

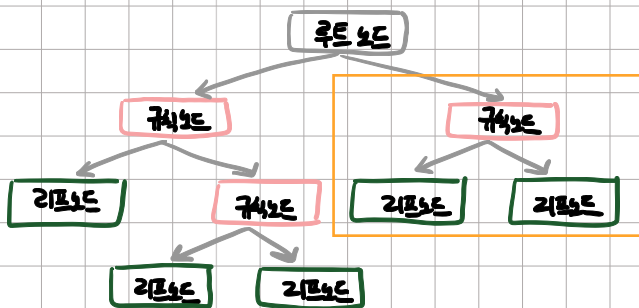
Week 2 - Ch4. 분류

01 분류의 개요

- 분류: 학습 데이터 주어진 데이터의 파와 레이블값을 이상형
있고리즘으로 학습해 모델 생성
→ 모델에 새로운 데이터 같이 주어진다면 미리 레이블값 예측

02 결정 트리

- 결정 트리 (Decision Tree)



- 내부노드: 규칙조각이 됨
- 리프노드: 정해진 클래스 값
- 샘플링: 새로운 규칙 조각마다 생성

- 트리를 어떻게 분할 할 것인가?

⇒ "균일도" 최대한 균일한 데이터 샘플을 구성할 수 있도록 분할

- 엔트로피 → 정보이득각
- 지니계수

- 결정 트리의 특징

- 장점
 - 쉽고, 직관적이다
 - 시간적 효율성도 ↓

- 단점: overfitting ~ 알고리즘 성능 ↓ → 트리크기 제한 ~ 특성

- 결정 트리 파라미터

CART - Decision Tree Classifier 클래스

- min_samples_split
- min_samples_leaf
- max_features
- max_depth
- max_leaf_nodes

03 앙상블 학습

- 앙상블 학습 (Ensemble Learning):

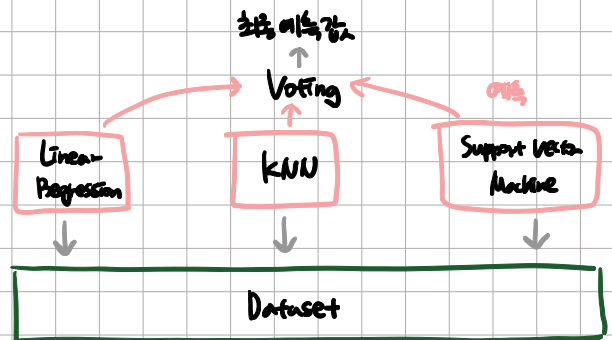
여러개의 분류기 (Classifier)를 생성, 그 예측 결합

→ 보다 강력한 최종 예측 도출하는 기법

- 정형 데이터 분류 ~ 뛰어난 성능.

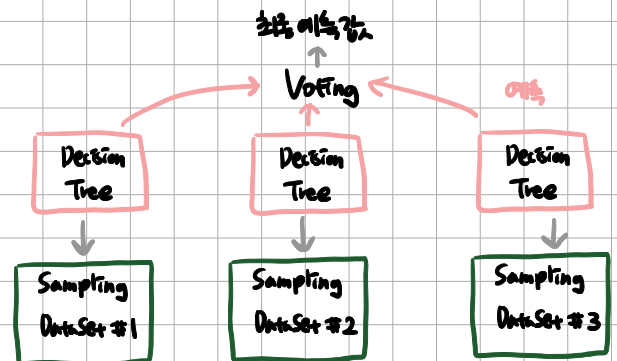
- 보팅 (Voting):

서로 다른 알고리즘을 가진 분류기 결합



- 배깅 (Bagging):

같은 알고리즘 기반의 분류기, 데이터 샘플링을 다르게



↳ 부트스트래핑 (Bootstrapping) 분할 방식.
Cv와 다르게 dataset 간 중복 허용

- 부스팅 (Boosting):

앞에서 학습한 분류기 ~ 예측 틀림 → 다음 분류기 기출치 보충한편 학습, 예측 전체

- 하드 보팅 (Hard Voting): 다수결로 결정

소프트 보팅 (Soft Voting): 분류기의 레이블값 결합
확률을 평균내서 가장 높은 것으로

04 랜덤 포레스트

- 랜덤 포레스트 : 다양한 대표적인 알고리즘. 앙상블을 기본으로.
- Random Forest Classifier
 - `n_estimators`
 - `max_features`
 - `max_depth`, `min_samples_leaf` 등 각 트리의 사용되는 파라미터