

# 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	KEY CONCEPT	.....
A.10	Hyperparameter	.....
A.11	k-fold Cross Validation	.....
A.12	k-fold Cross Validation	.....
A.13	k-fold Cross Validation	.....
A.14	5-fold CV	.....
A.15	Training MSE vs Validation MSE	.....
A.16	Final Test Fit	.....

- A.17 Bias-Variance Trade-Off . . . . .
- A.18 Training MSE vs Validation MSE . . . . .
- A.19 Prediction MSE . . . . .
- A.20 Assessing Model Prediction Accuracy . . . . .
- A.21 In the Classification Setting . . . . .
- A.22 Trade-off in the new approach . . . . .
- A.23 K-Nearest Neighbor . . . . .
- A.24 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

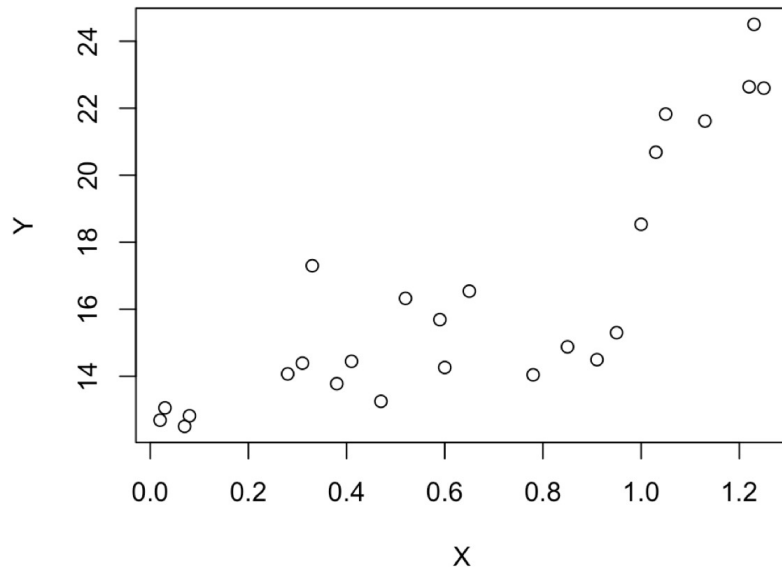
- Want to guess the next  $Y$  as accurate as possible

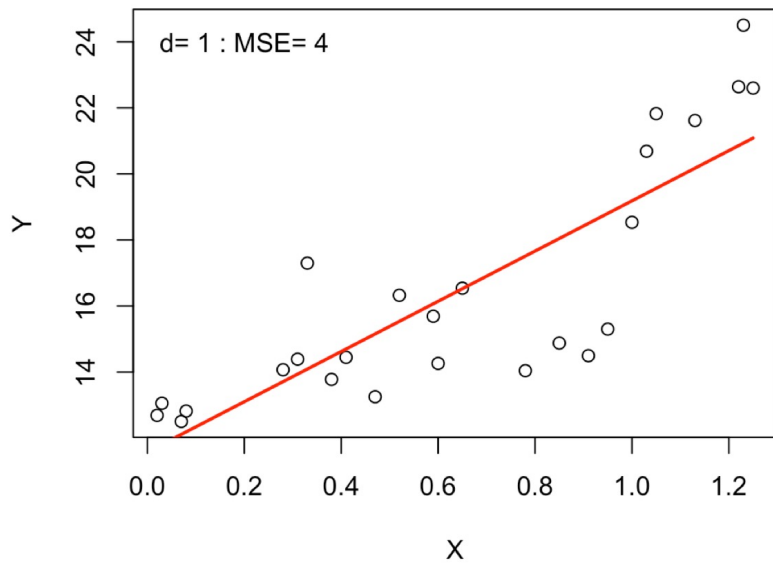
$$\hat{Y} = \hat{f}(X)$$

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

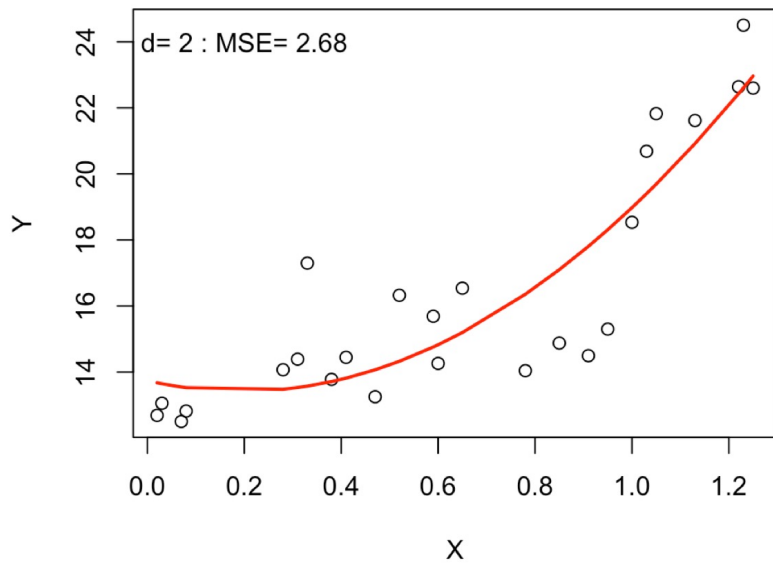
$$MSE = E(Y - \hat{Y})^2 = E(Y - \hat{f}(X))^2$$

