Non-Parametric models Parameters keeps on adding.
↓
Doesn't mean it has not Parameters.

$$P(x) = \int_a P(x|a) \, \phi(a) \, da$$

Conjugate priors:

$$\boxed{\text{Gaussian}} \leftrightarrow \boxed{\text{Gaussion}}$$
↓
Data comes from gaussion & mean of gaussion
also comes from gaussion.

Multinomial Distribution: models the probability of
counts for each of "K"-sided die rolled n-times.
→ generalization of binomial distribution.

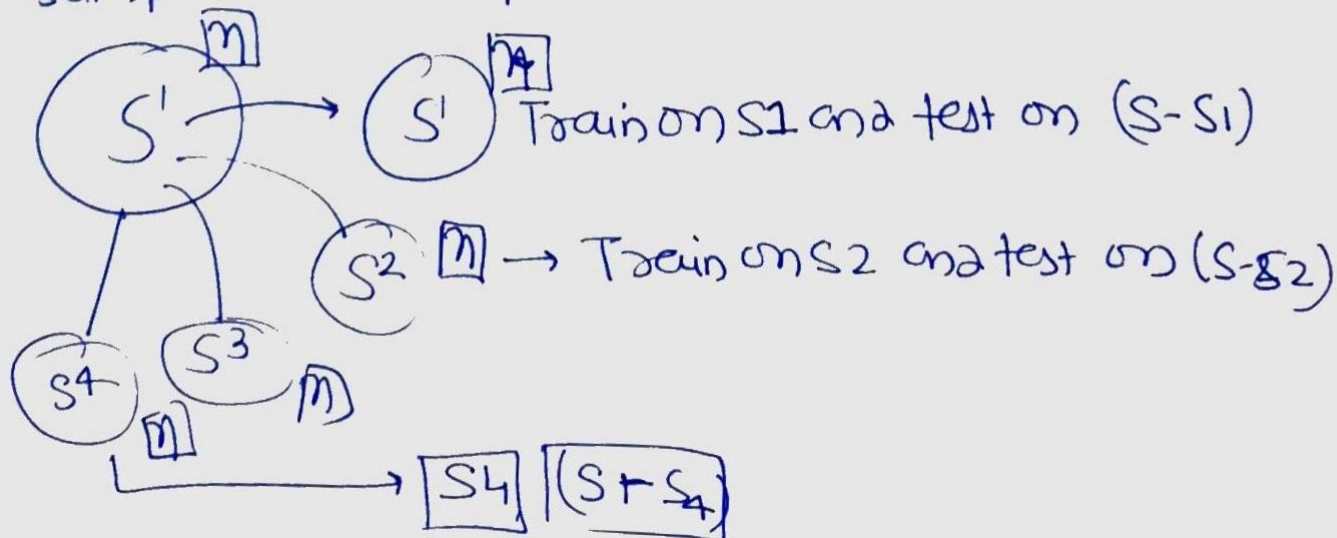## Lecture-42    Evaluation Measures 1

Why to divide Training datasets into multiple sets

- because - Variance of model decreases, if we train it
  on multiple datasets.

Given 'n' independent observations $Z_1, Z_2, \ldots Z_n$, each
with variance $\sigma^2$, the variance of mean $\bar{Z}$ of the
observation is given by $\dfrac{\sigma^2}{n}$.

## Bootstrapping & cross-validation
↓
sample with replacement.



Train on S1 and test on (S-S1)

Train on S2 and test on (S-S2)

S4 → (S + S4)

# K- FOLD CROSS VALIDATION

Which one gives better estimate? bootstrap/cross validation

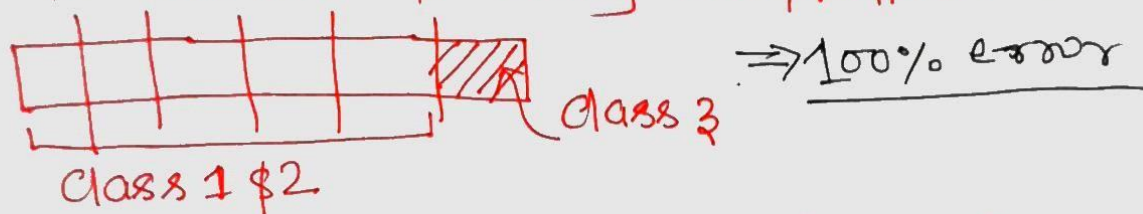No. of samples are sufficiently large : bootstrap. (67% Train)

" " " " " low : Cross validation

bootstrap is basically used to estimate the variance
of model. like (mean of distribution.)
variation of ⟩

$\boxed{K}$ , PROCESS of creating 'k' → (5 $\overset{or}{\text{to}}$ 10)

larger k ⇒ More reduction in variance

but increasing k ⇒ More biasness. , ⇒ we need
to choose 'k' using bias/variance reduction.

⇒ 100% error

class 3

class 1 & 2

STRATIFIED SAMPLIN : AFTER SAMPLING, Data
should have same proportion in all Folds.
before CV.

IMBALANCE DATASET : Fewer No. of folds.
since minority class population is very low.
our models needs to have good amount of minority
class in Training DATA.