

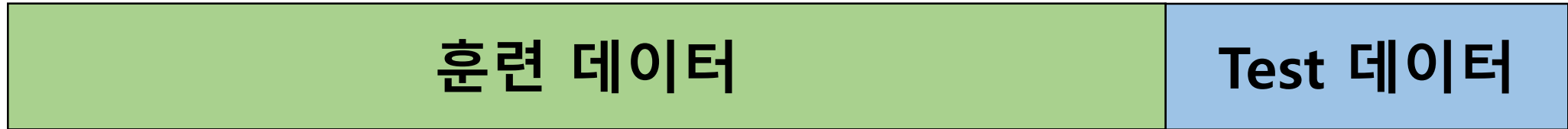
Resampling

Introduction to Statistical Learning

황성원

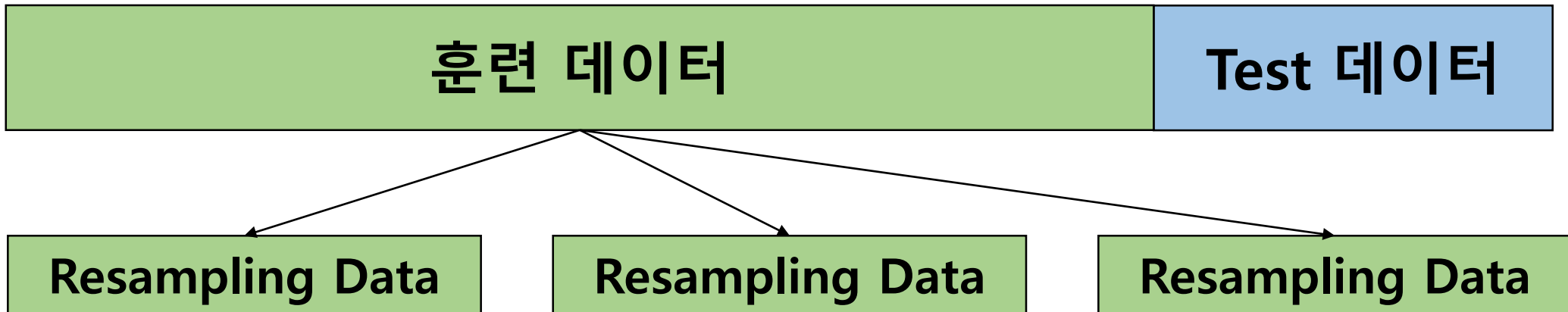
Resampling?

- 기존 방법 -



훈련데이터 → Fitting Model → Test (한 번의 Training 과정)

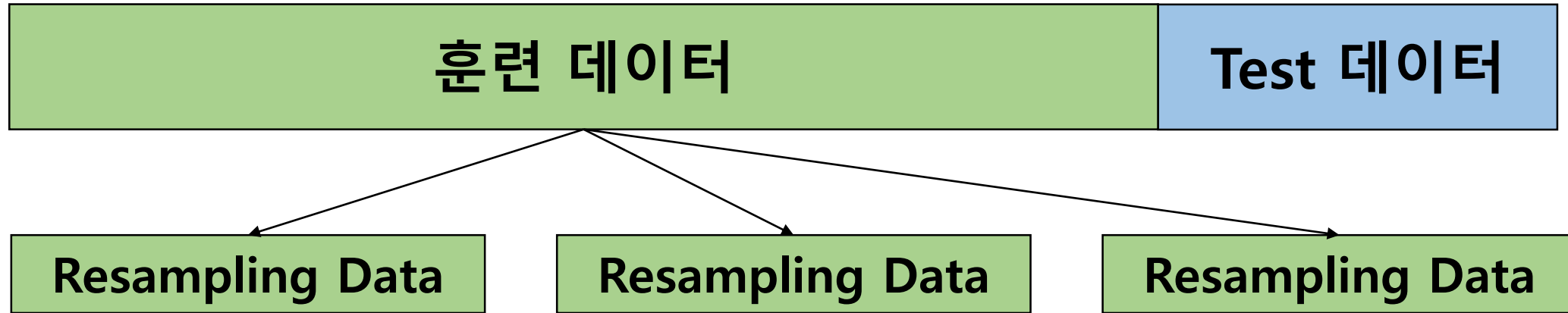
- Resampling 방법 -



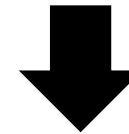
훈련데이터 → Resampling → Fitting Model → Test (여러 번의 Training 과정)_{SW}

Resampling?

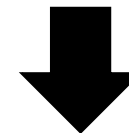
- Resampling 방법 -



훈련데이터 → Resampling → Fitting Model → Test (여러 번의 Training 과정)

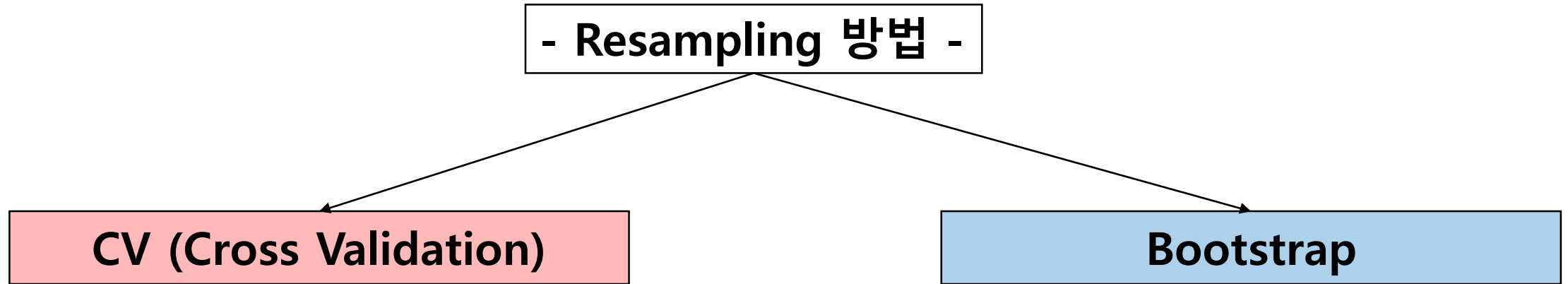


계산 부담 증가



컴퓨터 파워 증가로 해결!

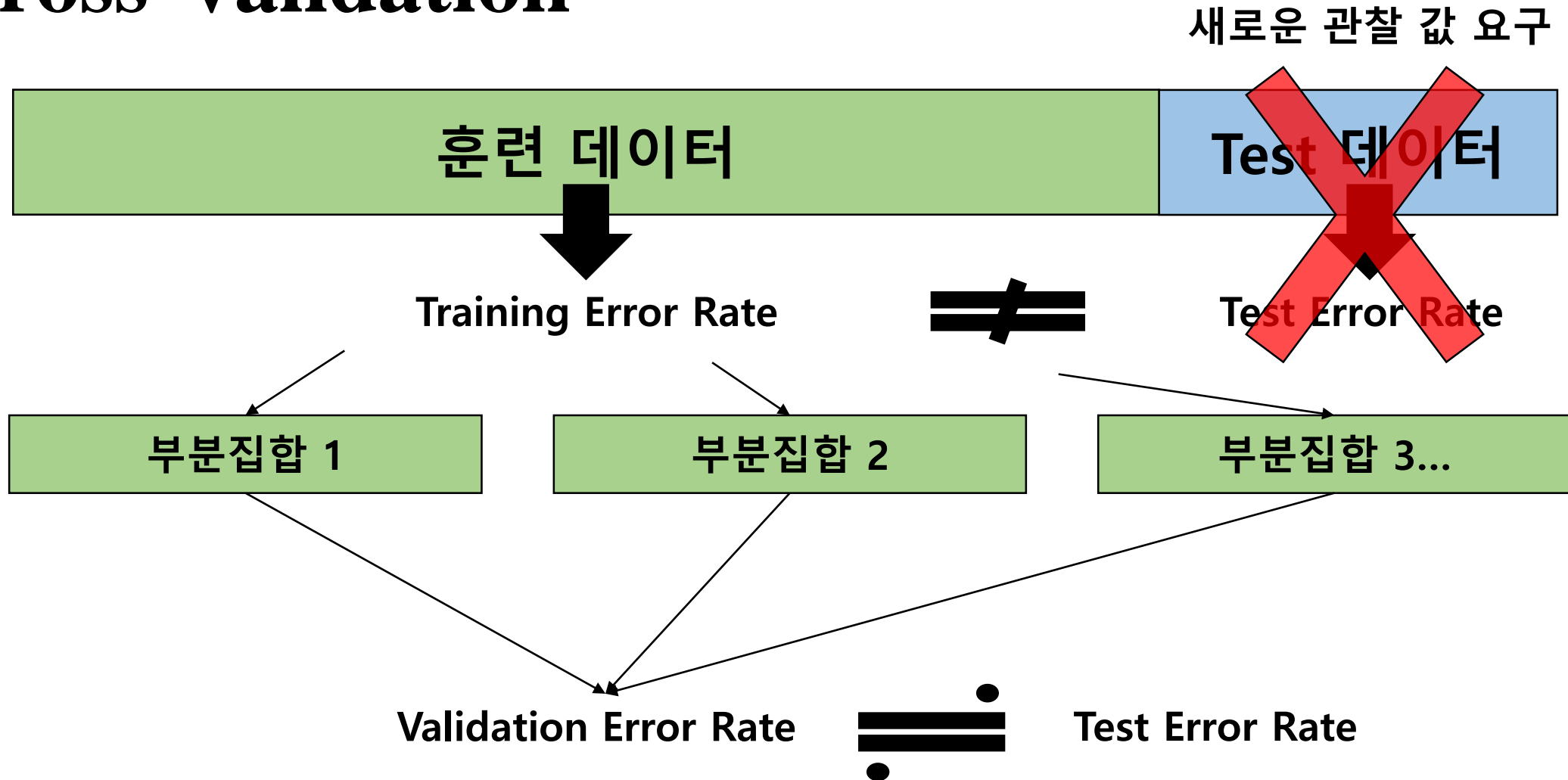
Resampling?



