

Major Exam

Total marks: 30 pts

Total time: 1 hr 40 mins

HONOUR CODE

As a student of IIT Delhi, you will not give or receive aid in examinations. You will do your share and take an active part in seeing to it that others as well as yourself uphold the spirit and letter of the Honour Code.

Instructions:

- **Write your name and roll number on answer script**
- With the exception of Question 1, all your answers must be clearly motivated! *A correct answer without a proper motivation will score zero points!*
- Vectors are denoted with a single underline \underline{a} , and matrices by double underline, $\underline{\underline{A}}$. Scalars appear without any underline. Please follow this rule through your answer book.

1. [8 pts] Answer **True** or **False**.

Each correct answer scores 1 point, each incorrect answer scores -1 point and each missing answer scores 0 point.

- T** (a) In ensemble methods, averaging over parameters (instead of predictions) can not be a valid method?
- T** (b) The `flatten` function is used to transform the activations of the convolution layers, which have channel, height, and width dimensions, into flat vectors that are input into the fully connected layers
- F** (c) The parameters of hard-margin support vector regression have a closed form solution
- F** (d) Neural Networks are non-parametric methods
- F** (e) To talk about the bias-variance trade-off only has meaning when learning by minimizing the mean squared error cost function.
- (f) Please select all that apply about k -NN in the following options. Assume a point can be its own neighbor.
- 0.5** ✓ (A) k -NN works great with a small amount of data, but struggles when the amount of data becomes large.
- 0.5** ✓ (B) k -NN is sensitive to outliers; therefore, in general we decrease k to avoid overfitting.
- (C) k -NN can be applied to classification problems but not regression problems.
- 0.5** ✓ (D) We can always achieve zero training error (perfect classification) with k -NN, but it may not generalize well in testing.
- F** (g) Logistic regression and quadratic discriminative analysis (QDA) will always produce the same decision boundary for binary classification problems
- T** (h) k -means clustering algorithm can be considered a special case of Gaussian Mixture Model

-0.5 for incorrect

2. [8 pts] A company is working on a machine learning problem with $p = 5000$ input variables $\underline{x} = \{x_1, \dots, x_{5000}\}^T$ with the aim to predict a single binary output $y \in \{1, -1\}$. They have collected $N = 200$ datapoints with equally many examples from each class.

(a) [3 pts] From previous experience, most of the features are useless for predicting y . To build a model, the company has executed the following steps:

i) Compute the linear correlation between each input and the output and keep only the 100 most correlated input variables (positive or negative correlation).

ii) Learn a logistic regression model.

iii) Evaluate the performance of the model by running 10-fold cross-validation to estimate the accuracy on new, unseen data.

The cross-validation indicated an accuracy of over 95% for the model. However, when they tried it in production, the performance was much worse.

Is the problem too hard or did the company do anything wrong? Give advice as to how the company should act.

(b) [4 pts] The company also tried LASSO. They learned the parameters of a logistic regression model by minimizing the cost function

$$J(\theta) = \sum_{i=1}^{200} \log(1 + e^{-y_i \cdot \theta^T \underline{x}_i}) + \frac{1}{\lambda} \sum_{j=1}^{5000} |\theta_j|$$

To learn the regularization parameter λ , they have produced the following graph.

