

What is "relative" entropy?

- Definition: given p, q , p and q both supported on X
and $\text{Supp}(q) = \{x \in X \mid q(x) > 0\} \subseteq \text{Supp}(p) = \{x \in X \mid p(x) > 0\}$
the relative entropy:

$$D(p \parallel q) = \mathbb{E}_{x \sim p} \left[\log \left(\frac{p(x)}{q(x)} \right) \right].$$

- also called: "Kullback-Leibler Divergence (KL)"

- Motivating questions:

- how can we interpret $D(p \parallel q)$? What does it mean?
(separate from a convenient object arising in analysis)
- in particular, why is the asymmetry $D(p \parallel q) \neq D(q \parallel p)$
a feature of the definition rather than a bug?

- What we already Know:

1. Properties:

- Nonnegativity: $D(p \parallel q) \geq 0$ and $= 0$ if and only if $p = q$ (almost surely for $x \sim p$)

• loosely, $D(p \parallel q)$ is a "distance" between p & q
we'll qualify this

- Asymmetry: $D(p \parallel q) \neq D(q \parallel p)$

- apparent from definition. p and q are not exchangeable
... average over p not q
- true distances (metrics) should be symmetric
- not a distance, if we treat as a pseudo-metric
Then this is just a bug
- why is this substantive / useful?

- Joint Convexity (see TBC): if $\lambda \in [0, 1]$ then

$$D(\lambda p + (1-\lambda)p' \parallel \lambda q + (1-\lambda)q') \leq \lambda D(p \parallel q) + (1-\lambda) D(p' \parallel q').$$

2. Example uses (so far):

(i) excess description length (cost of misspecification in compression):

$$(a) \text{ given code } C, \ell(x) = |c(x)|, L(C) = \mathbb{E}_{x \sim p} [\ell(x)]$$

- if we can only send n characters (d -ary) per unit time, then, using C , the channel capacity (max-average X sent per unit time) is $\approx 1/L(C)$

- wanted to minimize $L(C)$ (maximizes capacity for noiseless channel)

inefficiency from mismatch in code design and distribution of X

- saw $L(C) = \underbrace{H[X]}_{\text{Theoretical lower bound}} + D(p \parallel q) - \log(z)$

where $q(x) = \frac{1}{z} D^{-\ell(x)}, z = \sum_x D^{-\ell(x)}$.

(b) if we designed the code C using the assumption $X \sim q$ but $X \sim p$ then (ex: Shannon code)

$$L(C) \in H[X] + \underbrace{D(p \parallel q)}_{\text{excess description length}} + [0, 1]$$

- interpretation: $D(p \parallel q)$ is the cost of misspecification
 - true for other tasks
 - asymmetric since p is the true dist, generates samples
 - q is the approximation/proposal, does not generate samples
- average over p
not q

(ii) differences in entropies ("relative"):

- Ex: $H[X] - \log(1/X) = -D(p \parallel \text{Uniform}[X])$

* recall from limiting construction of differential entropy

$$\begin{aligned} (H[X] - \log(1/X)) &= H[X] + \log(1/X^{-1}) = \mathbb{E}_{x \sim p} [-\log(p(x)) + \log(1/X^{-1})] = -\mathbb{E}_{x \sim p} [\log(\frac{p(x)}{1/X^{-1}})] \\ &= -D(p \parallel \text{Uniform}[X]) \end{aligned}$$

$$\text{so: } H[X] = H[\text{Uniform}(X)] - D(p \parallel \text{Uniform}(x))$$

difference between max possible entropy (uniform) and entropy
of p ... information needed to go from $\text{Uniform}(x)$ to p

- this algebra does not work for generic q ($D(p \parallel q) \neq H(q) - H(p)$ generically)
- works when misspecified...

- Suppose $X \sim p$ but we think $X \sim q$
- measure our surprise on sampling X as $-\log(q(x))$
- our expected surprise is:

$$- \mathbb{E}_{x \sim p} [\log(q(x))] \neq H(q) \text{ or } H(p)$$

if we have many samples could estimate

let $\tilde{H}(p \parallel q)$ be our estimate to $H[X]$

using q for surprise

"cross entropy", usually $H(q, p)$

$$\cdot \text{Compare: } \tilde{H}(p \parallel q) = - \mathbb{E}_{X \sim p} [\log(q(x))] \text{ w/ } H[X] = H(p) \dots$$

$$\cdot \text{IE error in estimate: } \tilde{H}(p \parallel q) - H(p) = \mathbb{E}_X [-\log(q(x)) + \log(p(x))]$$

$$= D(p \parallel q) \geq 0$$

- so, if $X \sim p$, think $X \sim q$, and estimate $H[X]$ using samples
 - then the expected error in our estimate is $D(p \parallel q)$.
 - the error is > 0 if $p \neq q$. Always too/overly surprised (overestimate $H[X]$)
 - $D(p \parallel q)$ is our excess surprise (when misspecified)

$$3. \text{ Relation to Mutual Info: } I[X; Y] = D(p_{x,y} \parallel p_x p_y)$$

so: (a) mutual info is the expected excess description length of joint draws X, Y if we ignored their dependence

} contextualize
 I via D
not D ...

(b) mutual info is the excess surprise when observing joint samples X, Y if we thought they were independent

(c) we will see... if X is an unknown w/ prior p , observe Y , posterior $X|Y=y$, then $D(p_x \parallel p_{x|Y=y})$ = the information gained about X via observation.