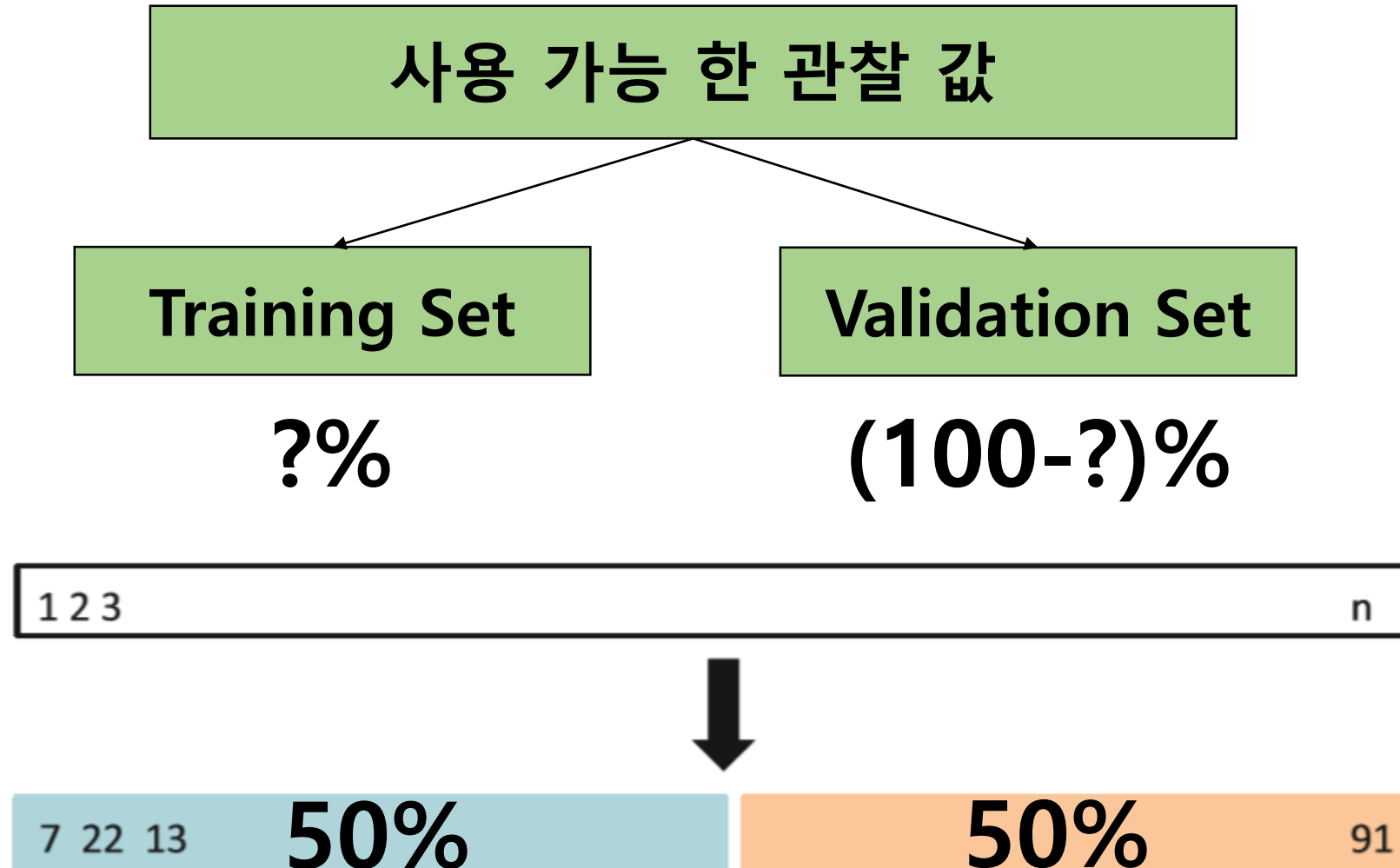
모델의 Performance를 평가 → Model Assessment

모델의 Flexibility를 선택 → Model Selection

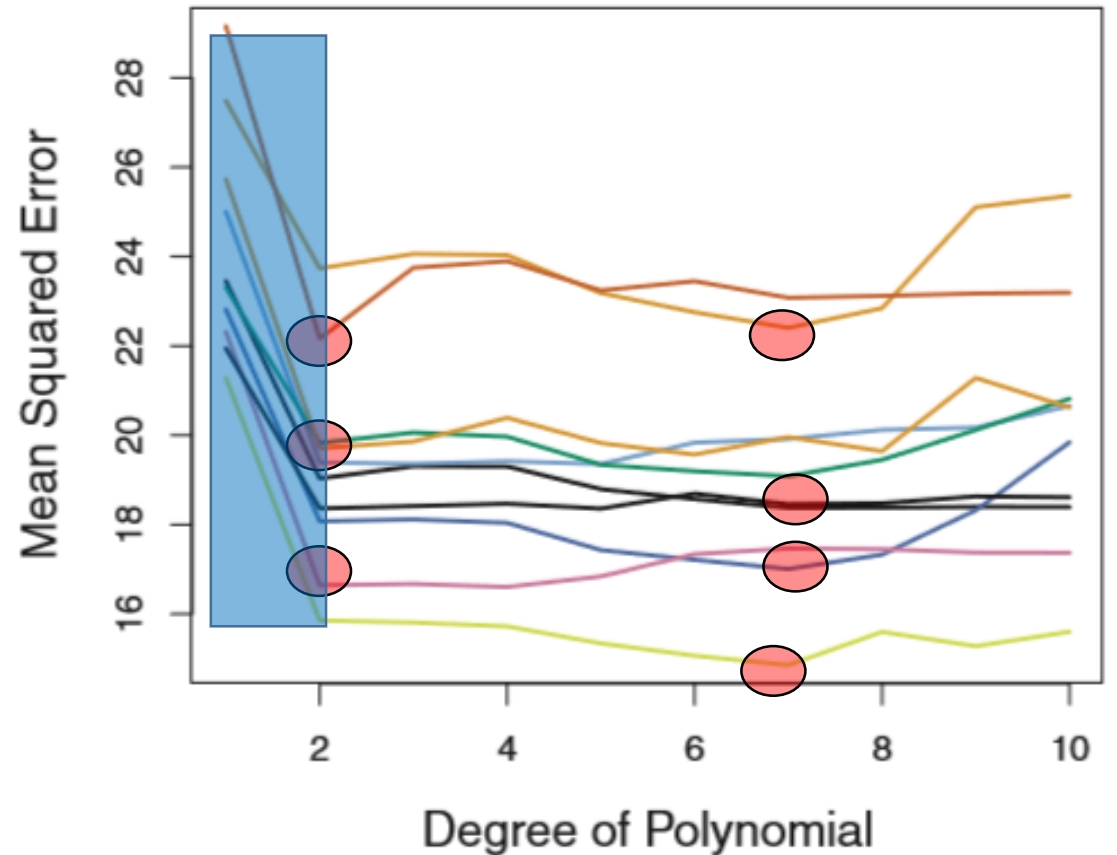
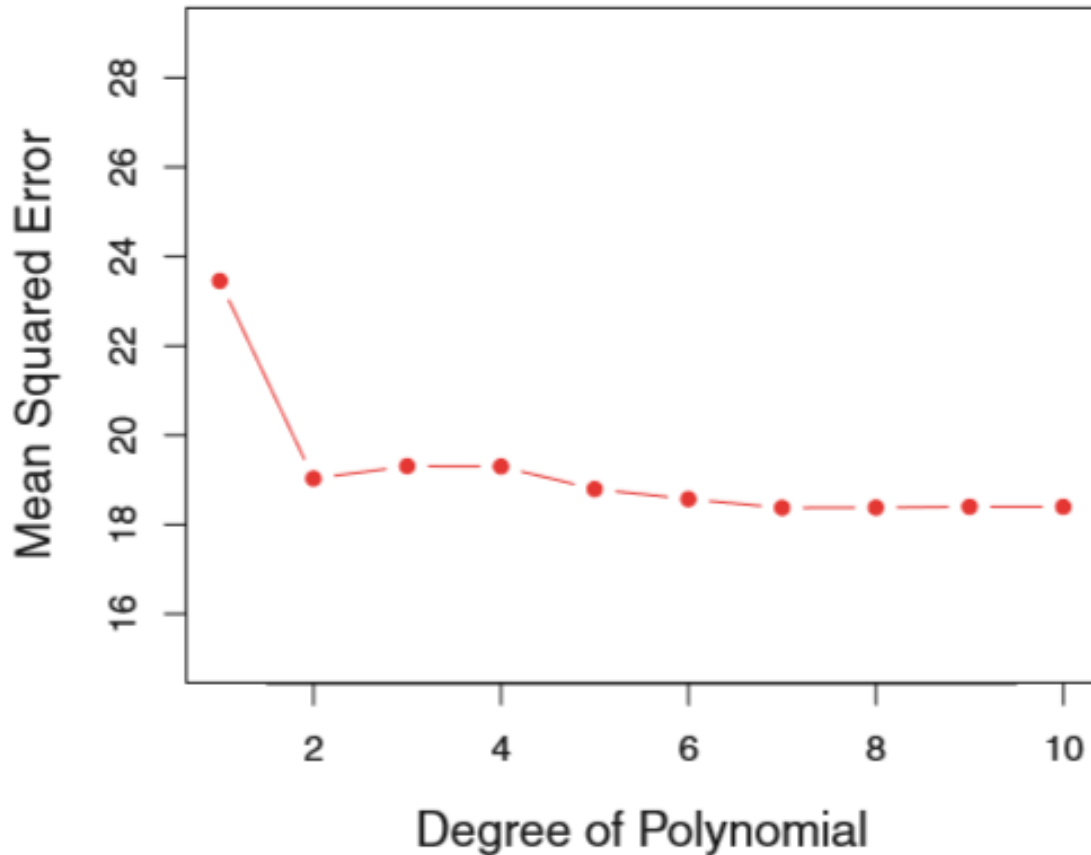
Cross-Validation



Validation Set Approach

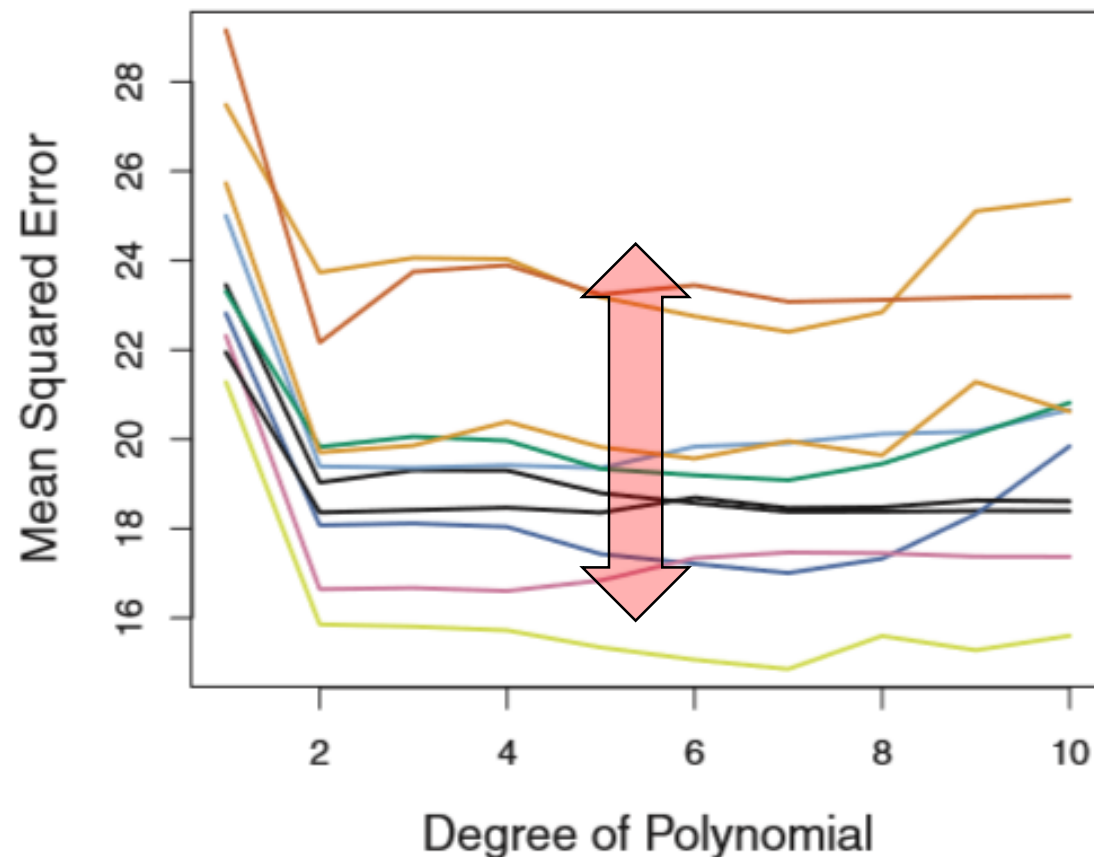
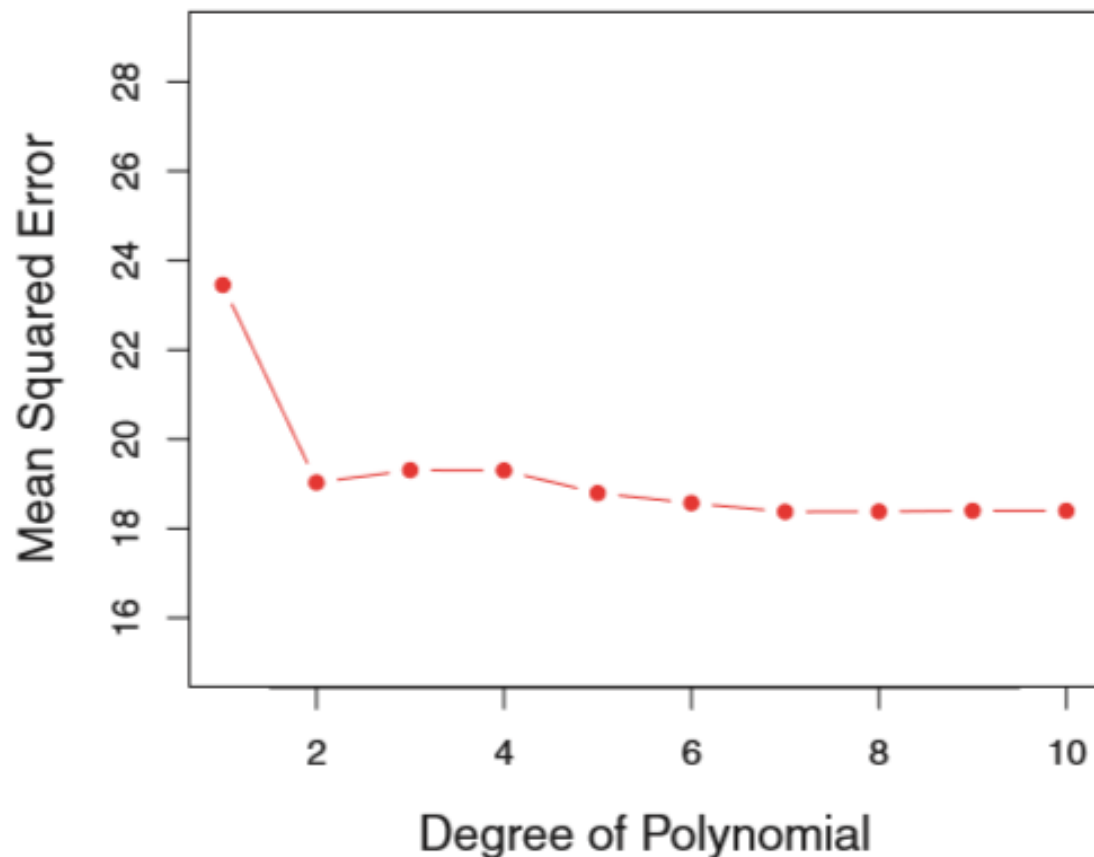


알 수 있는 것! (2 가지)



$$\text{mpg} = \beta_0 + \beta_1 \times \text{horsepower} + \beta_2 \times \text{horsepower}^2 + \epsilon$$

단점! (2 가지) – 경우에 따라 변화 심함



$$\text{mpg} = \beta_0 + \beta_1 \times \text{horsepower} + \beta_2 \times \text{horsepower}^2 + \epsilon$$

단점! (2 가지) – Training 사이즈 축소 = Performance 감소

사용 가능 한 관찰 값 = **Training Set**

VS.

사용 가능 한 관찰 값

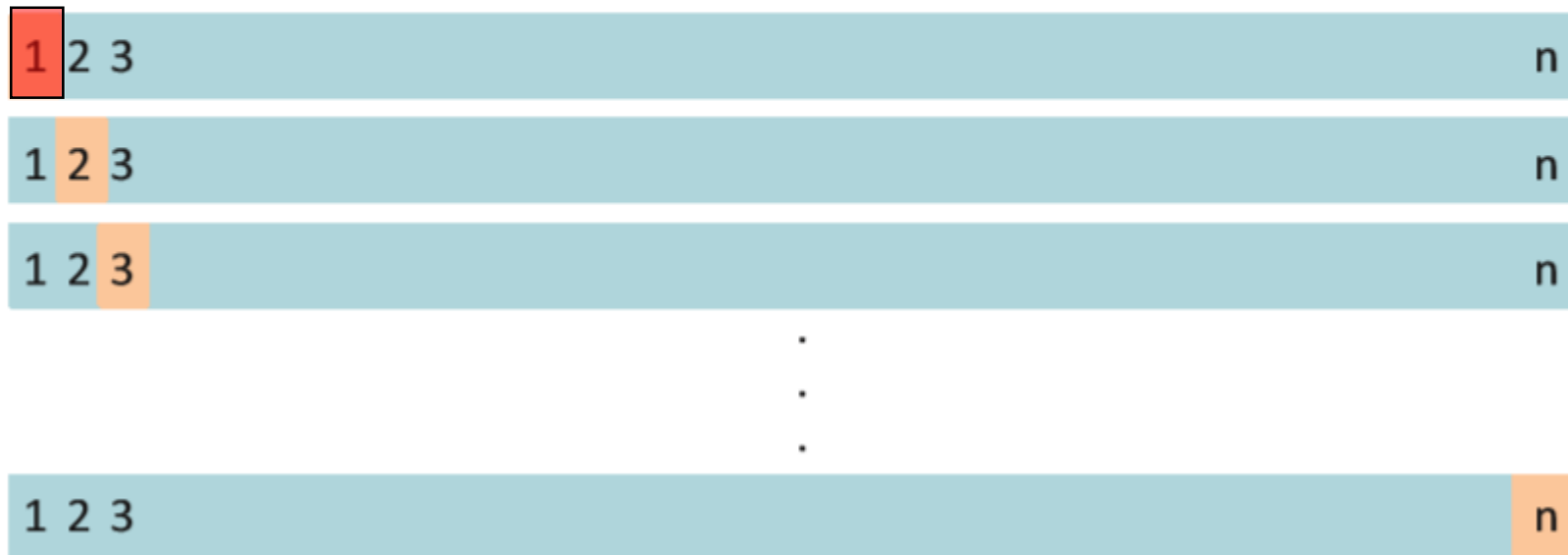
Training Set

Validation Set

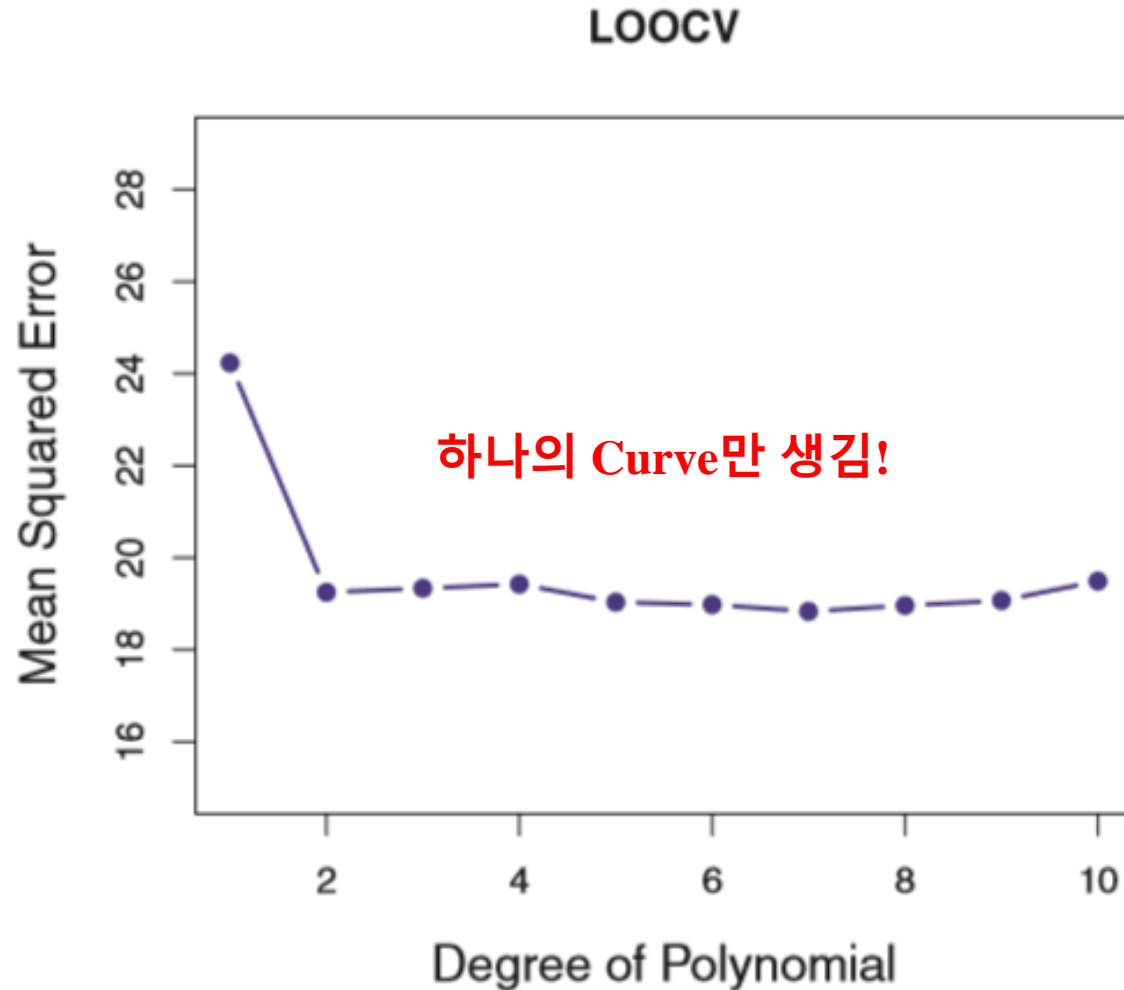
단점 제거 → LOOCV (Leave-One-Out Cross-Validation)



하나만 Validation Set으로 씀!



