

Week 5

Sunday, August 10, 2023 5:23 PM

1. Ensemble

: 어떤 Base 모델 예측을 통합하여 예측 정확성을 향상 / Better than random.

$$\text{Ensemble} = \sum_{i=1}^n \binom{N}{i} e^i (1-e)^{N-i} : \text{양자수 확률}$$

Decision Tree - low computation & Non-parametric \rightarrow Ensemble

데이터 블록 문제 X

2. Random Forest

: Diversity & Randomness

\downarrow \downarrow
특정 모집단에서 \rightarrow Random Subspace.

어떤 Training data \rightarrow Bagging

Bagging : Bootstrap Aggregating \rightarrow Bootstrap 샘플링에 상관을 보일 확률.

복수로 훈련을 통해 정규 데이터 크기로 샘플링.

Classification

- Majority Voting : $\text{Ensemble}(\hat{y}) = \arg \max_i \left(\sum_{j=1}^n I(y_j = i), i \in \{0, 1\} \right)$

- Weighted Voting : $\text{Ensemble}(\hat{y}) = \arg \max_i \left(\frac{\sum (\text{Train Acc}_j) \cdot I(y_j = i)}{\sum (\text{Train Acc}_j)}, i \in \{0, 1\} \right)$

or

$= \arg \max_i \left(\frac{1}{n} \sum P(y=i), i \in \{0, 1\} \right)$

Random Subspace

① 훈련 및 테스트 쪽 쓰는 임의 변수 추출 선택

② 선택된 임의 변수 중 카탈로그 변수 선택 \rightarrow 봉착 희석, 훈련 및 테스트 쪽의 선택하여 그 변수들만 고려

③ full-grown tree 를 떠올리기 때문.

④ Generalization Error $\leq \frac{\bar{P}(1-s^2)}{s^2}$: upper bound (\bar{P} : tree 개수의 평균 샘플링
 s : 훈련과 예측한 Tree와 맞은 예측한 Tree 수 차이 평균)

Random forest는 Tree \uparrow ~ 예측오류로
error가 수렴

① 개별 tree 정확도 $\propto s$

② Bagging + Random Subspace \rightarrow 훈련 오류율, 예측오류율 감소

② Bagging + Random Subspace → 훈련과 테스트에 대한 흐름을 확장하여 학습률을 줄여

③ Feature Importance.

원래 데이터에 따라 OOB Error 계산 → 평균으로 나온 데이터의 OOB Error 계산

$$r_i \quad x_i \quad c_i$$

$$\begin{aligned} d_i &= e_i - r_i \quad \left\{ \bar{d} = \frac{1}{t} \sum d_i \right. \\ &\quad \left. S_d^2 = \frac{1}{t-1} \sum (d - \bar{d})^2 \right. \end{aligned} \Rightarrow x_i \text{ importance } V_i = \frac{\bar{d}}{S_d} \begin{array}{l} (\text{기여도}) \\ (\text{스케일링}) \end{array}$$

4. Hyperparams.

• Decision Tree 개수 (정답률이 2000)

• 노드 분할의 주제와 선택되는 변수의 수

Classification : sort(변수의 수)

Regression : 정수 / 3