

Ch3 Regression to Machine Learning

Contents

3 Subsection

A.1	Statistical Learning
A.2	How do we find 'overall pattern'? - Inference
A.3	How do we find 'overall pattern'? - Prediction
A.4	Polynomial Regression 1
A.5	Problem
A.6	Measure of Quality of Fit
A.7	KEY CONCEPT
A.8	Leave-some-out Fitting Procedure 1
A.9	Leave-some-out Fitting 2
A.10	k-fold Cross Validation
A.11	k-fold Cross Validation
A.12	CV within the Training Set
A.13	Training MSE vs Validation MSE
A.14	Assessing Model Prediction Accuracy
A.15	Bias-Variance Trade-Off
A.16	Prediction MSE

A.17 In the Classification Setting	
A.18 Trade-off in the new approach	
A.19 K-Nearest Neighbor	
A.20 K-NN examples	

Textbook: James et al. ISLR 2ed.

3 Subsection

[\[ToC\]](#)

A.1 Statistical Learning

- General Model

$$Y = f(X) + \epsilon$$

- We don't want to assume that $f(X)$ is linear function.
- Two types of motivation:
 - Model Estimation
 - Prediction
- Pattern recognition

A.2 How do we find 'overall pattern'? - Inference

- Want to understand the relationship between X and Y
- Which predictors are associated with the response?
- What is the relationship between the response and each predictor?
- Can the relationship between Y and each predictor be adequately summarized using a linear equation, or is the relationship more complicated?

A.3 How do we find 'overall pattern'? - Prediction

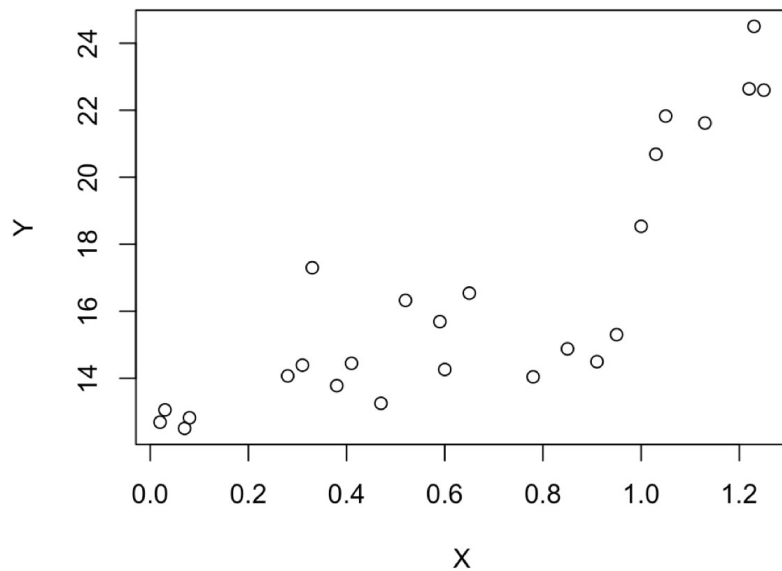
- Want to guess the next Y as accurate as possible

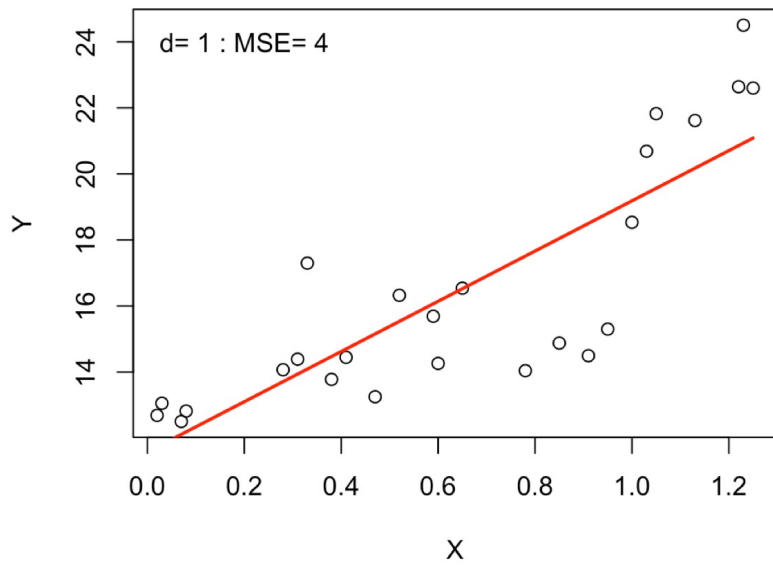
$$\hat{Y} = \hat{f}(X)$$

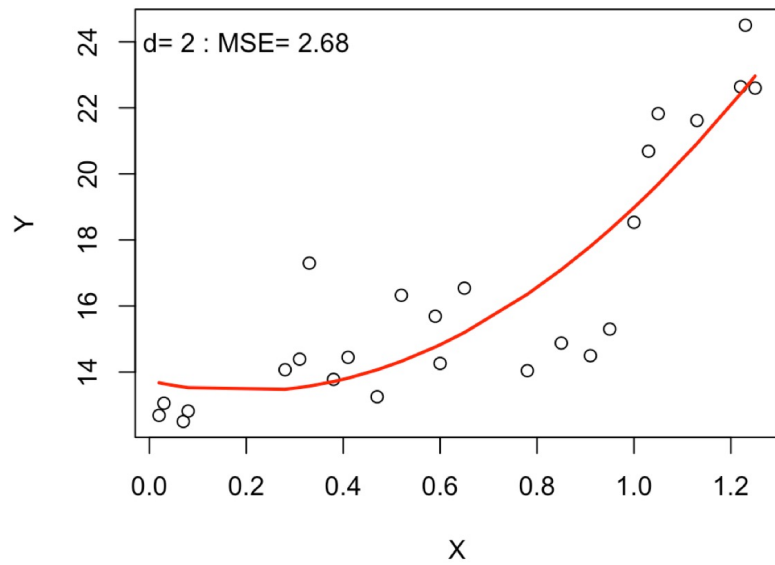
- f can be a black box
- reducible error and irreducible error in prediction
- Want to reduce prediction Mean Squared Error:

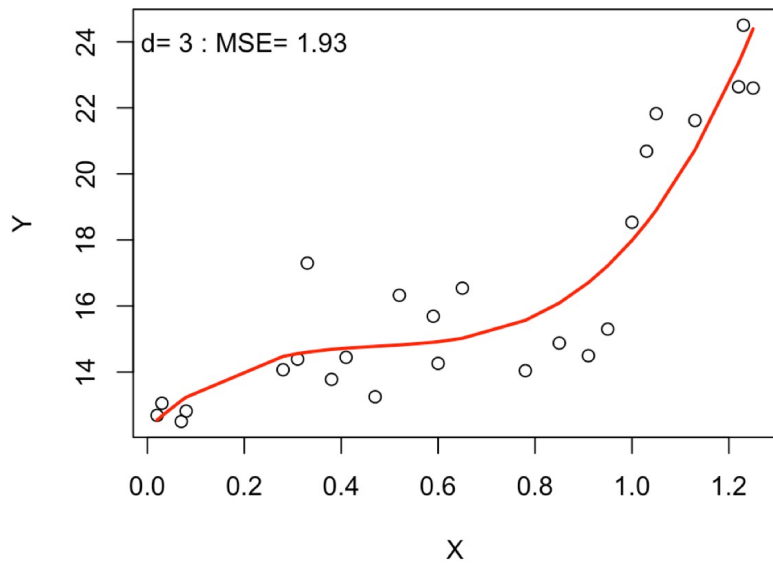
$$MSE = E(Y - \hat{Y})^2 = E(Y - \hat{f}(X))^2$$

