

Name: _____

Roll No.: _____

Dept.: _____

Instructions:

Total: 80 marks

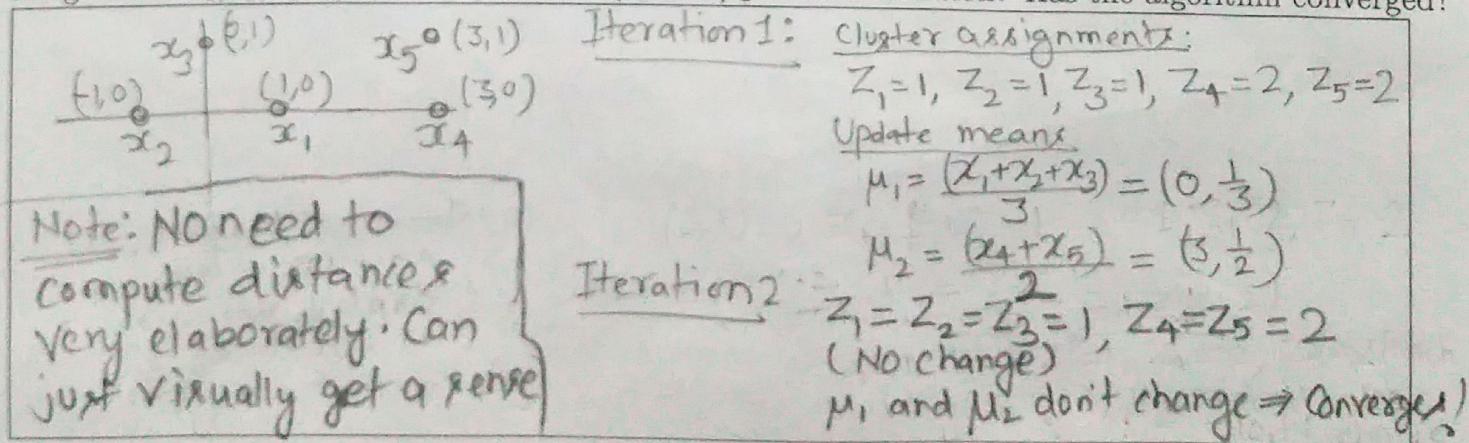
1. This question paper contains a total of 8 pages (8 sides of paper). Please verify.
2. Please write your name, roll number, department on **every side of every sheet** of this booklet.
3. Write final answers **neatly with a pen**. Pencil marks can get smudged and you may lose credit.
4. **Important:** 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. If needed, you may use the personal rough space for more detailed derivations.
5. The last page of the question paper lists some formulae if you need them.

Section 1 (True or False: $10 \times 1 = 10$ marks). For each of the following simply write **T** or **F** in the box.

1. T The prediction cost (time taken to predict the label for a test input) for 1-nearest neighbors classification is higher as compared to that of prototype based classification.
2. T A least squares regression problem with ℓ_1 norm regularizer on the weight vector \mathbf{w} will have a globally optimal solution.
3. F The softmax regression based discriminative model for classification with K classes would require learning $K - 1$ probability distributions.
4. T Changing the Perceptron update rule from $\mathbf{w}^{(t+1)} = \mathbf{w}^{(t)} + y_n \mathbf{x}_n$ to $\mathbf{w}^{(t+1)} = \mathbf{w}^{(t)} + \gamma y_n \mathbf{x}_n$ would effectively learn the same hyperplane separator.
5. F The size of kernel induced feature mapping ϕ of a polynomial kernel with degree $d \geq 2$ is the same for any value of d .
6. F If the MAP objective has a unique optima then the predictive distribution computed using the MAP estimate and computed by averaging over the full posterior will be the same.
7. F For linear/logistic regression, it is not possible to regularize different entries of the weight vector differently if using a zero mean Gaussian prior over the weight vector.
8. T Increasing the parameter C in soft-margin SVM objective $\frac{\|\mathbf{w}\|^2}{2} + C \sum_{n=1}^N \xi_n$ tends to increase the ℓ_2 squared norm of \mathbf{w} .
9. F Running a linear model on landmark based features or kernel random features would always be faster than a linear model on the original features.
10. F The SGD algorithm for a binary linear classification model would update the weight vector only when the current weight vector mispredicts the chosen training example.