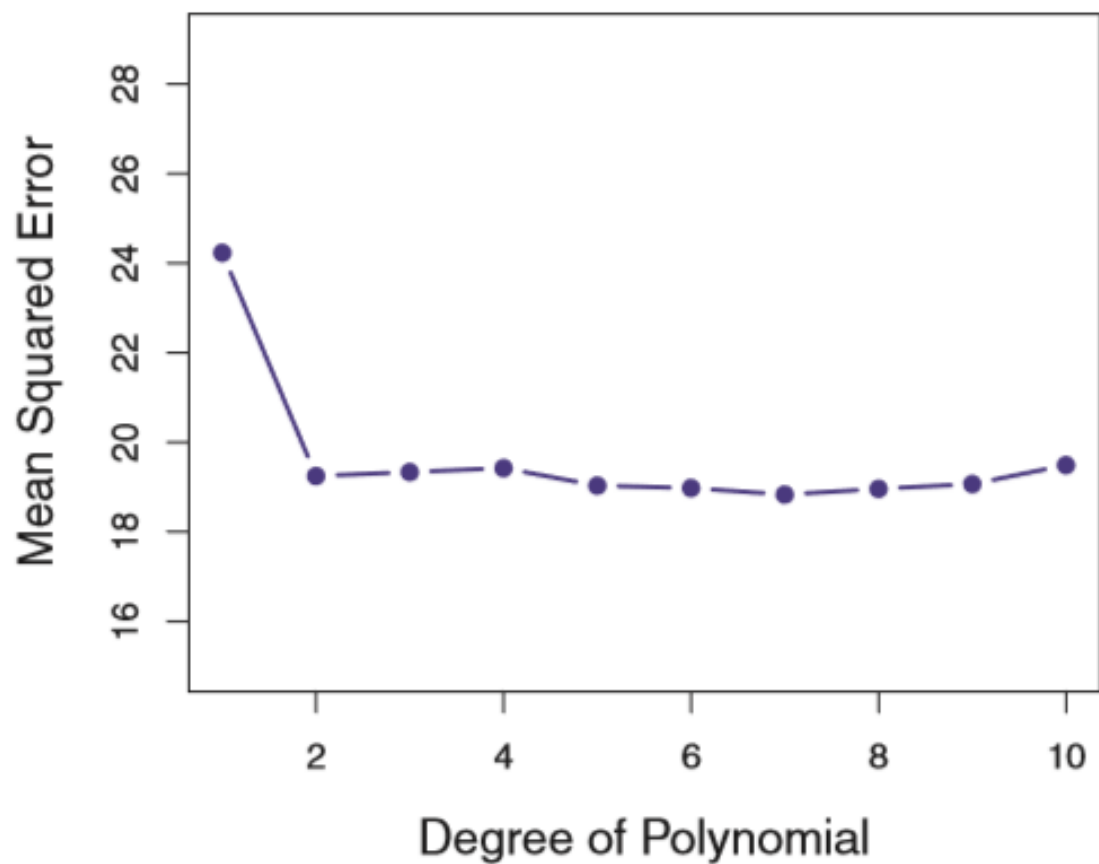
단점 제거 → LOOCV (Leave-One-Out Cross-Validation)



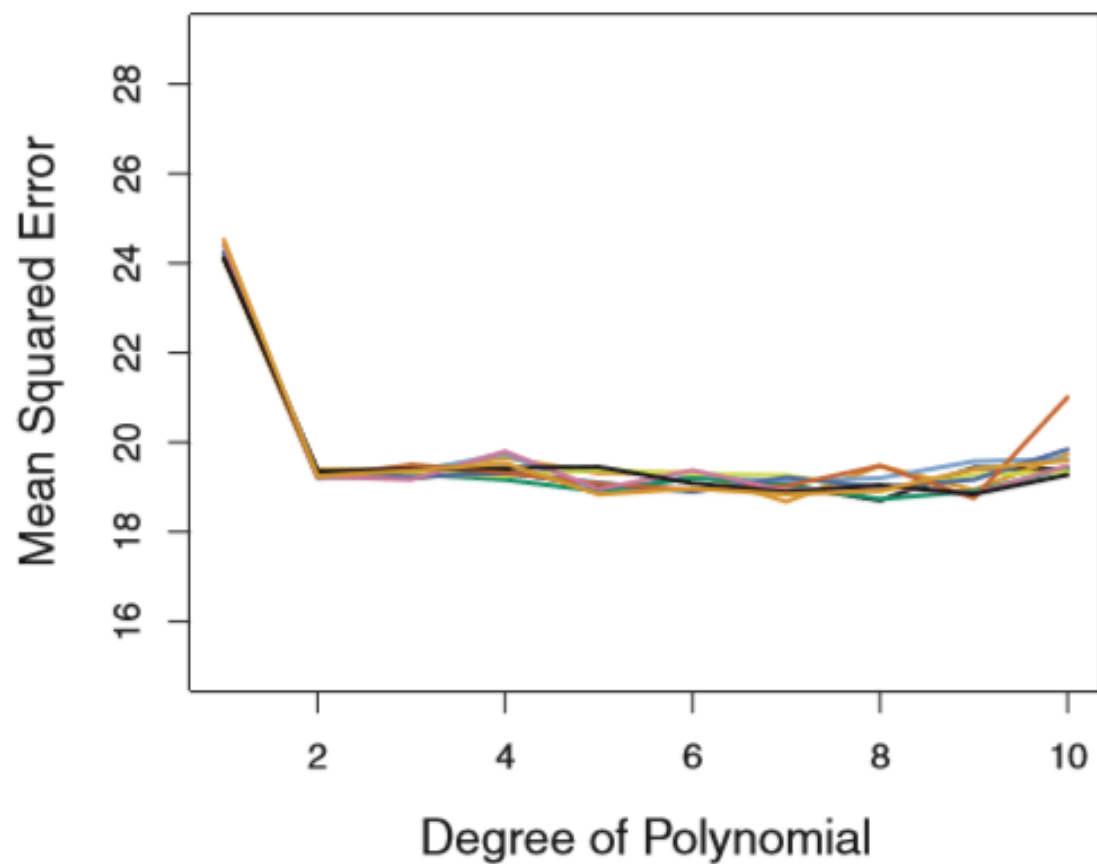
$$CV_{(n)} = \frac{1}{n} \sum_{i=1}^n \text{MSE}_i$$

정확도는?

LOOCV

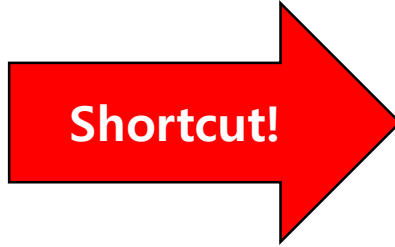


10-fold CV



특수한 경우의 추정 test MSE: Linear Model Fitting

$$CV_{(n)} = \frac{1}{n} \sum_{i=1}^n \text{MSE}_i$$



$$CV_{(n)} = \frac{1}{n} \sum_{i=1}^n \left(\frac{y_i - \hat{y}_i}{1 - h_i} \right)^2$$

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

||

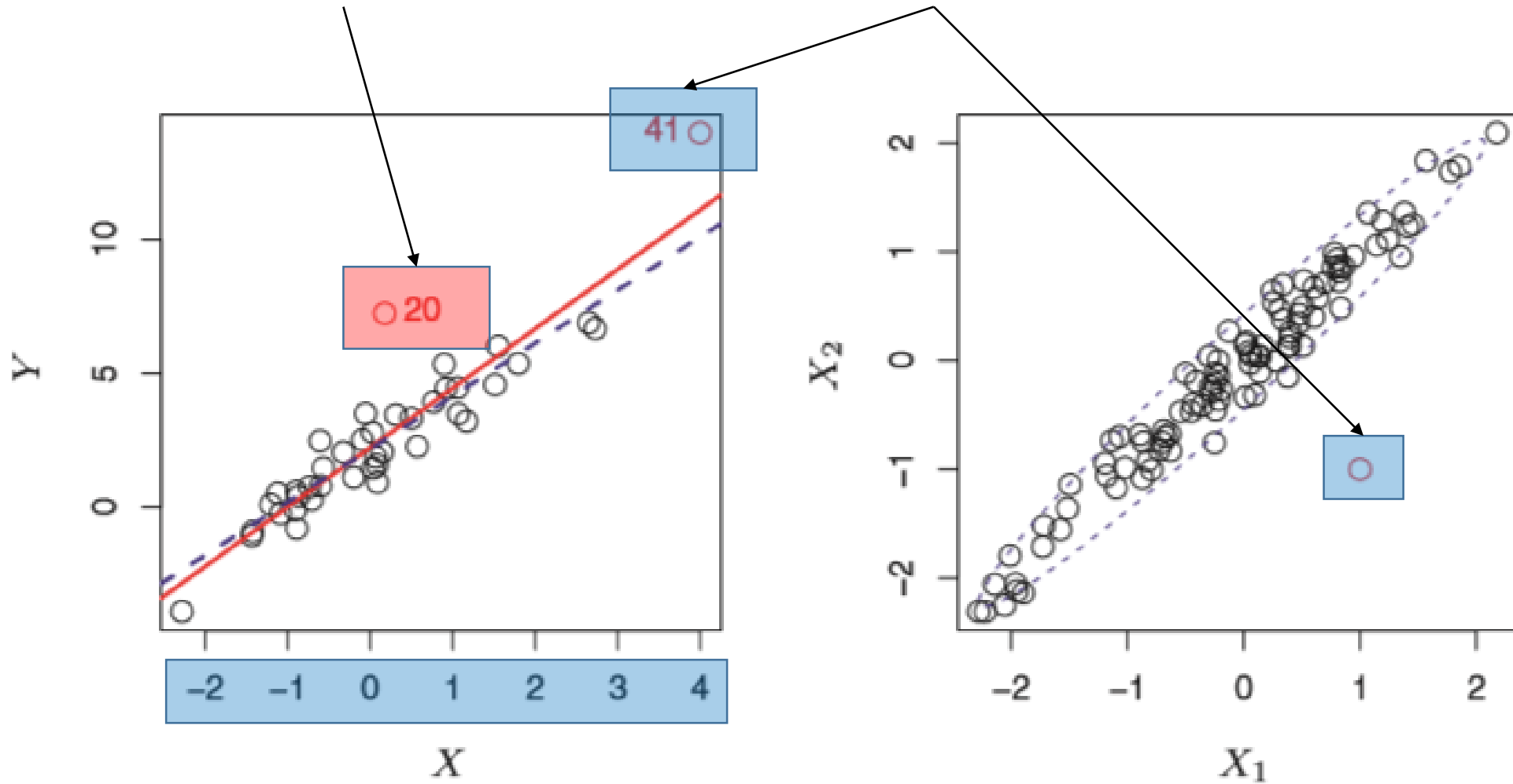
$$\text{MSE}_2 = (y_2 - \hat{y}_2)^2$$

(x_2, y_2) 를 제외한 나머지 데이터에서 Fitting한 Model에서 예측된 값!

$$h_i = \frac{1}{n} + \frac{(x_i - \bar{x})^2}{\sum_{i'=1}^n (x_{i'} - \bar{x})^2}$$

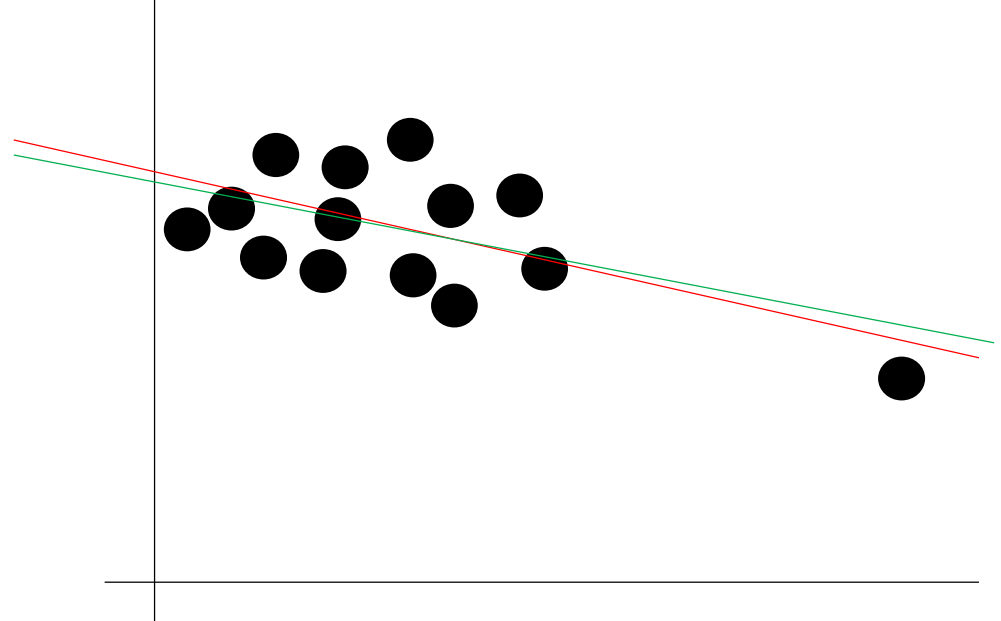
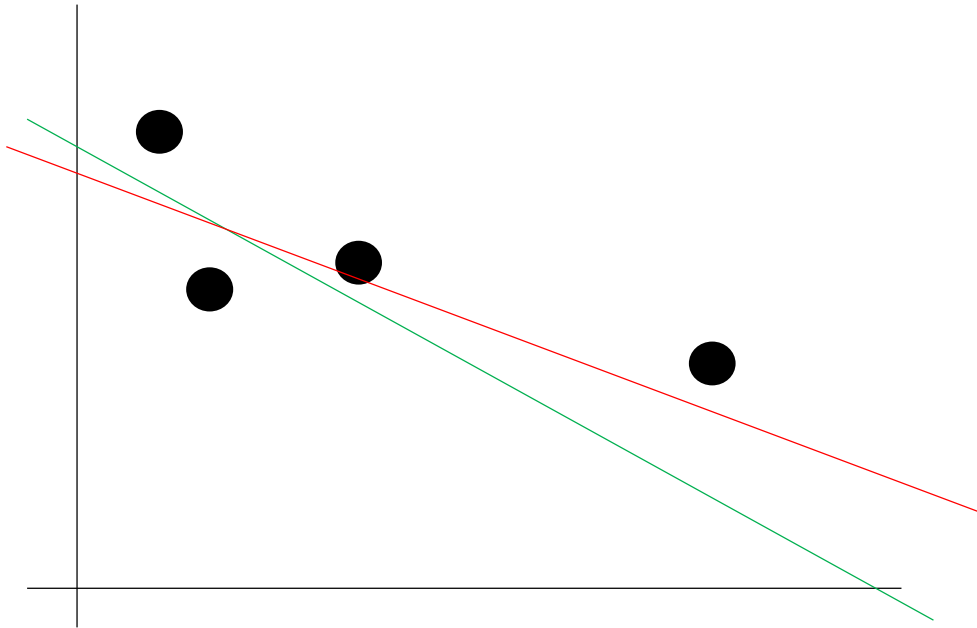
(x_2, y_2) 를 포함한 모든 데이터에서 Fitting한 Model에서 예측된 값!

복습: Outlier VS. Leverage point



Leverage statistic
$$h_i = \frac{1}{n} + \frac{(x_i - \bar{x})^2}{\sum_{i'=1}^n (x_{i'} - \bar{x})^2}$$

복습: Outlier VS. Leverage point



Leverage statistic $h_i = \frac{1}{n} + \frac{(x_i - \bar{x})^2}{\sum_{i'=1}^n (x_{i'} - \bar{x})^2}$

Why shortcut?

$$CV_{(n)} = \frac{1}{n} \sum_{i=1}^n \text{MSE}_i$$

Shortcut!

$$CV_{(n)} = \frac{1}{n} \sum_{i=1}^n \left(\frac{y_i - \hat{y}_i}{1 - h_i} \right)^2$$

$$\text{MSE}_2 = (y_2 - \hat{y}_2)^2$$

(x_2, y_2) 를 제외한 나머지 데이터에서 Fitting한 Model에서 예측된 값!

(x_2, y_2) 를 포함한 모든 데이터에서 Fitting한 Model에서 예측된 값!

