

2022 - 02 - 09

Likelihood: Observation: patient has yellow eyes

Explanation A: Patient has hepatitis B (HB) 乙肝

Explanation B: Patient has COVID

Explanation A is more reasonable because

$$P(\text{Yellow eyes} | \text{HB}) \gg P(\text{Yellow eyes} | \text{COVID})$$

observation E , Explanation H

$P(E|H)$ is a measure of

↑ ↑ how well H explains E .

Fixed chance

$L(H|E) = P(E|H)$ likelihood

↗ of H

function of H

Data : D (observations)

Model : m (explanation)

$$L(m | D) = P(D|m)$$

likelihood of model given Data

measures how well model explains data

fits

prior	$\frac{P(m D)}{P(D m)}$	$L(m D) = \frac{P(D m)}{P(m)}$ likelihood
-------	-------------------------	--

$$P(m|D) = \frac{P(m, D)}{P(D)} = \frac{P(D|m) P(m)}{P(D)} \rightarrow \text{Evidence partition function}$$

↓ ↓
 $\propto P(m) P(D|m)$

posterior \propto prior \times likelihood

Bayes' Theorem

parameter estimation

$$\theta = P(X=H)$$

↑

result of

thumb tack toss

θ :

A: 0

B: 0.5

C: 0.6 ✓ Explains / fits the data
 the best

D: 0.7 maximizes $L(\theta|D) = P(D|\theta)$

maximum likelihood Estimation (MLE)

$$D = \{H, T, H, H, T\} \quad \theta = P(X=H)$$

d_1, d_2, d_3, d_4, d_5

$$L(\theta | D) = P(D | \theta)$$

$$= P(d_1, d_2, \dots, d_5 | \theta)$$

Independence
assumption

$$= P(d_1 | \theta) P(d_2 | \theta) P(d_3 | \theta) P(d_4 | \theta) P(d_5 | \theta)$$

$$= \theta^3 (1-\theta)^2 \quad \theta \quad (1-\theta)$$

$$= \theta^3 (1-\theta)^2 \quad \text{identically distributed}$$

$\downarrow m_h \quad \downarrow m_t$

$$D' = \{H, H, H, T, T\} \rightarrow \text{(i.i.d)}$$

$$L(\theta | D') = \theta^3 (1-\theta)^2 \quad \text{sufficient statistics}$$

$$\max_{\theta} L(\theta | D) = \theta^3(1-\theta)^2$$

$$l(\theta | D) = \log L(\theta | D)$$

$$\uparrow = 3 \log \theta + 2 \log (1-\theta)$$

log likelihood : math simplicity

$$\frac{d l}{d \theta} = 0 \Rightarrow \theta = \frac{3}{3+2} = 0.6$$

$$\text{MLE: } \theta = \frac{m_h}{m_h + m_t}$$

Bayesian Estimation

$$p(\theta | D) \propto p(\theta) p(D | \theta)$$

$$\propto \theta^{\alpha_h - 1} (1 - \theta)^{\alpha_t - 1} \theta^{m_h} (1 - \theta)^{m_t}$$

$$= \theta^{\alpha_h + m_h - 1} (1 - \theta)^{\alpha_t + m_t - 1}$$

$$\Rightarrow p(d_{m+1} = H | D) = \frac{\alpha_h + m_h}{\alpha_h + m_h + \alpha_t + m_t}$$

α_h, α_t : Virtual Counts

$$\underline{\alpha_h = \alpha_t = 100}$$

Beta Distribution

200 : equivalent
Sample

$$P(d_{m+1} = H | D) = \frac{\alpha_h + m_h}{\alpha_h + m_h + \alpha_t + m_t}$$

thumbtack: $\alpha_h = \alpha_t = 0$

$$m_h = 3, \quad m_t = 7: \quad P(d_{m+1} = H | D) = \frac{3}{10} = 0.3$$

coin: $\alpha_h = \alpha_t = 1,000$

$$\textcircled{1} \quad M_h = 3, \quad m_t = 7 \quad P(d_{m+1} = H | D) = \frac{1,000 + 3}{1,000 + 3 + 1,000 + 7}$$

$$\textcircled{2} \quad M_h = 30,000, \quad M_t = 70,000 \quad P(d_{m+1} = H | D) = \frac{1,000 + 30,000}{1,000 + 30,000 + 1,000 + 70,000} \approx 0.3$$

f is concave

$$x_1, x_2, \dots, x_n$$

$$p_i \geq 0, \sum_i p_i = 1$$

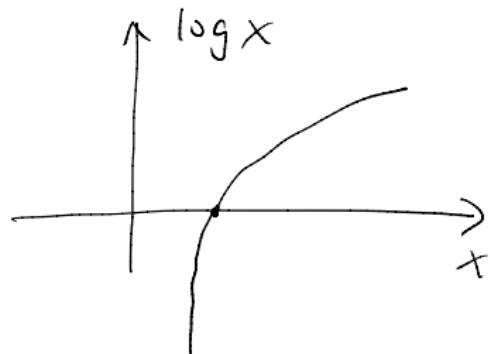
$$p_1, p_2, \dots, p_n$$

$$\underbrace{\sum_{i=1}^n p_i f(x_i)} \leq f\left(\sum_{i=1}^n p_i x_i\right) \quad \begin{matrix} \text{Jensen's} \\ \text{inequality} \end{matrix}$$

