

Decision Tree : not Good

Random Forests

Linear Reg

Logistic Reg

Naïve Bayesian \Rightarrow "basic" method (^{= not Good model}Weak Model)

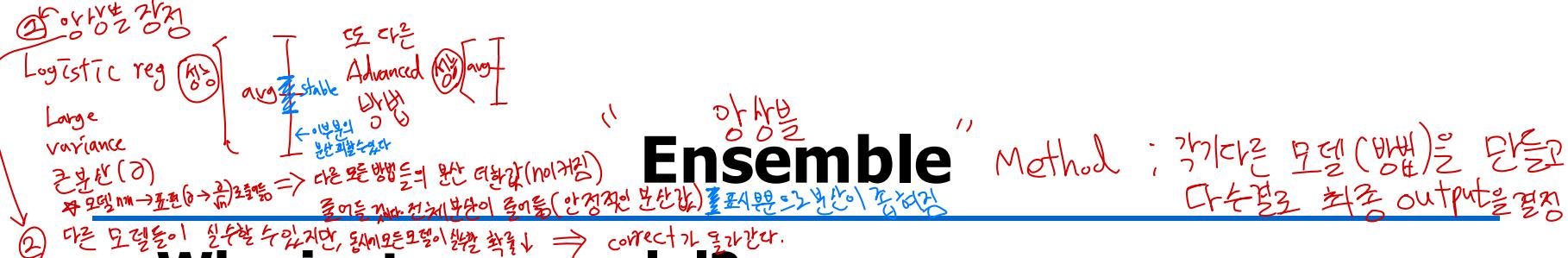
Decisiontree

\rightarrow 성능도 일반적 (Not Good)

