

Relative Entropy is a Divergence (not a distance)

- roughly, a divergence $p \parallel q$ measures how distinguishable q is from p when we draw data from p

- If $X \sim p$, but I thought $X \sim q$, how hard would it be to tell?
 - How much data would I need to find out?

... easiest to motivate divergences via hypothesis testing ...

- Problem: We observe a sequence of samples $\{X_i\}_{i=1}^n \sim_{\text{iid}} p$ where $p = p_0$ or p_1 , but we don't know which...

- Hypotheses: $H_0: p = p_0$, $H_1: p = p_1$

- Test:
 1. pick a test statistic $s(x_1, x_2, \dots, x_n): X^n \rightarrow \mathbb{R}$
 2. estimate the sampling distribution of $s(X_1, X_2, \dots, X_n)$ under $X \sim p_0$ and $X \sim p_1$
 3. compare observed s to its sampling distribution under p_0 and p_1
 4. define accept (H_1) and reject (H_0) regions
 - using 2. (e.g. a threshold on s for a one-sided test)
 5. reject/accept based on where the observed s is

(e.g. accept if $\Pr(\text{sampling } s(X_1, X_2, \dots, X_n) \text{ less likely than observed } s(x_1, x_2, \dots, x_n))$ under $H_0 \leq \alpha$ for some $\alpha \Rightarrow \alpha$ controls the FPR...)

i.e. chance we pick H_1 when H_0 is true

i.e. controls the significance of the test)

could exchange the
roles of H_0 and
 H_1 here, convention
is arbitrary...

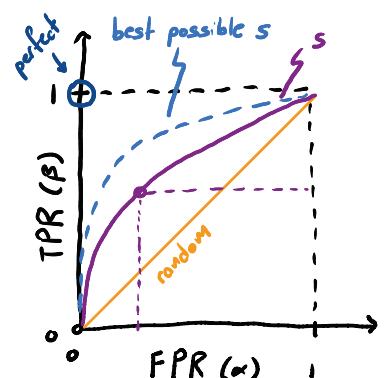
- Choosing s and the accept/reject region (α) determines the test
(fixes decision thresholds on s)

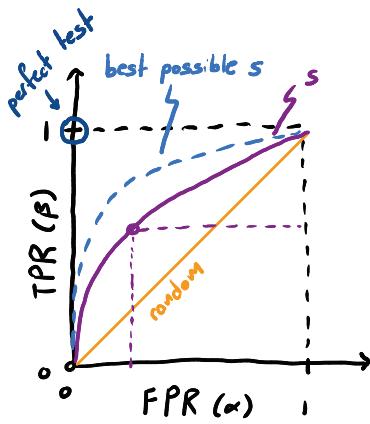
• thus, also fixes the FNR, i.e. prob. of choosing H_0 when H_1 is true
i.e. controls power/sensitivity

- generally want power & significance (sensitivity & specificity)

- for a given s , visualize trade-off w/ ROC

• plot best TPR (power) possible for a given FPR (significance)
equivalently, best TNR for a given FNR





- Want to choose our statistic s to:

(\cdot) maximize power (sensitivity) for all significances (specificity)
when using H_0 to control (fix accept/reject to control
 $FPR = \text{control significance}$

($\cdot\cdot$) maximize significance (specificity) for all powers (sensitivity)
when using H_1 to control (fix accept/reject to control
 $FNR = \text{control sensitivity}$

- What s should we choose?

Neyman-Pearson: if $f: \mathbb{R}^+ \rightarrow \mathbb{R}$ is injective (one-to-one) and monotonic
and we define:

$$\text{"likelihood ratio"} \Rightarrow LR(x_1, x_2, \dots, x_n) = \frac{\Pr(X_1=x_1, X_2=x_2, \dots, X_n=x_n | H_1)}{\Pr(X_1=x_1, X_2=x_2, \dots, X_n=x_n | H_0)}$$

then $s(x_1, x_2, \dots, x_n) = f(LR(x_1, x_2, \dots, x_n))$

will achieve (\cdot) and ($\cdot\cdot$).

if large, data more likely under H_1 , choose H_1
if small, ... under H_0 , choose H_0

idea: the best way to choose H_0 vs. H_1 is by comparing
the likelihood of the data under each model.

Ex: (i) LR test (ii) log(LR) test (iii) $H_0: X \sim \exp(\lambda_0)$, $H_1: X \sim \exp(\lambda_1)$ can use $s = \frac{1}{n} \sum_{j=1}^n X_j$.

to understand how distinguishable p_0 and p_1 are, we should study the
sampling dist. of $s(x_1, \dots, x_n)$

at least for large n ...

