

Roll No: MA22C043

Name: Shubham Singh

Collaborators (if any):

References/sources (if any):

- Use \LaTeX to write-up your solutions (in the solution blocks of the source \LaTeX file of this assignment), submit the resulting rollno.asst2.answers.pdf file at Crowdmark by the due date, and properly drag that pdf's answer pages to the corresponding question in Crowdmark (do this properly, otherwise we won't be able to grade!). (Note: **No late submissions** will be allowed, other than one-day late submission with 10% penalty or four-day late submission with 30% penalty.)
- Please upload to moodle a rollno.zip file containing three files: rollno.asst2.answers.pdf file mentioned above, and two code files for the programming question (rollno.ipynb file and rollno.py file). Do not forget to upload to Crowdmark your results/answers (including Jupyter notebook **with output**) for the programming question.
- Collaboration is encouraged, but all write-ups must be done individually and independently, and mention your collaborator(s) if any. Same rules apply for codes written for any programming assignments (i.e., write your own code; we will run plagiarism checks on codes).
- If you have referred a book or any other online material or LLMs (Large Language Models like ChatGPT) for obtaining a solution, please cite the source. Again don't copy the source *as is* - you may use the source to understand the solution, but write-up the solution in your own words (this also means that you cannot copy-paste the solution from LLMs!). Please be advised that *the lesser your reliance on online materials or LLMs for answering the questions, the more your understanding of the concepts will be and the more prepared you will be for the course exams*.
- Points will be awarded based on how clear, concise and rigorous your solutions are, and how correct your answer is. The weightage of this assignment is 12% towards the overall course grade.

1. (6 points) [A DIRECT/DISCRIMINANT APPROACH TO CLASSIFICATION] For the dataset below, we would like to learn a classifier, specifically a discriminant of the form: $\hat{y} = \text{sign}(wx)$ (assume $\text{sign}(u) = +1$ if $u \geq 0$, and -1 otherwise).

x	y
-1	-1
1	+1
20	+1

Let $z_i := z(x_i) := wx_i$. For a training dataset of size n , the parameter w of the classifier can be learnt by minimizing the

(L1) 0-1 loss function aka misclassification error $\sum_{i=1}^n (1 - \text{sign}(y_i z_i))/2$,

(L2) squared loss function $\sum_{i=1}^n (y_i - z_i)^2$, or

(L3) logistic loss function $\sum_{i=1}^n \log(1 + \exp(-y_i z_i))$.

- (a) (2 points) The 0-1 loss function is the most intuitive choice to build a good classifier. What value of w will lead to such a good classifier for this dataset: $w = 0$ or $w = 1$?

Solution: $z_i := z(x_i) := wx_i$, $w = 0$ gives +1 to every point and $w = 1$ gives -1,+1,+1 to point 1,2,3 respectively, which is more accurate, hence $w = 1$ will lead to a good classifier for this dataset.

- (b) (4 points) Between these values of w ($w = 0$ vs. $w = 1$), determine what value is preferred by the squared and logistic loss functions. Report the actual losses for these w values, and argue which loss function is better.

(Note: Optimizing squared loss is equivalent to applying linear regression methodology to solve this classification problem - did it work fine when there are outliers like $x = 20$ in the dataset?)

Solution:

Point	True Class	$z(x_i) := wx_i$	$z(x_i) := wx_i$	$\text{sign}(y_i z_i)$	$\text{sign}(y_i z_i)$
	y	$w = 0$	$w = 1$	$w = 0$	$w = 1$
-1	-1	0	-1	1	-1
1	1	0	1	1	1
20	1	0	20	1	1

squared loss for $w = 0 := 3$

Logistic loss for $w = 0 := 3\log(2) = 0.9030$

squared loss for $w = 0 := 19^2 = 361$

Logistic loss for $w = 0 := 2\log(2) + \log(1 + \exp(20)) = 9.2879$

$w = 0$, will be preferred by the squared and logistic loss functions. hence we can see that outliers really affects the actual solutions.

2. (6 points) [THINKING LOGISTIC-ALLY...] Consider the scenario in which a user maintains a dataset consisting of songs that he has downloaded over a period of time. He also tracks the likes (-1)/dislikes(+1) for each song along with a set of features X_1, X_2 . X_1 is a binary variable that takes value 1 if the song is sung by his favorite singer and X_2 corresponds to song duration in minutes. This dataset with 10 datapoints is given below:

$$[X_1 \ X_2] = \begin{bmatrix} 1 & 0 & 1 & 0 & 1 & 0 & 0 & 1 & 0 & 0 \\ 5 & 10 & 13 & 2 & 3 & 5 & 2 & 10 & 10 & 3 \end{bmatrix}^T$$

$$y = [-1 \ -1 \ -1 \ -1 \ -1 \ 1 \ 1 \ 1 \ 1 \ 1]^T$$

- (a) (3 points) Train a logistic regression model on this dataset by performing the gradient descent algorithm steps by setting the initial weights to 0. No bias (intercept term) is required and the step size $\eta = 1$. Report the updated weights at the end of two iterations.