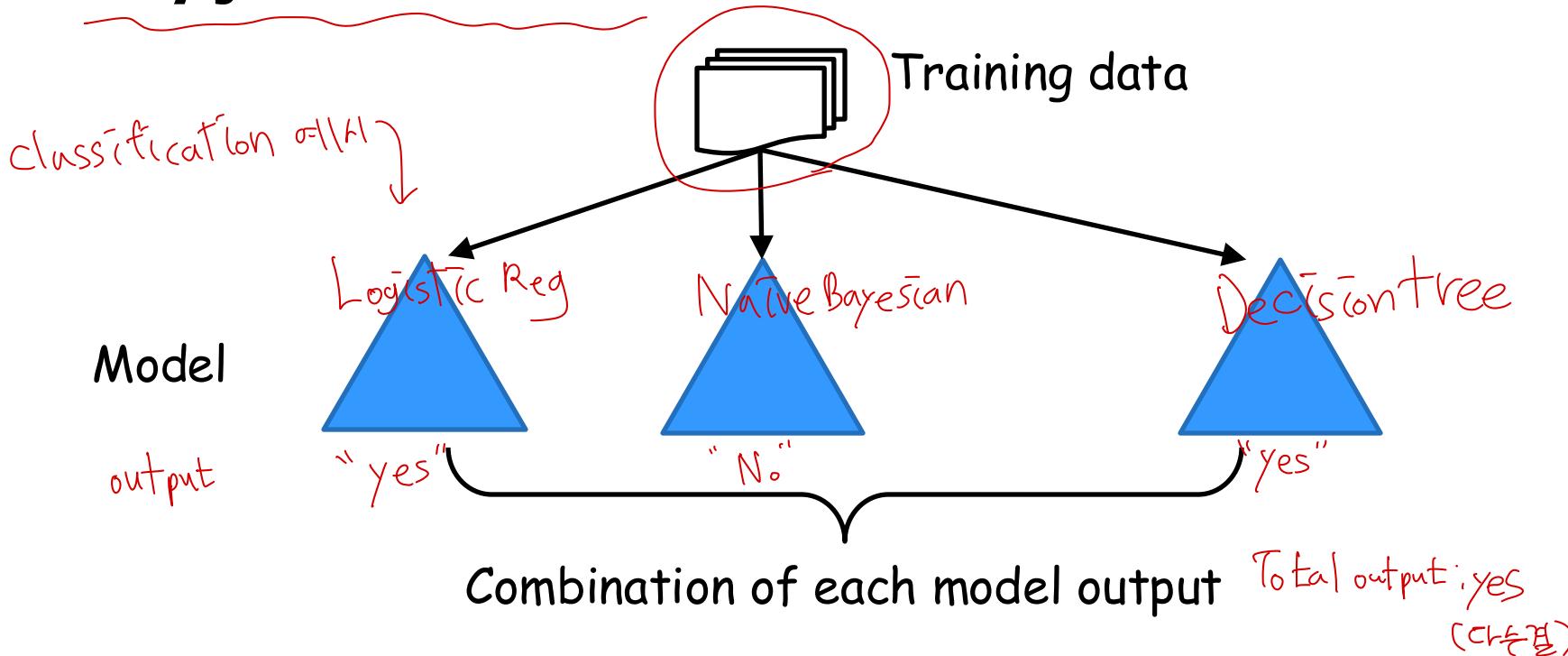
a set of Decision trees
+
Randomness

기타 \downarrow Random Forest
SVM
NN

Strong model



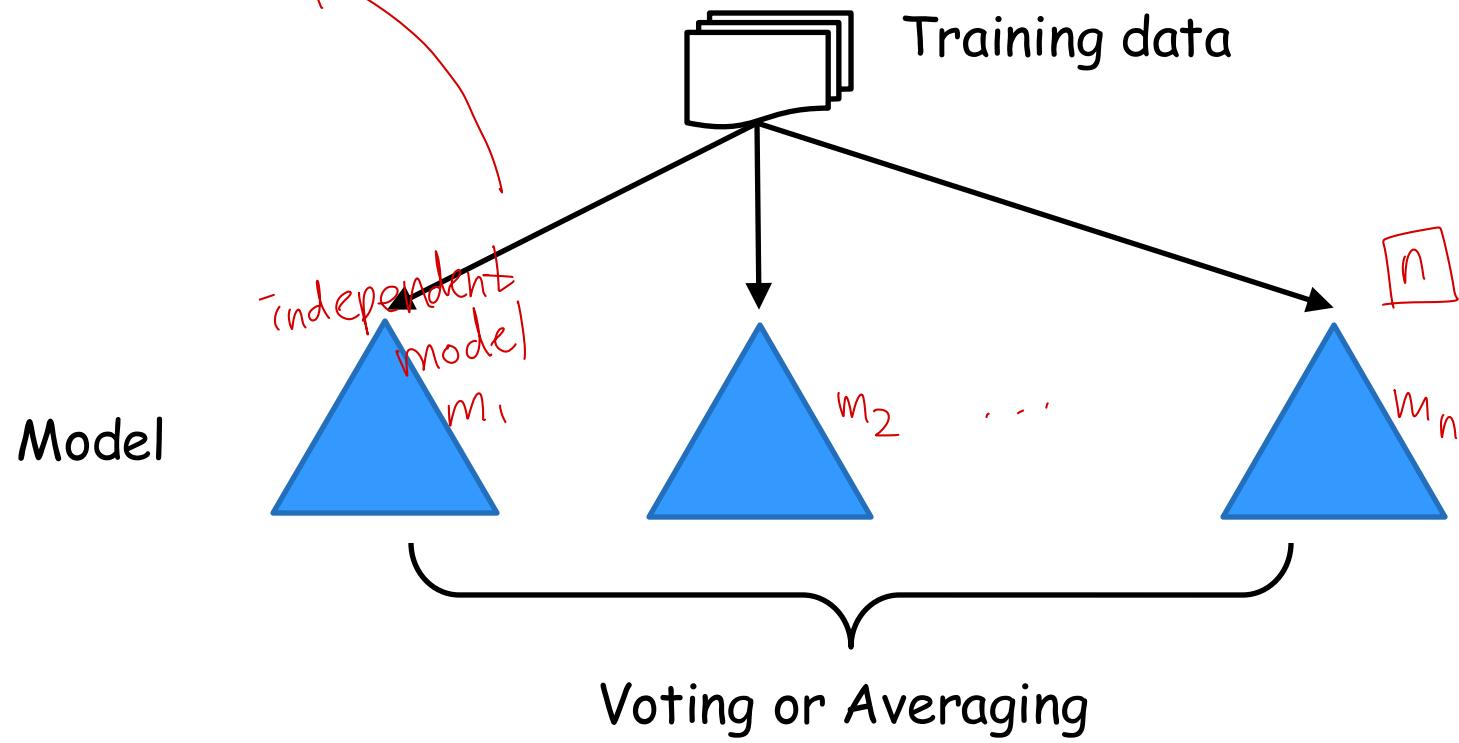
■ Why just one model?



- Accuracy of each model should be high enough
- Correlation of each model should be low enough

Ensemble

■ Bagging



- All models are equivalent

$$\begin{aligned}
 & m_1 \rightarrow \epsilon_1 \\
 & m_1 + m_2 \rightarrow \epsilon_2 \quad \therefore \epsilon_2 < \epsilon_1 \\
 & m_1 + m_2 + m_3 \rightarrow \epsilon_3
 \end{aligned}$$

