

Decision Tree : 의사결정 규칙을 나무구조로 나타내 전체 데이터를 선택적으로 분류하거나 예측하는 방법

↳ 좋은 Decision Tree : 좋은 정확도를 내면서, simple한 것을 선택한다! → 각 노드가 최대한 한 가지 클래스만 가지도록

↳ 좋은 기준이란 어떻게 정하는 것일까?
 분도 : Entropy, Gini index
 어떤 기준으로 노드를 분류 : ID3, CART 알고리즘

$$\text{Entropy} \left(- \sum_{k=1}^m p_k \log_2(p_k) \right)$$

: 데이터의 불확실성을 나타낸다, 즉, Entropy 값 = 분도 값 = 순도 증가

$$\begin{cases} \text{불확실성 최소} = \text{순도 최대} : \text{엔트로피 } 0 \\ \text{불확실성 최대} = \text{순도 최소} : \text{엔트로피 } 1 \end{cases}$$

ID3 알고리즘

상위 노드의 Entropy에서 하위 노드의 Entropy를 뺀 값

Entropy 차를 통해 Information Gain을 : 즉, Information Gain 값이 클수록, 엔트로피를 많이 줄였다는 의미 = 엔트로피가 작아졌다

$$\text{Gain}(S, A) = E(S) - \underline{I(S, A)} = E(S) - \sum_i \frac{|S_i|}{|S|} \cdot E(S_i)$$

$$\hookrightarrow I(S, A) = \sum_i \frac{|S_i|}{|S|} \cdot E(S_i)$$

$$\text{Gini index} \left(\sum_{j=1}^m \frac{|D_j|}{|D|} \cdot \text{Gini}(D_j) \right)$$

데이터의 불확실성 분산 정도를 정량화 해서 표현한 값, Gini index 값 = 분도 값 = 순도 증가

$$\text{Gini}(A) = \sum_{j=1}^m \frac{|D_j|}{|D|} \cdot \text{Gini}(D_j) \quad \rightarrow \quad \text{Gini}(D_i) = 1 - \sum_{j=1}^m p_j^2$$

CART 알고리즘

Gini index를 이용한 알고리즘, Binary split을 전체로 분석함

▷ feature가 어떻게 나눠질까?

1. 전체 데이터를 모두 기준으로 분할 후 분도를 계산한다.
2. 분할수, 선택기준을 결정한다.
3. Label의 class가 바뀌는 수를 결정한다.

STEP1. 각 Feature에 대해 2중첩으로 정렬 → STEP2. Label의 class가 변하는 지점을 찾기

→ STEP3. 전체의 평균값을 기준으로 잡기 → STEP4. 각 기준에 대해 분할 후, Gini index 혹은 Entropy 계산

↳ 위와같은 과정을 통해 첫번째 기준점을 정한다. 이 과정을 반복하면 좋은 Decision Tree model이 만들어진다.

▷ 가지치기

모든 terminal node의 순이 100%인 상태를 Full tree라고 한다. 분기가 너무 많아 복잡함 문제가 발생한다

- Pre pruning (사전 가지치기) : 미리 최대 depth나 분할의 횟수 등을 미리 지정
- Post pruning (사후 가지치기) : 모델을 만든 후 데이터 포인트가 작은 노드를 삭제 / 병합

Naive Bayes Classifier

▷ 확률 기초

- 확률 : 특정한 사건이 일어날 가능성을 나타내는 것
- 조건부 확률 : 어떤 사건이 일어난 조건 하에서, 다른 사건이 일어날 확률

$$P(B|A) = P(B \cap A) = P(A \cap B) = P(A|B)P(B)$$

- 독립과 조건부 독립 :
 - 독립 : 한 사건이 일어날 확률이 다른 사건이 일어날 확률에 영향을 미치지 않는 상태
 $P(A \cap B) = P(A)P(B)$
 - 조건부 독립 : 한 사건이 일어났다는 가정 하에서, 서로 다른 두 사건은 독립인 상황
 $P(A, B|C) = P(A|C) \cdot P(B|C)$

▷ 베이즈 정리

"두 확률 변수와 사건 확률과 사후 확률 사이의 관계를 나타내는 정리"

$$P(H|D) = \frac{P(D|H) \cdot P(H)}{P(D)}$$

- $P(H)$: 사전 확률, 과거 정보를 토대로 현재로 판단한 파라미터(H)의 확률
- $P(H|D)$: 사후 확률, 관측 결과 사건 D가 일어난 조건 하의 파라미터의 확률
- $P(D|H)$: 사전 확률의 과거 정보를 잘 설명하는 정도, 모델 파라미터를 바탕으로 한 관측 결과의 확률
- $P(D)$: 사건 D의 발생 가능성

▷ 계산의 한계

변수가 늘어남에 따라 계산량이 증가함 → 해결책 : 조건부 독립을 가정!!

가정 : 공통 변수가 주어졌을 때, 입력 변수들이 모두 독립이다!! (조건부 독립 가정)

결과가 주어졌을 때, 여러 변수 밖에서 정해진 공통 확률은 각 공통 확률의 곱으로 충분히 잘 표현할 수 있다는 단순한 가정을 기점으로 한다.

$$P(A \cap B|C) = P(A|C)P(B|C)$$

알려야 할 파라미터의 수가 대폭 줄어든다. $P(X=x|Y=y) \Rightarrow dk$

Feature들의 공통 배경에서 계산이 수월해진다

- 라플라스 수정 : likelihood가 0이 되는 것을 방지하기 위해 사전 확률의 확률을 조정해서

$$PLAP = \frac{C(x) + 1}{\sum C(x) + 1}, \quad P(x|c) = \frac{\text{count}(x, c) + 1}{\sum_{x \in V} \text{count}(x, c) + V}$$

세변나 관변의 다양성을 고려한 가산

세변나 관변의 다양성을 고려한 가산