

$G_h \Rightarrow 5$

## Resampling Methods

Data is currency of applied Machine Learning. Therefore, it is important that it is both collected & used safely.

Data resampling refers to method of economically using a collected dataset to improve the estimate of the population parameters & help to quantify the uncertainty of the estimate.

Resampling methods are an indispensable tool in modern statistics.

For example, in order to estimate the variability of a linear regression fit, we can repeatedly draw different samples from ~~training~~ data, fit a linear regression to each new sample, and thus examine the extent to which the resulting fits differ. Such an approach may allow us to obtain information that would not be available from fitting the model only once using training examples.

## Cross-Validation:

We had a simple approach to test our data on the test set by dividing data into 2 sets → One for training & another for testing. This though requires a lot of data & that is not the case a lot of times. In those times resampling methods came into the picture.

→ We hold out <sup>random</sup> some data of our training set & estimate our ~~from~~ testing errors on that subset of data & repeat this.

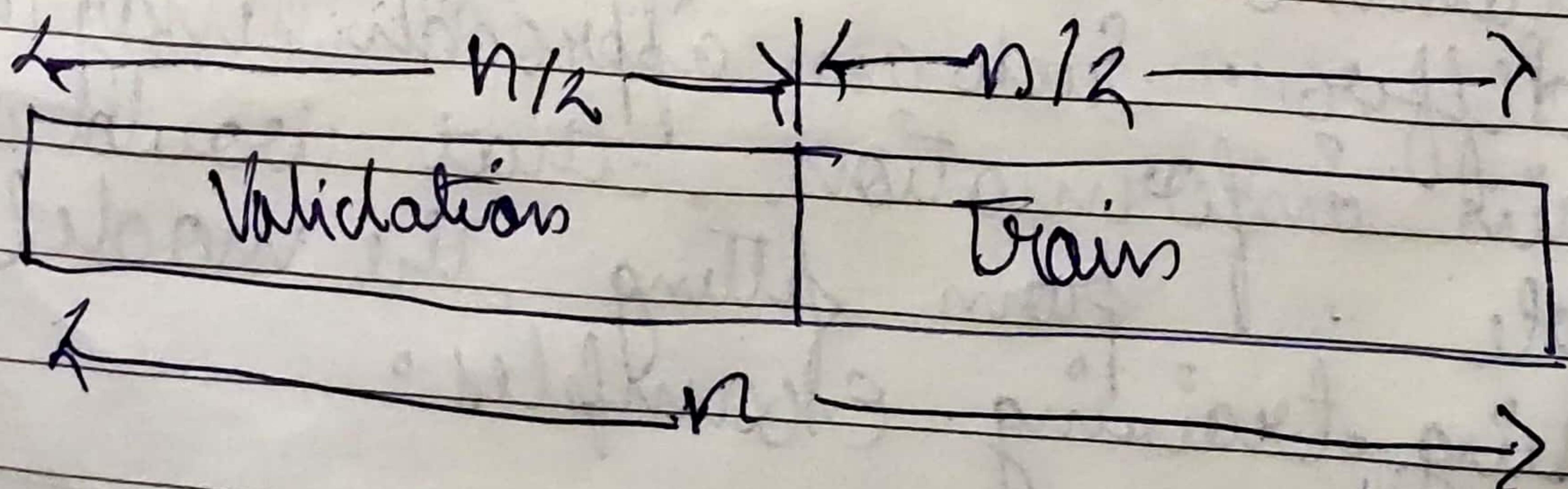
1.

### The validation Set Approach:

→ Here, we randomly divide available set of samples into 2 parts, training set & validation set.

→ The model is fit on training set & the fitted model is used to predict the responses for the observations of validation set.

→ The resulting validation error provides an estimate of the test error.



This has 2 potential drawbacks.

1. The validation set error estimates have high variability on the estimates of test errors.  
Because, the more the data, the better is the training & this approach kind of loses ~~its~~ because of this reason.

