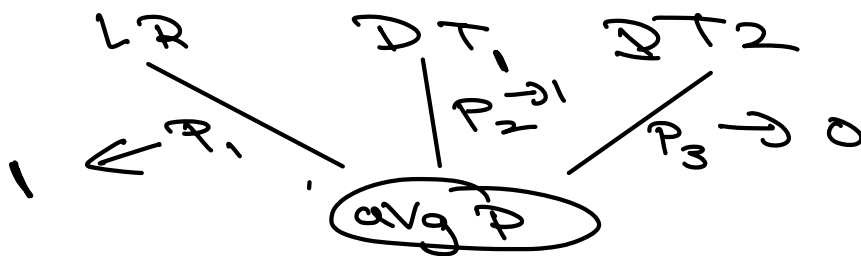


* Using DT classifier (max-depth=4)

Train-Score $\approx 85\%$

Test-Scores $\approx 78\%$

How to improve this?



Ensemble

Base-Learners + Combine/Aggregate Result
(as Unique as Possible)

Variants of Ensemble Technique

1) Bagging (Ex:- RF) *

2) Boosting* (Ex:- GBDT, XGBOOST, LightGBM)

3) Stacking

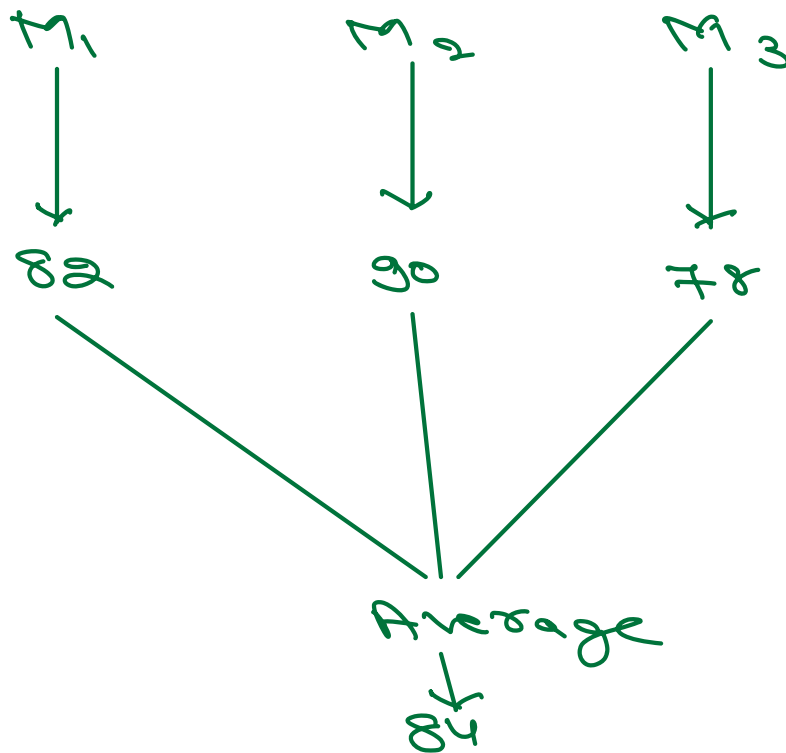
4) Cascading

} not so important from industry perspective
(Highly used ML Competitions)

Bagging

(Bootstrap Aggregation)

Bootstrap Sampling + Aggregation



Random Forest

(Bagging Model)

