

Statistical Machine Learning

7주차

담당: 15기 염윤석

1. What is Ensemble?

2. Ensemble Methods

3. Ensemble Models

1. What is Ensemble Learning?

Ensemble

Ensemble learning

- 다수의 기본 분류 모델(base classifier, weak classifier)의 예측 결과를 종합하여, 정확한 예측 성능을 얻도록 하는 방법론

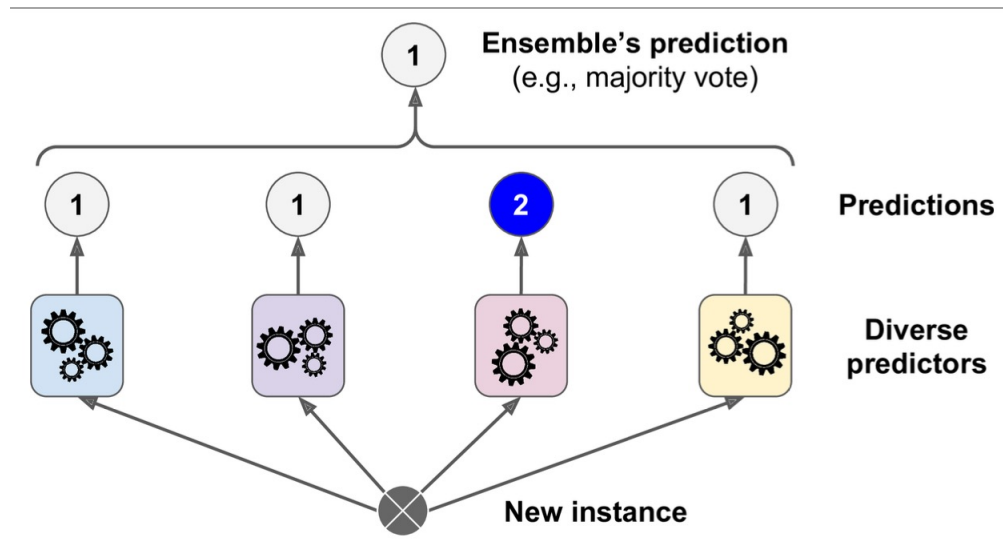
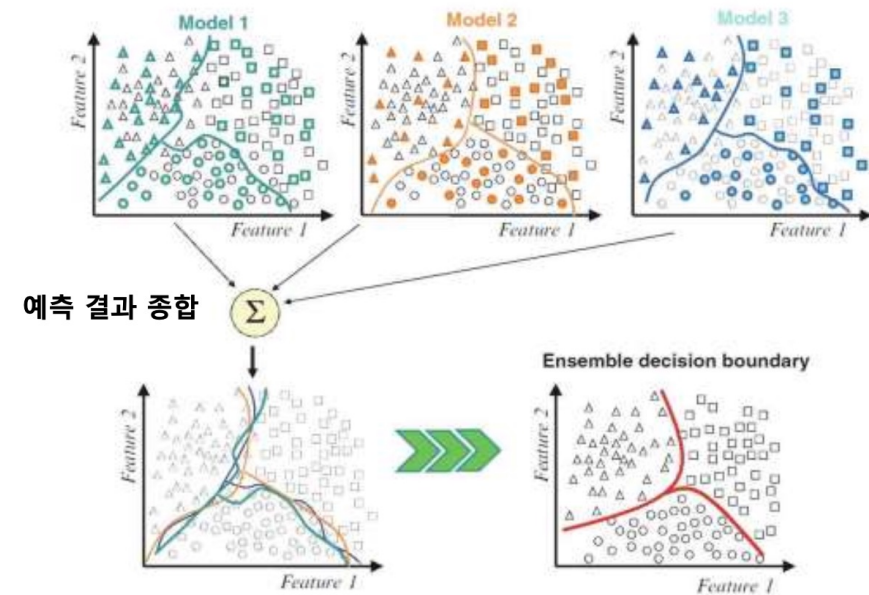


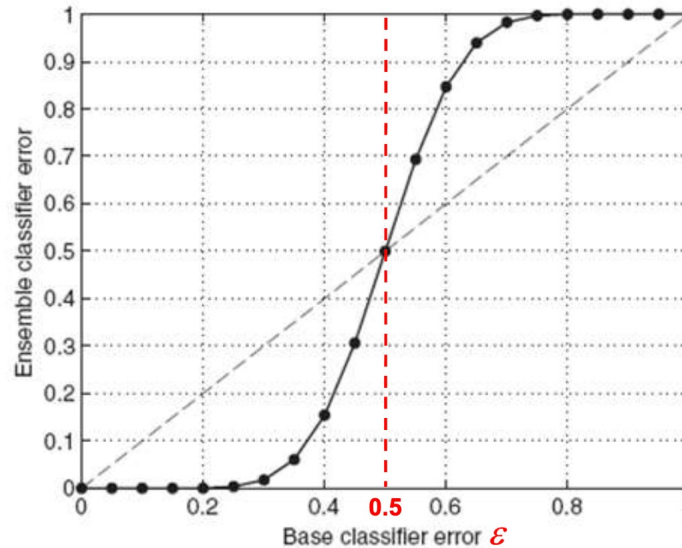
Figure 7-2. Hard voting classifier predictions



Ensemble

Example

- 25 base classifiers
- Error rate $\varepsilon = 0.35$
- Each independent
- Ensemble classifier : Majority vote