Section 2 (8 problems: $8 \times 3 = 24$ marks). Write your answers precisely and concisely in the provided box.

1. Consider a data set with 5 points $\{\mathbf{x}_1, \mathbf{x}_2, \mathbf{x}_3, \mathbf{x}_4, \mathbf{x}_5\}$ in two dimensions: $\{(1, 0), (-1, 0), (0, 1), (3, 0), (3, 1)\}$. Run two iterations of K -means with initial points at $\mu_1 = (-1, 0)$ and $\mu_2 = (3, 1)$. What are the assignments z_1, z_2, z_3, z_4, z_5 and the centers μ_1 and μ_2 at each iteration? Has the algorithm converged?



ame: _____

oll No.: _____

Dept.: _____

2. Rank the following methods in terms of the ~~time~~^{test} time prediction cost: linear SVM, kernelized Perceptron, and kernelized SVM (fastest first, slowest last), and briefly justify your answer.

① Linear SVM : can simply do $\text{sign}(w^T x) \leftarrow$ only ONE dot product

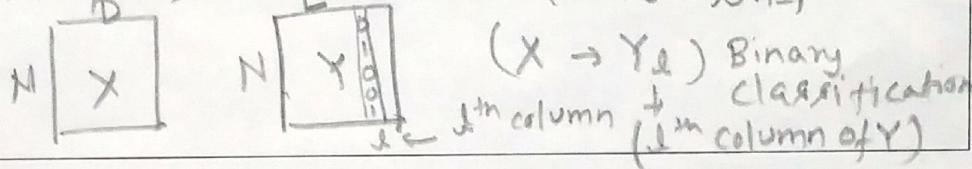
② Kernel SVM : have to do $\text{Sign}(\sum_{n=1}^N \alpha_n K(x_n, x))$ but very few α_n s nonzero. (Support vectors)

③ Kernel Perceptron: — Do — but α 's are not necessarily sparse as in SVM, so may need more dot products.

3. Consider multi-label classification given training data $\{(x_n, y_n)\}_{n=1}^N$. Here each output $y_n \in \{0, 1\}^L$, i.e., a binary vector of length L . Briefly describe an approach to learn a multi-label classification model using this data if you only have access to a *binary classification* algorithm B (e.g., a binary SVM).

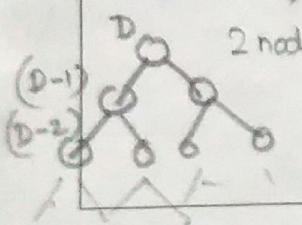
Learn L binary classifiers, one for predicting each of the L labels.

For j^{th} binary classifier, the data will be $\{(x_n, y_{nj})\}_{n=1}^N \xrightarrow{\text{Binary Classification}} 0/1$



4. Consider learning a decision tree, given some training data where each input has D binary-valued features. Let's assume that we will not test any feature that has been tested at one of the previous levels (but we can possibly test a feature at multiple nodes at the same level). How many information gain calculations would be needed to construct the full decision tree (i.e., assuming no pruning)? Just give the basic expression; no need to try simplifying it too much to get a more compact expression.

$$D + 2(D-1) + 2^2(D-2) + \dots + (2^{D-1} \times 1) \quad \text{Redundant}$$



↓ ↓
2 nodes, $D-1$ features to test ↓
4 nodes $D-2$ features to test

actually no test required
since only one feature left now

5. Consider a binary classification dataset with two-dimensional inputs, where each input is of the form $x = (x_1, x_2)$. Suppose the decision boundary is given by $\frac{(x_1-1)^2}{2} + \frac{(x_2-2)^2}{3} = 1$. Note that this is a nonlinear boundary (equation of an ellipse). Write down a mapping $\phi(x)$ that will make it possible to separate the two classes using a linear separator. Also write down the weights of this linear model.