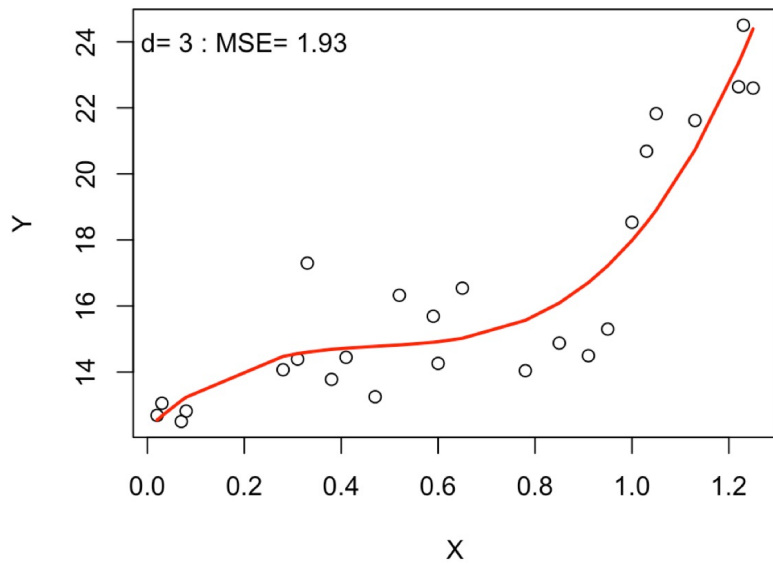
## A.4 Polynomial Regression 1

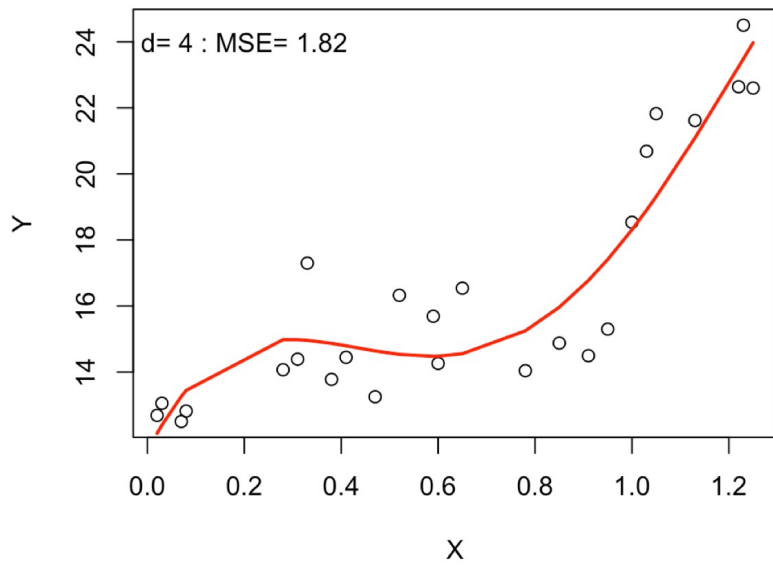


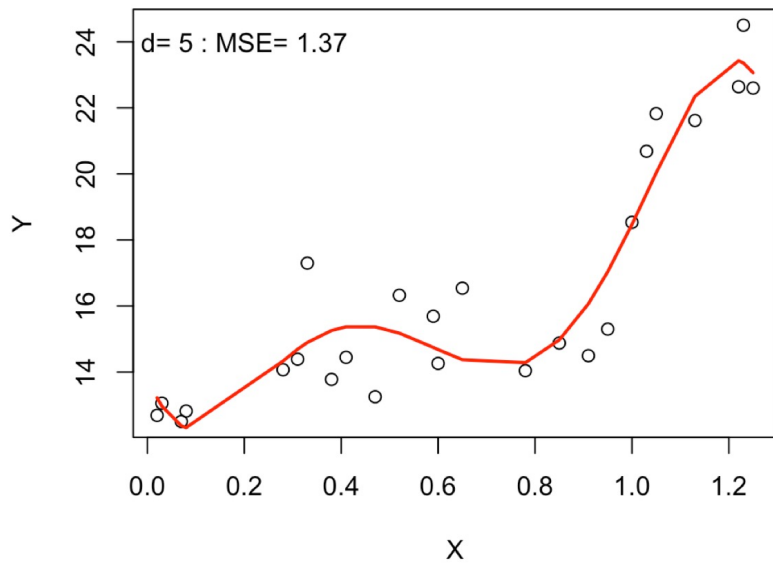


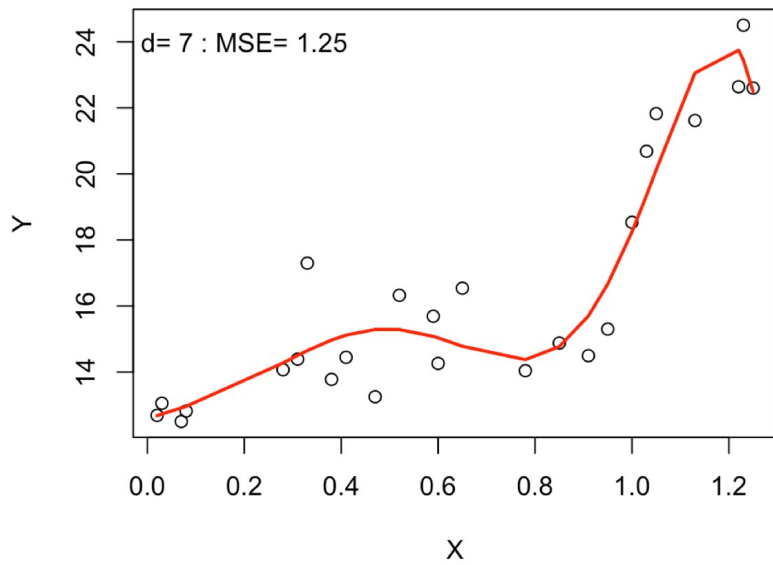


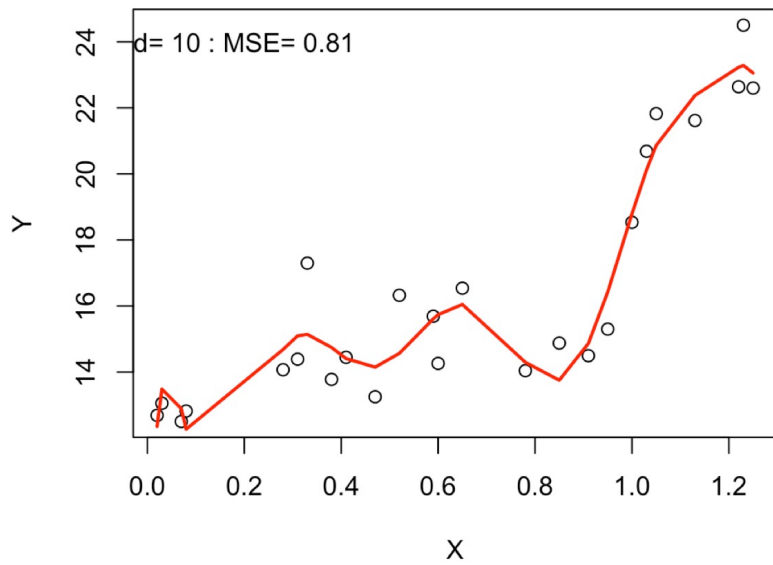


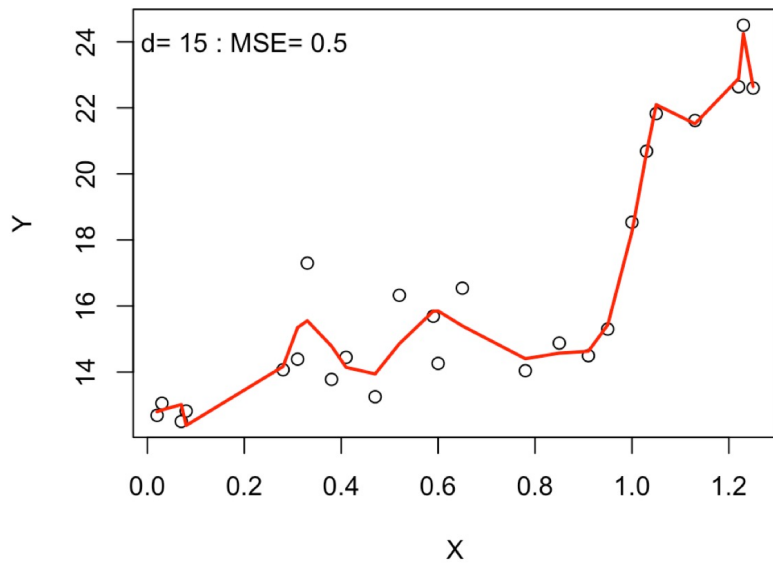


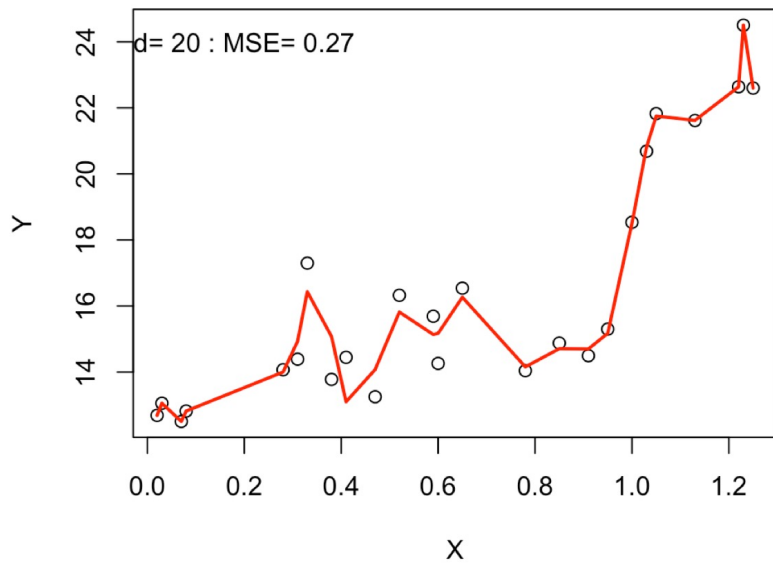














## A.5 Problem

- 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)