...
dependent

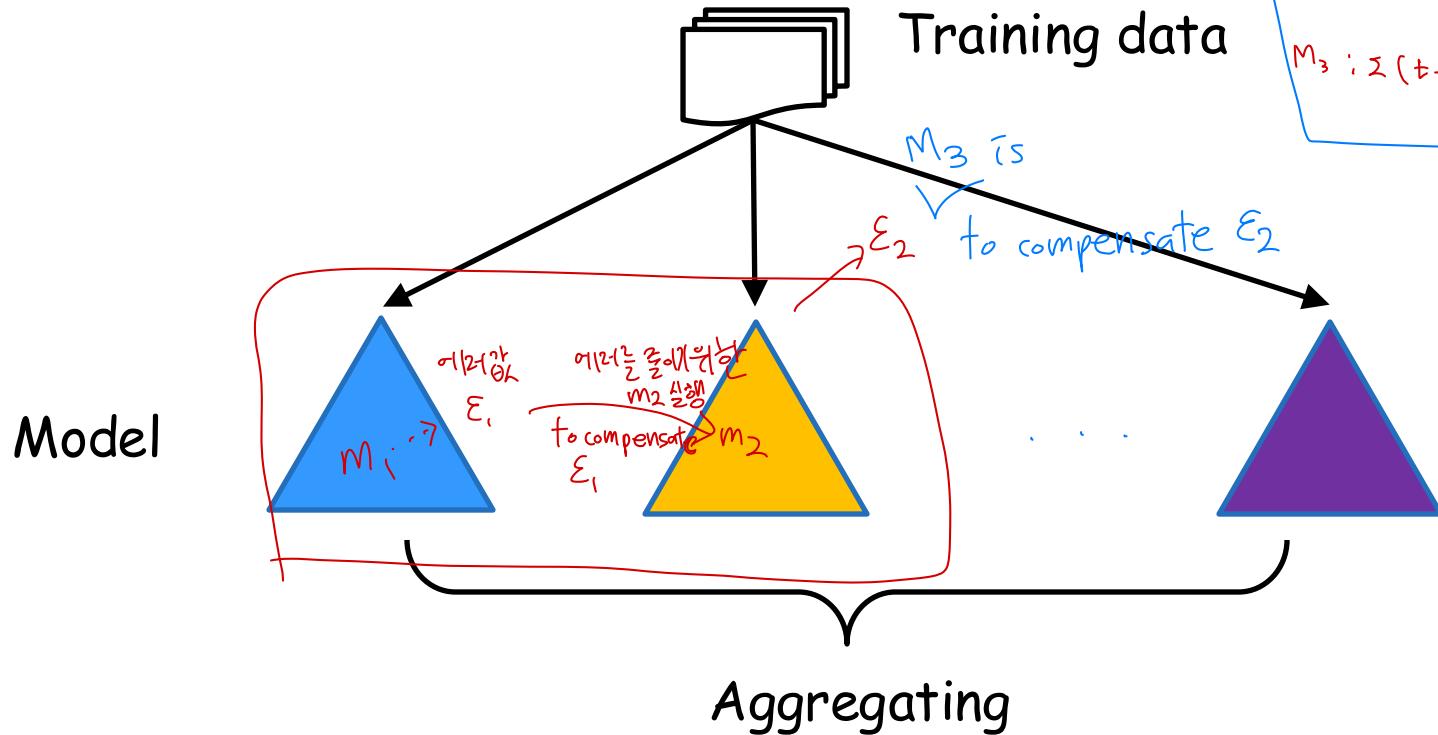
Adaboost : Adaptive boosting

옛날에 좋은 방법이 있다.

Ensemble

(boosting 방법)

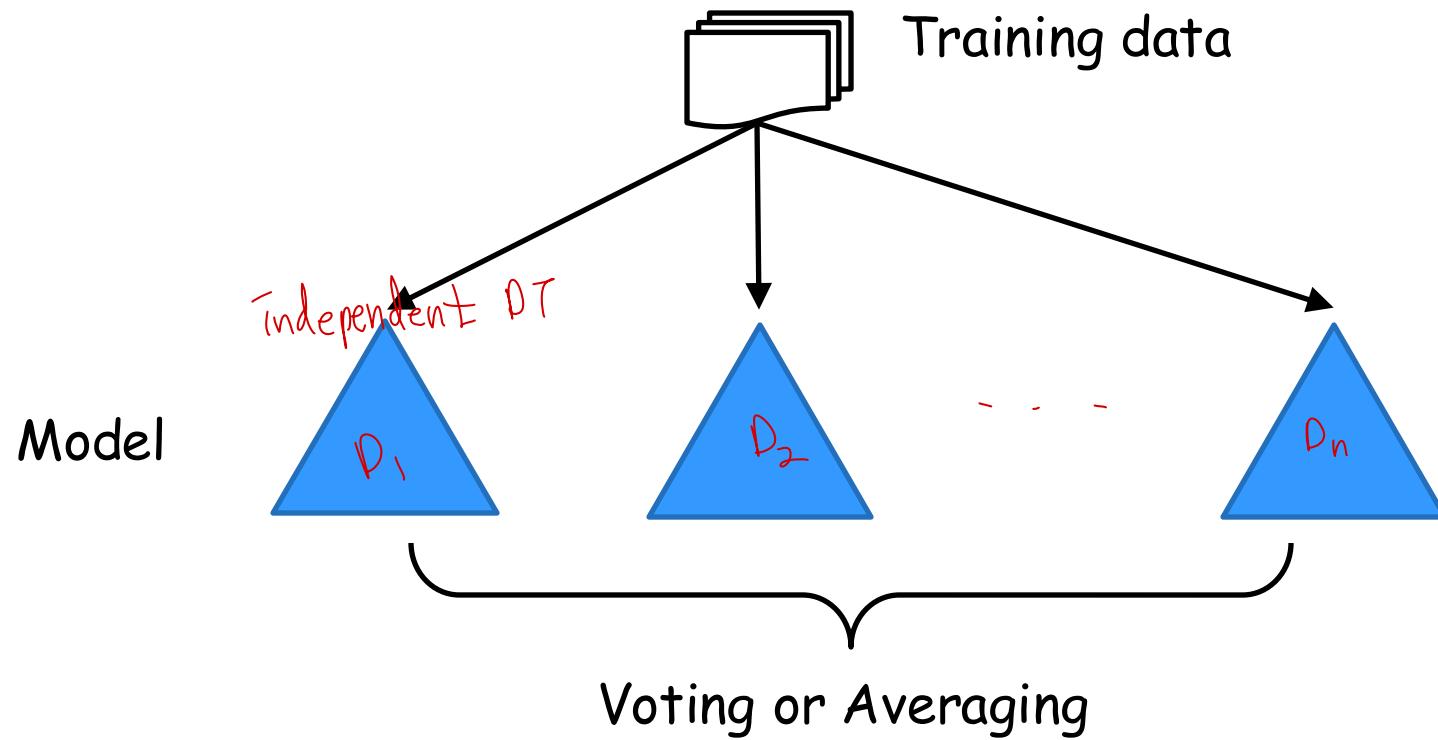
- Boosting** (Incremental Way)



