

Preliminaries of Probability:

- Random Variables are denoted by upper-case letters, e.g., X , & its support set denoted by calligraphic letters: \mathcal{X} .
- we write $X \sim P_x$ to denote X is distributed according the dist P_x
- In this course, we are only concerned with either discrete (X is a discrete set) continuous r.v (X is a continuum).
- Discrete r.v are usually characterized by their pmf: $X \sim P_x \Leftrightarrow \Pr(X=x) = P_x(x)$
- continuous r.v " " " pdf: $X \sim P_x \& P_x \text{ has density } p_x \Leftrightarrow \Pr(X \in A) = \int_A p_x(u) dx$
- cdf: $F_x(x) := \Pr(X \leq x)$
- Expected value: $X \sim P_x : E[X] := \sum x P_x(x)$
or $E[X] := \int x p_x(u) dx$ pdf

- k^{th} Moment: $E[X^k]$.
- For example: 2nd moment: $E[X^2] = \sum x^2 p_x(x)$ or $E[X^2] = \int x^2 p_x(x) dx$
- Variance: $\text{Var}(X) := E[(X - E[X])^2] = E[X^2] - (E[X])^2$

* Examples of discrete r.v:

1. Bernoulli: $X \sim \text{Bernoulli}(p)$

\uparrow
 $0 \leq p \leq 1$

pmf

$$P_x(x) := \Pr(X=x) \quad x \in \{0,1\}$$

$$= \begin{cases} p & x=1 \\ 1-p & x=0 \end{cases}$$

$X = \{0,1\}$

$E[X] = p$

$\text{Var}(X) = p \cdot (1-p)$

2. Uniform: $X \sim U([k])$ for some integer k : $P_x(x) = \begin{cases} \frac{1}{k} & x \in [k] \\ 0 & \text{o.w.} \end{cases}$

$E[X] = \frac{k+1}{2}$

$\text{Var}(X) = \frac{k^2-1}{12}$

$= \frac{1}{k} \cdot \mathbb{1}_{\{x \in [k]\}}$

3. Binomial : $X \sim \text{Bin}(n, p)$ ^{integer n}
 & $p \in [0, 1]$

Let $B_1, \dots, B_n \stackrel{\text{iid}}{\sim} \text{Ber}(p)$. Then $X = B_1 + B_2 + \dots + B_n$

It can be easily shown that

$$\Pr(X=k) = \begin{cases} \binom{n}{k} \cdot p^k \cdot (1-p)^{n-k} & 0 \leq k \leq n \\ 0 & \text{otherwise} \end{cases}$$

$$E[X] = np$$

$$\text{Var}(X) = n \cdot p \cdot (1-p)$$

Example of continuous dist:

1- Exponential : $X \sim \text{Exp}(\lambda)$ with $\lambda > 0$:

$$f_X(x) = \begin{cases} \lambda e^{-\lambda x} & x \geq 0 \\ 0 & \text{ow} \end{cases} \quad X = \mathbb{R}^+$$

Verify that $\int_0^\infty e^{-\lambda x} dx = 1 \Rightarrow$ the above is a valid pdf.

Also, verify that:

$$E[X] = \frac{1}{\lambda} \quad \& \quad \text{Var}(X) = \frac{1}{\lambda^2}$$

2- Laplace $X \sim L(\mu, b)$ ↑ scale parameter $\mu \in \mathbb{R} \& b > 0$

Pdf: $f(x) = \frac{1}{2b} e^{-\frac{|x-\mu|}{b}}$ $X = \mathbb{R}$

Verify that $\frac{1}{2b} \int_{-\infty}^{\infty} e^{-\frac{|x-\mu|}{b}} dx = 1$.

