

1. Exercise 1

- a. Better because a flexible approach will fit the data closer with such a large sample size.
- b. Worse because the flexible approach would overfit the small number of observations n .
- c. Better because a flexible approach obtains a better fit with more degrees of freedom.
- d. Worse because a flexible approach would fit to the increased variance of the error terms.

2. Exercise 2

- a. This scenario is a regression problem and we are most interested in inference.
 - i. $n = 500$ firms
 - ii. $p =$ profit, number of employees, industry
- b. This scenario is a classification problem and we are most interested in prediction.
 - i. $n = 20$ similar products
 - ii. $p =$ price charged, marketing budget, competition price, ten other variables
- c. This scenario is a regression problem and we are most interested in prediction.
 - i. $n = 52$ weeks of 2012 weekly data
 - ii. $p =$ % change in US market, % change in British market, % change in German market

3. Exercise 3

- a. Forgot to transfer my graph picture and ran out of time to add it in. Will be better prepared for the next problem set.
- b. Each of the five curves has the shape displayed because:
 - i. Variance – this increases monotonically because the increases in flexibility yield an overfit.
 - ii. Bias – this decreases monotonically because the increases in flexibility yield a closer fit.
 - iii. Test Error – this shows as a concave curve facing up because of increases in flexibility yielding a closer fit before it overfits.
 - iv. Training Error – this decreases monotonically, similarly to bias, because increases in flexibility yield a closer fit.
 - v. Bayes Error – has a shape that defines the lower limit and the test error is bounded by the Bayes error due to variance in the error in the output values.

4. Exercise 4

- a. Three real-life applications in which classification might be useful:
 - i. Car Part Replacement. The response would be “need to replace” or “in good standing”. The predictors would be “age of part” and “mileage used for”. The goal of this application would be prediction.
 - ii. Stock Market Price Prediction. The response would be “up”, “down”, or “flat”. The predictors would be “yesterday price % move”, “previous

week price % move", "order-flow", and many potential others. The goal of this application would be prediction.

- iii. Fed Rate Movement Prediction. The response would be "increase", "decrease", or "stable". The predictors would be "economic state", "stock-market previous month price % movement", "global political tensions", etc. The goal of this application would be prediction.

5. Exercise 5

- a. The advantages of a very flexible approach for regression or classification compared to a less flexible approach are obtaining a better fit for non-linear models and in general, decreasing bias. The disadvantages of a very flexible approach for regression or classification compared to a less flexible approach are that it requires estimating a greater number of parameters, it can overfit the data, and a general increase in variance. A flexible approach would be preferred to a less flexible approach if we are interested in prediction and not the interpretation of the results. The opposite is true as well though – a less flexible approach would be preferred to a more flexible approach if we are interested more in the inference and interpretation of the results compared the prediction of results.

6. Exercise 6

- a. A parametric statistical learning approach reduces the problem of estimating f down to estimating a set of parameters because it assumes a form for f . This is different from a non-parametric statistical learning approach where it does not assume a functional form for f and so requires a very high number of observations in order to accurately estimate f . The advantages of a parametric approach to regression or classification (as opposed to a non-parametric approach) are the simplification of modeling f to a few parameters and the reduction of observations that are required compared to a non-parametric approach. The disadvantages of a parametric approach to regression or classification are a potential to inaccurately estimate f if the form of f assumed is wrong or the potential to overfit the observations if more flexible models are used instead of inflexible models.