

Machine Learning

Supervised Learning

Data: X (n obs, p features), y (labels)

Regression, classification

Train/learn/fit f from data (model)

Score: for new x , get $f(x)$

Algos: LR, k-NN, DT, RF, GBM, NN/DL, SVM, NB...

Goal: max acc/min err new data

Metrics: MSE, AUC (ROC)

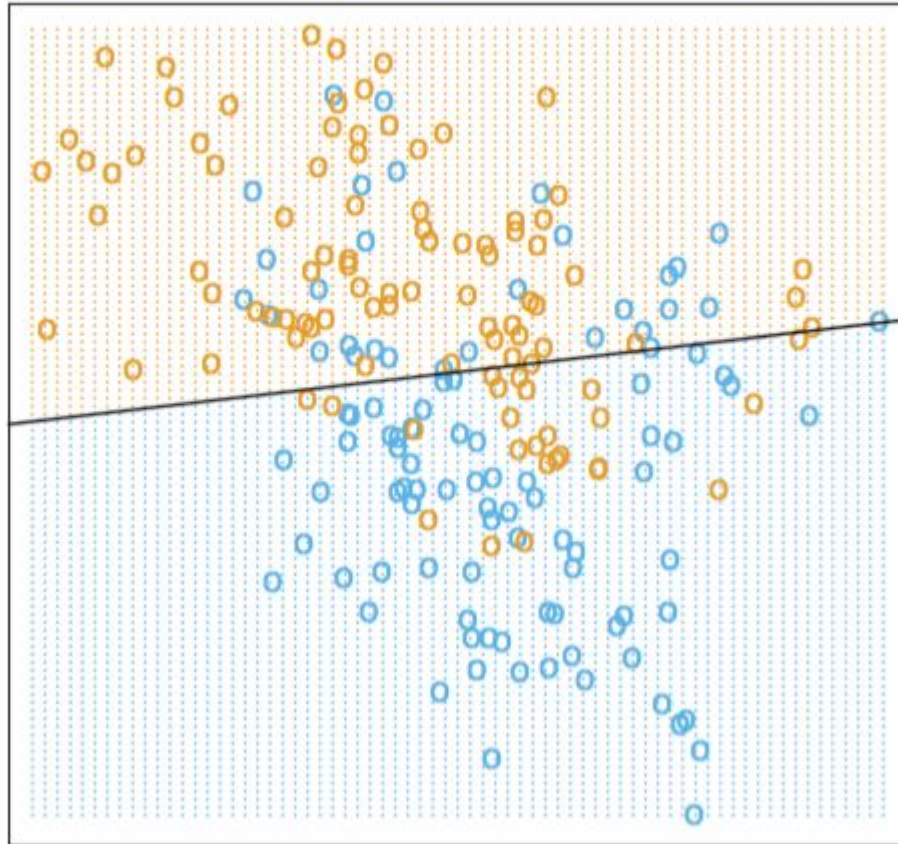
Bad: measure on train set. Need: test/CV

Hyperparams, model capacity, overfitting

Regularization

Model selection

Linear Regression of 0/1 Response



source:

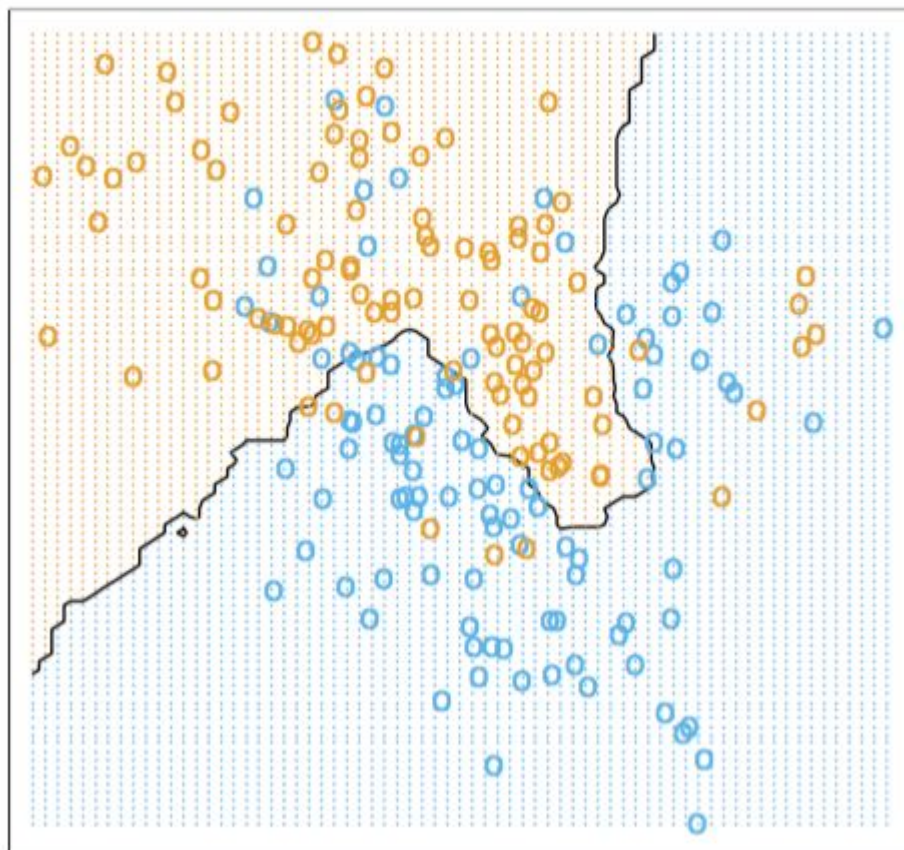
Trevor Hastie
Robert Tibshirani
Jerome Friedman

