

## Lec 3: $f$ -divergence

---


Yarjun Han

---

---

---

---



Defn. ( $f$ -divergence, Csiszár '63)

Let  $f: (0, \infty) \rightarrow \mathbb{R}$  be convex with  $f(1) = 0$ . The  $f$ -divergence between two distributions  $P$  and  $Q$  on the same space is

$$D_f(P \parallel Q) = \mathbb{E}_Q \left[ f\left(\frac{dP}{dQ}\right) \right].$$

Remark: 1. Some defn. additionally assumes that  $f'(1) = 0$ . This is WLOG:

$f(x)$  and  $f(x) + c(x-1)$  give the same  $f$ -divergence.

2. If  $\frac{dP}{dQ} = 0$ , define  $f(0) := f(0+)$ ;

If  $P \not\ll Q$ , define  $D_f(P \parallel Q) = \int_{q=0} q f\left(\frac{p}{q}\right) d\mu + f'(\infty) P(q=0)$ ,  
with  $f'(\infty) := \lim_{x \rightarrow \infty} \frac{f(x)}{x}$ .

Examples. ★ 1:  $f(x) = \frac{1}{2}|x-1|$ :  $D_f(P \parallel Q) = TV(P, Q) = \frac{1}{2} \int |dP - dQ|$   
(total variation (TV) distance)

★ 2:  $f(x) = (\sqrt{x}-1)^2$ :  $D_f(P \parallel Q) = H^2(P, Q) = \int (\sqrt{dP} - \sqrt{dQ})^2$   
(squared Hellinger distance)

★ 3:  $f(x) = x \log x$ :  $D_f(P \parallel Q) = D_{KL}(P \parallel Q) = \int dP \log \frac{dP}{dQ}$

★ 4:  $f(x) = (x-1)^2$ :  $D_f(P \parallel Q) = \chi^2(P \parallel Q) = \int \frac{(dP - dQ)^2}{dQ}$   
( $\chi^2$  divergence)

5.  $f(x) = \frac{1-x}{2(1+x)}$ :  $D_f(P \parallel Q) = LC(P, Q) = \frac{1}{2} \int \frac{(dP - dQ)^2}{dP + dQ}$   
(Le Cam distance)

6.  $f(x) = x \log x + (x+1) \log \frac{2}{x+1}$ :  $D_f(P \parallel Q) = JS(P, Q) = D_{KL}(P \parallel \frac{P+Q}{2}) + D_{KL}(Q \parallel \frac{P+Q}{2})$   
(Jensen-Shannon divergence)

Basic properties. ①  $D_f(P \parallel Q) \geq 0$

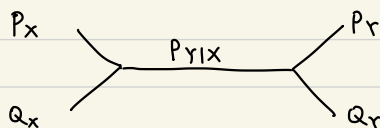
Pf.  $D_f(P \parallel Q) = \mathbb{E}_Q \left[ f\left(\frac{dP}{dQ}\right) \right] \geq f(\mathbb{E}_Q \left[ \frac{dP}{dQ} \right]) = f(1) = 0$  □

②  $(P, Q) \mapsto D_f(P \parallel Q)$  is jointly convex.

Pf. For convex  $f$ , the perspective transform  $\mathbb{R}_+^2 \ni (x, y) \mapsto y f(\frac{x}{y})$  is also convex.

Check Hessian: 
$$\begin{bmatrix} \frac{1}{y} f'(\frac{x}{y}) & -\frac{x}{y^2} f'(\frac{x}{y}) \\ -\frac{x}{y^2} f'(\frac{x}{y}) & \frac{x^2}{y^3} f''(\frac{x}{y}) \end{bmatrix} \succeq 0. \quad \square$$

③ Data processing inequality:  $D_f(P_X \parallel Q_X) \geq D_f(P_Y \parallel Q_Y)$



Pf. Follow from joint convexity (similar to the KL proof) □

Why f-divergence? Binary hypothesis testing.