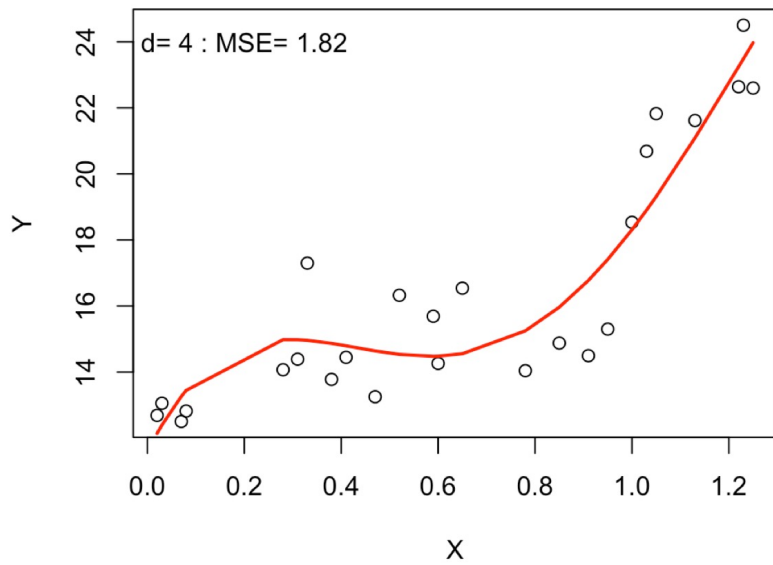
A.4 Polynomial Regression 1

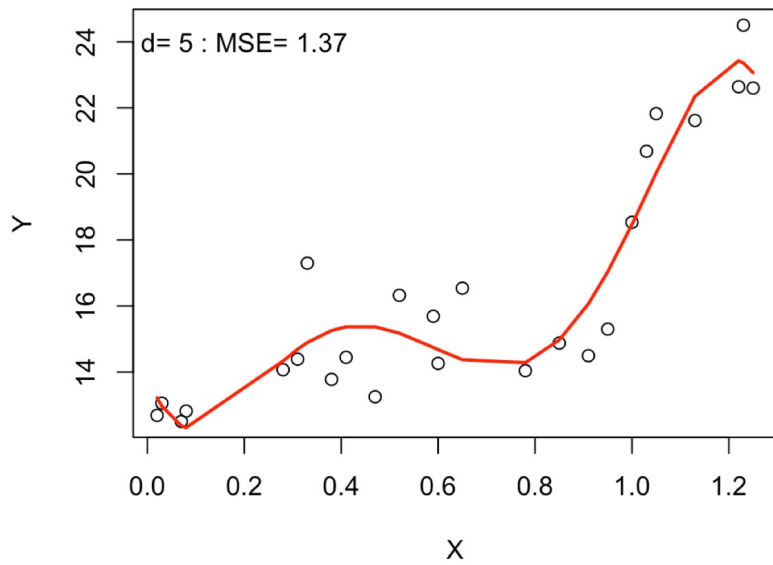


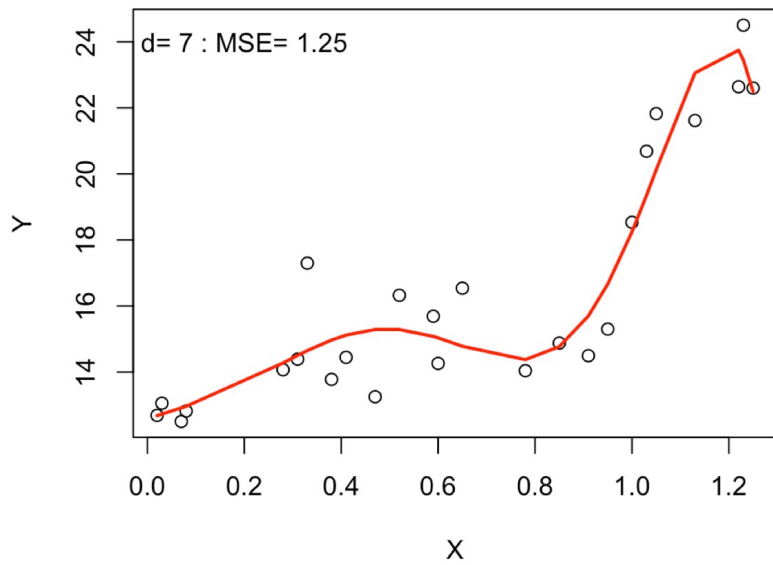


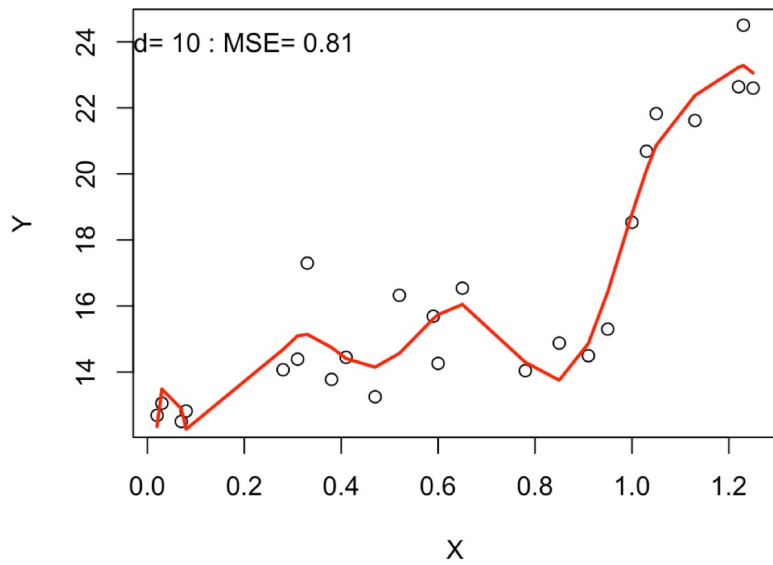


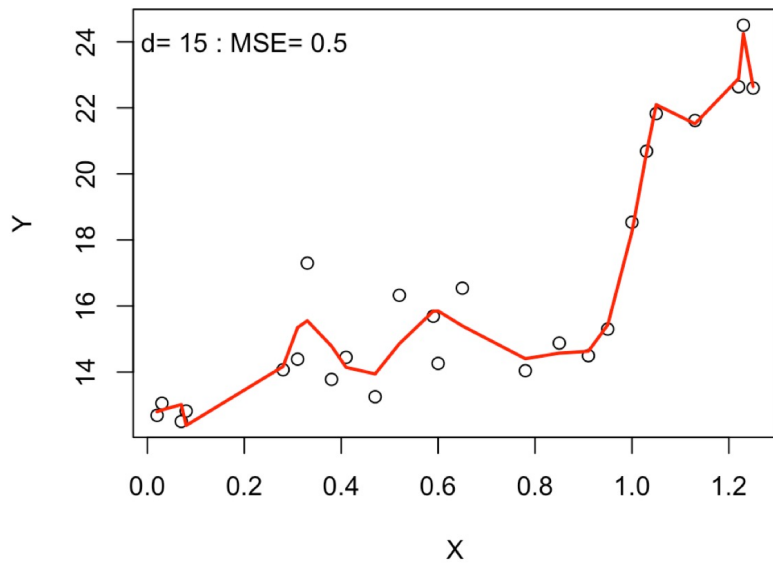


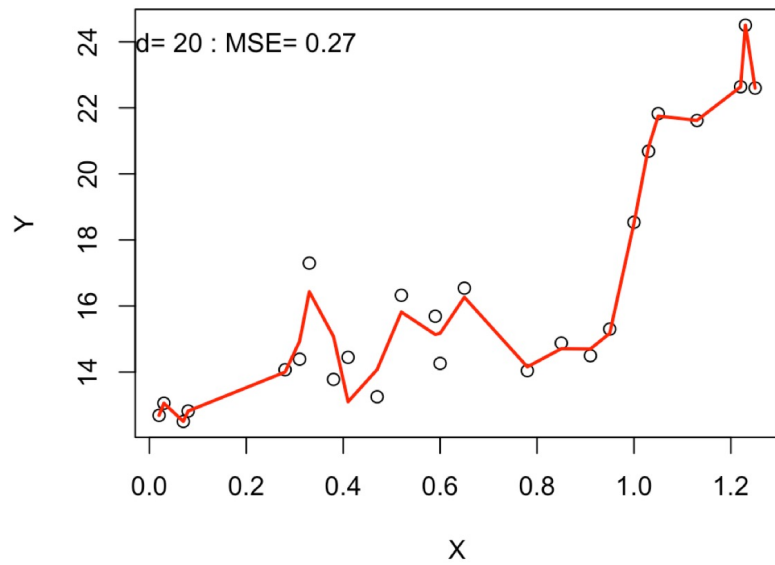












A.5 Problem

- More flexibility in the model is always going to result in better fit to the data.
- Better fitting model is not always inferential.
- Better fitting Leave some out and use it for 'validation' and 'testing'.
- Underlying mechanism:

$$Y = f(X) + \epsilon$$

