

General Properties for f -divergence:

1- $D_f(p||q) \geq 0$

2- IF f is strictly convex at 1, then $D_f(p||q) = 0 \Rightarrow p=q$.

\uparrow
* IF $\forall x \neq y$ & $\lambda \in (0,1)$ such that $\lambda x + \bar{\lambda} y = 1$,

then: $f(1) < \lambda f(x) + \bar{\lambda} f(y)$.

or equivalently

* Note that $p=q \Rightarrow D_f(p||q) = 0$ if $\sum_{i=1}^n \lambda_i \cdot \bar{\lambda}_i = 1$

proof: $D_f(p||q) = \sum q(x_i) \cdot f\left(\frac{p(x_i)}{q(x_i)}\right) \geq f(1)$ (*)

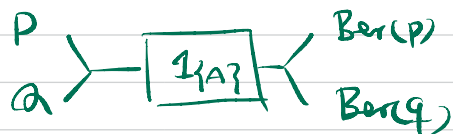
IF f is strictly convex at 1 $\Rightarrow f(1) < \sum q(x_i) \cdot f\left(\frac{p(x_i)}{q(x_i)}\right)$
(because $\sum q(x_i) \cdot \frac{p(x_i)}{q(x_i)} = 1$)

Thus, equality in (*) happens only if $\frac{p(x_i)}{q(x_i)} = c \quad \forall x$.

A more direct proof:

Suppose, for the sake of a contradiction, that $D_f(P||Q)=0$ for some $P \neq Q$. Then there exists some set A such that $P(A) \neq Q(A)$.

Let $p := P(A)$ & $q := Q(A)$.



By Df5 (which is about to be proved in Item 5), we have

$$D_f(\text{Ber}(p) || \text{Ber}(q)) \leq D_f(P||Q) = 0$$

$$\Rightarrow D_f(\text{Ber}(p) || \text{Ber}(q)) = 0$$

$$D_f(\text{Ber}(p) || \text{Ber}(q)) = q f(p/q) + \bar{q} f(\bar{p}/\bar{q}) = f(1)$$

↑
contradicts the
definition of
strict convexity.

3- $(p, q) \mapsto D_F(p \| q)$ is jointly convex

$$D_F(\lambda p_1 + \bar{\lambda} p_2 \| \lambda q_1 + \bar{\lambda} q_2) \leq \lambda D_F(p_1 \| q_1) + \bar{\lambda} D_F(q_2 \| q_2)$$

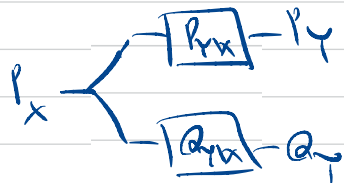
Define: $g(a, b) = b \cdot f(\frac{a}{b})$ [perspective of f]

If f is convex g is jointly convex

$$\begin{aligned} D_F(\lambda p_1 + \bar{\lambda} p_2 \| \lambda q_1 + \bar{\lambda} q_2) &= \underbrace{\sum (\lambda q_1 + \bar{\lambda} q_2) f\left(\frac{\lambda p_1 + \bar{\lambda} p_2}{\lambda q_1 + \bar{\lambda} q_2}\right)}_{\substack{\uparrow \\ \text{convex}}} \\ &= (\lambda b_1 + \bar{\lambda} b_2) f\left(\frac{\lambda a_1 + \bar{\lambda} a_2}{\lambda b_1 + \bar{\lambda} b_2}\right) \\ &= (\lambda b_1 + \bar{\lambda} b_2) f\left(\frac{\lambda b_1}{\lambda b_1 + \bar{\lambda} b_2} \frac{a_1}{b_1} + \frac{\bar{\lambda} b_2}{\lambda b_1 + \bar{\lambda} b_2} \frac{a_2}{b_2}\right) \\ &\leq (\lambda b_1 + \bar{\lambda} b_2) \left[\frac{\lambda b_1}{\lambda b_1 + \bar{\lambda} b_2} f\left(\frac{a_1}{b_1}\right) + \frac{\bar{\lambda} b_2}{\lambda b_1 + \bar{\lambda} b_2} f\left(\frac{a_2}{b_2}\right) \right] \\ &= \lambda \sum g(p_{i1}, q_{i1}) + \bar{\lambda} \sum g(p_{i2}, q_{i2}) \end{aligned}$$

