

1.1

$$Y = W^T \phi(x) + \epsilon \quad [\text{crossed out}] \quad \epsilon \sim N(0, \sigma^2)$$

$$\mathcal{L}(W) = \underset{W}{\operatorname{argmin}} \left[\|Y - XW\|^2 + \lambda \|W\|_1 \right]$$

Prior on $W \sim \text{Laplace}(0, b)$

$$\rightarrow \text{Laplace}(w_i | 0, b) = \frac{1}{2b} e^{-\frac{|w_i - 0|}{b}}$$

Now

$$P(\theta | D) \propto P(D | \theta) P(\theta)$$

$$\log(P(\theta | D)) \propto \log(P(D | \theta)) + \log(P(\theta))$$

$$\hat{\theta}_{\text{MAP}} = \underset{\theta}{\operatorname{argmax}} \left[\log(P(D | \theta)) + \log(P(\theta)) \right]$$

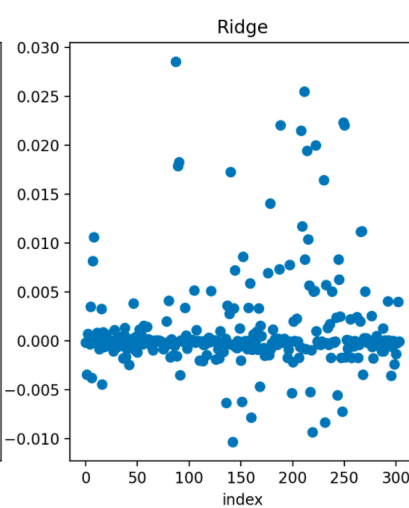
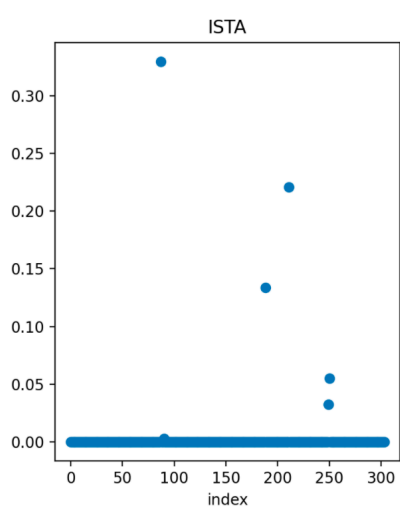
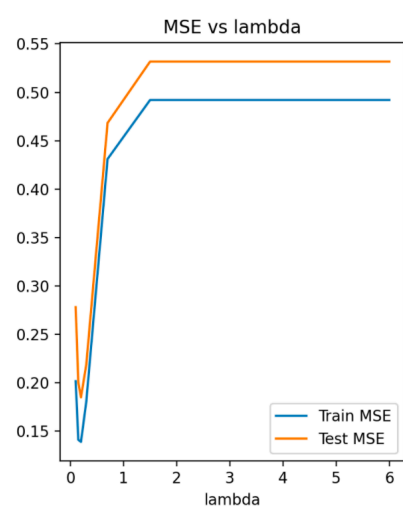
$$\hat{W}_{\text{MAP}} = \underset{W}{\operatorname{argmax}} \left[\log \left(\prod_{i=1}^N \frac{1}{\sigma \sqrt{2\pi}} e^{-\frac{(y_i - W^T x_i)^2}{2\sigma^2}} \right) + \log \left(\prod_{i=1}^M \frac{1}{2b} e^{-\frac{|w_i|}{2b}} \right) \right]$$

$$\hat{W}_{\text{MAP}} = \underset{W}{\operatorname{argmax}} \left[-\sum_{i=1}^N \frac{(y_i - W^T x_i)^2}{2\sigma^2} - \frac{\|W\|_1}{2b} \right]$$

$$\hat{W}_{\text{MAP}} = \underset{W}{\operatorname{argmin}} \left[\sum (y_i - W^T x_i)^2 + \lambda \|W\|_1 \right]$$

$$\lambda = \frac{2\sigma^2}{2b}$$

\hat{W}_{MAP} is same as $\mathcal{L}(W)$ - hence proved



1.2 (b) PLOT Explanation

When λ is in range 0.1 to 0.2

Ideally if we let the model converge we should observe increasing train error but since maxiter is capped at 10,000 we observe higher error then we should

but test error is fine as increase in lambda reduces overfitting and therefore we observe decrease in test mse initially

Further increasing value of lambda leads to increase in error as sparsity in weights is increased i.e. features with non zero weights reduce. Error later on becomes constant as model is not learning.

1.2 (c)

For weight vector of ISTA the number of components with non zero value is very less compared to that of Ridge. This is expected and justified as loss function of ~~both~~ LASSO is such that it only selects very few features.