A.6 Measure of Quality of Fit

- Training MSE (sample)

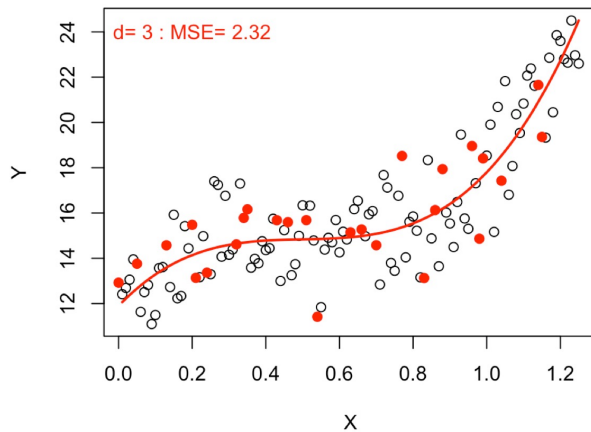
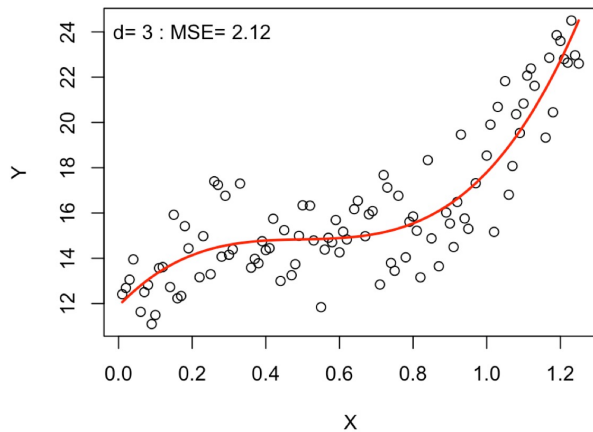
$$\text{MSE}_{tr} = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{f}(x_i))^2$$

- But we want minimum Prediction MSE

$$\text{MSE} = E(Y - \hat{f}(X))^2$$

- Solution: look at Test MSE (sample) as estimator

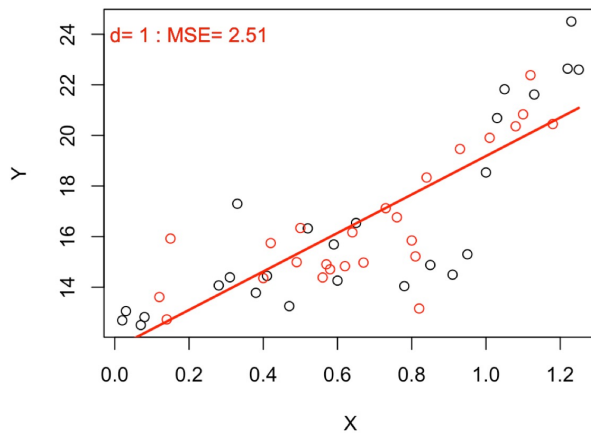
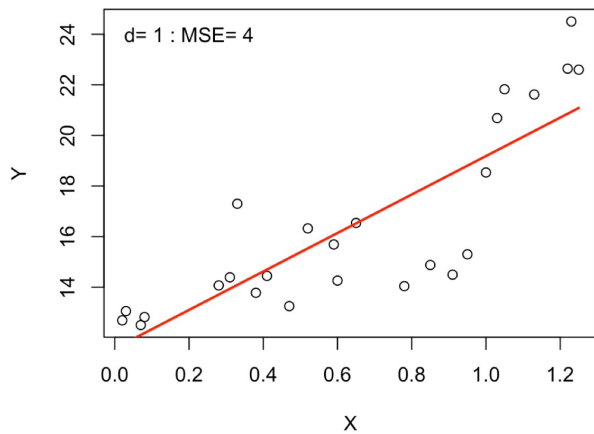
$$\text{MSE}_{test} = \frac{1}{m} \sum_{j=1}^m (y_j - \hat{f}(x_j))^2$$

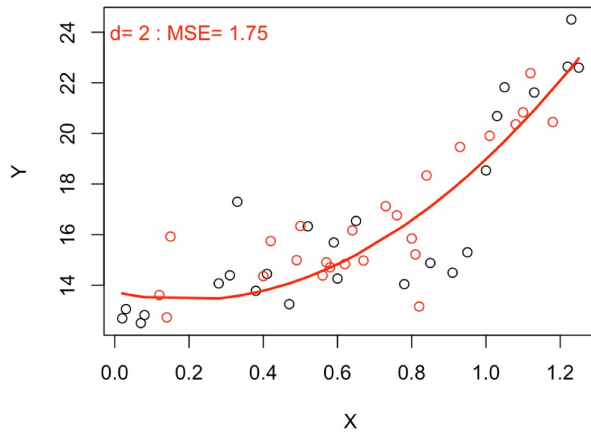
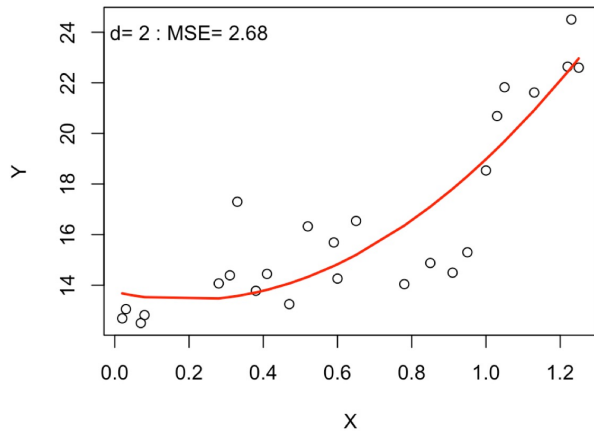


