

Ensemble Method

① Simple Ensemble Techniques

1.1 Man Voting

- used for classification problem mostly.
- multiple models are used to make prediction for each data points.
- Prediction \Rightarrow vote "GO BY MAJORITY"

Eg Input ($x_1, x_2, x_3, \dots, x_n$)

Outputs

model(DTC)

5

model(KNN)

4

model(LR)

5

model(SVM)

4

model(Ridge Regress)

4

~~final pred~~
~~13~~
~~final pred~~

final predⁿ

man voting

4

averaging

4.40

Weighted avg

4.41

| | model1 | model2 | model3 | model4 | model5 | final |
|--------|--------|--------|--------|--------|--------|-------|
| act | 0.23 | 0.23 | 0.18 | 0.18 | 0.18 | |
| rating | 5 | 4 | 5 | 4 | 4 | 4.41 |

$$5 \times 0.23 + 4 \times 0.23 + 5 \times 0.18 + 4 \times 0.18 + 4 \times 0.18 = 4.41$$

1.2 Averaging

1.3 Wt Average

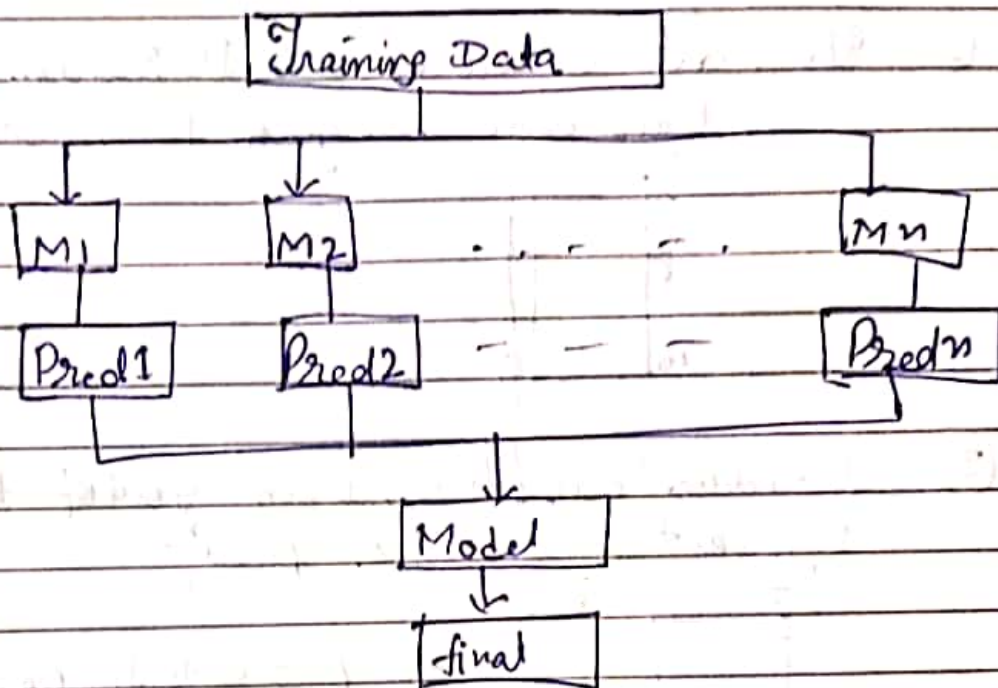
Can be
used for
Regression

2. Advanced Ensemble Techniques

- Stacking
- Blending
- Bagging
- Boosting

2.1 Stacking

Uses the predⁿ from multiple models to build a new model. This model is used to make predⁿ of test set.



We take the training data and run it through multiple models M_1 to M_n , called base models, to generate predⁿ.

Pred1 to Predⁿ are prediction & this is the input sent to (Model), instead of Max Voters.

Instead of model takes this input & gives final predⁿ.

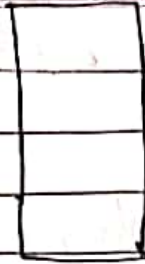
Right model is chosen acc to the problem.
i.e. whether it is classification or regression.

"Overfitting" is a measure

Reducing Overfitting

(1)

Data

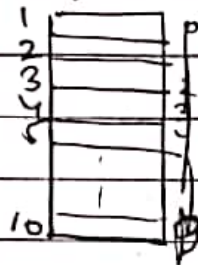


Training

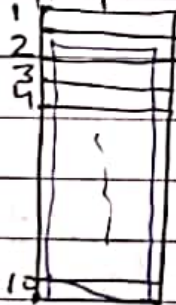


Test

(2) Train data is divided into 10 random parts.
So 1 train data has given 10 smaller Datasets.