**The Elements of
Statistical Learning**

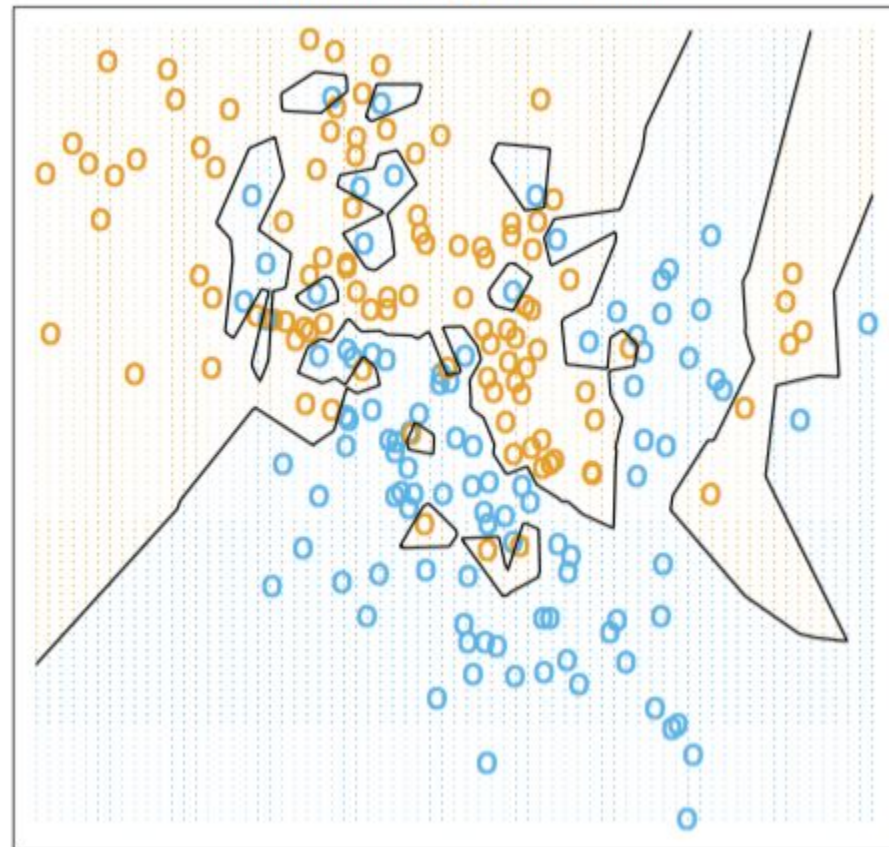
Data Mining, Inference, and Prediction

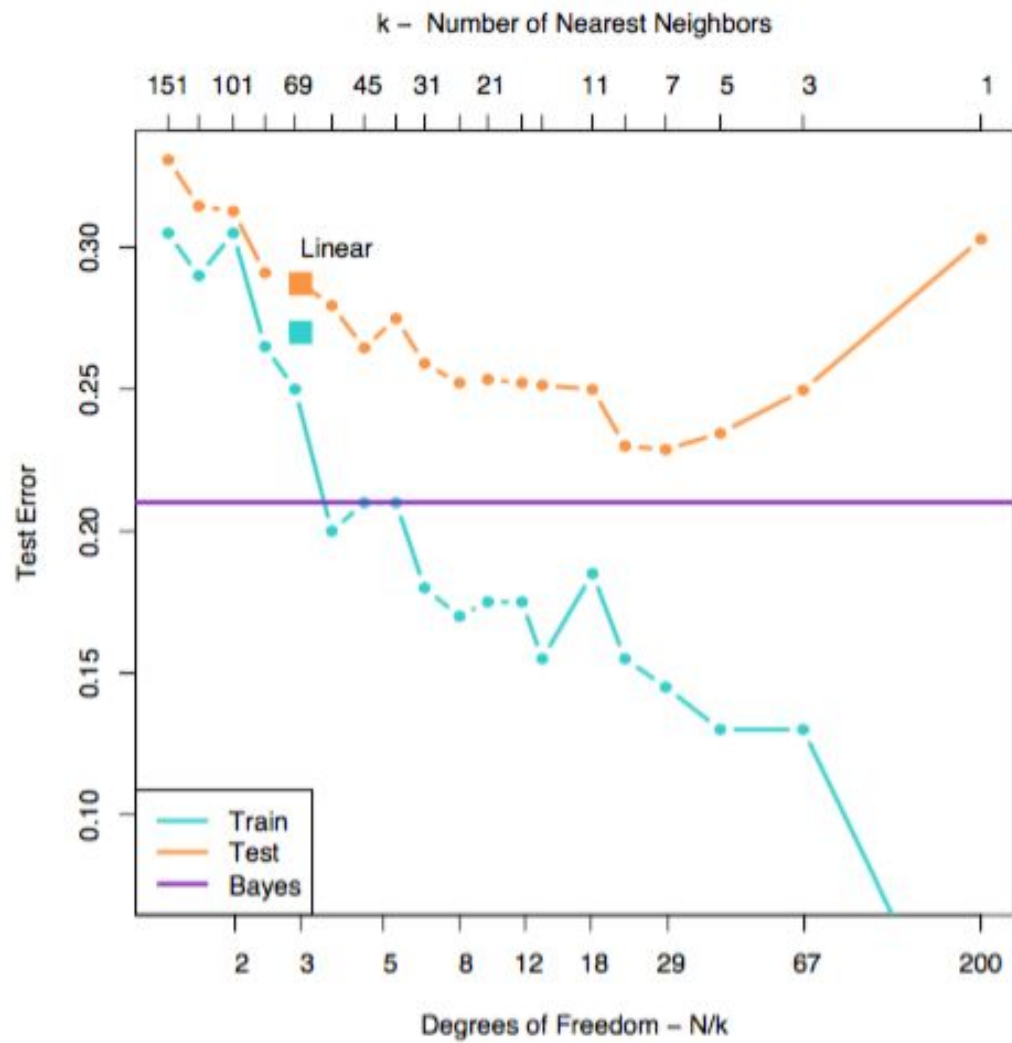
Second Edition

15-Nearest Neighbor Classifier

