

# PRESENTATION

**4팀**

김동영  
강수경  
김재희  
유경민  
최윤희

# INDEX

---

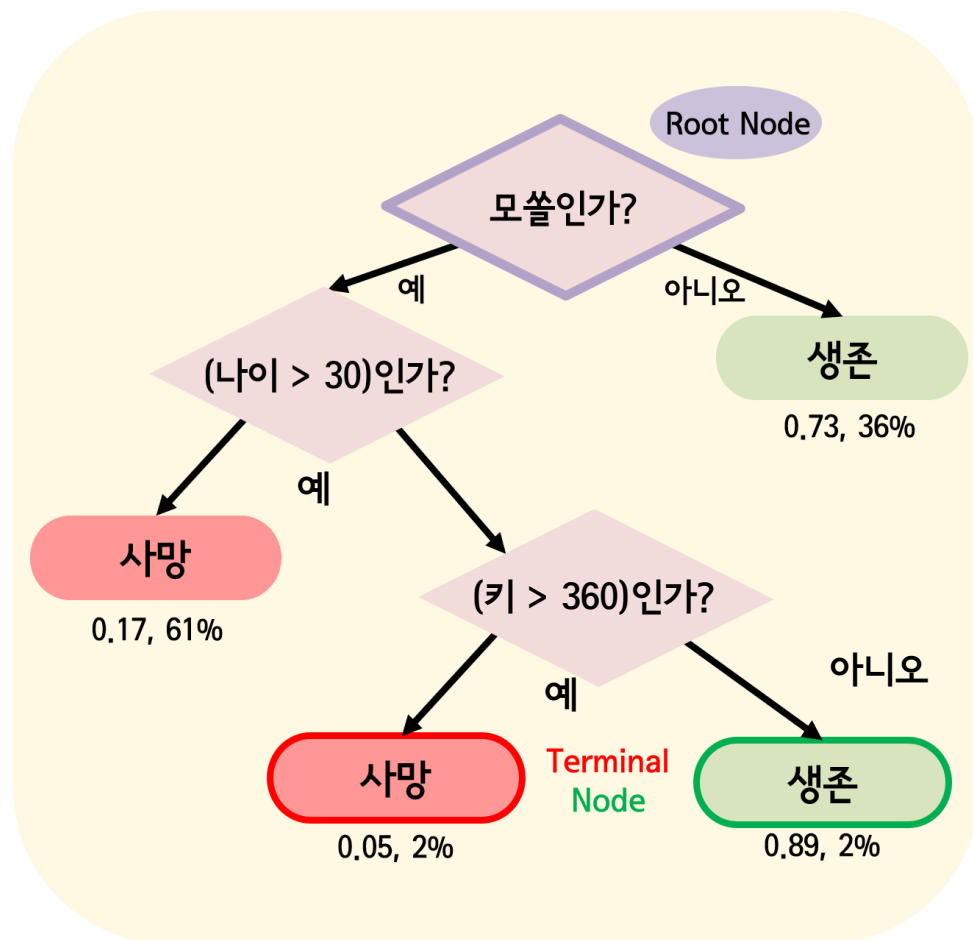
1. CART

2. Bootstrap

3. Bagging

4. Random Forest

## CART - 분류모델



각 node는 서로 배타적이다

Terminal node 들의 합  
=  
Root node의 데이터 수

중복되어 분할 X

“모든 데이터는 분할에 사용”

