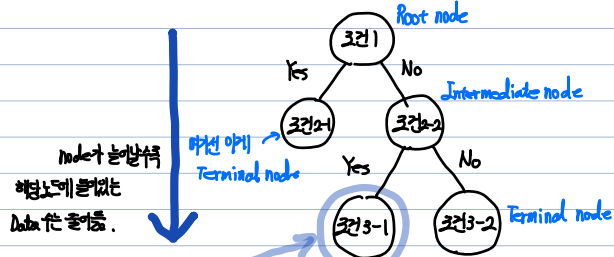


덱스 3주와 정리 Decision Tree

Unit 1. 의사결정 나무란?

- 의사결정 규칙을 나무구조로 나타내어 데이터를 분류/예측.

→ Data 세에 존재하는 패턴을 예측 가능한 규칙들이 조합으로 나타냄.
결국 이 조합이 나무와 함께 의사결정 나무라는 이름이 붙여짐.



What is good decision tree? → ① (목표한 정답을 가지는 문제해.) ② 이런 노드를 봤을 때 현존하는 Data가 더 단련에 좋은 Decision Tree! ③ 불분했다? 좋은 Decision Tree!

좋은 decision tree란 무엇인가?

Unit 2. ID3 알고리즘

Unit 2. ID3 알고리즘

Unit 2. ID3 알고리즘

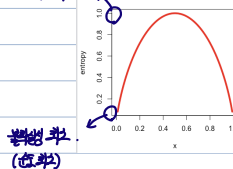
Entropy → (in Data,) 데이터의 불확실성.

→ entropy가 높다: 집단의 특성을 찾기가 어렵다. (불확실한 것인가...?)

→ ∴ entropy를 낮추는 방향으로 분류해야 한다.

$$-\sum_{k=1}^m P_k \log_2(P_k)$$

불확실성의 정도 (엔트로피)



← 엔트로피는 분포에 따라 달라진다.
(P(확률 X)가 쉬울수록 entropy가 낮아짐)

이런 entropy를 이용해 불확실성을 줄여주는 것을 찾는 알고리즘...

ID3 알고리즘 → entropy가 → Information Gain을 → 제일 크게 값이 나온 것을 기준으로 선택.

Information Gain = 정보 엔트로피 - 하위 엔트로피.

∴ Information Gain 값이 크면: 엔트로피를 많이 줄였다.

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

각 변수의 고위값 구하는 방법이고, 각 class 값에 대해 entropy 구하기

→ Information Gain 하기

→ Information Gain은 제일 크게 값을 갖는 변수를 기준으로 선택.

→ 조건에 Data에 대해 분류를 반복하기.

Unit 3. cart 알고리즘

Gini Index - Data의 불확실성을 측정하여 선택한다.

∴ gini index를 낮추는 방향으로 분류하기.

$$Gini(A) = \sum_{i=1}^2 \frac{|D_i|}{|D|} * Gini(D_i)$$

$$Gini(D_i) = 1 - \sum_{j=1}^2 P_{ij}^2$$

이런 Gini를 이용해 분류하는 알고리즘...

Cart 알고리즘

→ 데이터의 특성값을 골라나 분류할 것을 계산한다.

→ Gini Index를 각 변수마다 계산하여 가장 작은 Gini index 값을 갖는 변수가 분류를 선택한다.

→ 대신 ID3와 달리 Binary split을 선호한다.

Unit 4. feature가 연형이라면?

step 1. 각 feature에 대해 중요도로 정렬

step 2. label의 class가 변하는 지점을 찾기

step 3. 정렬된 feature를 기준으로 정렬

step 4. 각 feature에 대해 Gini index 혹은 Entropy 계산.

→ 과정 반복!!

Unit 5. 가지치기

왜 가지치기를 하는가? → Full tree 인 경우, 분기가 너무 많아 과적합 위험이 발생함. (분기가 너무 많으면 일반화 능력이 떨어진다.)

→ ∴ 이를 방지하기 위해 적절한 수준에서 terminal node를 한정한다.

→ 가지치기 종류 - 사전 가지치기 (pre pruning): Tree의 최대 깊이를 제한하여 leaf node의 최대 개수를 제한. 노드를 분할시키기 위해 필요한 최소한의 Data 개수를 강제해야 함.

사후 가지치기 (post pruning): 트리를 만든 후 하위 노드를 제거 or 병합. / 하위 노드의 불순도의 값이 특정 값 이하인 경우 하위 노드를 병합함.

→ 이 밖에 Cross-Validation 이나 Independent Validation Set 등의 방법을 사용하기도 한다.!