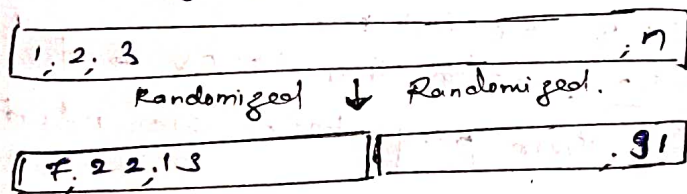# Resampling Methods.

## * Introduction:

- **Resampling method:** repeatedly drawing samples from a training set and fitting a model of interset on each sample in order to obtain additional Information about the model.

- **cross validation:** can be used to estimate the test error associated with a given statistical learning method in order to evaluate its performance or to select appropriate level of flexibility.

- **model assessment:** process of evaluating a model's performance

- **model selection:** process of selecting the proper level of flexibility.

- **bootstrap:** used in several contexts, most commonly to provide a measure of accuracy of a parameter estimate or of a given statistical learning method.

## 1. Cross validation.

### 5-1-1 : The validation set approach.

- it involves randomly dividing the available set of observations into two parts: a training set( used for model fitting) & a validation set/hold-out set (prediction) & calculating estimate of test error)

| 1, 2, 3 | ; n |
|---|---|

Randomized ↓ Randomized.

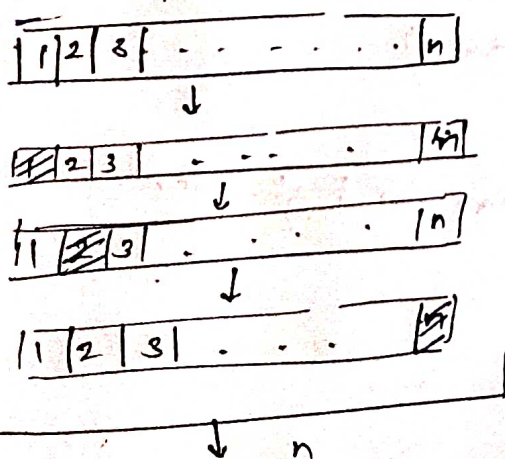| 7, 22, 13 | | .91 |
|---|---|---|

may tend to overestimate the test error rate.

**Drawbacks**

* validation estimate of the test error rate can be highly variable, depending on which observations included in train & test.

* if training set small - model may perform bad & validation set error

### 5.1.2 Leave One out cross validation. (LOOCV)

- attempts to addren drawbacks of cross validation set approach.

| 1 | 2 | 3 | . . . . . | n |

↓

| 1 | 2 | 3 | . . . | n |

↓

| 1 | 2 | 3 | . . . | n |

↓

| 1 | 2 | 3 | . . | n |

↓

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

$MSE_i := (y_i - \hat{y}_i)^2 \to$ not included in validation set the $y_i$th observation.

→ consider all $\binom{n}{1} =$ n possible leave one out cross validation sets then get the average : that is the required estimation.

* **Advantages** - less bias (on repetitions).
  - does not over estimate test error.
  - no randomness involved in the train-validation split.

Note: LOOCV has potential to be expensive for implementation.
For least squares linear/polynomial regression:

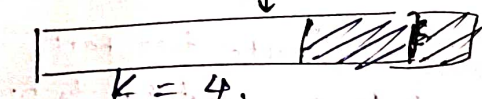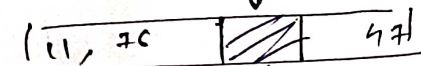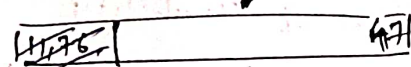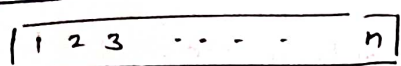$$CV_{(n)} = \frac{1}{n} \sum_{i=1}^{n} \frac{(y_i - \hat{y}_i)^2}{(1 - h_i)} \longrightarrow$$ This formulae does not hold true in general!

where: $\hat{y}_i$ : ith fitted value from the original least squares fit.

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

**5.1.3 : k-fold cross validation:** (1) Randomly split into k many non overlapping groups of equal size.
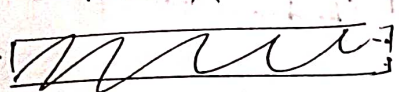


k = 4.

(2) $i^{th}$ group hold out as validation set. remaining 'k-1' used for training model.