Solution:

$$R(w) = \sum_{i=1}^N \log (1 + \exp (-y_i w^T x_i))$$

$$\sigma(y_i w^T x_i) = \frac{1}{1 + \exp (-y_i w^T x_i)}$$

$$\nabla_w R(w) = \sum_{i=1}^N \sigma(-y_i w^T x_i) (-y_i x_i)$$

$$w_{t+1} \leftarrow w_t + \eta \sum_{i=1}^N \sigma(-y_i w_0^T x_i) (y_i x_i)$$
$$\eta = 1$$

for $w_0 = [0, 0]$ $\sigma(-y_i w^T x_i) = 1/2$ for all i ,

$$w_1 = \sum_{i=1}^N \frac{1}{2} (y_i x_i)$$

$$w_1 = [-1, -1.5]$$

similarly for w_2

$$w_2 \leftarrow w_1 + \sum_{i=1}^N \sigma(-y_i w_1^T x_i) (y_i x_i)$$

$$w_2 = [-2.99572655, -4.52249828]$$

- (b) (1 point) What will be the prediction for a new song with features $[0, 20]$ using the trained logistic regression model?

Solution:

$$z = w_2^T x$$

$$z = [-2.99572655, -4.52249828]^T [0, 20]$$

$$z = -90.448$$

$$\text{sign}(z) = -1$$

Logistic regression will predict that, he will dislike the track.

- (c) (2 points) Discuss if logistic regression is a good choice for addressing this specific problem. If not, what are other better options?

Solution:

logistic regression is not a good choice for addressing this specific problem as any datapoint will lie in first quadrant as first component is 0 or 1, and second component is positive as it represents song duration and we can see for any datapoint x_i , $\text{sign}(w_2^T x_i)$ is -1.

Kernel Regression and Support Vector Machines are better options than logistics regression.

3. (12 points) [SVM'S TO THE RESCUE] A Gaussian or Radial Basis Function (RBF) kernel with inverse width $k > 0$ is

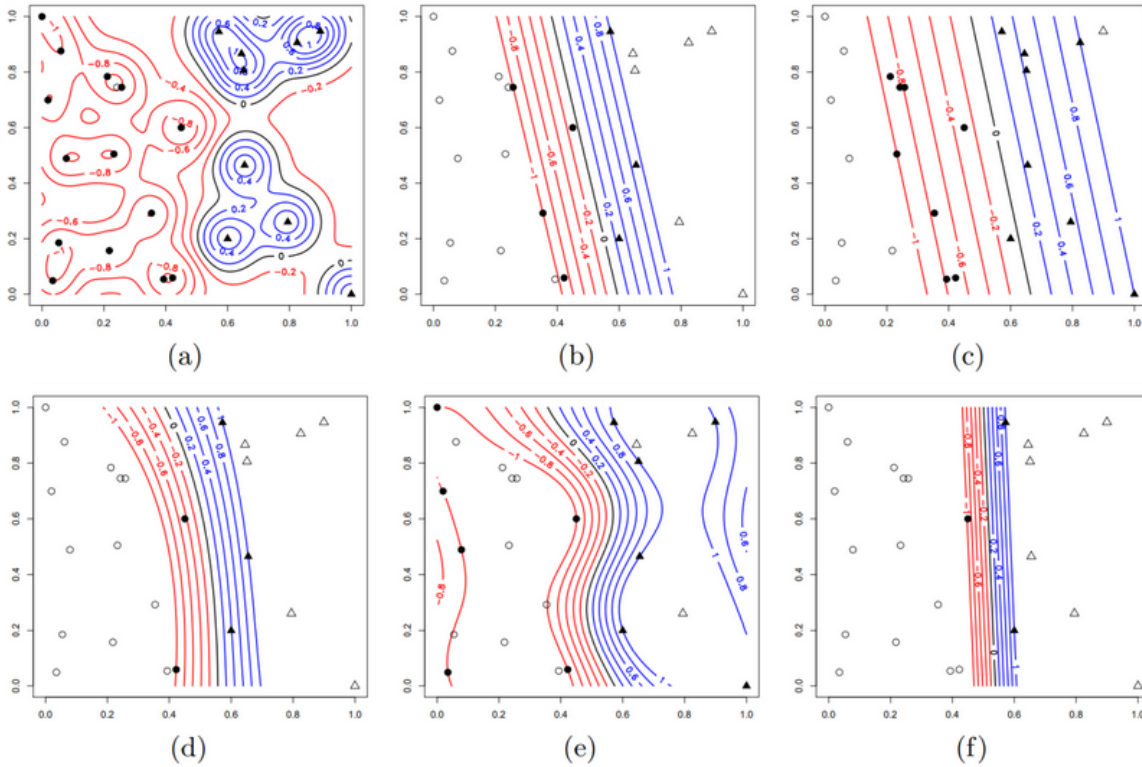
$$K(u, v) = e^{-k||u-v||^2}.$$

Below figures show decision boundaries and margins for SVMs learned on the exact same dataset. Parameters used for the different runs are as follows:

- (i) Linear Kernel with $C = 1$
- (ii) Linear Kernel with $C = 10$
- (iii) Linear Kernel with $C = 0.1$
- (iv) RBF Kernel with $k = 1, C = 3$
- (v) RBF Kernel with $k = 0.1, C = 15$
- (vi) RBF Kernel with $k = 10, C = 1$

Find out which figure plot would have resulted after each run mentioned above. Justify your answer.