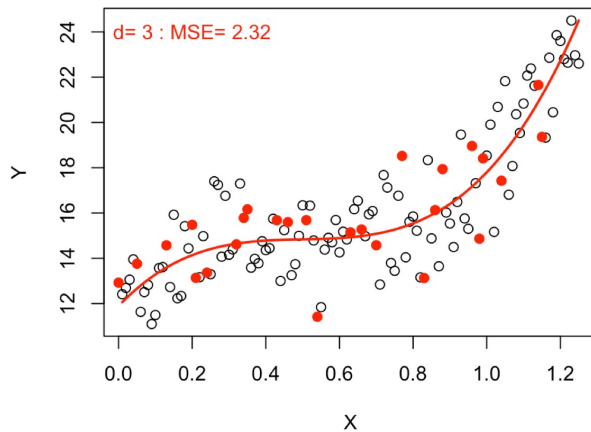
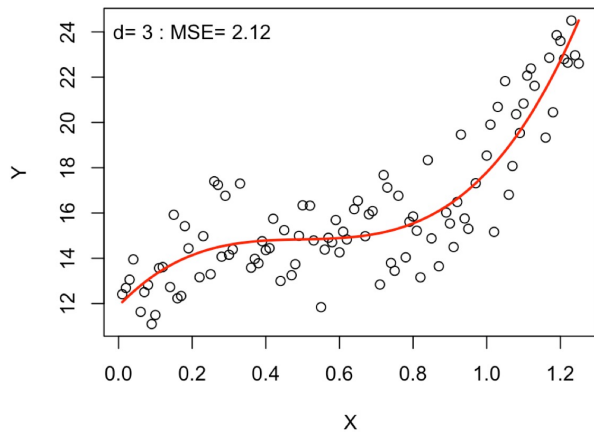
$$\text{MSE}_{tr} = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{f}(x_i))^2$$

- But we want minimum Prediction MSE

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

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

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

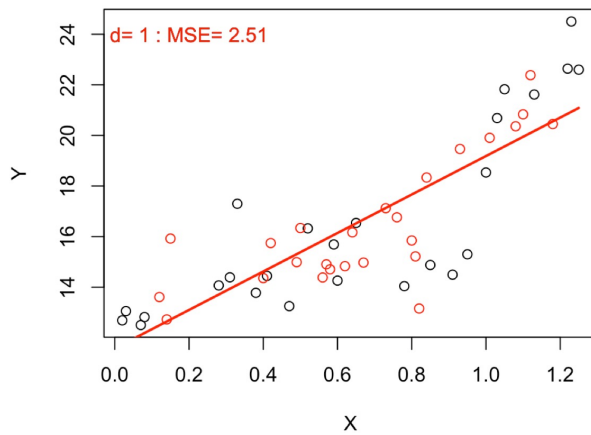
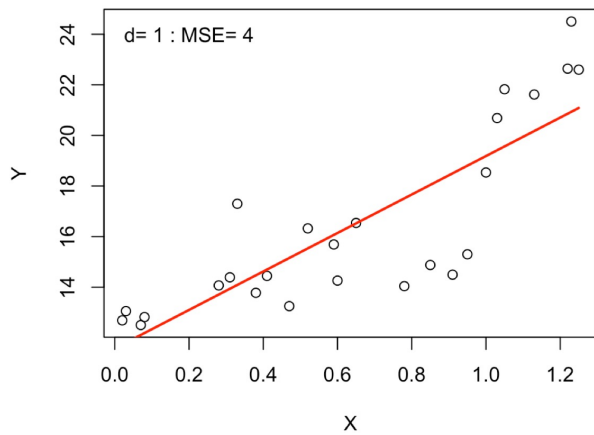