- Each model compensates the error of the previous models

Bagging of Decision Trees

■ Bagging

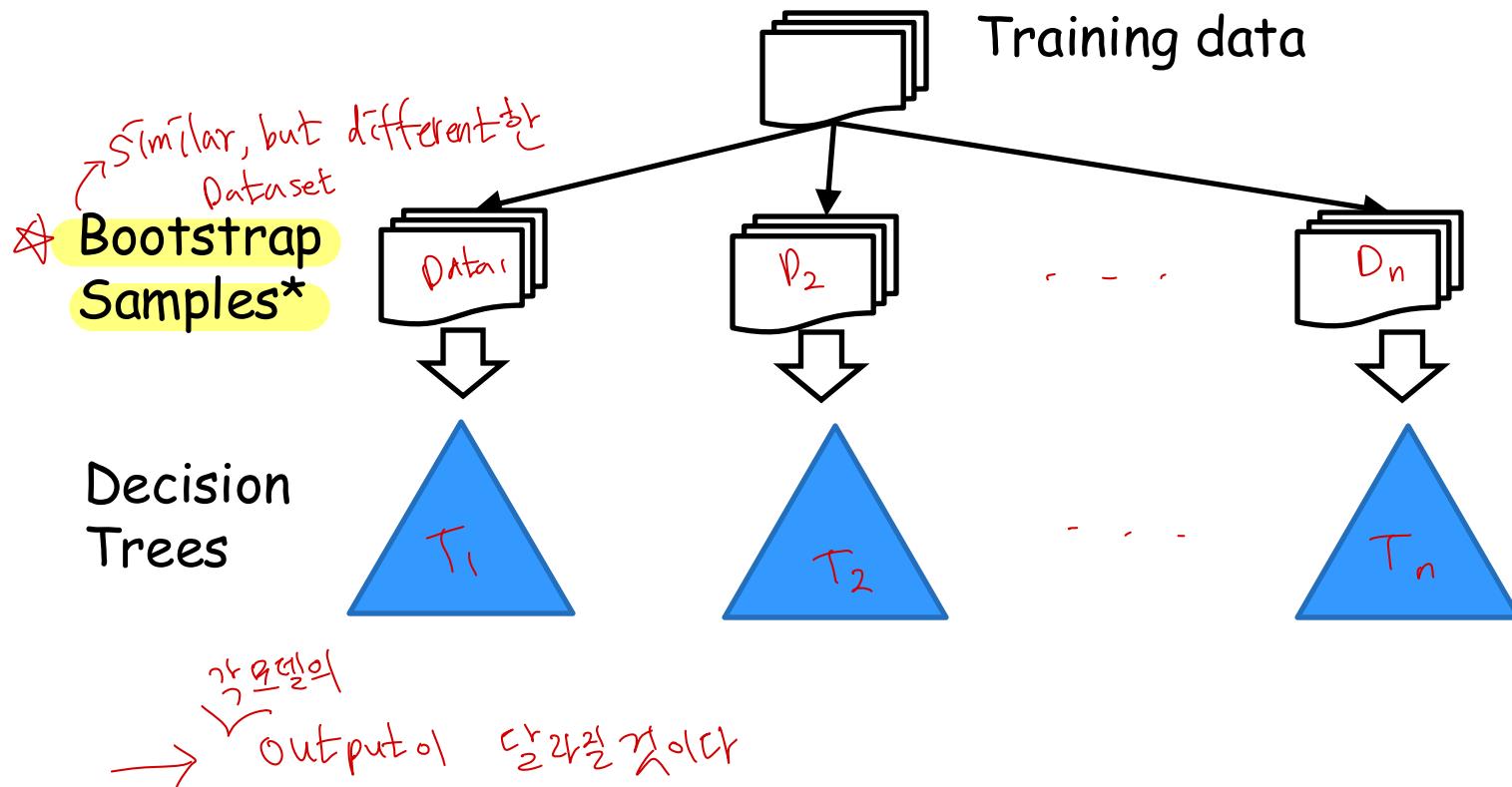


bagging의 문제점

- All models must be the same!! \Rightarrow 같은 답을 내보내면 meaning less 하다.
: 티타이타늄, algorithm이 같기 때문에

Bagging of Decision Trees

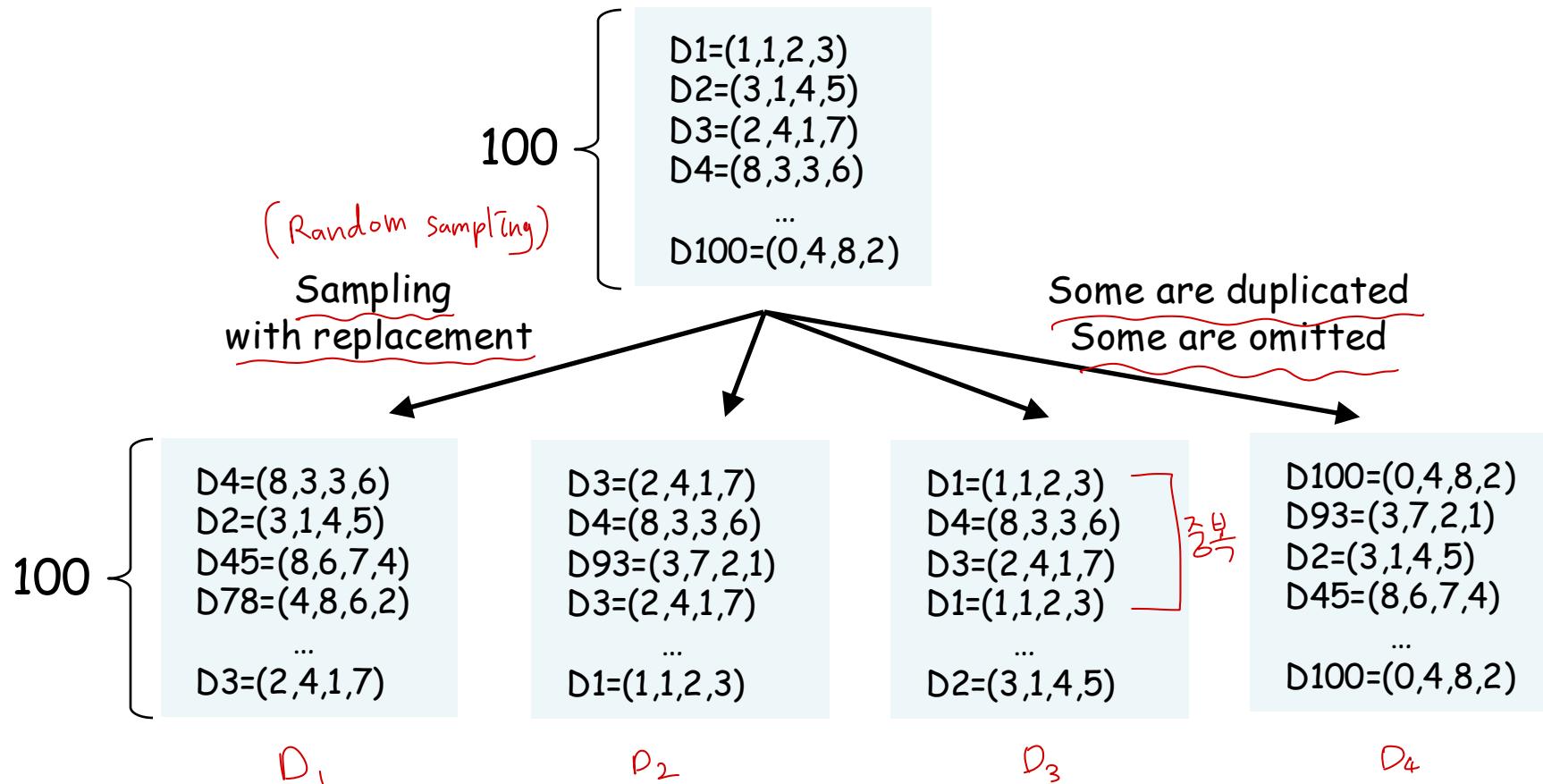
- Why just one tree?



