

# 크로스 엔트로피 (cross entropy) $H[p, q]$

- 두 확률 분포  $p(y)$ ,  $q(y)$ 의 크로스 엔트로피  $H[p, q]$ 는 다음과 같이 정의한다.

$$H[p, q] = - \sum_{k=1}^K p(y_k) \log_2 q(y_k)$$

$$H[p, q] = - \int p(y) \log_2 q(y) dy$$

↳ "크로스 엔트로피"는 주로 분류 문제 (classification)의 "실측 분포"와 "예측 분포"를 비교하는데 사용된다.

↳ "크로스 엔트로피"는 언제 사용하나요?

⇒  $p$  하고  $q$ 가 얼마나 닮았는지 측정할 때 사용된다.

⇒ 일종의 Error Function 일수도.

⇒  $p$  하고  $q$ 가 비슷하면 크로스 엔트로피 값은 ↓  
 $p$  하고  $q$ 가 다르면 크로스 엔트로피 값은 ↑

# 클러크-라이블러 분산 (Kullback-Leibler divergence) ⇒ 상대 엔트로피

↳ 두 확률 분포  $p(y)$ ,  $q(y)$ 의 차이를 정량화하는 방법이다. (relative entropy)

$$KL(p || q) = H[p, q] - H[p] = \int p(y) \log_2 \left( \frac{p(y)}{q(y)} \right) dy$$

⇒ 값은 항상 양수이며, 두 확률 분포  $p(x)$ ,  $q(x)$ 가 완전히 같을 때만 0이 됨.

↳ 크로스 엔트로피, 엔트로피, 클러크-라이블러 분산 모두  $\log$  값을 계산하므로 계산량이 많아서 다른 방식을 찾아보니 "지니 불순도" (Gini Impurity) 사용.

# 지니 불순도 (Gini Impurity)

↳ 엔트로피처럼 확률 분포가 어느 쪽으로 치우쳐 있는가를 계는 척도이지만 로그를 사용하지 않으므로 계산량이 적어서 엔트로피 대신 사용됨.

$$G[Y] = \sum_{k=1}^K p(y_k) (1 - p(y_k))$$