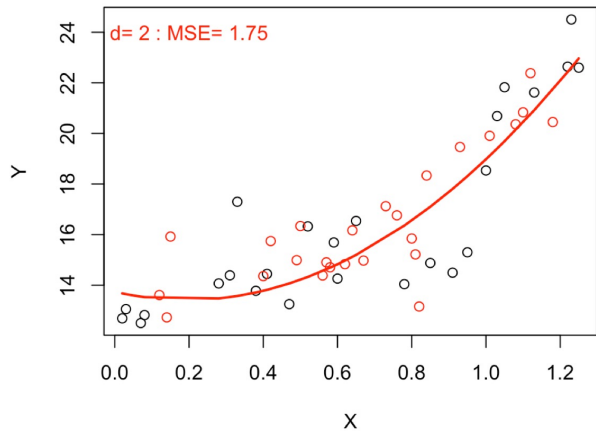
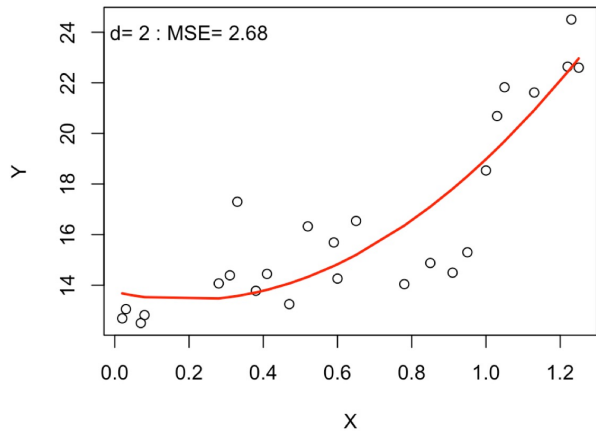


## A.7 KEY CONCEPT

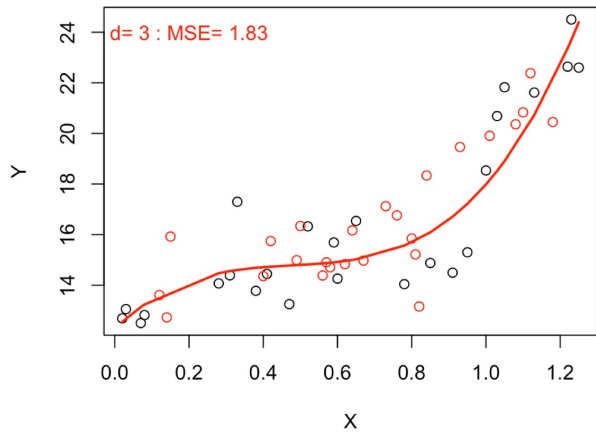
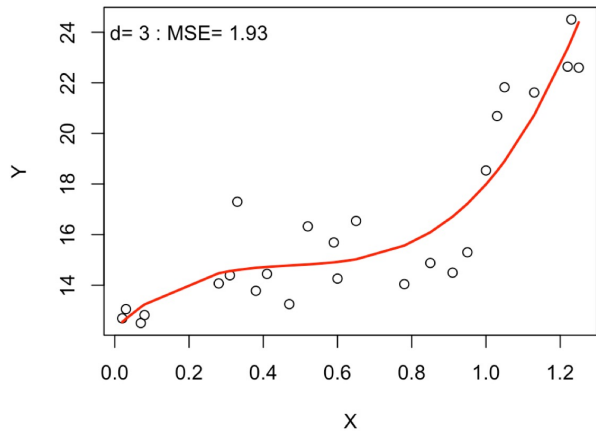
- Cross Validation
- Don't use all data when you are fitting a model
- Leave some out and use it for 'validation' and 'testing'.

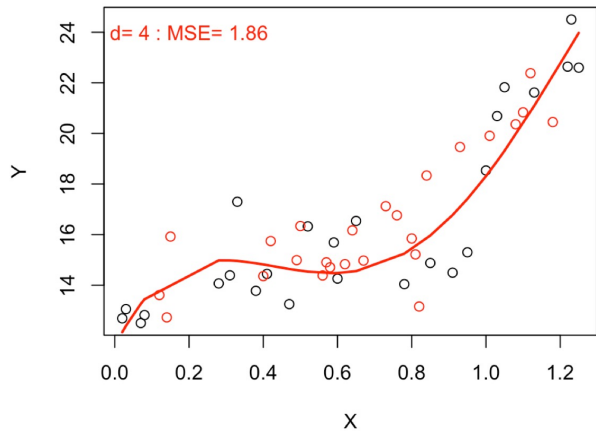
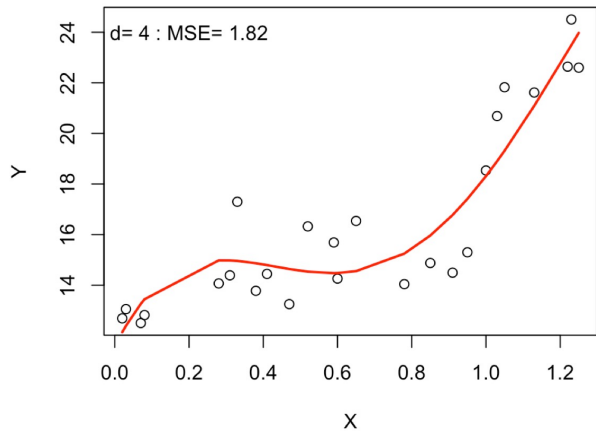
## A.8 Leave-some-out Fitting Procedure 1

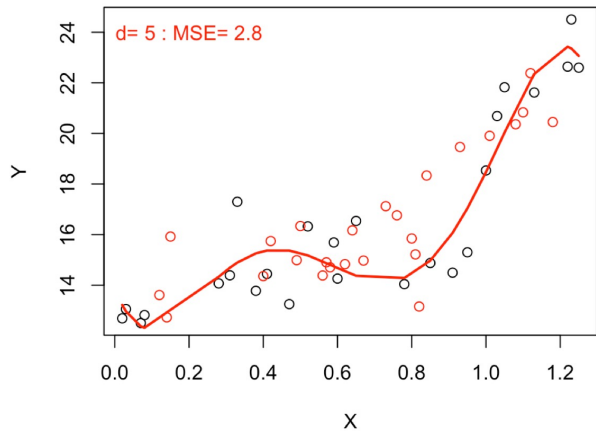
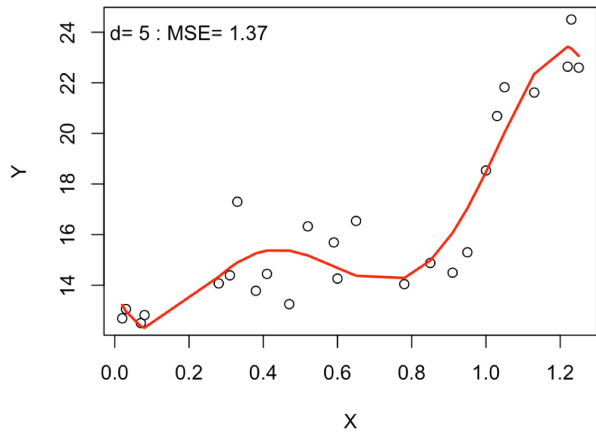


