# Module 5: Resampling Methods

Lecture 11
Feb 8th, 2023
Ch 5.1.3-4:

- LOOCV

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

- Leverage — *outlier*

- K-fold CV

$$CV_{(K)} = \sum_{k=1}^{K} \frac{n_k}{n} MSE_k \quad \Longleftarrow$$

$$K = n$$

- Hybrid between validation set and LOOCV

$$\text{CV}_{(K)} = \sum_{k=1}^{K} \frac{n_k}{n} \text{MSE}_k$$

LOOCV

10-fold CV

Blue — True error

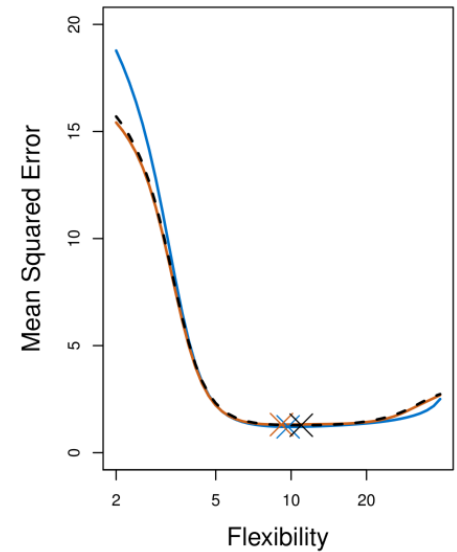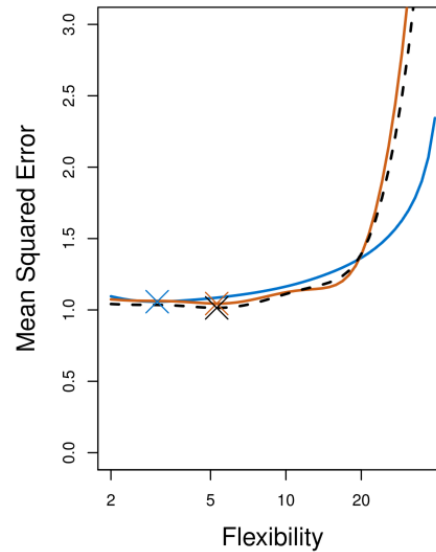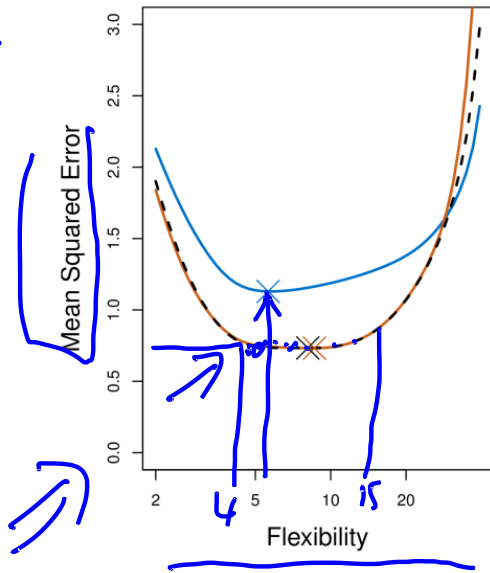Dashed — LOOCV estimate

Orange — 10-fold CV

- We divide the data into *K* roughly equal-sized parts

Compute

$$CV_K = \sum_{k=1}^{K} \frac{n_k}{n} Err_k$$

where $Err_k = \sum_{i \in C_k} I(y_i \neq \hat{y}_i)/n_k.$

# Example: GWAS

$5000 \gg 50$

- Sample size $n = 50$, number of predictors (SNPs) $p = 5000$; Predictor heart attach after age of 60.

- Step 1: Select the top 100 predictors having the largest correlation with the class labels

- Step2: We then apply a classifier (logistic regression) using only these 100 predictors

How do we estimate the test set performance of this classifier?

ATCG → Grace

Clark

Terence Tao

Alan Turing

P1 | ATCGATTA
P2 | ATCGTTTA

SNP

5000

Selected set
of predictors

Predictors

Outcome

Y

Samples

CV folds

Tests

100

No！！！

Selected set of predictors

Predictors

Outcome

Samples

CV folds

information leaky

*CV*

- Bootstrap is used to quantify the uncertainty associated with an estimator or machine learning method.