Why shortcut?

$$CV_{(n)} = \frac{1}{n} \sum_{i=1}^n \text{MSE}_i$$

Shortcut!

$$CV_{(n)} = \frac{1}{n} \sum_{i=1}^n \left(\frac{y_i - \hat{y}_i}{1 - h_i} \right)^2$$

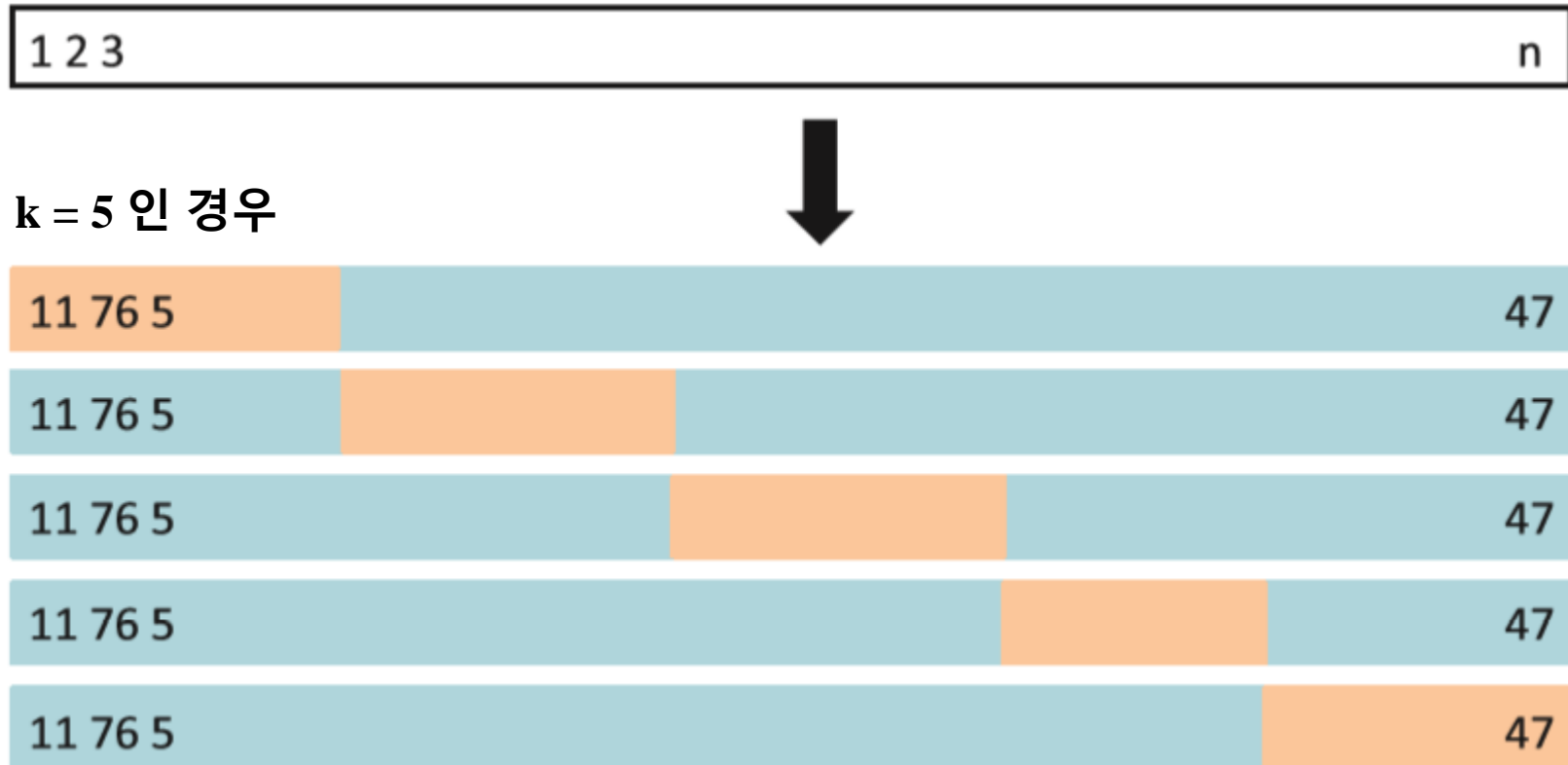
1 2 3 한 Set만 Fitting 하면 된다! n



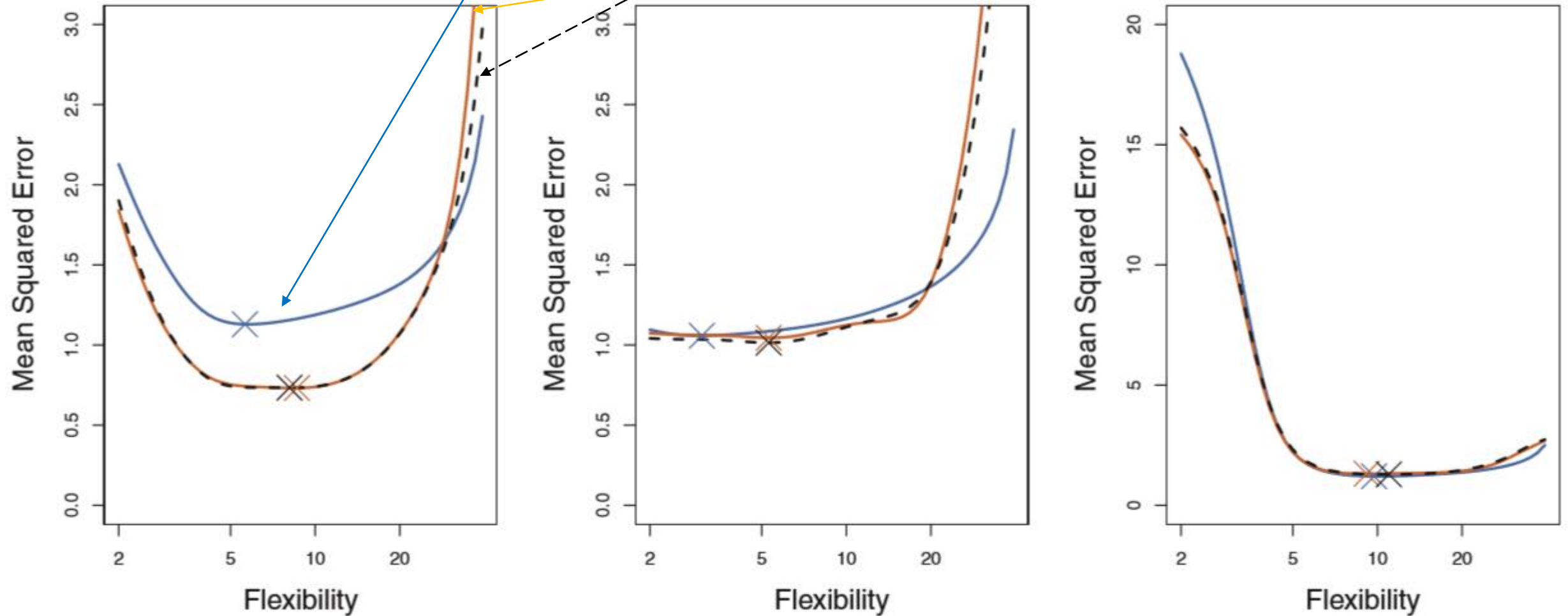
LOOCV는 Linear Model 아닌 경우, 계산 량 너무 크다!

따라서, 더욱 일반적인 형태가 필요!

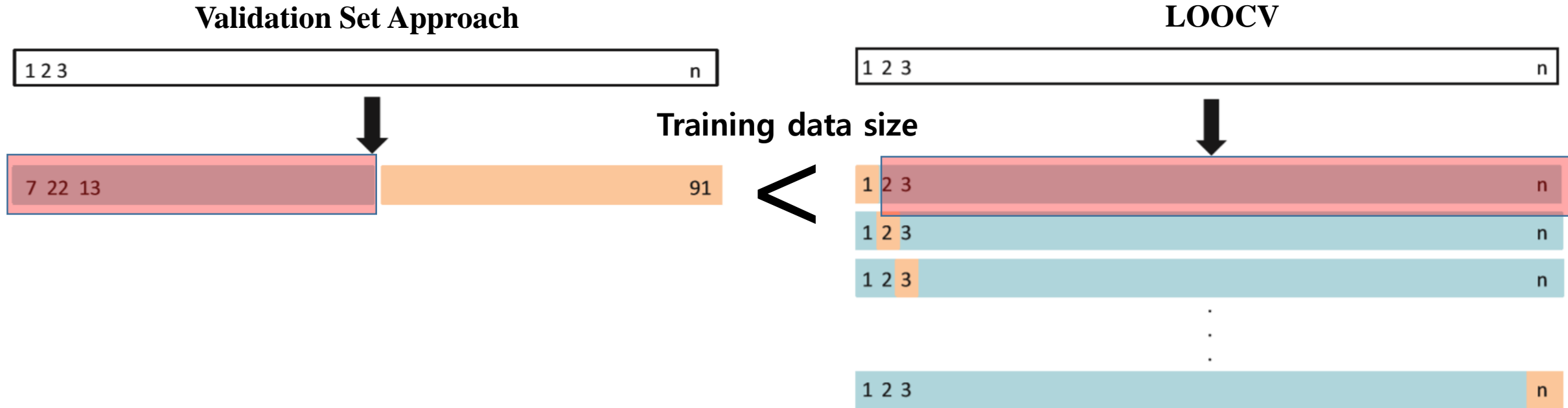
→ k-Fold Cross-Validation



Comparison: True VS. LOOCV VS. 10-fold CV

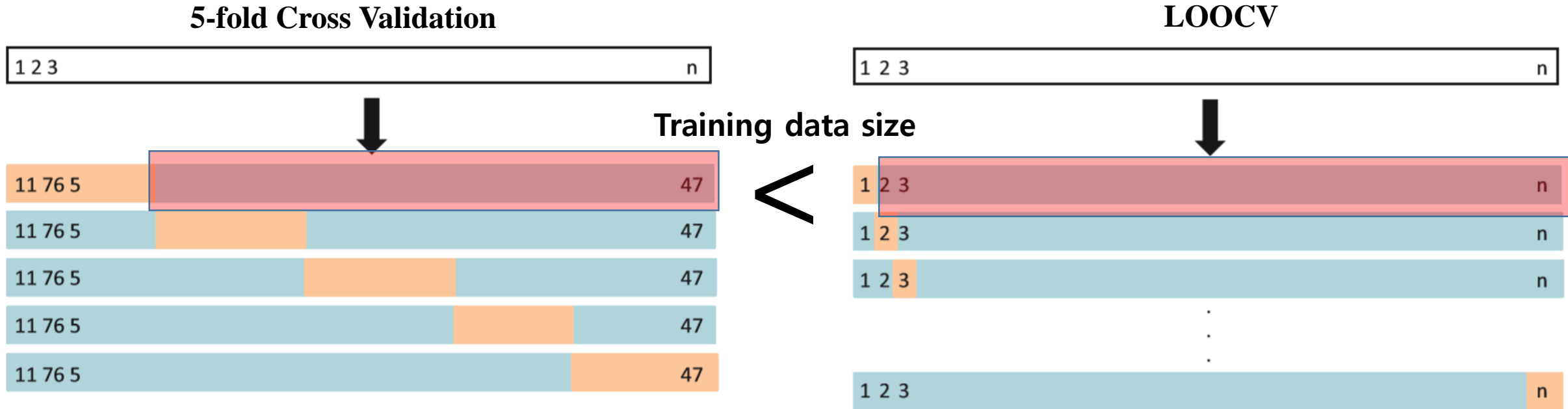


Bias / Variance Trade-off



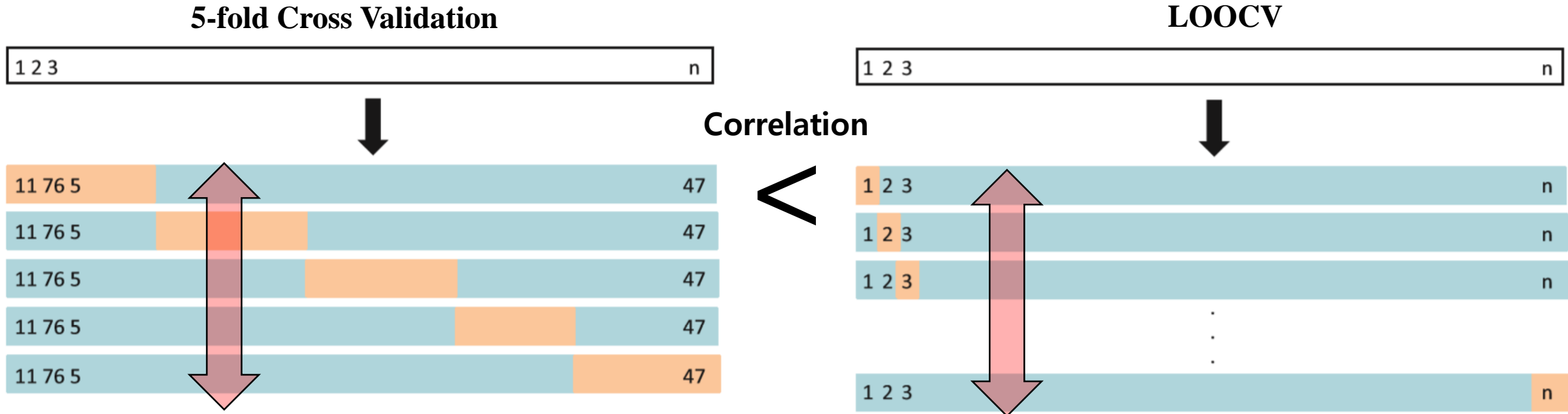
Training에 사용되는 데이터의 수가 많아지면서 Bias는 적어지게 된다.

Bias / Variance Trade-off



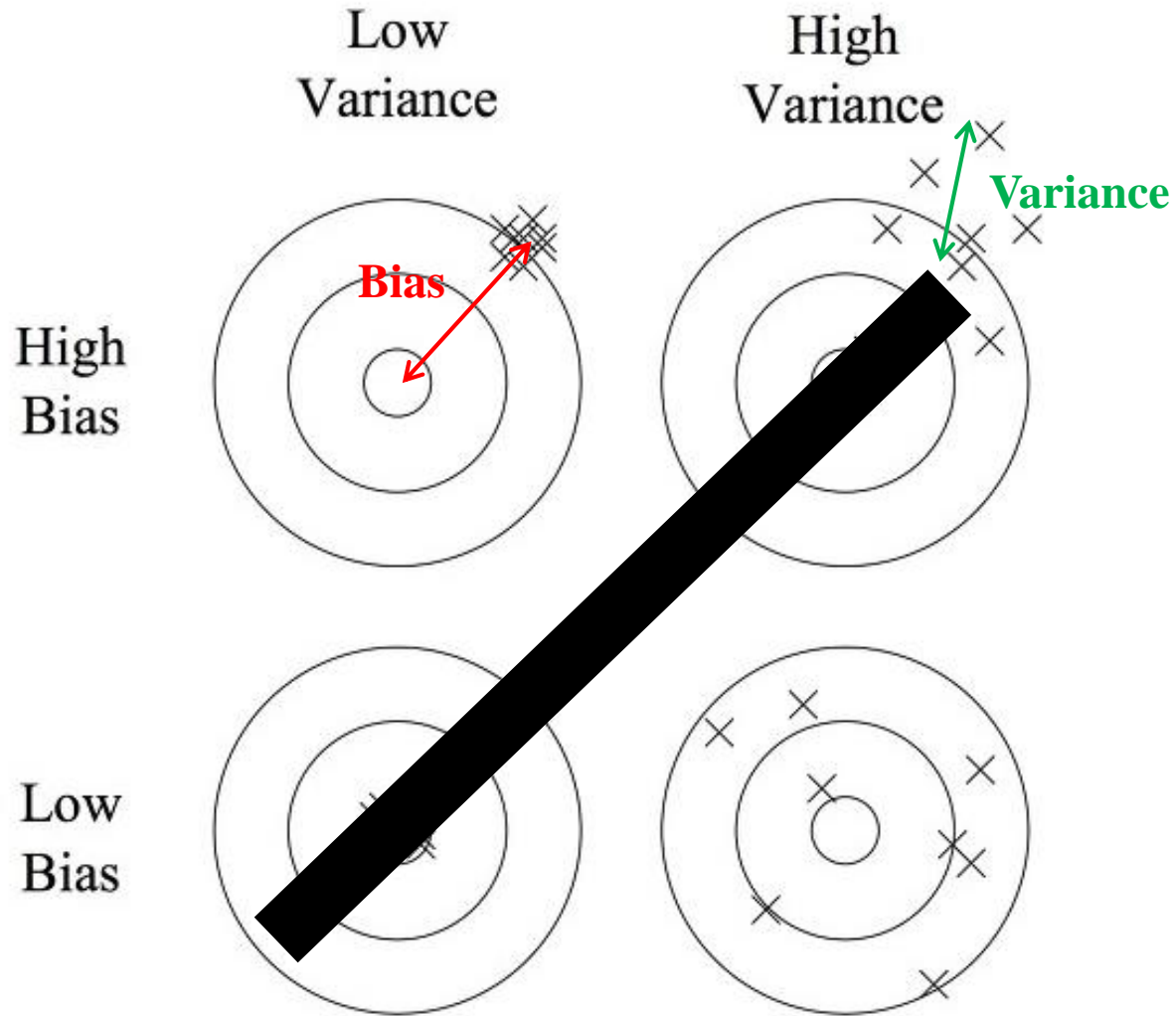
LOOCV가 Bias는 가장 적게 Fitting한다!