Each tree shall be
trained on Randomly
Sampled subset of Data

① Sample Rows (Row and Col Sampling)
② Sample Columns



RFB \Rightarrow DT's + RS + C.S + Aggregation

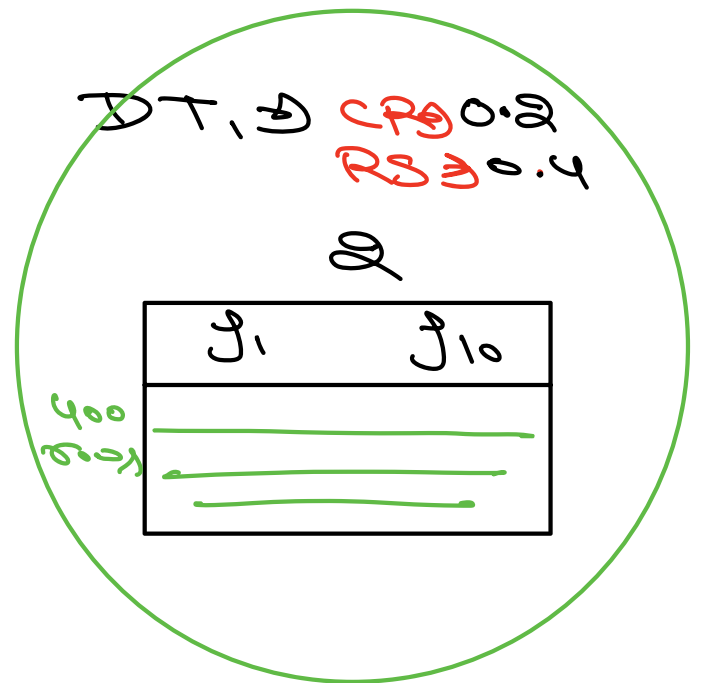
Unique and different Trees

Prediction

$Q = 10$