Recall the simple hypothesis testing problem:

Null  $H_0$ :  $X \sim P$

Alternative  $H_1$ :  $X \sim Q$

Test:  $T: X \rightarrow \{0, 1\}$

Type I error:  $P(T(X)=1)$

Type II error:  $Q(T(X)=0)$

Thm.

$$\inf_T (P(T(X)=1) + Q(T(X)=0)) = 1 - TV(P, Q)$$

Pf. Easy to show  $TV(P, Q) = \sup_A P(A) - Q(A)$

( $\leq$ ) Take  $T(X) = 1(X \in A)$  for  $A$  attaining the supremum;

( $\geq$ ) Take  $A = \{T(X)=1\}$ . □

Remark: ①  $TV(P, Q) = 0$  :  $P = Q$ , totally indistinguishable

②  $TV(P, Q) = 1$  :  $P \perp Q$ , perfectly distinguishable

③  $TV(P, Q) < 1$  : partially indistinguishable

(Important quantity for establishing minimax lower bounds later)

Why not just TV?

①  $TV(P, Q)$  can be hard to compute

② TV does not tensorize; e.g.  $TV(P^{\otimes n}, Q^{\otimes n}) \leq n TV(P, Q)$  is the best possible inequality in general, but is often loose.

Example. How large is  $TV(\text{Bern}(\frac{1}{2})^{\otimes n}, \text{Bern}(\frac{1}{2} + \delta)^{\otimes n})$ ?

Using  $TV(P^{\otimes n}, Q^{\otimes n}) \leq n TV(P, Q)$  :  $n\delta$  upper bound

Using Pinsker's inequality:  $TV(P^{\otimes n}, Q^{\otimes n}) \leq \sqrt{\frac{1}{2} D_{KL}(P^{\otimes n} \| Q^{\otimes n})}$   
 $= \sqrt{\frac{n}{2} D_{KL}(P \| Q)} = O(\sqrt{n\delta})!$

Popular  $f$ -divergences that tensorize:

①  $H^2$ :  $1 - \frac{1}{2} H^2(\prod_i P_i, \prod_i Q_i) = \prod_i (1 - \frac{1}{2} H^2(P_i, Q_i))$

② KL:  $D_{KL}(\prod_i P_i \| \prod_i Q_i) = \sum_i D_{KL}(P_i \| Q_i)$

③  $\chi^2$ :  $\chi^2(\prod_i P_i \| \prod_i Q_i) + 1 = \prod_i (\chi^2(P_i \| Q_i) + 1)$ .

Remark (optional): All of them follow from the tensorization of Rényi divergences

i.e.  $D_\lambda(\prod_i P_i \| \prod_i Q_i) = \sum_i D_\lambda(P_i \| Q_i)$ , with

$$D_\lambda(P \| Q) \triangleq \frac{1}{\lambda - 1} \log \mathbb{E}_Q \left[ \left( \frac{dP}{dQ} \right)^\lambda \right].$$

For  $\lambda = \frac{1}{2}, 1, 2$ ,  $D_\lambda$  corresponds to  $H^2$ , KL and  $\chi^2$ .

## Similarities and differences between f-divergences

Locally  $\chi^2$ -like: when  $f''(1)$  exists and  $P \approx Q$ :

$$\begin{aligned} D_f(P \parallel Q) &= \mathbb{E}_Q \left[ f\left(\frac{dP}{dQ}\right) \right] \\ &\approx \mathbb{E}_Q \left[ \underbrace{f(1)}_{=0} + \underbrace{f'(1)}_{\mathbb{E}_Q[\cdot]=0} \left(\frac{dP}{dQ} - 1\right) + \frac{f''(1)}{2} \left(\frac{dP}{dQ} - 1\right)^2 \right] \\ &= \frac{f''(1)}{2} \chi^2(P \parallel Q). \end{aligned}$$

In parametric models: Fisher information: if  $(P_\theta)_{\theta \in \Theta}$  is a "regular" parametric model with  $\theta \in \mathbb{R}^d$ , then for  $h \in \mathbb{R}^d$  and  $t \approx 0$ :

$$\begin{aligned} \chi^2(P_{\theta+th} \parallel P_\theta) &= \int \frac{(f_{\theta+th} - f_\theta)^2}{f_\theta} \mu(dx) \quad (\text{assume } \frac{dP_\theta}{d\mu} = f) \\ &\approx t^2 h^T \int \frac{(\dot{f}_\theta)^2}{f_\theta} \mu(dx) h \quad (\dot{f}_\theta(x) = \frac{\partial f}{\partial \theta}(x)) \\ &=: t^2 h^T I(\theta) h, \end{aligned}$$

where  $I(\theta) \in \mathbb{R}^{d \times d}$  is the Fisher information:

$$\begin{aligned} I(\theta) &= \int \frac{(\dot{f}_\theta)^2}{f_\theta} d\mu = \mathbb{E}[(\nabla_\theta \log f_\theta(X))(\nabla_\theta \log f_\theta(X))^T] \\ &= \mathbb{E}[-\nabla_\theta^2 \log f_\theta(X)]. \end{aligned}$$

f-divergence as "average statistical information".