$$\frac{(x_1-1)^2}{2} + \frac{(x_2-2)^2}{3} = 1 \Rightarrow 3(x_1-1)^2 + 2(x_2-2)^2 - 6 = 0$$

$$3(x_1^2 + 1 - 2x_1) + 2(x_2^2 + 4 - 4x_2) - 6 = 0$$

$$3x_1^2 - 6x_1 + 2x_2^2 - 8x_2 + 3 + 8 - 6 = 0$$

$$\phi(x) = (x_1^2, x_1, x_2^2, x_2)$$

$$w^T \phi(x) + b = 0$$

$$w = (3, -6, 2, -8), b = 5$$

Name: _____

Roll No.: _____

Dept.: _____

6. Consider solving multi-class classification by reducing it to binary SVM classification problems using One-vs-All (OvA) and All-Pairs (note: all-pairs is also called "All vs All" or AvA). Suppose we have K classes and M examples per class ($N = KM$ examples total). Typical binary SVM takes roughly N^2 time to train on N examples. For what values of K is it computationally cheaper to use OVA instead of AVA?

OVA: Learns K classifiers, each binary. Each uses all the examples

$$\text{Cost} = K \times (KM)^2 = K^3 M^2$$

AVA: Learns $\frac{K(K-1)}{2}$ classifiers, each binary. Each uses examples from 2 classes

$$\text{Cost} = \frac{K(K-1)}{2} \times (2M)^2 = 2K(K-1)M^2$$

$$\text{We want } K^3 M^2 \leq 2K(K-1)M^2 \Rightarrow K^2 \leq 2(K-1) \Rightarrow \text{No value of } K \geq 2 \text{ satisfying,}$$

so OVA is never cheaper.

7. Consider the hard-margin versions of the SVDD problem (left) and the one-class SVM problem (right).

$$\min_{c \in \mathbb{R}^D, R \in \mathbb{R}} R^2$$

$$\min_{c \in \mathbb{R}^D, \tau \in \mathbb{R}} \frac{\|c\|^2}{2} - \tau$$

$$\text{subject to } \|x_n - c\|^2 \leq R^2$$

$$\text{subject to } c^T x_n \geq \tau + \frac{\|x_n\|^2}{2}$$

Show that these two problems are equivalent for some value of τ . What is that value τ ?

$$\|x_n - c\|^2 \leq R^2 \quad (\text{SVDD constraint})$$

$$\|x_n\|^2 + \|c\|^2 - 2c^T x_n \leq R^2$$

$$\Rightarrow c^T x_n \geq \frac{1}{2}(\|c\|^2 - R^2) + \frac{\|x_n\|^2}{2}$$

$$\text{Comparing with OSVM constraint: } \tau = \frac{1}{2}(\|c\|^2 - R^2)$$

$$\text{plugging in this } \tau \text{ in } \min_{c, \tau} \frac{\|c\|^2}{2} - \tau \Rightarrow \min_{c, R} \frac{R^2}{2} = \min_{c, R} R^2$$

equivalent

8. Consider a least squares regression problem with N examples $(\mathbf{X}, \mathbf{y}) = \{(x_n, y_n)\}_{n=1}^N$ and regression weight vector $w \in \mathbb{R}^D$. Assume no regularization on w . However, suppose we add another M "fake" examples $(\tilde{\mathbf{X}}, \tilde{\mathbf{y}}) = \{(\tilde{x}_m, \tilde{y}_m)\}_{m=1}^M$. Show that using these additional fake examples is equivalent to using an ℓ_2 regularizer. For what values of $(\tilde{\mathbf{X}}, \tilde{\mathbf{y}})$, it will correspond to an ℓ_2 regularizer $\lambda w^\top w$?

$$L(w) = \sum_{n=1}^N (y_n - w^\top x_n)^2 + \sum_{m=1}^M (\tilde{y}_m - w^\top \tilde{x}_m)^2$$

using the least square solution, it is easy to see that

$$\hat{w} = (\tilde{\mathbf{X}}^\top \tilde{\mathbf{X}})^{-1} (\tilde{\mathbf{X}}^\top \tilde{\mathbf{y}})$$

comparing it with $\hat{w} = (\mathbf{X}^\top \mathbf{X} + \lambda \mathbf{I}_D)^{-1} \mathbf{X}^\top \mathbf{y}$, we see that

$$\tilde{\mathbf{X}} = \sqrt{\lambda} \mathbf{I}_D \quad \text{and} \quad \tilde{\mathbf{y}} = \mathbf{0} \text{ (vector)}$$

+ $M=D$ in this case

[Note: Any $\tilde{\mathbf{X}}$ that satisfies $\tilde{\mathbf{X}}^\top \tilde{\mathbf{X}} = \lambda \mathbf{I}_D$ + $\tilde{\mathbf{y}} = \mathbf{0}$ vector would correspond to an ℓ_2 reg.]