$= \lambda b_1 f\left(\frac{a_1}{b_1}\right) + \bar{\lambda} b_2 f\left(\frac{a_2}{b_2}\right)$

4. Conditioning increase f -divergence

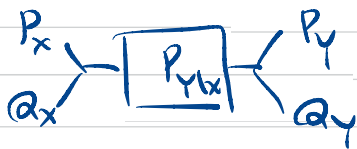


$$D_f(P_Y || Q_Y) \leq \sum_x P_X D_f(P_{Y|X=x} || Q_{Y|X=x})$$

$$(P, Q) \mapsto D_f(P || Q)$$

$$\begin{aligned} \sum_x P_X(x) D_f(P_{Y|X=x} || Q_{Y|X=x}) &\stackrel{f}{\geq} D_f\left(\sum_x P_X(x) \cdot P_{Y|X=x} \parallel \sum_x P_X(x) \cdot Q_{Y|X=x}\right) \\ &= D_f(P_Y || Q_Y) \end{aligned}$$

5. DPI



$$D_f(P_Y || Q_Y) \leq D_f(P_X || Q_X)$$

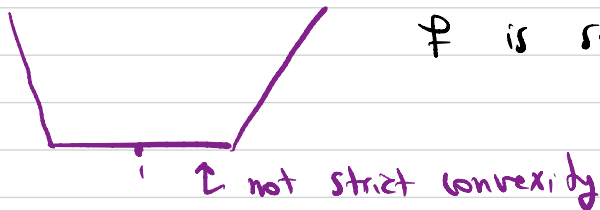
* Recall PROOF of DPI for KL depends on chain rule which doesn't hold for f -divergence!

$$\begin{aligned}
D_f(p_x \| q_x) &= \sum_x q_x(x) \cdot f\left(\frac{p_x(x)}{q_x(x)}\right) \\
&= \sum_{x(y)} p_{Y|X}(y|x) \cdot q_x(x) \cdot f\left(\frac{p_x(x)}{q_x(x)} \cdot \frac{p_{Y|X}(y|x)}{p_{Y|X}(y|x)}\right) \\
&= \sum_{x(y)} q_Y(y) \cdot \overset{\substack{\text{backward} \\ \text{channel}}}{\downarrow} q_{X|Y}(x|y) \cdot f\left(\frac{p_{XY}(x,y)}{q_Y(y) \cdot q_{X|Y}(x|y)}\right) \\
&= \sum_y q_Y(y) \cdot \underbrace{\sum_x q_{XY}(x,y)}_{\text{Jensen}} \cdot f\left(\frac{p_{XY}(x,y)}{q_Y(y) \cdot q_{X|Y}(x|y)}\right) \\
&\geq \sum_y q_Y(y) \cdot f\left(\sum_x q_{X|Y}(x|y) \cdot \frac{p_{XY}(x,y)}{q_Y(y) \cdot q_{X|Y}(x|y)}\right) \\
&= \sum_y q_Y(y) \cdot f\left(\sum_x \frac{p_{XY}(x,y)}{q_Y(y)}\right) \\
&= \sum_y q_Y(y) \cdot f\left(\frac{p_Y(y)}{q_Y(y)}\right) \checkmark
\end{aligned}$$

Remark:

1) As we proved: $D_f(p||q) = 0 \Rightarrow p=q$ Only if

f is strictly convex at 1



For example: For hockey-stick divergence,
for $\lambda > 1$



Hence, in general $E_r(p||q) = 0 \Rightarrow p=q$.