In binary hypothesis testing, if  $P(H_1) = \pi \in (0, 1)$ , then the Bayes error is

$$\begin{aligned} B_\pi(P, Q) &= \inf_T (\pi P(T(X)=1) + (1-\pi) Q(T(X)=0)) \\ &= \int (\pi dP \wedge (1-\pi) dQ) \quad (x \wedge y := \min\{x, y\}) \end{aligned}$$

The statistical information is the difference between "a priori" and "a posteriori" Bayes losses:

$$I_{\pi}(P, Q) = \pi \wedge (1 - \pi) - B_{\pi}(P, Q),$$

which is a  $f$ -divergence with  $f_{\pi}(t) = \pi \wedge (1 - \pi) - (\pi t) \wedge (1 - \pi)$ .

Thm (Liese & Vajda '06). For any  $f$ -divergence,  $\exists$  a measure  $\Gamma_f$  on  $(0, 1)$  s.t.

$$D_f(P \parallel Q) = \int_0^1 I_{\pi}(P, Q) \Gamma_f(d\pi) \quad \forall P, Q.$$

Remark: every  $f$ -divergence is an "average" statistical information, with different weights on  $\pi$ .

Pf.  $f(1) = 0$ , and WLOG assume  $f'(1) = 0$ . Then

$$\begin{aligned} f(t) &= \int_1^t (t-x) f''(dx) & (\text{For } f \in C^2, f''(dx) = f''(x) dx; \\ &\stackrel{\text{check}}{=} \int_0^1 (x - t \wedge x) f''(dx) & \text{in general, any convex function gives} \\ &\quad + \int_1^{\infty} (t - t \wedge x) f''(dx). & \text{rise to a "measure" } f''(dx)) \end{aligned}$$

Define  $\tilde{f}(t) = \int_0^1 (x - t \wedge x) f''(dx) + \int_1^{\infty} (1 - t \wedge x) f''(dx)$ , then

$$\mathbb{E}_a[(f - \tilde{f}) \left( \frac{dP}{dQ} \right)] = \mathbb{E}_a \left[ \int_1^{\infty} \left( \frac{dP}{dQ} - 1 \right) f''(dx) \right] = 0.$$

On the other hand,

$$1 \wedge x - t \wedge x = (1+x) \left( \frac{1}{1+x} \wedge \frac{x}{1+x} - \frac{t}{1+x} \wedge \frac{x}{1+x} \right) = (1+x) f_{\frac{1}{1+x}}(t).$$

so

$$\begin{aligned} \int_0^{\infty} (1+x) I_{\frac{1}{1+x}}(P, Q) f''(dx) &= \mathbb{E}_a \left[ \int_0^{\infty} (1+x) f_{\frac{1}{1+x}} \left( \frac{dP}{dQ} \right) f''(dx) \right] \\ &= \mathbb{E}_a \left[ \tilde{f} \left( \frac{dP}{dQ} \right) \right] = \mathbb{E}_a \left[ f \left( \frac{dP}{dQ} \right) \right] = D_f(P \parallel Q), \end{aligned}$$

and  $\Gamma_f(\pi)$  is the pushforward measure of  $(1+x) f''(dx)$  by the map

$$x \in (0, \infty) \mapsto \frac{1}{1+x} \in (0, 1)$$



## Different guarantees on contiguity

Def (contiguity)  $\{P_n\}$  is contiguous w.r.t.  $\{Q_n\}$  (written as  $\{P_n\} \triangleleft \{Q_n\}$ )  
if  $Q_n(A_n) \rightarrow 0$  implies  $P_n(A_n) \rightarrow 0$ .

Clearly,  $TV(P_n, Q_n) \rightarrow 0$  implies  $\{P_n\} \triangleleft \{Q_n\}$ .

