

Name:

IIT Kanpur

CS771 Intro to ML

Roll No.: Dept.:

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. ☐ F The kernel SVM weight vector 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. ☐ F Learning a single hidden layer neural network with infinite many hidden units is equivalent to learning a kernelized model with an RBF kernel.
3. ☐ T Both alternating optimization (ALT-OPT), as well as the expectation maximization (EM) algorithm, are sensitive to initialization.
4. ☐ T It is possible to get closed form solutions for all the parameters of a fully supervised generative classification model with Gaussian class-conditionals.
5. ☐ F A K -nearest neighbors classifier that uses Euclidean distances can only learn a linear decision boundary regardless of the value of K .
6. ☐ T 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. ☐ F 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. ☐ F Iteration $t + 1$ of Adaboost is faster than iteration t because iteration $t + 1$ trains only using the misclassified examples from iteration t .
9. ☐ F 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. ☐ T 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. ☐ F 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. ☐ F 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 ℓ_1 regularizer, ☒ B Hinge loss with ℓ_2 reg., ☒ C Hinge loss with ℓ_1 reg., ☒ D Huber loss with ℓ_2 reg., ☒ E Huber loss with ℓ_1 reg.

Name: Roll No.: Dept.:

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 ℓ_1 norm is non-convex, ☒ C ℓ_2 norm is convex, ☐ D ℓ_1 norm promotes non-negativity, ☐ E Using a Laplace prior is equivalent to using an ℓ_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 ℓ_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.

$$X = [x_1, x_2]$$

$$P(x_1|y=0) = \text{multinoulli}(\pi_1^{(0)}, \pi_2^{(0)}, \pi_3^{(0)}, \pi_4^{(0)}, \pi_5^{(0)}) \Rightarrow 5 \text{ params}$$

$$P(x_2|y=0) = \text{Bern}(\mu^{(0)}) \Rightarrow 1 \text{ param}$$

$$\text{Likewise for } P(x_1|y=1) \text{ and } P(x_2|y=1) \Rightarrow 5+1$$

$$P(y) = \text{Bern}(\theta) \Rightarrow 1 \text{ param}$$

$$\text{Thus total} = 2 \times 6 + 1 = 13$$

$$\text{Also acceptable: } 2 \times (4+1) + 1 = 11$$

(actually 4 are sufficient since $\sum \pi_k = 1$)

Name:

IIT Kanpur

CS771 Intro to ML

Roll No.: Dept.:

End-semester Examination

Date: November 29, 2018

2. You are given N inputs $\{x_1, x_2, \dots, x_N\}$. Suppose, for each $x_n \in \mathbb{R}^D$, you want to obtain a K dimensional and *non-sparse* feature vector z_n , where the sum of the K features is one. Briefly describe how you would compute such feature vectors $\{z_1, z_2, \dots, z_N\}$, using a K -means clustering algorithm on this data?

Compute soft assignment for each x_n (just like soft K-means)
 e.g.
$$z_{nk} = \frac{\exp(-\|x_n - \mu_k\|^2)}{\sum_{j=1}^K \exp(-\|x_n - \mu_j\|^2)} \Rightarrow \sum_k z_{nk} = 1$$

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

It will promote w 's entries to not just be small but also similar to each other. Note: Suppose someone gave us a graph A ($D \times D$ binary matrix) with $A_{dd'} = 1$ denoting that features d and d' are similar then we can use $R(w) = \|w\|^2 + \sum_{d,d'} A_{dd'} (w_d - w_{d'})^2$

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.

We can use something like EM or ALT-OP to learn the parameters where, in each iteration, we compute a guess for the label of each unlabeled example and use these guesses when updating the parameters (we did it in one of the HW).

5. Can we compute the squared ℓ_2 norm $\|w\|^2$ of the kernel ridge regression weights $w = \sum_{n=1}^N \alpha_n \phi(x_n)$, assuming ϕ 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 ℓ_1 norm of w .

$$\|w\|^2 = w^T w = \left(\sum_{n=1}^N \alpha_n \phi(x_n) \right)^T \left(\sum_{m=1}^N \alpha_m \phi(x_m) \right) = \sum_{n=1}^N \sum_{m=1}^N \alpha_n \alpha_m \phi(x_n)^T \phi(x_m)$$

 Thus
$$\|w\|^2 = \sum_{n=1}^N \sum_{m=1}^N \alpha_n \alpha_m K(x_n, x_m) = \alpha^T K \alpha \text{ (if matrix-vector form)}$$

