

10.3 교차엔트로피와 쿨백-라이블러 발산

- 교차엔트로피 \rightarrow 분류기 성능 평가 지표
- 쿨백-라이블러 발산 \rightarrow 교차엔트로피를 응용한 것으로,
두 확률분포의 모양이 얼마나 유사한지 평가.

▣ 교차-엔트로피

두 확률분포 p, q 의 교차-엔트로피 (cross entropy) $H[p, q]$ 는

이산 확률분포의 경우,

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

연속 확률분포의 경우,

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

* 교차엔트로피는 지금까지 공부한 엔트로피, 결합엔트로피, 조건부엔트로피와 다르게
확률변수가 아닌 확률분포를 인수로 받는다는 점에 유의!

▣ 예제 ▣ 다음 두 분포의 교차엔트로피 = 0.25

$$p = [1/4, 1/4, 1/4, 1/4]$$

$$q = [1/2, 1/4, 1/8, 1/8]$$

$$= -\frac{1}{4} \cdot \log_2 \frac{1}{2} - \frac{1}{4} \log_2 \frac{1}{4} - \frac{1}{4} \cdot \log_2 \frac{1}{8} - \frac{1}{4} \cdot \log_2 \frac{1}{8}$$

$$= 1/4$$

교차엔트로피를 사용한 분류 성능 측정

"Y가 0 또는 1이라는 값만 갖는 이진분류문제"

- p 는 X 값이 정해졌을 때 정답인 Y 의 확률분포다.
- 이진분류문제에서 Y 는 0 또는 1이다.

따라서 p 는

- 정답이 $Y=1$ 일 때,

$$p(Y=0)=0, \quad p(Y=1)=1$$

- 정답이 $Y=0$ 일 때,

$$p(Y=0)=1, \quad p(Y=1)=0$$

분포 q 는 X 값이 정해졌을 때의 예측값의 확률분포다. 모수가 μ 인 베르누이분포로 가정한다

$$q(Y=0) = 1 - \mu$$

$$q(Y=1) = \mu$$

따라서 확률분포 p 와 q 의 교차엔트로피는,

$$\begin{aligned} H[p, q] &= -\overset{=0}{p(Y=0)} \log_2 \overset{=1}{q(Y=0)} - \overset{=1}{p(Y=1)} \log_2 \overset{=0}{q(Y=1)} \\ &= -\log_2 \mu \end{aligned}$$

$$\begin{aligned} \text{• 정답이 } Y=0 \text{ 일 때, } &\overset{=1}{p(Y=0)} \log_2 \overset{=0}{q(Y=0)} - \overset{=0}{p(Y=1)} \log_2 \overset{=1}{q(Y=1)} \\ H[p, q] &= -\overset{=1}{p(Y=0)} \log_2 \overset{=0}{q(Y=0)} - \overset{=0}{p(Y=1)} \log_2 \overset{=1}{q(Y=1)} \\ &= -\log_2 (1 - \mu) \end{aligned}$$

교차엔트로피는 분류성능이 좋을수록 작아지고, 분류성능이 나쁠수록 커진다.

이유는 다음과 같다.

→ $Y=1$ 일 때는 μ 가 작아질수록, 즉 예측이 틀릴수록 $-\log_2 \mu$ 의 값도 커진다.

→ $Y=0$ 일 때는 μ 가 커질수록, 즉 예측이 틀릴수록 $-\log_2 (1-\mu)$ 의 값도 커진다.

따라서, 교차엔트로피 값은 예측의 틀린 정도를 나타내는 오차함수의 역할을 할 수 있다.
 N 개의 학습 데이터 전체에 대해, 교차엔트로피 평균을 구하면 다음 식으로 표현할 수 있다. 이 값을 로그 손실 (log loss) 이라고 한다.

$$\text{(binary)} \quad \log \text{ loss} = -\frac{1}{N} \sum_{i=1}^N \left(\overset{\text{정답}(Y=1)}{y_i \log_2 \mu_i} + \overset{\text{정답}(Y=0)}{(1-y_i) \log_2 (1-\mu_i)} \right)$$

같은 방법으로, 이진분류가 아닌 다중분류에서도 교차엔트로피를 오차함수로 사용할 수 있다.
다중분류 문제의 교차엔트로피 손실함수를 카테고리 로그 손실 (categorical log-loss) 이라고 한다.

$$\text{categorical log loss} = -\frac{1}{N} \sum_{i=1}^N \sum_{k=1}^K \left(\underbrace{\mathbb{I}(y_i=k)}_{\text{지시함수}} \log_2 \underbrace{p(y_i=k)}_{\text{분류기가 계산한 } y_i=k \text{ 확률}} \right)$$

* 위 식에서 $\mathbb{I}(y_i=k)$ 는 $y_i=k$ 일 때만 1이고, 그렇지 않으면 0이 되는 지시함수 (indicator function) 다. $p(y_i=k)$ 는 분류모델이 계산한 $y_i=k$ 확률이다.