Bias / Variance Trade-off



LOOCV가 Variance는 CV 보다 크게 Fitting한다!

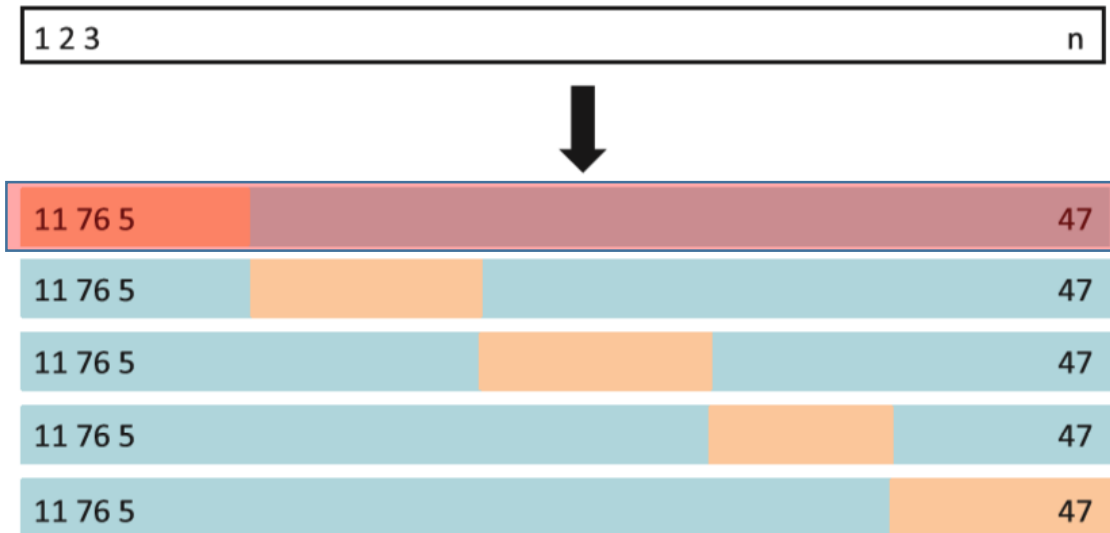
Bias / Variance Trade-off



Bias / Variance Trade-off

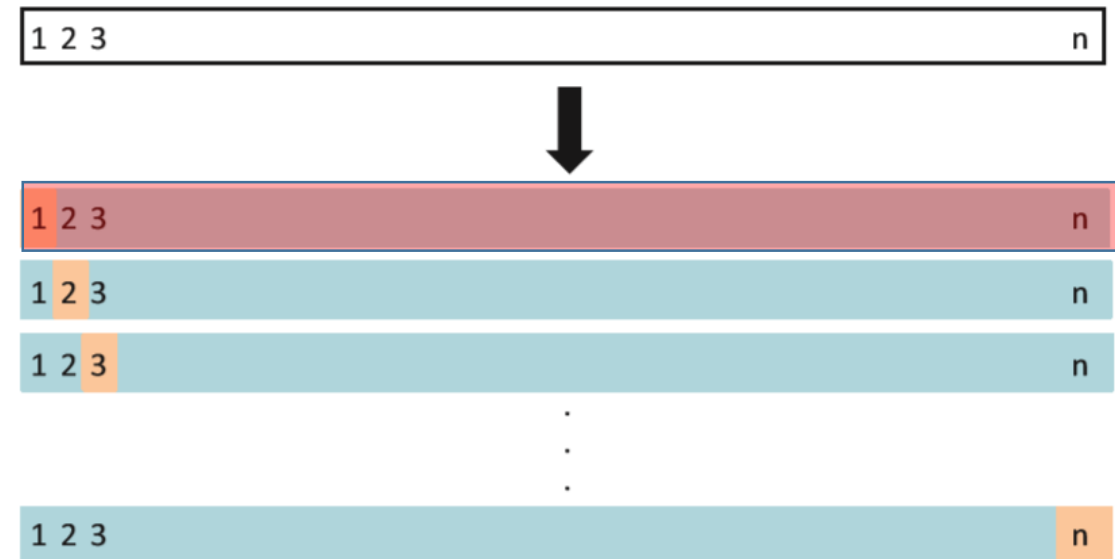
5-fold Cross Validation

$$CV_{(k)} = \frac{1}{k} \sum_{i=1}^k \text{MSE}_i$$



LOOCV

$$CV_{(n)} = \frac{1}{n} \sum_{i=1}^n \text{MSE}_i$$



Bias / Variance Trade-off

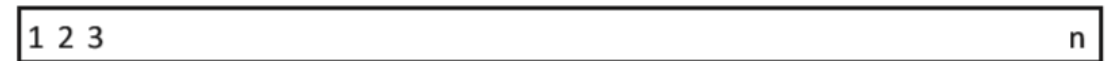
5-fold Cross Validation

$$CV_{(k)} = \frac{1}{k} \sum_{i=1}^k \text{MSE}_i$$



LOOCV

$$CV_{(n)} = \frac{1}{n} \sum_{i=1}^n \text{MSE}_i$$

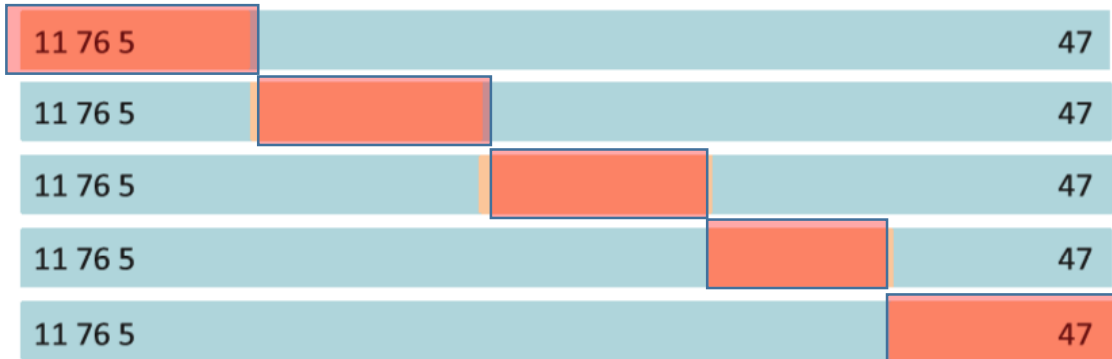


Bias / Variance Trade-off

5-fold Cross Validation

$$CV_{(k)} = \frac{1}{k} \sum_{i=1}^k \text{MSE}_i$$

1 2 3 n



LOOCV

$$CV_{(n)} = \frac{1}{n} \sum_{i=1}^n \text{MSE}_i$$

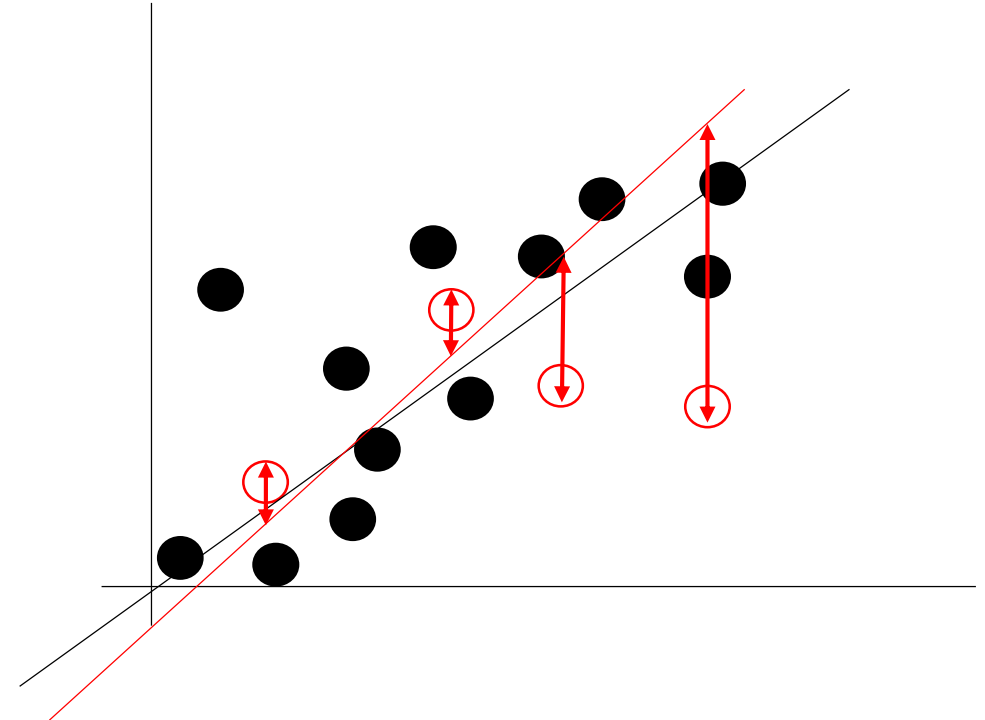
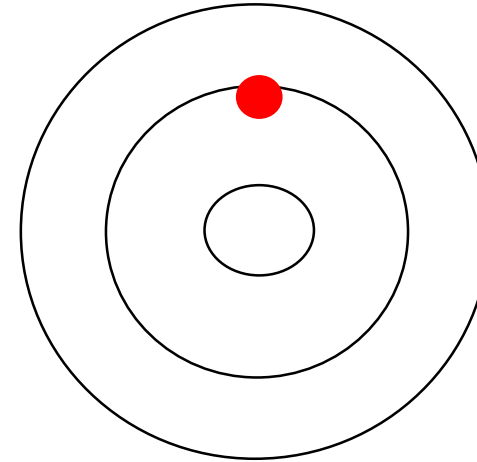
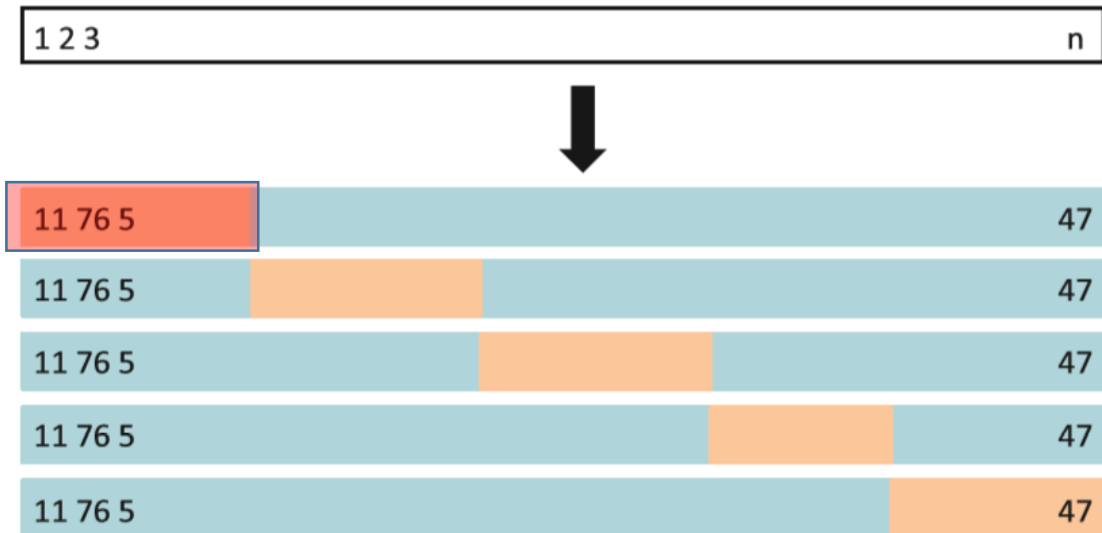
1 2 3 n



Bias / Variance Trade-off

5-fold Cross Validation

$$CV_{(k)} = \frac{1}{k} \sum_{i=1}^k \text{MSE}_i$$

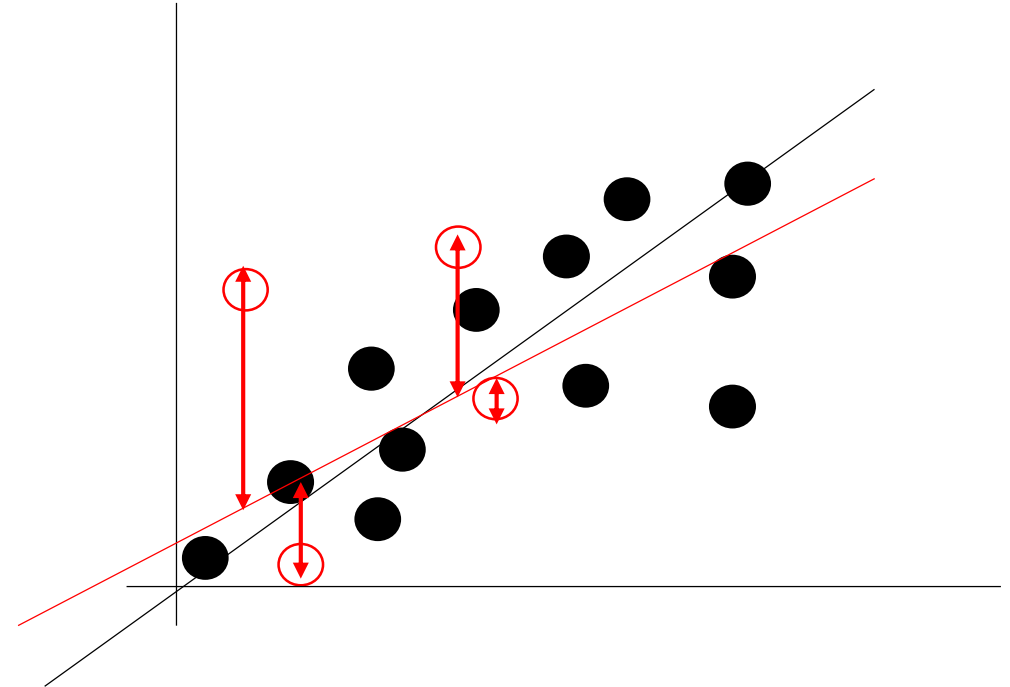
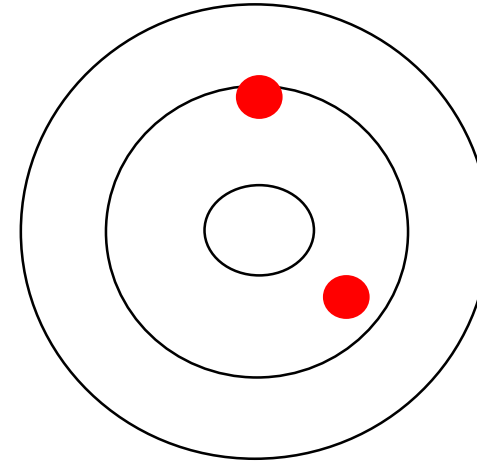
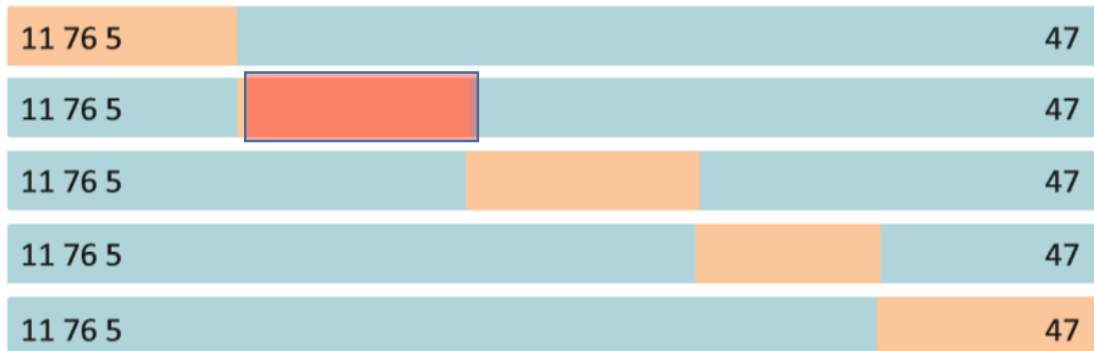