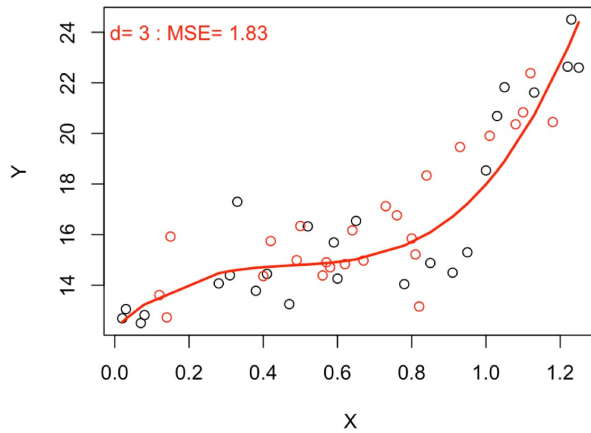
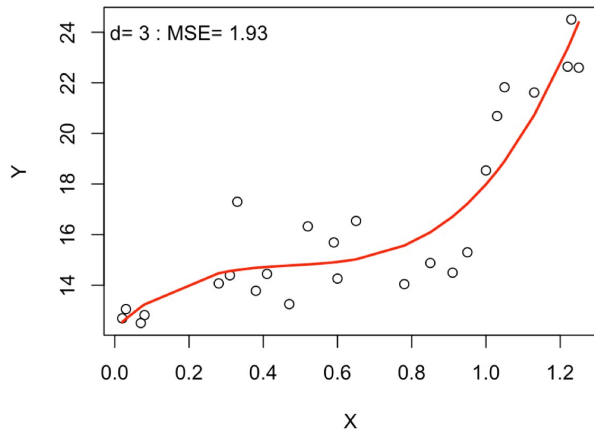
A.7 KEY CONCEPT

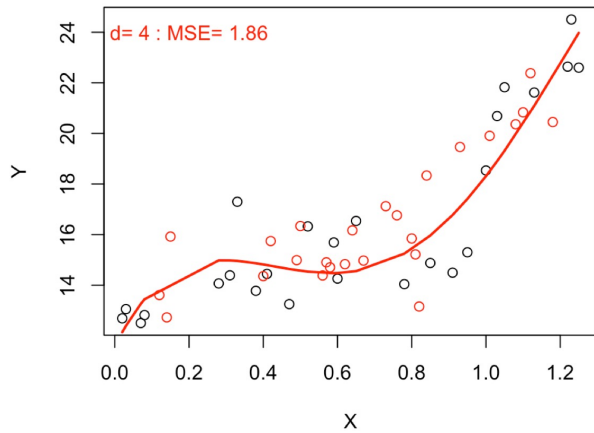
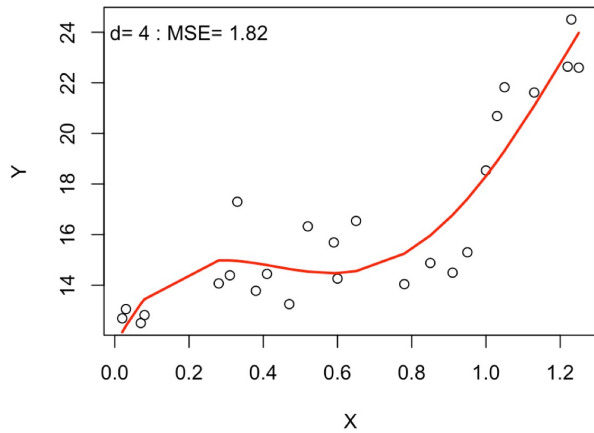
- Cross Validation
- Don't use all data when you are fitting a model
- Leave some out and use it for 'validation' and 'testing'.

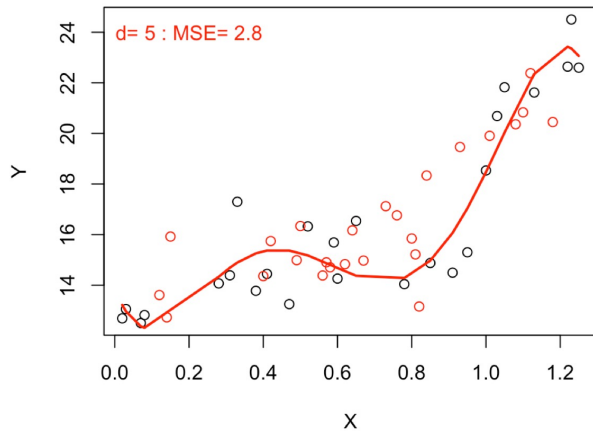
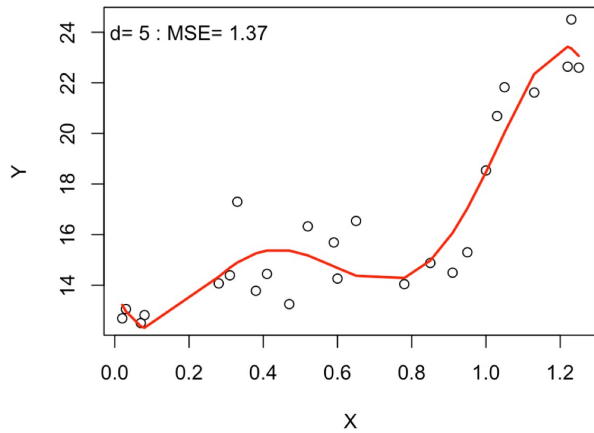
A.8 Leave-some-out Fitting Procedure 1

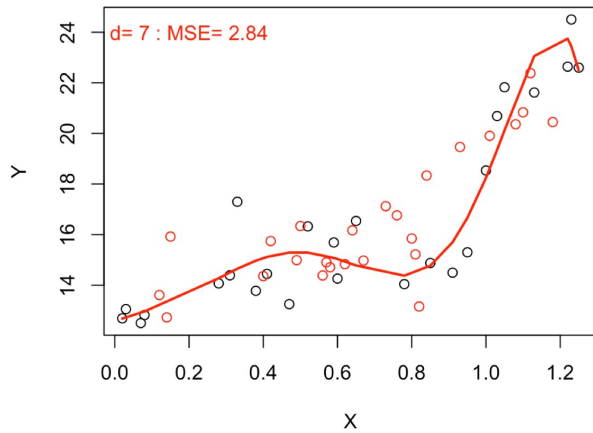
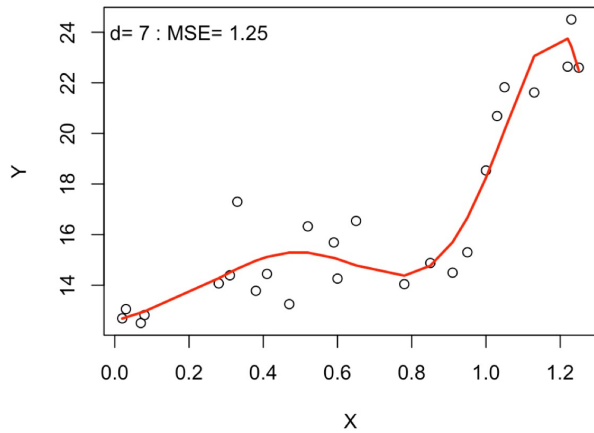


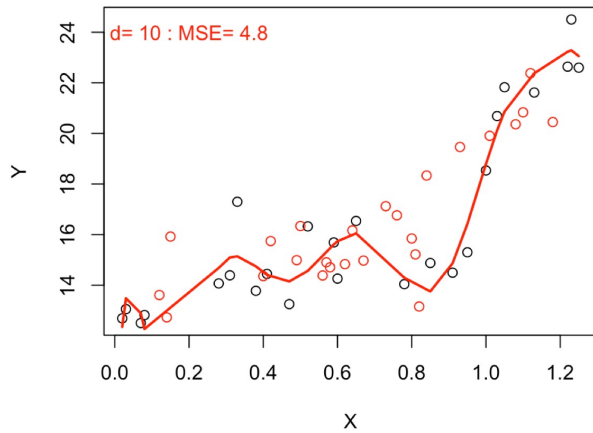
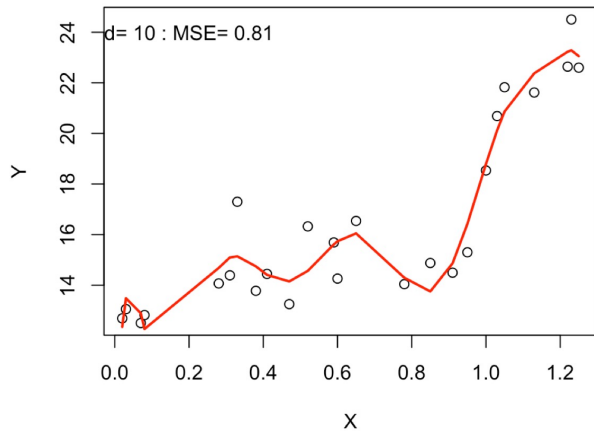