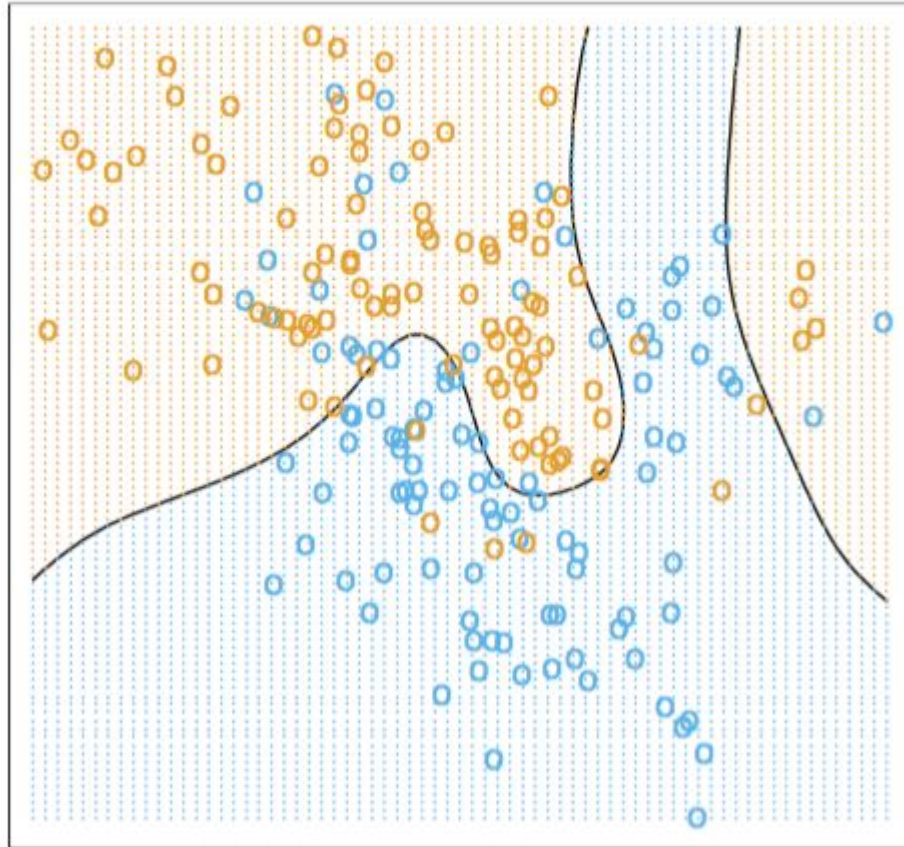


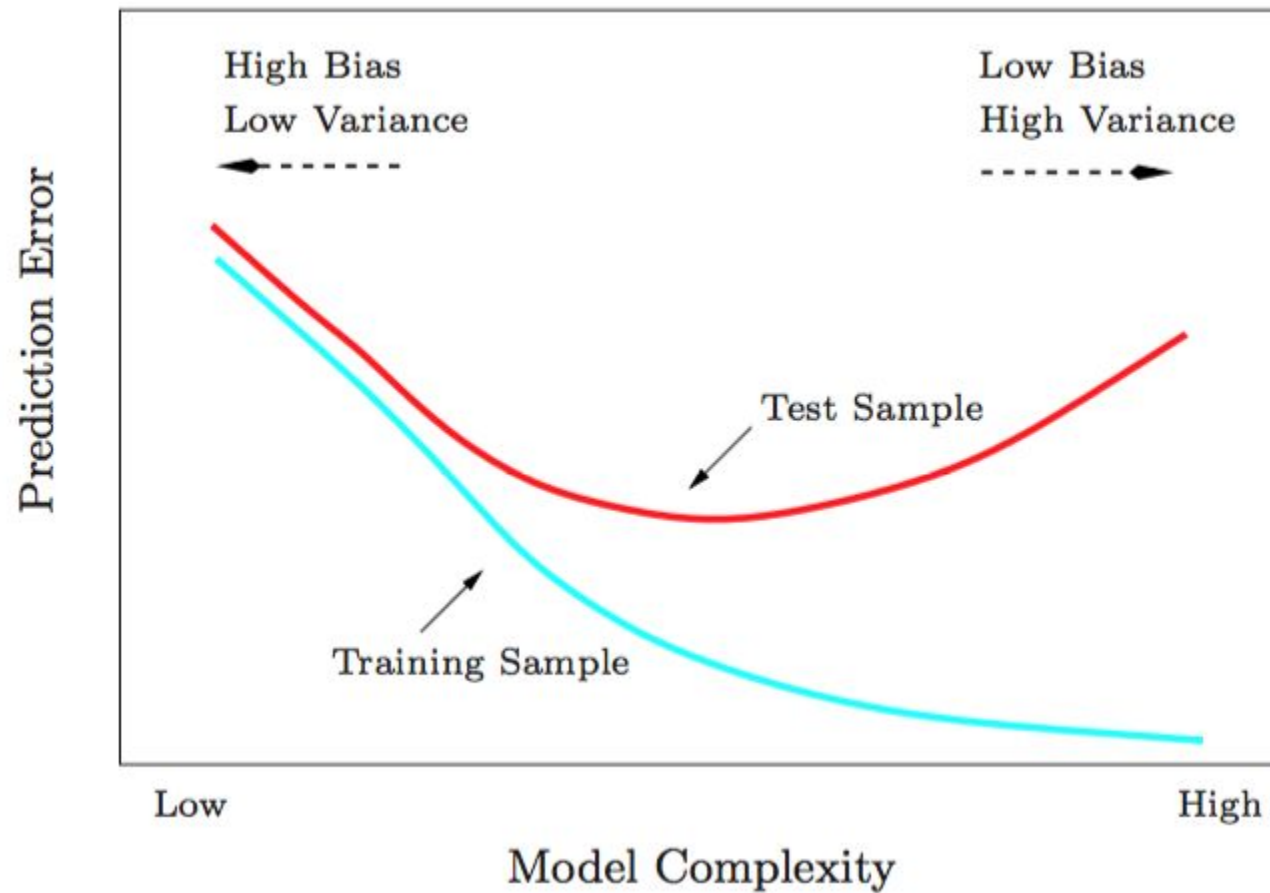
1-Nearest Neighbor Classifier

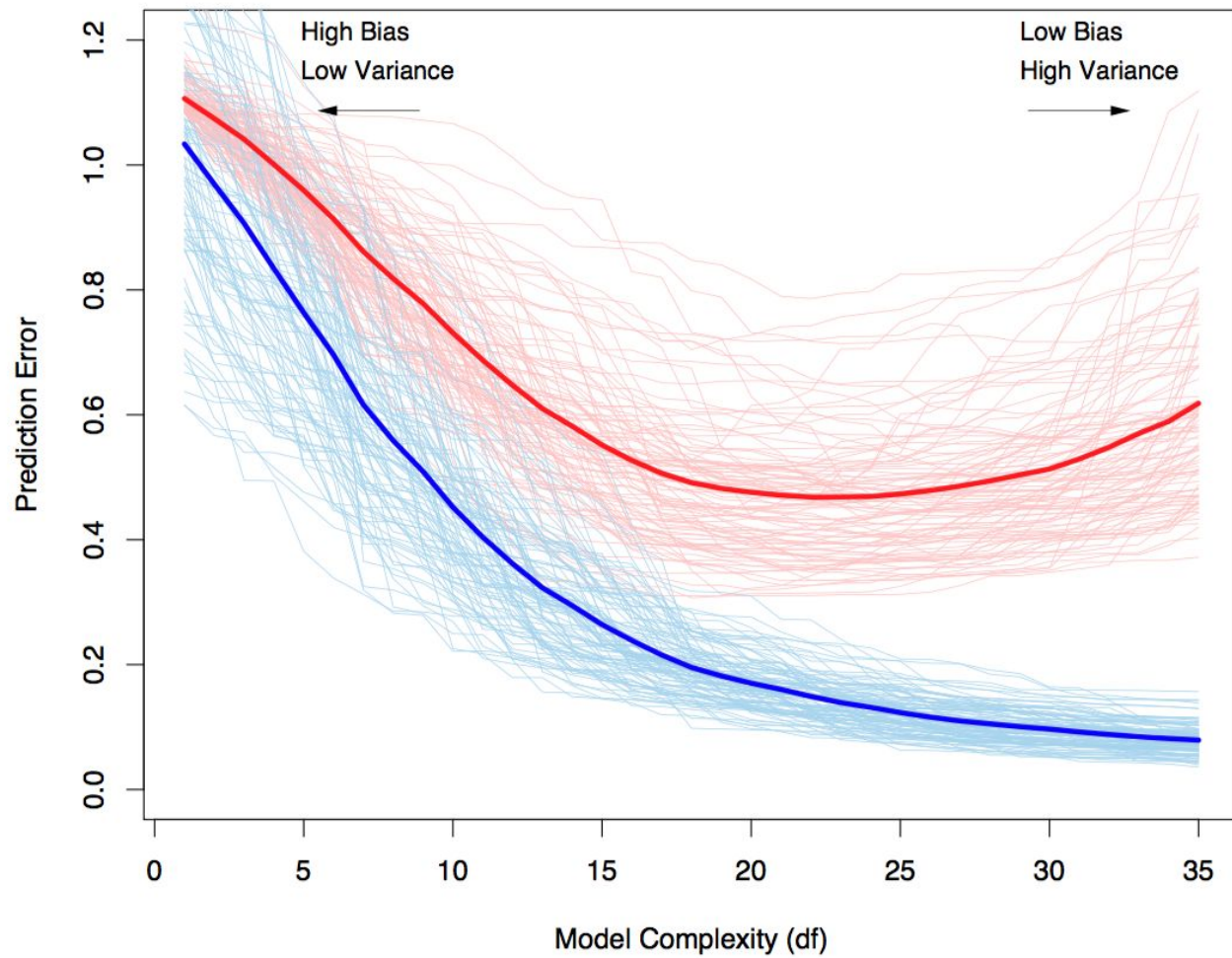


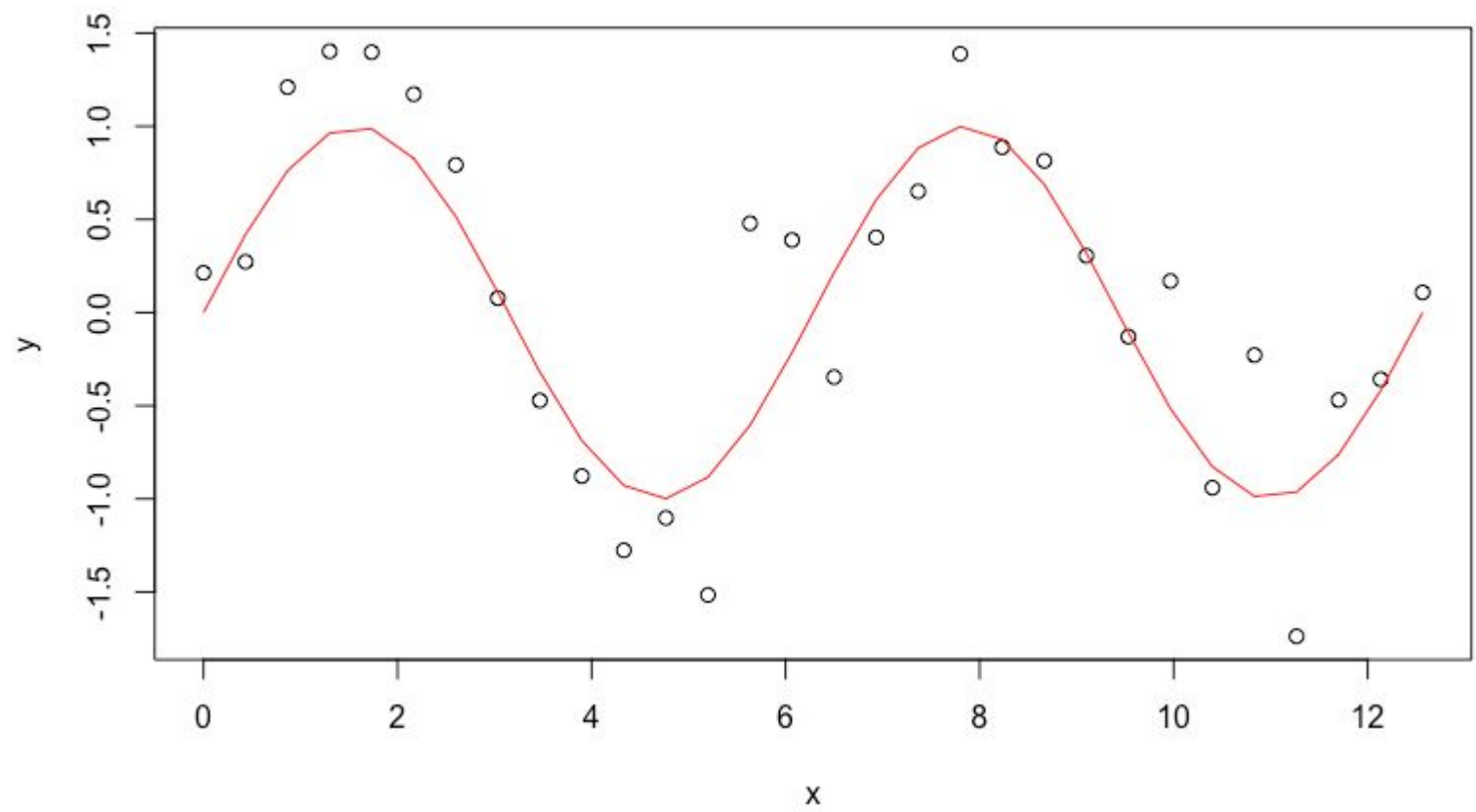