## 불순도 / 불확실성

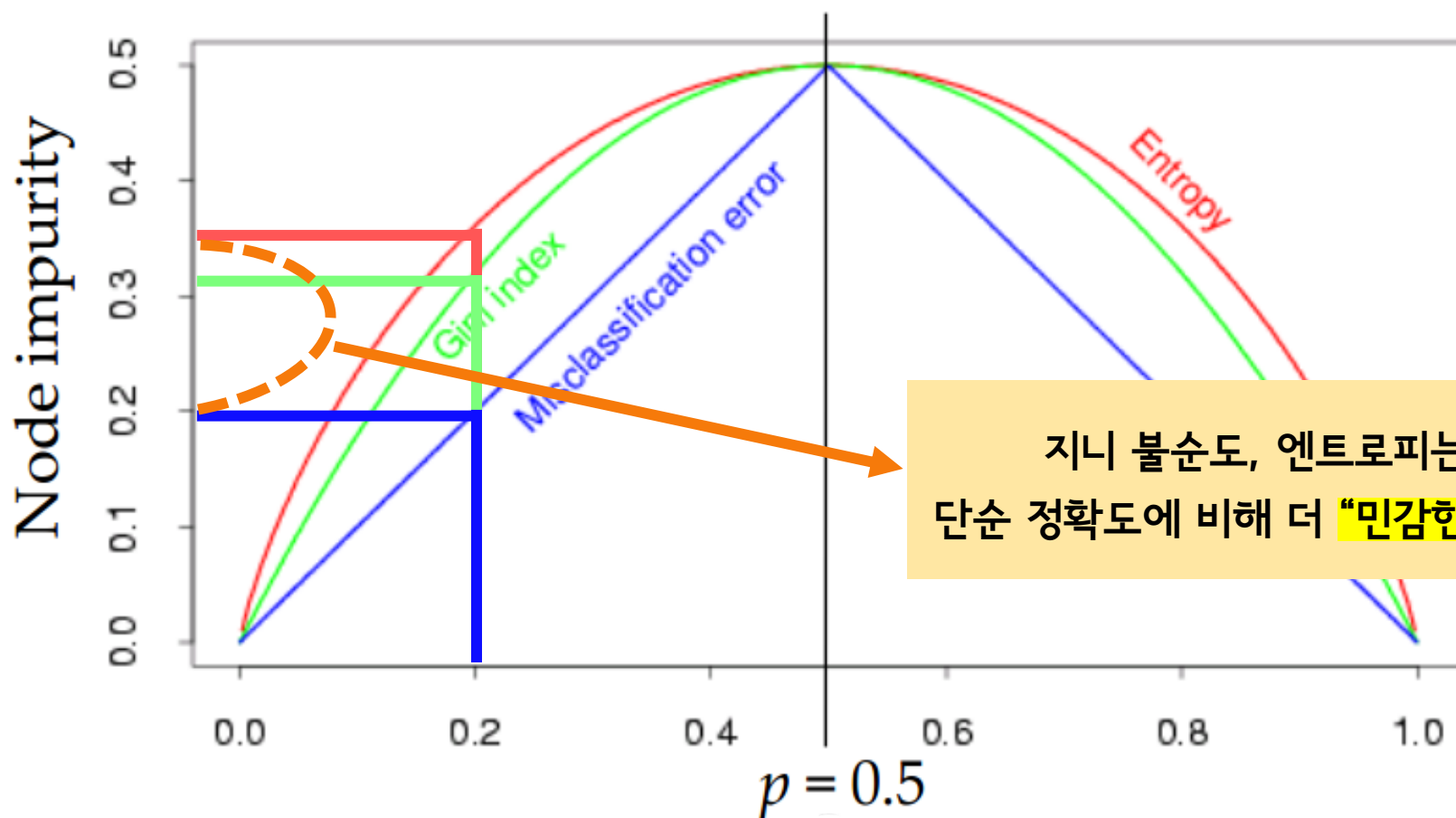


이 모델은 각 **노드의 순도가 증가하는 방향**  
(**불순도/불확실성이 감소하는 방향**) 으로 학습을 진행한다.



단순 정확성 (accuracy) 보다는  
“**지니 불순도**” 사용!