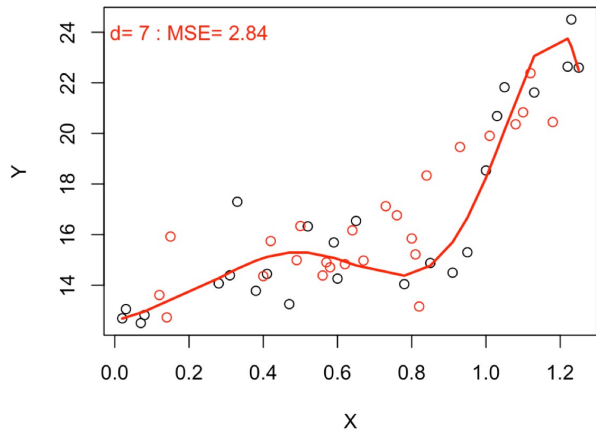
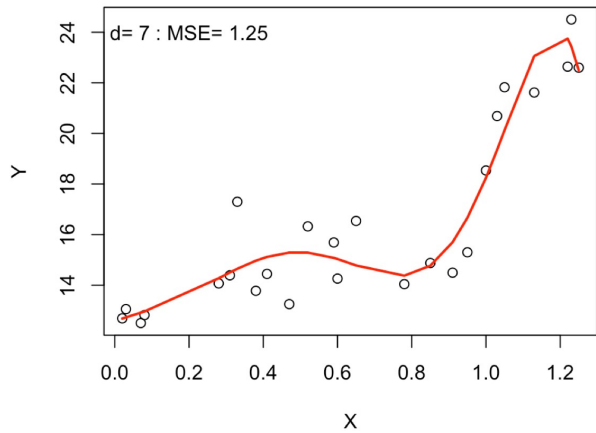


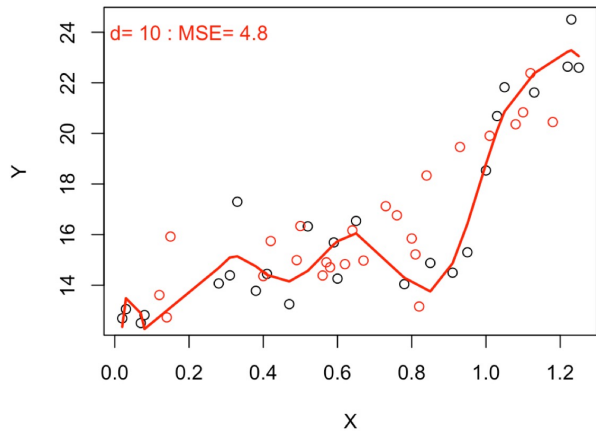
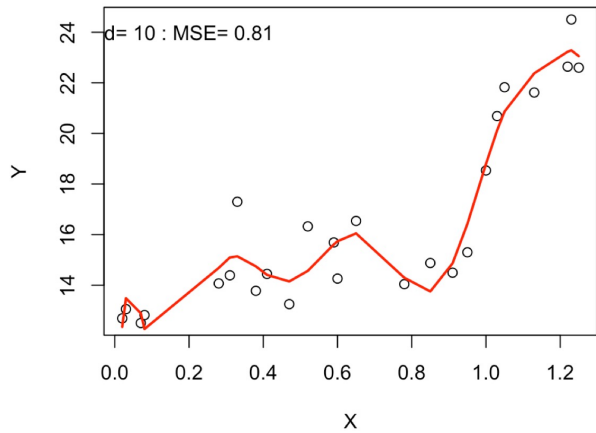


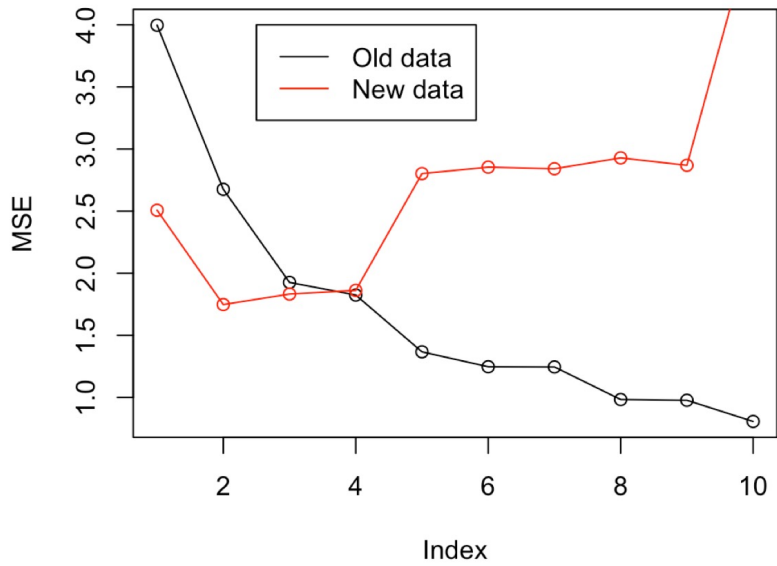












## A.9 KEY CONCEPT

- Divide data into

[Training] vs [Testing]

[Training] vs [Validation]

[In-sample] vs [Out-sample]

- Fit the model using [Training] data
- See if the model actually fit [Testing] data (the model hasn't seen these observations yet)
- Has one problem:

## A.10 Hyperparameter

- 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.
- Tuning Set, Training Set, Validation Set, and Testing Set



## A.11 k-fold Cross Validation

- Usually  $k = 5$  or  $k = 10$ . We use  $k = 5$  in this class.
- Randomly assign data into  $k + 1$  groups.
- For example, if  $n = 155$  and  $k = 5$ ,

n=155

```
[-----Tuning Set 125-----] [Test Set]
[fold 1][fold 2][fold 3][fold 4][fold 5]
[ 25  ][ 25  ][ 25  ][ 25  ][ 25  ] [ 30  ]
```