(3) To reduce "overfitting", train 9 out of these 10 parts,
10th part is used for predⁿ.

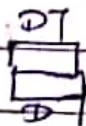
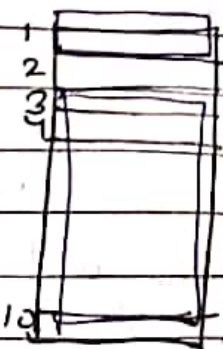


DT

- part 2 to 10 for train
- part 1 for predⁿ

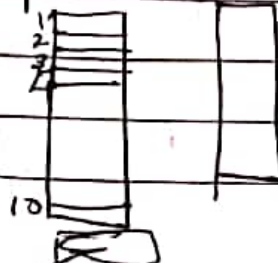
(4)

Repeat this



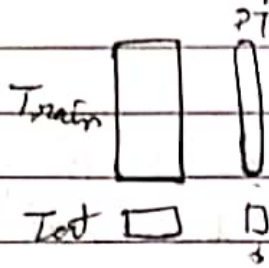
(5)

finally predⁿ for all
10 parts. Each of these
predⁿ is coming from
a model that has not
seen the same data
points DT

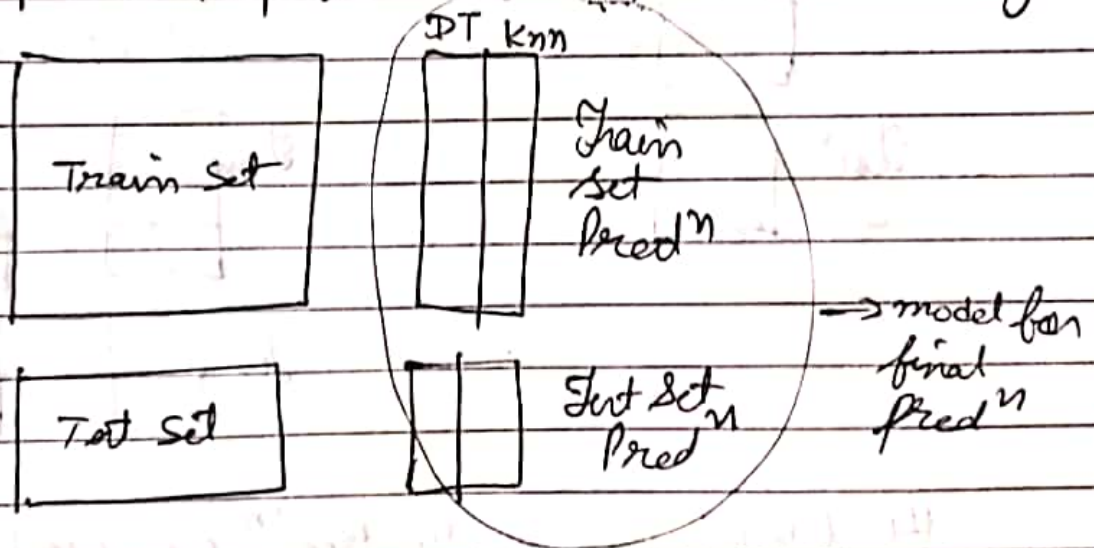


⑤ We created 10 models to get

⑥ for creating a test dataset, we can use the entire train dataset..
Again train the models on entire training datasets & make predⁿ on test



⑦ Step 1-6 repeated with another base model (say Knn)

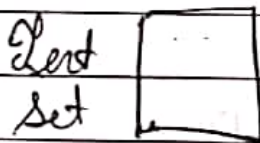
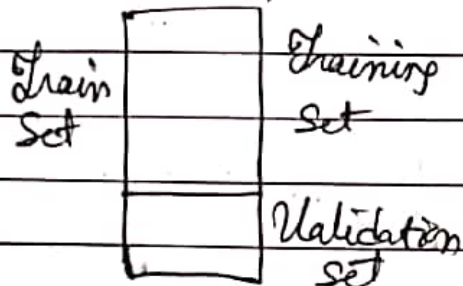


Eg process 1 to 6 is used to gen DT, Knn dataset
then the third model LR is used on
predⁿ of DT & Knn models