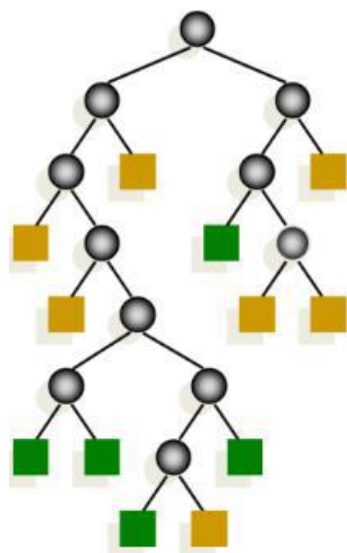
## 불순도 / 불확실성



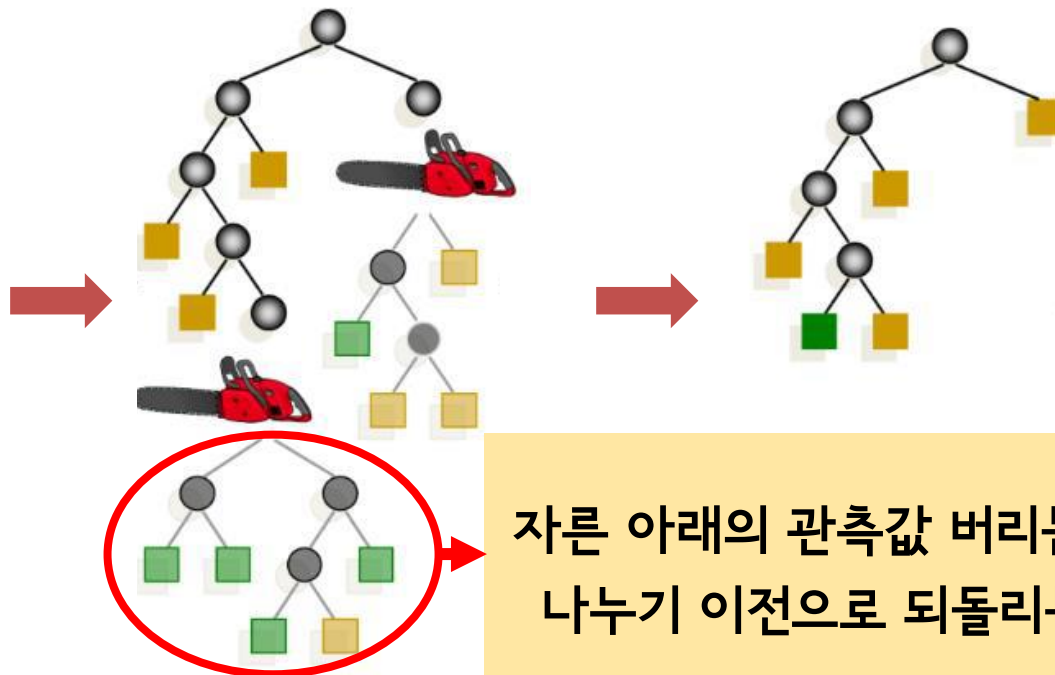
\*Entropy 부록참조\*

## 가지치기 (pruning)

트리를 **최대로 나누기!**