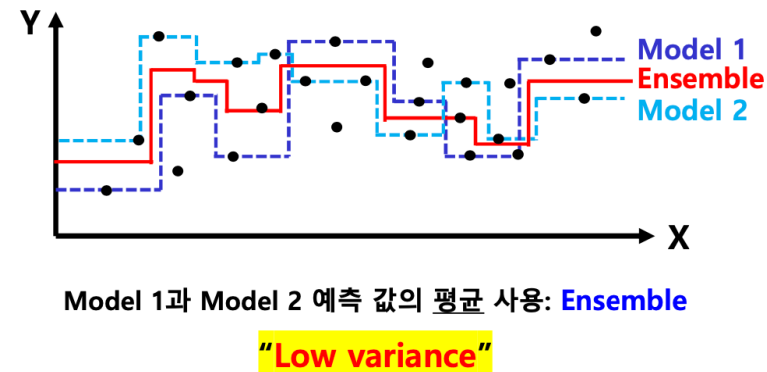
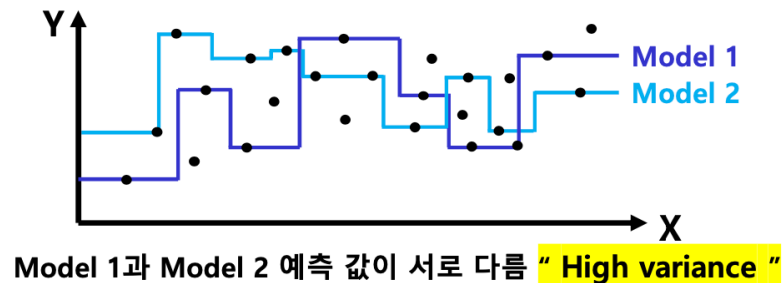
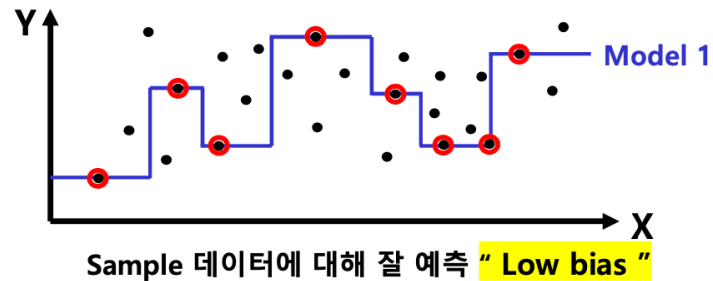
$$e_{ensemble} \sim \text{Binomial}(25, 0.35)$$

$$P(\text{incorrect classifier} \geq 13) = e_{ensemble} = \sum_{i=13}^{25} \binom{25}{i} \varepsilon^i (1 - \varepsilon)^{25-i} = 0.06$$

