

0.1 엔트로피 ENTROPY

$Y=0$ 또는 $Y=1$ 인 두 가지 값은 확률변수의 확률분포가
다양과 같이 세 종류가 있다고 하자.

- 확률분포 $Y_1 : P(Y=0)=0.5, P(Y=1)=0.5$
- 확률분포 $Y_2 : P(Y=0)=0.8, P(Y=1)=0.2$
- 확률분포 $Y_3 : P(Y=0)=1.0, P(Y=1)=0.0$

(베이지안 관점에서 Y_1, Y_2, Y_3 의 의미)

- $Y_1 \Rightarrow y$ 값에 대해 아무것도 모른다 0.5/0.5
- $Y_2 \Rightarrow y$ 값이 0 이라고 믿지만 아닐 가능성도 있다는 것을 아는 상태 0.8/0.2
- $Y_3 \Rightarrow y$ 값이 0 이라고 100% 확신하는 상태 1.0/0.0

확률분포가 가지는 이러한 차이를 하나의 숫자 3로 나타낸 것

\Rightarrow 엔트로피

엔트로피의 정의

✓ 확률분포가 가지는 정보의 확실히 또는 정보량을 수치로 표현한 것.

- 확률분포에서 특정한 값에 나올 확률 상승 \rightarrow 엔트로피 감소.
- 확률분포에서 여러가지 값이 비슷한 확률로 나올 \rightarrow 엔트로피 증가.
- 확률 또는 확률분포가 특정한 값에 몰려있음 (편향함) \rightarrow 엔트로피 감소.
- 확률 또는 확률분포가 골고루 퍼져있음 (완전함) \rightarrow 엔트로피 증가.

* 물리학에서의 엔트로피

• 물질의 상태가 분산되는 정도 = 엔트로피.

↳ 고주 분산 \Rightarrow 엔트로피 \uparrow

특정 상태로 몰려있음 \Rightarrow 엔트로피 \downarrow

수학적으로, 엔트로피는 확률분포함수를 입력으로 받아 숫자를 출력하는
함수 (functional)로 정의한다. $H[\cdot]$ 기호로 표시한다.

binary

확률변수 Y 가 가변부분과 같은 이산확률변수이면, 다음처럼 정의한다.

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

가변부분 없음.

for $k = X$ 가 가질 수 있는 클래스의 수
 $p(y) =$ 확률질량함수 (pmf)

(* 확률의 로그값은 항상 음수이므로,
음수기호를 붙여서 양수로 만들어줌.)

확률변수 Y 가 정변분과 같은 연속확률변수이면 다음처럼 정의한다.

$$H[Y] = - \int_{-\infty}^{\infty} p(y) \log_2 p(y) dy$$

for $p(y) =$ 확률밀도함수 (pdf)

엔트로피 값에서 \log base = 2로 정의된 이유

\rightarrow 정보량과 관련이 있는 역사적 배경!

* 엔트로피 계산에서 $p(y)=0$ 인 경우 로그값이 정의되어있지 않음,
 다음과 같은 극한값을 사용한다.

$$\lim_{p \rightarrow 0} p \log_2 p = 0$$

(python 계산)

$$\left\{ \begin{array}{l} 0 \cdot \log_2(0) \\ \downarrow \\ \text{eps} * \log_2(\text{eps}) \end{array} \right\}$$

↳ ^{미적분} 값은 로피탈의 정리 (l'Hôpital's rule) 에서 구할 수 있다.

• Y_1, Y_2, Y_3 의 엔트로피 계산

$$H[Y_1] = -\frac{1}{2} \log_2 \frac{1}{2} - \frac{1}{2} \log_2 \frac{1}{2} = 1$$

$$H[Y_2] = -\frac{8}{10} \log_2 \frac{8}{10} - \frac{2}{10} \log_2 \frac{2}{10} \approx 0.72$$

$$H[Y_3] = \underbrace{-1 \cdot \log_2 1}_{=0} - \underbrace{0 \cdot \log_2 0}_{=0} = 0$$

극고주
퍼짐

특정한
값

엔트로피의 ~~특성~~ 성질

확률변수가 결정론적이면, 확률분포에서 특정항 하의 값이 나온 확률이 1이다. 이때 엔트로피는 0이 되고, 이 값은 엔트로피가 가질 수 있는 최솟값이다.

반대로, 엔트로피의 최대값은 이산확률변수의 클래스 개수에 따라 달라진다. 만약 이산확률분포가 가질 수 있는 값이 2^k 개면, 엔트로피의 최대값은 각 값에 대한 확률이 모두 같은 값인 $\frac{1}{2^k}$ 이다.

엔트로피의 값은

$$H = -2^k \times \left(\frac{1}{2^k} \cdot \log_2 \frac{1}{2^k} \right) = k$$

↓
2^k개만큼 더함.