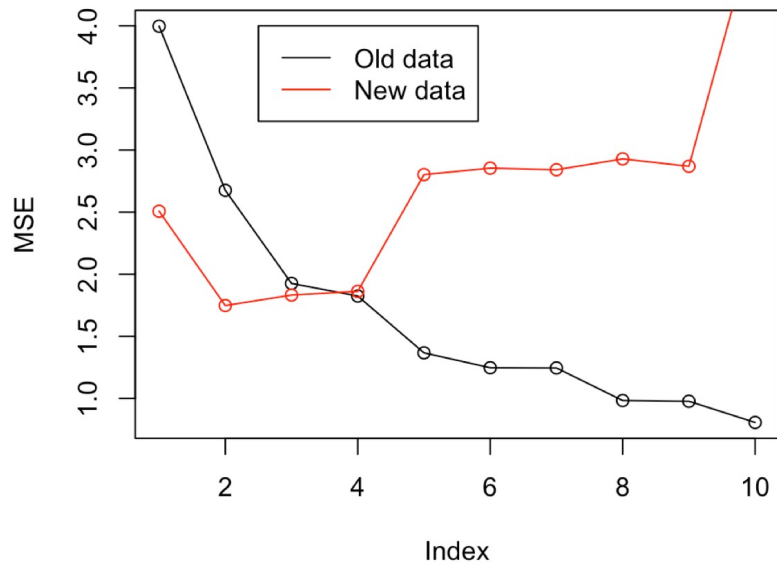






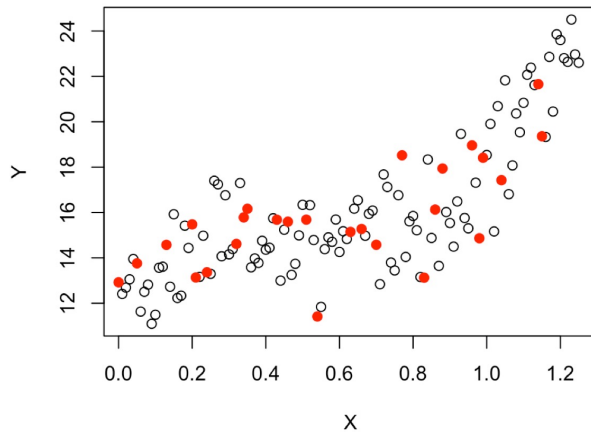
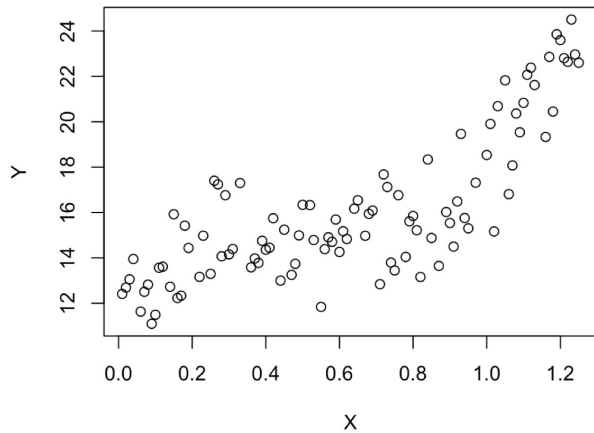


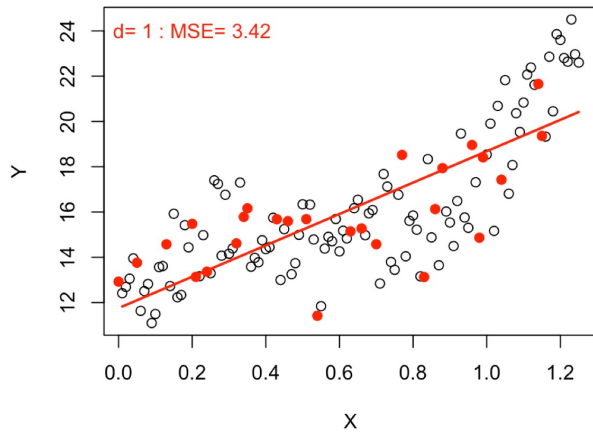
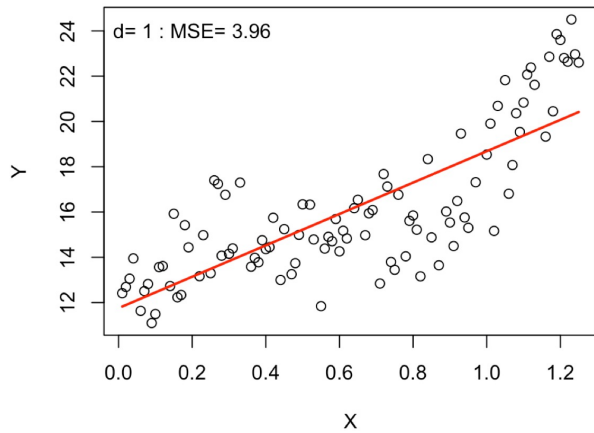


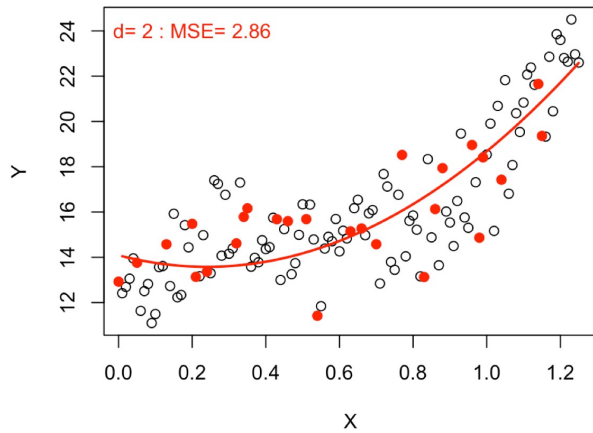
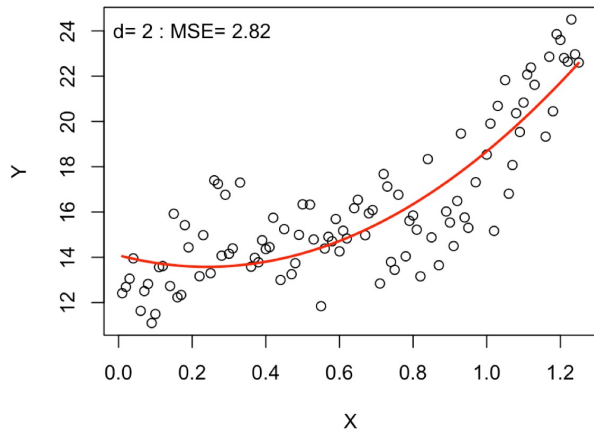