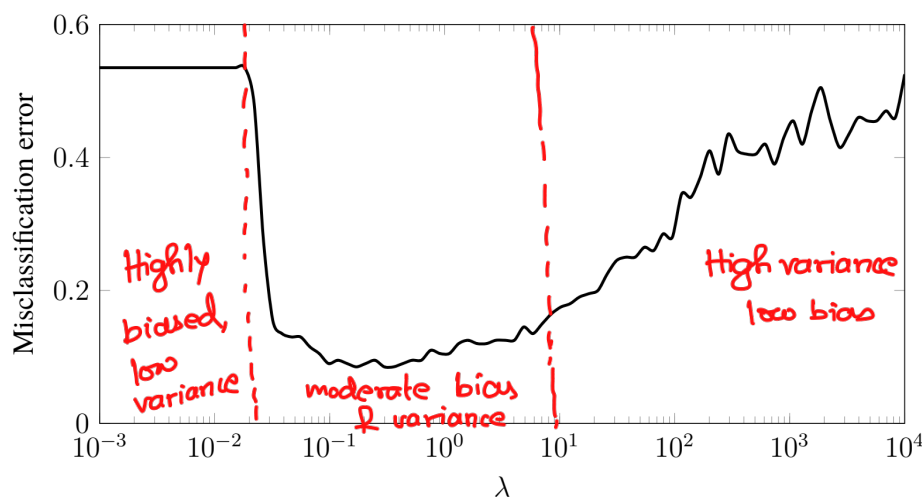


Figure 1: Misclassification error estimated using 10-fold cross-validation for different values of the regularization parameter λ

(2) Explain why the graph looks the way it does. Specifically, explain the shape in terms of the bias-variance trade-off and which term dominates where. Explain also what bias and variance mean in terms of model complexity and how this relates to λ .

(c) [1 pt] The left part of the graph is completely flat. Why?

$\underline{0} \rightarrow \underline{0}$
All parameters are almost equal to zero
So there is nothing learned when λ is very low, so the prediction remains same.

More non-zero parameters imply higher model complexity. Lower λ implies more sparsification of parameters \Rightarrow lower model compl.

① To clear out the gradients of all parameters that optimizer is tracking

3. [2 pts] What is the purpose of zero_grad method in PyTorch training of neural networks? What happens if the method is not called?

① Gradients from previous minibatches would add up leading to incorrect gradient computation

4. [2 pts] In the linearly separable case of SVM if one of the training samples is removed, will the decision boundary shift toward the point removed or shift away from the point removed or remain the same? Justify your answer.

If training sample == support vector, then yes } ②
 " " ≠ support vector, then no

5. [6 pts] In the lecture, we covered the soft-margin SVM. over the training set $\{(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)\}$.

Here, we are interested in a modified version of the soft-margin SVM where we have a custom margin ξ_i for each of the n data points. In the standard soft-margin SVM, we pay a penalty of ξ_i for each of the data point. In practice, we might not want to treat each training point equally – for example, we might know that some data points are more important than the others.

We formally define the following optimization problem:

$$\begin{aligned} \min_{\underline{\theta}, \{\xi_i\}} \quad & \frac{1}{2} \|\underline{\theta}\|_2^2 + C \sum_{i=1}^n \phi_i \xi_i \\ \text{s.t.} \quad & y_i(\underline{\theta}^T \underline{x}_i - b) \geq 1 - \xi_i \quad \forall i \rightarrow \xi_i \geq 1 - \gamma_i (\underline{\theta}^T \underline{x}_i - b) \\ & \xi_i \geq 0 \quad \forall i \rightarrow \xi_i \geq 0 \end{aligned}$$

Since $\sum \xi_i$ is being minimized therefore

Note that the only difference is that we have a custom weighting factor $\phi_i > 0$ for each of the slack variables ξ_i in the objective function. These ϕ_i are some constants given by the prior knowledge, and thus they can be treated as known constants in the optimization problem. Intuitively, this formulation weighs each of the violations (ξ_i) differently according to the prior knowledge (ϕ_i).

- (a) [2 pts] For the standard soft-margin SVM, we have shown that the constrained optimization problem is equal to the following unconstrained optimization problem:

$$\min_{\underline{\theta}, b} \frac{1}{2} \|\underline{\theta}\|_2^2 + C \sum_{i=1}^n \max(0, 1 - y_i(\underline{\theta}^T \underline{x}_i - b))$$

← $\xi_i \geq \max(0, 1 - \gamma_i (\underline{\theta}^T \underline{x}_i - b))$
 Motivation has partial marks

What's the corresponding unconstrained optimization problem for the SVM with custom margins?