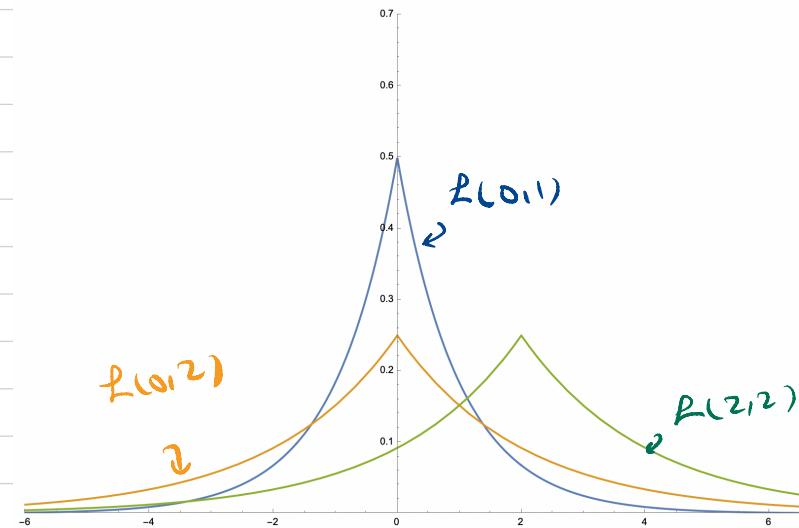
$$E[X] = \mu$$

$$\text{Var}(X) = 2b^2$$

Exercise: Let $X \sim f(\mu, b)$

Show $P(|X - \mu| \geq t) = e^{-t/b}$

Approach I: Directly



Approach II: First verify that:

$$|X - \mu| \sim \text{Exp}(b^{-1})$$

3 - Gaussian $X \sim N(\mu, \sigma^2)$ $\mu \in \mathbb{R}$

$$\text{pdf: } f_X(x) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$

$\sigma > 0$
 $x \in \mathbb{R}$

Verifying that $\int_{-\infty}^{\infty} \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(x-\mu)^2}{2\sigma^2}} dx = 1$?

↑ Note that Gaussian has a faster decay than Lap

$$E[X] = \mu \quad \& \quad \text{Var}(X) = \sigma^2.$$

Exercise: $\Pr[|X-\mu| > t] = 2Q(t/\sigma)$

where

$$Q(a) := \frac{1}{\sqrt{2\pi}} \int_a^{\infty} e^{-\frac{x^2}{2}} dx$$

$$E[(X-\mu)^k]$$

central moment
 k odd

$$= \begin{cases} 0 & k \text{ odd} \\ \sigma^k (k-1)!! & k \text{ even} \end{cases}$$

k even

double factorial

$$n!! = n \cdot (n-2) \cdot n \cdot \dots$$

Next, we need to define a family of distance measure between distributions. But before that, we need to discuss Jensen's inequality

Detour: Jensen's inequality:

Definition: * let f be a real-valued function. We say that f is convex

$$\text{if } f(\alpha x + \bar{\alpha} y) \leq \alpha f(x) + \bar{\alpha} f(y) \quad \forall x, y, \forall \alpha \in [0, 1] \quad (*)$$

* Strictly convex if \uparrow strict inequality.

* Concave if \geq

* Strictly concave $>$

We can generalize (*):

(**)

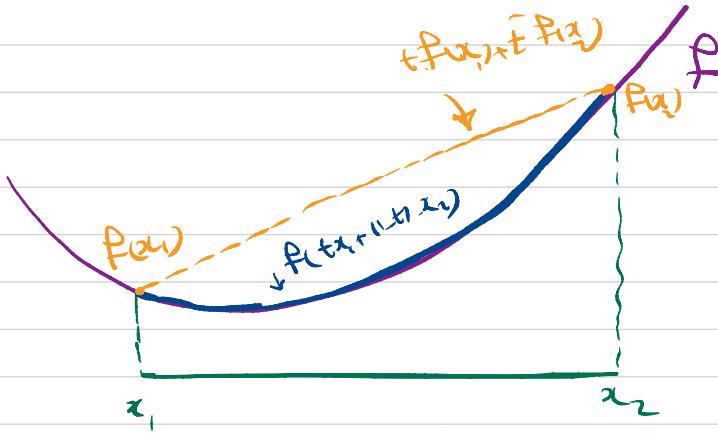
$$f(\sum p_i x_i) \leq \sum p_i f(x_i) \quad \text{for all } p_i \text{ such that } \sum p_i = 1.$$

For example:

$$\begin{aligned} f(a_1 x_1 + a_2 x_2 + a_3 x_3) &= f(a_1 x_1 + (1-a_1)(\underbrace{\frac{a_2}{1-a_1} x_2 + \frac{a_3}{1-a_1} x_3}_{\alpha})) \\ &\leq a_1 f(x_1) + (1-a_1) \cdot f\left(\frac{a_2}{1-a_1} x_2 + \frac{a_3}{1-a_1} x_3\right) \\ &\leq a_1 f(x_1) + (1-a_1) \cdot \left[\frac{a_2}{1-a_1} \cdot f(x_2) + \frac{a_3}{1-a_1} f(x_3) \right] \\ &= a_1 f(x_1) + a_2 f(x_2) + a_3 f(x_3) \end{aligned}$$

More general statement (**) can be proved by induction.