Name: _____

Roll No.: _____

Dept.: _____

Section 3 (6 problems: $6 \times 6 = 36$ marks). Write your answers precisely and concisely in the provided box.

1. Assuming hard-margin SVM, show that, given the solution for the dual variables α_n 's, the bias term $b \in \mathbb{R}$ can be computed as $b = y_s - t_s$ where s can denote the index of *any* of the support vectors, and t_s is a term that requires computing a summation defined over all the support vectors. (Hint: Use KKT conditions)

Using KKT conditions: $\partial_s [1 - y_s(w^T x_s + b)] = 0$

Since $\alpha_s \neq 0$ for support vectors, we must have

$$y_s (w^T x_s + b) = 1$$

Using $w = \sum_{n \in SV} \alpha_n y_n x_n$, we get

$$y_s \left(\sum_{n \in SV} \alpha_n y_n x_n \right)^T x_s + y_s b = 1$$

Multiplying both sides by y_s , we get (using $y_s^2 = 1$)

$$b = y_s - \underbrace{\left(\sum_{n \in SV} \alpha_n y_n x_n \right)^T x_s}_{t_s}$$

2. Show that we can rewrite regression with absolute loss function $|y_n - w^T x_n|$ as a reweighted least squares objective where the squared loss term for each example (x_n, y_n) is multiplied by an importance weight $s_n > 0$. Write down the expression for s_n , and briefly explain why this expression for s_n makes intuitive sense. Given N examples $\{(x_n, y_n)\}_{n=1}^N$, briefly outline the steps of an optimization algorithm that estimates the unknowns (w and the importance weights $\{s_n\}_{n=1}^N$) for this reweighted least squares problem.

$$|y_n - w^T x_n| = \frac{1}{|y_n - w^T x_n|} (y_n - w^T x_n)^2 = s_n (y_n - w^T x_n)^2$$

Thus $s_n = \frac{1}{|y_n - w^T x_n|}$. It makes sense because it reduces

the importance of "outliers" (which are examples that deviate too much from the model's fit; note that large $|y_n - w^T x_n|$ would imply small importance).

$$\text{Thus our objective: } L(w, \{s_n y_n^H\}_{n=1}^N) = \sum_{n=1}^N s_n (y_n - w^T x_n)^2$$

We can estimate w and $\{s_n\}_{n=1}^N$ in alternating fashion.

- Initialize $\{s_n\}_{n=1}^N = S$ (NN vector)

- Solve for w given S (a simple reweighted least squares)

$$w = (X^T \text{diag}(S) X)^{-1} X^T \text{diag}(S) y \leftarrow \text{It is okay}$$

- Given w , set $s_n = \frac{1}{|y_n - w^T x_n|} \forall s$

- Repeat until convergence

even if you didn't give this equation!

Name: _____

Roll No.: _____

Dept.: _____

3. Consider the following way to generate a binary random number y : Draw K count-valued random variables $m_k \sim \text{Poisson}(\lambda_k)$, $k = 1, \dots, K$. Define $m = \sum_{k=1}^K m_k$, and generate y as $y = \mathbb{I}[m > 0]$, where $\mathbb{I}[\cdot]$ is indicator function. Derive the expression for $p(y=1|\lambda_1, \dots, \lambda_K)$.