이다.

#클래스 2³개
→ H=3
#클래스 2⁵개
→ H=5

엔트로피의 추정

이론적인 확률밀도함수가 없고, 실제 데이터가 주어진 경우,

- ① 데이터에서 확률밀도함수를 추정한다
- ② 추정 결과를 기반으로 엔트로피 계산

(ex) 데이터가 모두 00개 있고, 그 중 Y=0인 데이터가 40개, Y=1인 데이터가 40개 있을 경우 엔트로피는 1이다.

$$P(y=0) = \frac{4}{8} = \frac{1}{2}$$

$$P(y=1) = \frac{4}{8} = \frac{1}{2}$$

$$H[Y] = -\frac{1}{2} \log_2 \left(\frac{1}{2}\right) - \frac{1}{2} \log_2 \left(\frac{1}{2}\right) = \frac{1}{2} + \frac{1}{2} = 1.$$

stats 함수의 stats 함수와 - entropy 함수 제공
(base 인수로 2)

$$\text{Sp. stats. entropy}([0.5, 0.5], \text{base}=2)$$

지니불순도

엔트로피와 유사한 개념으로 지니불순도 (Gini impurity) 라는 것이 있다.

• 엔트로피와 유사하게, 불순도가 치우쳐 있는가를 재는 척도

• \log_2 사용하지 않으므로 계산량이 더 적어서

엔트로피 대신 사용됨.

* 경제학에서의 '지니계수'라는 것도.

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



엔트로피 최대화 ~ (생각적으로 맞고 있을 것!)

기댓값 0, 분산 σ^2 이 주어졌을 때 엔트로피 $H[p(x)]$ 를 가장 크게 만드는 확률밀도함수 $p(x)$ 는 정규분포가 된다. 이는 다음처럼 증명한다.

우선 확률밀도함수가 지켜야 할 (제한조건)은 다음과 같다.

(1) 확률밀도함수 총면적 = 1

$$\int_{-\infty}^{\infty} p(x) dx = 1$$

(2) 기댓값 = 0

$$\int_{-\infty}^{\infty} x p(x) dx = 0$$

(3) 분산 = σ^2

$$\int_{-\infty}^{\infty} x^2 p(x) dx = \sigma^2$$

엔트로피를 목적함수 (objective functional) \rightarrow entropy.

$$H[p(x)] = - \int_{-\infty}^{\infty} p(x) \log_2 p(x) dx$$

라그랑주 승수법으로 제한조건을 추가하면 다음과 같다.
 $+ \lambda_1, \lambda_2, \lambda_3$

$$\begin{aligned}
H[p(x)] &= - \int_{-\infty}^{\infty} p(x) \log p(x) dx \\
&\quad + \lambda_1 \left(\int_{-\infty}^{\infty} p(x) dx - 1 \right) \\
&\quad + \lambda_2 \left(\int_{-\infty}^{\infty} x p(x) dx - 0 \right) \\
&\quad + \lambda_3 \left(\int_{-\infty}^{\infty} x^2 p(x) dx - \sigma^2 \right) \\
&= \int_{-\infty}^{\infty} (-p(x) \log p(x) + \lambda_1 p(x) + \lambda_2 x p(x) + \lambda_3 x^2 p(x) \\
&\quad - \lambda_1 - \lambda_3 \sigma^2) dx
\end{aligned}$$

변분법에서 도함수를 다음과 같이 계산하여,

$$\frac{\delta H}{\delta p(x)} = -\log p(x) - 1 + \lambda_1 + \lambda_2 x + \lambda_3 x^2 = 0$$

따라서 확률밀도함수의 형태는,

$$p(x) = \exp(-1 + \lambda_1 + \lambda_2 x + \lambda_3 x^2)$$

적분을 통해 규 형태의 확률밀도함수의 변칙, 기댓값, 분산을 계산하고
주어진 제약을 만족하도록 변칙계수를 풀면, 각각의 상수를 다음과 같이 구할 수 있다.
(풀이과정은 생략)

$$\lambda_1 = 1 - \frac{1}{2} \log 2\pi \sigma^2$$

$$\lambda_2 = 0$$

$$\lambda_3 = -\frac{1}{2\sigma^2}$$

위 값을 대입하면, 정규분포

$$\phi(x) = \frac{1}{\sqrt{2\pi}\sigma^2} \exp\left(-\frac{x^2}{2\sigma^2}\right)$$

"정규분포 \rightarrow Entropy $\uparrow\uparrow$ "

따라서, 정규분포는

- 가장 값과 표준편차를 안고있는 확률분포중에서 가장 엔트로피가 크고
따라서 가장 정보와 적은 확률분포

☆ 정규분포는 베이지 추정에 있어서

"사실상의 무정보" 사전 확률분포로 자주 사용한다.