(3) repeat process for holding out each group i : i = 1(1) k-1

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

generally : k = 5 or 16 (computation perspective)

→ very less variability compared to validation set approach

**5.1.4. Bias - Variance Trade off for k-fold Cross-Validation.**

- k fold cv gives more accurate estimate of the test error.
- no. of observations in k fold $\approx \frac{(k-1)n}{k}$ . → more then LOO CV but less than validation set approach ⇒ from perspective of bias reduction LOOCV is preferred to k-fold CV.
- LOOCV has higher variance than does k-fold cv with k<n.

→ LOOCV :- training sets almost identical - highly +vely correlated.
   - averaging output of these n fitted outputs.

k fold : - ten overlap in training sets ⇒ less correlated

var (mean of non corr. data) < var (mean of highly correlated data)

∴ Var (k fold) < Var (LOOCV)

⇒ * for k fold CV with k=5 or 10 : no high bias nor high variance

## 5.1.5 Cross validation on classification Problems.

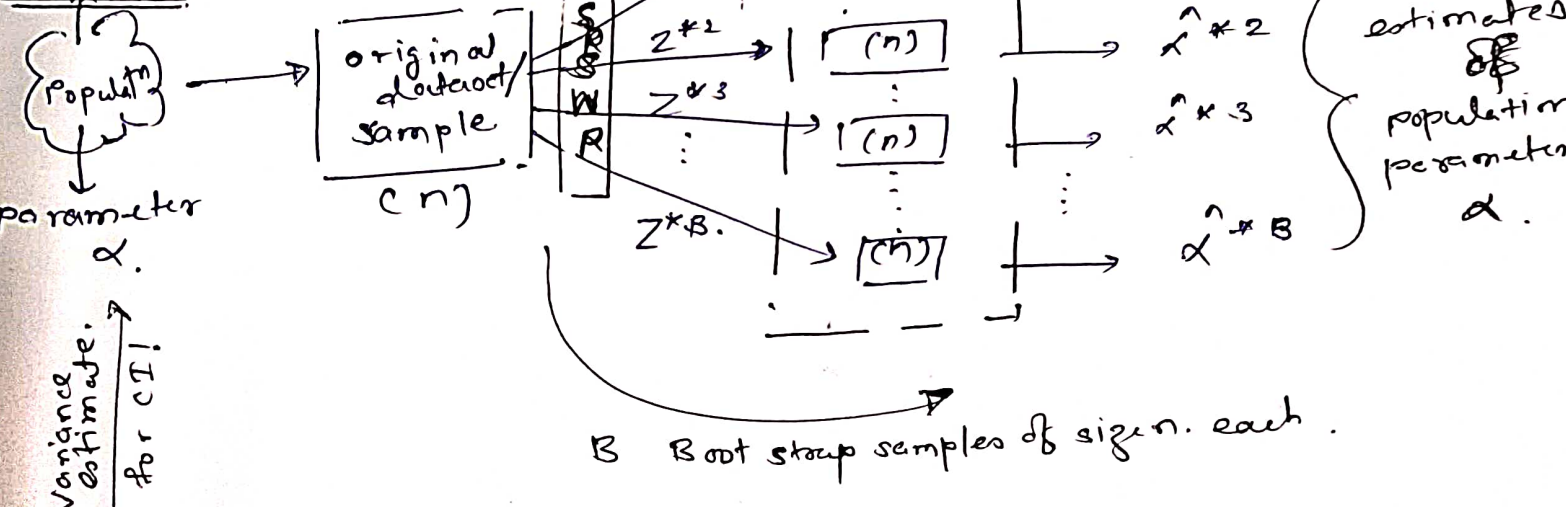- Regression: MSE , classification: error rate/ ~~part of~~ fraction of misclassified objects.

For LOOCV :
$$CV_{(n)} = \frac{1}{n} \sum_{i=1}^{n} Err_i \quad \text{where } Err_i = I(y_i \neq \hat{y}_i)$$

For K fold CV :
$$CV_{(k)} = \frac{1}{k} \sum_{i=1}^{k} Err_i$$

## * 2. Boot strap.

- widely applicable & extremely powerful statistical tool that can be used to quantify the uncerteinldy associated with a given estimator or statistical learning method.
- used when measure of variability is hard to obtain.
- rather than repeatedly obtaining independent data sets from the population we instead obtain distinct data sets by repeatedly sampling observations from the original data set.

Unknown!



B Boot strap samples of size n. each.

variance estimate. for CI!

$$SE_B(\hat{\alpha}) = \sqrt{\frac{1}{B-1} \sum_{r=1}^{B} \left( \hat{\alpha}^{*r} - \frac{1}{B} \hat{\alpha}^{*r} \right)^2}$$

* Advantage : - Can be used in almost all situations.
- No complicated mathematical ~~situal~~ formulae required.