Bias / Variance Trade-off

5-fold Cross Validation

$$CV_{(k)} = \frac{1}{k} \sum_{i=1}^k MSE_i$$

1 2 3 n

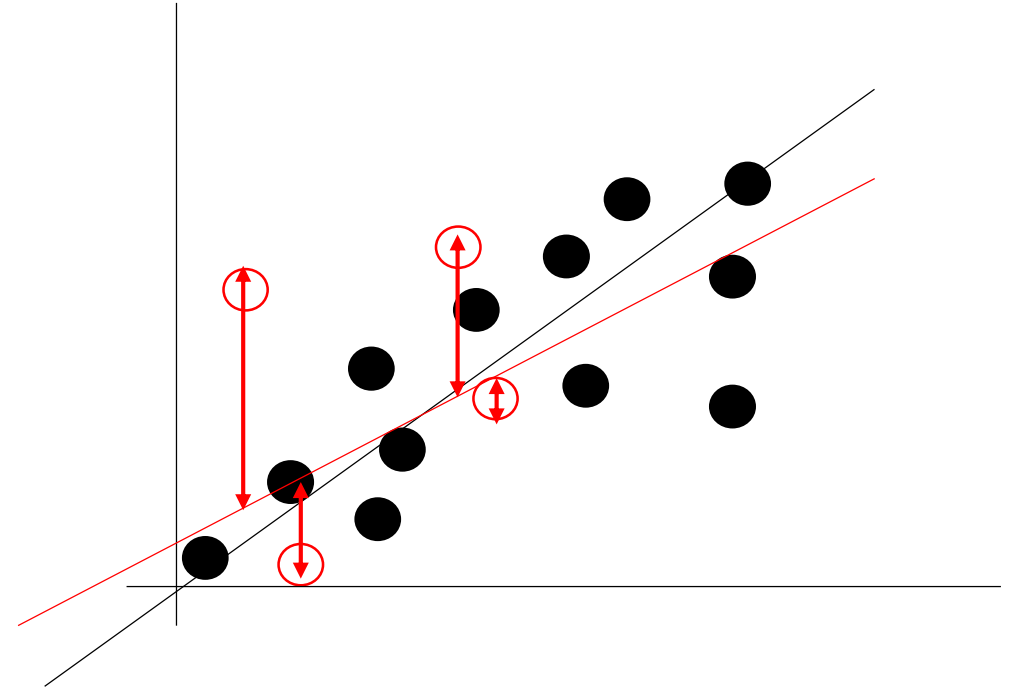
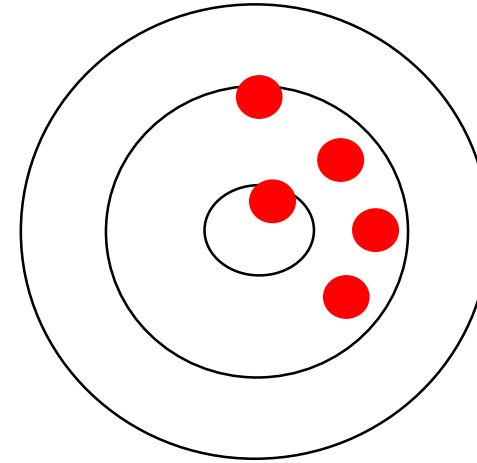
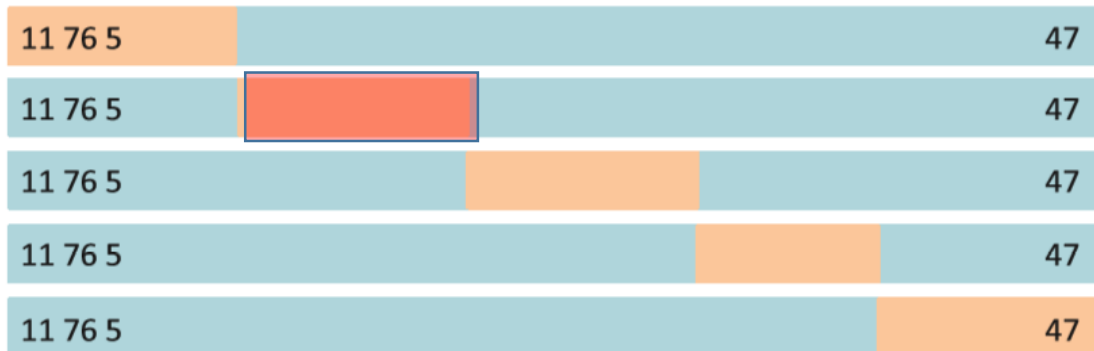