In comparison,  $KL(P_n \parallel Q_n) \leq C$  already establishes contiguity, as

$$P_n(A_n) \log \frac{P_n(A_n)}{Q_n(A_n)} \leq KL(P_n \parallel Q_n) \leq C \quad (\text{see Lec 2})$$

$\chi^2(P_n \parallel Q_n) \leq C$  leads to an even stronger guarantee:

$$\frac{(P_n(A_n) - Q_n(A_n))^2}{Q_n(A_n)(1 - Q_n(A_n))} \stackrel{\text{DPI}}{\leq} \chi^2(P_n \parallel Q_n) \leq C$$
$$\Rightarrow P_n(A_n) \leq Q_n(A_n) + \sqrt{C \cdot Q_n(A_n)}.$$

Therefore, different  $f$ -divergences have different powers in establishing contiguity results, due to different growth of  $f(t)$  as  $t \rightarrow \infty$ . In this context, a popular choice is to upper bound  $\chi^2(P_n \parallel Q_n)$ , known as the "second moment method" in random graph theory & property testing (Lec 8).

## Dual representations of $f$ -divergence.

Similar to KL,  $f$ -divergences also admit dual representations.

Def (convex conjugate): for a convex function  $f$  on  $\mathbb{R}$ , its convex conjugate is defined as

$$f^*(y) = \sup_x (xy - f(x)).$$

- Properties:
- ①  $f^*$  is convex;
  - ②  $f^{**} = f$ ;
  - ③ Young's inequality:  $f(x) + f^*(y) \geq xy$ .

The following result is then immediate:

Thm.  $D_f(P \parallel Q) = \sup_{g: \mathbb{E}_Q[f^* \circ g] < \infty} \mathbb{E}_P g - \mathbb{E}_Q[f^* \circ g].$

Pf.  $D_f(P \parallel Q) = \mathbb{E}_Q \left[ f \left( \frac{dP}{dQ} \right) \right] = \mathbb{E}_Q \left[ \sup_y y \frac{dP}{dQ} - f^*(y) \right]$

$$= \sup_{g: X \rightarrow \mathbb{R}} \mathbb{E}_P g - \mathbb{E}_Q[f^* \circ g]. \quad \square$$

Example 1 (TV): When  $f(x) = \frac{1}{2}|x-1|$ ,  $f^*(y) = \begin{cases} y & \text{if } |y| \leq \frac{1}{2} \\ +\infty & \text{if } |y| > \frac{1}{2} \end{cases}$ , so

$$TV(P, Q) = \sup_{\|g\|_\infty \leq \frac{1}{2}} \mathbb{E}_P g - \mathbb{E}_Q g = \frac{1}{2} \sup_{\|g\|_\infty \leq 1} |\mathbb{E}_P g - \mathbb{E}_Q g|.$$

Example 2 (KL): When  $f(x) = x \log x$ ,  $f^*(y) = e^{y-1}$ , so

$$D_{KL}(P \parallel Q) = \sup_g \mathbb{E}_P g - \mathbb{E}_Q e^{g-1} = \sup_g \mathbb{E}_P g - (\mathbb{E}_Q e^g - 1).$$

As  $\mathbb{E}_Q e^g - 1 \geq \log \mathbb{E}_Q e^g$ , this is weaker than Donsker-Varadhan.

A way to recover Donsker-Varadhan is

$$\begin{aligned} D_{KL}(P \parallel Q) &= \sup_g \sup_{a \in \mathbb{R}} \mathbb{E}_P [g+a] - \mathbb{E}_Q e^{g+a-1} \\ &= \sup_g \left( \mathbb{E}_P [g] - \underbrace{\inf_{a \in \mathbb{R}} (\mathbb{E}_Q e^{g+a-1} - a)}_{= \log \mathbb{E}_Q e^g} \right) \\ &= \log \mathbb{E}_Q e^g, \text{ by taking } a = 1 - \log \mathbb{E}_Q e^g. \end{aligned}$$



Example 3 ( $\chi^2$ ): When  $f(x) = (x-1)^2$ ,  $f^*(y) = y + \frac{y^2}{4}$ , so

$$\begin{aligned}\chi^2(P\|Q) &= \sup_g \mathbb{E}_P[g] - \mathbb{E}_Q[g + \frac{g^2}{4}] \\ &= \sup_g \sup_{\lambda, c \in \mathbb{R}} \mathbb{E}_P[\lambda(g+c)] - \mathbb{E}_Q[\lambda(g+c) + \frac{\lambda^2(g+c)^2}{4}] \\ &= \sup_g \frac{(\mathbb{E}_P g - \mathbb{E}_Q g)^2}{\text{Var}_Q g}.\end{aligned}$$