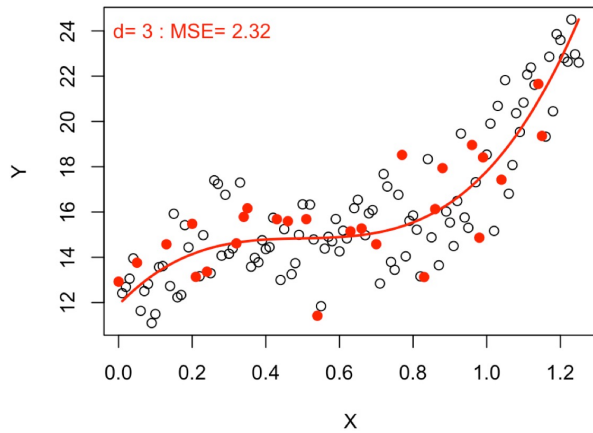
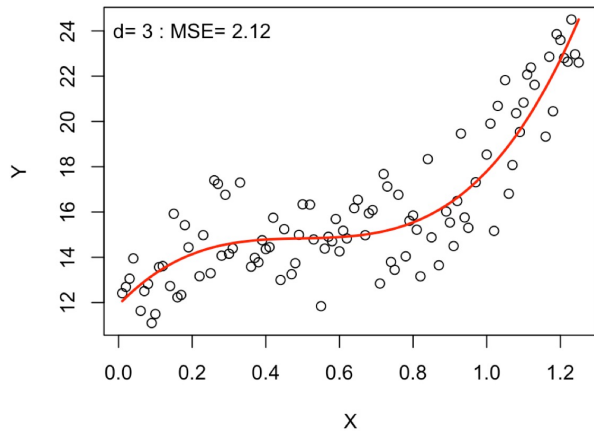
A.9 Leave-some-out Fitting 2

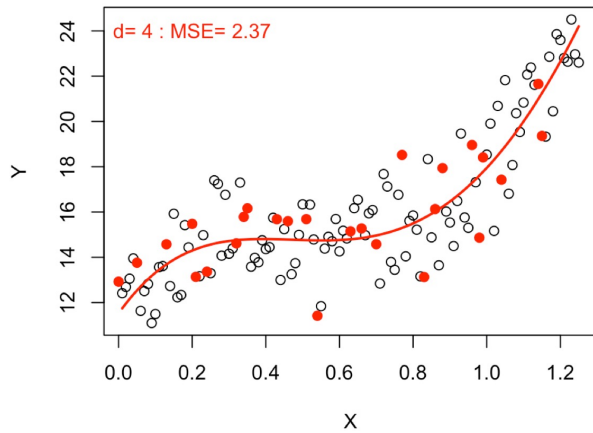
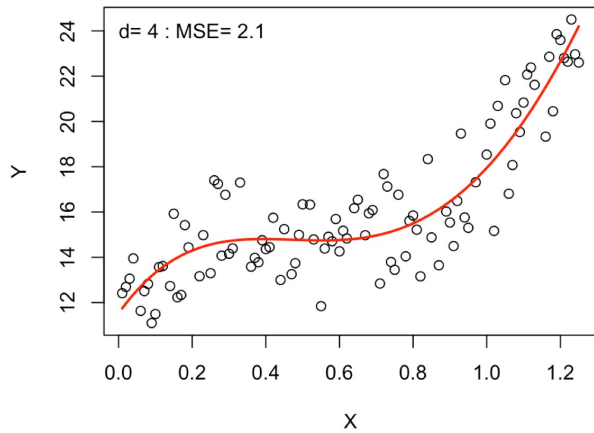
Larger dataset. $n = 100$ and $m = 26$.

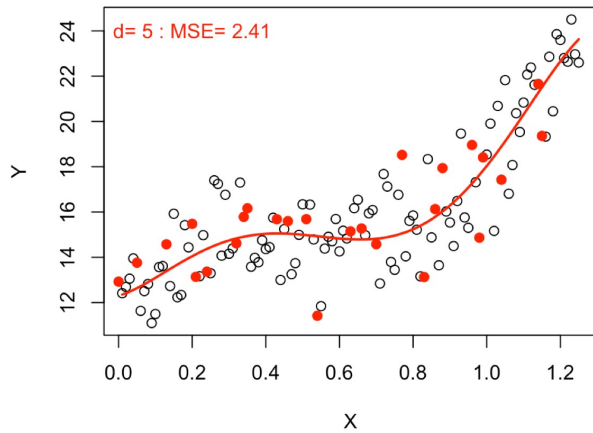
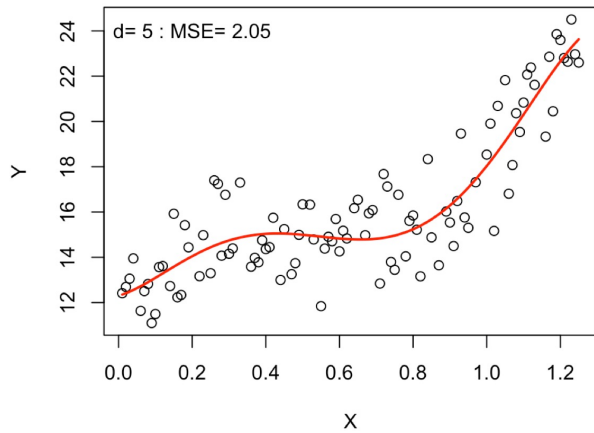


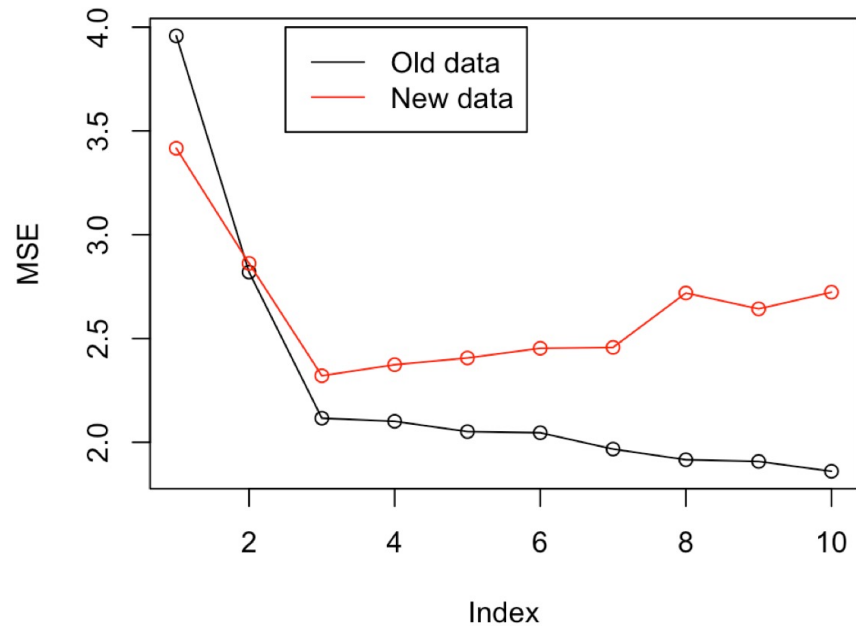












A.10 k-fold Cross Validation

- Training, Validation, and Testing Set
- Usually $k = 5$ or $k = 10$. We use $k = 5$ in this class.
- Randomly divide n into 6 groups.
- For example, if $n = 150$ and $k = 5$,

n=150

```
[-----Training Set 125-----]    [Test Set]
[fold 1][fold 2][fold 3][fold 4][fold 5]
[  25  ][  25  ][  25  ][  25  ][  25  ]    [   25   ]
```

A.11 k-fold Cross Validation

- Hyperparameter - parameter in the model that controls flexibility.
- e.g. Polynomial Regression $\rightarrow d$.
- Use Cross-Validation within the training set to tune the hyperparameter.
- Once you pick your hyperparameter, test the final model using [Training Set] to fit the model, and use [Test Set] to find testing MSE.
- Cross-Validation of the training set can be used repeatedly.
- Test Set should be used only once per method at the end.