Bias / Variance Trade-off

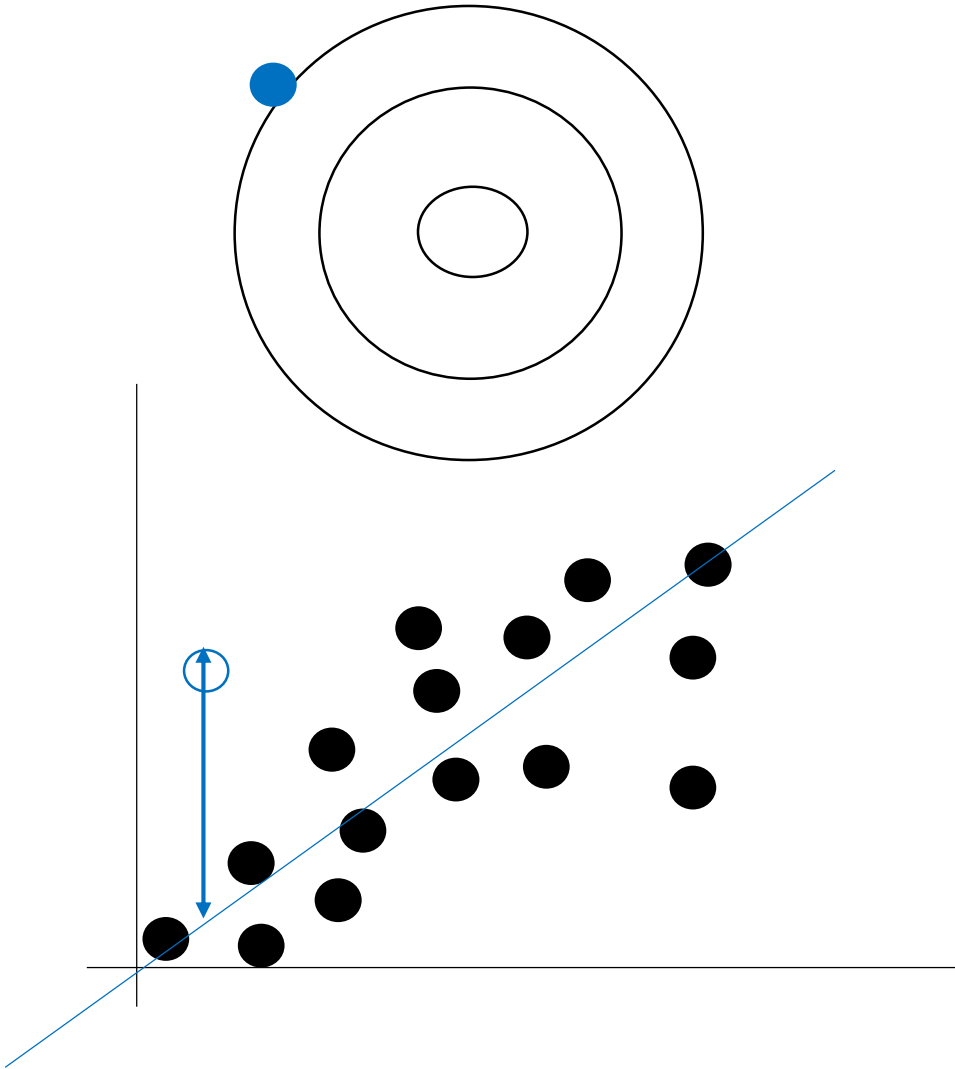
5-fold Cross Validation

$$CV_{(k)} = \frac{1}{k} \sum_{i=1}^k \text{MSE}_i$$

1 2 3 n

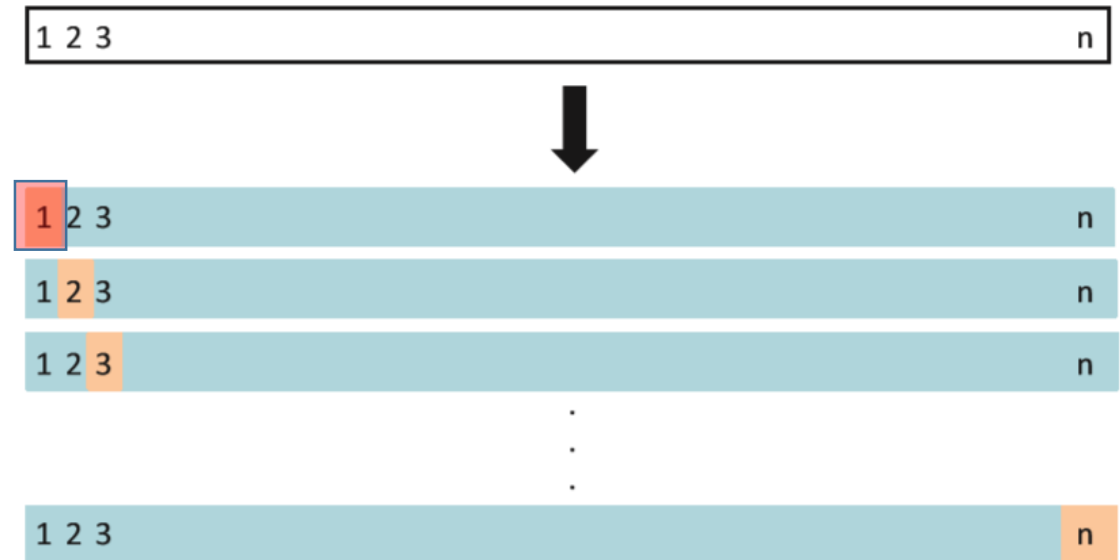


Bias / Variance Trade-off

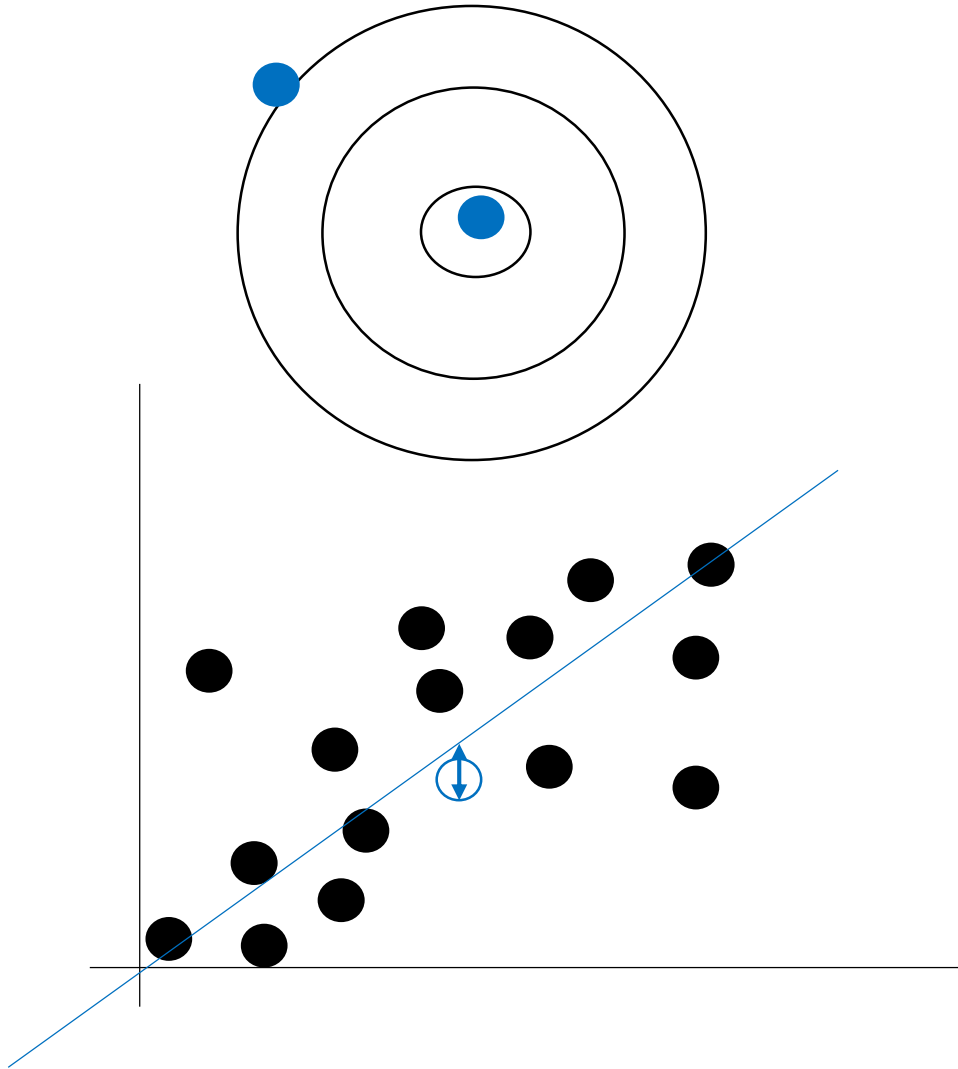


LOOCV

$$CV_{(n)} = \frac{1}{n} \sum_{i=1}^n MSE_i$$

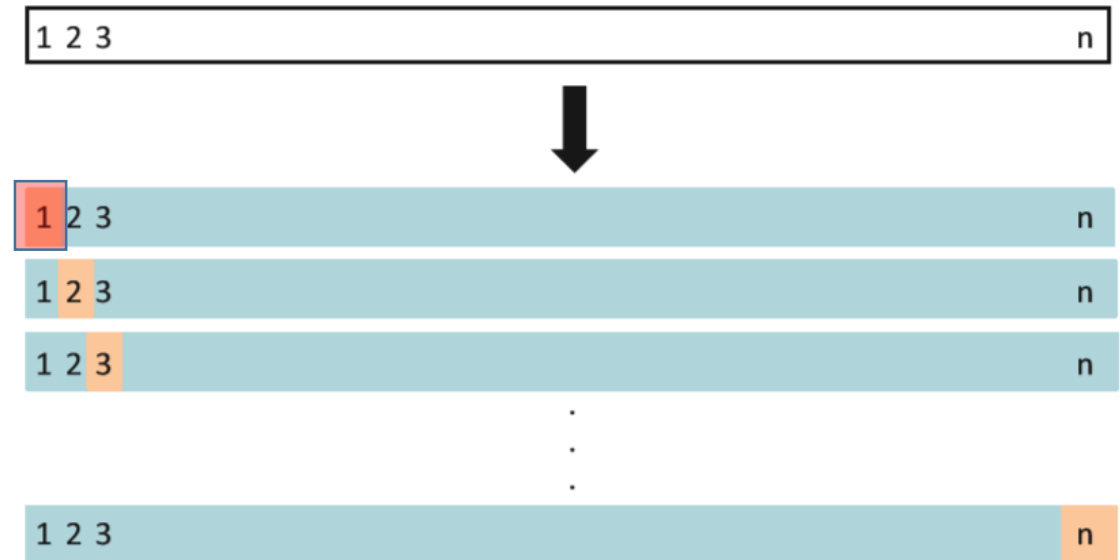


Bias / Variance Trade-off

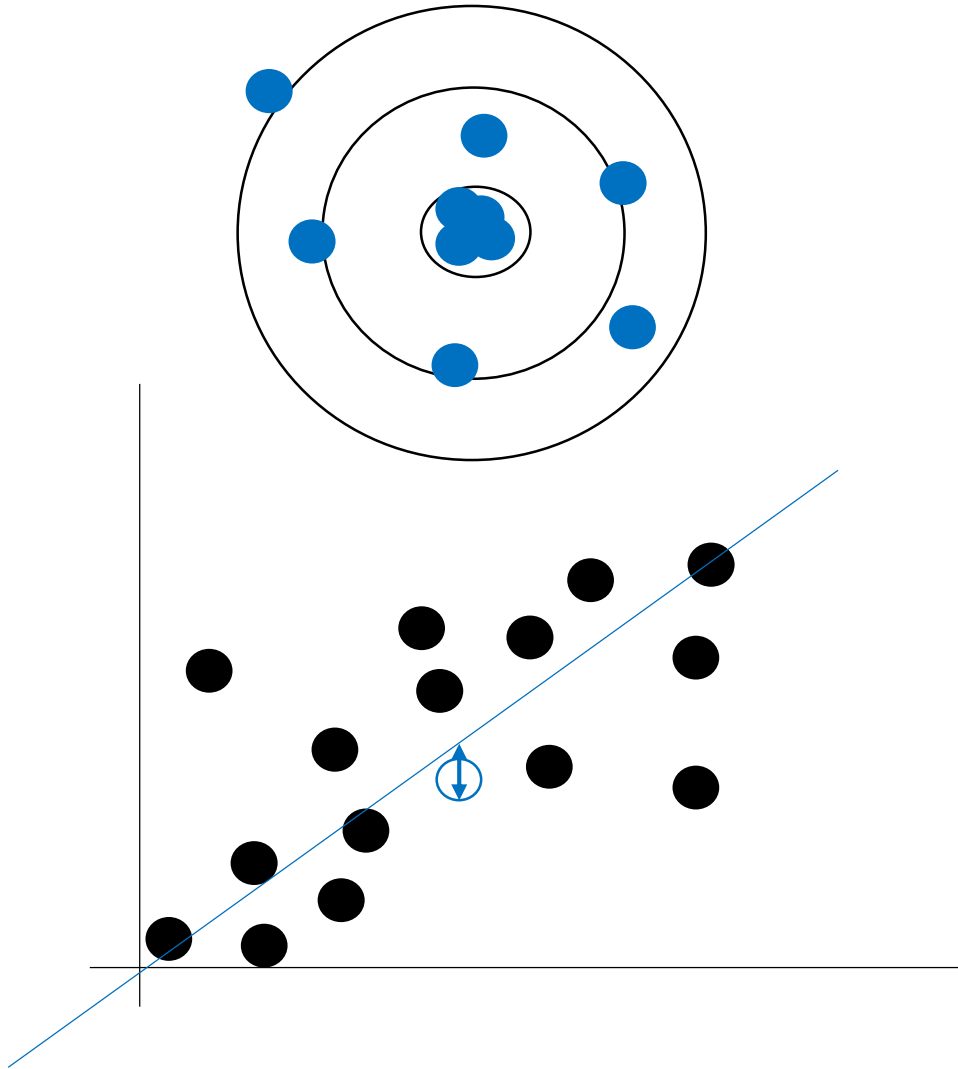


LOOCV

$$CV_{(n)} = \frac{1}{n} \sum_{i=1}^n \text{MSE}_i$$

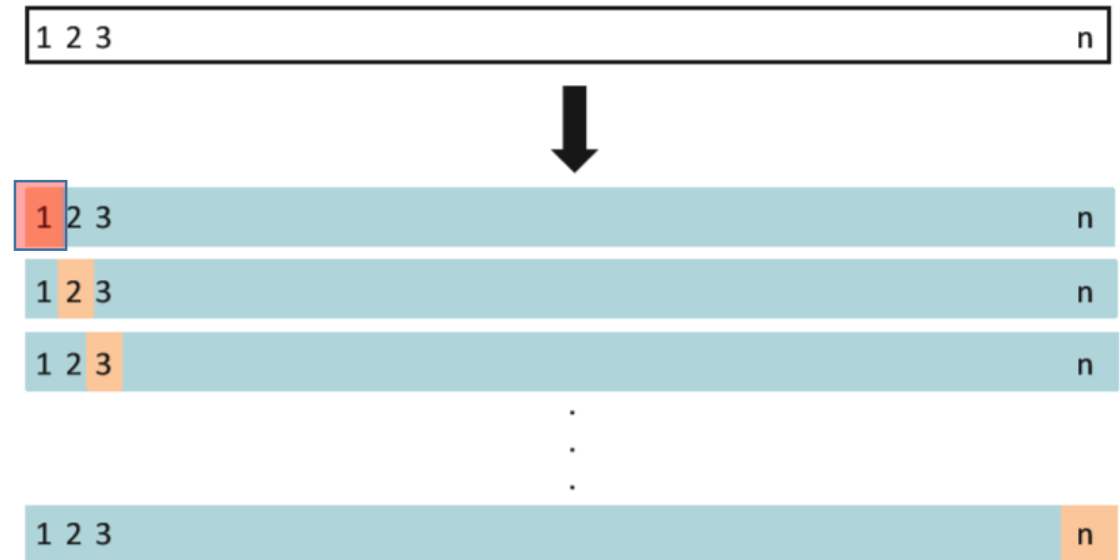


Bias / Variance Trade-off



LOOCV

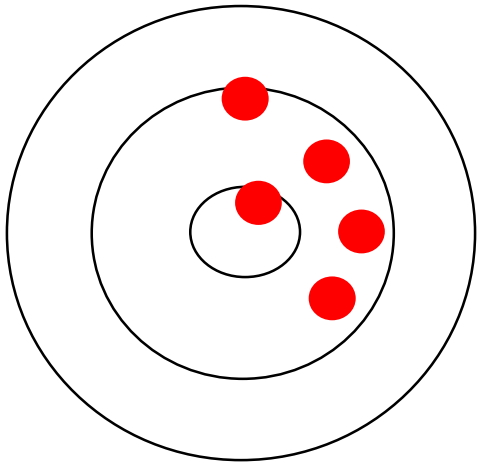
$$CV_{(n)} = \frac{1}{n} \sum_{i=1}^n \text{MSE}_i$$



Bias / Variance Trade-off

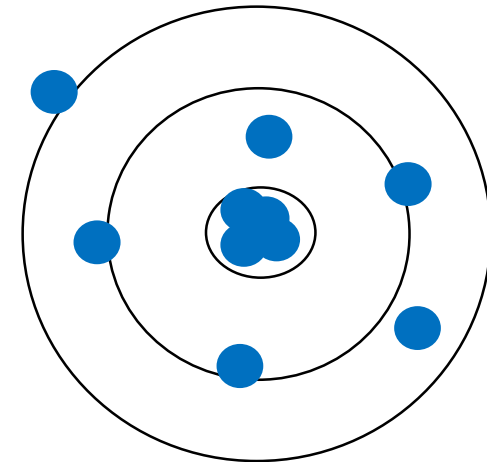
5-fold Cross Validation

$$CV_{(k)} = \frac{1}{k} \sum_{i=1}^k \text{MSE}_i$$

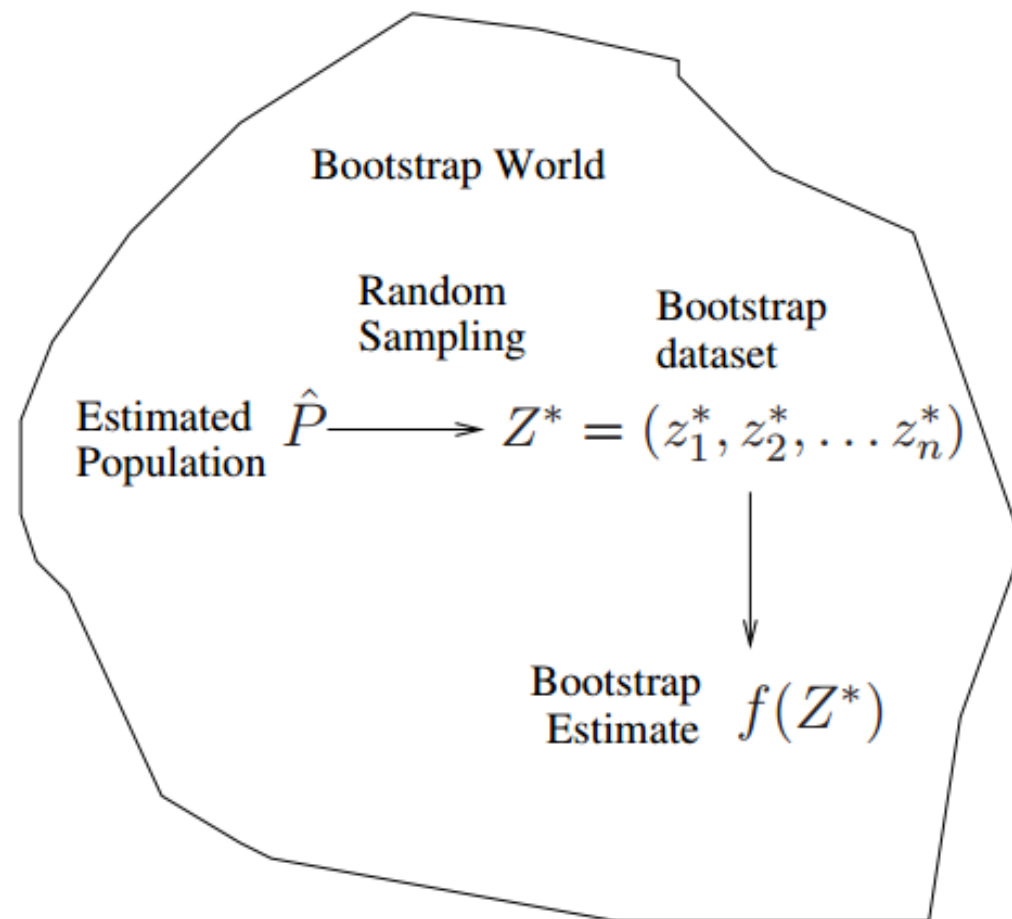
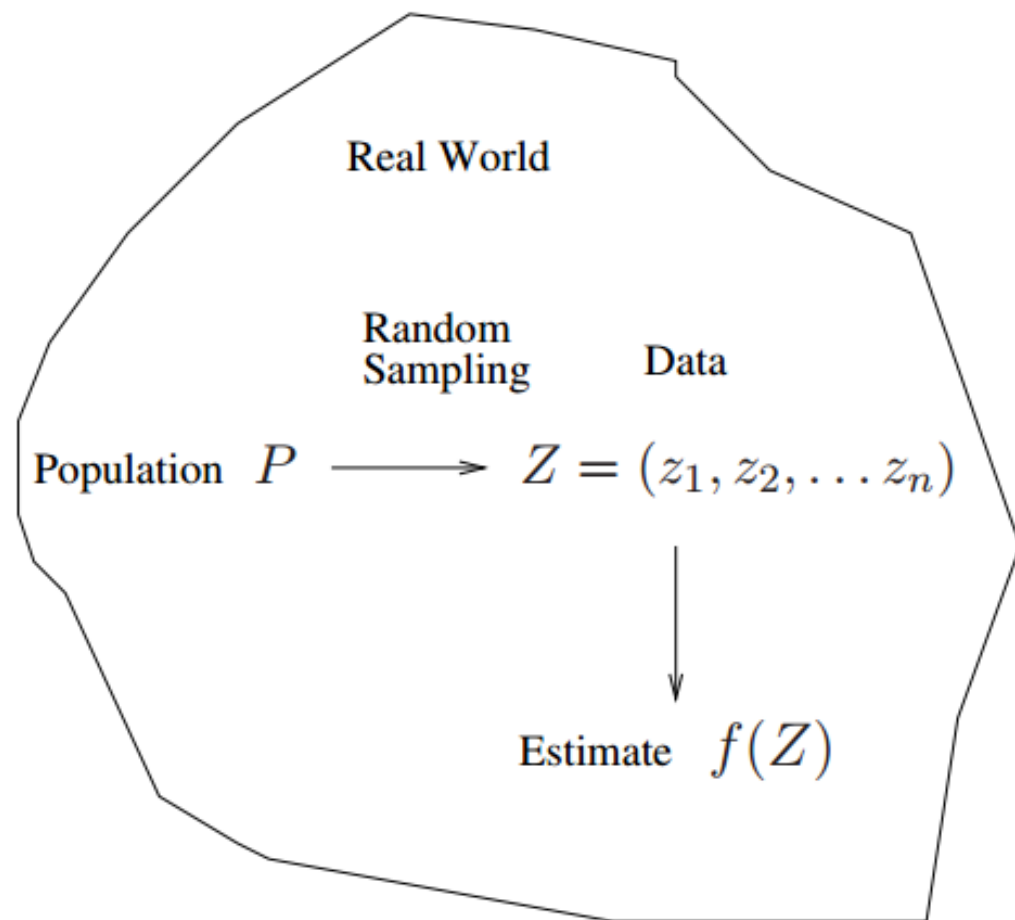


LOOCV

$$CV_{(n)} = \frac{1}{n} \sum_{i=1}^n \text{MSE}_i$$



Bootstrap



Bootstrap

반복을 허용하고,

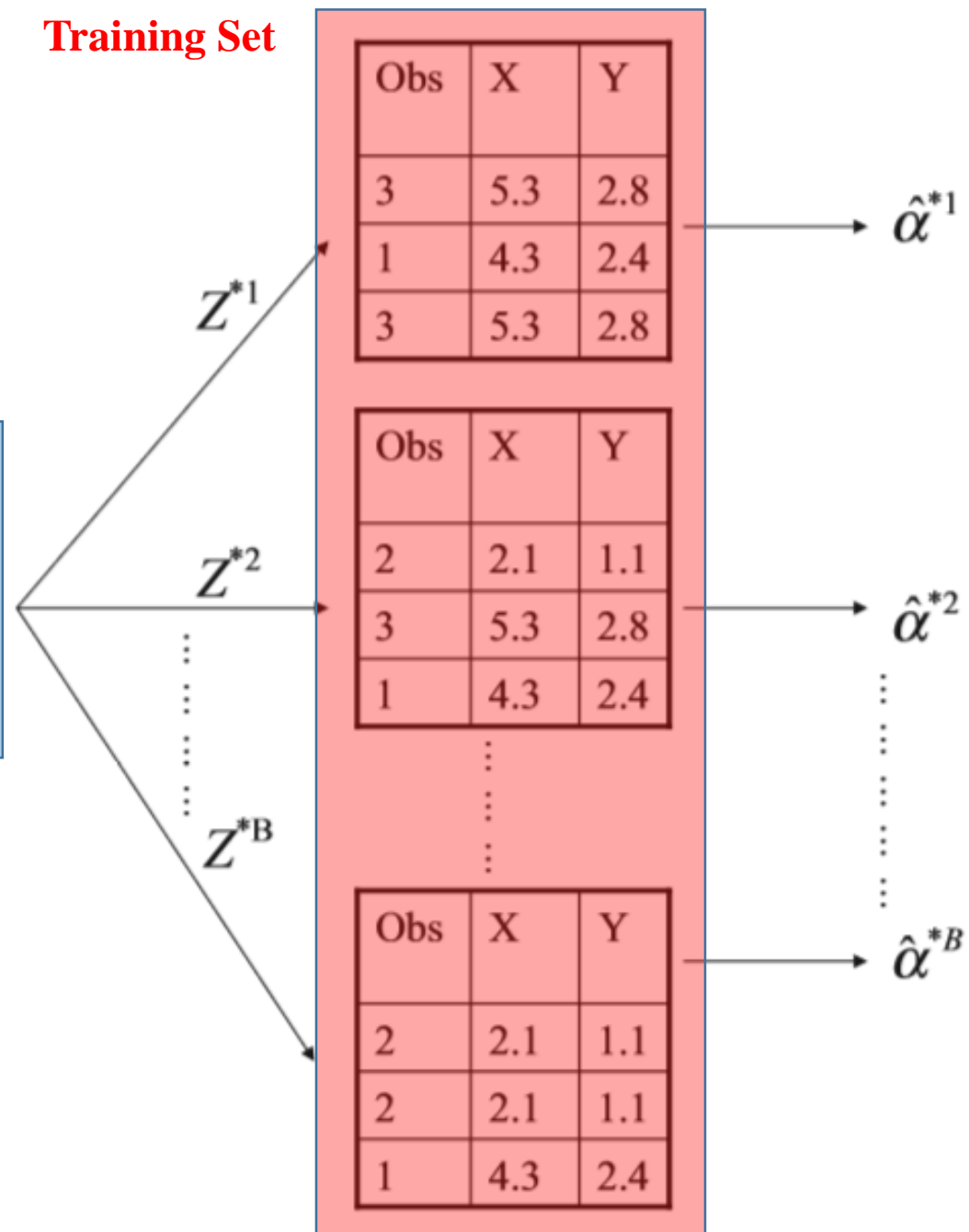
관찰 값에서 Resampling 한다!

