

## 의사결정 나무 (Decision Tree)

2) 데이터를 분류하는 방법은 선형적이지 않은 선형 공간을 분할하는 분류 문제.  
i) Classification ii) Regression

i) Classification ii) Regression 이 모두 사용 가능해서 CART (Classification And

이사 결정 사유를 이용한 분류

Regression Tree) 2가지의 이.

① 의사결정 사무를 임용한 부근에서 다음과 같다.

- ① 여러가지 독립 변수 중 하나의 독립 변수를 선택하고, 그 독립 변수에 대한 가절 (threshold)을 정한다.  $\Rightarrow$  이를, "분류기준"이라고 함.
  - ② 전체 학습 데이터 집합 (학습 노드)을 해당 독립 변수의 값이 가절보다 작은 데이터 그룹 (노드1)과 해당 독립 변수의 값이 가절보다 큰 데이터 그룹 (자식 노드)로 나눈다.
  - ③ 각각의 자식 노드에 대해 1~2 단계를 반복하여 하나의 자식 노드를 만든다.  
만, 자식 노드에 한 가지 클래스의 Data만 존재한다면, 더 이상 자식 노드를 나누지 않고 중지한다.
- $\Rightarrow$  이렇게, 자식 노드 나누기를 연속적으로 적용하면, 노드가 계속 증가하는 나무와 같은 형태로 표현할 수 있다.

이러한 결론을 바탕으로 사물론의 본질에 대해

- 의사 결정 나무의 전체 training data를 모두 적용해보면 각 Data는 특정한 Node로 가고 내려가게 된다. 각 노드는 그 노드를 선택한 데이터 샘플을 가진다.
- 이 때, 노드에 속한 data의 클래스의 비율을 구하여 이를 그 노드의 조건부 확률 바로  $P(Y=k|X)_{node}$  라고 정의한다.

$$P(F=k | X)_{node} \approx \frac{N_{node, k}}{N_{node}}$$

- Test Data ( $X_{test}$ )의 클래스를 예측할 때는, 가장 상위의 노드를 시작으로  
 왼쪽 규칙을 차례대로 적용하여, 마지막으로 도달하는 노드의 조건부 클래스를  
 이용하여 클래스를 예측한다.

$$\hat{y} = \arg \max_k P(r=k | X_{\text{test}}) \text{ last node}$$