Corollary (Hammersley-Chapman-Robbins (HCR) lower bound)

In a parametric family  $(P_\theta)_{\theta \in \mathbb{R}}$ , if an estimator  $\hat{\theta}$  is unbiased, then

$$\text{Var}_\theta(\hat{\theta}) \geq \sup_{\theta' \neq \theta} \frac{(\theta - \theta')^2}{\chi^2(P_{\theta'}\|P_\theta)}.$$

In particular, by taking  $\theta' \rightarrow \theta$ , it recovers the Cramér-Rao bound

$$\text{Var}_\theta(\hat{\theta}) \geq \frac{1}{I(\theta)}.$$

Example 4 (JS): When  $f(x) = x \log x + (x+1) \log \frac{2}{x+1}$ ,  $f^*(y) = \begin{cases} -\log(2-e^y), & y \leq \log 2 \\ +\infty, & y \geq \log 2 \end{cases}$

$$\begin{aligned}\text{JS}(P, Q) &= \sup_{g \leq \log 2} \mathbb{E}_P g + \mathbb{E}_Q [\log(2-e^g)] \\ &\stackrel{h = \frac{e^g}{2}}{=} \sup_{0 < h < 1} \mathbb{E}_P [\log h] + \mathbb{E}_Q [\log(1-h)] + \log 2.\end{aligned}$$

So generative adversarial networks (GAN) aim to minimize

$$\min_G \text{JS}(P, P_{G(Z)}) = \min_G \sup_D \mathbb{E}_{X \sim P} [\log D(X)] + \mathbb{E}_{Z \sim N} [\log(1 - D(G(Z)))]$$

$\uparrow$   
generator
 $\uparrow$   
data distribution
 $\uparrow$   
noise
 $\uparrow$   
discriminator

Joint range: given two  $f$ -divergences, how to prove inequalities between them?

(For example, is there a general paradigm to prove Pinsker's inequality

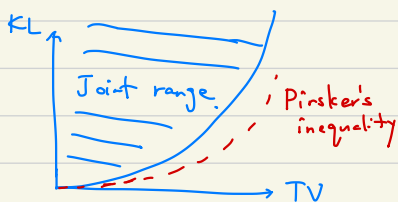
$$2TV(P, Q)^2 \leq D_{KL}(P \parallel Q)?)$$

Def (Joint range): Fix two  $f$ -divergences  $D_f(P \parallel Q)$  and  $D_g(P \parallel Q)$ .

Define:  $R = \{(D_f(P \parallel Q), D_g(P \parallel Q)) : P, Q \text{ general prob. measures}\}$

$R_K = \{(D_f(P \parallel Q), D_g(P \parallel Q)) : P, Q \text{ prob. measures on } [K]\}$ .

Example (TV vs. KL):



Thm (Harremoës-Vajda'11)  $R = \text{conv}(R_2) = R_4$ .

Implication: to establish inequalities between  $D_f$  and  $D_g$ , suffices to prove them for  $P = (p, 1-p)$  and  $Q = (q, 1-q)$ !

Pf (of a simpler case  $P \ll Q$ )

①  $R \subseteq \text{conv}(R_2)$ : Fix any point  $(D_f(P \parallel Q), D_g(P \parallel Q)) \in R$ .

Then  $L = \frac{dP}{dQ}$  is a RV in  $[0, \infty)$  with  $\mathbb{E}_Q[L] = 1$ , and

$$(D_f(P \parallel Q), D_g(P \parallel Q)) = (\mathbb{E}_Q[f(L)], \mathbb{E}_Q[g(L)]).$$

Next consider the set  $C$  of all prob. measures on  $[0, \infty)$  with mean 1.

For  $\mu \in C$ , we associate a point  $(\mathbb{E}_\mu f(L), \mathbb{E}_\mu g(L)) \in \mathbb{R}^2$ .

Clearly  $C$  is convex, and

extremal points of  $C = \left\{ \text{distributions with mean 1 and support size } \leq 2 \right\}$   
(i.e. all points  $x$  that cannot be expressed as  $x = \lambda y + (1-\lambda)z$  with  $y, z \in C$ ,  $\lambda \in (0, 1)$ )

In fact, if  $A_1, A_2, A_3$  form a partition of  $[0, \infty)$ , and

$$\mu = \lambda_1 \mu_1 + \lambda_2 \mu_2 + \lambda_3 \mu_3, \quad \lambda_i > 0, \quad \text{supp}(\mu_i) \subseteq A_i.$$

Then the probability and mean constraints only require

$$\begin{cases} \lambda_1 + \lambda_2 + \lambda_3 = 1, \\ \lambda_1 m(\mu_1) + \lambda_2 m(\mu_2) + \lambda_3 m(\mu_3) = 1, \end{cases}$$

which is a line containing  $(\lambda_1, \lambda_2, \lambda_3)$ . So  $\mu$  cannot be an extremal point.