*samples equal in size to the original dataset, but selected with replacement

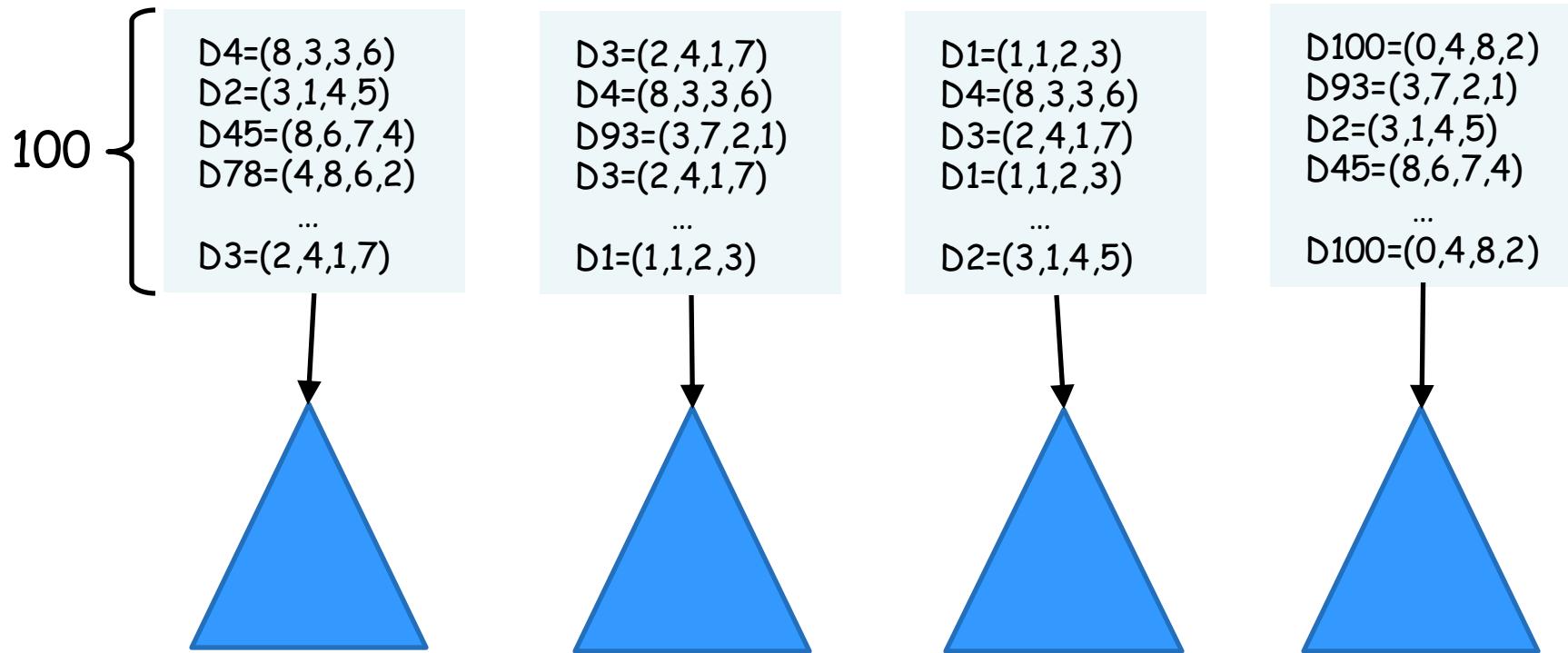
Bagging of Decision Trees

■ Bootstrap Samples



Bagging of Decision Trees

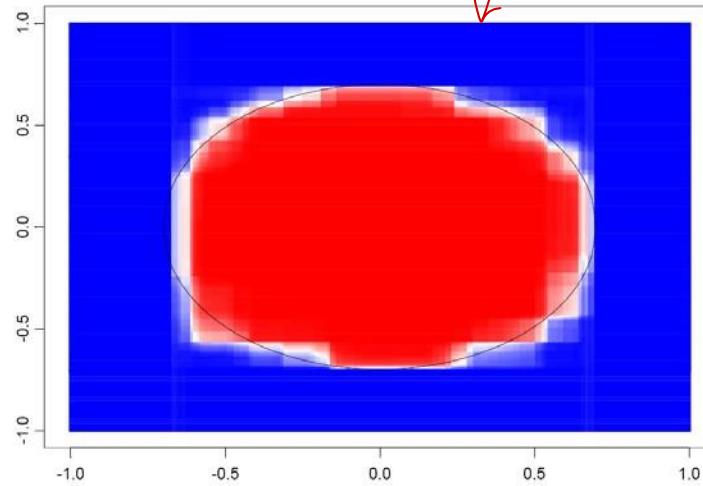
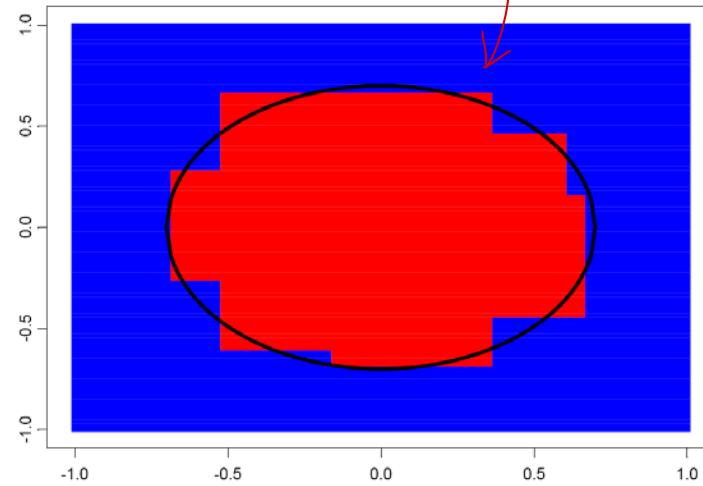
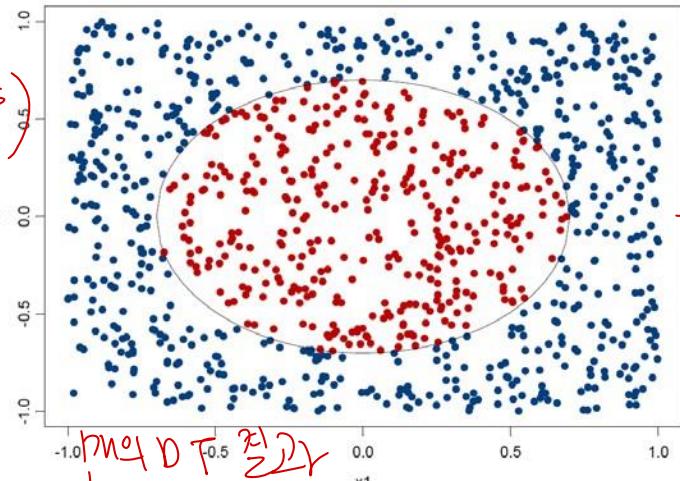
■ Bootstrap Samples



Bagging of Decision Trees

XGBoost
Tree based
ensemble learner.

- Deep Learning
 - Lots of Data samples
 - Time
- Tree-based ensemble
 - Good performance
 - less data
 - less time & computing



Bagging of Decision Trees

- **Bagging = Bootstrap Aggregation** → 배깅

- Create multiple models with bootstrap
 - Create new training sets with bootstrap samples
 - Create models with each samples
 - Averaging or major voting of the models

- **Advantage of Bagging**