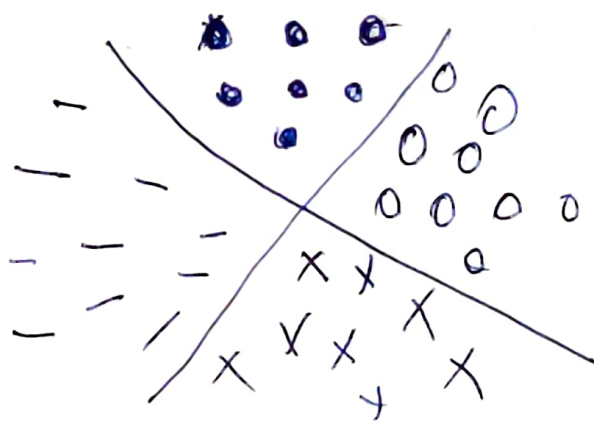
2.1

1 vs 1

→ Splits a multiclass classification into one binary classification problem per each pair of class

→ ~~to~~ Can classify even when classes are mutually linearly separable but not 1 vs rest linearly separable

example →



1 vs rest

→ Splits a multiclass classification into one binary classification problem per class

→ If not linearly separable as 1 vs rest, ~~then~~ Then it can't classify correctly

3.1

(a) increase in λ reduces overfitting ~~with~~ which has direct relation to decrease in variance and increase in bias. To be clear overfitting was reduced as in lasso increase of λ increases sparsity in weight vector i.e. components with zero weight increase hence reduced overfitting

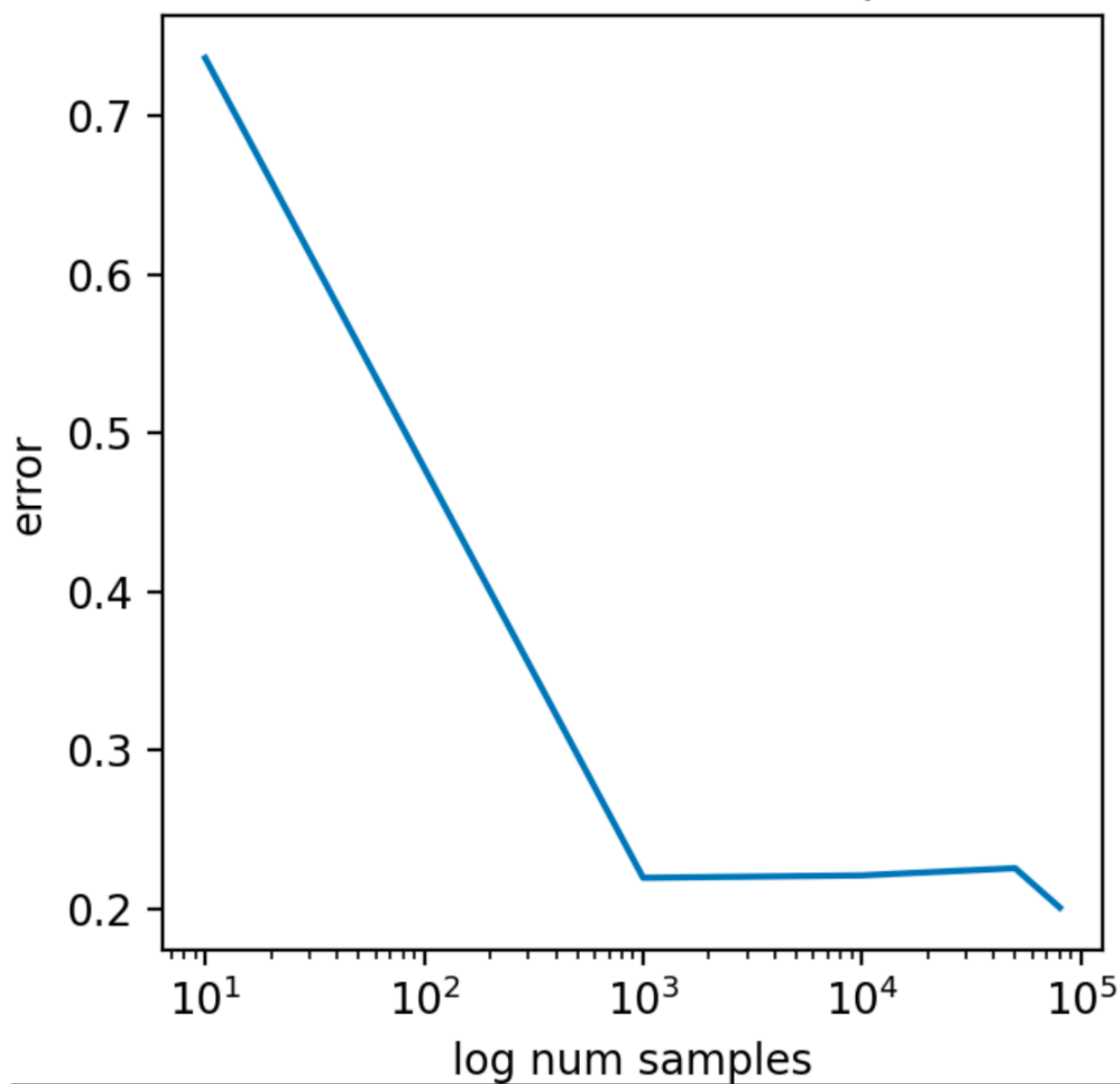
(b) The variance decreases and bias remains the same as the classification planes are more defined