Geometric Interpretation:



Any segment connecting any two points
lies above the function.

Theorem. If f is twice differentiable, then:

$$f \text{ is convex} \Leftrightarrow f''(x) \geq 0$$

$$f \text{ is st. convex} \Leftrightarrow f'(x) \geq 0$$

$$f \text{ is concave} \Leftrightarrow f'(x) \leq 0$$

$$f \text{ is st. concave} \Leftrightarrow f'(x) \leq 0$$

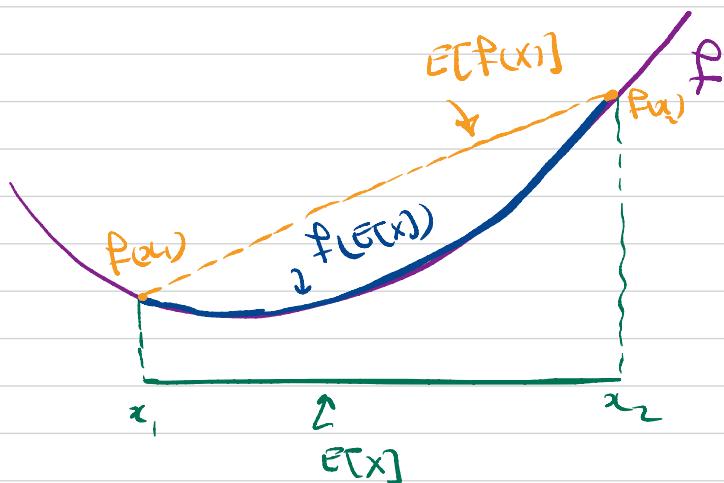
Jensen's Inequality:

Let f be a convex function & X be a random variable.

Then:

$$E[f(X)] \geq f(E[X])$$

"=" iff
 $X = E[X]$



Distance measure among distributions

Question: Given two distributions P & Q , how to quantify the "distance" between them?

In math, any appropriate measure of "distance" between any pair of objects, must satisfy the following :

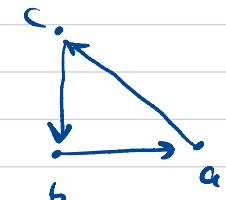
can be distributions

1. $d(a, b) = 0 \iff a = b$

2. $d(a, b) = d(b, a)$ Symmetry

3. Triangle inequality: $d(a, c) \leq d(a, b) + d(b, c)$

Example:
Euclidean
distance



We often ignore some of these properties to define "divergence" instead of distance.

I. Total Variation distance (TV)

$$TV(P, Q) := \frac{1}{2} \sum |P(x) - Q(x)|$$

↑
↑
pmf pmf

discrete

$$:= \frac{1}{2} \int |P(x) - Q(x)| dx$$

↑
↑
pdf pdf

continuous

Fact : TV is in fact a distance metric, meaning it satisfies:

$$1. \ TV(P, Q) = 0 \iff P = Q$$

$$2. \ TV(P, Q) = TV(Q, P)$$

$$3. \ TV(P, R) \leq TV(P, Q) + TV(Q, R)$$

why? $TV(P, R) = \frac{1}{2} \sum |P(x) - R(x)| \leq$

$+ Q(x) - Q(x)$

triangle inequality of 1.1

Example: $TV(Bernoulli(p), Bernoulli(q)) = \frac{1}{2} |p(0) - q(0)| + \frac{1}{2} |p(1) - q(1)|$

$$= \frac{1}{2} |1-p - (1-q)| + \frac{1}{2} |p-q|$$

$$= |p-q|$$

$TV(L(0,b), L(\mu,b)) = 1 - e^{-\frac{|\mu|}{2b}}$ assume $\mu > 0$



$$\geq \frac{1}{2} \int \left| \frac{1}{2b} e^{-\frac{|x|}{b}} - \frac{1}{2b} e^{-\frac{|x-\mu|}{b}} \right| dx = \frac{1}{4b} \left[\int_{-\infty}^{\mu/2} \cdot + \int_{\mu/2}^{\infty} \cdot \right]$$