2.2 Blending

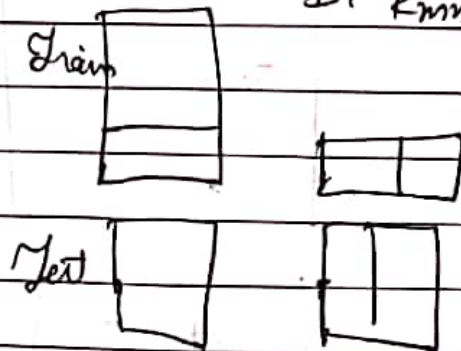
Similar to stacking except that train data is splitted into training set and validation set. ~~and the predⁿ set~~
Validation set & predⁿ are used to ~~re~~ build a model to ~~re~~ on test set

① Train set split



② models are fitted on training set.

③ Predⁿ are made on Validation & Test set

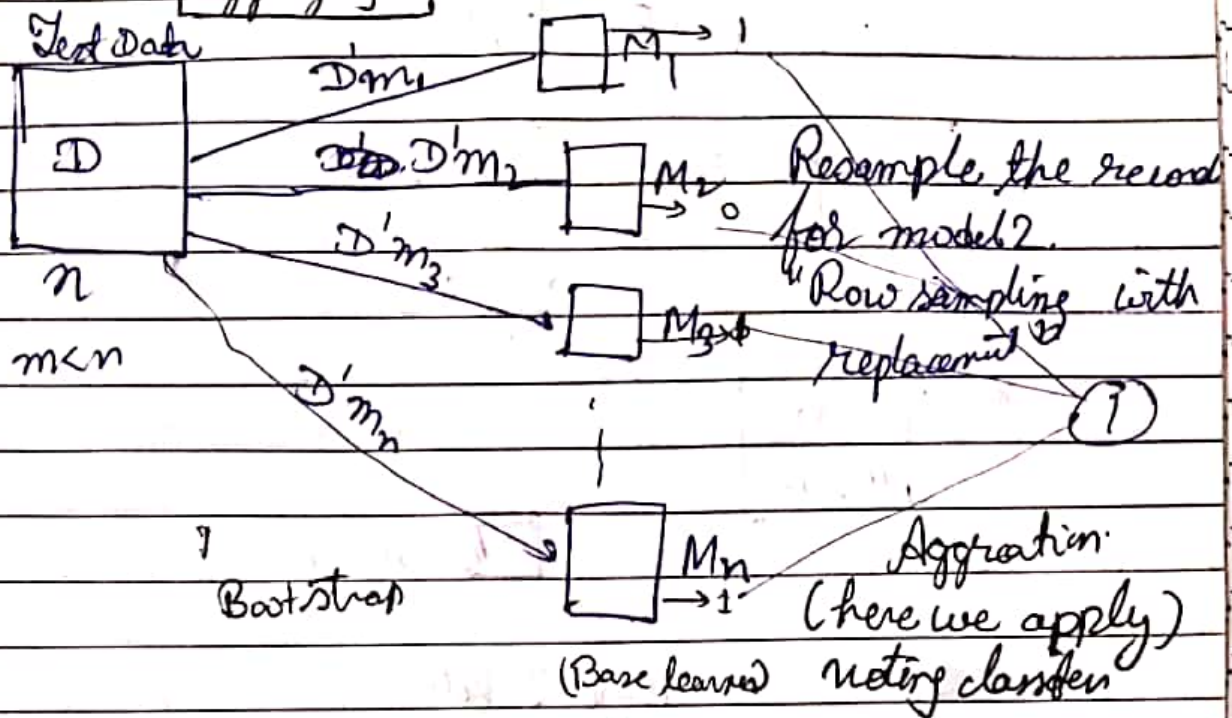


④ Validation Set & its predⁿ are used as features to build a new model.