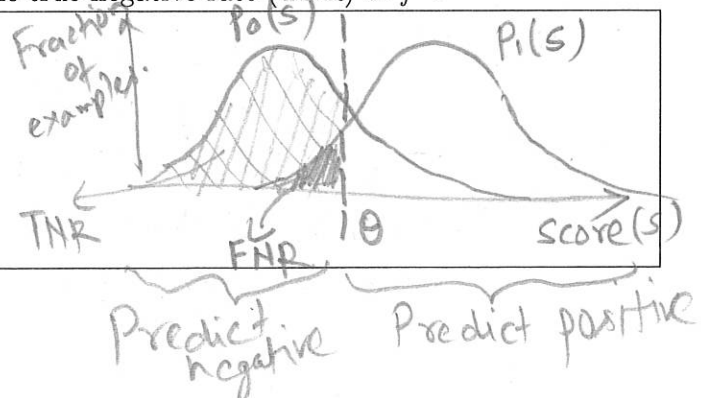
 However, we can't compute ℓ_1 norm since $\phi(x_n)$ is infinite-dim.

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 ?

For a threshold θ

$$\text{FNR} = F_1(\theta)$$

$$\text{TNR} = F_0(\theta)$$



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.

Mixture of two Bernoulli distributions

$$p(x) = \mu \times \text{Bern}(x|\pi_1) + (1-\mu) \text{Bern}(x|\pi_2)$$

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.

Without nonlinearity in hidden layer nodes: $h_n = Wx_n$
 $y_n = V^T h_n = V^T W x_n = \hat{W}^T x_n$ where $\hat{W} = (V^T W)^T$

Therefore we get a linear model if there is no nonlinearity.

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 $\sum_{n=1}^N (y_n - w^T x_n)^2$. 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 (x_n, y_n) .

$$\nabla \mathcal{L}(w) = -2(y_n - w^T x_n)x_n \quad (\text{using a single example } (x_n, y_n))$$

$$\text{SGD update: } w^{(t+1)} = w^{(t)} + \eta(y_n - w^{(t)T} x_n)x_n \quad (\eta \text{ subsumes the factor of 2, also } \eta > 0)$$

Suppose $y_n > w^{(t)T} x_n$, then note that

$$w^{(t+1)T} x_n = w^{(t)T} x_n + \underbrace{\eta(y_n - w^{(t)T} x_n)x_n^T x_n}_{> 0}$$

$$\Rightarrow w^{(t+1)T} x_n > w^{(t)T} x_n$$

which is an improvement since $w^T x_n$ value is moving closer to y_n .

can use a similar argument for the case when

$y_n < w^{(t)T} x_n$. In this case, $w^{(t+1)T} x_n < w^{(t)T} x_n$,

which is again an improvement since $w^T x_n$ value is moving closer to y_n .

Name:

IIT Kanpur

CS771 Intro to ML

End-semester Examination

Date: November 29, 2018

Roll No.: Dept.:

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

In prototype based classification (unkernelized case)

$$y_* = \text{Sign} \left[\underbrace{\|x_* - \mu_- \|^2}_{f(x_*)} - \|x_* - \mu_+ \|^2 \right]$$

In the kernelized case

$$f(x_*) = \underbrace{\|\phi(x_*) - \phi(\mu_-)\|^2}_{f(x_*)} - \|\phi(x_*) - \phi(\mu_+)\|^2$$

where $\phi(\mu_-)$ and $\phi(\mu_+)$ are means of negative and positive examples in the ϕ space, e.g. $\phi(\mu_-) = \frac{1}{N_-} \sum_{n: y_n = -1} \phi(x_n)$

substituting and expanding, we get

$$f(x_*) = \left\| \phi(x_*) - \frac{1}{N_-} \sum_{n: y_n = -1} \phi(x_n) \right\|^2 - \left\| \phi(x_*) - \frac{1}{N_+} \sum_{n: y_n = +1} \phi(x_n) \right\|^2$$

Common mistake of you have made: $\frac{2}{N_+} \sum_{n: y_n = +1} \phi(x_*)^T \phi(x_n) - \frac{2}{N_-} \sum_{n: y_n = -1} \phi(x_*)^T \phi(x_n) + \frac{1}{N_+^2} \sum_{n, m: y_n = y_m = +1} \phi(x_n)^T \phi(x_m) - \frac{1}{N_-^2} \sum_{n, m: y_n = y_m = -1} \phi(x_n)^T \phi(x_m)$

Annotations: $k(x_*, x_n)$, $k(x_*, x_m)$, $k(x_n, x_m)$

Important: $\phi(\mu_-)$ is not ϕ simply applied to $\frac{1}{N_-} \sum_{n: y_n = -1} x_n$ (likewise for $\phi(\mu_+)$) but

3. Consider K -means clustering where we are trying to learn K means μ_1, \dots, μ_K , given N observations $\{x_1, \dots, x_N\}$, with each $x_n \in \mathbb{R}^D$. Suppose we have some *a priori* information that the K means are "close" to known vectors μ_1^*, \dots, μ_K^* , respectively. Propose a suitable prior for each mean μ_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.

