

② 배깅 (Bagging)

- 앞서 모형 결함에 사용자는 독립적인 모형의 수가 많을 수록 "성능 향상"이 일어난 가능성이 높을 것이다.
↳ 하지만, 각각 다른 모형을 사용하는데 한계가 있으므로 보통은 "배깅" 방법을 사용하며
같은 학습 모형을 쓰지만, 서로 다른 결과를 출력하는 다수의 모형을 만든다. \Rightarrow What is bagging?

- "배깅" (bagging)은 i) 동일한 모형과 ii) 동일한 모수를 사용하는 개인 부트스트래핑 (Bootstrap)과 유사하게 트레이닝 Data를 랜덤하게 선택해서 다수결 모형을 적용.

- Trainy Data를 선택하는 방법에 따라 다음과 같이 부른다.

- i) 같은 Sample Data를 중복하여 사용하자 \Rightarrow Pasting
 - ii) 같은 Sample Data를 중복하여 사용하면 \Rightarrow Bagging
 - iii) 데이터가 아니라 "특성 독립 변수" 중에 "일부 자원"을 선택하면 \Rightarrow Random Subspaces
 - iv) Data Sample과 독립 변수 자원 모두 일부를 random하게 사용하면 \Rightarrow Random Patches
- ↳ 성능 평가 시에는 트레이닝 Data가 아닌 다른 Data를 사용할 수 있다.

↳ \Rightarrow OOB (Out-of-Bag)

③ 랜덤 포레스트 (Random Forest)

↳ 의사결정 모형 (Decision Tree)를 개별 모형으로 사용하는 모형 결함 방법을 알린다.

↳ 배깅과 마찬가지로 Data 샘플의 일부를 선택하여 사용한다.

↳ 하지만, 노드 분할 시 모든 독립 변수들을 비교하여 최선의 독립 변수를 선택하는 것이 아니라, 독립 변수의 자원을 랜덤하게 감소시켜 그 중에서 독립 변수를 선택한다.

\Rightarrow 이렇게 함으로써, 개별 모형들 사이의 상관관계가 줄어들기 때문에, 모형 성능의 변동이 감소하는 효과가 있다.

↳ 이러한 방법을 극단적으로 이용한 것이

"Extremely Randomized Trees" 모형

↳ 각 노드에서 랜덤하게 독립 변수를 선택한다.

↳ 랜덤 포레스트의 "중요도"는 각 독립 변수의 중요도 (Feature Importance)를 계산할 수 있다.

\Rightarrow 포레스트 안에서 사용된 모든 노드에 대해 어떤 독립 변수를 사용했고

그 노드에서 얻은 정보량 (Information Gain)을 구할 수 있으므로 각각의 독립 변수들이 얻어낸 Information Gain의 평균을 비교하면, 어떤 독립 변수가 중요한지 비교 가능