$P(y=0|\lambda_1, \dots, \lambda_K)$ is easier. Note that $y=0$ if $m=0$, i.e.

$$m_1 = m_2 = \dots = m_K = 0$$

$$\text{Thus } P(y=0|\lambda_1, \dots, \lambda_K) = \prod_{k=1}^K P(m_k=0|\lambda_k) = \prod_{k=1}^K \frac{\lambda_k^0 e^{-\lambda_k}}{0!} = e^{-\sum_{k=1}^K \lambda_k}$$

$$\text{Thus } P(y=1|\lambda_1, \dots, \lambda_K) = 1 - e^{-\sum_{k=1}^K \lambda_k}$$

4. Suppose we wish to do an “online” K -means by performing an SGD-style optimization on the K -means objective $\sum_{n=1}^N \|\mathbf{x}_n - \mu_{z_n}\|^2 = \sum_{n=1}^N \sum_{k=1}^K z_{nk} \|\mathbf{x}_n - \mu_k\|^2$. Assume mini-batch size = 1, i.e., you get one randomly chosen example \mathbf{x}_n . Assume learning rate η . What will be the SGD update equations for the cluster means μ_1, \dots, μ_K ? In which direction does each mean move as a result of these SGD updates?

When doing SGD, we will first assign \mathbf{x}_n to its closest mean:

$$z_n = \underset{k \in \{1, \dots, K\}}{\operatorname{argmin}} \|\mathbf{x}_n - \mu_k\|^2$$

Note that if $z_n=k$, it means $z_{nk}=1$ and $z_{nl}=0 \forall l \neq k$ (in one-hot notation)

Thus when doing SGD, the loss function for the n^{th} example will be $\|\mathbf{x}_n - \mu_k\|^2$, and its gradient will be $-2(\mathbf{x}_n - \mu_k)$. Thus, μ_k will be updated as

$$\mu_k^{(t+1)} = \mu_k^{(t)} + 2\eta (\mathbf{x}_n - \mu_k^{(t)})$$

$$= (1-2\eta)\mu_k^{(t)} + 2\eta \mathbf{x}_n \leftarrow \text{thus } \mu_k \text{ moves toward } \mathbf{x}_n$$

None of the other means will be updated, so they won't move.

ame:

oll No.:

Dept.:

5. Consider N scalar-valued observations x_1, \dots, x_N from a Gaussian $\mathcal{N}(\mu, \lambda^{-1})$. Suppose the mean μ is known and precision λ is unknown. Assume a gamma prior on λ , i.e., $p(\lambda) = \text{Gamma}(\lambda|a, b) = \frac{1}{\Gamma(a)} b^a \lambda^{a-1} \exp(-b\lambda)$. Compute the posterior distribution of λ , i.e., $p(\lambda|x_1, \dots, x_N)$. Is the posterior available in closed form? If yes, why, and what's the name of this distribution? If no, why not? (Hint/suggestion: Try computing the posterior first, before answering yes or no).

$$\begin{aligned} P(\lambda|x_1, \dots, x_N) &\propto P(\lambda)P(x_1, \dots, x_N|\mu, \lambda) \\ &= P(\lambda) \prod_{n=1}^N P(x_n|\mu, \lambda) \\ &\propto \lambda^{a-1} \exp(-b\lambda) \times \prod_{n=1}^N \sqrt{\frac{1}{2\pi}} \exp\left[-\frac{\lambda}{2}(x_n - \mu)^2\right] \\ &\propto \lambda^{a+\frac{N}{2}-1} \exp\left[-(b + \frac{1}{2} \sum_{n=1}^N (x_n - \mu)^2)\lambda\right] \end{aligned}$$

Note that the above distribution is again Gamma, with shape param $a + \frac{N}{2}$ and rate param $b + \frac{1}{2} \sum_{n=1}^N (x_n - \mu)^2$

Thus we do get a closed form posterior

This is because Gaussian and Gamma are also conjugate if Gaussian's mean is fixed and its precision has a gamma prior.

6. Consider linear regression with squared loss $\mathcal{L}(\mathbf{w}) = \sum_{n=1}^N (y_n - \mathbf{w}^\top \mathbf{x}_n)^2$ (no regularizer), and apply Newton's method to find the optimal \mathbf{w} . Write down the expression for the weight update at each iteration. What is the minimum number of iterations that Newton's method will take to converge? Is that what you would expect Newton's method to do for this problem? If yes, why? If no, why not?