정규분포 \sim 자연계 현상 \sim 자연계는 엔트로피가 높아져야 하는 경향

가변길이 인코딩

엔트로피는 현재 통신분포에서 데이터가 차지하는 정보량을 계산하기 위해 고안되었다.

예를 들어, 특히 문자 A, B, C, D로 쓰여진 다음 문자가 있다고 하자.

'BDA BABACBABAACBBAAACBBCAB...' (200자)

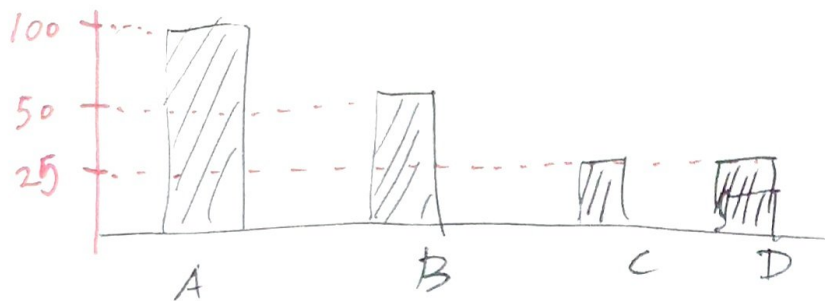
이 문자를 0과 1로 이루어진 이진수로 변환할 때, 다음과 같이 인코딩한다.

• A = "00"	• C = "10"
• B = "01"	• D = "11"

\rightarrow (400자)

① 문자수를 더 줄일 수 있는가?

→ 문자 수의 분포 확인.



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

(* 지프의 법칙 (Zipf's law)에 따르면,
이러한 분포는 현실의 글자 빈도수에서 흔히 나타난다.



특정 분포와 차가 없을 때, 다음처럼 인코딩하면 인코딩된 후의 이진수 수를 줄일 수 있다.

more frequent ↑

- A = "0"
- B = "10"
- C = "110"
- D = "111"

문자마다 인코딩하는 이진수 숫자가 다르다
= 가변길이 인코딩 (variable length encoding)

가장 많이 출현하는 'A' : 1글자
'C', 'D'는 3글자, but 빈도수 적어서 영향 ↓

문자의 분포 = 위에서 정정한 분포일 때,
인코딩된 이진수의 숫자는 350개가 된다.

$$(200 \times \frac{1}{2}) \cdot 1 + (200 \times \frac{1}{4}) \cdot 2 + (200 \times \frac{1}{8}) \cdot 3 + (200 \times \frac{1}{8}) \cdot 3 = 350$$

따라서, 알파벳 한 글자를 인코딩하는데 필요한 평균 비트수는 $350/200 = 1.75$
이 값은 확률변수의 엔트로피 값과 같다.

$$H = -\frac{1}{2} \log_2 \frac{1}{2} - \frac{1}{4} \log_2 \frac{1}{4} - \frac{2}{8} \log_2 \frac{1}{8} = 1.75$$

• 연습문제 10.1.4

A, B, C, D, E, F, G, H의 8글자로 이루어진 문서가 있고,

각각의 글자가 나올 확률이 다음과 같다고 가정하자.

$$\left\{ \frac{1}{2}, \frac{1}{4}, \frac{1}{8}, \frac{1}{16}, \frac{1}{64}, \frac{1}{64}, \frac{1}{64}, \frac{1}{64} \right\}$$

① 이 문서를 위한 가변길이 인코딩 방식을 서술하고,

② 한 글자를 인코딩하는데 필요한 평균 비트수를 계산하라.

$$2^3 = 8$$

NON-가변길이 인코딩	가변길이 인코딩	
000	0	$-\frac{1}{2}$
001	10	$-\frac{1}{4}$
010	110	$-\frac{1}{8}$
011	1110	$-\frac{1}{16}$
100	11100	} $\frac{1}{64}$
101	11101	
110	11110	
111	11111	

(avg. 3 bit)

② 한 글자를 인코딩하는데 필요한 평균 비트수 = 엔트로피

$$\begin{aligned}
 H &= -\frac{1}{2} \log_2 \frac{1}{2} - \frac{1}{4} \log_2 \frac{1}{4} - \frac{1}{8} \log_2 \frac{1}{8} - \frac{1}{16} \log_2 \frac{1}{16} \\
 &\quad - \frac{1}{64} \log_2 \frac{1}{64} - \frac{1}{64} \log_2 \frac{1}{64} - \frac{1}{64} \log_2 \frac{1}{64} - \frac{1}{64} \log_2 \frac{1}{64} \\
 &= \boxed{2 \text{ bit}}
 \end{aligned}$$