Example: $E_r(\text{Ber}(\frac{r-1}{2r}) || \text{Ber}(1/2)) = \left(\frac{r-1}{2r} - \overbrace{r(1/2)^r}^{(r-1)} \right)_+ + \left(\frac{r+1}{2r} - \overbrace{r(1/2)^r}^{1} \right)_+$

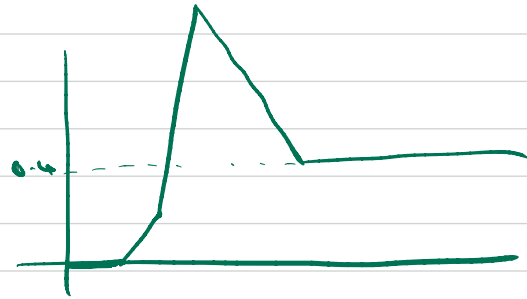
for $r > 1$

$$= \left(\underbrace{\frac{r-1}{2r} - (r-1)}_{< 0} \right)_+ + \left(\underbrace{\frac{1+r}{2r}}_{< 0} \right)_+ = 0$$

Example: $E_r(\text{Ber}(p) \parallel \text{Ber}(q)) = (p - rq)_+ + (\bar{p} - r\bar{q})_+ - (1-r)_+$

Example 2: $P = [0.1, 0.2, 0.3, 0.4]$ $Q = [0.4, 0.5, 0.1, 0]$

$$\begin{aligned} E_r(P \parallel Q) &= (0.1 - r \times 0.4)_+ + (0.2 - r \times 0.5)_+ + (0.3 - r \times 0.1)_+ \\ &\quad + (0.4 - r \times 0)_+ \\ &= (0.1 - 0.4r)_+ + (0.2 - 0.5r)_+ + (0.3 - 0.1r)_+ + 0.4 \end{aligned}$$



Exercise: $E_r(N(0, \sigma^2) \parallel N(\mu, \sigma^2)) = Q\left(\frac{\log r}{\beta} - \frac{\beta}{2}\right) - rQ\left(\frac{\log r}{\beta} + \frac{\beta}{2}\right)$

where $\beta = \frac{\|\mu\|}{\sigma}$

the same result holds even

$$N(\mu_1, \sigma^2 I_d) \parallel N(\mu_2, \sigma^2 I_d)$$

with $\beta = \|\mu_1 - \mu_2\| / \sigma$

Properties of E_γ -divergence

$$E_\gamma(p||q) := \sum_x (p(x) - \gamma q(x))_+ - (1-\gamma)_+$$

what is $E_1(p||q)$?

- $E_\gamma(p||q) \neq E_\gamma(q||p)$

Equivalent forms:

$$1. E_\gamma(p||q) = \frac{1}{2} \sum |p(x) - \gamma q(x)| - \frac{1}{2} |1-\gamma|$$

$$2. (a)_+ = |a| + a$$

$$\begin{aligned} \text{so: } E_\gamma(p||q) &= \sum (p(x) - \gamma q(x))_+ - (1-\gamma)_+ \\ &= \sum \frac{|p(x) - \gamma q(x)| + (p(x) - \gamma q(x))}{2} - \frac{|1-\gamma| + (1-\gamma)}{2} \\ &= \frac{1}{2} \sum |p(x) - \gamma q(x)| + \cancel{\frac{1-\gamma}{2}} - \frac{1}{2} |1-\gamma| - \cancel{\frac{1-\gamma}{2}} \end{aligned}$$

$$2- E_r(p||q) \stackrel{\text{For } r \geq 1}{=} \sup_A (p(A) - r q(A)) = p(A^*) - r q(A^*)$$

$$\text{For } \underline{r \leq 1}: \sup_A (r q(A) - p(A))$$

where

$$A^* = \{x: p(x) \geq r q(x)\}$$

proof requires two steps:

$$\textcircled{1} \text{ For any } A: p(A) - r q(A) \leq E_r(p||q).$$

$$\textcircled{2} \exists A^* \text{ such that } p(A^*) - r q(A^*) = E_r(p||q).$$

$$\begin{aligned} \text{L } p(A) - r q(A) &= \sum_{x \in A} (p(x) - r q(x)) \leq \sum_{x \in A} (p(x) - r q(x))_+ \\ &\quad \uparrow \text{Because we ignore all negative terms} \\ &\leq \sum_{x \in X} (p(x) - r q(x))_+ \\ &= E_r(p||q) \end{aligned}$$

because on A^* , we always $p(x) \geq r q(x)$

$$\begin{aligned}
 2- \quad p(A^*) - r q(A^*) &= \sum_{x \in A^*} (p(x) - r q(x)) = \sum_{x \in A^*} (p(x) - r q(x))_+ \\
 &= \sum_{x \in A^*} (p(x) - r q(x))_+ + \underbrace{\sum_{x \in A^{*c}} (p(x) - r q(x))}_{{=0}} \\
 &= \sum_x (p(x) - r q(x))_+ = E_r(p \| q).
 \end{aligned}$$