$$\text{Gradient } g = -2 \sum_{n=1}^N (y_n - \mathbf{w}^\top \mathbf{x}_n) \mathbf{x}_n$$

$$\text{Hessian } H = \frac{\partial g}{\partial \mathbf{w}} = -2 \sum_{n=1}^N \frac{\partial}{\partial \mathbf{w}} (y_n - \mathbf{w}^\top \mathbf{x}_n) \mathbf{x}_n^\top = 2 \sum_{n=1}^N \mathbf{x}_n \mathbf{x}_n^\top$$

$$\text{Thus } \mathbf{w}^{(t+1)} = \mathbf{w}^{(t)} - H^{(t)} g^{(t)} = \mathbf{w}^{(t)} + \frac{1}{2} \left(\sum_{n=1}^N \mathbf{x}_n \mathbf{x}_n^\top \right)^{-1} \sum_{n=1}^N (y_n - \mathbf{w}^{(t)} \mathbf{x}_n) \mathbf{x}_n$$

$$\Rightarrow \mathbf{w}^{(t+1)} = \mathbf{w}^{(t)} + \left(\sum_{n=1}^N \mathbf{x}_n \mathbf{x}_n^\top \right)^{-1} \sum_{n=1}^N \mathbf{x}_n y_n - \left(\sum_{n=1}^N \mathbf{x}_n \mathbf{x}_n^\top \right)^{-1} \sum_{n=1}^N (\mathbf{w}^{(t)} \mathbf{x}_n) \mathbf{x}_n$$

$$\Rightarrow \mathbf{w}^{(t+1)} = \mathbf{w}^{(t)} + (\mathbf{x}^\top \mathbf{x})^{-1} \mathbf{x}^\top \mathbf{y} - \mathbf{w}^{(t)} = (\mathbf{x}^\top \mathbf{x})^{-1} \mathbf{x}^\top \mathbf{y}$$

Thus Newton's update converges in ONE iteration.

This is NOT surprising. Squared error is quadratic function, and the quadratic approx. made by Newton's method is the same as the function being optimized.

inverse
cancels out

$\left(\sum_{n=1}^N \mathbf{x}_n \mathbf{x}_n^\top \right) \mathbf{w}^{(t)}$

Section 4 (1 problem: 10 marks). Write your answers precisely and concisely in the provided box.

1. The mistake-driven Perceptron update rule, after Perceptron makes a mistake on (\mathbf{x}_n, y_n) , updates the weight vector as $\mathbf{w}^{(t)} = \mathbf{w}^{(t-1)} + y_n \mathbf{x}_n$. This update will modify the weight vector (hyperplane separator) in such a way that it becomes *less incorrect* on the example (\mathbf{x}_n, y_n) , but not necessarily *correct*. For example, if $y_n = 1$ then $\mathbf{w}^{(t)} \mathbf{x}_n$ will become less negative than $\mathbf{w}^{(t-1)} \mathbf{x}_n$ (but not necessarily positive). Likewise, if $y_n = -1$ then $\mathbf{w}^{(t)} \mathbf{x}_n$ will become less positive than $\mathbf{w}^{(t-1)} \mathbf{x}_n$ (but not necessarily negative). Let's try to design a variant of Perceptron that guarantees that, after the update, the new weight vector $\mathbf{w}^{(t)}$ will ~~be~~ definitely be correct on the example (\mathbf{x}_n, y_n) , i.e., $y_n \mathbf{w}^{(t)} \mathbf{x}_n \geq 0$. At the same time, let's not make the new weight vector $\mathbf{w}^{(t)}$ drift too far from the current weight $\mathbf{w}^{(t-1)}$, i.e., we want to keep the squared ℓ_2 distance $\|\mathbf{w}^{(t)} - \mathbf{w}^{(t-1)}\|^2$ as small as possible.

Formulate the above as an optimization problem and solve it to derive the updates for this variant of the Perceptron. Also verify that your obtained expression for $\mathbf{w}^{(t)}$ does satisfy the constraint $y_n \mathbf{w}^{(t)} \mathbf{x}_n \geq 0$.

$$\begin{aligned} & \underset{\mathbf{w}^{(t)}}{\operatorname{argmin}} \frac{\|\mathbf{w}^{(t)} - \mathbf{w}^{(t-1)}\|^2}{2} \quad \text{s.t. } -y_n \mathbf{w}^{(t)} \mathbf{x}_n \leq 0 \\ & \text{Lagrangian } L(\mathbf{w}^{(t)}, \alpha) = \frac{\|\mathbf{w}^{(t)} - \mathbf{w}^{(t-1)}\|^2}{2} - \alpha_n y_n \mathbf{w}^{(t)} \mathbf{x}_n \end{aligned}$$