Ensemble

Ensemble Learning

- Reduce Learning error
- Reduce Bias
- Reduce Variance



Ensemble

Voting

Hard Voting
Soft Voting
Weighted Voting

Stacking

Meta level Learning
Blending

Bagging

Bootstrap + Aggregating

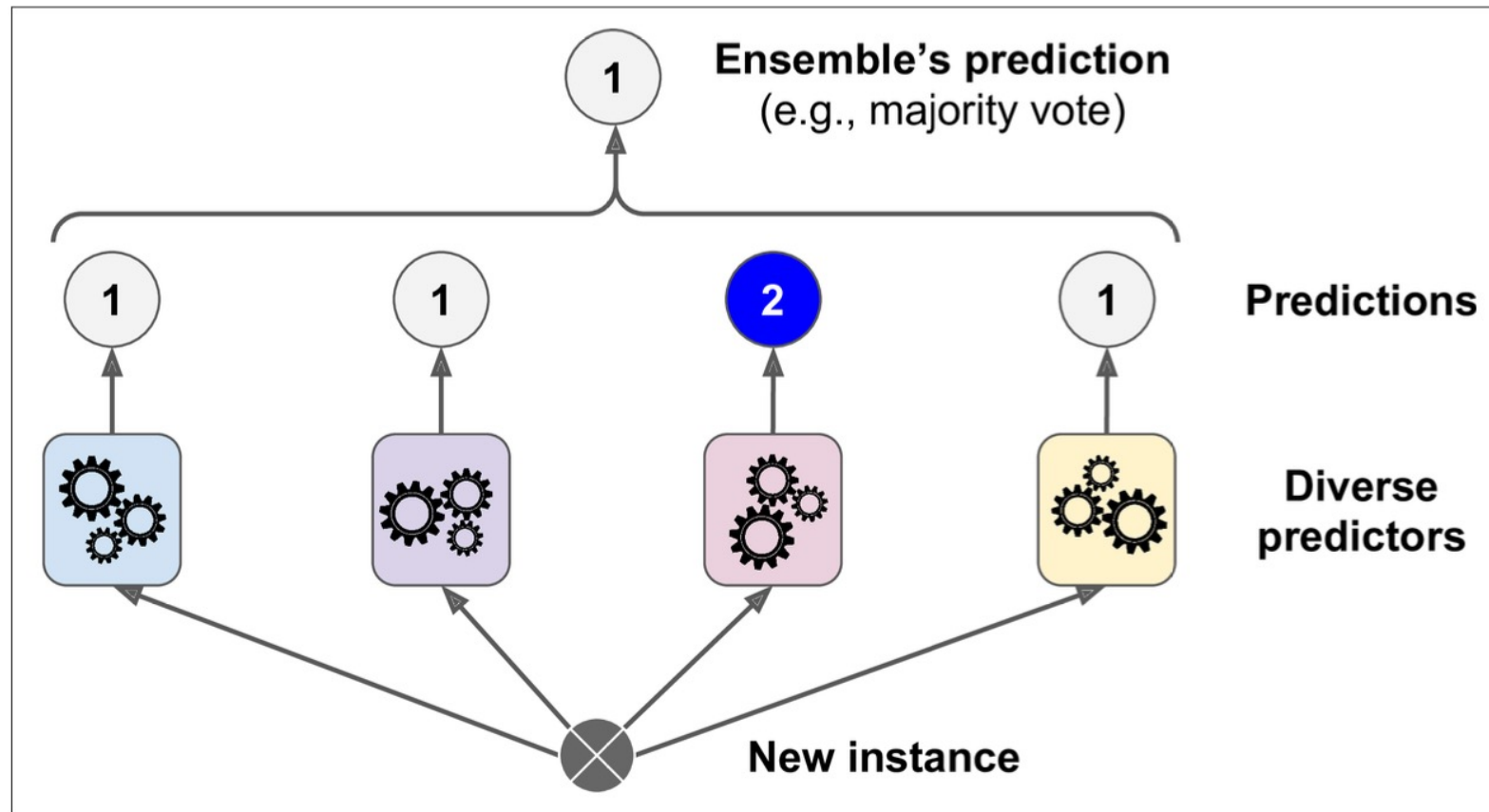
Boosting

Error learner

2. Ensemble Methods