- Decreasing test error by lowering prediction variance while leaving bias unchanged

= with weak model 와 같은 품질

Random Forests

Bootstrap으로 충분하지 않다
= Bagging of DTs
+ Randomness

■ Definition

- An ensemble method consisting of a bagging of un-pruned decision tree learners
- with a randomized selection of features at each split.

■ Difference with Bagging

- Bagging: selecting the best feature in the full set of feature at each split
- Random Forests: selecting the best feature in a randomly selection subset of feature at each split
(Additional Randomness)

같은 데이터
 ↓
 같은Alg ← 랜덤
 ↓
 이전으로 다른DT를 만들게 됩니다.

if (random)
 else 이런 느낌

Random Forests

= "Algorithm of DT" + Randomness

Bagging of DT + Additional Randomness

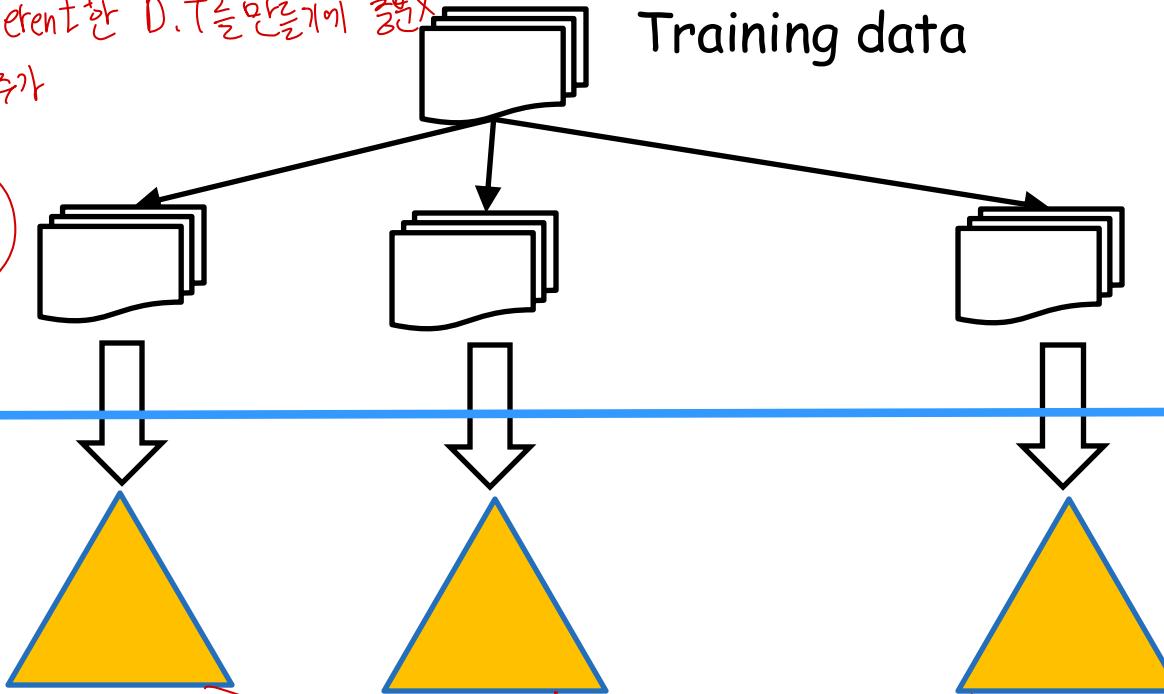
Similar, but Different 한 DT를 만들기에 충분히
 → randomness 추가