Min. w.r.t. $\mathbf{w}^{(t)}$ we get

$$\begin{aligned} (\mathbf{w}^{(t)} - \mathbf{w}^{(t-1)}) - \alpha_n y_n \mathbf{x}_n &= 0 \\ \Rightarrow \mathbf{w}^{(t)} &= \mathbf{w}^{(t-1)} + \alpha_n y_n \mathbf{x}_n \end{aligned}$$

Now what's the optimal α_n ? To get that we plug-in $\mathbf{w}^{(t)}$ into the Lagrangian and then max. over α_n .

$$\begin{aligned} L(\alpha_n) &= \frac{\|\alpha_n y_n \mathbf{x}_n\|^2}{2} - \alpha_n y_n (\mathbf{w}^{(t-1)} + \alpha_n y_n \mathbf{x}_n) \mathbf{x}_n \\ &= \alpha_n^2 \frac{\|\mathbf{x}_n\|^2}{2} - \alpha_n y_n \mathbf{w}^{(t-1)} \mathbf{x}_n - \alpha_n^2 \|y_n \mathbf{x}_n\|^2 \quad (\text{note: } y_n^2 = 1) \end{aligned}$$

Taking deriv. w.r.t α_n and setting to zero gives:

$$\begin{aligned} \alpha_n \|\mathbf{x}_n\|^2 - 2 \alpha_n \|y_n \mathbf{x}_n\|^2 - y_n \mathbf{w}^{(t-1)} \mathbf{x}_n &= 0 \\ \Rightarrow \alpha_n &= -\frac{y_n \mathbf{w}^{(t-1)} \mathbf{x}_n}{\|\mathbf{x}_n\|^2} \end{aligned}$$

$$\text{Thus } \mathbf{w}^{(t)} = \mathbf{w}^{(t-1)} - \frac{y_n^2 \mathbf{w}^{(t-1)} \mathbf{x}_n}{\|\mathbf{x}_n\|^2} \mathbf{x}_n = \mathbf{w}^{(t-1)} - \frac{(\mathbf{w}^{(t-1)} \mathbf{x}_n) \mathbf{x}_n}{\|\mathbf{x}_n\|^2}$$

Also, note that $\mathbf{w}^{(t)} \mathbf{x}_n = 0 \Rightarrow y_n \mathbf{w}^{(t)} \mathbf{x}_n = 0$ so the constraint does get satisfied.

Name:

IIT Kanpur

CS771 Intro to ML

Mid-semester Examination

Date: September 20, 2018

Roll No.: Dept.:

Some formulae you might need

- Gaussian PDF: $\mathcal{N}(x|\mu, \lambda^{-1}) = \sqrt{\frac{\lambda}{2\pi}} \exp(-\frac{\lambda}{2}(x - \mu)^2)$
- Poisson PMF: $\text{Poisson}(x|\lambda) = \frac{\lambda^x \exp(-\lambda)}{x!}$
- Hessian of a scalar-valued function $\mathcal{L}(\mathbf{w})$ w.r.t. $\mathbf{w} \in \mathbb{R}^D$ is $\mathbf{H} = \frac{\partial^2 \mathcal{L}(\mathbf{w})}{\partial \mathbf{w} \partial \mathbf{w}^\top} = \frac{\partial}{\partial \mathbf{w}} \left[\frac{\partial \mathcal{L}(\mathbf{w})}{\partial \mathbf{w}} \right]^\top = \frac{\partial \mathbf{g}^\top}{\partial \mathbf{w}}$
- Derivatives - linear form: $\frac{\partial \mathbf{x}^\top \mathbf{a}}{\partial \mathbf{x}} = \frac{\partial \mathbf{a}^\top \mathbf{x}}{\partial \mathbf{x}} = \mathbf{a}$, quadratic form: $\frac{\partial}{\partial \mathbf{x}} (\mathbf{x} - \mathbf{s})^\top \mathbf{W} (\mathbf{x} - \mathbf{s}) = 2\mathbf{W}(\mathbf{x} - \mathbf{s})$