$P(\mu_k) = N(\mu_k^*, \lambda^{-1} I_D) \Rightarrow$ corresponds to a regularizer of the form $\lambda \|\mu_k - \mu_k^*\|^2$ for each μ_k .

The K -mean objective using this additional regularizer can be written as (assuming z given)

$$L(\mu_1, \dots, \mu_K) = \sum_{k=1}^K \left[\sum_{n: z_n = k} \|x_n - \mu_k\|^2 + \lambda \|\mu_k - \mu_k^*\|^2 \right]$$

For each μ_k , we need to solve

$$\hat{\mu}_k = \underset{\mu_k}{\text{argmin}} \sum_{n: z_n = k} \|x_n - \mu_k\|^2 + \lambda \|\mu_k - \mu_k^*\|^2$$

Taking derivative w.r.t μ_k and setting to zero:

$$\sum_{n: z_n = k} [-2(x_n - \mu_k) + 2\lambda(\mu_k - \mu_k^*)] = 0 \Rightarrow \hat{\mu}_k = \frac{\sum_{n: z_n = k} x_n + \lambda \mu_k^*}{N_k + \lambda}$$

so μ_k^* is like an additional data point with weight λ .

Name: Roll No.: Dept.:

4. Consider learning a linear regression model by minimizing the squared loss function $\sum_{n=1}^N (y_n - \mathbf{w}^T \mathbf{x}_n)^2$. 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: $\sum_{n=1}^N (y_n - \mathbf{w}^T \tilde{\mathbf{x}}_n)^2$. 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).

This is basically what the DROPOUT technique of neural net does! REGULARIZATION (in each layer)

Let's denote $\tilde{L}(\mathbf{w}) = \sum_{n=1}^N (y_n - \mathbf{w}^T \tilde{\mathbf{x}}_n)^2$. Its expectation will be

$$\mathbb{E}[\tilde{L}(\mathbf{w})] = \sum_{n=1}^N \mathbb{E}[y_n^2 + (\mathbf{w}^T \tilde{\mathbf{x}}_n)^2 - 2y_n \mathbf{w}^T \tilde{\mathbf{x}}_n] = \sum_{n=1}^N y_n^2 + \mathbb{E}[(\mathbf{w}^T \tilde{\mathbf{x}}_n)^2] - 2y_n \mathbb{E}[\mathbf{w}^T \tilde{\mathbf{x}}_n]$$

Note that $\mathbb{E}[(\mathbf{w}^T \tilde{\mathbf{x}}_n)^2] = \text{Var}[\mathbf{w}^T \tilde{\mathbf{x}}_n] + \mathbb{E}[\mathbf{w}^T \tilde{\mathbf{x}}_n]^2$. Plugging it in expression above,

$$\mathbb{E}[\tilde{L}(\mathbf{w})] = \sum_{n=1}^N y_n^2 - 2y_n \mathbb{E}[\mathbf{w}^T \tilde{\mathbf{x}}_n] + \mathbb{E}[(\mathbf{w}^T \tilde{\mathbf{x}}_n)^2] + \text{Var}[\mathbf{w}^T \tilde{\mathbf{x}}_n]$$

$$= \sum_{n=1}^N (y_n - \mathbb{E}[\mathbf{w}^T \tilde{\mathbf{x}}_n])^2 + \text{Var}[\mathbf{w}^T \tilde{\mathbf{x}}_n] \leftarrow \text{We're pretty much at the answer now!}$$

Note that $\mathbb{E}[\mathbf{w}^T \tilde{\mathbf{x}}_n] = \mathbf{w}^T \mathbb{E}[\tilde{\mathbf{x}}_n] = p \mathbf{w}^T \mathbf{x}_n$ (Since $\tilde{\mathbf{x}}_n = \mathbf{x}_n \circ \mathbf{m}_n$ and $\mathbb{E}[m_{nd}] = p$)

