# CS-337/335 || Assignment-2

Niraj Mahajan
Roll No- 180050069

## 1    Lasso and ISTA:

1.1   Assume the weights $w_i$ $\forall$ $i \in [0, d)$ be IID from Laplace Distribution, with 0 mean.

$$P(w_i | \sigma) = \frac{1}{2\sigma} \cdot e^{\frac{-|w_i|}{\sigma}}$$

$$\therefore P(w|\sigma) = \prod_i P(w_i|\sigma) = \frac{1}{(2\sigma)^d} e^{\frac{-\sum_{i=0}^{d} |w_i|}{\sigma}}$$

$$\boxed{P(w|\sigma) = \frac{1}{(2\sigma)^d} e^{\frac{-||w||_1}{\sigma}}} \quad \sim ①$$

Now, For $P(D|w)$,

→ we know for any sample $x_i, y_i$

$$P(y_i | x_i, w) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{\left(\frac{-[y_i - w^T \phi(x_i)]^2}{2\sigma^2}\right)}$$

Hence, $\boxed{P(D|w) = \prod_{i=0}^{n} \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(\frac{-[y_i - w^T \phi(x_i)]^2}{2\sigma^2}\right)} \quad \sim ②$

Hence, our Posterior.

$$P\left(w/_D\right) \propto P\left(D/w\right) \cdot P(w)$$

Hence forth, ~~once~~ ~~will~~ we will ignore constants.

$$P\left(w/_D\right) \propto \exp\left(\sum_{i=0}^{n}\left\{\frac{[y_i - w^T\phi(x_i)]^2}{2\sigma^2}\right\}\right) \cdot \exp\left(\frac{-||w||_1}{\sigma}\right)$$

$$P\left(w|_{D,\sigma}\right) \propto \exp\left(-\sum_{i=0}^{n}\left\{\frac{[y_i - w^T\phi(x_i)]^2}{2\sigma^2}\right\} \rightarrow \frac{||w||_1}{\sigma}\right).$$

Taking the Log-likelihood:-

$$LL\left(D|w,\sigma\right) = -\sum_{i=0}^{n}\left\{\frac{[y_i - w^T\phi(x_i)]^2}{2\sigma^2}\right\} - \frac{||w||_1}{\sigma}$$

Hence our MAP estimate $\hat{w} = \underset{w}{\text{argmax}}\ LL\left(D|w,\sigma\right)$

$$= \underset{w}{\text{argmin}}\ -LL\left(D|w,\sigma\right).$$

$$\therefore \hat{w} = \underset{w}{\text{argmin}}\ \sum_{i=0}^{n}\left\{[y_i - w^T\phi(x_i)]^2\right\} + \lambda||w||_1$$

This is exactly what's LASSO Regression.

**Q 1.2]** Plots for part b, part c attached at the end. (Thanks for your cooperation :) )

**b)**

Trends observed in $\frac{train}{test}$ plot :- (Test MSE $\frac{Train\ MSE\ vs\ lambda}{OR}$ vs Lambda

→ Decreases till $\lambda = 0.2$, then increases, then remains constant.

⇒ <u>Explanation for Test MSE</u>.

  ↳ Increase in $\lambda$ for Lasso leads to $w_i$'s becoming '0'. Hence, initially, the Test MSE decreases as the $w_i$'s that were resulting into overfitting are ~~get~~ forced to 0'.

  → Then, the Test MSE increases as after a moment, excess number of $w_i$'s are forced to zero, leading to underfitting.

  → Then, after a particular lambda, all $w_i$'s become '0' leading to constant Test MSE.

⇒ ~~For Train MSE, the initial degree of freedom of the model is so high that the.~~

⇒ For Train MSE, the initial decrease in MSE due to increase in lambda can be explained as follows.

  → Initially, the degree of freedom of the model is very high and the iterations are not sufficient to train the data. As lambda increases, the $w_i$'s get forced to zero and hence [This can be proved by re plotting the graphs with increasing max_iter

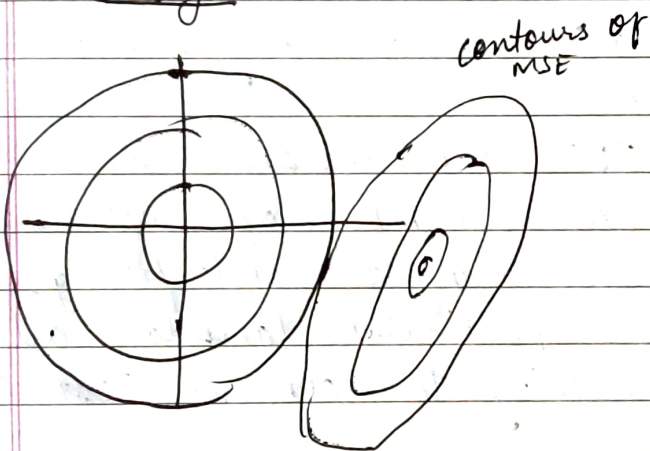⇒ The explanations for the increase and the subsequent plateau of the train MSE is the same as that for Test MSE.