(사이킷런 (skit-learn) 패키지) metrics submodule은 로그 손실을 계산하는 log-loss 함수를 제공한다. `sklearn.metrics.log_loss`

쿨백-라이블러 발산

→ 변분법적 추론 (Variational inference)

↳ Deep learning → VAE (variational auto encoder)

쿨백-라이블러 발산 (Kullback-Leibler divergence)은 두 확률분포

$p(y)$, $q(y)$ 의 분포모양이 얼마나 다른지를 숫자로 계산한 값이다.

$KL(p||q)$ 로 표기한다.

?

GAN!

(이산확률분포)

$$KL(p||q) = H[p, q] - H[p]$$

$$= \sum_{i=1}^K p(y_i) \log_2 \left(\frac{p(y_i)}{q(y_i)} \right)$$

(연속확률분포)

$$KL(p||q) = H[p, q] - H[p]$$

$$= \int p(y) \log_2 \left(\frac{p(y)}{q(y)} \right) dy$$

쿨백-라이블러 발산 = 2차원트랜스에서 가중치 되는 p분포의 엔트로피값을 빼 값

= 상대엔트로피 (relative entropy)

(값이 항상 양수이며,
 $p(x)$, $q(x)$ 두 확률분포가 완전히 같을 경우에만 0이 된다.)

when $p(x) = q(x)$,

$$\begin{aligned} KL(p||p) &= H[p, p] - H[p] \\ &= \int p(y) \log_2 \left(\frac{p(y)}{p(y)} \right) dy \\ &= 0 \end{aligned}$$

$$KL(p||q) = 0 \iff p = q$$

콜백-라이블러 발산은 거리 (distance) 가 아니라, 확률분포 p 가 다른 확률분포 q 와 얼마나 다른지를 나타내는 값이다. 두 확률분포의 위치가 달라지면 일반적으로 값도 달라진다.

$$KL(p||q) \neq KL(q||p)$$

\therefore 콜백-라이블러 발산 = 상대 엔트로피.

예제

가변길이 인코딩과 콜백-라이블러 발산

4개의 문자 A, B, C, D로 쓰여진 다음과 같은 문서를
가변길이 인코딩하는 경우를 생각하자.

'DBCADDAAAAA...'

이 문서를 구성하는 문자의 확률분포는 다음과 같다.

$$p(Y=A) = \frac{1}{2}, \quad p(Y=B) = \frac{1}{4}, \quad p(Y=C) = \frac{1}{8}, \quad p(Y=D) = \frac{1}{8}$$

이 때, 한 문자당 인코딩된 문자수는 분포 q 의 엔트로피인 1.75가 된다.

$$\begin{aligned} \sum_{k=1}^K p(y_k) \log_2 p(y_k) &= -\frac{1}{2} \cdot \log_2 \frac{1}{2} + \left(-\frac{1}{4} \cdot \log_2 \frac{1}{4}\right) + \left(-\frac{1}{8} \cdot \log_2 \frac{1}{8}\right) \\ &\quad + \left(-\frac{1}{8} \cdot \log_2 \frac{1}{8}\right) = \boxed{1.75} \end{aligned}$$

그런데, 가변길이 인코딩을 사용하지 않고 고정길이 인코딩을 사용한다는 것은 다음과 같은 분포를 가정한 것과 같다.

$$q(Y=A) = \frac{1}{4}, \quad q(Y=B) = \frac{1}{4}, \quad q(Y=C) = \frac{1}{4}, \quad q(Y=D) = \frac{1}{4}$$

실제로 한 문자당 인코딩된 글자수는 다음과 같이 계산할 수 있다.

$$\begin{aligned} \sum_{i=1}^K p(y_i) \log_2 q(y_i) &= -\frac{1}{2} \cdot \log_2 \frac{1}{4} + \left(\frac{1}{4} \cdot \log_2 \frac{1}{4} \right) + \left(-\frac{1}{8} \cdot \log_2 \frac{1}{4} \right) \\ &\quad + \left(-\frac{1}{8} \cdot \log_2 \frac{1}{4} \right) \\ &= \frac{1}{2} \cdot 2 + \frac{1}{4} \cdot 2 + \frac{1}{8} \cdot 2 + \frac{1}{8} \cdot 2 \\ &= 2 \end{aligned}$$

클록-라이블러 법칙은 잘못된 분포 q 로 인코딩했을 때 한 문자당 인코딩된 글자수와 원래의 분포 p 를 사용하였을 때 한 문자당 인코딩된 글자수의 차이를 0.25라 한다.

$$\begin{aligned} KL(p||q) &= \sum_{i=1}^K p(y_i) \log_2 \left(\frac{p(y_i)}{q(y_i)} \right) \\ &= -\sum_{i=1}^K p(y_i) \log_2 q(y_i) - \left(-\sum_{i=1}^K p(y_i) \log_2 p(y_i) \right) \\ &= H[p, q] - H[p] \\ &= 2.0 - 1.75 = 0.25 \end{aligned}$$

즉, 확률분포 q 의 모양이 확률분포 p 와 다른 정도를 정량화한 값

= (클록-라이블러 법칙)

scipy.stats.entropy 함수 → 두 개의 확률분포를 인자로 넣으면
클록-라이블러 법칙을 계산 (base=2)