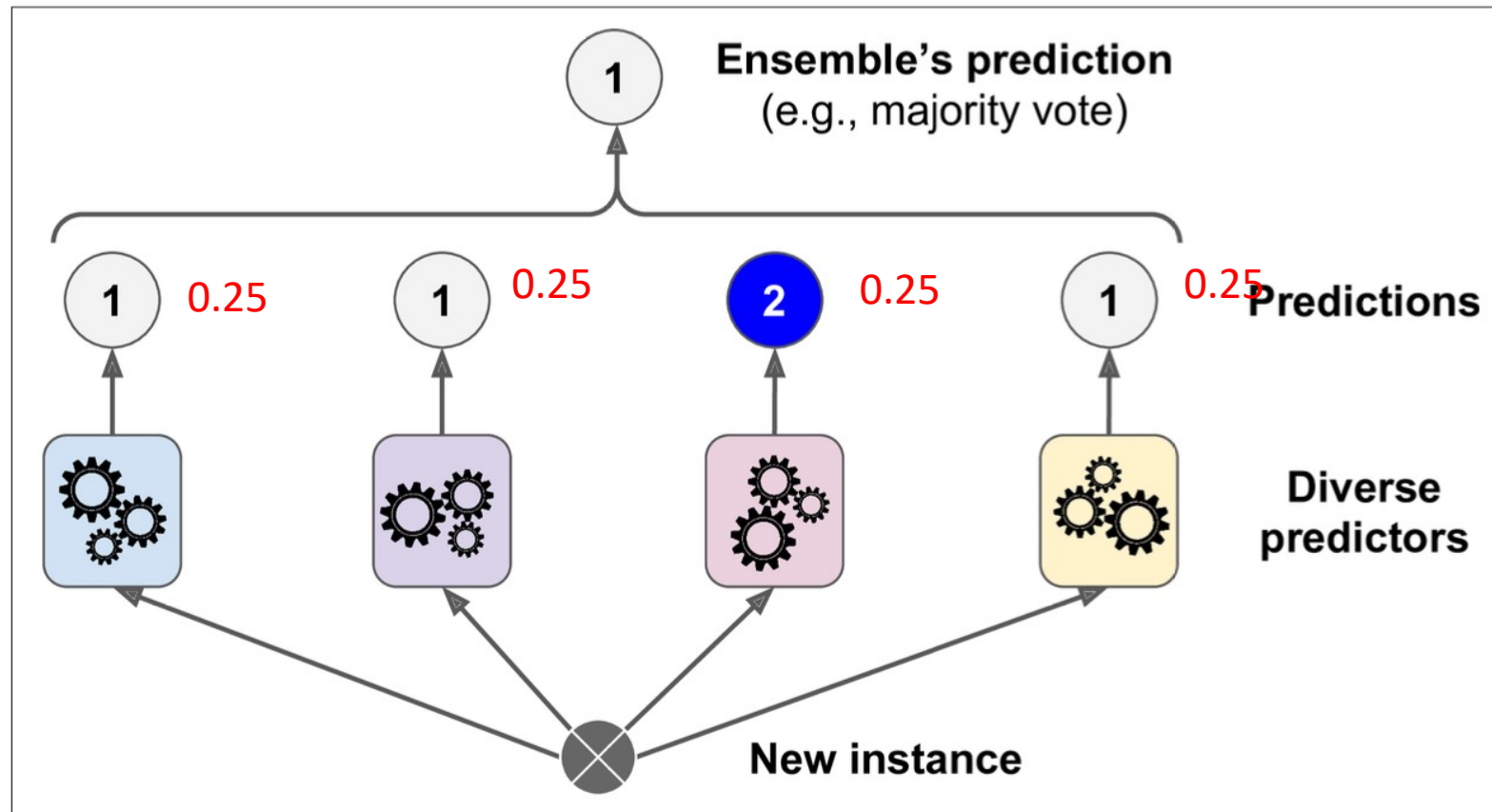
Voting

Hard Voting : Majority Voting



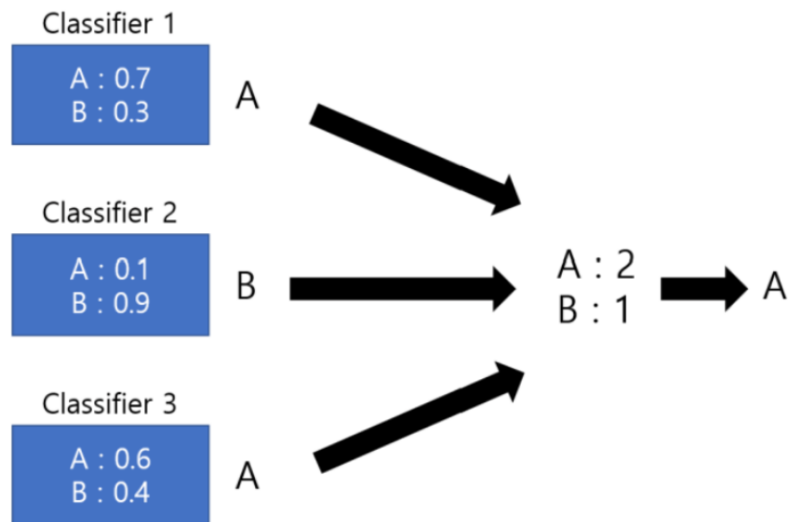
Voting

Soft Voting : Average Voting

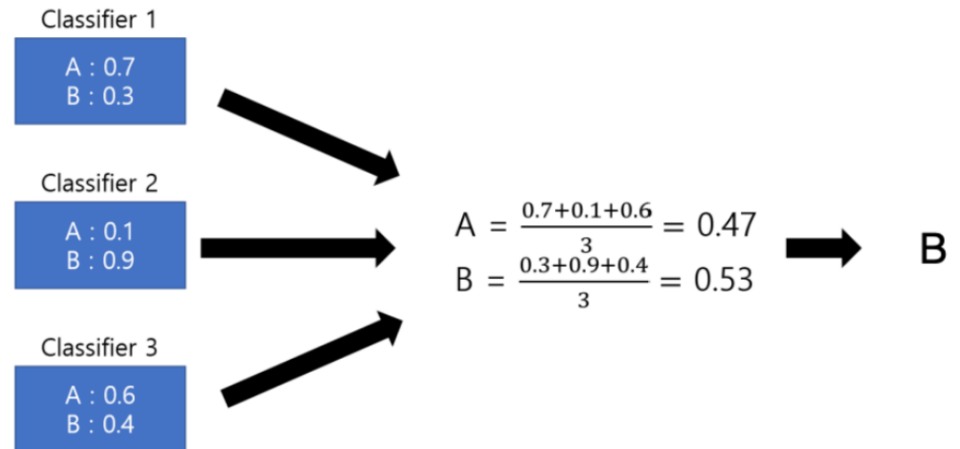


Voting

Hard Voting

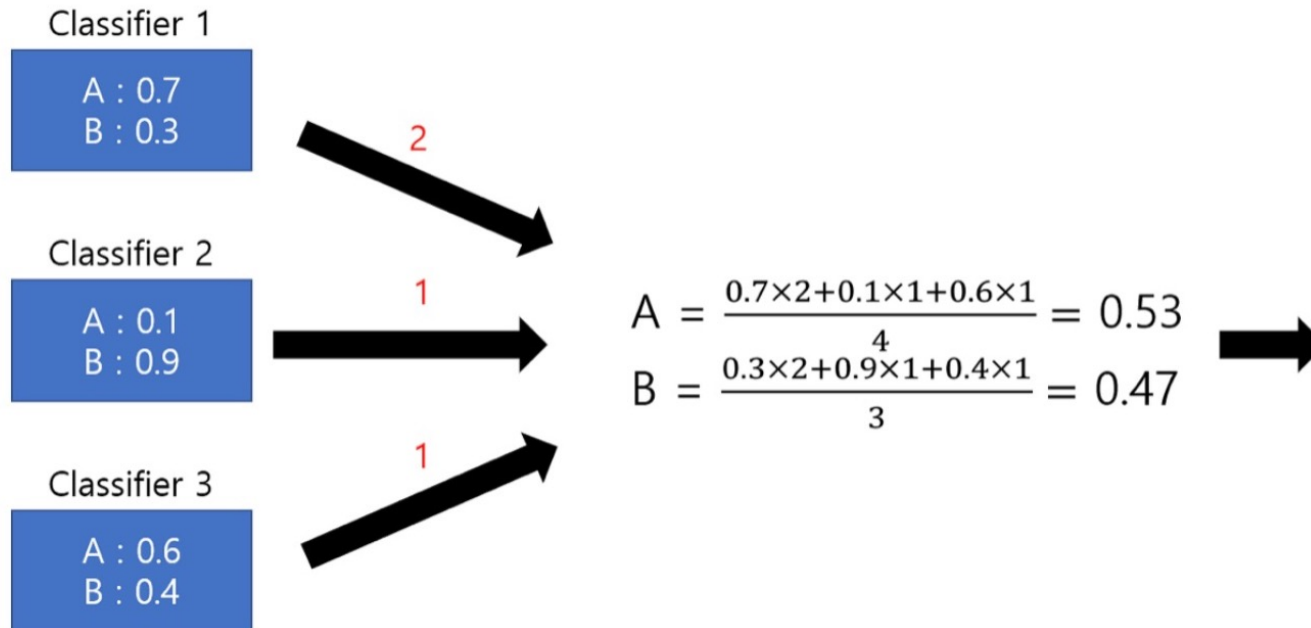


Soft Voting



Voting