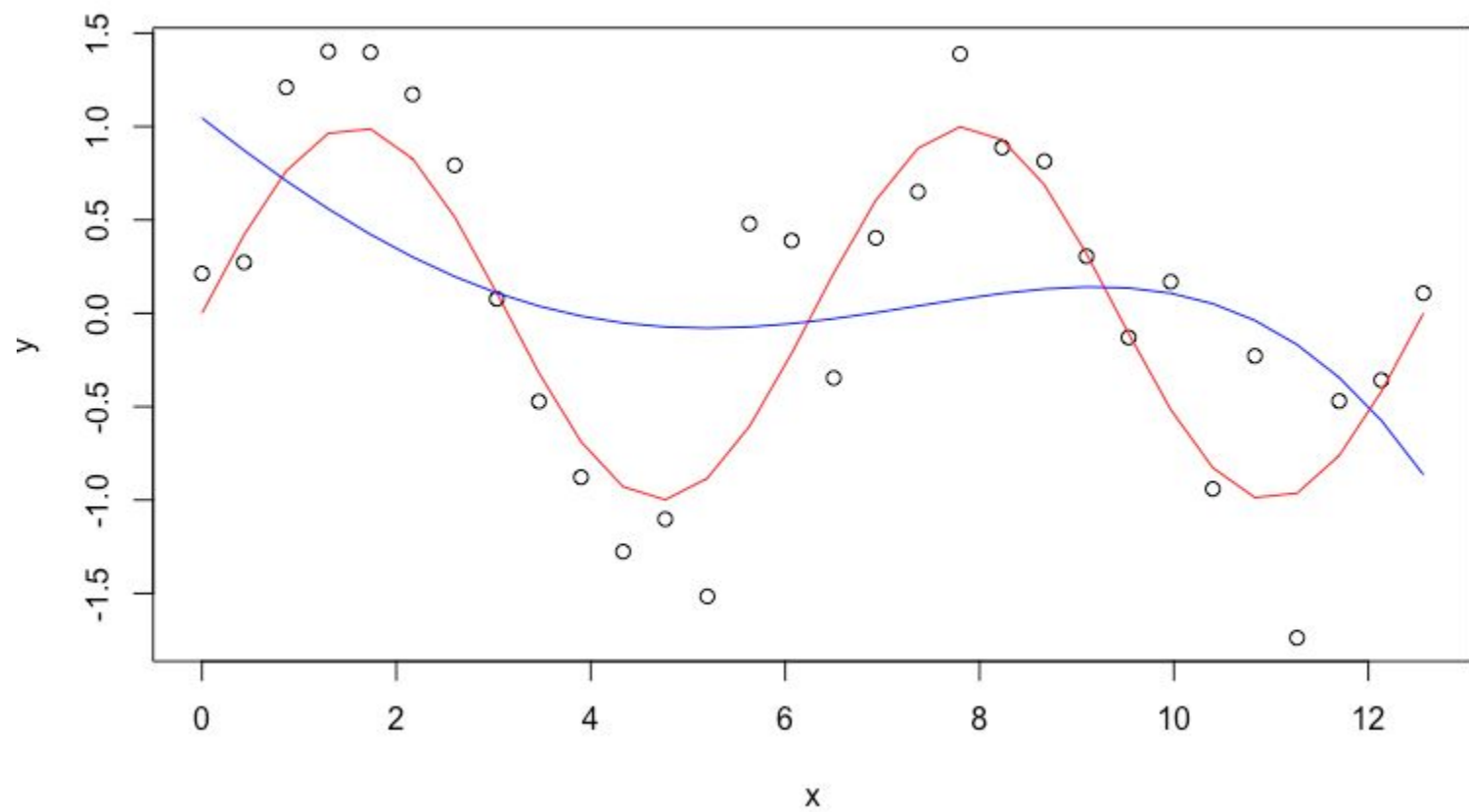
Bayes Optimal Classifier



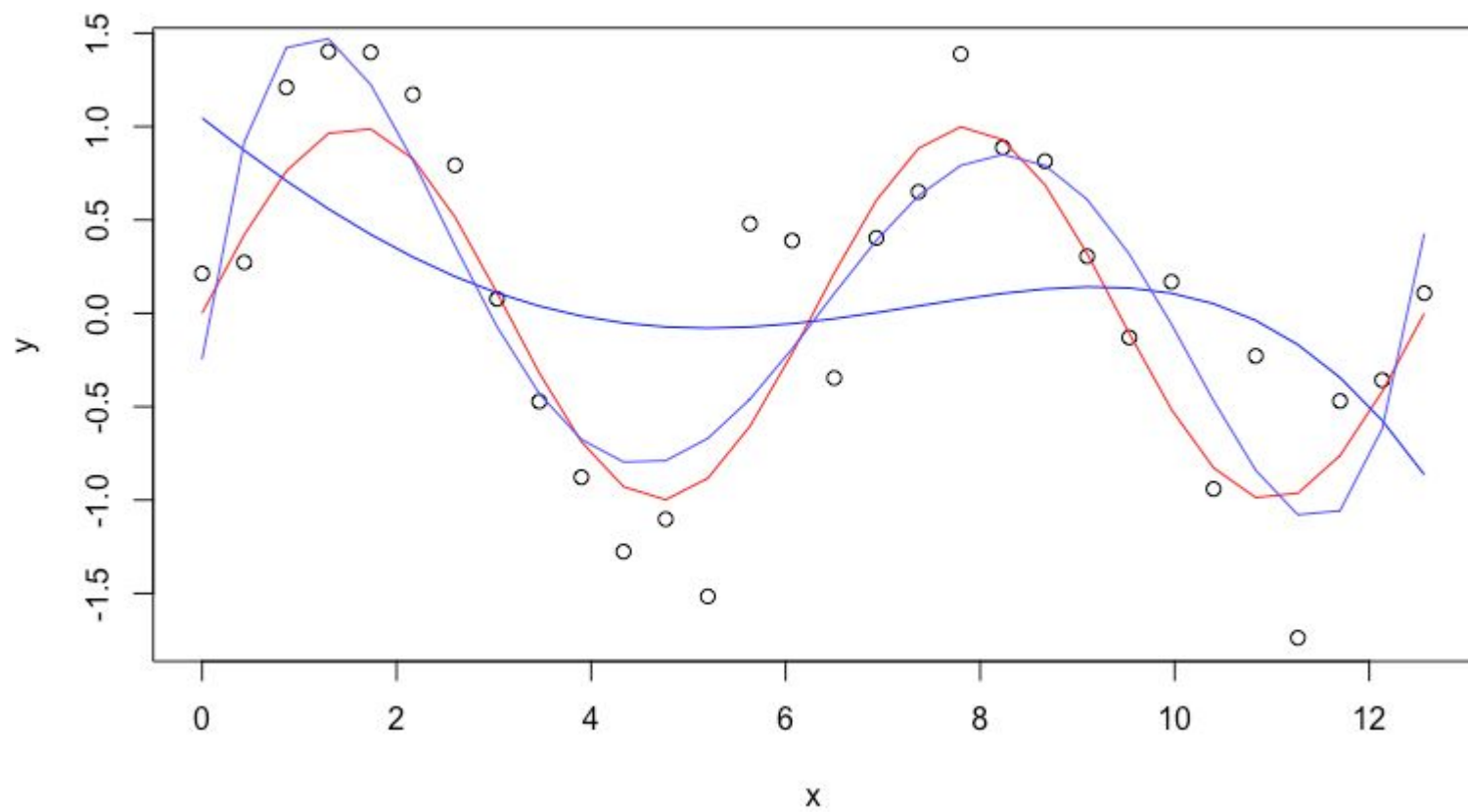




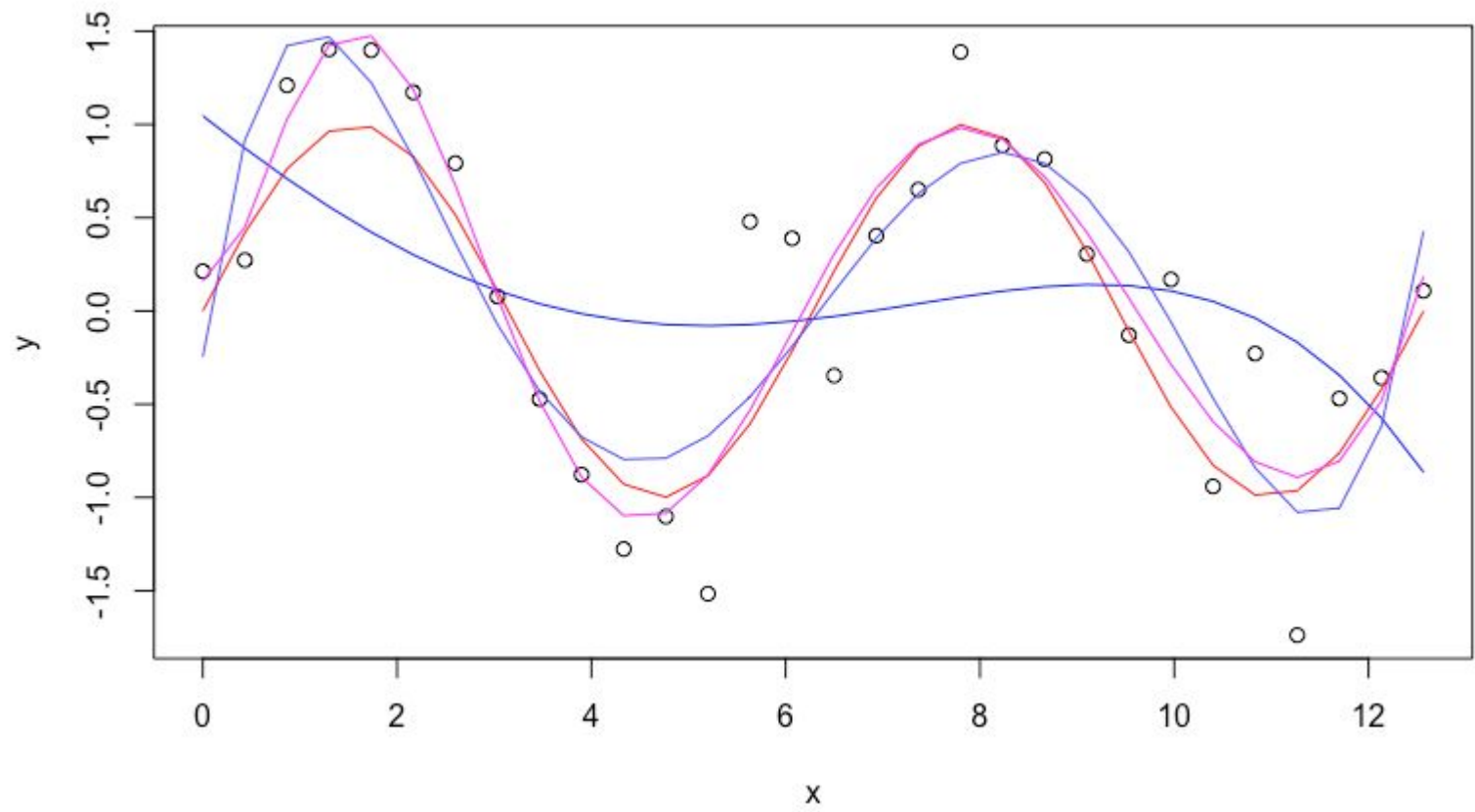




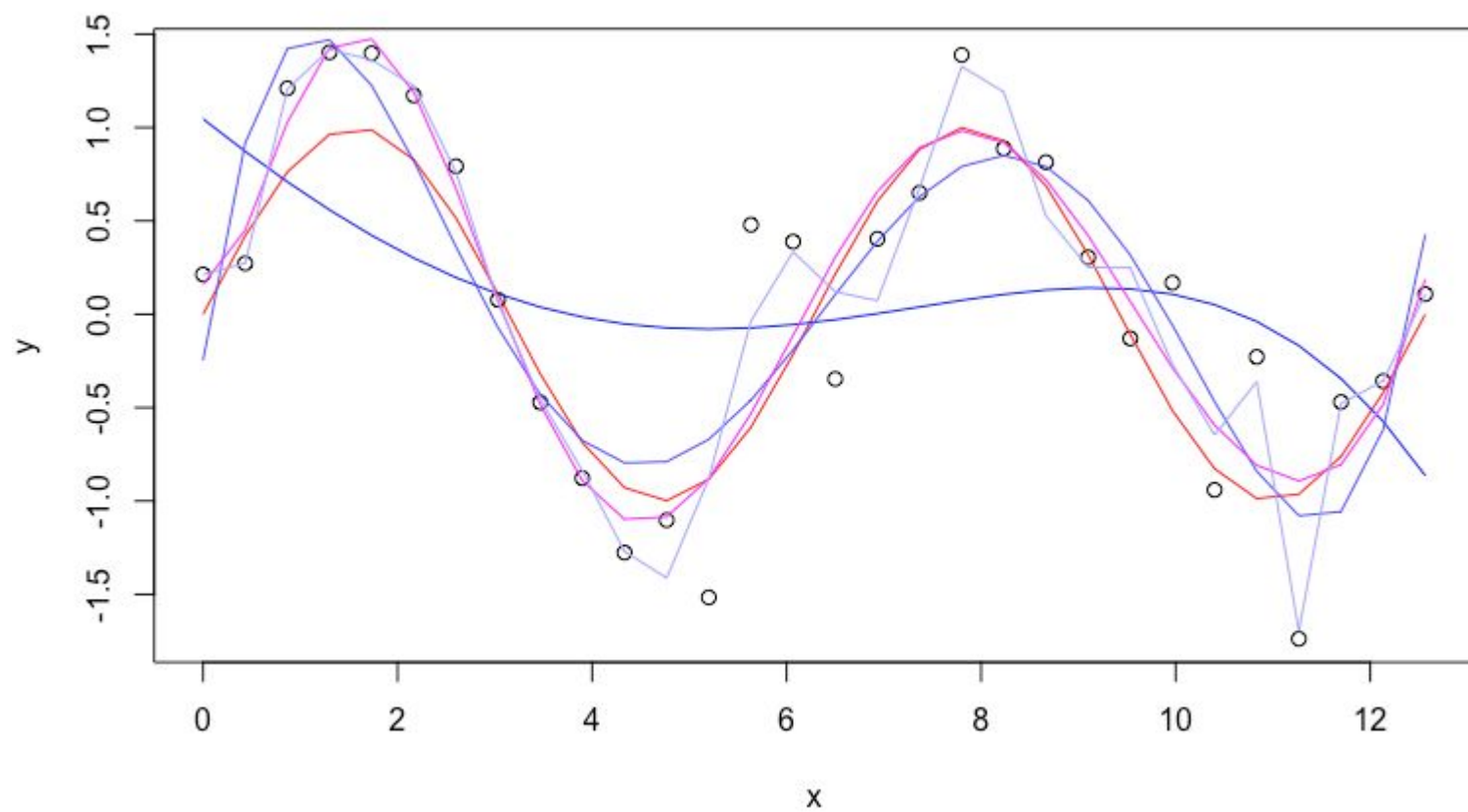
d = 4