Q-1.5

c) Observation:- Most of the values of $w_i$'s in Lass Regression are `0` while in Ridge Regression, most of the values are non zero (though they are near 0).
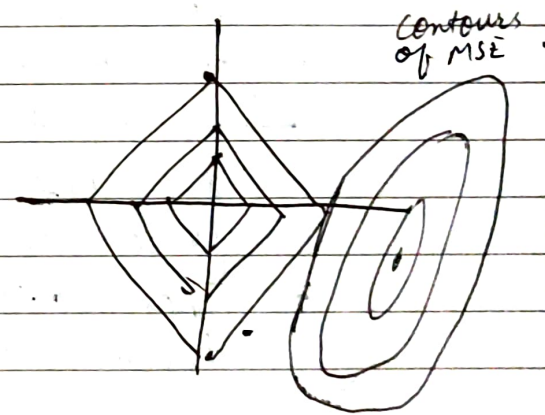
⇒ To explain this, consider Lasso and Ridge $\Omega(w)$ functions in a 2-d space.

Plotting $\Omega(w) = C$.

| Ridge | Lasso |
|---|---|

contours of MSE                    Contours of MSE

As we can see, since the Lass plots are diamonds in the n-D space, it is most likely that the MSE contours first touch the Lasso at the corners, (where $w_i$'s are zero). But, since Ridge Regression contours are hyper spheres, the MSE contours are most likely to touch tangentially at some other point.

→ Hence, the weights in the Lasso solution are mostly 0's but the weights in the Ridge Solution are, although small numbers, non zero.

**Q2.1]** In one vs one perceptron, we have a classifying plane for each pair of classes, while, in case of one vs many perceptron, we have a single classifying plane for each class.

<u>1 vs rest.</u>

<u>Advantages.</u>
↳ Computationally efficient as compared to 1 vs 1 as only $n$ planes are to be trained.

<u>Disadvantages.</u>
↳ Consider the following Distribution:
Here each circle is a cluster of points belonging to class $c_i$
It is clear that a single plane will not be able to seperate $c_5$ from the rest classes.

$C_1$   $C_2$

$C_5$

$C_3$   $C_4$

fig-1.

→ <u>1 vs 1</u>.

<u>Advantages.</u>
↳ If pairwise classes are linearly separable, then this method guarantees good results, irrespective of the distribution of the clusters in the subspace.

<u>Disadvantages.</u>
↳ Computationally costly as $^nC_2$ planes are to be trained.
↳ Even after training, for a input $x$, and class $c$, we have $(n-1)$ scores corresponding to whether $x$ belongs to class $c$ ($\forall c \in c_i$) { $(n-1)$ ~~because one~~ corresponding to each 1 vs 1 planes }

**Q-3.]1]**     Definations:-

Bias :- Error arrising due to erroneous assumption made in model hyperparameters and architecture.

Variance :- ~~test~~ Variability in the test error arising due to small variations in train data.

(Source :- Wikipedia)

**a)** Increase in $\lambda$ will force more weights to 0 in Lasso Regression. Hence, the degree of polynomial will decrease.

Hence, ~~decrease~~ Increase in $\lambda$ will lead to :

→ Increase in Bias.

→ Decrease in Variance.

**b)** Adding higher number of training examples in perceptron :-

→ No change in Bias.

→ Decrease in variance.

Reason :- ~~If~~ we look at the defination of bias, no assumption " or hyperparameter" is changing on adding more test data. Hence, there is no reason for bias to change.

On adding more data, Using Law of Large Numb our estimated solution will start coinciding with the true solution and hence, variability in the test accuracy will decrease as our predic become more and more ideal. Hence, there is a decrease in variance.

c) • Adding more features (Redundant).
→ No change in Bias.
→ No change in Variance.
Assuming :- Data is normalised.
→ On adding redundant dimensions in Lasso,
the algorithm will reduce the weight of the
~~corresponding~~ redundant dimension to 0. Hence the
solution obtained before and after are the
exact same. Hence, there is no change in bias
or variance.


Q.3.2] Plots attached to the end. Thank you for
your co-operation ☺ !


a) → Plot I on test error vs number of samples.
As explained in par 3.1] b, Adding more
training examples will decrease variance but
have no change on $\underset{bias}{\sout{Lasso}}$ ~~Because~~
Reason in [3.1 b].

Hence, we can see that the test error is
decreasing ~~upto~~ ~~some~~ extent.
But after a point, the variance has decreased
greatly and hence the variation in test error
plateaus down.