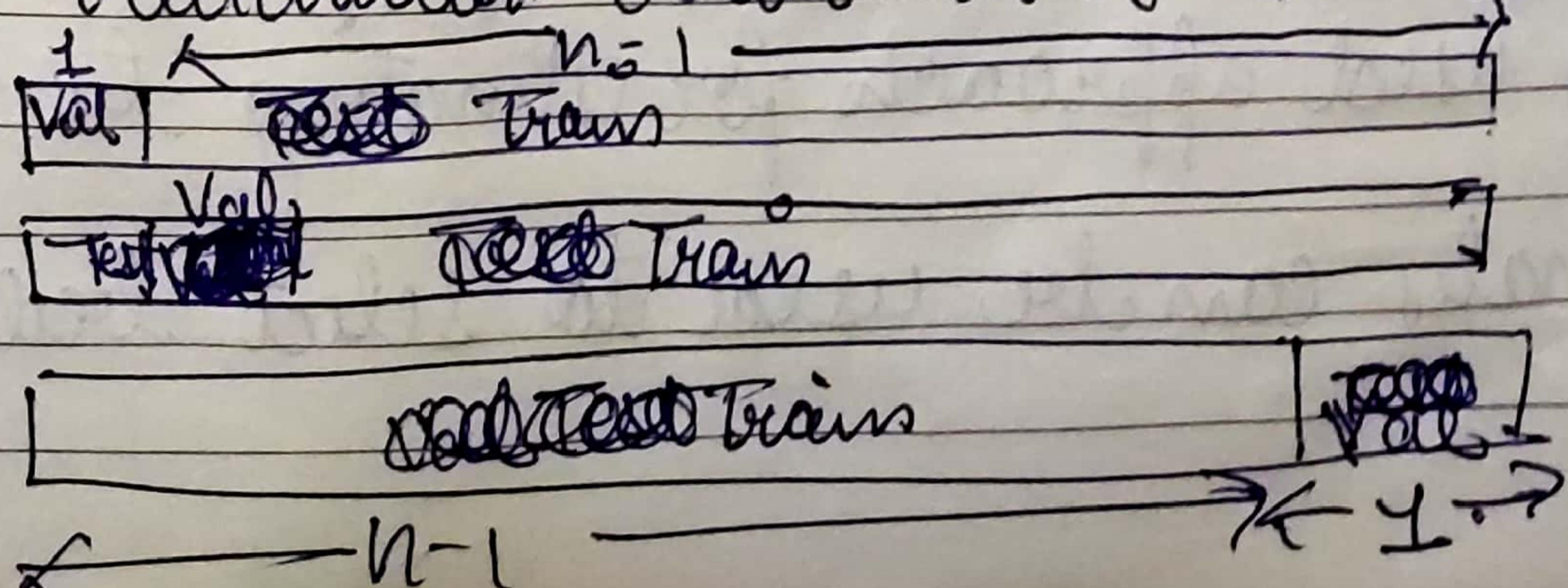
a.) ~~Leave One Out Cross-Validation~~

This is closely related to validation set approach.

→ Here we split data into 2 set, one set contains  $n-1$  values while others one contains only one value.

→ we apply this methods in two different times, on a validation set.

The validation error is calculated by averaging the validation errors over  $N = 1$



## Advantages Over Validation Set Approach

1. Firstly, there is far less bias because it fits data in over a dataset of "n-1" values.
2. It doesn't overestimate the test error due to lack of training or test examples.  
 This will almost yield same results every time.

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

where  $h_i^*$  is , high leverage points learnt in  
ch-3. ~~Ch-3~~ Pg 98

### Disadvantage:

1. High Computational expensive.
2. Estimate of each fold is highly correlated hence their average can have high variance.

### (3) K-Fold Cross Validation :

An alternative to LOOGV is K-fold CV.

- Widely used approach for estimating test error.
- Estimates can be used to select best model.