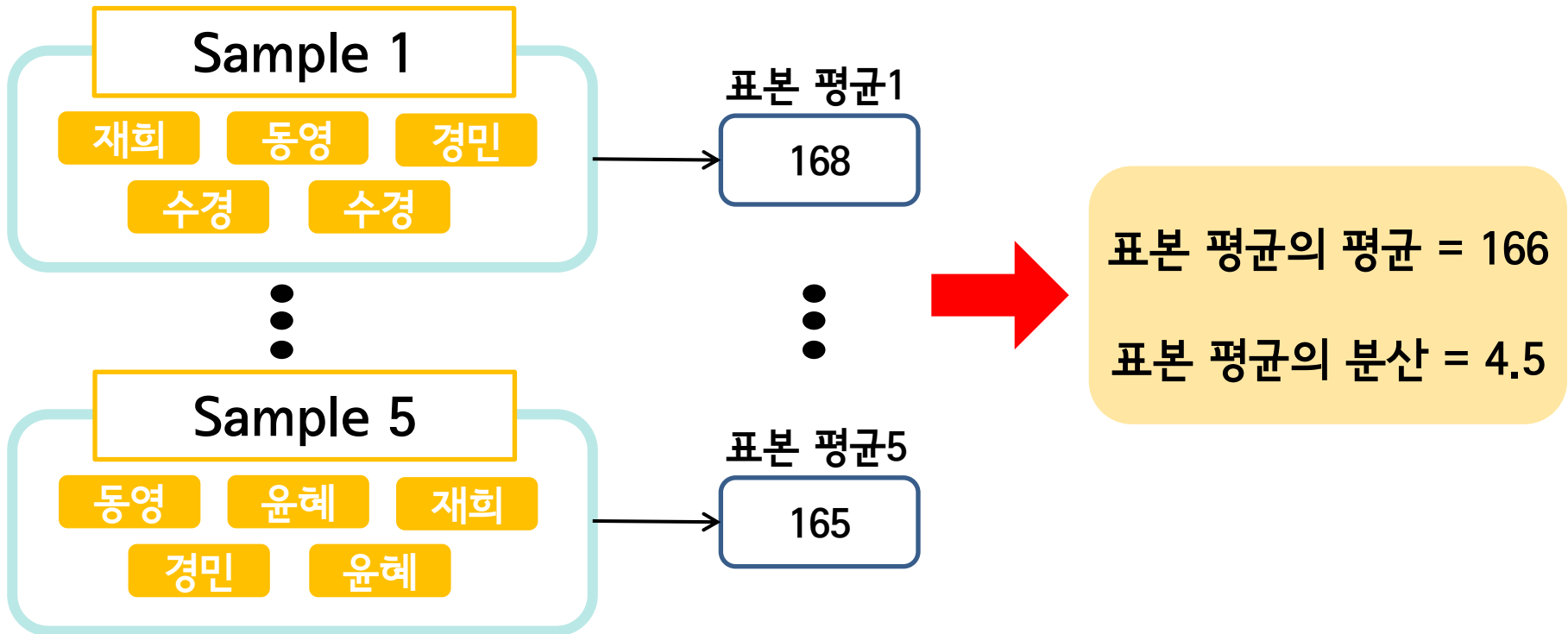


**순도가 높아지는 방향으로 트리 자르기!**



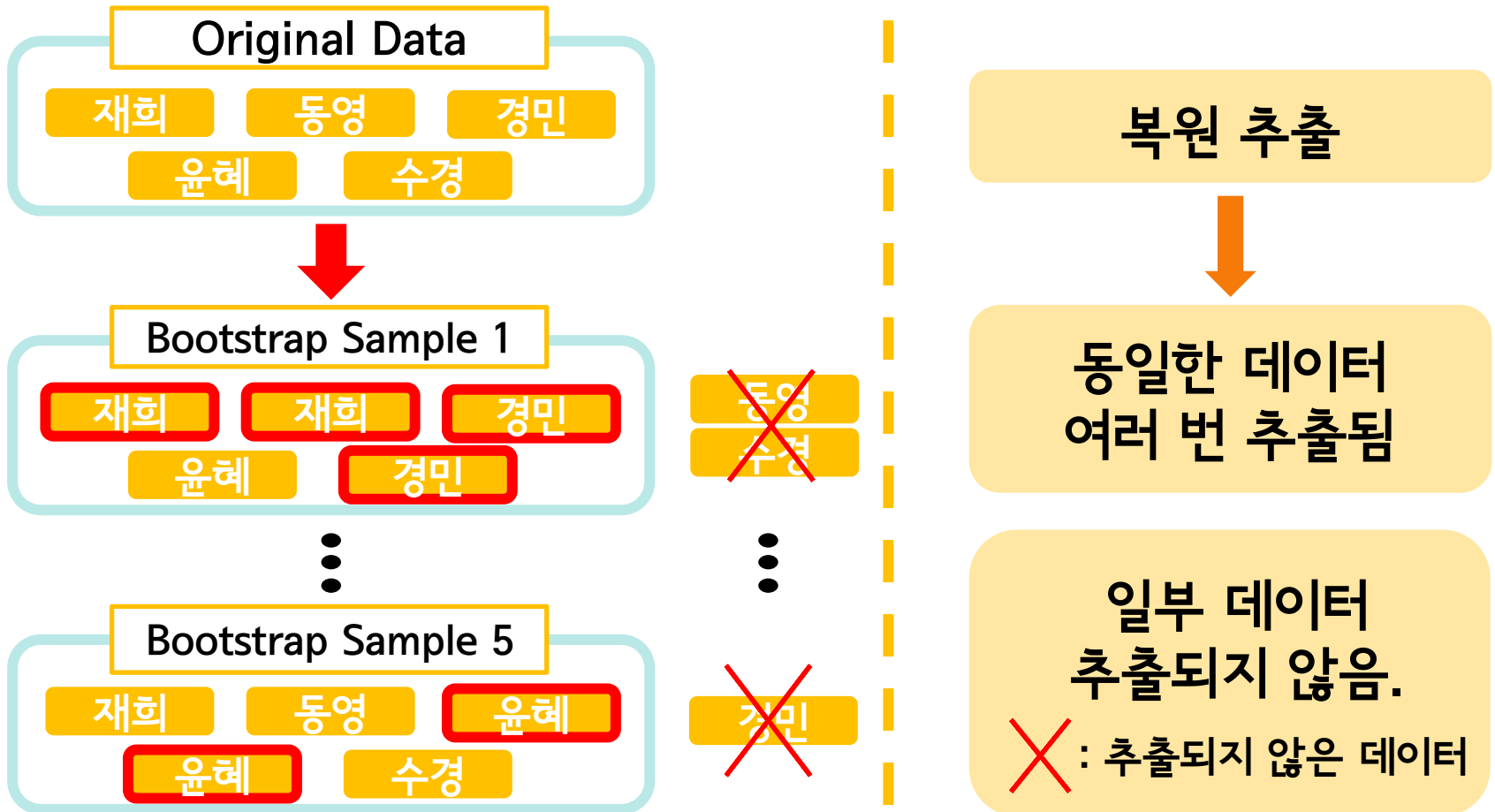
자른 아래의 관측값 버리는 것 X  
나누기 이전으로 되돌리는 것!