- (b) [4 pts] As seen in lecture, the dual form of the standard soft-margin SVM is:

$$\begin{aligned} \min_{\underline{\alpha}} \quad & \frac{1}{2} \underline{\alpha}^T \underline{K} \underline{\alpha} - \underline{\alpha}^T \underline{1} \\ \text{s.t.} \quad & \sum_{i=1}^n \alpha_i y_i = 0 \\ & 0 \leq \alpha_i \leq C \quad i = 1, \dots, n \end{aligned}$$

where $\underline{K} = (\text{diag}(\underline{y})) \underline{X} \underline{X}^T (\text{diag}(\underline{y}))$.

What is the dual form of the SVM with custom margins? To start, you are provided with the Lagrangian, which is given by

Step marking is followed!

$$\mathcal{L}(\theta, b, \xi, \alpha, \gamma) = \frac{1}{2} \|\theta\|_2^2 + \sum_{i=1}^n \alpha_i - \sum_{i=1}^n \alpha_i y_i (\theta^T \underline{x}_i - b) + \sum_{i=1}^n (C\phi_i - \alpha_i - \gamma_i) \xi_i$$

The derivation is very similar to one done in lecture
Only difference will be in $0 \leq \alpha_i \leq \Phi_i C$

6. [4 pts] Recall that in k -means clustering, we attempt to minimize the objective $\rightarrow \mathcal{J}$

$$\min_{r_{ik}, \underline{\mu}_k} \sum_{i=1}^N \sum_{k=1}^K r_{ik} \|\underline{x}_i - \underline{\mu}_k\|_2^2$$

where it turned out that: if the assignment r_{ik} is known for each \underline{x}_i , the mean of the k th cluster is given by

$$\underline{\mu}_k = \frac{1}{\sum_{i=1}^N r_{ik}} \sum_{i=1}^N r_{ik} \underline{x}_i$$

The samples are $\{\underline{x}_1, \dots, \underline{x}_N\}$, where $\underline{x}_i \in \mathbb{R}^d$.

- (a) [1 pt] What is the minimum value of the objective when $K = N$ (the number of clusters equals the number of sample points)? **Each datapoint will be a cluster $\Rightarrow \mathcal{J} = 0$** — ①
- (b) [3 pts] Suppose we add a regularization term to the above objective, as follows:

$$\min_{r_{ik}, \underline{\mu}_k} \sum_{i=1}^N \sum_{k=1}^K r_{ik} \|\underline{x}_i - \underline{\mu}_k\|_2^2 + \lambda \|\underline{\mu}_k\|_2^2$$

Determine the expression for $\underline{\mu}_k$, assuming known r_{ik} from a previous assignment step.

$$\mathcal{J}(\underline{\mu}) = \sum_{i=1}^N \sum_{k=1}^N r_{ik} (\underline{x}_i - \underline{\mu}_k)^T (\underline{x}_i - \underline{\mu}_k) + \lambda \underline{\mu}_k^T \underline{\mu}_k$$

$$\frac{\partial \mathcal{J}}{\partial \underline{\mu}_k} = 0 \Rightarrow \frac{\partial \mathcal{J}}{\partial \underline{\mu}_k} = - \sum_{i=1}^N r_{ik} (\underline{x}_i - \underline{\mu}_k) + \lambda \underline{\mu}_k = 0 \quad \text{--- ①}$$

$$\Rightarrow \left(\sum_{i=1}^N r_{ik} + \lambda \right) \underline{\mu}_k = \sum_{i=1}^N r_{ik} \underline{x}_i$$

$$\Rightarrow \underline{\mu}_k = \frac{\sum_{i=1}^N r_{ik} \underline{x}_i}{\sum_{i=1}^N r_{ik} + \lambda} \quad \text{--- ①}$$