Exercise: $TV(N(0, \sigma^2), N(\mu, \sigma^2)) = Q(-\frac{\mu}{\sigma}) - Q(\frac{\mu}{\sigma})$

$$= 1 - Q(\frac{\mu}{\sigma})$$

Theorem. For any dist. $P \neq Q$ on support set A , we have:

$$TV(P, Q) = \sup_{A \subseteq A} [P(A) - Q(A)]$$

Equivalent expression
for TV

Proof- we need to prove:

$$\textcircled{1} \quad P(A) - Q(A) \leq TV(P, Q) \quad \forall A \in A$$

$$\textcircled{2} \quad \exists A^* \in A \text{ s.t. } P(A^*) - Q(A^*) = TV(P, Q).$$

↓ Note $P(A) - Q(A) = \frac{1}{2} (P(A) - Q(A)) + \frac{1}{2} (Q(A) - P(A))$

$$\textcircled{1} \quad \text{Take arbitrary } A: \quad P(A) - Q(A) = \frac{1}{2} (P(A) - Q(A)) + \frac{1}{2} (Q(A) - P(A)) = \frac{1}{2} \sum_{x \in A} (P(x) - Q(x)) + \frac{1}{2} \sum_{A^c} (Q(x) - P(x))$$

$$\leq \frac{1}{2} \sum_A |P(x) - Q(x)| + \frac{1}{2} \sum_{A^c} |Q(x) - P(x)|$$

$$= \frac{1}{2} \sum |P(x) - Q(x)| = TV(P, Q)$$

What is the maximizing A^* ? $A^* = \{x : p(x) \geq q(x)\}$

$$\begin{aligned} TV(p, q) &= \frac{1}{2} \sum_x |p(x) - q(x)| = \frac{1}{2} \left[\sum_{x \in A^*} |p(x) - q(x)| + \sum_{x \in A^{*c}} |p(x) - q(x)| \right] \\ &= \frac{1}{2} \left[\sum_{x \in A^*} (p(x) - q(x)) + \sum_{x \in A^{*c}} (q(x) - p(x)) \right] \\ &= P(A^*) - Q(A^*) \end{aligned}$$

□

Another equivalent expression for TV:

$$TV(p, q) = \frac{1}{2} \sup_{f: |f(x)| \leq 1 \forall x} E_p[f(x)] - E_q[f(x)]$$

Proof is similar!

There are two other equivalent expressions for TV:

Supported on finite set

* Theorem: For any $P \otimes Q^V$, we have:

$$TV(P, Q) = \inf_{P_{XY}} \Pr(X \neq Y)$$

$\begin{cases} P_X = P \\ Q_Y = Q \end{cases}$ ← minimization is over all "coupling" of P & Q .

* Theorem: For any P, Q , we have:

$$\begin{aligned} TV(P, Q) &= 1 - \int \min\{p(x), q(x)\} dx \\ &= 1 - \sum \min\{p(x), q(x)\} \end{aligned}$$

II. kullback - Leibler KL Divergence :

$$D(P \parallel Q) := E_p \left[\log \frac{P(x)}{Q(x)} \right]$$

↑
 natural
 log.

discrete
 $\Rightarrow \sum_x p(x) \cdot \log \frac{P(x)}{Q(x)}$
 cont.
 $\Rightarrow \int_X p(x) \cdot \log \frac{P(x)}{Q(x)}$

Example:

$$\begin{aligned}
 & - D(\text{Bernoulli}(p) \parallel \text{Bernoulli}(q)) \\
 &= P(1) \cdot \log \frac{P(1)}{Q(1)} + P(0) \cdot \log \frac{P(0)}{Q(0)} \\
 &= p \log \frac{p}{q} + (1-p) \cdot \log \frac{1-p}{1-q}.
 \end{aligned}$$

$$\begin{aligned}
 & - D(F(0, b) \parallel F(\mu, b)) = b e^{-|\mu|/b} - b - |\mu|
 \end{aligned}$$