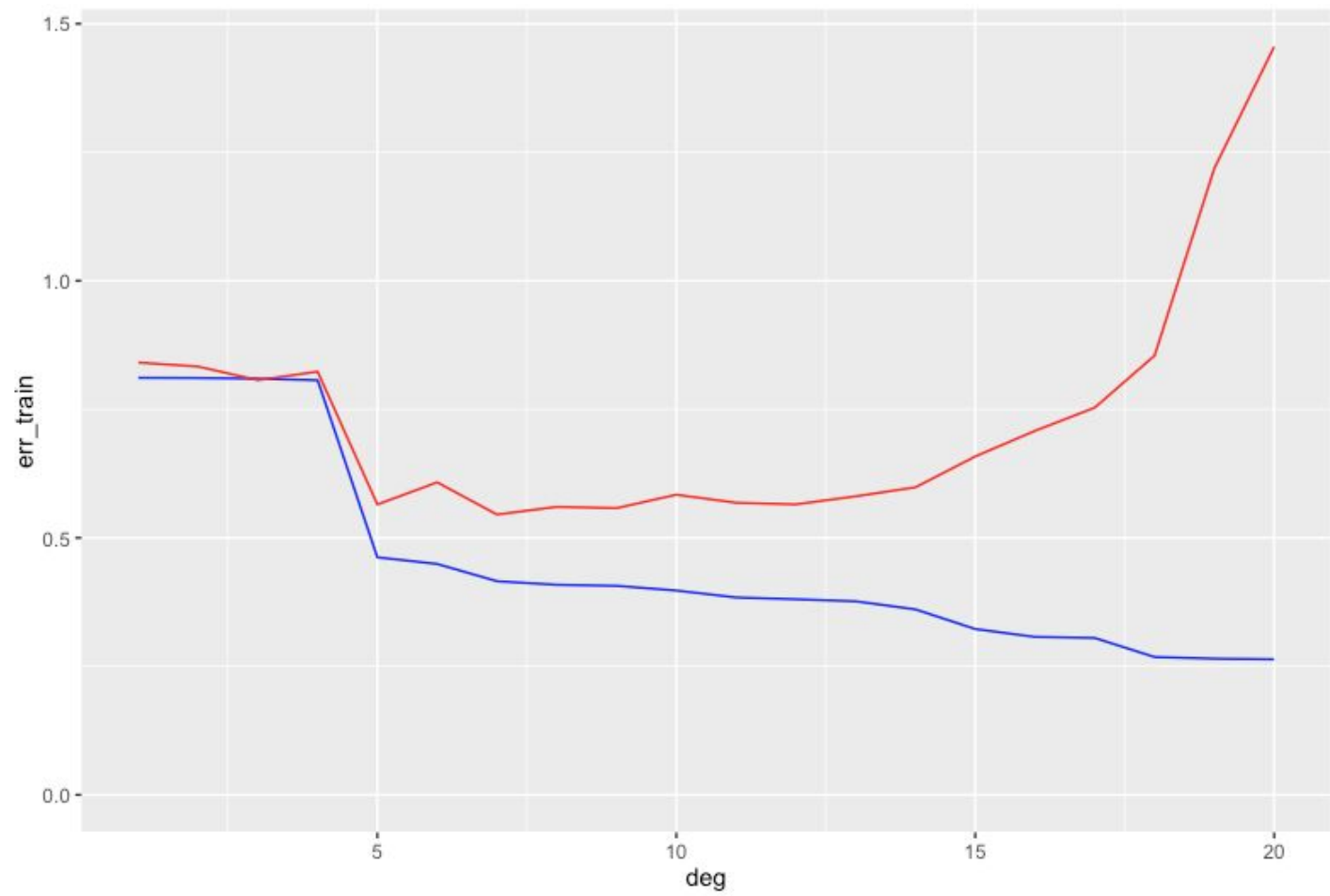
$d = 5$

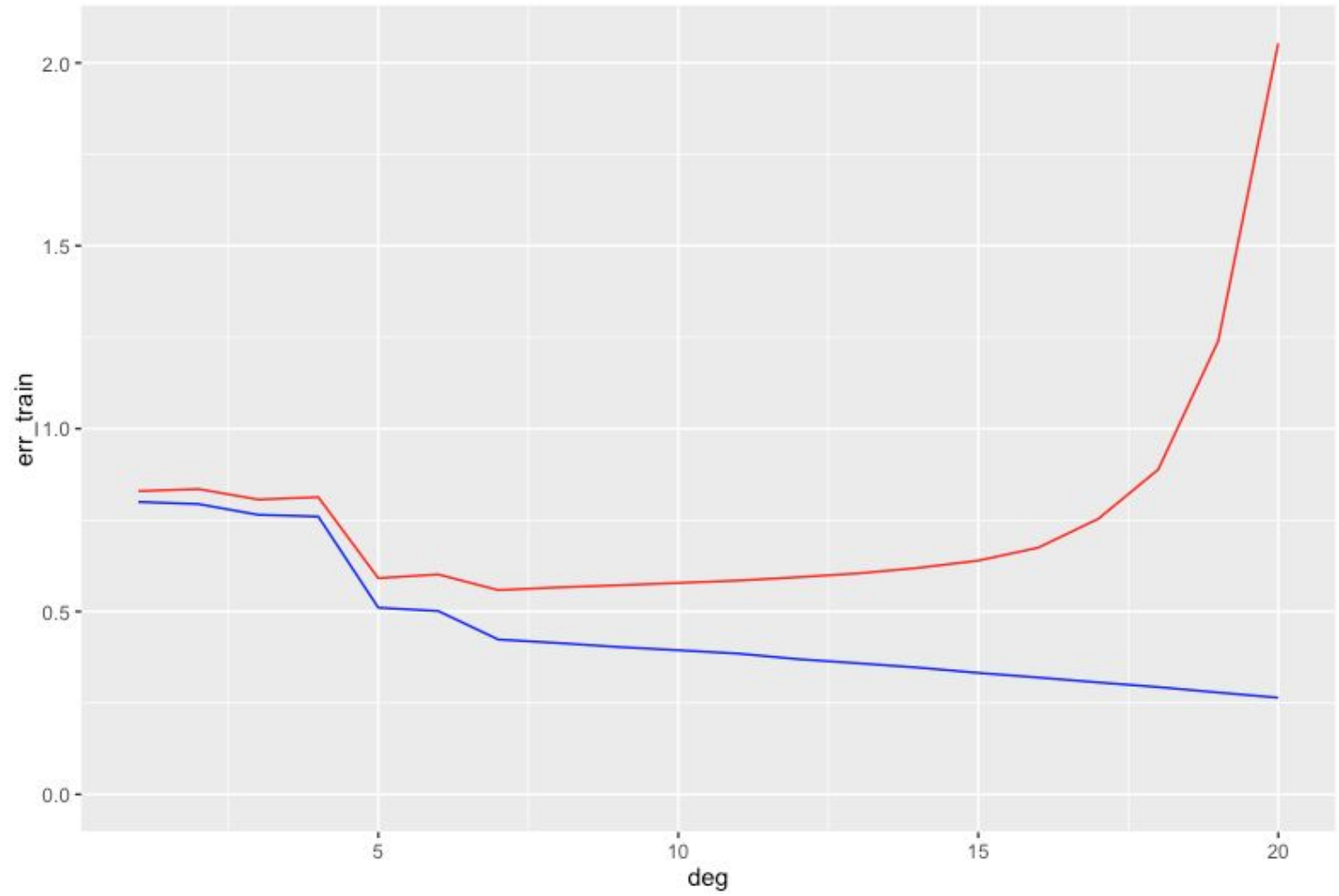


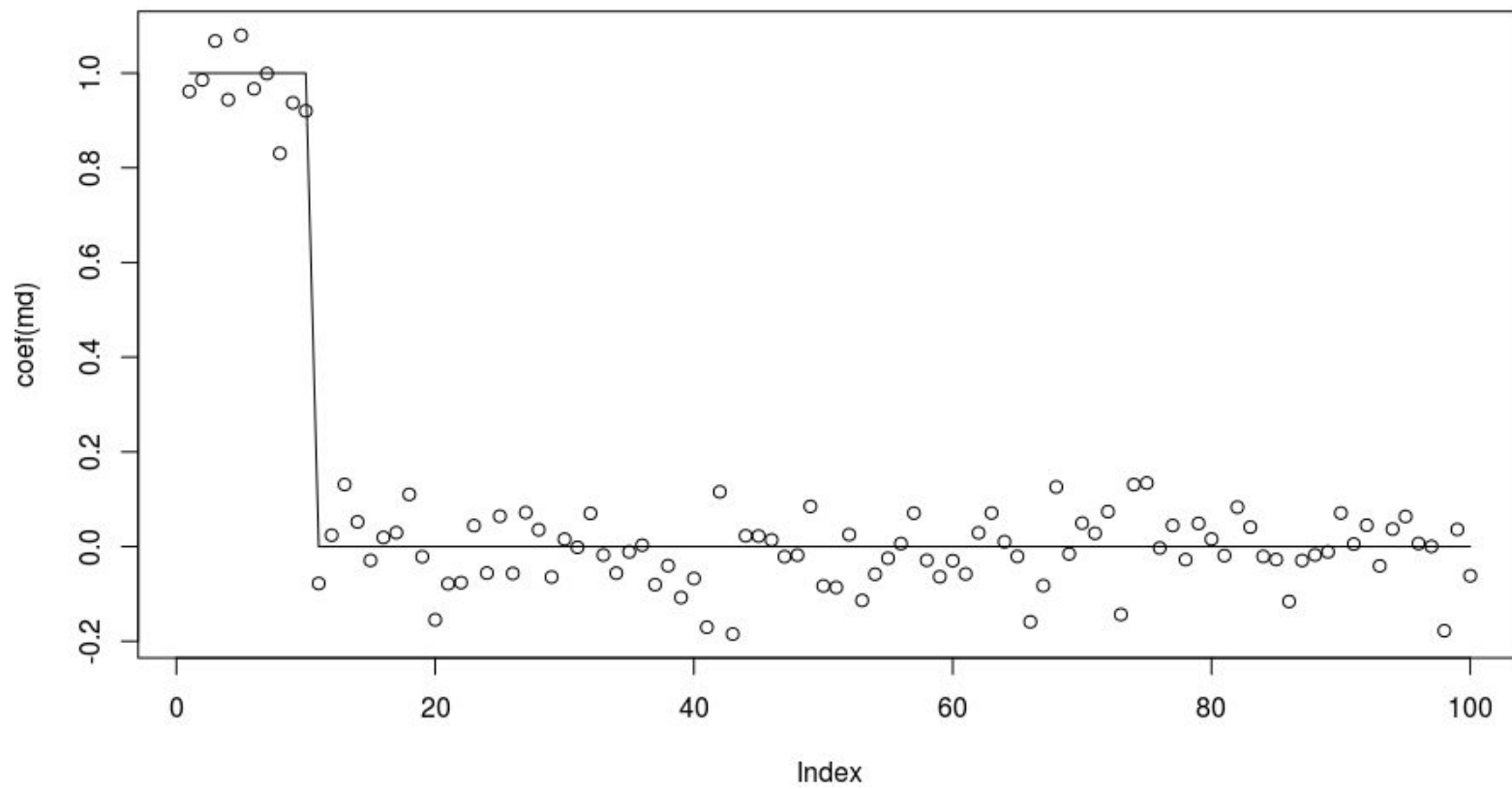
$d = 8$

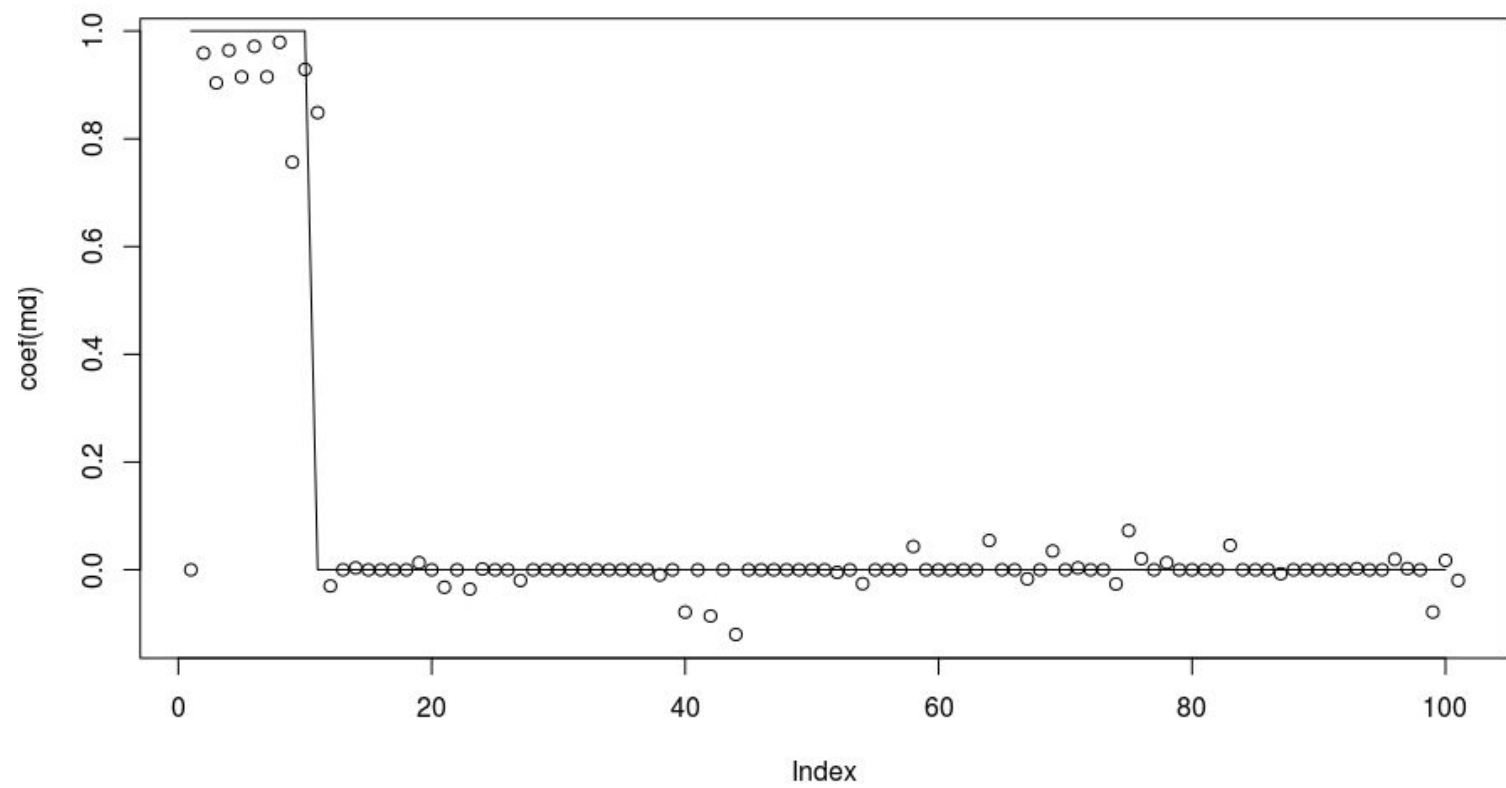


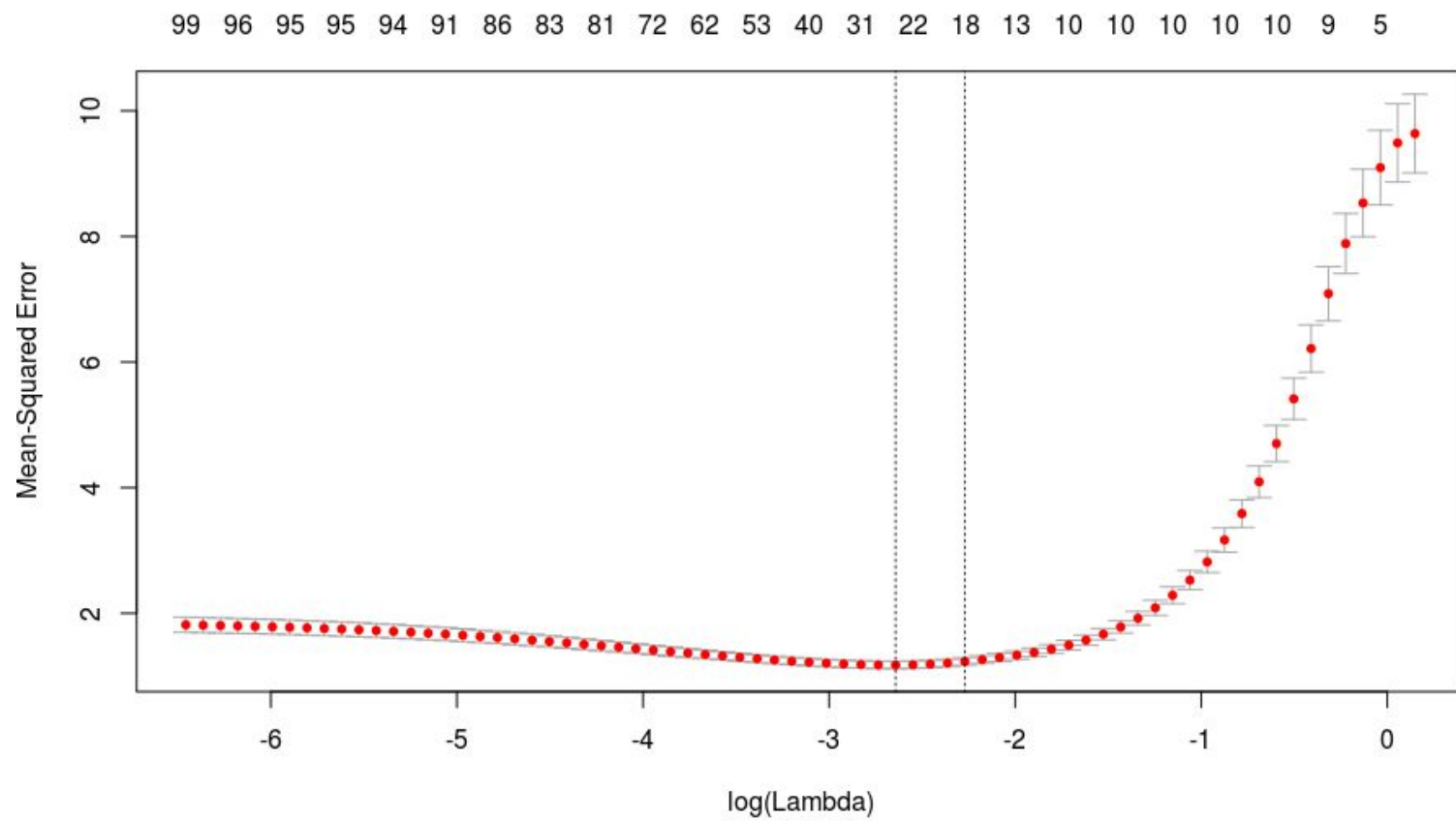