## Resampling



여러 개의 샘플을 추출하면 추리 통계량을 알아낼 수 있다!

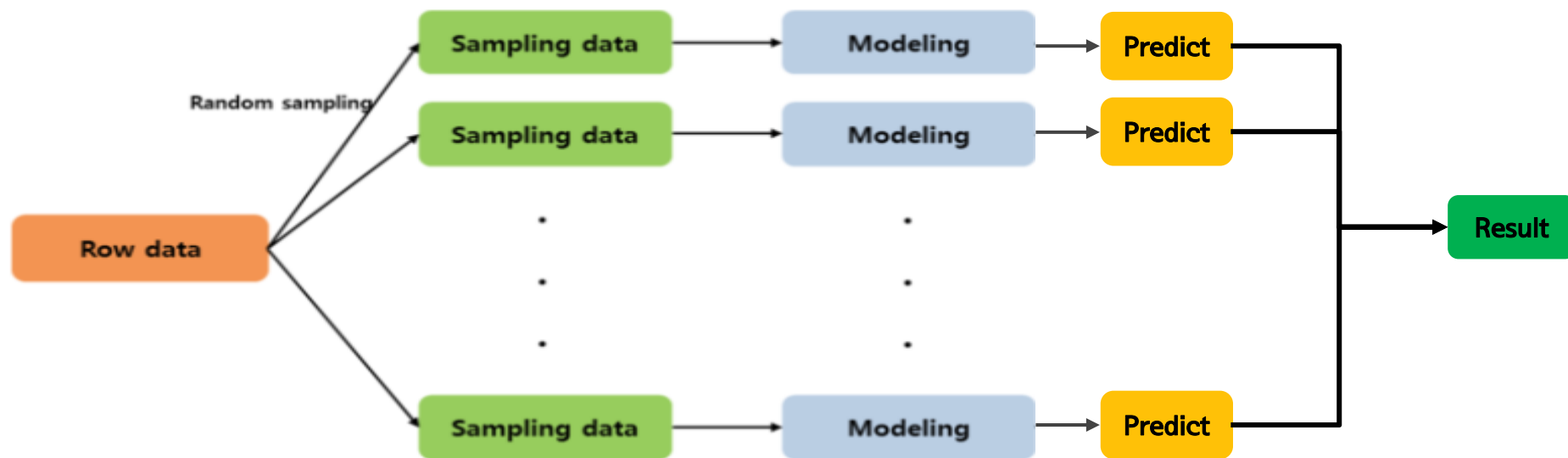
## Bootstrap





## Bagging

Bagging = Bootstrap + Aggregating



Bootstrap으로 얻은 표본에 대해 모델링한 결과를 합침

## Bagging 예측값 산출 방법

### 회귀 문제

$$\hat{f}_{bag}(x) = \frac{1}{B} \sum_{b=1}^B \hat{f}^{*b}(x)$$

B개 트리의 예측값의 **평균**

### 분류 문제

$$\arg \max_k \sum_{b=1}^B I(\hat{f}^{*b}(x) = k)$$

B개 트리의 예측값 중  