Bootstrap Samples

Creating with additional randomness

Decision Trees

Training data



상관관계

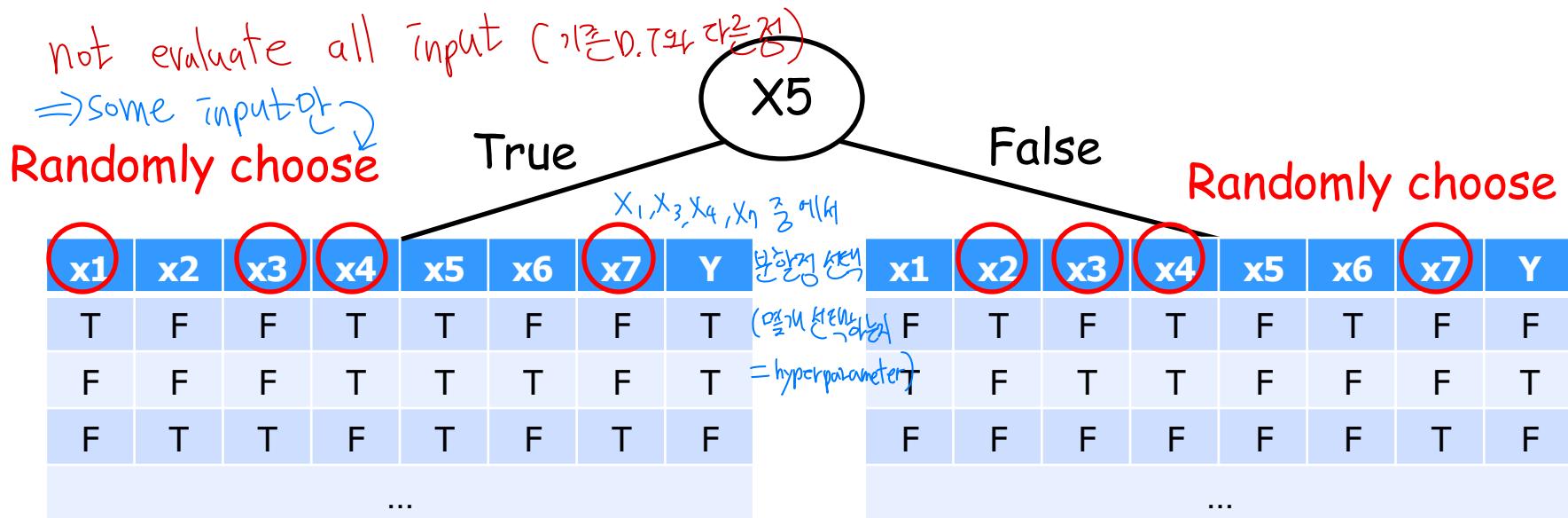
Corelation between trees are reduced by randomness

"Similar but different enough"

Random Forests

▪ Additional Randomness

- Further reduces variance even on smaller sample set sizes, improving accuracy
- Subset of features much faster to search than all features



Random Forests

■ Steps

Let N_{trees} be the number of trees to build

Let m_{try} be the number of features to be selected at each split
=inputs

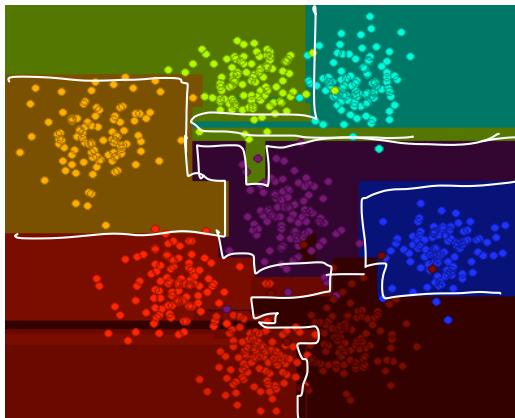
For each of N_{trees} iterations

1. Select a new bootstrap sample from training set
2. Grow an un-pruned tree on this bootstrap.
3. At each internal node, randomly select m_{try} features and determine the best split using only these features.
4. Do not perform cost complexity pruning. Save tree as is.

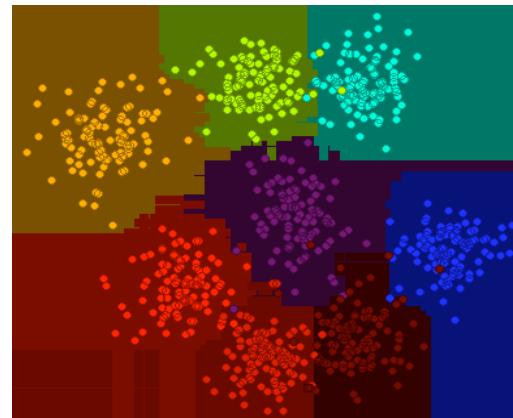
Output overall prediction as the average response (regression) or majority vote (classification) from all individually trained trees