↑
 assume $\mu > 0$

$$\begin{aligned}
 & \frac{1}{2b} \int e^{-\frac{|x|}{b}} \cdot \log \frac{e^{-|x|/b}}{e^{-|x-\mu|/b}} dx = \int_{-\infty}^0 + \int_0^\mu + \int_\mu^\infty
 \end{aligned}$$

$$\begin{aligned}
 - D(N(a, \sigma_1^2) || N(b, \sigma_1^2)) &= E_{N(a, \sigma_1^2)} \left[\log \frac{e^{-\frac{(x-a)^2}{2\sigma_1^2}}}{e^{-\frac{(x-b)^2}{2\sigma_1^2}}} \right] \\
 &= \frac{1}{2\sigma_1^2} \cdot E[(a-b)(2x-a-b)] \log e \\
 &= \frac{1}{2\sigma_1^2} (a-b) \cdot \log e \underbrace{E[(2x-a-b)]}_{a-b} \\
 &= \frac{(a-b)^2}{2\sigma_1^2}
 \end{aligned}$$

Exercise:

$$D(N(a, \sigma_1^2) || N(b, \sigma_2^2)) = \frac{1}{2} \log \frac{\sigma_2^2}{\sigma_1^2} + \frac{\sigma_1^2 + (a-b)^2}{2\sigma_2^2} - \frac{1}{2}$$

Properties of KL divergence:

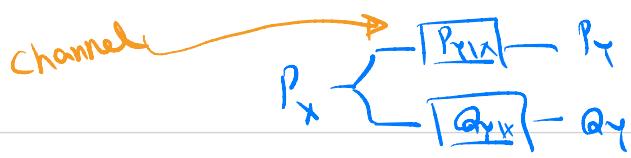
1) $D(P||Q) \geq 0 \quad \forall P, Q$, " = " if & only if $P=Q$

2) $D(P_{xy}||Q_{xy}) = D(P_x||Q_x) + \sum P_x(x) D(P_{y|x=x}||Q_{y|x=x})$
 ↑ $P_{x|x}$ $Q_{x|x}$ marginal

3) $D(P_{xy}||Q_{xy}) \geq D(P_x||Q_x)$
 or $D(Q_y||Q_y)$
 monotonicity

4) Let P_y & Q_y be output dist. of $P_{x|x}$ & $Q_{x|x}$ with common input P_x . Then

$$D(P_y||Q_y) \leq \sum p(x) D(P_{y|x=x}||Q_{y|x=x})$$



5) Data Processing Inequality:

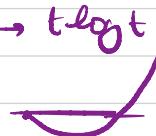


$$D(P_y \| Q_y) \leq D(P_x \| Q_x)$$

→ Outputs of a channel are closer than the corresponding inputs.

Proofs. 1. $D(P \| Q) = \sum_x p(x) \log \frac{p(x)}{q(x)}$

$$= \sum_x q(x) \cdot \frac{p(x)}{q(x)} \log \frac{p(x)}{q(x)} = \sum q(x) \cdot f\left(\frac{p(x)}{q(x)}\right)$$



Jensen's ineq $\geq f\left(\sum_n q(x_i) \frac{p(x_i)}{q(x_i)}\right) = f(1) = 0$

$$\begin{aligned}
 2. \quad D(P_{xy} || Q_{xy}) &= \sum_{x,y} P_{xy} \cdot \log \frac{P_{xy}}{Q_{xy}} = \sum_{x,y} p(x, y) \cdot P_{xy|x} \cdot \log \frac{P(x, y)}{Q(x, y)} \\
 &= \sum_{x,y} p(x, y) \cdot P_{xy|x} \cdot \log \frac{P(x)}{Q(x)} + \sum_{x,y} p(x, y) \cdot P_{xy|x} \cdot \log \frac{P(y|x)}{Q(y|x)} \\
 &= \sum_x p(x) \cdot \log \frac{P(x)}{Q(x)} + \sum_x p(x) \cdot \sum_y p(y|x) \cdot \log \frac{P(y|x)}{Q(y|x)}
 \end{aligned}$$

$$\begin{aligned}
 3. \text{ Note } D(P_{xy} || Q_{xy}) &= D(P_x || Q_x) + \underbrace{\sum_{y|x} D(P_{y|x=x} || Q_{y|x=x})}_{\geq 0} \\
 \Rightarrow D(P_{xy} || Q_{xy}) &\geq D(P_x || Q_x)
 \end{aligned}$$

4. consider joint dist. $P_{xy} = P_x \cdot P_{y|x}$ & $Q_{xy} = P_x \cdot Q_{y|x}$

$$D(P_{xy} || Q_{xy}) = \underbrace{D(P_x || P_x)}_{\substack{\uparrow \\ \text{chain rule}}} + \sum_{x} p(x) \cdot D(P_{y|x=x} || Q_{y|x=x})$$