$d = 20$

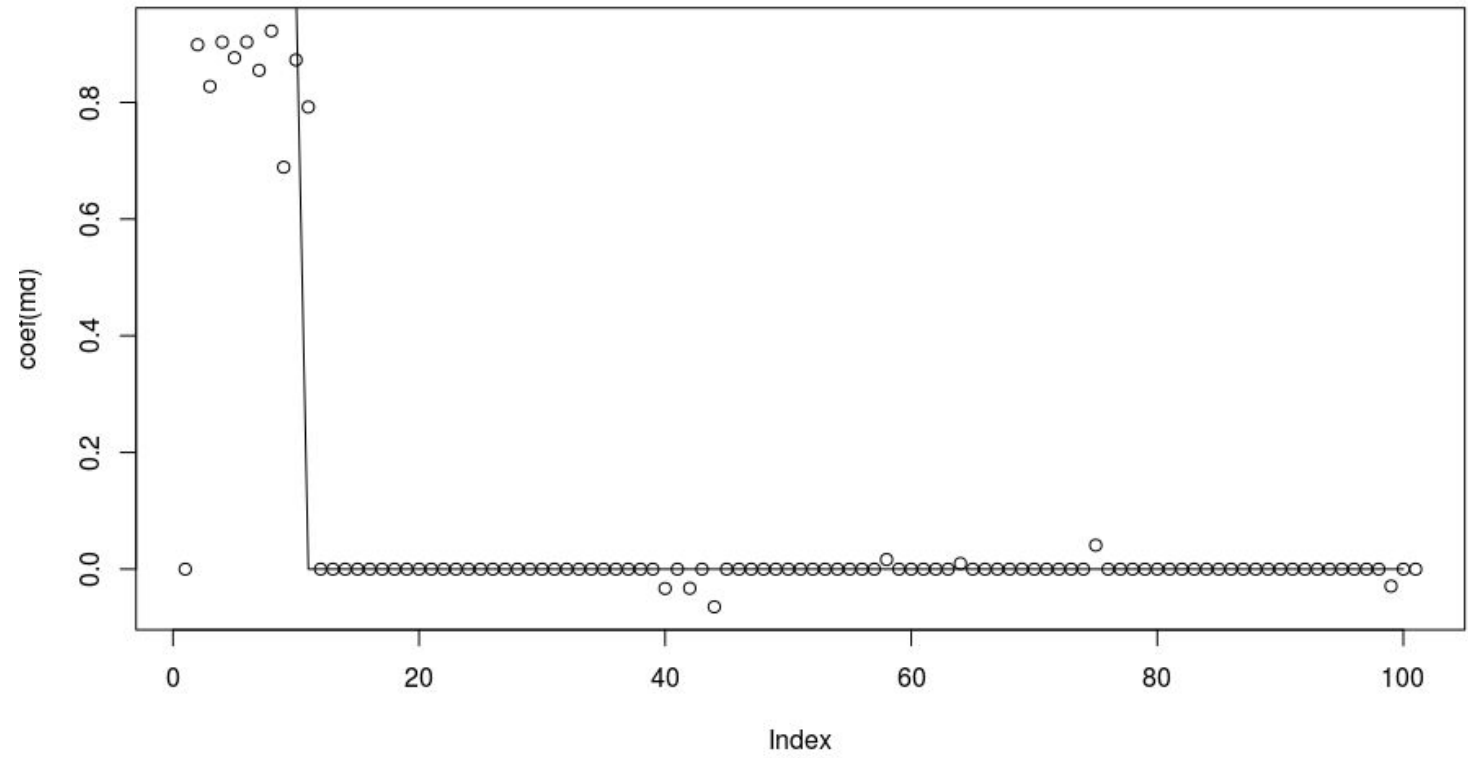










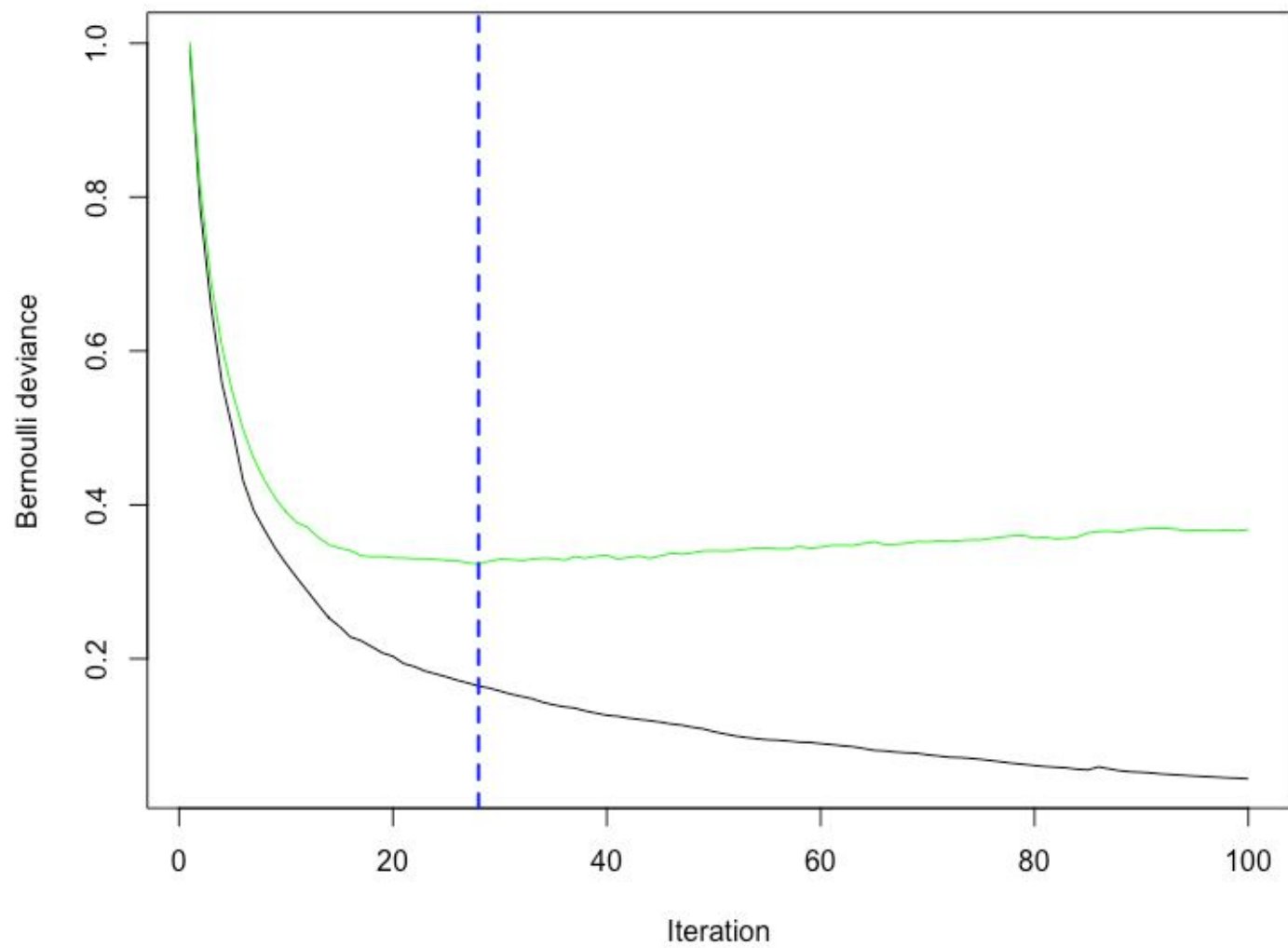


1	2	3	4	5
Train	Train	Validation	Train	Train

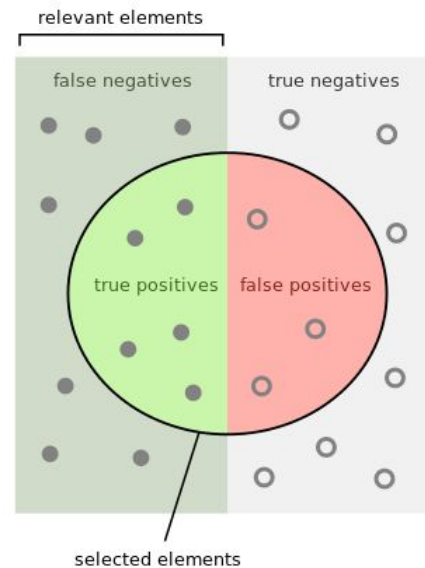
Train

Validation

Test



		True Class	
		Positive (P)	Negative (N)
Predicted Class	Positive (+)	True Positive Count (TP)	False Positive Count (FP)
	Negative (-)	False Negative Count (FN)	True Negative Count (TN)



$$\text{Precision} = \frac{\text{true positives}}{\text{true positives} + \text{false positives}}$$

$$\text{Recall} = \frac{\text{true positives}}{\text{true positives} + \text{false negatives}}$$