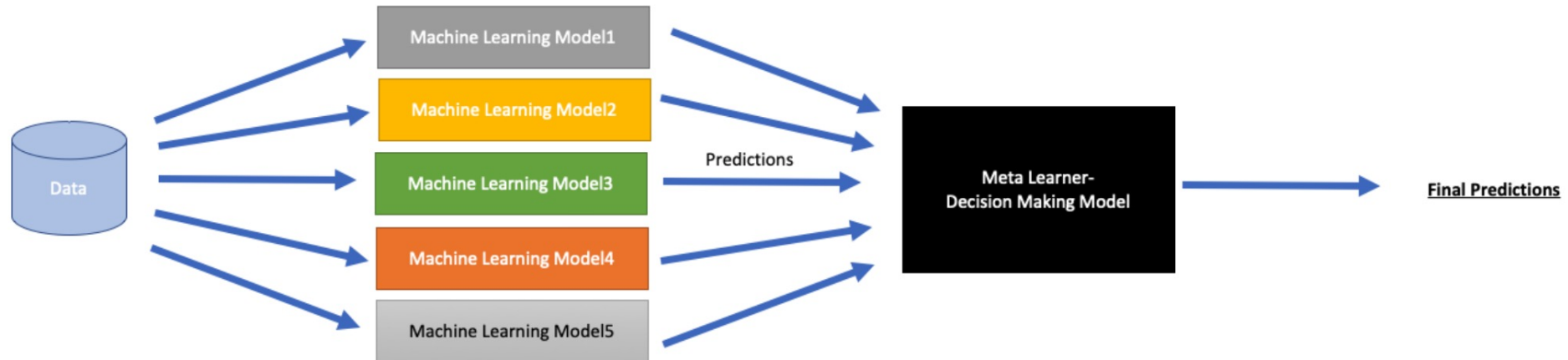
Soft + Weighted



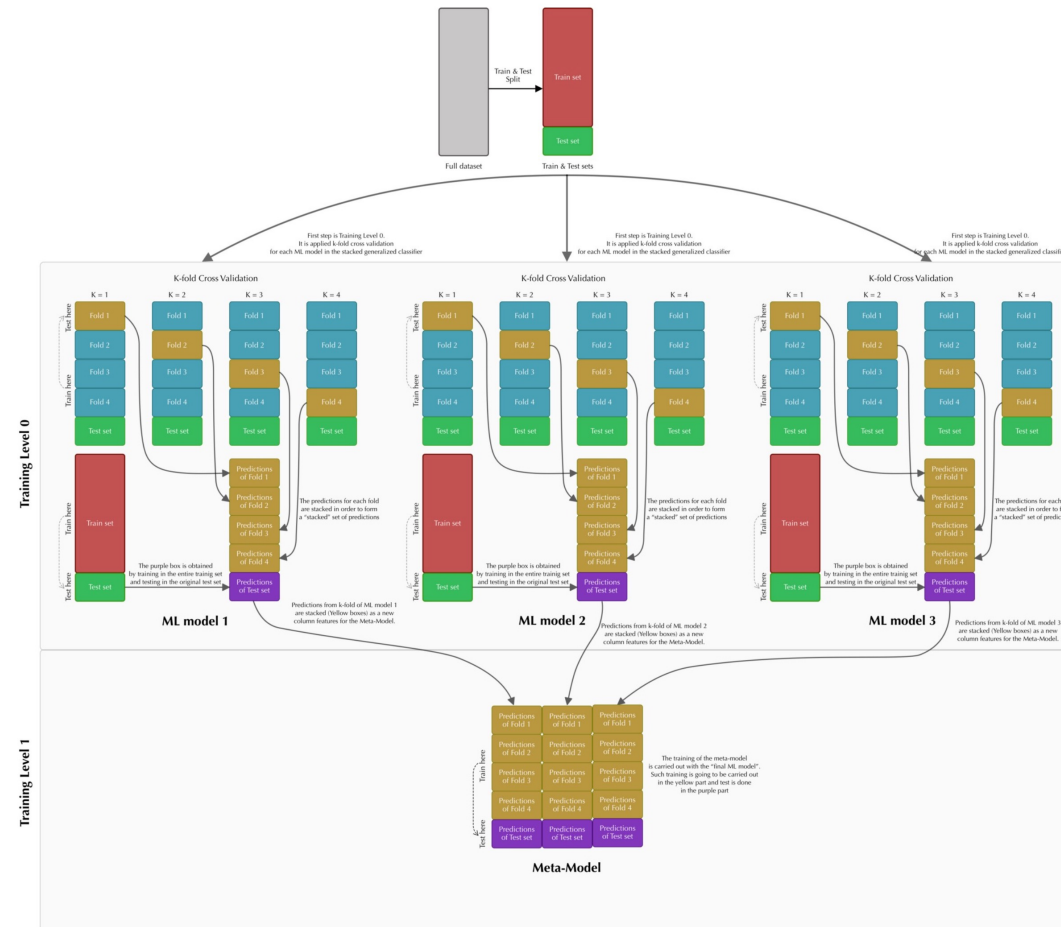
Stacking

Stacking Generalization

- Meta-learning model
- 개별 모델의 예측값을 다시 input으로 사용
- K-fold cv
- Step 0 : 각 weak model에 k-fold cv를 적용하여 예측 데이터를 형성
- Step 1 : step 0에서 만든 예측 데이터를 stack.하여 meta-model을 train 및 예측