$n = 1000$

j_1	j_2	j_{10}

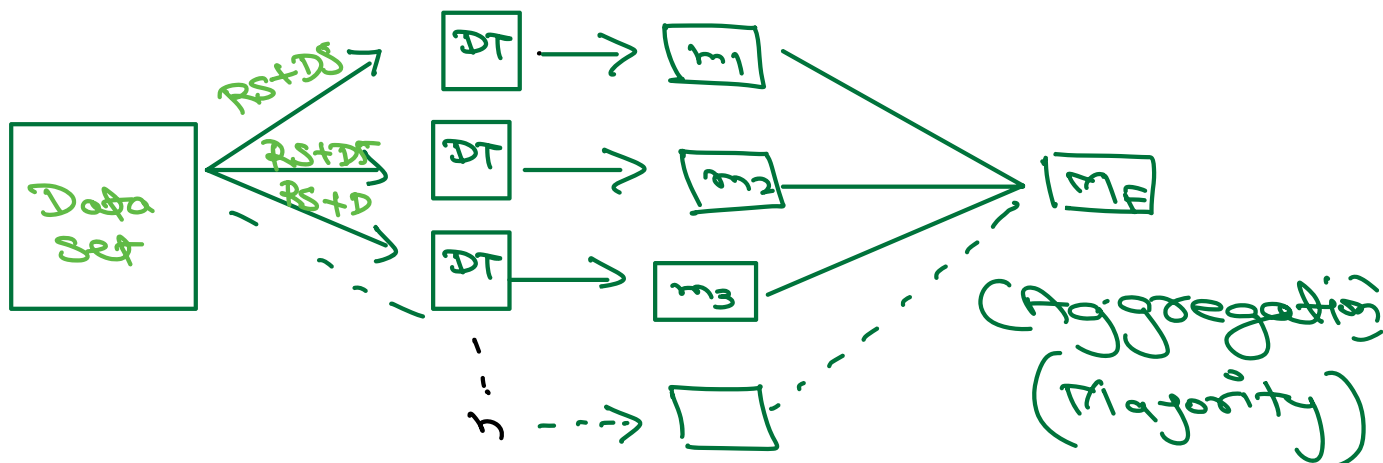


Repeat this

n -Estimators $n = 100$

- * RSR \Rightarrow rows
- * CSR \Rightarrow columns

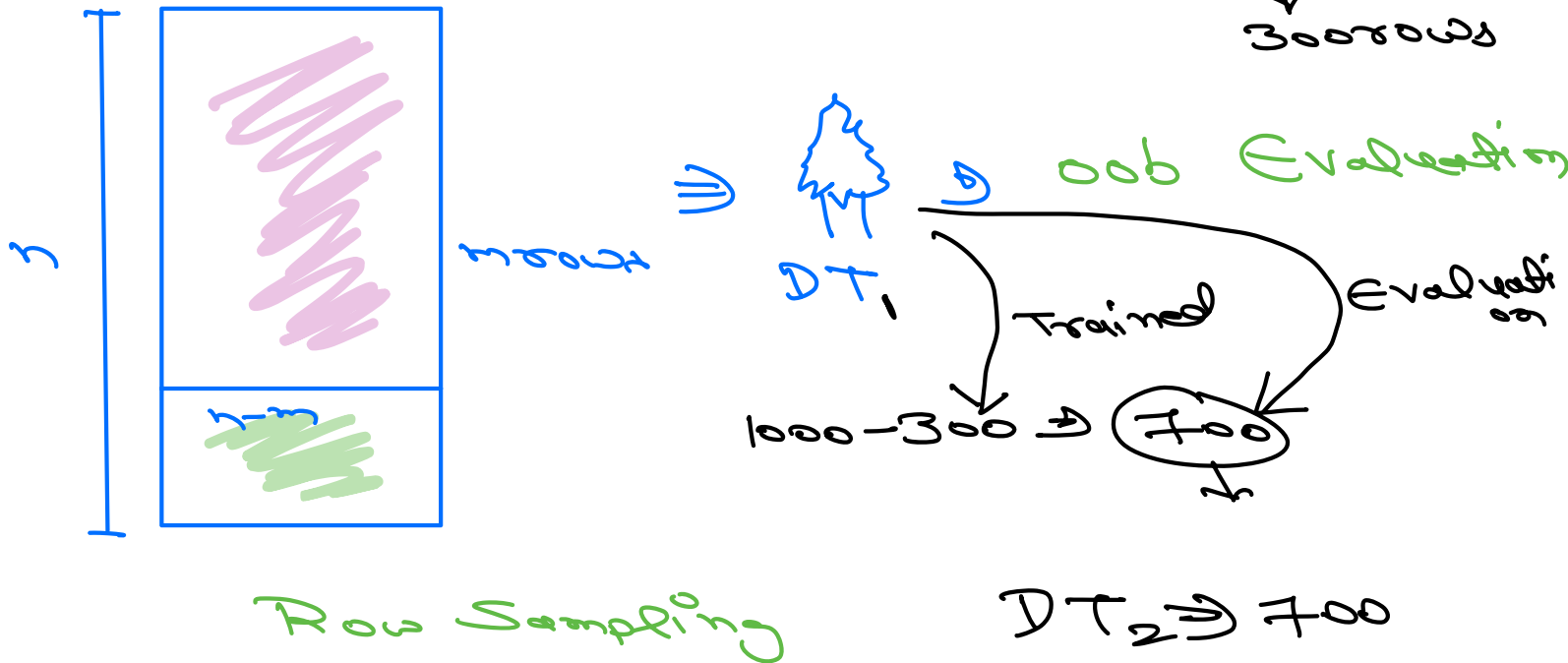
Bootstrap Sampling \Rightarrow Sampling with Replacement