A.12 CV within the Training Set

Let $k=4$

Data = [A,B,C,D,E,F,G,H,I,J,K,L,M,N,O,P,Q,R,S,T,U,V,W,X,Y,Z]

$n=26$

Training = $5 \times 4 = 20$ obs.

Test = 6 obs.

```
[ fold 1 ][ fold 2 ][ fold 3 ][ fold 4 ]   [   Test   ]  
[J G M T V][L B Q E I][N U X A P][Z K S D F]   [O C W Y H R]   # random assignment
```



```
CV rd 1 Training:      [fold 2][fold 3][fold 4] -> Training MSE
CV rd 1 Validation: [fold 1]                    -> Validation MSE
```

```
CV rd 2 Training:  [fold 1]          [fold 3][fold 4] -> Training MSE
CV rd 2 Validation:      [fold 2]                    -> Validation MSE
```

```
CV rd 3 Training:  [fold 1][fold 2]          [fold 4] -> Training MSE
CV rd 3 Validation:      [fold 3]                    -> Validation MSE
```

```
CV rd 4 Training:  [fold 1][fold 2][fold 3]          -> Training MSE
CV rd 4 Validation:      [fold 4]                    -> Validation MSE
```

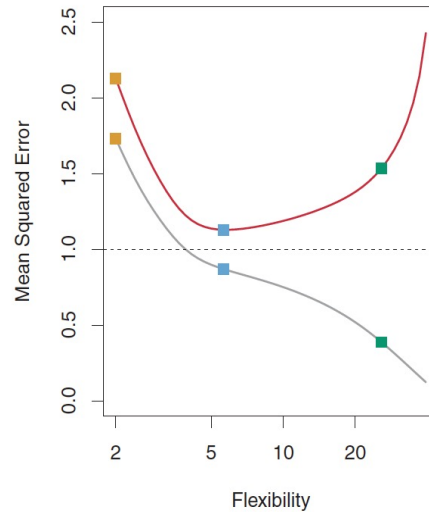
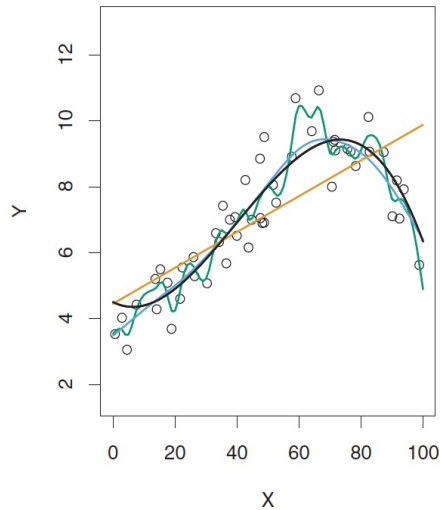
```
#--- Pick hyperparameter based on Average Validation MSE ---
```

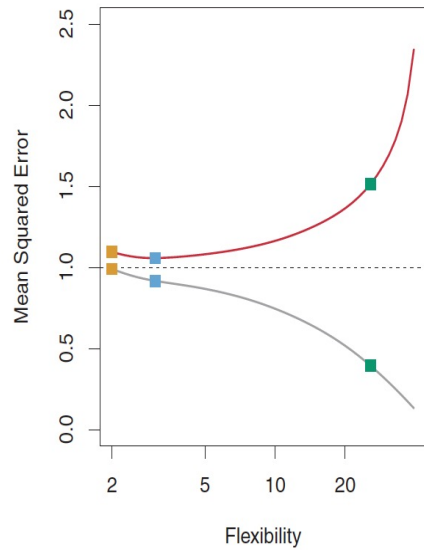
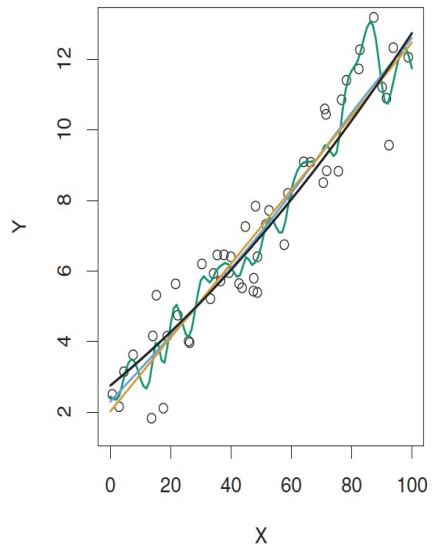
```
Final Fit Training: [  -- All Training Set ---  ]          -> Training MSE
Final Fit Test:      [  Test  ] -> Test MSE
```

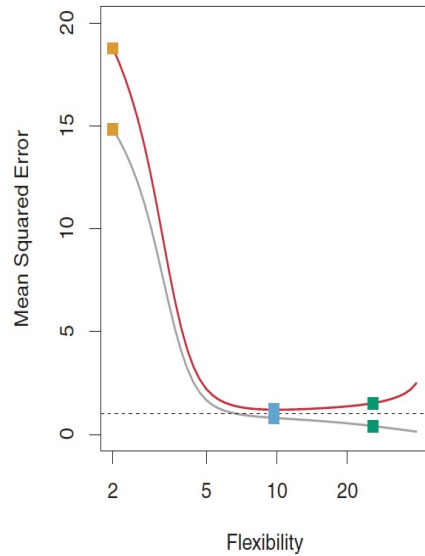
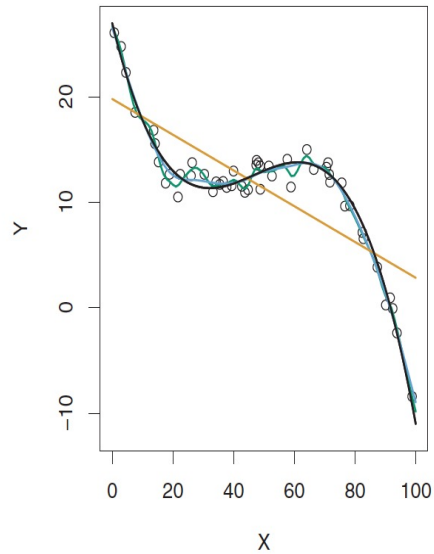
```
#--- Method should be compared to other Method using Test MSE
```


A.13 Training MSE vs Validation MSE

Also called in-sample vs out-sample







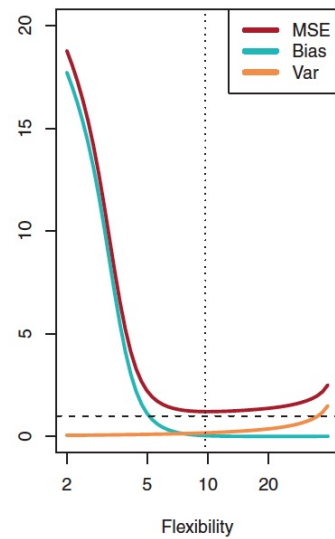
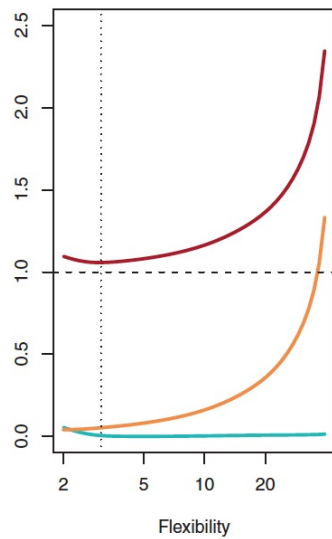
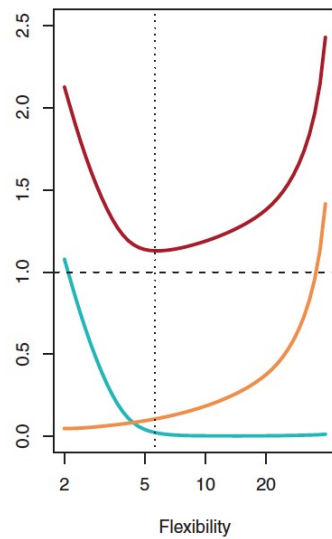
A.14 Assessing Model Prediction Accuracy