Stacking

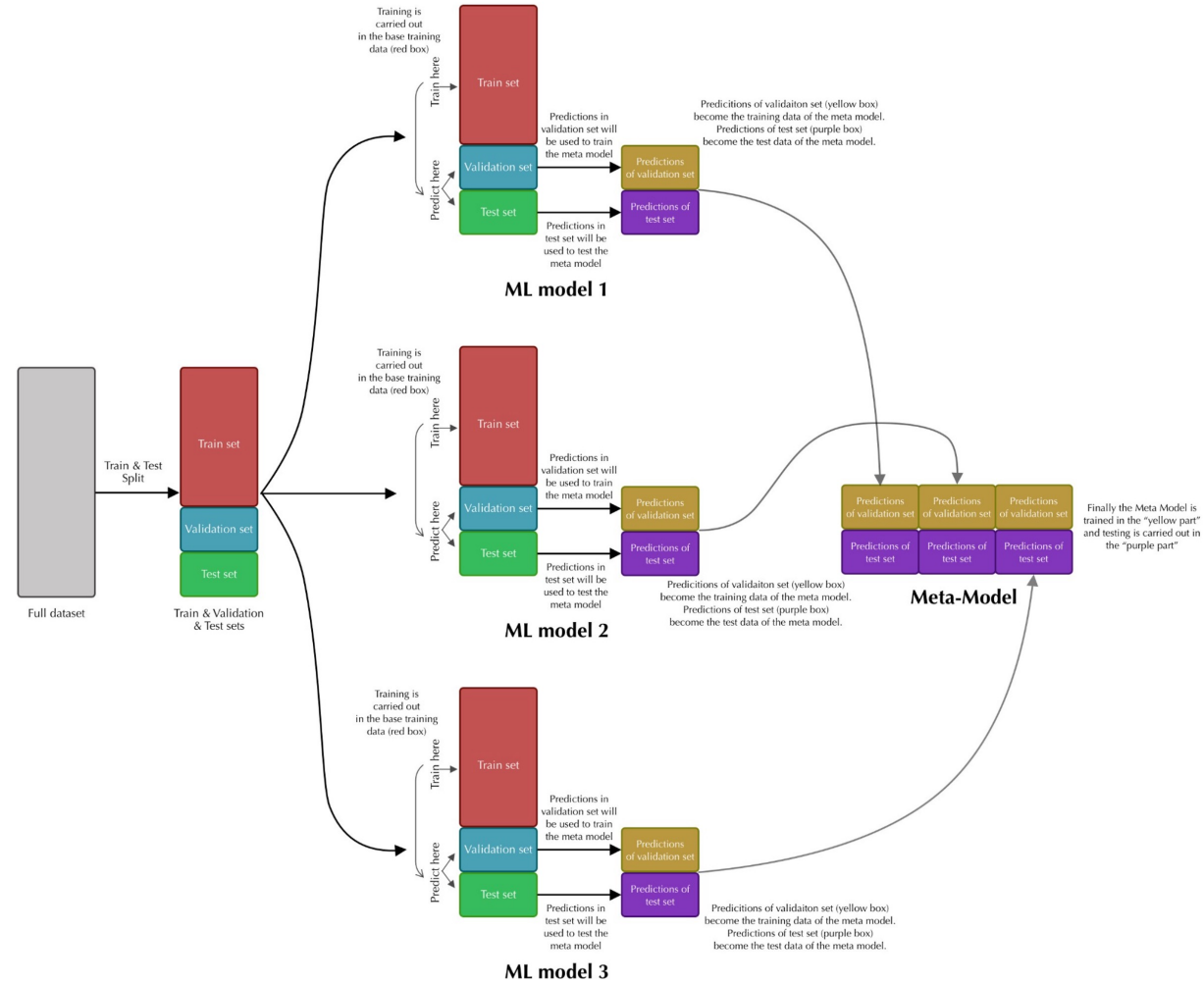


Blending

Blending Generalization

- Meta-learning model
- 개별 모델의 예측값을 다시 input으로 사용
- One-hold out

Blending



Bagging

Bagging = Bootstrap + Aggregating(Average)

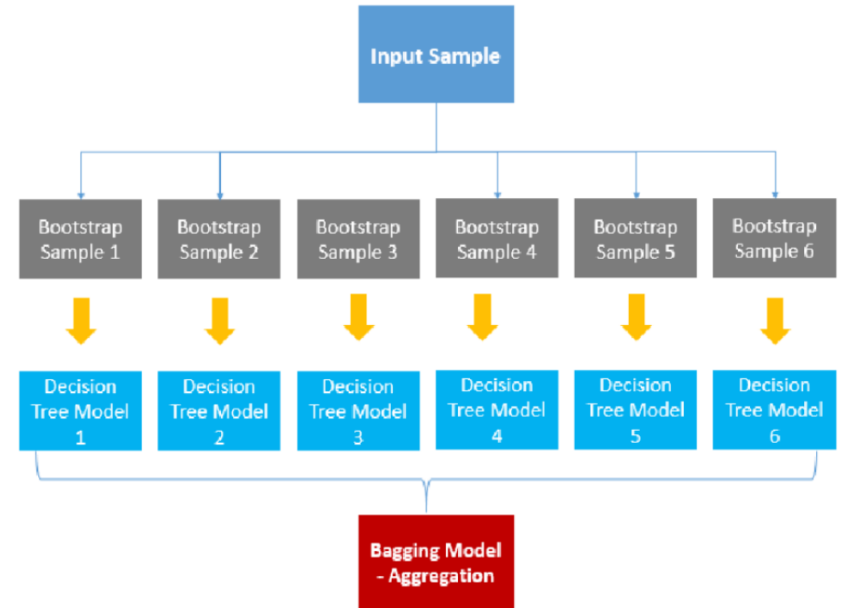
Bootstrap : sampling with Replacement → Variance 개선

Probability that a record is chosen by bootstrap
(N records, N sample size)

$$= 1 - \left(1 - \frac{1}{N}\right)^N$$

If N is large enough, then $\lim_{N \rightarrow \infty} 1 - \left(1 - \frac{1}{N}\right)^N = 1 - e^{-1} = 0.632$
63.2% of original train dataset

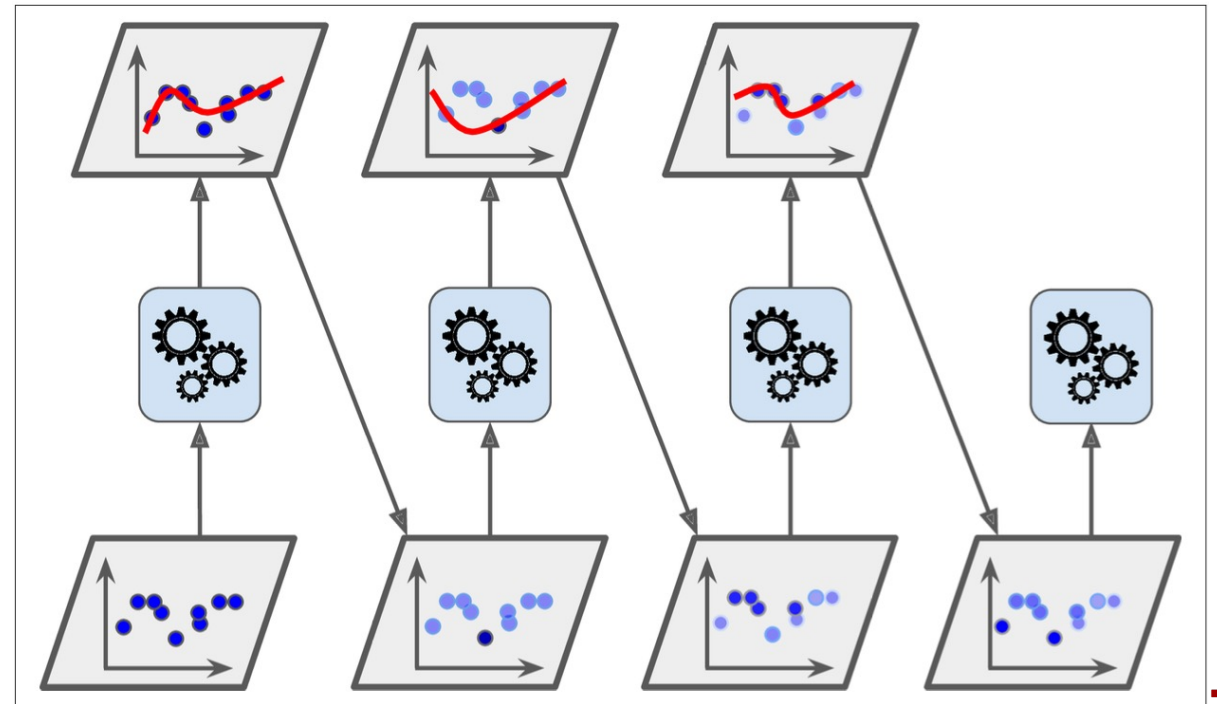
Aggregating : Majority Voting, Weighting, Soft voting