→ Idea is to randomly divide the data into  $K$ -equal sized parts. We leave out part  $K$ , fit the model to other  $K-1$  parts (Combined) & then obtains for left out  $K^{\text{th}}$  part.

→ This is done in two ways & their results are combined

$\frac{n}{K}$	$\frac{n}{K}$	$\frac{n}{K}$	$\frac{n}{K}$	$\frac{n}{K}$
Val	Test			

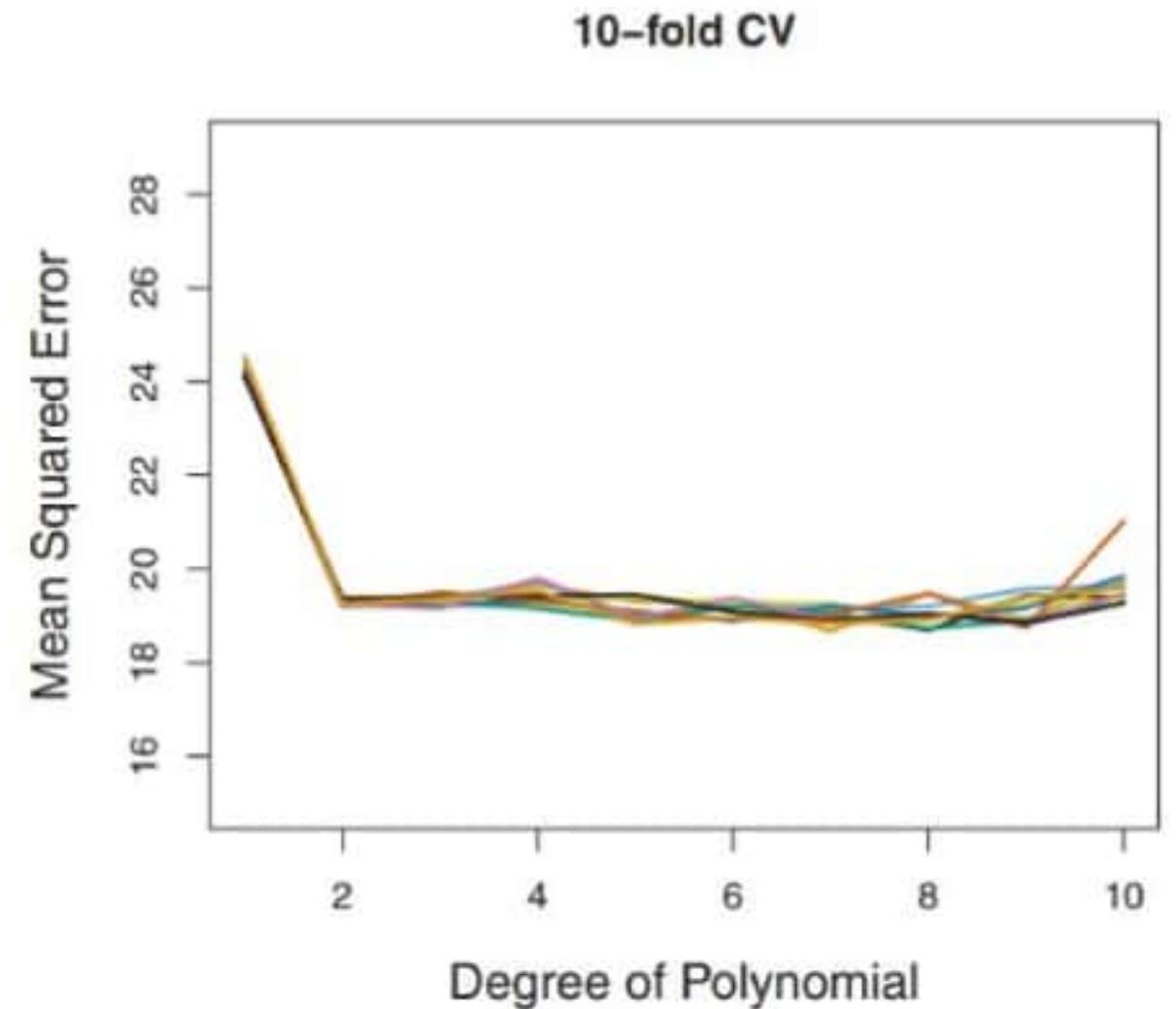
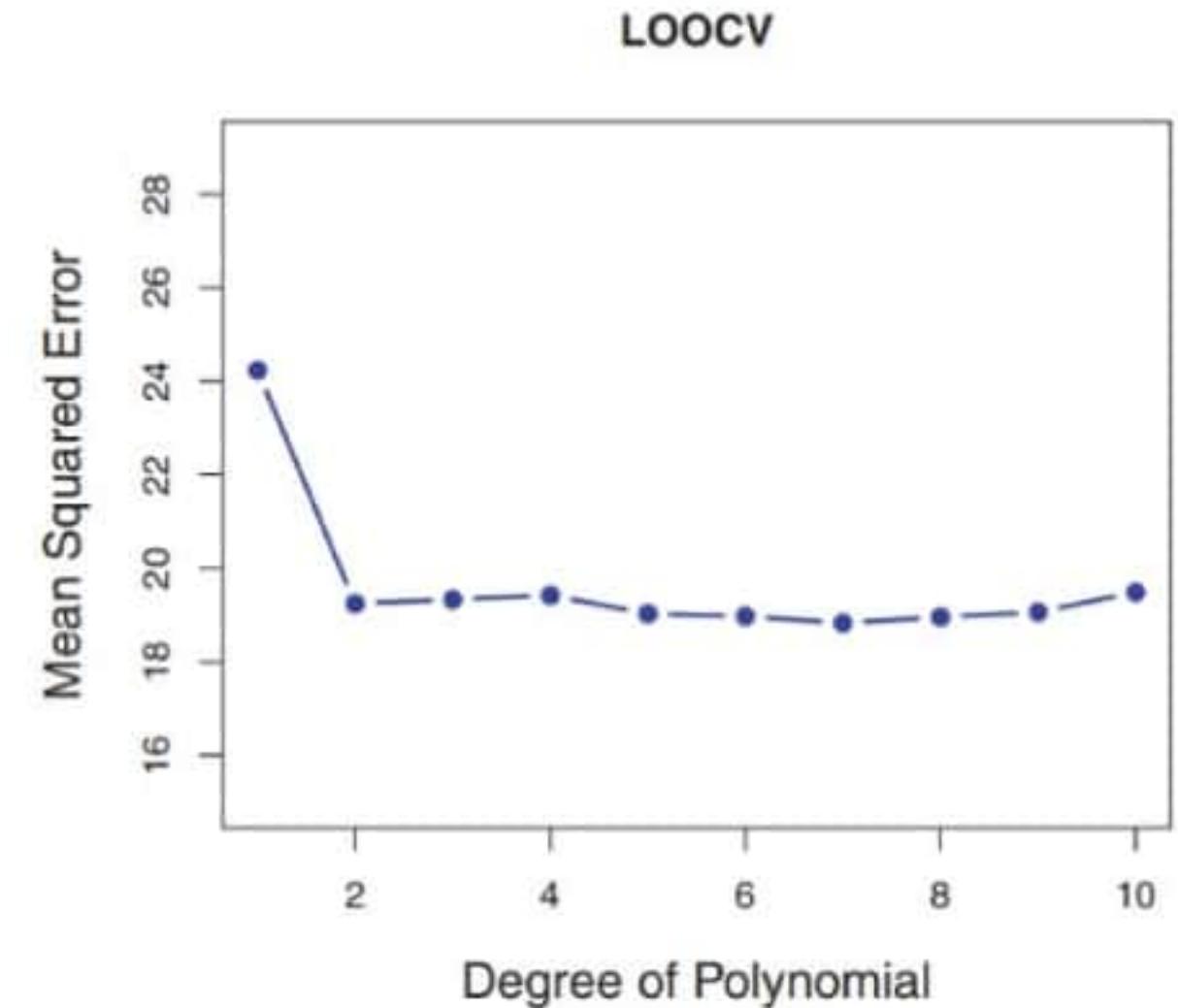
	Val			