Also, $\text{Var}[\mathbf{w}^T \tilde{\mathbf{x}}_n] = \mathbf{w}^T \text{Cov}(\tilde{\mathbf{x}}_n) \mathbf{w} = \mathbf{w}^T \begin{bmatrix} p(1-p)x_{n1}^2 & 0 & \dots & 0 \\ 0 & \ddots & & \\ 0 & & p(1-p)x_{nD}^2 & \\ 0 & & & 0 \end{bmatrix} \mathbf{w} = \mathbf{w}^T \Sigma_n \mathbf{w}$ (Like an L2 regularizer summed over all n)

Thus $\mathbb{E}[\tilde{L}(\mathbf{w})] = \sum_{n=1}^N (y_n - p \mathbf{w}^T \mathbf{x}_n)^2 + \mathbf{w}^T \Sigma \mathbf{w}$ where Σ is the above diagonal matrix

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

Note that $\mathbf{X} \mathbf{X}^T = \mathbf{U} \mathbf{\Lambda} \mathbf{V}^T (\mathbf{U} \mathbf{\Lambda} \mathbf{V}^T)^T = \mathbf{U} \mathbf{\Lambda} \mathbf{V}^T \mathbf{V} \mathbf{\Lambda} \mathbf{U}^T = \mathbf{U} \mathbf{\Lambda}^2 \mathbf{U}^T$ (using $\mathbf{V}^T \mathbf{V} = \mathbf{I}$)

$$(\mathbf{X} \mathbf{X}^T) \mathbf{u}_n = (\mathbf{U} \mathbf{\Lambda}^2 \mathbf{U}^T) \mathbf{u}_n = \lambda_n^2 \mathbf{u}_n \quad (\text{using orthonormality of } \mathbf{U} = [\mathbf{u}_1, \mathbf{u}_2, \dots, \mathbf{u}_N])$$

$\Rightarrow \mathbf{u}_n$ is an eigenvector of $\mathbf{X} \mathbf{X}^T$

Likewise $\mathbf{X}^T \mathbf{X} = (\mathbf{U} \mathbf{\Lambda} \mathbf{V}^T)^T (\mathbf{U} \mathbf{\Lambda} \mathbf{V}^T) = \mathbf{V} \mathbf{\Lambda} \mathbf{U}^T \mathbf{U} \mathbf{\Lambda} \mathbf{V}^T = \mathbf{V} \mathbf{\Lambda}^2 \mathbf{V}^T$ (using $\mathbf{U}^T \mathbf{U} = \mathbf{I}$)

$$(\mathbf{X}^T \mathbf{X}) \mathbf{v}_d = (\mathbf{V} \mathbf{\Lambda}^2 \mathbf{V}^T) \mathbf{v}_d = \lambda_d^2 \mathbf{v}_d$$

$\Rightarrow \mathbf{v}_d$ is an eigenvector of $\mathbf{X}^T \mathbf{X}$.

Name:

IIT Kanpur

CS771 Intro to ML

Roll No.: Dept.:

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 $\{(\mathbf{x}_n, \mathbf{y}_n)\}_{n=1}^N$, with each $\mathbf{x}_n \in \mathbb{R}^D$ and $\mathbf{y}_n \in \mathbb{R}$. Assume the following generative story for each $(\mathbf{x}_n, \mathbf{y}_n)$: (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 $\mathbf{y}_n \sim \mathcal{N}(\mathbf{w}_{z_n}^T \mathbf{x}_n, \beta^{-1})$.

Your goal is to estimate the parameters $\Theta = \{\pi_k, \mu_k, \Sigma_k, \mathbf{w}_k\}_{k=1}^K$ of this model. Assume β 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 Θ . 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 Θ . Also give a brief sketch of the overall ALT-OPT algorithm.

This is basically a mixture of experts model where \mathbf{x}_n and \mathbf{y}_n are also modeled.

$$\begin{aligned}
 \text{CLL}(\Theta) &= \log \prod_{n=1}^N P(\mathbf{x}_n, \mathbf{y}_n, z_n | \Theta) \\
 &= \log \prod_{n=1}^N \prod_{k=1}^K \left[P(\mathbf{x}_n | z_n=k) P(\mathbf{y}_n | z_n=k, \mathbf{x}_n) P(z_n=k) \right] \quad (\text{ignoring } \Theta \text{ from the notation}) \\
 &= \sum_{n=1}^N \sum_{k=1}^K z_{nk} \left[\log P(\mathbf{x}_n | z_n=k) + \log P(\mathbf{y}_n | z_n=k, \mathbf{x}_n) + \log P(z_n=k) \right] \\
 &\quad \downarrow \qquad \qquad \qquad \downarrow \qquad \qquad \qquad \downarrow \\
 &\quad \mathcal{N}(\mathbf{x}_n | \mu_k, \Sigma_k) \qquad \mathcal{N}(\mathbf{y}_n | \mathbf{w}_k^T \mathbf{x}_n, \beta^{-1}) \qquad \pi_k \\
 &= \sum_{n=1}^N \sum_{k=1}^K z_{nk} \left[-\frac{1}{2} \log |\Sigma_k| - \frac{1}{2} (\mathbf{x}_n - \mu_k)^T \Sigma_k^{-1} (\mathbf{x}_n - \mu_k) - \frac{\beta}{2} (\mathbf{y}_n - \mathbf{w}_k^T \mathbf{x}_n)^2 + \log \pi_k \right]
 \end{aligned}$$