In these plots, circles are Class 1, triangles are Class 2, and solid points are support vectors.



Solution:

- (i) figure(b) will corresponds to Linear Kernel with $C = 1$, as it allows some datapoints to be inside margins which means it is not penalising points inside margins too much or too less.
- (ii) figure(f) will corresponds to Linear Kernel with $C = 10$, here margin is small compare to other linear margins, which means it peanlise data too much when they are inside the margins.
- (iii) figure(c) will corresponds to Linear Kernel with $C = 0.1$, as it has the larger margins than other it is safe to say it doesn't penalise data much.
- (iv) figure(d) will corresponds to RBF Kernel with $k = 1$, $C = 3$, as it allows some datapoints to be inside margins which means it is not penalising points inside margins too much or too less.
- (v) figure(e) will corresponds to RBF Kernel with $k = 0.1$, $C = 15$, here no datapoints are inside margins which is equivalent to hard margin SVM problem, which means it peanlise data too much when they are inside the margins.
- (vi) figure(a) will corresponds to RBF Kernel with $k = 10$, $C = 1$, here too many datapoints are inside margins compare to other RBF Kernel which means it doesn't penalise data much.

4. (12 points) [GUESS-TIMATING BIAS AND VARIANCE] You are given a dataset consisting of 100 datapoints in [this folder](#). You have to fit a polynomial ridge regression model to this data.

As seen in class, a model's error can be decomposed into bias, variance, and noise. A "learning curve" provides an opportunity to determine the bias and variance of machine learning models, and to identify models that suffer from high bias (underfitting) or high variance (overfitting). The "learning curve" typically shows the training error and validation/testing error on the y-axis and the model complexity on the x-axis.

- (a) (2 points) Read the last Section (Section 4) on "Bias and Variance in practice" in this [document](#), and summarize briefly how you will heuristically find whether your model suffers from (i) high bias, or (ii) high variance, using only the train and validation/test errors of the model.

Solution: Heuristic Evaluation for Bias and Variance:

High Bias (Underfitting):

- *Training Error:* If the training error is high, it indicates that the model is not fitting the training data well. This could be a sign of a too simplistic model (high bias).
- *Validation/Test Error:* If the validation/test error is also high, it suggests that the model doesn't generalize well to new data.

High Variance (Overfitting):

- *Training Error:* If the training error is low (model fits the training data very well), but...
- *Validation/Test Error:* ... the validation/test error is high, it suggests that the model is too sensitive to the training data and doesn't generalize well to new data.

Balanced Model: A good model typically has low training error and low validation/test error, indicating good generalization to new, unseen data.

- (b) (2 points) Start with the code for polynomial regression from the tutorial (the code without in-built package functions in Tutorial #8), and add quadratic regularization functionality to the code. That is, your code should do polynomial regression with quadratic regularization that takes degree d and regularization parameter λ as input. **Do not use any inbuilt functions from python packages (except for plotting functions and functions to compute polynomial features for each datapoint).**

Solution:

- (c) (3 points) Run your code on the provided dataset for degree $d = 24$ and each λ in the set:

$$\{10^{-15}, 10^{-9}, 10^{-6}, 10^{-3}, 10^{-2}, 10^{-1}, 1, 10^1, 10^2, 10^3, 10^6, 10^9, 10^{15}\}$$

- Perform 5-fold cross-validation on the 100 data points (20 datapoints in each fold). For each validation fold, compute both training (4-folds-based) and validation (1-fold-based) errors using mean squared error measure.
- Calculate the average training and validation errors across the 5 folds.

Solution:

- (d) (3 points) Construct a learning curve by plotting the average training and validation errors against the model complexity ($\log_{10} \lambda$). Based on this learning curve, identify the (i) model with the highest bias, (ii) model with the highest variance?, and (iii) the model that will work best on some unseen data.

Solution:

- (e) (2 points) Plot the fitted curve to the given data (\hat{y} against x curve) for the three models reported in part (e), and superimposed with the training and validation datapoints for any one fold.

Solution:

Please use the template.ipynb file in the [same folder](#) to prepare your solution. Provide your results/answers in the pdf file you upload to Crowdmark named rollno.asst3.answers.pdf, and submit your pdf and code separately also in [this](#) moodle link. The pdf+code submitted should be a rollno.zip file containing three files: rollno.asst3.answers.pdf, rollno.ipynb file (including your code as well as the exact same results/plots uploaded to Crowdmark) and the associated rollno.py file.