**3.2](b)** <u>Plot at the end</u> (. — . — —) :)

⇒ Theoritically, the least train error should be obtained on degree = 6, because we are giving (largest possible degree)

✓

our model more and more features to overfit on. (And indeed this is observed in the plots).

⇒ On analysing the plots, we see that the test error decreases till degree = 5, and then increases. Hence degree = 5 is the optimal degree. Here, as we increas degree, we are reducing bias and increasing variance, and at degree = 5 we attain a maxima. After degree = 5, the model starts overfitting.

⇒ I also have added plots of how the data model fits on the train data for various degrees. As we can see, the fit gets better and better as the degree increases.

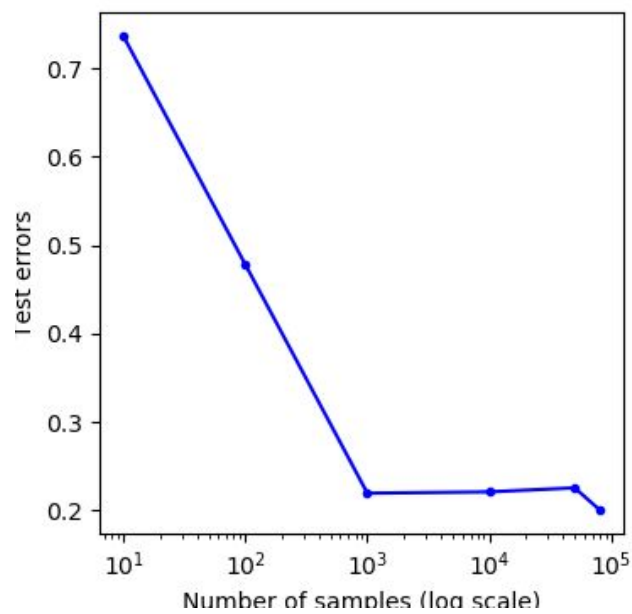# Assignment 2 Graphs: Niraj Mahajan: 180050069

Q1.2 b)



The optimal lambda in the above figure is 0.2. The explanation is provided in the written in the scans.
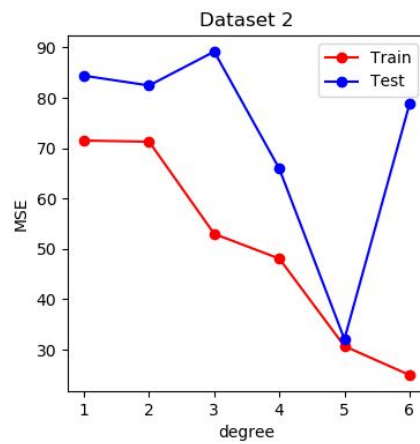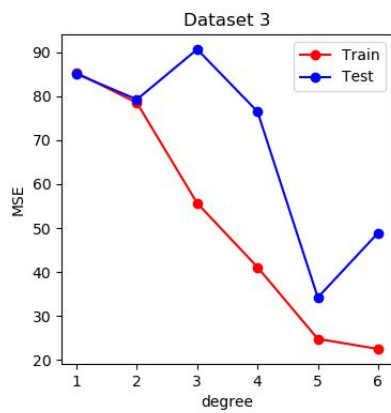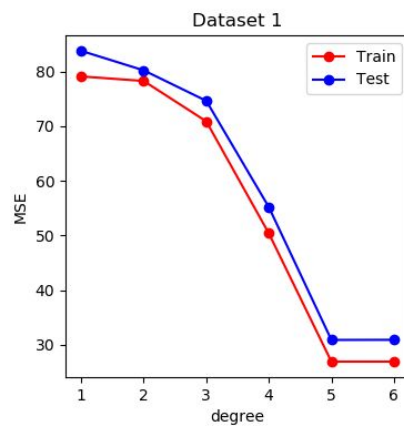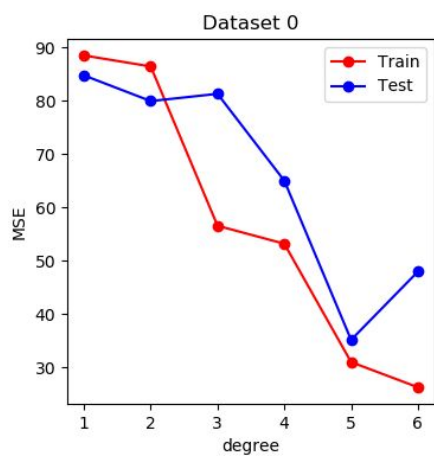
Q1.2 c]

Q 3.2] a]



Q3.2]b]

Plots of fitting model for various degrees. (order is as follows)
1,2
3,4
5,6