⑤ Model is used to make predⁿ on test & meta features

2.3 (8)

Bagging Bootstrap Aggregation



$D'm_1 \neq D'm_2$ (although some of the records may be same)

Row sampling with Replacement \Rightarrow Bootstrap

R sklearn Bagging hyperparameters

base-estimator: DT by default else specify

n -estimator: no of base-estimators \leftarrow small no: error prone.
large no: computationally expensive

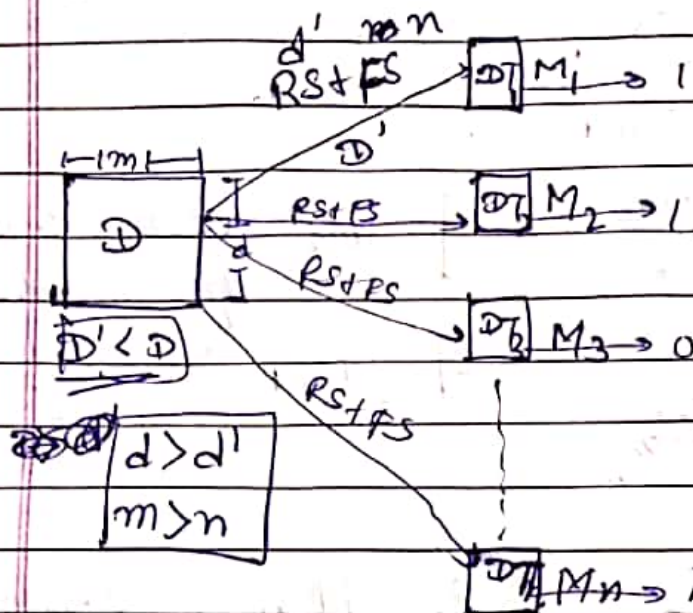
max-sample: size of subset.

max no of samples to train each base estimator

max-features: no of features from dataset.

~~n jobs~~

RANDOM FOREST



Majority Vote

DT

- DT to its complete depth \rightarrow ① Low Bias & ② High Variance
- Low Bias \Rightarrow training error v. low
- high variance \Rightarrow test error high
- Overfitting

On combine all DTs, i.e., majority vote, high variance converts to low variance.

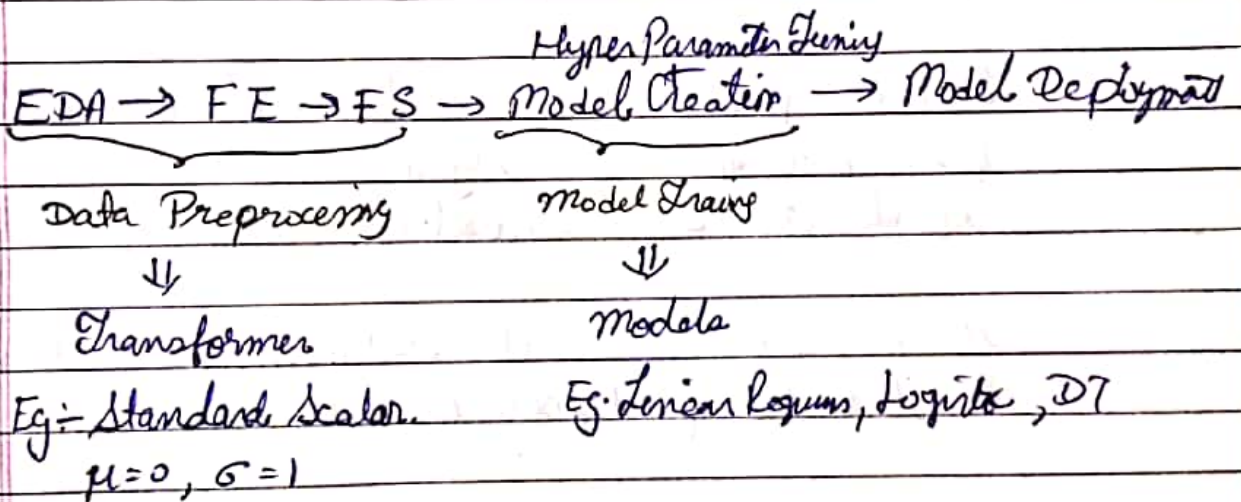
~~1000 records~~

- If we have 1000 records, change of 200 records won't affect the result much as this 200 records will be properly splitted away.

If it is a regression problem, we take the mean or median of the output.

hyperparameter \Rightarrow how many DTs for RFs.

Diff^{nce} $fit()$, $transform()$, $fit_transform()$ and $predict()$ methods in scikit learn.



MinMax Scaler, PCA, Imputer (handles NaN).