## A.12 k-fold Cross Validation

- Round 1  
[-----Training Set 100-----] [validation set 25]  
[fold 2][fold 3][fold 4][fold 5] [fold 1]  
[ 25 ][ 25 ][ 25 ][ 25 ] [ 25 ]
- Round 2  
[-----Training Set 100-----] [validation set 25]  
[fold 1][fold 3][fold 4][fold 5] [fold 2]  
[ 25 ][ 25 ][ 25 ][ 25 ] [ 25 ]
- Round 3  
[-----Training Set 100-----] [validation set 25]

```
[fold 1][fold 2][fold 4][fold 5] [fold 3]
[ 25 ][ 25 ][ 25 ][ 25 ] [ 25 ]
```

- Round 4

```
[-----Training Set 100-----] [validation set 25]
[fold 1][fold 2][fold 3][fold 5] [fold 4]
[ 25 ][ 25 ][ 25 ][ 25 ] [ 25 ]
```

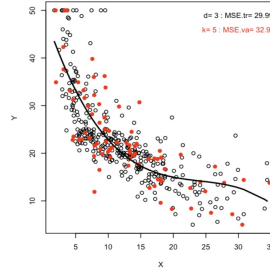
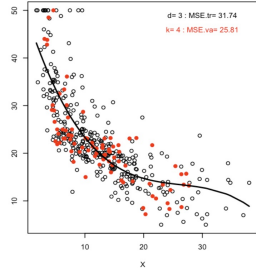
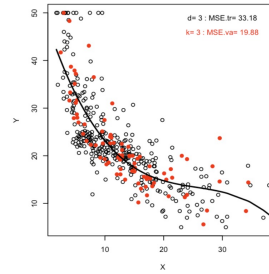
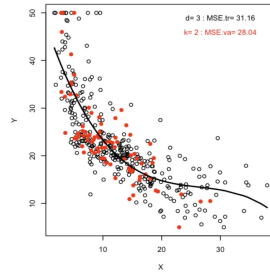
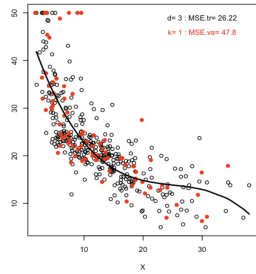
- Round 5

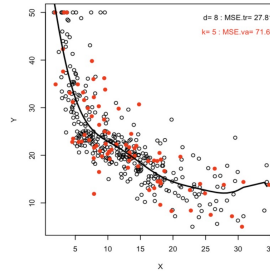
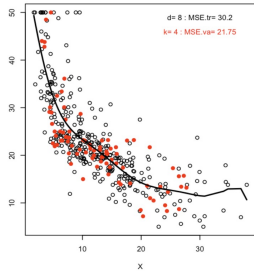
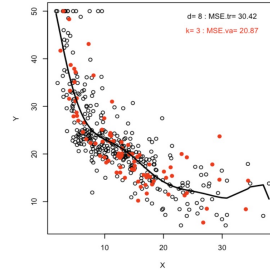
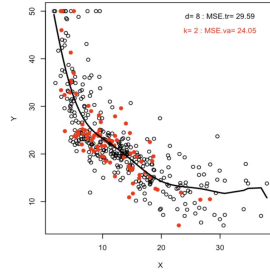
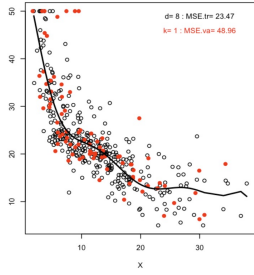
```
[-----Training Set 100-----] [validation set 25]
[fold 1][fold 2][fold 3][fold 4] [fold 5]
[ 25 ][ 25 ][ 25 ][ 25 ] [ 25 ]
```

## A.13 k-fold Cross Validation

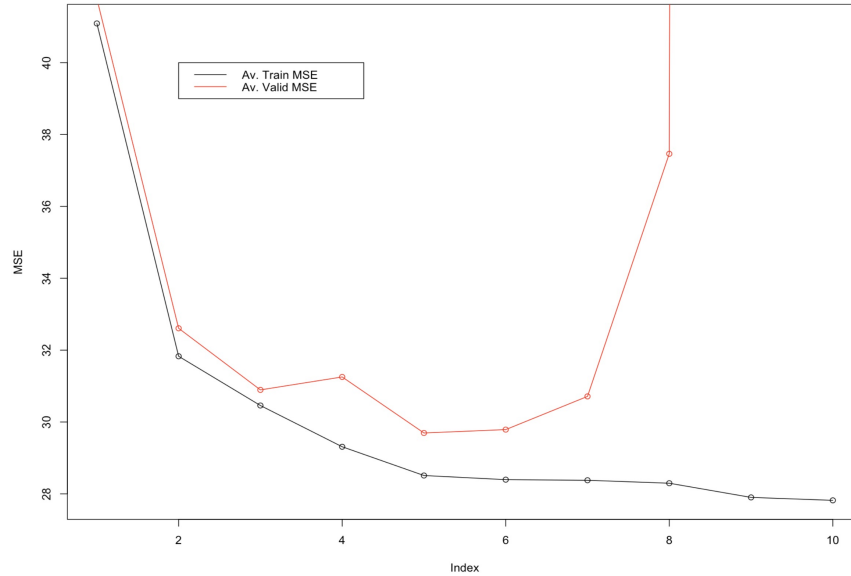
- Take [Tuning Set], and use 5-fold CV and fit 5 times using 5 different [Training Set] and [Validation Set]
- Use average validation MSE to decide on the best value of the hyperparameter.
- Now use the chosen hyperparameter, and fit entire [Tuning Set]. Then test it on [Test Set].
- Test Set should be used only once per method.
- Test MSE is the measure of performance for the method.