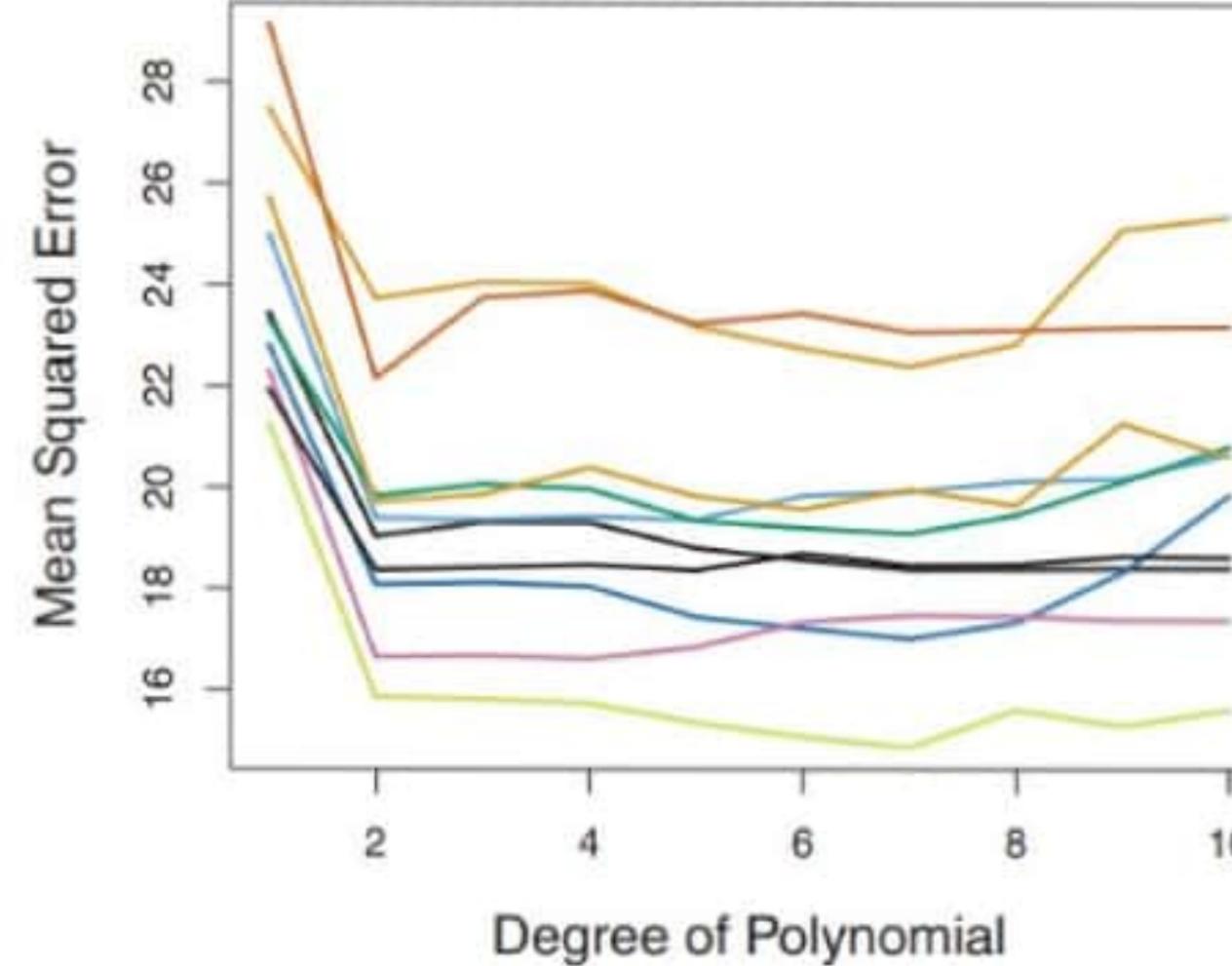