Random Forests

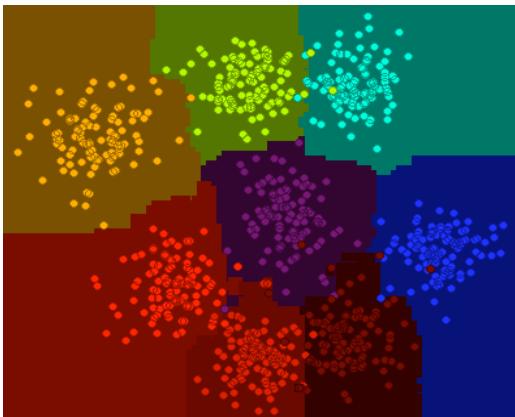
- Illustration



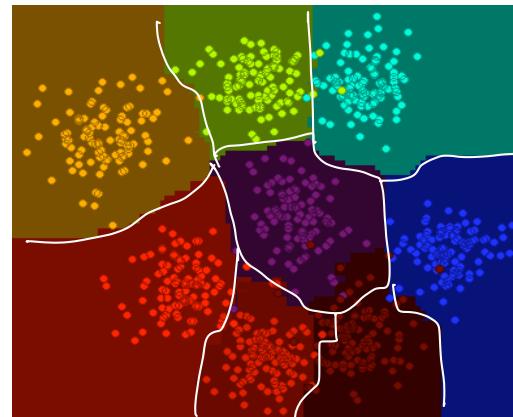
1 Tree



10 Trees



100 Trees



500 Trees

Comparison

- Intrinsically multiclass



- Handles Apple and Orange features

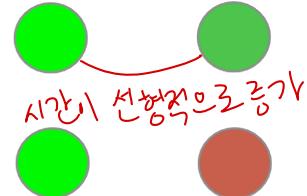
"Categorical"



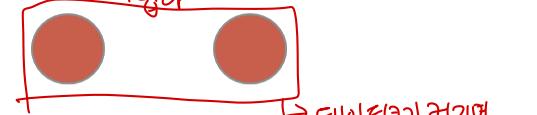
이미지를 각
개별로 정의해
여기다가

- Scalability (large learning set)

데이터 크기가 커져도 시간이 확 늘어나지 않는 정도

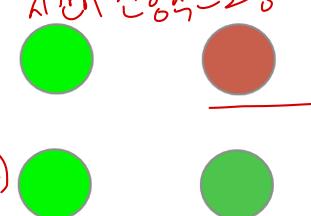


시간이 선형적으로 증가



데이터크기 커지면
시간이 지수급으로
늘어남

- Prediction accuracy



Not Good

- Parameter tuning (= sensitivity to Hyperparameter)

얼마나 많은 수의 파라미터를 설정해야 하는가?



"k"

17

Advantages

- **Overall**

- For many data sets, it produces a highly accurate classifier.
(Almost same as SVMs and NNs)
- It is faster to train and has fewer parameters than SVMs and NNs (good for large databases)
- It is interpretable

- **Training**

- It generates an internal unbiased estimate of the generalization error (Cross validation is unnecessary)
- Resistance to over training

New concept

Out-Of-Bag

R.F 일반적방법 → Train | Validation | Test

Outofbag → Train
Outofbag → Test

= Boot strap | Outofbag Test
≠ Validation

Bootstrap Sampling

$D_1 = (1, 1, 2, 3)$
 $D_2 = (3, 1, 4, 5)$
 $D_3 = (2, 4, 1, 7)$
 $D_4 = (8, 3, 3, 6)$

...
 $D_{100} = (0, 4, 8, 2)$

chosen
100
2
3
 $D_4 = (8, 3, 3, 6)$
 $D_2 = (3, 1, 4, 5)$
 $D_45 = (8, 6, 7, 4)$
 $D_{78} = (4, 8, 6, 2)$
...
 $D_3 = (2, 4, 1, 7)$

$D_3 = (2, 4, 1, 7)$
 $D_4 = (8, 3, 3, 6)$
 $D_{93} = (3, 7, 2, 1)$
 $D_3 = (2, 4, 1, 7)$
...
 $D_1 = (1, 1, 2, 3)$

not chosen $(\frac{1}{3}) 33\%$
out of bag

또 다른 33%

how many samples
may not be chosen? ⇒ 다음장

What is the prob. that
 D_1 is not include in Bag 1?

Validation sample은 600개의 D.T의 일부로 구성되는 전체 R.F에서 뽑은 샘플

• OOD accuracy → $D_1 \sim D_{100}$ 까지 D.T를
각각의 D_1, D_2, \dots, D_{100} 에 존재하는 Sub-R.F에
choose D_1 as Validation Dataset (avg)

accuracy를 측정하여 평균

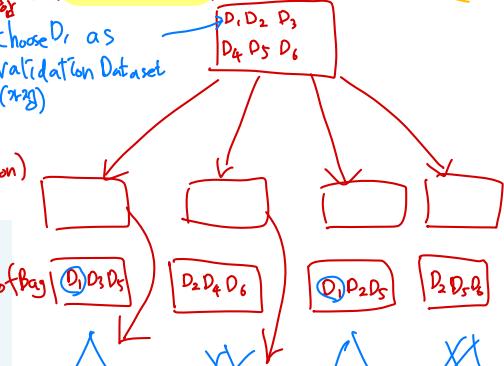
(\approx Good Estimation of Validation)

= OOD evaluation
(Indirect evaluation)

장점
: Training 데이터로
Build, evaluate
동시에 가능하
... Data utilization
...!

$D_2 = (3, 1, 4, 5)$

33%



out of D_1 이 있는 이 2개의 D.T만 평균
만약 D.T가 600개 있다면 D_1 이 존재하는 D.T는 약 200개($\frac{1}{3}$)

Out-Of-Bag

- Bootstrap sample set from learning set D
 - Remaining samples called out-of-bag samples (OOB)
 - About 33% of original data are not selected in bootstrap

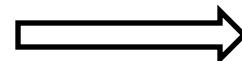
$$\lim_{n \rightarrow \infty} \left(1 - \frac{1}{n}\right)^n = e^{-1} = \frac{1}{2.718\ldots} \approx 33\%$$

선택되지 않은 확률
= not chosen의 확률 (데이터 1개)

Data }
n
= Data Amount

D1=(1,1,2,3)
D2=(3,1,4,5)
D3=(2,4,1,7)
D4=(8,3,3,6)
...
D100=(0,4,8,2)

Sampling
with replacement



Some are duplicated
Some are omitted

D4=(8,3,3,6)
D2=(3,1,4,5)
D45=(8,6,7,4)
D78=(4,8,6,2)
...
D3=(2,4,1,7)

Bootstrap
sample