Boosting

Boosting

- 오답을 다시 학습
 - 예측이 틀린 데이터가 다시 뽑힐 가중치가 높아진다.
 - 이전 모델이 잘못 예측한 부분을 집중적으로 학습
- Bias 개선



Bagging & Boosting

<Normal>

10개년 6,9 수능 문제 1회독

<Bagging>

10개년 6,9, 수능 전체 문제에서
랜덤 복원 추출
→ 10번 반복

<Boosting>

10개년 6,9, 수능 전체 문제에서
랜덤 복원 추출
이때, 틀린 문제는 반드시 포함해서 추출
→ 10번 반복

10개년 6,9, 수능 문제 1회독
→ 틀린 문제만 뽑아서 다시 1회독
→ 다시 틀린 문제만 뽑아서 1회독
→ 다시 틀린 문제만 뽑아서 1회독
→ (반복)

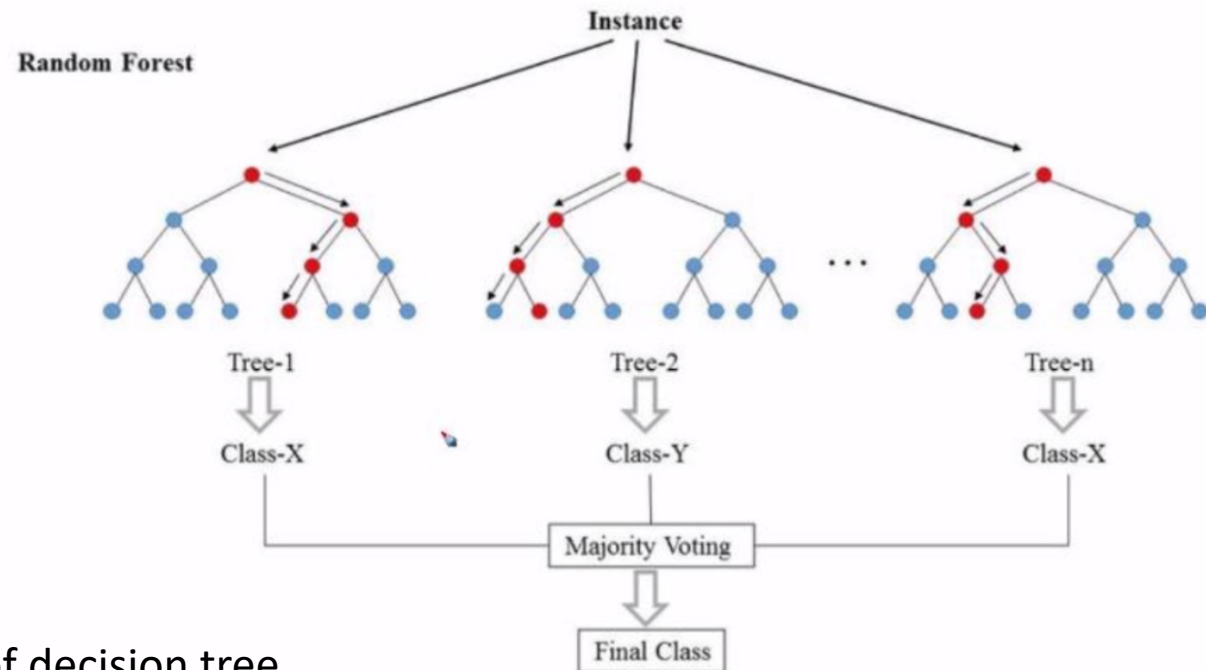
3. Ensemble Models

RandomForest

Feature Bagging → RandomForest

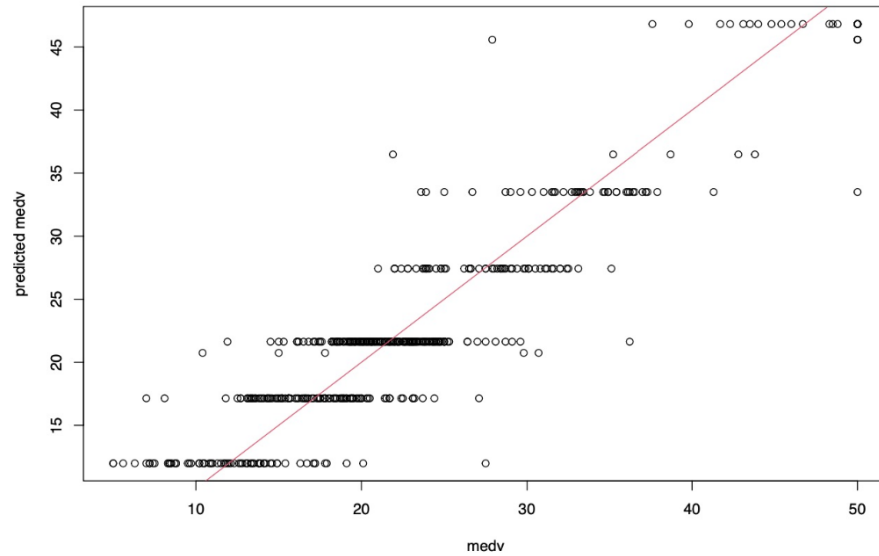
RandomForest Decision Tree Generation

- **Forest-RI(random input)**
randomly select F features
to split each node of decision tree
- **Forest-RC(randomly combined)**
 F randomly combined new features
(F linear combination)
- **Randomly select**
one of the F best splits at each node of decision tree