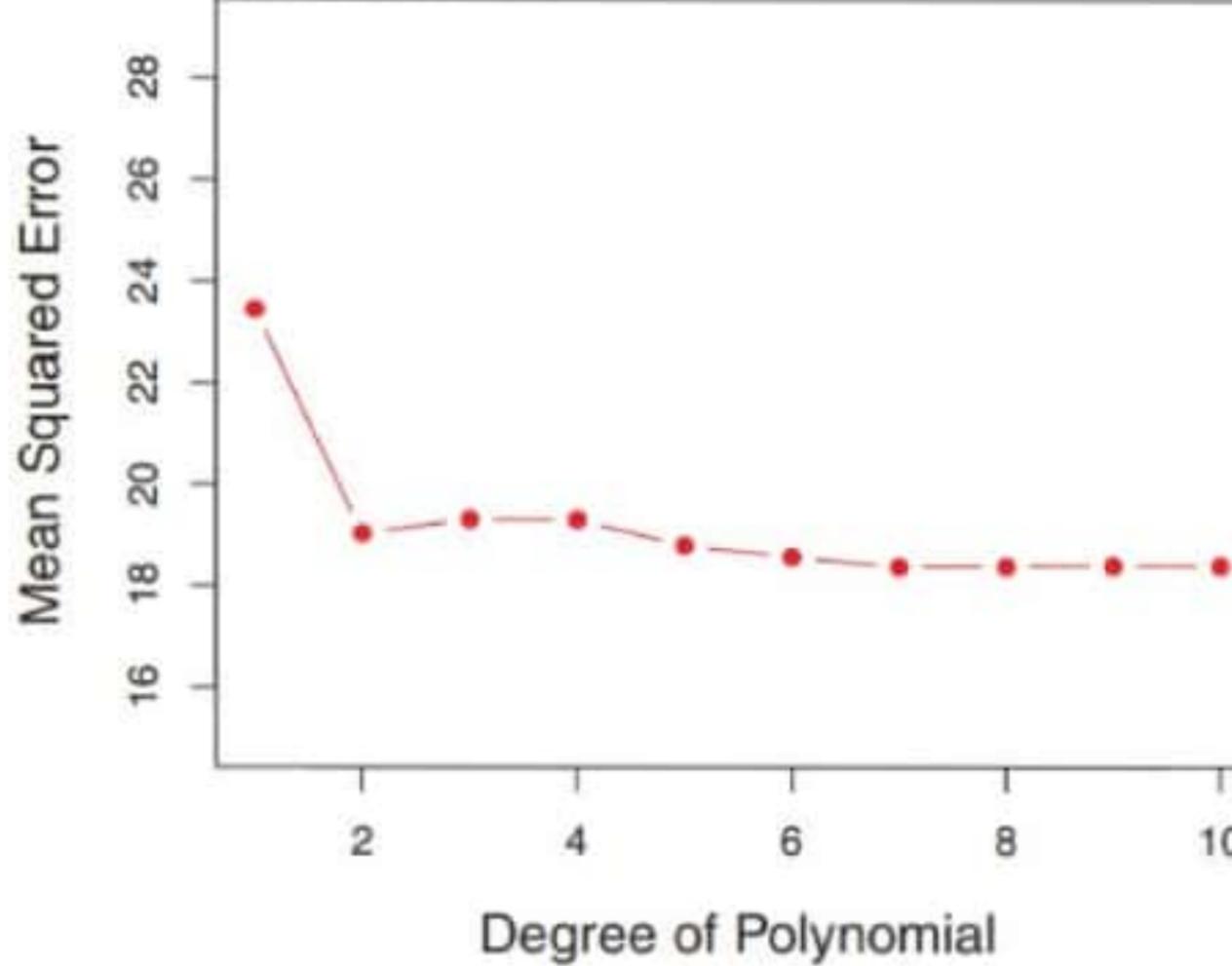
				Val