It can provide an estimate of the standard error of a coefficient. Then we can do hypothesis test and confidence interval.

$$\boxed{\hat{\beta}}\, ? \qquad \boxed{\hat{f}}\, ?$$

← TSLA

← KO

- To invest two stocks that yield return of X and Y, where X and Y are random.
- We will invest a fraction $\alpha$ of our money in X and the rest in Y.
- Assume both stock have the same average return over the years. What criteria should we use for allocate the investment?

$$\frac{E(X) = E(Y)}{\alpha}$$

$$Var(X) \gg Var(Y)$$

# Bootstrap： Finance Example

- We want to minimize the total risk or variance of our investment.

$$\text{Var}(\alpha X + (1-\alpha)Y).$$

- The solution is

$$\alpha = \frac{\sigma_Y^2 - \sigma_{XY}}{\sigma_X^2 + \sigma_Y^2 - 2\sigma_{XY}}$$

where $\sigma_X^2 = \text{Var}(X), \sigma_Y^2 = \text{Var}(Y),$ and $\sigma_{XY} = \text{Cov}(X,Y).$

*Handwritten annotations:*

$X = Y$

$\text{Var}(X+Y) \neq \text{Cov}(X,Y)$

$\text{Var}(2X) = 2^2 \text{Var}(X)$

$\text{Var}(x)$

$\text{Cov}(X,y)$

$(X+Y) = (1\ 1)\begin{pmatrix} X \\ Y \end{pmatrix}$

$\text{Var}\left( (1\ 1)\begin{pmatrix} X \\ Y \end{pmatrix} \right)$

$= (1\ 1) \text{Var}\left(\begin{pmatrix} X \\ Y \end{pmatrix}\right)\begin{pmatrix} 1 \\ 1 \end{pmatrix}$

$= (1\ 1)\begin{bmatrix} \text{Var}(X) & \text{Cov}(X,Y) \\ \text{Cov}(X,Y) & \text{Var}(Y) \end{bmatrix}\begin{pmatrix} 1 \\ 1 \end{pmatrix}$

$= Var(X) + Var(Y) + 2\, Cov(X, Y)$

- But the values of $\sigma_X^2$, $\sigma_Y^2$, and $\sigma_{XY}$ are unknown.
- We can compute estimates for these quantities, $\hat{\sigma}_X^2$, $\hat{\sigma}_Y^2$, and $\hat{\sigma}_{XY}$, using a data set that contains measurements for $X$ and $Y$.
- We can then estimate the value of $\alpha$ that minimizes the variance of our investment using

$$\hat{\alpha} = \frac{\hat{\sigma}_Y^2 - \hat{\sigma}_{XY}}{\hat{\sigma}_X^2 + \hat{\sigma}_Y^2 - 2\hat{\sigma}_{XY}}.$$
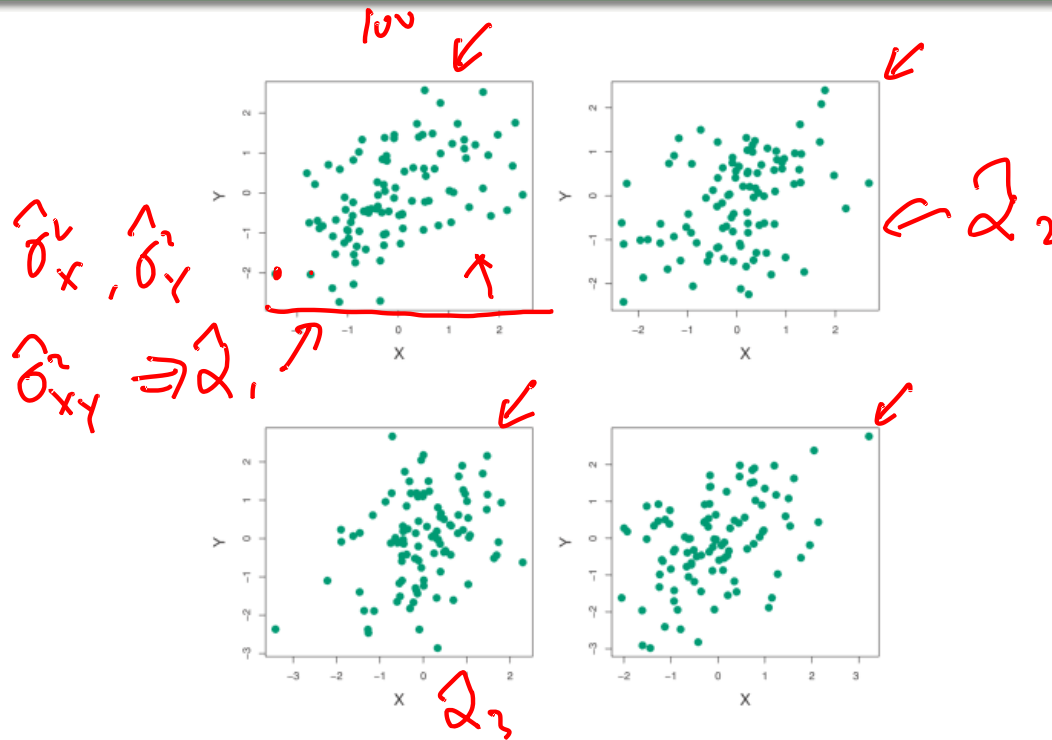
FIGURE 5.9. *Each panel displays* 100 *simulated returns for investments* X *and* Y. *From left to right and top to bottom, the resulting estimates for* α *are* 0.576, 0.532, 0.657, *and* 0.651.