□

Properties of E_r -divergence:

- 1) Given p, q , $r \mapsto E_r(p \| q)$ is non-increasing & convex for $r \geq 1$
- 2 is non-decreasing & convex

- monotonicity is obvious:

Note that for $r < 1$
use the other def.

$$E_r(p||q) = \sum_{r < 1} (r q_i - p_i)_+$$

- convexity: $E_r(p||q) = \sup \underbrace{p(A) - r q(A)}_{\text{linear}} \rightarrow$ Supremum of linear function is convex

Direct approach:

$$\begin{aligned} E_{t r_1 + \bar{t} r_2}(p||q) &= \sup [p(A) - (t r_1 + \bar{t} r_2) q(A)] \\ &= \sup [p(A) - t r_1 q(A) - \bar{t} r_2 q(A)] \\ &= \sup [t (p(A) - r_1 q(A)) \\ &\quad + \bar{t} (p(A) - r_2 q(A))] \\ &\leq t \sup [p(A) - r_1 q(A)] \\ &\quad + \bar{t} \sup [p(A) - r_2 q(A)] \end{aligned}$$

2) Reciprocity: $E_r(p||q) = r \cdot E_{r^{-1}}(q||p)$

suppose $r \geq 1$: $E_r(p||q) = \sum (p - rq)_+$

& for $r < 1$: $E_r(p||q) = \sum (rq - p)_+$

$r \cdot E_{r^{-1}}(q||p) \stackrel{r^{-1}}{\downarrow} = r \cdot \left[\sum \left(\frac{1}{r} p - q \right)_+ \right] = \sum (p - rq)_+$

3) $\lim_{r \rightarrow 0} E_r(p||q) = 0$

$\lim_{r \rightarrow \infty} E_r(p||q) = 0$

\downarrow If p & q are supported on the same alphabet

Strong DPI

We saw that all f -divergences satisfies DPI. But in many cases & for many channels, DPI is strict.

$$D_f(P_Y \| Q_Y) \leq D_f(P_X \| Q_X)$$

↑ we wish to improve this

So, we wish to find smallest $\lambda \leq 1$ such that

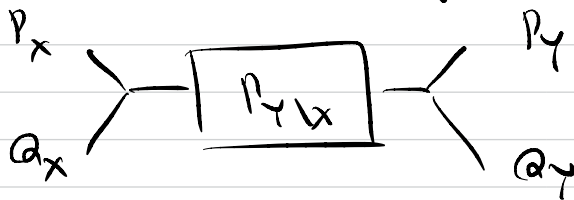
$$D_f(P_Y \| Q_Y) \leq \lambda D_f(P_X \| Q_X) \quad \forall P_X \& Q_X$$

What's the best λ ?

Smallest λ is given by $\eta_F(P_{Y|X}) = \sup_{P_X \neq Q_X} \frac{D_F(P_Y \| Q_Y)}{D_F(P_X \| Q_X)}$

We call this number: contraction coefficient of channel $P_{Y|X}$ under F -divergence,

By definition



$$D_F(P_Y \| Q_Y) \leq \eta_F(P_{Y|X}) \cdot D_F(P_X \| Q_X)$$

If $\eta_F(P_{Y|X}) < 1$ then this inequality is strictly

better than DpF \Rightarrow Stronger DPI

we then call the channel $P_{Y|X} \Rightarrow$ Contractive.

Computing $\eta_f(P_{Y|X})$ is not easy: Infinite-dimensional optimization.

But we know a lot about them for specific Dp:

1- $\eta_{X^2}(P_{Y|X}) := \sup_{P_X \neq Q_X} \frac{X^2(P_Y || Q_Y)}{X^2(P_X || Q_X)}$ is related to

Rényi maximal correlation: $\sup_{f,g} \rho(f(X), g(Y))$

2- $\eta_{KL}(P_{Y|X})$ is related propagation of information over network.