5 Folds

### Disadvantages:

- Since each training set is  $(K-1)/K$  in size, the estimates of predictions errors will be higher & have high bias.
- When  $K=n$ , bias decreases, but the estimates have high variance.
- $K=5$  to  $10$  provides a good compromise for bias variance tradeoff.





## Cross Validation: right and wrong

→ Consider a simple classifier applied to some two class data:

1. Starting with 5000 predictors & 50 samples, find 100 predictors with largest correlation with class labels.

2. We then apply a classifier using only these 100 predictors!

\* Can we apply CV in step 2, forgetting about Step-1?

NO!

→ This would ignore the fact that in Step 1, the procedure has already seen the labels of training data & made use of them. This is the form of training & must be used in validation process.

→ It is easy to stimulate realistic data without the class labels independent of outcomes so true test error = 50%, but if we ignore Step-1, CV error estimate will be

Right Way:

Apply CV to both steps.

# THE BOOTSTRAP

The bootstrap is a flexible and powerful statistical tool that can be used to quantify the uncertainty associated with a given estimator or statistical learning method.

for example, it can provide an estimate of the Standard Error of a coefficient, or a confidence interval for that coefficient.

EXAMPLE:

- Suppose that you wish to invest a fixed sum of money in 2 finances that yields  $X$  &  $Y$ .
- let  $\alpha$  be fraction of money invested in  $X$ .
- so, we'll try to minimize  $(\text{Var}(\alpha X + (1-\alpha)Y))$

To minimize the risk  $\alpha$  is given by:  $\frac{\partial^2 - \sigma_{XY}}{\partial \alpha^2 + \sigma_X^2 - 2\sigma_{XY}}$

But,  $\sigma_X$ ,  $\sigma_Y$  &  $\sigma_{XY}$  are unknown. So, we'll compute estimates for them.

$$\hat{\alpha} = \frac{\hat{\sigma}_X^2 - \hat{\sigma}_{XY}^2}{\hat{\sigma}_X^2 + \hat{\sigma}_Y^2 - 2\hat{\sigma}_{XY}}$$

- To estimate  $\hat{\alpha}$ , will simulate 100 paired samples of  $X$  &  $Y$  1000 times.

Now,  $\hat{\alpha}$  will be very close to  $\alpha$ .