- To estimate the standard deviation of $\hat{\alpha}$, we repeated the process of simulating 100 paired observations of $X$ and $Y$, and estimating $\alpha$ 1,000 times.
- We thereby obtained 1,000 estimates for $\alpha$, which we can call $\hat{\alpha}_1, \hat{\alpha}_2, \ldots, \hat{\alpha}_{1000}$.

  - For these simulations the parameters were set to $\sigma_X^2 = 1, \sigma_Y^2 = 1.25$, and $\sigma_{XY} = 0.5$, and so we know that the true value of $\alpha$ is 0.6 (indicated by the red line).

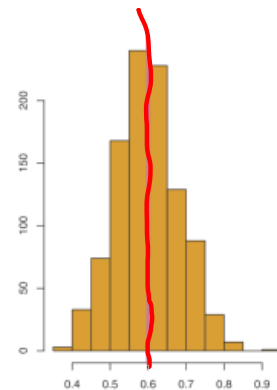- The mean over all 1,000 estimates for $\alpha$ is

$$\bar{\alpha} = \frac{1}{1000} \sum_{r=1}^{1000} \hat{\alpha}_r = 0.5996,$$

very close to $\alpha = 0.6$, and the standard deviation of the estimates is

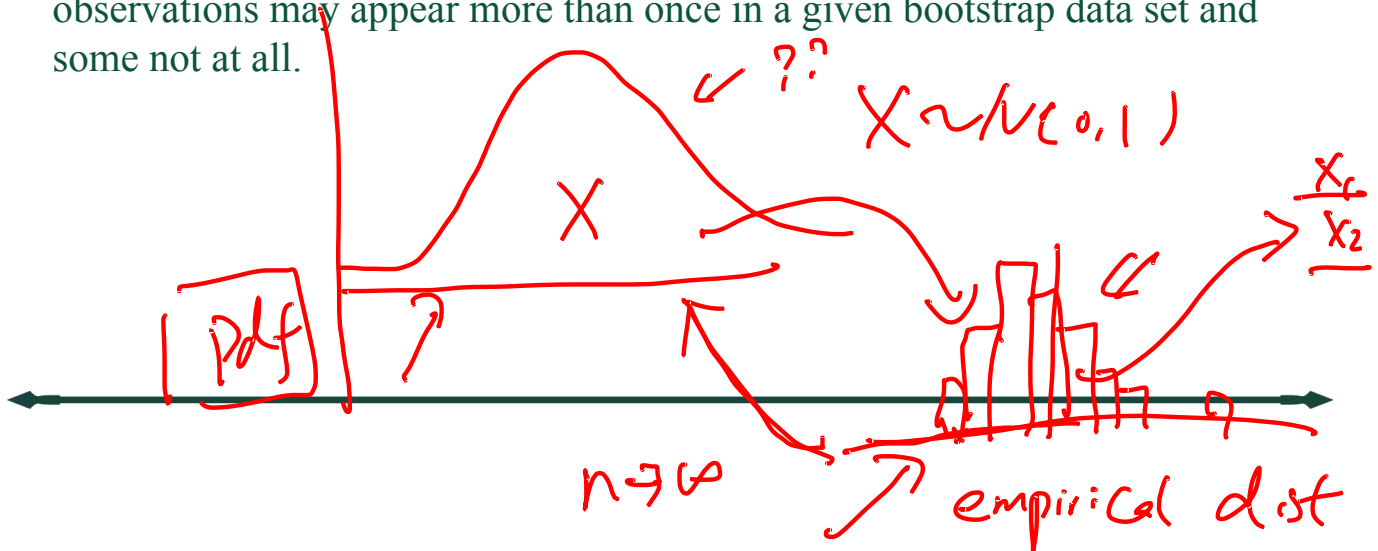$$\sqrt{\frac{1}{1000 - 1} \sum_{r=1}^{1000} (\hat{\alpha}_r - \bar{\alpha})^2} = 0.083.$$