- CS + RS Leads to
- 1) Set of underfitted DT's
 - 2) Different opinions

OOB

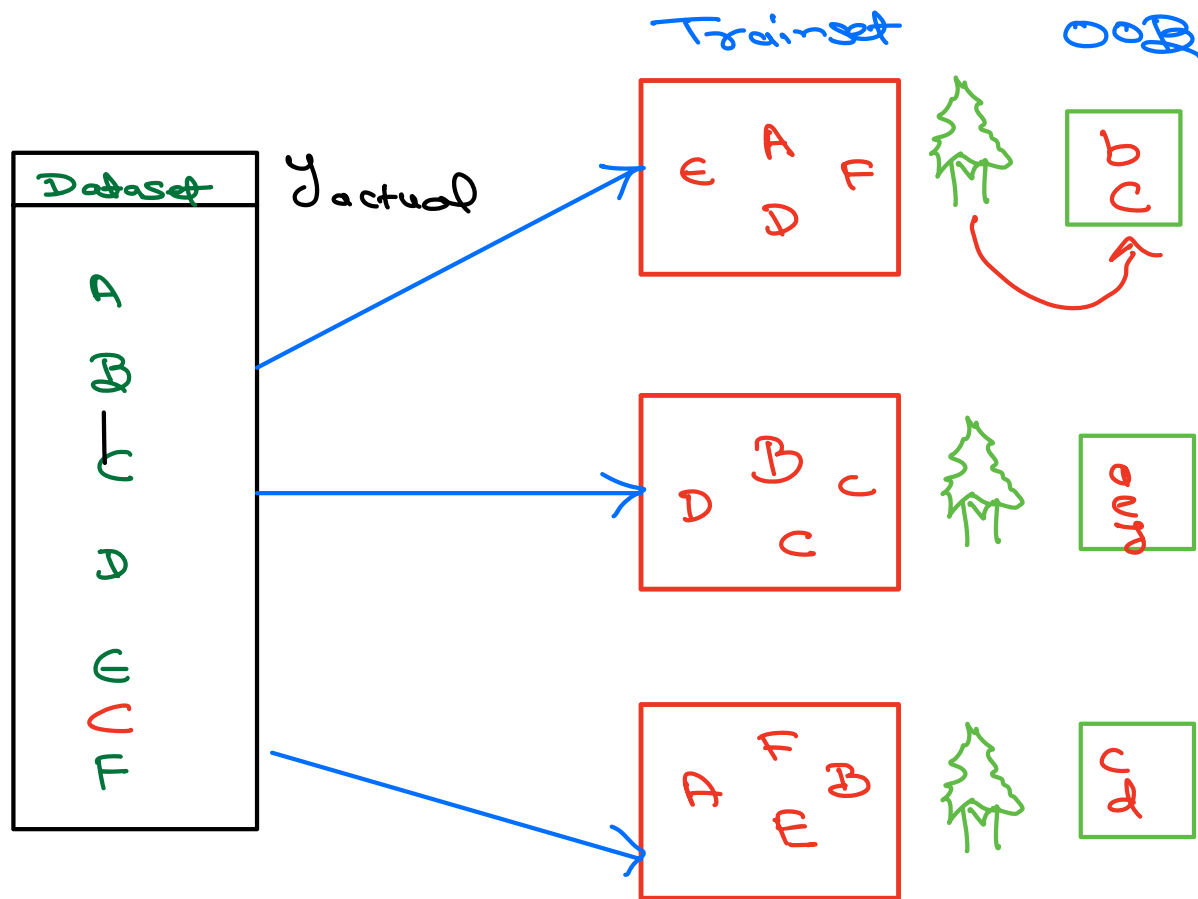
DE 1000 row
R.S.D 0.3
↓
300 rows



Left out Rows after Row Sampling from each tree can be used for Evaluating the performance of Model

OOB Score ⇒ Variation of Cross-Val Score

110



Bootstrap Sampling

OOB-score

$$P_c \ni \underbrace{Q_{t_1}(c) + Q_{t_3}(c)}_Q$$

$$y_c - P_c \ni \text{Error}_c$$

$$\sum_{i=1}^M \text{Error}_i$$

OOB-error

Pending

Implement RF with Data

Bias Variance in RA

HyperParam Tuning with Grid Search or Random Search ?

Bias Variance Tradeoff

DT \Rightarrow High Variance + Low Bias (overfit)

Ideal Model \Rightarrow Low Bias and Low Var

Base Learner (DT) \Rightarrow Random Forest
(Low Bias, Low Variance)
+
Bootstrap aggregation
Bagging \Rightarrow

Error \Rightarrow Bias + Variance + ϵ
(Simple) Models (Underfitted) (Complex) Models (Overfit)

Training Parallelization

* Since Every tree is independent, we can train using Multi-processing Jobs.