$x_1 \quad x_2 \quad x_3 \quad x_4 \quad \dots \quad x_n$

x_1

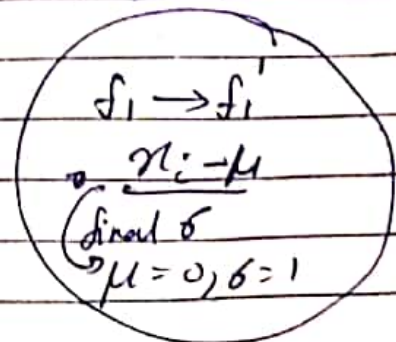
x_2

\vdots

x_i

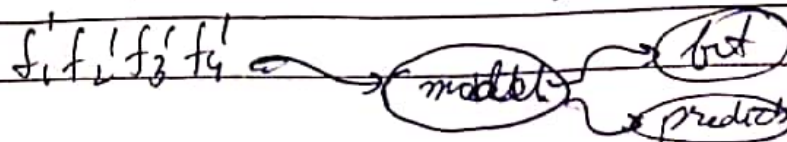
\vdots

Standard Scaler



fit \rightarrow compute μ, σ

transform \rightarrow applies formula



2

~~Bootstrap~~ B₂

Bootstrapping

Sampling With replacement

Resampling method that uses random sampling with replacement

| | | | | | | | | |
|----|----|----|---|---|----|----|---|---|
| 1 | 2 | 3 | | 1 | 10 | 3 | 8 | 3 |
| 4 | 5 | 6 | | 4 | 2 | 11 | 9 | 3 |
| 7 | 8 | 9 | ⇒ | | | | | |
| 10 | 11 | 12 | | | | | | |

initial dataset

5 12 1 5 6

Bootstrap sample of size 5

Suppose you have an initial sample with 3 observations
Using Bootstrap method, you'll create a new
sample with 3 obs as well.

In this case, second observation was chosen randomly
& will be the first observation in our new sample

| X | Y |
|---|----|
| 5 | 10 |
| 4 | 8 |
| 2 | 4 |

⇒

| X | Y |
|---|---|
| 4 | 8 |
| 2 | 4 |
| 4 | 8 |

It is very much
possible for an
already chosen sample
to be chosen again

2.3.1 Boosting Algorithm.

If a data point is predicted incorrectly by first model and then the next (probably all models), combining them will not give a very good result.

Such situation is taken care by ~~them~~ boosting.

It is a ~~seper~~ sequential process where each subsequent model attempt to correct the errors of previous model. Next model is dependent on previous model.

Steps

- ① A subset is created from the original dataset.
- ② Initially all the points are given equal weights.
- ③ Base model is created on this subset.
- ④ This model makes a predⁿ on entire dataset.

| | | | | |
|---|--|---|---|---|
| | | + | + | |
| + | | - | | + |
| | | | - | |
| + | | | | - |

- ⑤ Errors are calculated using actual & predicted values.
- ⑥ Incorrectly prediction have higher weight.
- ⑦ (Another model) will predict on entire dataset.
(is created)

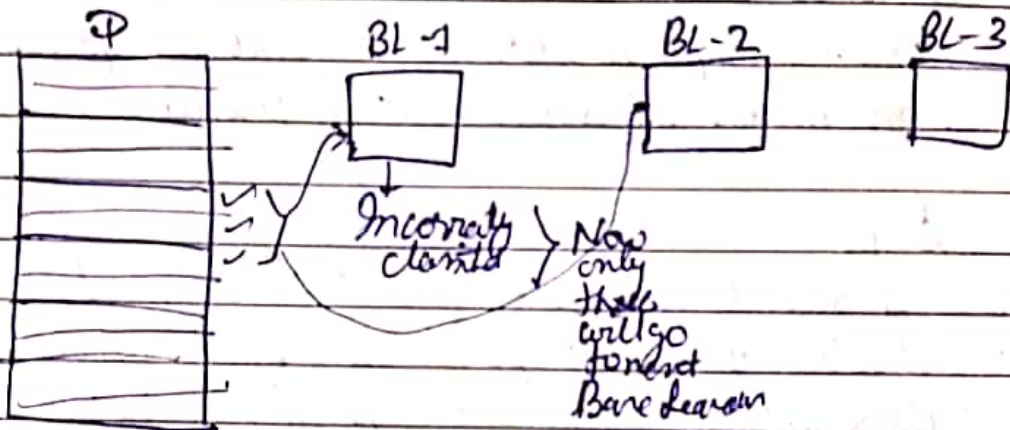
| | | | | |
|---|--|---|---|---|
| | | + | + | |
| + | | - | | + |
| | | | - | |
| + | | | | - |

- ⑧ Multiple models are created, each corrects prev model.
- ⑨ final model / strong learner is weighted mean of all models.