$$E_p [f(x)] \leq f(E_p [x])$$

$\log x$ is concave

$$E_p[\log x] \leq \log E_p[x]$$



Example 2: Thumb tack $\theta = P(X=h)$

Guess X : result of Thumbtack toss

$$\textcircled{1} \quad \theta = 0.2$$

$$\textcircled{2} \quad \underline{\theta = 0.5}$$

$$\textcircled{3} \quad \theta = 0.8$$

$$\begin{aligned} H(x) &= P(X=h) \log \frac{1}{P(X \neq h)} + P(X=t) \log \frac{1}{P(X \neq t)} \\ &= \theta \log \frac{1}{\theta} + (1-\theta) \log \frac{1}{1-\theta} \\ &= -\theta \log \theta - (1-\theta) \log (1-\theta) \end{aligned}$$

$$H(x) = \sum_x p(x) \log \frac{1}{p(x)} \quad x \in \{1, 2, \dots, 6\}$$

$$\textcircled{1} \quad H(x) \geq 0$$

$$\leq \log \sum_x p(x) \frac{1}{p(x)}$$

$$\textcircled{2} \quad H(x) \leq \log |x|$$

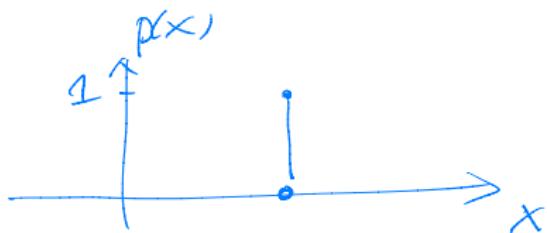
$$= \log |x|$$

$$H(x) = \log |x| \text{ when } p(x)$$

is uniform

↗ cardinality
of x

$$\textcircled{3} \quad H(x) = 0 \quad \text{when } p(x=a) = 1 \text{ for some } a$$



Conditional Entropy

$$H(Y) = \sum_Y p(y) \log \frac{1}{p(y)}$$

Y : initial dist $p(Y)$, $\underline{H(Y)}$

Evidence $X=a$: $p(Y|X=a)$

$$\underline{H(Y|X=a)} = \sum_Y p(Y|X=a) \log \frac{1}{p(y|X=a)}$$

How uncertainty about Y will change
if we test on X ?

$$\begin{aligned} \underline{H(Y|X)} &= \sum_a H(Y|X=a) p(X=a) \\ &= \sum_a \sum_Y p(Y|X=a) \log \frac{1}{p(y|X=a)} \underline{p(X=a)} \\ &= \sum_{X,Y} \underline{\underline{p(X,Y)}} \log \frac{1}{\underline{\underline{p(Y|X)}}} = -E[\log p(y|x)] \end{aligned}$$

KL Divergence

$$KL(P \parallel Q) = \sum_x P(x) \log \frac{P(x)}{Q(x)}$$

$$\textcircled{1} \quad KL(P \parallel Q) \geq 0 = - \sum_x P(x) \underbrace{\log}_{\frac{Q(x)}{P(x)}}$$

$$\textcircled{2} \quad KL(P \parallel Q) \geq 0 \\ \text{when } P=Q \quad \geq - \log \sum_x P(x) \frac{Q(x)}{P(x)}$$

$$\textcircled{3} \quad KL(P \parallel Q) = - \log \sum_x Q(x)$$

$$\neq KL(Q \parallel P) = - \log 1$$

$$= 0$$

$$\begin{aligned}
 0 &\leq KL(P \parallel Q) = \sum_x P(x) \log \frac{P(x)}{Q(x)} \\
 &= \sum_x P(x) \log P(x) + \sum_x P(x) \log \frac{1}{Q(x)} \\
 &= - \underbrace{\sum_x P(x) \log \frac{1}{P(x)}}_{\text{cross entropy}} + \sum_x P(x) \log \frac{1}{Q(x)} \\
 &= -H(P) + H(P, Q)
 \end{aligned}$$

$$\min KL(P \parallel Q) \equiv \min H(P, Q)$$

$$H(P, Q) \geq H(P)$$

$$-E_P[\log Q(x)] \geq -E_P[\log P(x)]$$

$$E_P[\log \underline{Q(x)}] \leq E_P[\log P(x)]$$

Gibbs Inequality

Mutual Information

0.9

0.2

Info Y has
about X

$$I(X;Y) = H(X) - H(X|Y)$$

↑

Uncertainty
about X

H(X|Y)

↑

Expected Uncertainty
about X if test Y

$$= \underline{KL(P(x,y) || P(x)P(y))}$$

≥ 0

$I(Y|X)$

\Rightarrow if $P(x,y)$

$$= P(x)P(y)$$

① $I(X;Y) = I(Y;X)$

② $I(X;Y) \geq 0$

③ $I(X;Y) = 0$ when $X \& Y$ independent

④ $H(X) \geq H(X|Y)$

$$E_p [\log \underline{Q}(x)] \leq E_p [\log P(x)] \quad \text{Gibbs Inequality}$$

Thimble tack: Data m_h heads, m_t tails

$$I(\theta | D) = m_h \log \theta + m_t \log (1-\theta)$$

$$\frac{dI}{d\theta} = 0 \Rightarrow \theta^* = \frac{m_h}{m_h + m_t} \quad m = m_h + m_t$$

$$\max_Q m_h \log Q(x=h) + m_t \log Q(x=t) \quad \theta \quad Q(x=h)$$

$$\Leftrightarrow \max_Q \frac{m_h}{m} \log Q(x=h) + \frac{m_t}{m} \log Q(x=t)$$

$$\tilde{P}(x=h) \quad \tilde{P}(x=t)$$

$$\Leftrightarrow Q = \tilde{P}, \quad Q(x=h) = \frac{m_h}{m}$$

empirical
distribution

