

# Resampling Methods

**Resampling Methods:** tool in modern statistics, involving repeatedly drawing samples from a training set & refitting a model on each sample

**Example:**

We can draw samples from training & repeatedly fit a linear model & see how each differs

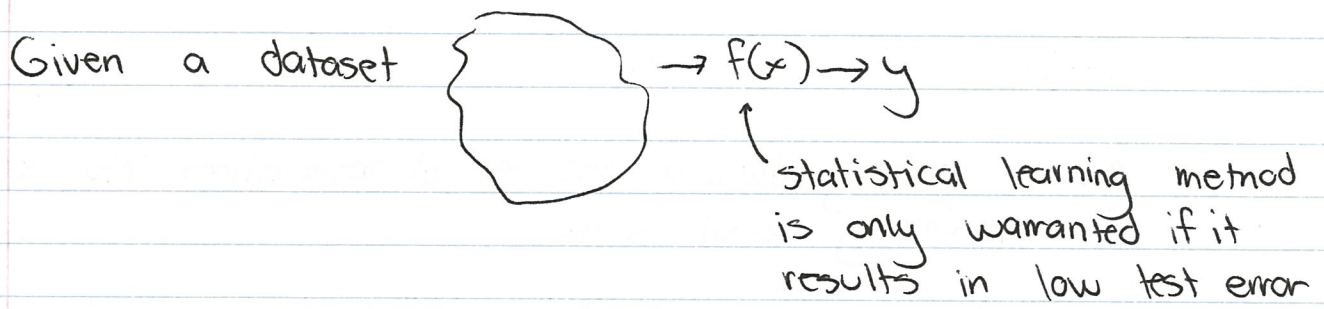
two most common resampling

1. cross validation
2. bootstrap

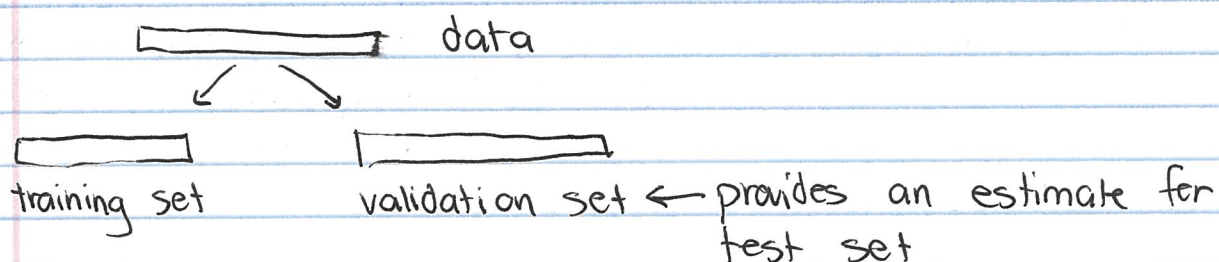
**Model Assessment:** the process of evaluating a model's performance

**Model Selection:** the process of selecting the proper level of flexibility for a model

## Cross-Validation



## Validation Set Approach



you can do multiple splits that are all different  
- can be highly variable

## Leave-One-Out Cross-Validation

- single observation is used for validation

$$CV(n) = \frac{1}{n} \sum MSE$$

Benefits:

- no randomness in training/validation
- doesn't overestimate test error

## k-Fold Cross Validation

involves randomly dividing the set of observations into  $k$  groups of approximately equal size

1st fold is treated as the validation set, & model is fit on  $k-1$  folds

repeat procedure  $k$  times

this procedure results in  $k$  estimates of test errors  $mse_1, mse_2, \dots, mse_k$

$$CV(k) = \frac{1}{k} \sum MSE$$

When we perform cross-validation our goal might be to determine how well a given statistical learning procedure can be expected to perform on independent data

↓ we want to eventually choose the model with the lowest test error

## Bias Variance Tradeoff for $k$ Fold Cross Validation

$k$  Fold CV gives an accurate estimate of test error

when  $k$

has bias

when your test set is small  
you can have overestimates of  
test error rate



## The Bootstrap

extremely powerful statistical tool that can be used to quantify uncertainty associated with a given estimator or statistical learning method

↓ gives you a standard error for your estimates

Example:

value that minimizes risk

$$\alpha = \frac{\sigma^2_y - \sigma_{xy}}{\sigma^2_x + \sigma^2_y - 2\sigma_{xy}}$$

we get estimates of parameter from data

we estimate  $\alpha$  on simulated data

we get 1000 estimates for  $\alpha$ ,  $\hat{\alpha}_1, \hat{\alpha}_2, \dots, \hat{\alpha}_{1000}$

$$\bar{\alpha} = \frac{1}{1000} \sum_{i=1}^{1000} \hat{\alpha}_i$$

- get  $\mu$  &  $\sigma$

from different samples we expect  $\alpha$  to differ by 0.08

$z_1 \quad n=3$

	x	y
1	4.3	2.4
2	2.1	1.1
3	5.3	2.8



3  
1  
3

2  
3  
1

2  
2  
1