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

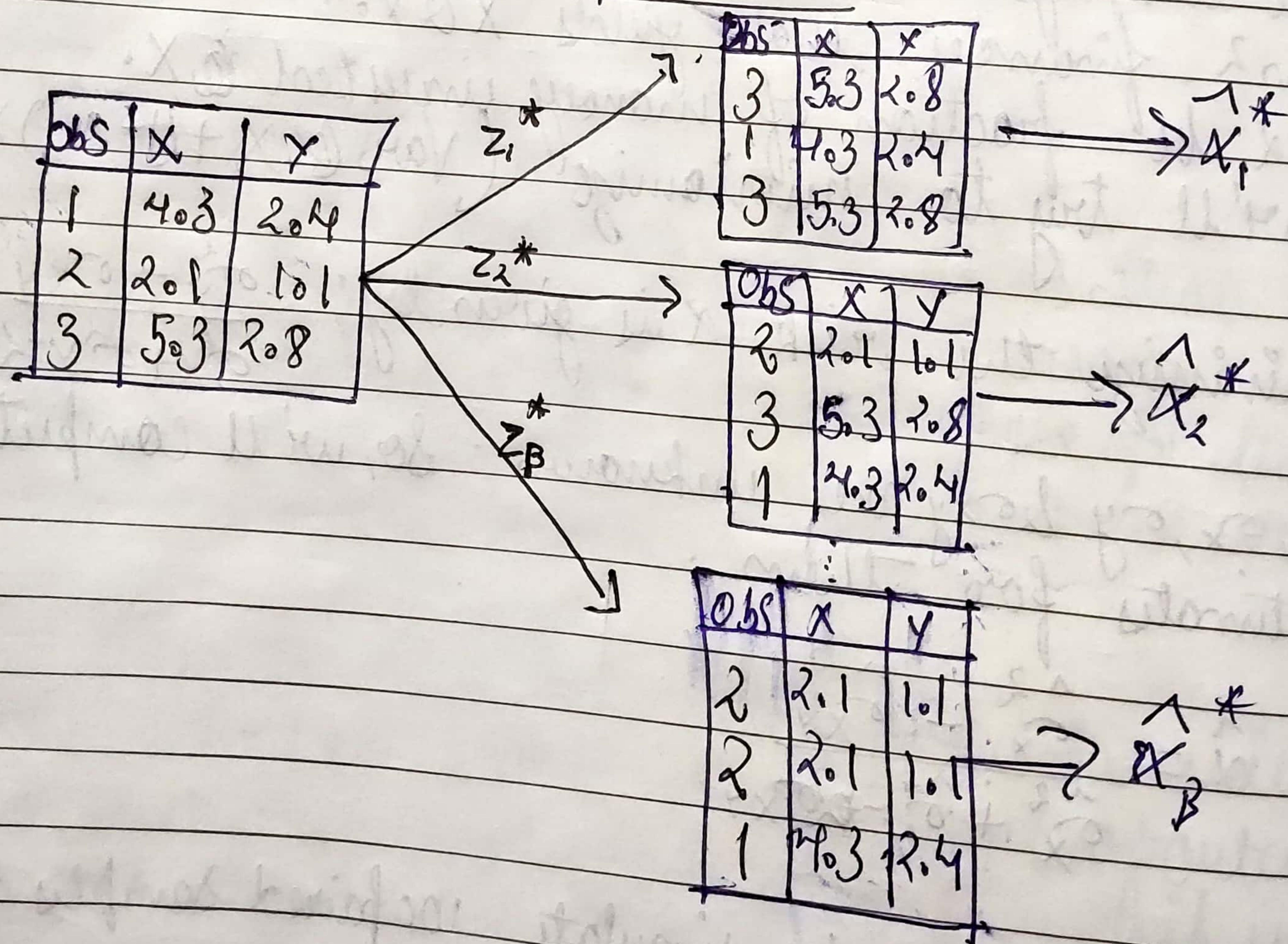
We also perform Bootstrap to calculate the standard error of these estimates.

$$SE = \sqrt{\frac{1}{1000-1} \sum_{j=1}^{1000} (\bar{x}_j - \bar{x})^2} = 0.0853$$

So, this gives very good estimate of accuracy.

HOW DO WE DO THIS IN REAL WORLD?

Rather than repeatedly sampling independent datasets from population we instead obtain distinct datasets by repeated sampling of original data with replacement.



## BOOTSTRAP IN GENERAL:

In complex data situations, figuring out appropriate way to generate bootstrap can require some thought.

For Example: if a data is in time-series, we can simply sample the population with replacement.  
Because, we ~~are~~ our data points are dependent on previous datapoints & we expect them to be correlated.

\* So, rather than doing that, we can divide data in blocks & then sample that.

## OTHER USES OF BOOTSTRAP:

- Primary use is to calculate Standard Error of the estimate.
- Can be used to calculate the confidence interval of our data. (By estimates & Random Sampling we can create an empirical Sampling Distribution & then problem solved.)

Can Bootstrap estimate predict error?

No, because of a lot of overlap in train & validation data. We can rather use k-fold CV.