- Ada Boost
- GBM & XGBoost
- XGBoost
- Light GBM
- Cat GBM

Ada Boost

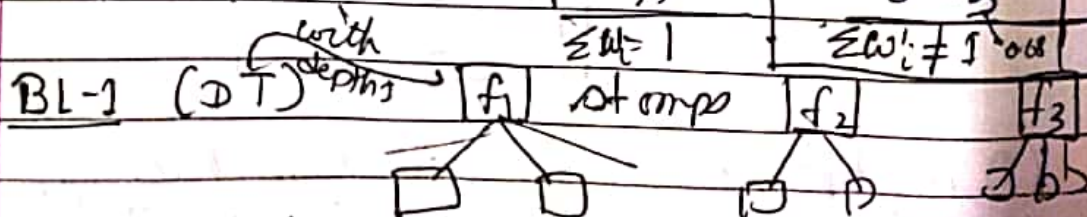
Boosting



Ada Boosting

$\omega = 1/n(\text{no of records})$

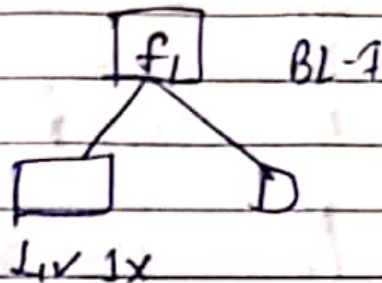
| | f_1 | f_2 | f_3 | O/p | Sample wt | Updated wt | Normalized |
|---|-------|-------|-------|-----|-----------|----------------------|-------------|
| 1 | | | | Y | $1/7$ | 0.05 | $0.05/0.68$ |
| 2 | | | | N | $1/7$ | 0.05 0.34 | $0.34/0.68$ |
| 3 | | | | Y | $1/7$ | 0.05 | |
| 4 | | | | N | $1/7$ | 0.05 | |
| 5 | | | | Y | $1/7$ | 0.05 | |
| 6 | | | | | $1/7$ | 0.05 | |
| 7 | | | | | $1/7$ | 0.05 | |



for each feature create stump

Step-1

from there we have to select BL-1 based on entropy & ginni index
(Lower entropy selected)



Step-2

• for this (1x) find total error.

$$\text{Total error} = 1/7$$

Step-3

$$\begin{aligned} \text{Performance of step} &= \frac{1}{2} \log_e \left(\frac{1 - TE}{TE} \right) \\ &= \frac{1}{2} \log_e \frac{1 - 1/7}{1/7} = \frac{1}{2} \log_e 6 \\ &= 0.896 \end{aligned}$$

Step-4

$$\begin{aligned} \text{New sample cut} &= \text{weight} \times e^{\text{Performance}_{\text{step}}} = \frac{1}{7} \times e^{0.895} \\ \text{(incorrect ones)} &= \cancel{\frac{1}{7} \times e^{0.895}} = 0.349 \end{aligned}$$

$$\begin{aligned} \text{Correctly Classified} &= \text{cut} \times e^{-\text{performance}_{\text{step}}} \\ &= \frac{1}{7} \times e^{-0.895} = 0.05 \end{aligned}$$

S-5 Sum of updated cut $\neq 1$ \therefore normalize
 $\therefore w'_i = w_i / 0.68$

S-6 Create a new dataset based on updated values will mostly select wrong records.

| f_1 | f_2 | f_3 | f_4 | O/P | Normalized | f_1 | f_2 | f_3 | O/P |
|-------|-------|-------|-------|-----|------------|-------------|-------|-------|-----|
| | | | | Y | 0.07 | 0 - 0.07 | | | |
| | | | | N | 0.51 | 0.07 - 0.58 | | | |
| | | | | Y | 0.07 | 0.58 - 0.65 | | | |
| | | | | N | 0.07 | 0.65 - 0.72 | | | |
| | | | | | " | | | | |
| | | | | | " | | | | |
| | | | | | " | | | | |

S-7 We'll create bucket

To select data, algorithm 8 iterations to select data. In first iteration, Algo selects a random value of 0.43 and check in which bucket it falls

1st
2nd

0.31