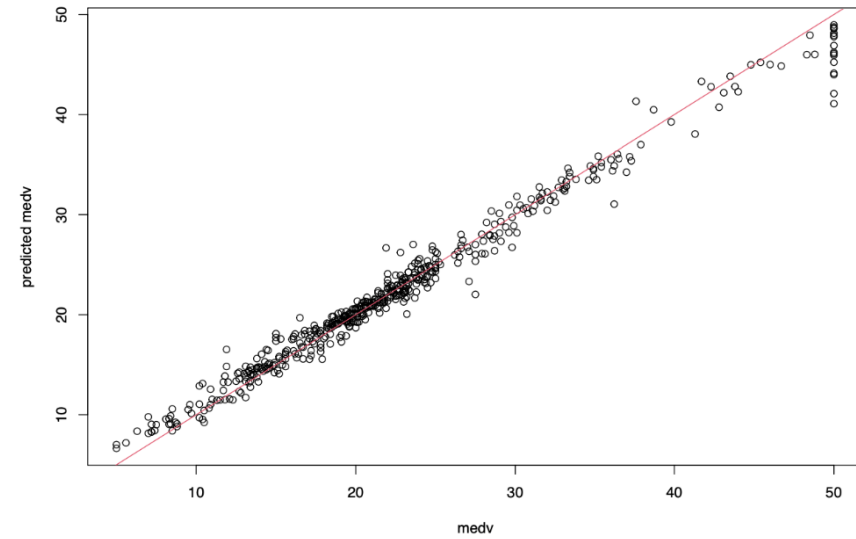


RandomForest

Single Tree



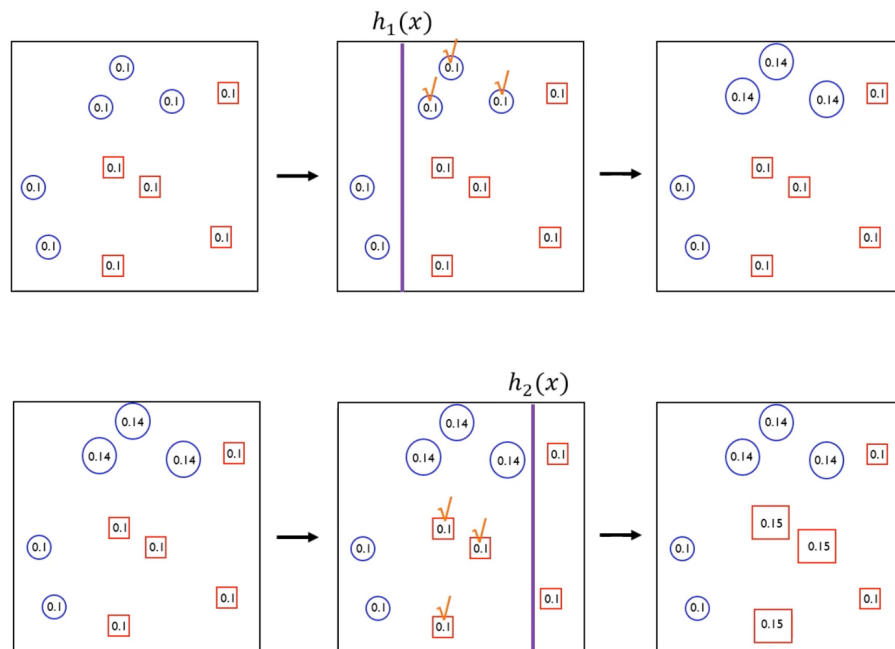
Random Forest



Adaboost

Adaboost : Adaptive + Boosting

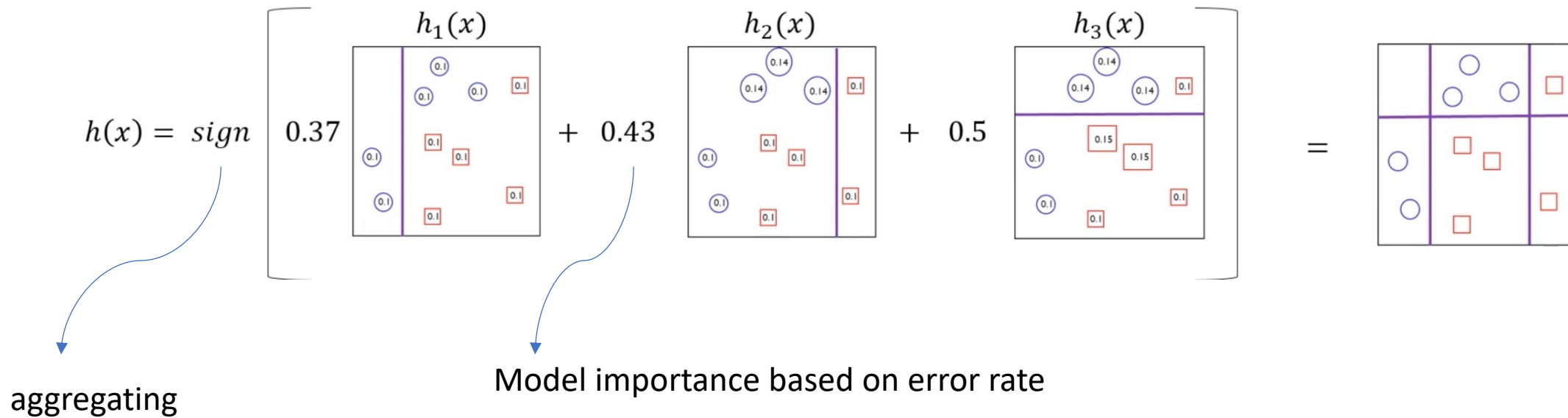
- **Adaptive** : 이전 모델이 잘못 분류한 데이터의 가중치를 adaptive하게 변경
- **Boosting** : 이전 모델이 잘못 분류한 데이터들을 중심으로 학습



정분류 sample : 그대로
오분류 sample : 가중치 \uparrow

Adaboost

$$h(x) = \text{sign} \left[\sum_{i=1}^{m=3} \alpha_j h_j(x) \right]$$



Ensemble models

- RandomForest
- ExtraTrees
- Adaboost
- GradientBoost
- XGBoost
- LightGBM
- CatBoost

수고하셨습니다!

해당 세션자료는 KUBIG Github에서 보실 수 있습니다!
Ensemble 및 autoML 패키지 소개가 있습니다.