By monotonicity: $D(P_y || Q_y) \leq D(P_{xy} || Q_{xy}) = \sum p_{xy} D(P_{y|x=x} || Q_{y|x=x})$

5. $\sum_{x,y} p_{xy} D(P_{y|x=x} || Q_{y|x=x})$

$$D(P_{xy} || Q_{xy}) = D(P_x || Q_x) + \sum_{x,y} p_{xy} D(P_{y|x=x} || Q_{y|x=x})$$

By monotonicity: $D(P_y || Q_y) \leq D(P_{xy} || Q_{xy}) = D(P_x || Q_x)$

* Note that

1) $D(\rho || \sigma) \neq D(\sigma || \rho)$ * Not symmetric *

2) $D(\rho || \tau) \leq D(\rho || \sigma) + D(\sigma || \tau)$ No triangle inequality

* So KL-divergence is not a distance metric

Theorem. (Pinsker's Inequality.)

$$TV(P, Q) \leq \sqrt{\frac{D(P||Q)}{2}}$$

* You may see different constant in RHS, but 2 is optimal. Note that this doesn't mean that

* Note that $0 \leq TV(P, Q) \leq 1$
 \uparrow
 $= P = Q$
 \downarrow
 if P & Q are mutually singular; i.e.,

However;
 $0 \leq D(P||Q) \leq \infty$ they have
 \uparrow non-overlapping support.

this is the best inequality connecting TV & D .

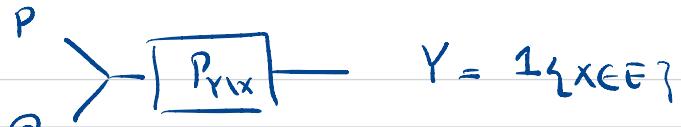
if & only if
 $\exists x_0 \in \text{Supp}(P)$

Proof: We first show that it suffices to prove this inequality but $x_0 \notin \text{Supp}(Q)$.

for Bernoulli dist. Then, we prove it assuming P & Q are Bernoulli.

Equivalently
 $\text{Supp}(Q) \subseteq \text{Supp}(P)$

I. consider the following channel:



Note that when $X \sim p \Rightarrow Y \sim \text{Ber}(p(E))$

where $E = \{x : P(x) > Q(x)\}$

$$X \sim q \Rightarrow Y \sim \text{Ber}(q(E))$$

$$\begin{aligned} \Pr[Y=1] &= \sum_x P(X=x) \cdot P(Y=1|X=x) \\ &= \sum_x p(x) \cdot \underbrace{P_{Y|X}(1|x)}_{1_{\{x \in E\}}} = p(E) \end{aligned}$$

$$\stackrel{\text{DPI}}{\geq} D(p||q) \geq D(\text{Ber}(p(E)) || \text{Ber}(q(E)))$$

Pinsker's inequality
for Bernoulli

$$\begin{aligned} 2 \text{TV}^2(\text{Ber}(p(E)), \text{Ber}(q(E))) &= 2 \cdot (p(E) - q(E))^2 \\ &= 2 \text{TV}^2(p, q) \end{aligned}$$

So we only need to prove it for binary:

II. Suppose $p = \text{Ber}(p)$ & $q = \text{Ber}(q)$. WLOG, we assume $p \geq q$.

[Note that $D(\text{Ber}(p) || \text{Ber}(q)) = D(\text{Ber}(\bar{p}) || \text{Ber}(\bar{q}))$ so we can always assume $p \geq q$]

$$D(p||q) = p \log \frac{p}{q} + \bar{p} \log \frac{\bar{p}}{\bar{q}} \quad \& \quad TV(p,q) = (p-q)$$

Define,

$$g(p,q) := p \log \frac{p}{q} + \bar{p} \log \frac{\bar{p}}{\bar{q}} - 2(p-q)^2. \quad \text{we need to show that} \quad g(p,q) \geq 0 \quad \forall p,q.$$

$$\frac{\partial g}{\partial q} = -\frac{p}{q} + \frac{\bar{p}}{\bar{q}} - 4(p-q) = \frac{q-p}{q(1-q)} - 4(p-q) \leq 0$$

so $q \mapsto g(p,q)$ is decreasing for $q \leq p \Rightarrow g(p,q) \geq g(1,p) = 0$



* Important Observation: a relationship between $TV \approx KL$

was obtained just by looking at the binary case!