3- $\eta_{TV}(P_{Y|X})$ has been a central quantity in studying Markov chain, graph theory, dynamical system:

$$\eta_{TV}(P_{Y|X}) := \sup_{P_Y = Q_X} \frac{TV(P_Y, Q_Y)}{TV(P_X, Q_X)}$$

Theorem: (Dobrushin's two-point characterization)

For any channel $P_{Y|X}$, we have:

$$\eta_{TV}(P_{Y|X}) = \max_{x_1, x_2} TV(P_{Y|X=x_1}, P_{Y|X=x_2})$$

Proof: Strategy: ① $TV(P_{Y|X=x_1}, P_{Y|X=x_2}) \leq \eta_{TV}(P_{Y|X})$ for any pair of x_1, x_2 ↑ Dobrushin coefficient

$$② \sup_{x_1, x_2} TV(P_{Y|X=x_1}, P_{Y|X=x_2}) \geq \eta_{TV}(P_{Y|X})$$

degenerate
 ① Pick P_x & Q_x distribution on some arbitrary points:

$$P_x(x_1) = 1 \quad \text{with } x_1 \neq x_2 \quad \Rightarrow \quad TV(P_x, Q_x) = 1$$

$$Q_x(x_2) = 1$$

corresponding $P_Y = P_{Y|X=x_1}$ & $Q_Y = P_{Y|X=x_2}$

So: $TV(P_Y, Q_Y) = TV(P_{Y|X=x_1}, P_{Y|X=x_2})$

thus: $\eta_r(P_{Y|X}) = \sup_{P_x \neq Q_x} \frac{TV(P_Y, Q_Y)}{TV(P_x, Q_x)} \geq TV(P_{Y|X=x_1}, P_{Y|X=x_2})$

② Fix two dist. P_x & Q_x & define $A = \{x: P_x(x) \geq Q_x(x)\}$ $\checkmark x_1, x_2$

Now, we can write: for any $B \subseteq \mathcal{Y}$
↑
output support

$$\begin{aligned}
 P_Y(B) - Q_Y(B) &= \sum_x P_{Y|X}(B|x) P_X(x) - P_{Y|X}(B|x) \cdot Q_X(x) \\
 &= \sum_x P_{Y|X}(B|x) \cdot [P_X(x) - Q_X(x)]
 \end{aligned}$$

$$= \sum_{x \in A} P_{Y|X}(B|x) \cdot [P_X(x) - Q_X(x)] + \sum_{x \in A^c} P_{Y|X}(B|x) \cdot [P_X(x) - Q_X(x)]$$

$$= \sum_{x \in A} P_{Y|X}(B|x) \cdot [P_X(x) - Q_X(x)] - \sum_{x \in A^c} P_{Y|X}(B|x) \cdot [Q_X(x) - P_X(x)]$$

$$= \text{TV}(P, Q) \sum_{x \in A} P_{Y|X}(B|x) \frac{[P_X(x) - Q_X(x)]}{\text{TV}(P, Q)} - \sum_{x \in A^c} P_{Y|X}(B|x) \frac{[Q_X(x) - P_X(x)]}{\text{TV}(P, Q)}$$

$U(x) = \frac{P_X(x) - Q_X(x)}{\text{TV}(P, Q)}$ $V(x) = \frac{Q_X(x) - P_X(x)}{\text{TV}(P, Q)}$

note that $\sum_{x \in A} U(x) = \text{TV}(P_X, Q_X)$

$$\sum_{x \in A^c} V(x) = \text{TV}(P_X, Q_X)$$

Thus U is a pmf on A

& V is a pmf on A^c

$$= \text{TV}(P, Q) \left[\sum_A U(x) \cdot P_{Y|X}(B|x) - \sum_{x \in A^c} V(x) \cdot P_{Y|X}(B|x) \right]$$

$$= TV(P, Q) \left[\sum_{x \in A} \sum_{x' \in A^c} u(x) \cdot v(x') P_{Y|X}(B|x) - \sum_{x \in A} \sum_{x' \in A^c} v(x') \cdot u(x) P_{Y|X}(B|x') \right]$$

$$= TV \left[\sum_{x \in A} \sum_{x' \in A^c} u(x) \cdot v(x') \underbrace{(P_{Y|X}(B|x) - P_{Y|X}(B|x'))}_{\leq TV(P_{Y|X=x}, P_{Y|X=x'})} \right]$$

$$\leq TV(P_X, Q_X) \cdot \sup_{x, x'} TV(P_{Y|X=x}, P_{Y|X=x'})$$

$$\Rightarrow P_Y(B) - Q_Y(B) \leq TV(P_X, Q_X) \cdot \sup_{x, x'} TV(P_{Y|X=x}, P_{Y|X=x'})$$