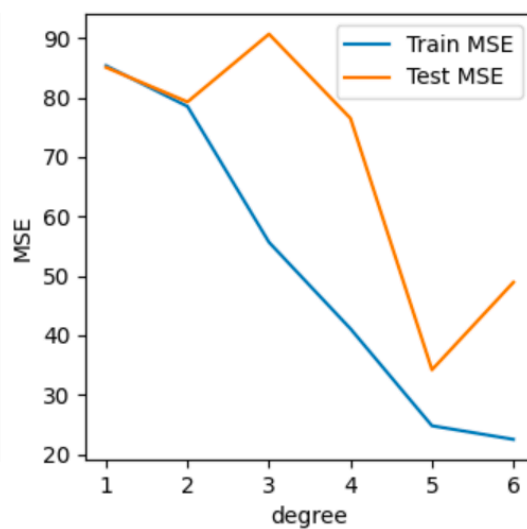
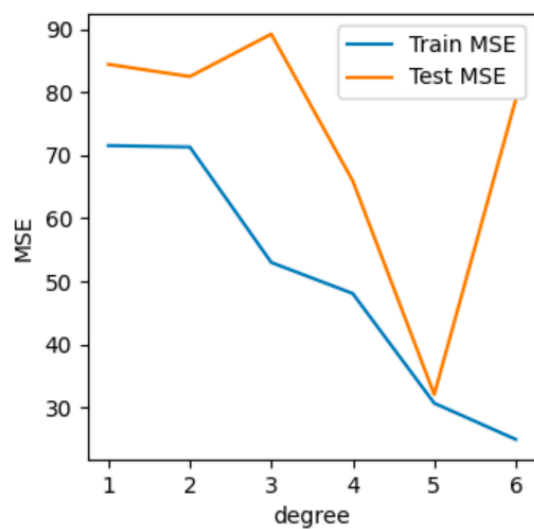
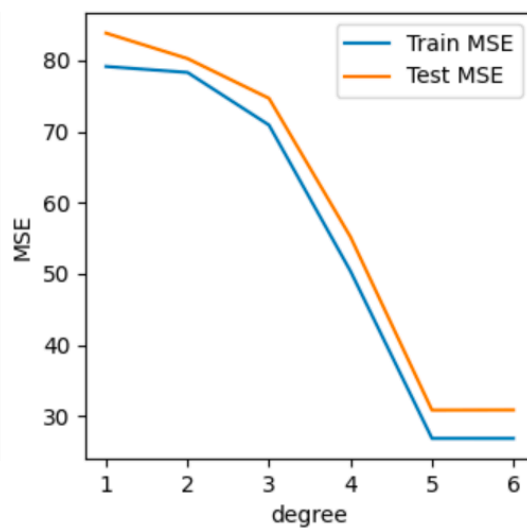
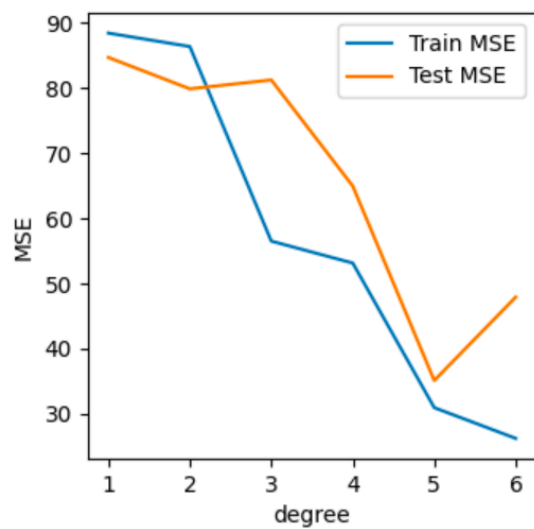
~~(c) Such a modification to data will lead to overfitting faster~~

~~there is increase in variance and decrease in bias~~

(c) both bias and variance remain the same as the newly added features are redundant and not have any effect on output

test error vs num samples





3.2

The test error will decrease as observed in graph as variance is reduced and bias remains the same.

Theoretically at degree = 6 we should obtain the least train error as it has the least bias of all

As we can see empirically in the plots that test mses seems to minimise at degree = 5. The bias decreases with increase in degree but variance increases as well so degree 5 is optimal for test mse. degree 6 is overfitting the data