Q: What is $E[s(x_1, \dots, x_n)]$? What is the sampling dist. for large n ?

Single Sample: $E[s(x)] = \begin{cases} E_{x \sim p_0}[f(LR(x))] & \text{if } H_0 \text{ true} \\ E_{x \sim p_1}[f(LR(x))] & \text{if } H_1 \text{ true} \end{cases}$

sign flips if use p_0/p_1 for LR

in particular: $\log(LR)$:

$$E_{x \sim p_0}[\log\left(\frac{p_1(x)}{p_0(x)}\right)] = -D(p_0 || p_1) \quad \text{if } H_0 \text{ true}$$

$E[s(x)] = \begin{cases} E_{x \sim p_0}\left[\log\left(\frac{p_1(x)}{p_0(x)}\right)\right] = -D(p_0 || p_1) & \text{if } H_0 \text{ true} \\ E_{x \sim p_1}\left[\log\left(\frac{p_1(x)}{p_0(x)}\right)\right] = +D(p_1 || p_0) & \text{if } H_1 \text{ true} \end{cases}$

X order depends which is true.

• interpretation: $D(p \parallel q)$ = expected value of log likelihood ratio (best test stat for distinguishing $p \neq q$) on a single draw, when p is true

• multi-sample: $\mathbb{E}[s(x_1, \dots, x_n)] = \mathbb{E}_{X \sim P_0 \text{ or } P_1} \left[f\left(\frac{P_1(x_1, x_2, \dots, x_n | H_1)}{P_0(x_1, x_2, \dots, x_n | H_0)}\right) \right] = \mathbb{E}_{X \sim P_0 \text{ or } P_1} \left[f\left(\prod_{j=1}^n \frac{P_1(x_j | H_1)}{P_0(x_j | H_0)}\right) \right]$

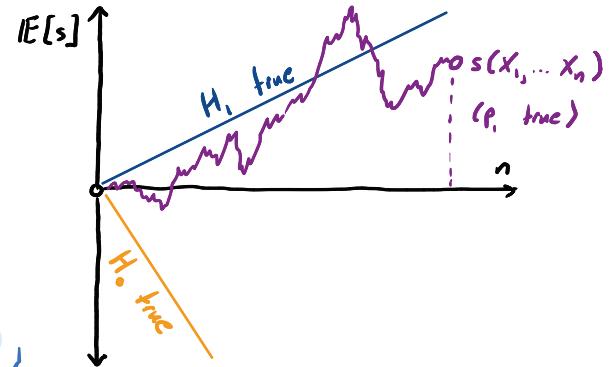
• ideally, choose f s.t. the terms here separate ... log is natural (separate product over each sample)

• log(LR): $\mathbb{E}[s(x_1, \dots, x_n)] = \mathbb{E}_{X \sim P_0 \text{ or } P_1} \left[\log\left(\prod_{j=1}^n LR(x_j)\right) \right] = \sum_{j=1}^n \mathbb{E}_{X \sim P_0 \text{ or } P_1} [\log(LR(x_j))]$

$$= \begin{cases} -n D(p_0 \parallel p_1) & \text{if } H_0 \\ n D(p_1 \parallel p_0) & \text{if } H_1 \end{cases}$$

• the larger $D(p_1 \parallel p_0)$ the faster we expect to acquire evidence distinguishing p_1 and p_0 when p_1 is true

• ... $D(p_0 \parallel p_1)$... when p_0 is true



$D(p \parallel q)$ is the rate we gain evidence to distinguish $p \neq q$ when p is true

• If we use $f(t) = \frac{1}{n} \log(t) = \log(t^{1/n})$

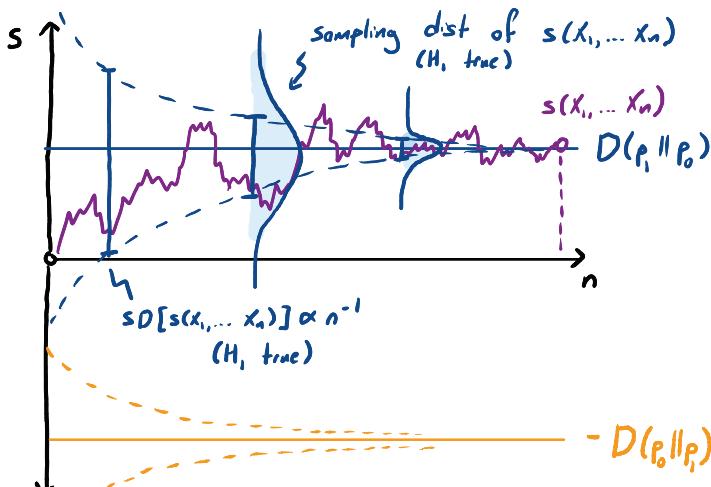
then:

$$\mathbb{E}[s(x_1, \dots, x_n)] = \begin{cases} -D(p_0 \parallel p_1) & \text{if } H_0 \text{ true} \\ +D(p_1 \parallel p_0) & \text{if } H_1 \text{ true} \end{cases}$$

i.i.d. draws

$$\text{and, } s(x_1, \dots, x_n) = \frac{1}{n} \sum_{j=1}^n \log(LR(x_j))$$

so, by the CLT, as $n \rightarrow \infty$, the sampling distribution of $s(x_1, \dots, x_n)$ converges to a N distribution w/ $\mathbb{E} = D(\cdot \parallel \cdot)$ and variance = $\Theta(n^{-1})$.

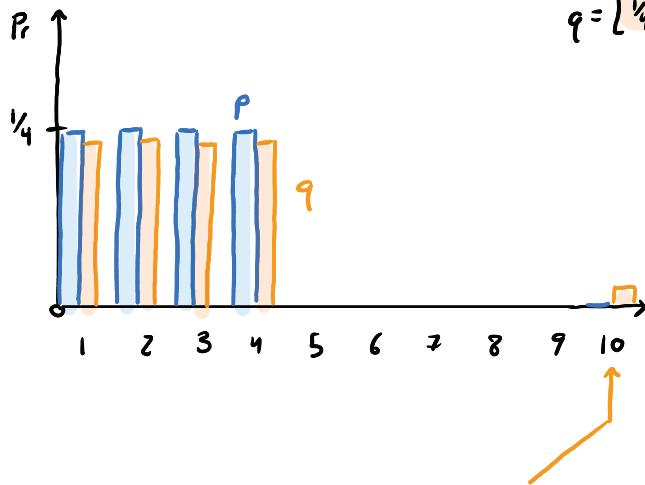


• asymptotically, $D(p \parallel q)$ and $D(q \parallel p)$ measure how distinguishable p is from q when drawing from p , or when drawing from q

- Interpretation: $D(p \parallel q) = \mathbb{E}$ of test statistic ($s(x_1, \dots, x_n) = \log(LR(x_1, \dots, x_n)^{1/n})$) for the $\log(LR)$ test for distinguishing $p \neq q$ when p is true.
- loosely, "distinguishability given $\text{data} \sim p$ "

- Changing f in $s(x) = f(LR(x))$ induces a family of "divergences" generated by $f\dots$ (see next section)
- This construction/interpretation explains the asymmetry of $D\dots$

• Ex: $X = \{1, 2, \dots, 10\}$ $p = [\frac{1}{4}, \frac{1}{4}, \frac{1}{4}, \frac{1}{4}, 0, 0, 0, 0, 0, 0]$ $q = [\frac{1}{4} \cdot (1-\varepsilon) \text{ (x4)}, 0, 0, 0, 0, 0, \varepsilon]$ } only differ by a rare event that is possible under q



• if p is true:

$$D(p \parallel q) = \mathbb{E}_{X \sim p} \left[\log \left(\frac{p(x)}{q(x)} \right) \right] = \sum_{x=1}^4 \log \left(\frac{1/4}{1/4(1-\varepsilon)} \right) \cdot \frac{1}{4}$$

$$= -\log(1-\varepsilon) \approx \varepsilon \xrightarrow{\varepsilon \gg 0} 0$$

so, when p is true, it's hard to distinguish q from p

rare event: $X=10$, if $X=10$, we know w/ 100% certainty $X \sim q$

• if q is true,

$$D(q \parallel p) = \mathbb{E}_{X \sim q} \left[\log \left(\frac{q(x)}{p(x)} \right) \right] = \sum_{x=1}^4 \log(1/(1-\varepsilon)) \cdot \frac{1}{4}(1-\varepsilon)$$

$$+ \varepsilon \log(\varepsilon)$$

$\log(\infty) = +\infty$

= +∞

we can build examples where

$$D(p \parallel q) \rightarrow 0 \text{ but } D(q \parallel p) \rightarrow \infty !!!$$

very asymmetric.

so, when q is true $p \neq q$ are easy to distinguish...

- Why? if q is true, we could sample $X=10$ (prob. ε) if we sample $X=10$ we know w/ complete certainty that $X \sim q$ not p since prob $X=10$ if $X \sim p$ is zero!