 5** (1 problem: 12 marks). Write your answers precisely and concisely in the provided box.

1. Assume you are given  $N$  examples  $\{(\mathbf{x}_n, y_n)\}_{n=1}^N$ , with each  $\mathbf{x}_n \in \mathbb{R}^D$  and  $y_n \in \mathbb{R}$ . Assume the following generative story for each  $(\mathbf{x}_n, y_n)$ : (1) Generate  $z_n \sim \text{multinoulli}(\pi_1, \dots, \pi_K)$ , (2) Generate the inputs  $\mathbf{x}_n \sim \mathcal{N}(\mu_{z_n}, \Sigma_{z_n})$ , and (3) Generate the outputs as  $y_n \sim \mathcal{N}(\mathbf{w}_{z_n}^\top \mathbf{x}_n, \beta^{-1})$ .

Your goal is to estimate the parameters  $\Theta = \{\pi_k, \mu_k, \Sigma_k, \mathbf{w}_k\}_{k=1}^K$  of this model. Assume  $\beta$  to be fixed.

- You have to derive an EM algorithm to compute the posterior distribution over unknowns  $\mathbf{Z} = \{z_1, \dots, z_N\}$  and point estimate (MLE) of unknowns  $\Theta$ . To do so, first write down the expression for the complete-data log-likelihood (CLL) for the model, and simplify it (ignore the constants).
- Now derive the necessary expressions that you would need for the EM algorithm for this model. If some of these derivations are obvious/familiar to you, you can skip those and directly write down the final expressions (but these expressions better be correct; no partial marks can be given for incorrect expressions in such a case :)). Also give a brief sketch of the overall EM algorithm.
- Assuming  $\pi_k = 1/K, \forall k$ , derive the ALT-OPT algorithm for this model (you may use the results from the above EM algorithm to get the ALT-OPT algorithm directly, without deriving from scratch). The ALT-OPT algorithm will compute point estimates for both  $\mathbf{Z}$  and  $\Theta$ . Also give a brief sketch of the overall ALT-OPT algorithm.

(if needed, you may continue the answer in the box on the next page)

Name:

Roll No.:

Dept.:

**IIT Kanpur**  
**CS771 Intro to ML**  
**End-semester Examination**  
*Date: November 29, 2018*

---



Name: Roll No.: Dept.: 

IIT Kanpur  
 CS771 Intro to ML  
 End-semester Examination  
 Date: November 29, 2018

### Some formulae you might need

- Bernoulli:  $\text{Bernoulli}(x|p) = p^x(1-p)^{1-x}$ . Expectation  $\mathbb{E}[x] = p$ , Variance  $\text{var}[x] = p(1-p)$
- Univariate Gaussian PDF:  $\mathcal{N}(x|\mu, \lambda^{-1}) = \sqrt{\frac{\lambda}{2\pi}} \exp(-\frac{\lambda}{2}(x-\mu)^2)$ ,  $\mathcal{N}(x|\mu, \sigma^2) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\{-\frac{(x-\mu)^2}{2\sigma^2}\}$
- Multivariate Gaussian PDF:  $\mathcal{N}(\mathbf{x}|\boldsymbol{\mu}, \boldsymbol{\Sigma}) = \frac{1}{\sqrt{(2\pi)^D|\boldsymbol{\Sigma}|}} \exp\{-\frac{1}{2}(\mathbf{x}-\boldsymbol{\mu})^\top \boldsymbol{\Sigma}^{-1}(\mathbf{x}-\boldsymbol{\mu})\}$ . Trace-based representation:  $\mathcal{N}(\mathbf{x}|\boldsymbol{\mu}, \boldsymbol{\Sigma}) = \frac{1}{\sqrt{(2\pi)^D|\boldsymbol{\Sigma}|}} \exp\left\{-\frac{1}{2}\text{trace}[\boldsymbol{\Sigma}^{-1}\mathbf{S}]\right\}$ ,  $\mathbf{S} = (\mathbf{x}-\boldsymbol{\mu})(\mathbf{x}-\boldsymbol{\mu})^\top$ .
- For  $x_k \in \{0, N\}$  and  $\sum_{k=1}^K x_k = N$ , multinomial( $x_1, \dots, x_K|N, \boldsymbol{\pi}$ ) =  $\frac{N!}{x_1! \dots x_K!} \pi_1^{x_1} \dots \pi_K^{x_K}$ , where  $\boldsymbol{\pi} = [\pi_1, \dots, \pi_K]$ , s.t.  $\sum_{k=1}^K \pi_k = 1$ . The multinoulli is the same as multinomial with  $N = 1$ .
- $\frac{\partial \mathbf{x}^\top \mathbf{a}}{\partial \mathbf{x}} = \frac{\partial \mathbf{a}^\top \mathbf{x}}{\partial \mathbf{x}} = \mathbf{a}$ , quadratic form:  $\frac{\partial}{\partial \mathbf{x}}(\mathbf{x}-\mathbf{s})^\top \mathbf{W}(\mathbf{x}-\mathbf{s}) = 2\mathbf{W}(\mathbf{x}-\mathbf{s})$
- $\frac{\partial}{\partial \boldsymbol{\mu}}[\boldsymbol{\mu}^\top \mathbf{A} \boldsymbol{\mu}] = [\mathbf{A} + \mathbf{A}^\top] \boldsymbol{\mu}$ ,  $\frac{\partial}{\partial \mathbf{A}} \log |\mathbf{A}| = \mathbf{A}^{-\top}$ ,  $\frac{\partial}{\partial \mathbf{A}} \text{trace}[\mathbf{A}\mathbf{B}] = \mathbf{B}^\top$
- For a random variable vector  $\mathbf{x}$ ,  $\mathbb{E}[\mathbf{x}\mathbf{x}^\top] = \mathbb{E}[\mathbf{x}]\mathbb{E}[\mathbf{x}]^\top + \text{cov}[\mathbf{x}]$
- For a random scalar  $x$ ,  $\text{var}[x] = \mathbb{E}[x^2] - \mathbb{E}[x]^2$

FOR ROUGH WORK ONLY

Name:

Roll No.:

Dept.:

**IIT Kanpur**  
**CS771 Intro to ML**  
**End-semester Examination**  
*Date: November 29, 2018*

---

FOR ROUGH WORK ONLY