* We discussed two measures of "distance": TV & KL .

We can construct a family of other divergences:

Def. Let $f: (0, \infty) \rightarrow \mathbb{R}$ be a convex function satisfying $f(1)=0$. Given two distributions P & Q , their f -divergence is:

$$D_f(P||Q) := E_Q \left[f \left(\frac{P(x)}{Q(x)} \right) \right]$$

dis $\stackrel{\cong}{=} \sum_x Q(x) \cdot f \left(\frac{P(x)}{Q(x)} \right)$ pmfs of P & Q

$$= \int q(x) \cdot f \left(\frac{p(x)}{q(x)} \right) dx$$

pdf of P & Q

Examples:

$$1. \quad F(t) = \frac{1}{2}(t-1) \quad : \quad D_F(P||Q) = E_Q \left[\frac{1}{2} \left| \frac{P(x)}{Q(x)} - 1 \right| \right]$$

$$\begin{aligned} & \text{discrete} \\ & \geq \frac{1}{2} \cdot \sum Q(\alpha_i) \cdot \left| \frac{P(\alpha_i)}{Q(\alpha_i)} - 1 \right| \\ & = \frac{1}{2} \sum |P(\alpha_i) - Q(\alpha_i)| \\ & = TV(P, Q). \end{aligned}$$

$$2. \quad F(t) = t \log t \quad : \quad D_F(P||Q) = D(P||Q)$$

$$3 - f(t) = (\sqrt{t} - 1)^2 : D_f(p||q) = \sum q(x) \cdot (\sqrt{\frac{p(x)}{q(x)}} - 1)^2$$

$$\begin{aligned} &= 2 \left(p(x) + q(x) - 2 \sqrt{p(x)q(x)} \right) \\ &= 2 - 2 \sum_x \sqrt{p(x)q(x)} \\ &= H(p, q) \end{aligned}$$

Squared Hellinger
distance

Can be obtained by $f(t) = 2(1 - \sqrt{t})$

$$4 - f(t) = (t - 1)^2 \quad \chi^2 \text{- divergence}$$

Same thing
for $f(t) = t^2 - 1$.

$$5 - f(t) = \frac{t^\alpha - 1}{\alpha - 1} \quad \text{for } \alpha \in (0, 1) \cup (1, \infty)$$

$$\text{or } f(t) = \frac{1}{t} - t$$

χ^2 divergence

$$\chi^\alpha(p||q) = \frac{1}{\alpha-1} \left[\sum_x P(x) Q(x)^{\frac{1-\alpha}{\alpha}} - 1 \right]$$

$\alpha=1/2 \rightarrow$ Squared Hellinger distance:

We are often interested in a function of χ^α :

$$D_\alpha(p||q) := \frac{1}{\alpha-1} \left[\log \left(1 + (\alpha-1) \cdot \chi^\alpha(p||q) \right) \right]$$

* Rényi divergence *

$$6 \quad f(t) = t \log t - (1+t) \cdot \log \left(\frac{1+t}{2} \right)$$

$$D_F(p||q) = D(p||\frac{p+q}{2}) + D(q||\frac{p+q}{2})$$

↑ Symmetric

7. Hockey-Stick divergence:
or \mathbb{E}_r -divergence

$$\mathbb{E}_r(p \parallel q) := D_{f_r}(p \parallel q)$$

$$\text{Hence, } \mathbb{E}_r(p \parallel q) = \sum_{x=1}^{r \geq 1} q(x) \cdot \left(\frac{p(x)}{q(x)} - r \right)_+ = \sum (p(x) - r q(x))_+$$

&

$$\mathbb{E}_r(p \parallel q) = \sum_{x=1}^{r \leq 1} q(x) \left(\frac{p(x)}{q(x)} - r \right)_+ - 1 \cdot r = \sum (p(x) - r q(x))_+ - \underbrace{(1-r)}_{\sum (p(x), r q(x))}$$

$$= \sum q(x) \left(r - \frac{p(x)}{q(x)} \right)_+ = \sum (r q(x) - p(x))_+$$

we could equivalently define: $f_r(t) = \begin{cases} (t-r)_+ & r \geq 1 \\ (r-t)_+ & r < 1 \end{cases}$

$$\text{Fix } r \geq 0 : \quad f_r(t) = (t-r)_+ - (1-r)_+$$