Hyperparameters

* n-estimators \Rightarrow No of trees

* Max-samples \Rightarrow R.S.R (0.1)

* Max-features \Rightarrow Column Sampling (Sort, log, Num)

* CCP-alpha \Rightarrow Cost Complexity pruning

Regression \Rightarrow Error + $\lambda ||w||$

Tree \Rightarrow Error + λ (No of Terminal Leaf Node)

Grid Search

\Rightarrow Used for finding Best Combination of Hyperparameters

① No of trees \Rightarrow [10, 100, 1000, 5000]

② Max-depth \Rightarrow [5, 10, 20,]

CCP

10, 5	10, 10	10, 20
100, 5	100, 10	100, 20
1000, 5	1000, 10	1000, 20
5000, 5	5000, 10	5000, 20

No of trees

Max depth


$4 \times 3 \Rightarrow 12$

Find Best model among 12

4 Notree

3 Max-dep $\rightarrow 60$

5 CCP-alka

C.V \rightarrow Cross-validation \rightarrow 2 \rightarrow 120


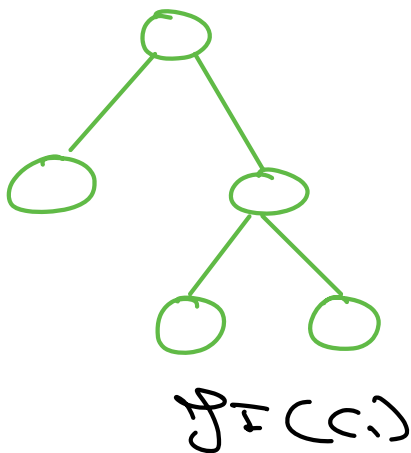
* Randomized Search \rightarrow

Randomly Select a subset
of Combination

Feature Importance

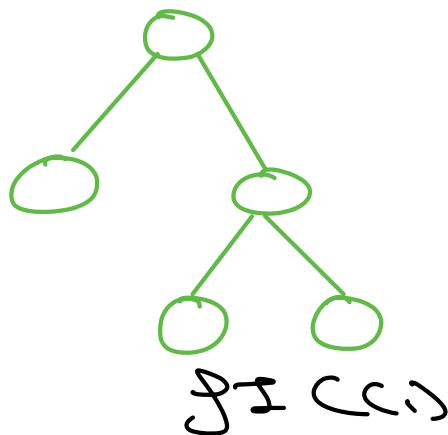
$D \in C_1, C_2, C_3$

Tree 1
 C_1, C_2

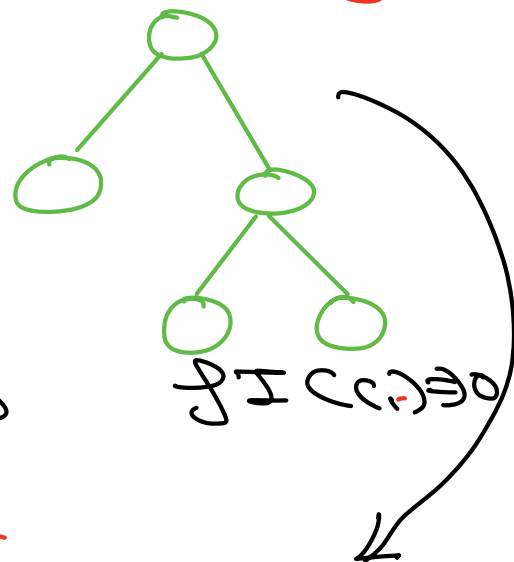


Random Forest

Tree 2
 C_1, C_3



Tree 3
 C_2, C_3



66% CS.R

$$FI_{Total} \Rightarrow \frac{FI(C_1) + FI(C_2) + FI(C_3)}{3} \Rightarrow 0$$