## A.14 5-fold CV

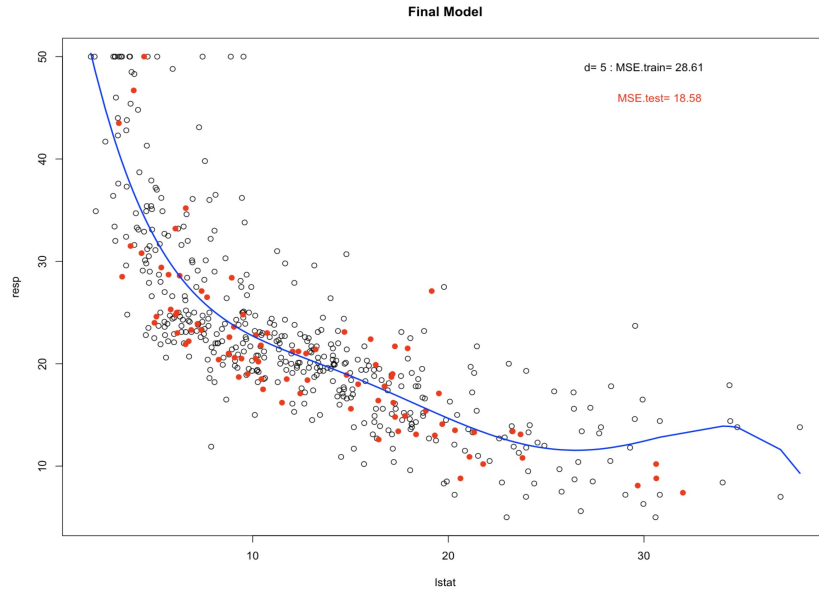




## A.15 Training MSE vs Validation MSE



## A.16 Final Test Fit



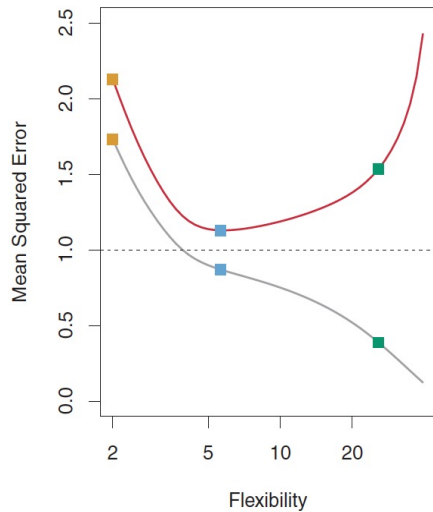
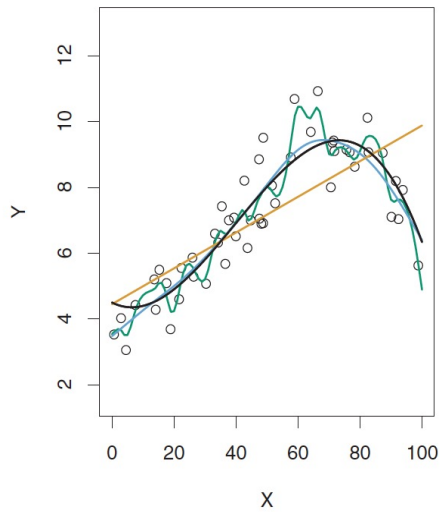


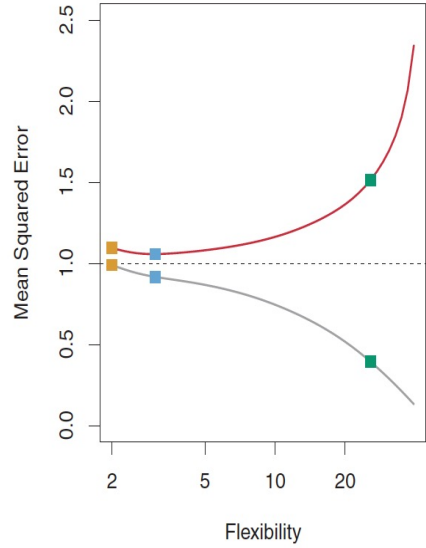
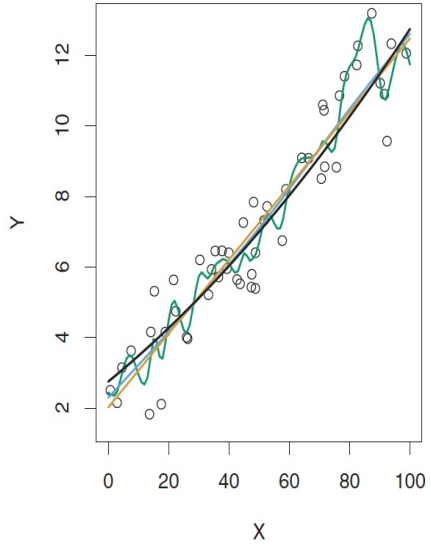


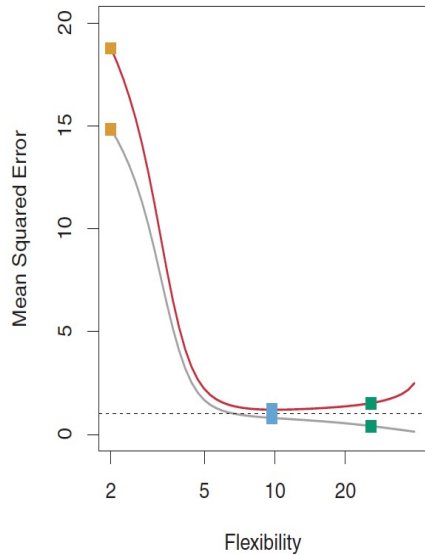
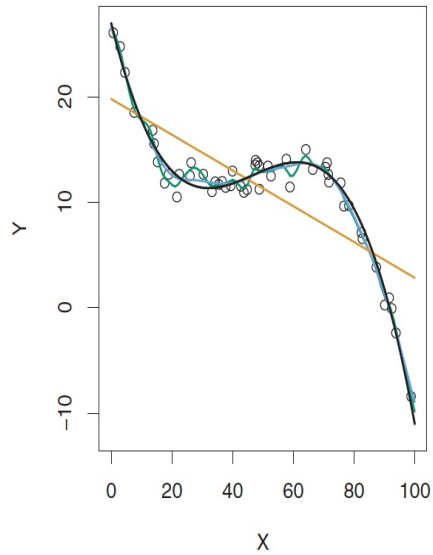
## A.17 Bias-Variance Trade-Off

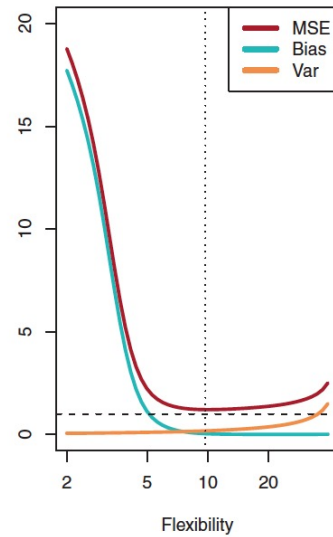
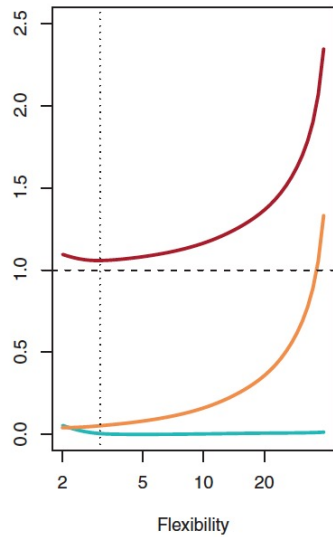
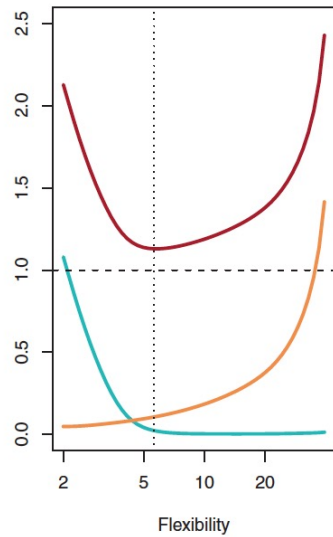
Prediction MSE can be decomposed as