- This gives us a very good idea of the accuracy of $\hat{\alpha}$: $\text{SE}(\hat{\alpha}) \approx 0.083.$

- The procedure outlined above cannot be applied, because for real data we cannot generate new samples from the original population.
- The bootstrap approach allows us to use a computer to mimic the process of obtaining new data sets.
- Rather than repeatedly sampling from the population, we instead obtain distinct data sets by repeatedly sampling observations from the original data set with replacement.
- Each of these "bootstrap data sets" is created by sampling with replacement, and is the same size as our original dataset. Some observations may appear more than once in a given bootstrap data set and some not at all.

$n$

$X \sim N(0,1)$

$\frac{X_1}{X_2}$

pdf

$X$

$n \to \infty$

empirical dist

$n = 3$

$n = 3$

| Obs | X | Y |
|-----|-----|-----|
| 3 | 5.3 | 2.8 |
| 1 | 4.3 | 2.4 |
| 3 | 5.3 | 2.8 |

$Z^{*1}$ → $\hat{\alpha}^{*1}$

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

$Z^{*2}$ → $\hat{\alpha}^{*2}$

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

Population

Original Data (Z)

$Z^{*B}$

| Obs | X | Y |
|-----|-----|-----|
| 2 | 2.1 | 1.1 |
| 2 | 2.1 | 1.1 |
| 1 | 4.3 | 2.4 |

→ $\hat{\alpha}^{*B}$

# A General Picture for the Bootstrap

Real World

Random Sampling — Data

Population $P$ — $Z = (z_1, z_2, \ldots z_n)$

Estimate $f(Z)$

Bootstrap World

Random Sampling — Bootstrap dataset

Estimated Population $\hat{P}$ — $Z^* = (z_1^*, z_2^*, \ldots z_n^*)$

Bootstrap Estimate $f(Z^*)$

app

**MICHIGAN STATE**
**U N I V E R S I T Y**

*Bonus Quiz*

- Can we use one boostrap dataset as training and the *?*
  original data set as test set?

*B₁*

*n 100*   *Boostrap* →   *n = 100*

*training*

*?*

*testing*

(20 pts) For a classification problem with $K = 2$ ($Y \in \{0, 1, 6\}$), we know the oracle classifier is

$$C(x) = j, \quad \text{if } p_j(x) = \max\{p_0(x), p_1(x)\},$$

which is based on the loss with equal weight for Type I and II error. If we know Type I erorr will cost $1000 while Type II error will cost $3000. Derive the new oracle classifier which minimizes this cost.

classification     logistic   regression

MICHIGAN STATE
U N I V E R S I T Y

- When sample size $n$ is large, we know a bootstrap dataset will contain $1 - e\text{-}1 = 63.2\%$ of original data. Write a code to demonstrate it using $n = 1000000$

Loss function

① Supervised ⎡ regression ⎡ ① KNN
                           ⎣ ② Linear regression

Classification ① KNN
                ② Logistic.

② unsupervised

→ Test the best model?

Bias-Variance trade-off

$\hat{\beta} = ?$    $SE(\hat{\beta})$? CI, HT

Linear  regression ←    X is qualitative, interaction terms
                         beyond linearity?