Since it holds for all $B \subseteq \mathcal{Y}$, we can take supremum over B . \square

For more than 60 years, TV was a only F -divergence whose contraction coefficient was known in closed-form.

However, it was proved recently, that a similar two-point characterization can be proved for hockey-stick divergence

Notation: let $\eta_r(P_{Y|X})$ denote the contraction coefficient of $P_{Y|X}$ under E_r -divergence. That is, we have:

$$\eta_r(P_{Y|X}) := \sup_{\substack{P_X, Q_X \\ E_r(P_X \| Q_X) \neq 0}} \frac{E_r(P_{Y|X} \| Q_X)}{E_r(P_X \| Q_X)}$$

Theorem: For any channel $P_{Y|X}$ & $r \geq 1$, we have

$$\eta_r(P_{Y|X}) = \sup_{x_1, x_2} E_r(P_{Y|X=x_1} \| P_{Y|X=x_2})$$

* You'll prove this result in HW 1.

* Note that setting $r=1$ in this result, we can obtain Dobrushin's two-point characterization result; thus this result is a natural generalization of Dobrushin's.

* Remark: This result gives $\eta_r(P_{Y|X})$ only for $r \geq 1$.

How should we compute $\eta_r(P_{Y|X})$ for $r < 1$?

Lemma. Let $P_{Y|X}$ be a channel & $r < 1$. Then we have

$$\eta_r(P_{Y|X}) = \eta_{1/r}(P_{Y|X})$$

Proof. Use reciprocity property of E_r .

Example: consider channel $P_{Y|X}$ with binary input:

$$P_{Y|X=0} = \text{Bernoulli}(\delta) \quad P_{Y|X=1} = \text{Bernoulli}(1-\delta)$$

for some $\delta \in (0, 1/2)$.

Compute $\eta_r(P_{Y|X})$ for $r \geq 1$ & $r < 1$.

We know from the above theorem that:

$$\eta_r(P_{Y|X}) = \sup_{x_1, x_2} E_r(P_{Y|X=x_1} \| P_{Y|X=x_2}) \quad \text{for } r \geq 1.$$

* this channel is often referred to as
Binary Symmetric Channel
with crossover probability
 δ ; denoted by
BSC(δ).

Since input of channel is binary, this supremum is over all binary x_1 & x_2 . Thus:

$$\eta_r(P_{Y|X}) = \max \{ E_r(P_{Y|X=0} \| P_{Y|X=1}), E_r(P_{Y|X=1} \| P_{Y|X=0}) \}$$

$$= \max \left\{ \underbrace{\bar{Q}(\text{Bern}(\delta) \parallel \text{Bern}(\bar{\delta}))}_{= (\bar{\delta} - \tau \delta)_+}, \underbrace{\bar{E}(\text{Bern}(\bar{\delta}) \parallel \text{Bern}(\delta))}_{= (\bar{\delta} - \tau \delta)_+} \right\}$$

note that we assume $\tau \geq 1$
 $\bar{\delta} > \delta \rightarrow \bar{\delta} - \tau \delta \leq \bar{\delta} - \delta$
 $\rightarrow (\bar{\delta} - \tau \delta)_+ = (\bar{\delta} - \delta)_+$

Thus $\eta_\tau(P_{Y|X}) = (\bar{\delta} - \tau \delta)_+$ for $\tau \geq 1$. This implies that:

- $\eta_{TV}(P_{Y|X}) = 1 - 2\delta$
- $\eta_\tau(P_{Y|X}) = \frac{1}{\tau} (\tau \bar{\delta} - \delta)_+$ for $\tau < 1$.