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









## A.18 Prediction MSE

- 

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

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

## A.19 Assessing Model Prediction Accuracy

- If the model is fitting well, your

[Av. Training MSE]

[AV. Validation MSE]

[Test MSE]

should be all comparable.

- Your [Test MSE] is the best estimate for the true Prediction MSE.



## A.20 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.21 Trade-off in the new approach

- Classical Statistics (Probabilistic Model)

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

- Assume parametric model for  $f(\cdot)$  and  $\epsilon$ .
- Sampling Probability of  $(y_1, \dots, y_n)$ , which are realizations of r.v.  $Y$ .
- Estimate parameters for  $f(\cdot)$  and  $\epsilon$ .
- Because the model distinguish the mechanism  $f(\cdot)$  vs noise  $\epsilon$ , 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.

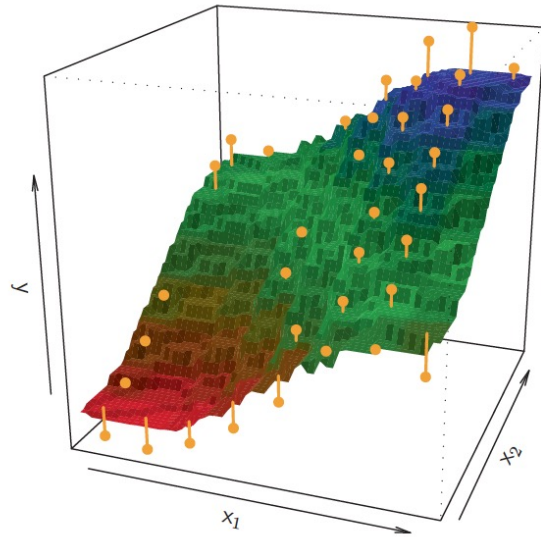
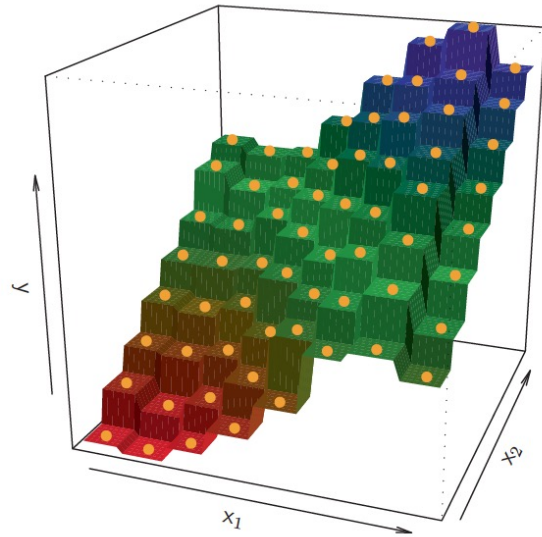
## A.22 K-Nearest Neighbor

- 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:

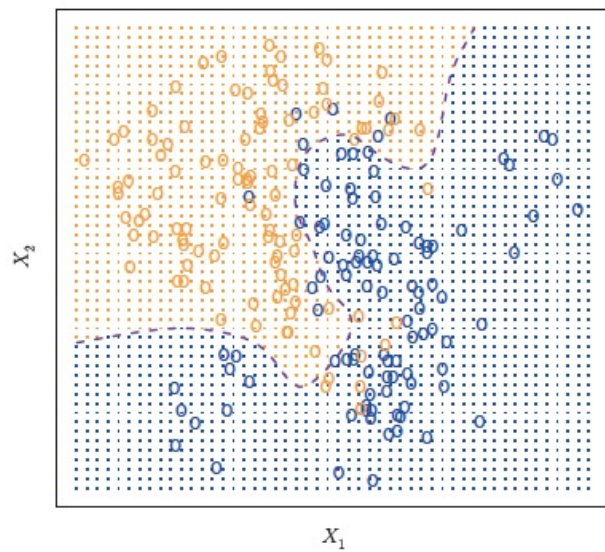
$$\hat{f}(x_0) = \frac{1}{K} \sum y_i.$$

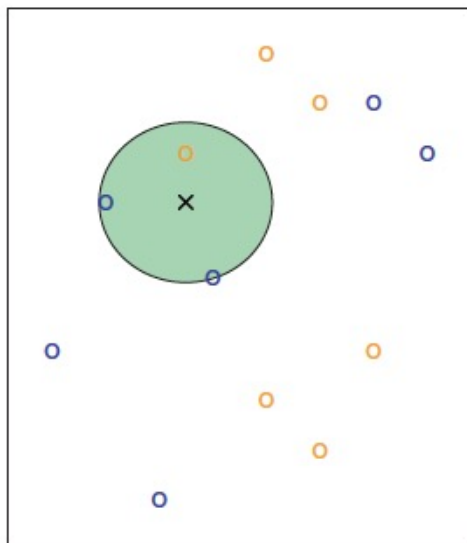
- $K$  is the hyperparameter.
- Can be used for Regression or Classification

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

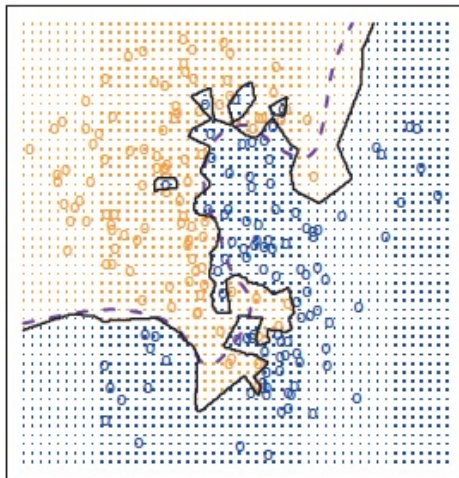


## A.23 K-NN examples





KNN: K=1



KNN: K=100