Assessing Model Prediction Accuracy

A.15 Bias-Variance Trade-Off

Prediction MSE can be decomposed as

$$\begin{aligned} E(Y - \hat{f}(X))^2 &= E\left(f(X) + \epsilon - \hat{f}(X)\right)^2 \\ &= E\left(f(X) - E(\hat{f}(X)) + E(\hat{f}(X)) - \hat{f}(X) + \epsilon\right)^2 \\ &= E\left(f(X) - E(\hat{f}(X))\right)^2 + E\left(E(\hat{f}(X)) - \hat{f}(X)\right)^2 + E(\epsilon^2) \\ &= Var(\hat{f}(X)) + Bias(\hat{f}(X))^2 + Var(\epsilon) \end{aligned}$$



A.16 Prediction MSE

-

$$E(Y - \hat{f}(X))^2 = Var(\hat{f}(X)) + Bias(\hat{f}(X))^2 + Var(\epsilon)$$

- can't have low variance and low bias
- has lower bound

A.17 In the Classification Setting

- Instead of MSE, work with Error Rate:

$$ER = \frac{1}{n} \sum_{i=1}^n I(y_i \neq \hat{y}_i)$$

A.18 Trade-off in the new approach

- Classical Statistics (Probabilistic Model)

$$Y = f(X) + \epsilon$$

- Assume parametric model for $f(\cdot)$ and ϵ .
- Sampling Probability of (y_1, \dots, y_n) , which are realizations of r.v. Y .
- Estimate parameters for $f(\cdot)$ and ϵ .
- Because the model distinguish the mechanism $f(\cdot)$ vs noise ϵ , looking at in-sample fit was enough (if the assumption is correct).
- Predict future Y using the estimated model.

- Pros and Cons
 - Model is interpretable.
 - Future effect of the model is easier to calculate.
 - No need for out-sample validation (test set), if assumption is correct.
 - Popular models are mathematically optimized already, to save the computational task.
 - Theory on prediction interval. Based on the assumption, often distribution on the prediction error is available.

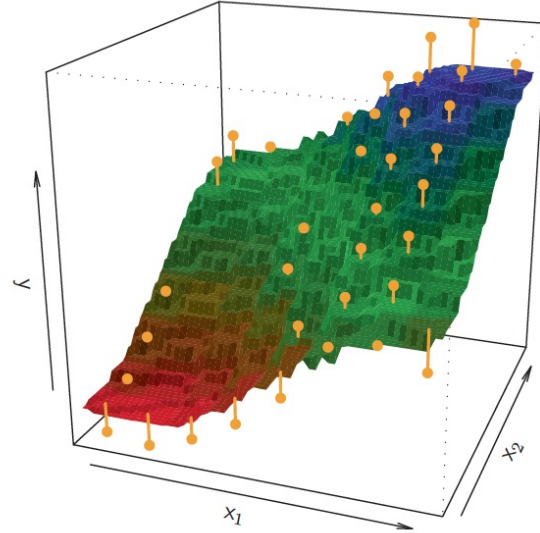
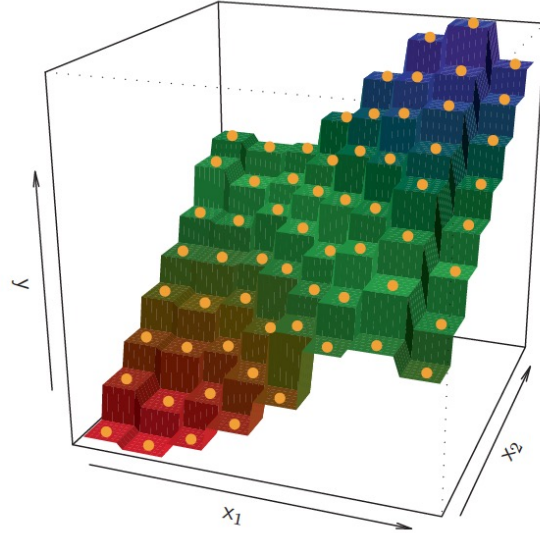
A.19 K-Nearest Neighbor

- One of elementary supervised learning model.
- Pick a point x_0 , find K nearest observations.
- $f(x_0)$ is estimated by the average of all K neighbors:

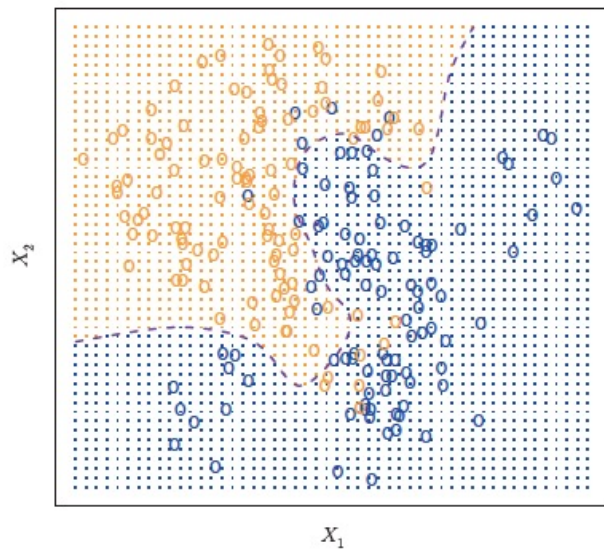
$$\hat{f}(x_0) = \frac{1}{K} \sum y_i.$$

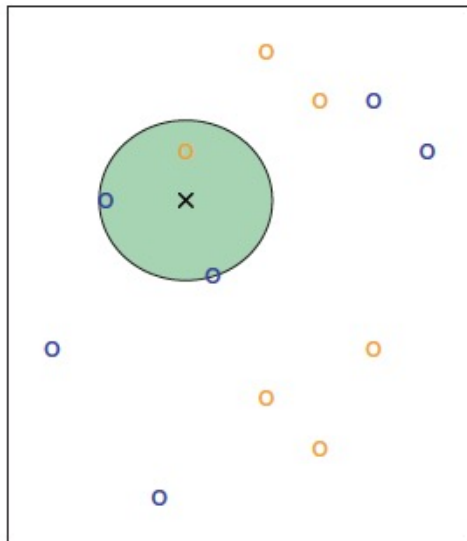
- K is the hyperparameter.
- Can be used for Regression or Classification

$K=1$ (left) and $K=9$ (right)

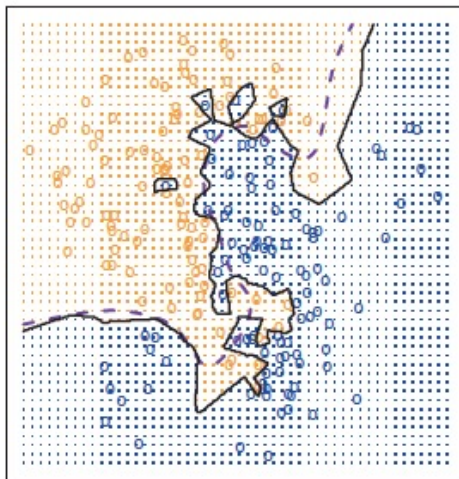


A.20 K-NN examples





KNN: $K=1$



KNN: $K=100$