Validation Set

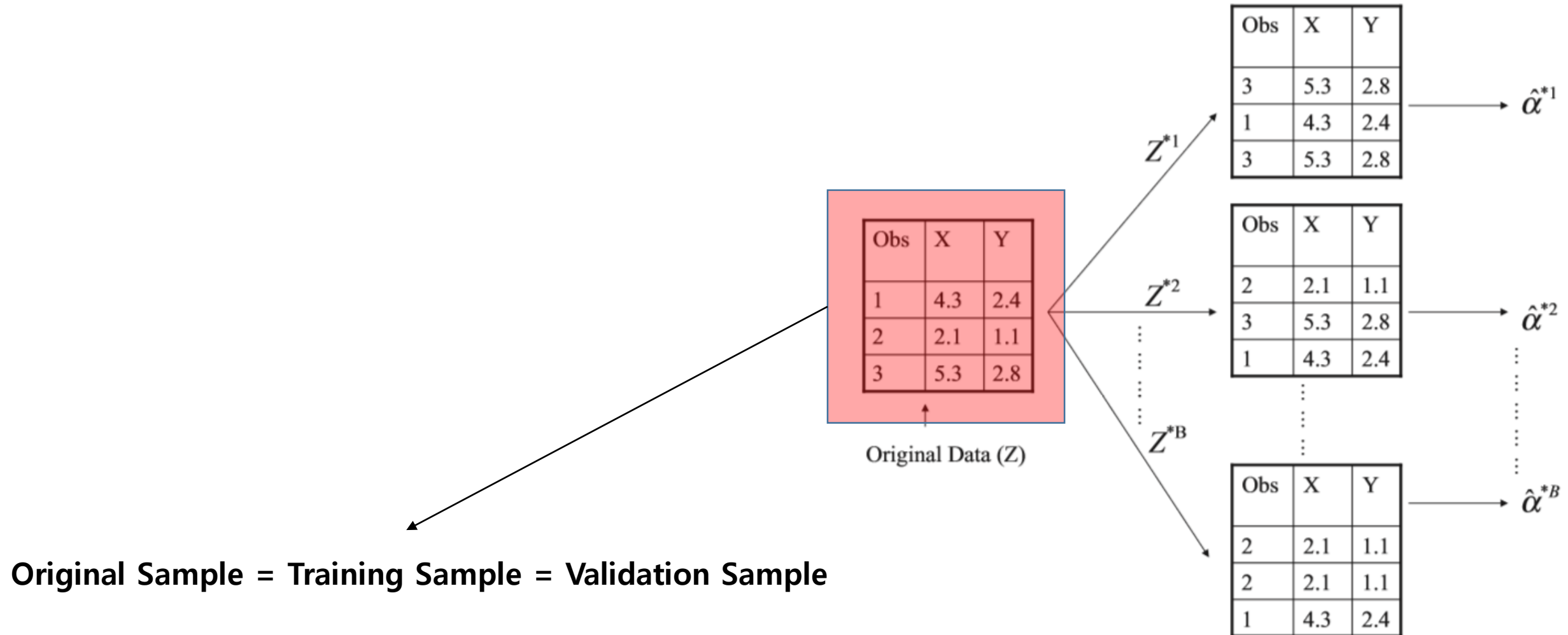
Obs	X	Y
1	4.3	2.4
2	2.1	1.1
3	5.3	2.8

Original Data (Z)

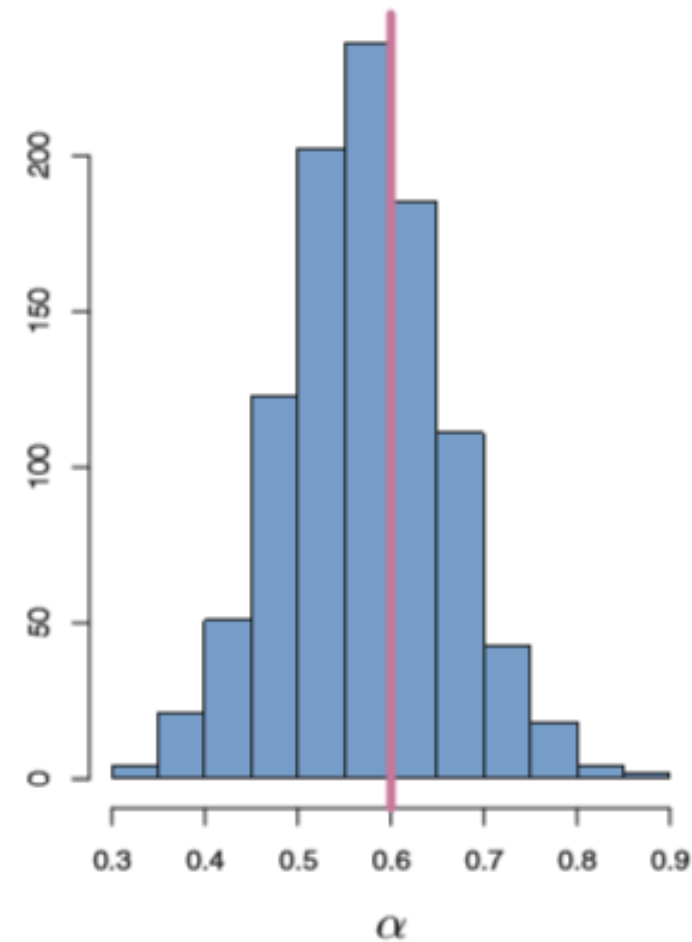
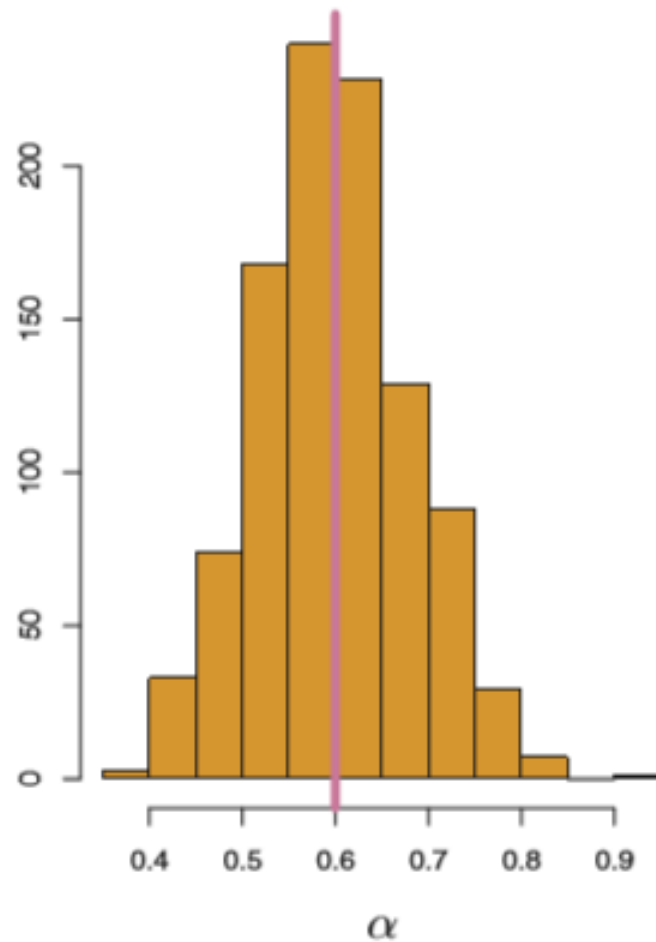
Training Set



Bootstrap – test error 추정이 가능? X



Bootstrap – Variability with parameter



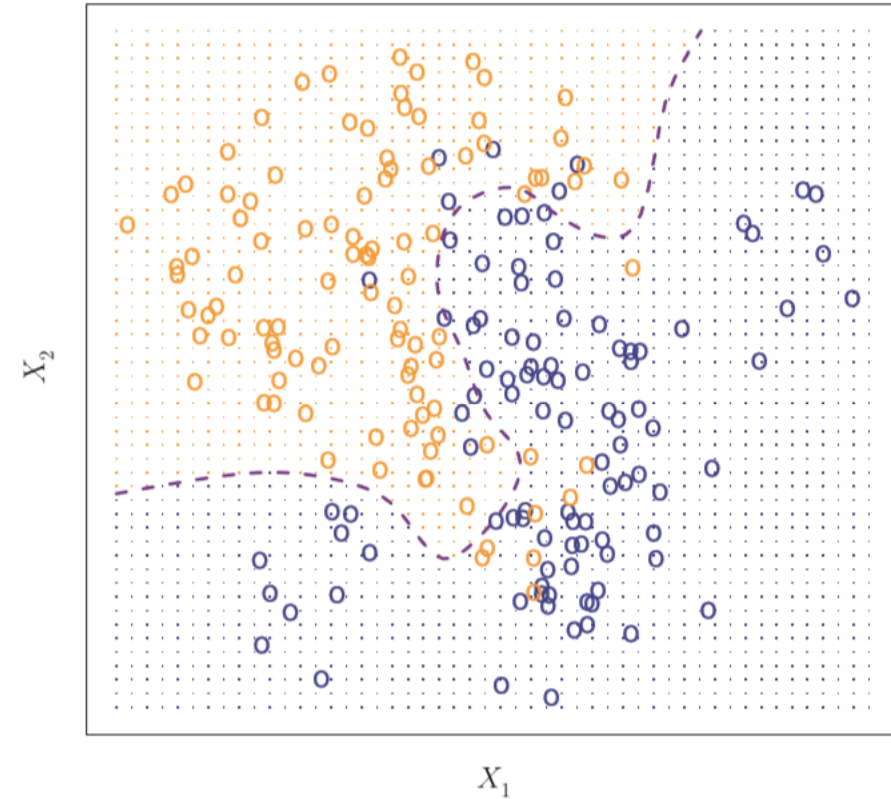
지난 주 이슈 1 – Bayes Classifier

$$P(Y = k | X = x) = \frac{f_k(x) \times \pi_k}{\sum_{l=1}^K \pi_l f_l(x)}$$

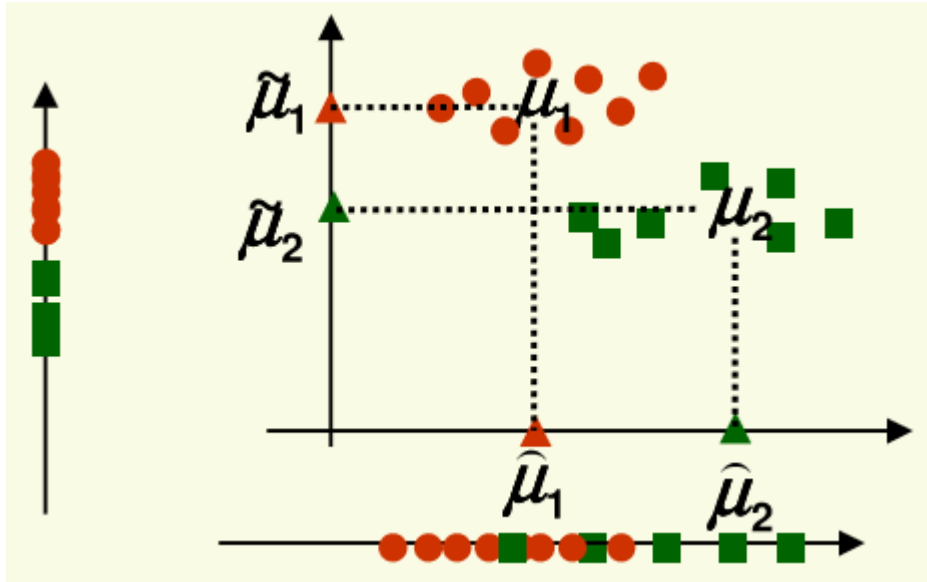
Naive Bayes Classifier

$$f_k(x) = P(X = x | Y = k)$$

$$P(X = x | Y = k) \approx \prod_{i=1}^m P(X = x_i | Y = k)$$



지난 주 이슈 2 – LDA for dimension reduction



Define their **scatter** as

$$s = \sum_{i=1}^n (z_i - \mu_z)^2$$

samples z_1, \dots, z_n

Sample mean is $\mu_z = \frac{1}{n} \sum_{i=1}^n z_i$

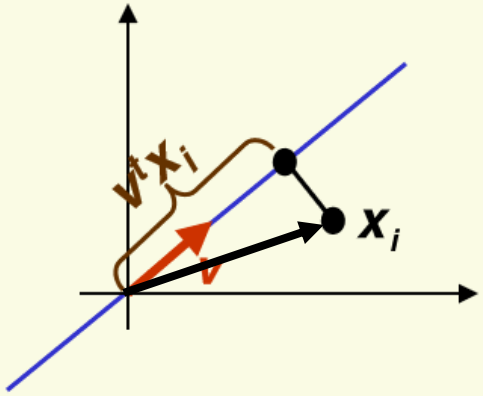
larger scatter:



smaller scatter:



지난 주 이슈 2 – LDA for dimension reduction



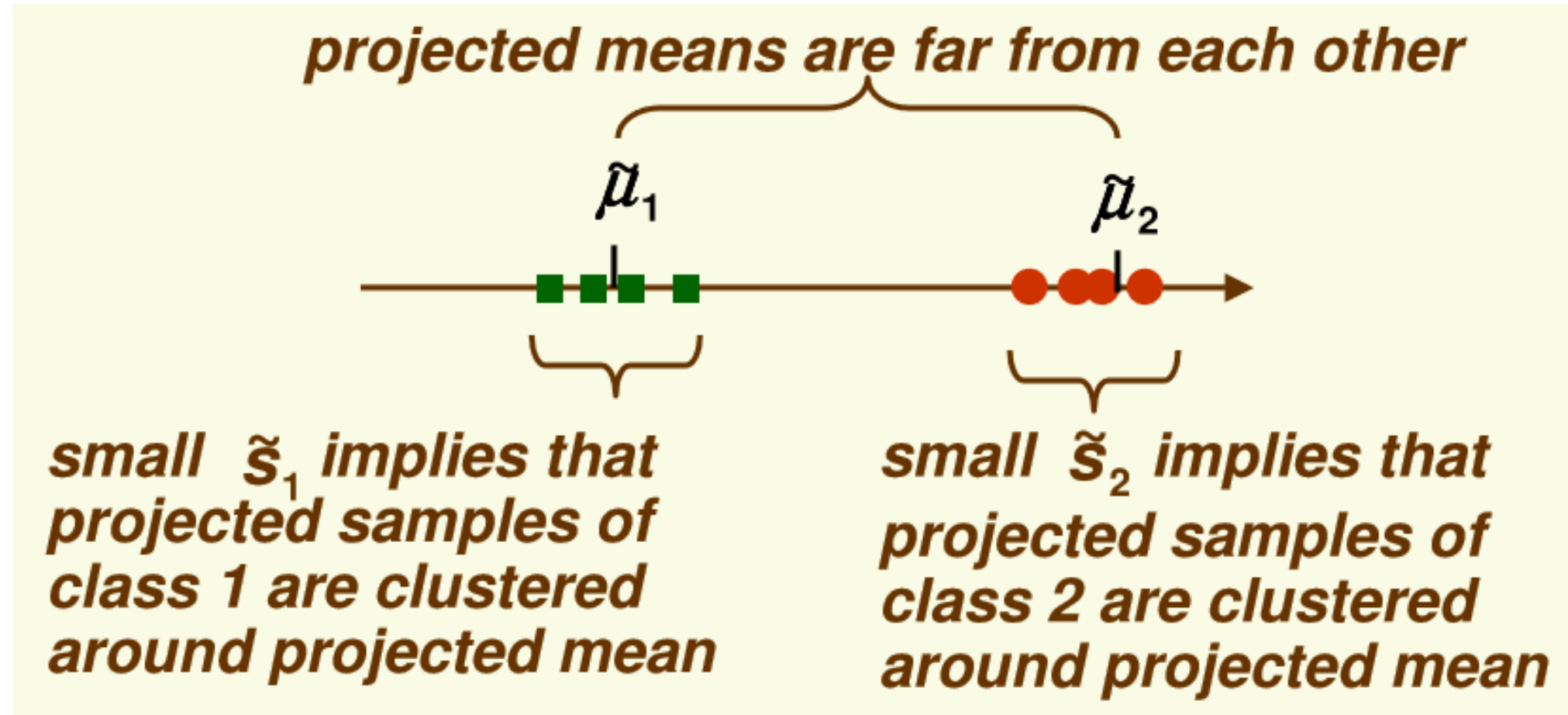
want projected means are far from each other

$$J(v) = \frac{(\tilde{\mu}_1 - \tilde{\mu}_2)^2}{\tilde{s}_1^2 + \tilde{s}_2^2}$$

want scatter in class 1 is as small as possible, i.e. samples of class 1 cluster around the projected mean $\tilde{\mu}_1$

want scatter in class 2 is as small as possible, i.e. samples of class 2 cluster around the projected mean $\tilde{\mu}_2$

지난 주 이슈 2 – LDA for dimension reduction



지난 주 이슈 2 – LDA for dimension reduction

$$J(\mathbf{v}) = \frac{(\tilde{\mu}_1 - \tilde{\mu}_2)^2}{\tilde{\mathbf{s}}_1^2 + \tilde{\mathbf{s}}_2^2} = \frac{\mathbf{v}^t \mathbf{S}_B \mathbf{v}}{\mathbf{v}^t \mathbf{S}_W \mathbf{v}}$$

$J(\mathbf{v})$ 를 \mathbf{v} 에 대해 미분 = 0 이 후, 수식 정리했을 때,
아래와 같이 고유 값 문제가 된다.!

$$\Rightarrow \underbrace{\mathbf{S}_B \mathbf{v} = \lambda \mathbf{S}_W \mathbf{v}}_{\text{generalized eigenvalue problem}}$$

Thank you!