— We need expected CLL. The only expectation we need is $E[z_{nk}]$ which is the same as the posterior probability of $z_{nk}=1$, $P(z_{nk}=1 | \mathbf{x}_n, \mathbf{y}_n, \Theta) \propto P(z_{nk}=1) P(\mathbf{x}_n | z_{nk}=1) P(\mathbf{y}_n | z_{nk}=1, \mathbf{x}_n)$ ← Unlike GMM or standard mixture of experts we have two likelihood terms, one for \mathbf{x}_n , one for \mathbf{y}_n

Thus $E[z_{nk}] \propto \pi_k \mathcal{N}(\mathbf{x}_n | \mu_k, \Sigma_k) \mathcal{N}(\mathbf{y}_n | \mathbf{w}_k^T \mathbf{x}_n, \beta^{-1})$ (which can be easily normalized). Suppose $\gamma_{nk} = E[z_{nk}]$, $N_k = \sum_{n=1}^N \gamma_{nk}$

Given γ_{nk} from the E step, the M step is straightforward. μ_k, Σ_k, π_k updates are just like GMM updates of these, e.g.

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

$$\pi_k = \frac{N_k}{N}, \quad \mu_k = \frac{1}{N_k} \sum_{n=1}^N \gamma_{nk} \mathbf{x}_n, \quad \Sigma_k = \frac{1}{N_k} \sum_{n=1}^N \gamma_{nk} (\mathbf{x}_n - \mu_k)(\mathbf{x}_n - \mu_k)^T$$

Name:

IIT Kanpur

CS771 Intro to ML

Roll No.: Dept.:

End-semester Examination

Date: November 29, 2018

The update of W_k is also like the standard mixture of experts with $\hat{W}_k = \left(\sum_{n=1}^N \gamma_{nk} X_n X_n^T \right)^{-1} \left(\sum_{n=1}^N \gamma_{nk} y_n X_n \right)$

ALT-OPT with $\pi_k = \frac{1}{K}$ will be identical, except that (1) we don't need to estimate π_k

(2) Z_n will be computed as follows

$$\begin{aligned} Z_n &= \underset{k \in \{1, \dots, K\}}{\operatorname{argmin}} \left[-\log p(X_n | Z_n = k) - \log p(y_n | Z_n = k, x_n) \right] \\ &= \underset{k \in \{1, \dots, K\}}{\operatorname{argmin}} \left[\frac{1}{2} \log |\Sigma_k| + \frac{1}{2} (x_n - \mu_k)^T \Sigma_k^{-1} (x_n - \mu_k) + \frac{1}{2} (y_n - W_k^T x_n)^2 \right] \end{aligned}$$

(Note that this is the same as

$$Z_n = \underset{k}{\operatorname{argmax}} p(Z_n = k | x_n, y_n, \theta)$$

$$= \underset{k}{\operatorname{argmax}} p(x_n, y_n | Z_n = k, \theta) \underbrace{p(Z_n = k | \theta)}_{\frac{1}{K}}$$

$$= \underset{k}{\operatorname{argmin}} -\log p(x_n, y_n | Z_n = k, \theta)$$

We can use the optimal value of Z_n to create a one-hot vector $\gamma_n = [\gamma_{n1}, \dots, \gamma_{nK}]$

(3) Updates of μ_k, Σ_k, W_k will have an identical form as in the EM case but since γ_n is a one-hot vector, only the points with $\gamma_{nk} = 1$ will contribute to the update of μ_k, Σ_k, W_k .

(Recall HW 3 problem on mixture of experts)

Skipping the ALT-OPT & EM sketch (obvious!),

Name:

Roll No.:

Dept.:

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 PRACTICE WORK ONLY

Name:

Roll No.:

Dept.:

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

ROUGH WORK ONLY