**다수결로 가장 많이 나온 값**

## Bagging

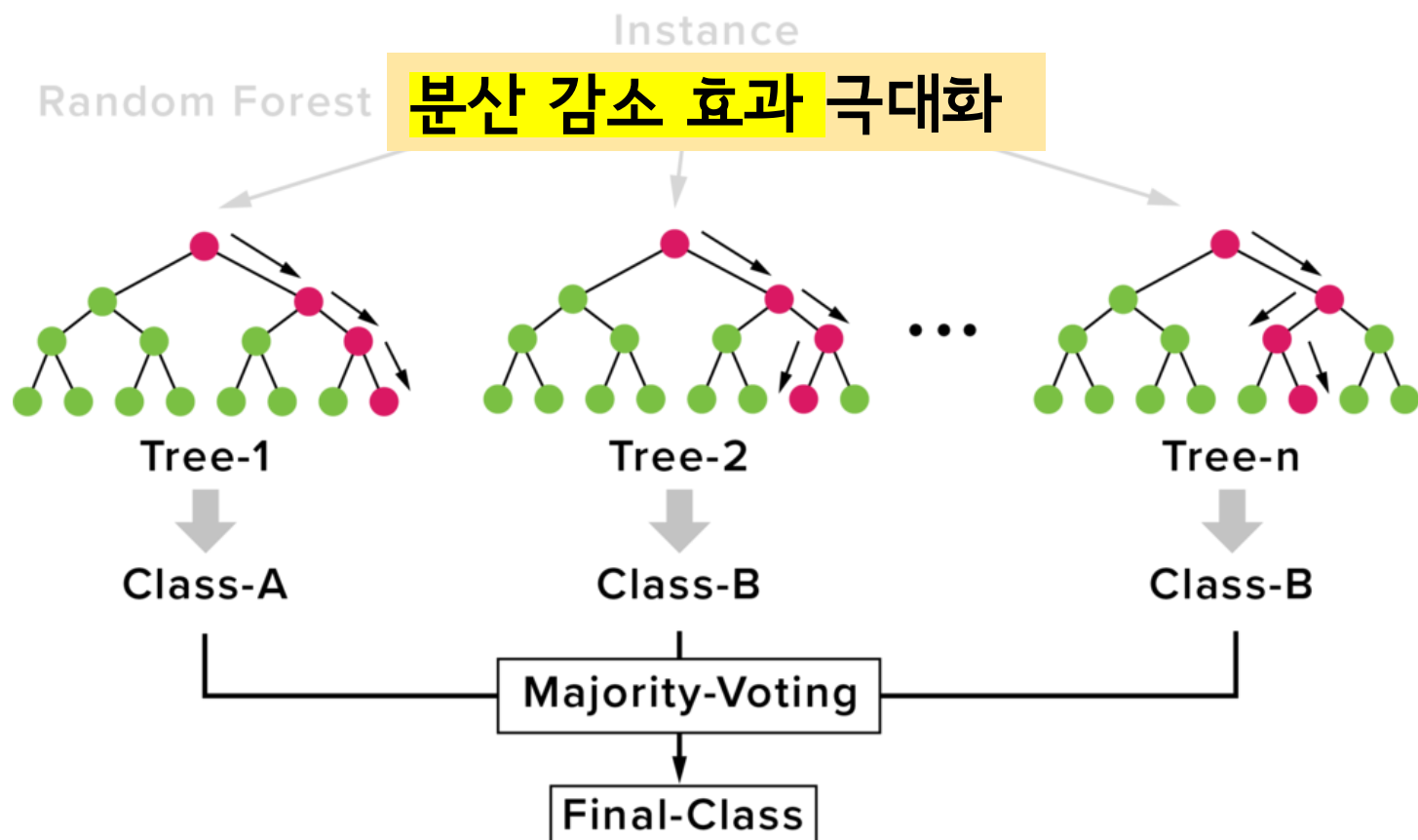


분산  편향 

분산  편향 

Bagging 결과 예측력 UP, 해석력 DOWN 됨

## Random Forest



## Random Forest Parameter

변수의 개수



mtry



max\_features

# PARAMETER



ntree



n\_estimators

배깅한 샘플 개수

## 평가지표

회귀 모델에서 사용되는 대표적 평가지표

**MSE**

Mean Square Error

**MAE**

Mean Absolute Error

## 평가지표

분류 모델에서 사용되는 대표적 평가지표

**Accuracy**

**f1 score**



범주팀의 클린업을 기대해주세요~!