Now by Choquet-Bishop-de Leeuw, any  $\mu \in C$  can be written as a convex combination of extremal points of  $C$ , i.e.  $R \subseteq \text{conv}(R_2)$ .

Thm (Choquet-Bishop-de Leeuw): if  $C$  is a metrizable convex compact subset of a locally convex topological vector space, then  $C = \text{conv}(\text{extremal}(C))$ .

②  $\text{conv}(R_2) \subseteq R_4$ : by Carathéodory theorem below, any point of  $\text{conv}(R_2) \subseteq \mathbb{R}^4$  (which is connected) can be written as a convex combination of 2 points of  $R_2$ , which belongs to  $R_4$ .

Thm (Carathéodory): Let  $S \subseteq \mathbb{R}^d$  and  $x \in \text{conv}(S)$ . Then there exists  $S' = \{x_1, \dots, x_k\}$  s.t.  $x \in \text{conv}(S')$ , with

①  $k \leq d+1$  in general;

②  $k \leq d$  if  $S$  has at most  $d$  connected components. [m]

Examples of inequalities:

① TV vs.  $H^2$ :  $\frac{H^2}{2} \leq \text{TV} \leq \sqrt{H^2(1 - \frac{H^2}{4})}$  (also the joint range)

② TV vs. KL:  $\text{TV}^2 \leq \frac{1}{2} \text{KL}$   
 $\text{TV} \leq 1 - \frac{1}{2} \exp(-\text{KL})$

③ KL vs.  $\chi^2$ :  $\text{KL} \leq \log(1 + \chi^2)$  (also the joint range)

## Special topic: chain rule for $H^2$

Thm (Jayram '09) For all  $P_{X^n}, Q_{X^n}$ :

$$H^2(P_{X^n}, Q_{X^n}) \leq C \sum_{i=1}^n \mathbb{E}_P [H^2(P_{X_i | X^{i-1}}, Q_{X_i | X^{i-1}})],$$

with  $C = \prod_{i=1}^{\infty} \frac{1}{1-2^{-i}} \approx 3.46$ .

The proof is surprisingly combinatorial. First, it suffices to prove the case  $n = 2^k$ ; for general  $2^{k-1} < n \leq 2^k$ , can consider  $P_{2^k} = P_{X^n} \otimes P_0^{2^k-n}$ ,  $Q_{2^k} = Q_{X^n} \otimes P_0^{2^k-n}$ . The proof uses several properties of  $H^2$ .

Lemma 1 ( $L^2$  geometry). For arbitrary distributions  $P_0, \dots, P_m$ :

$$\frac{1}{m} \sum_{1 \leq i < j \leq m} H^2(P_i, P_j) \leq \sum_{i=1}^m H^2(P_i, P_0).$$

Pf. This result holds for all  $L^2$  distance:

$$\frac{1}{m} \sum_{1 \leq i < j \leq m} \|P_i - P_j\|^2 \leq \sum_{i=1}^m \|P_i - P_0\|^2.$$

In fact,  $2 \cdot \text{LHS} = \frac{1}{m} \sum_{i,j=1}^m \|P_i - P_j\|^2$

$$\begin{aligned} &= \frac{1}{m} \sum_{i,j=1}^m \|P_i - P_0 - (P_j - P_0)\|^2 \\ &= \frac{1}{m} \sum_{i,j=1}^m (\|P_i - P_0\|^2 + \|P_j - P_0\|^2 - 2 \langle P_i - P_0, P_j - P_0 \rangle) \\ &= 2 \cdot \text{RHS} - \frac{2}{m} \left\| \sum_{i=1}^m (P_i - P_0) \right\|^2 \leq 2 \cdot \text{RHS}. \end{aligned}$$

Finally, note that

$$H^2(P, Q) = \int (\sqrt{p} - \sqrt{q})^2$$

is indeed on  $L^2$  distance.



Now for  $A \subseteq [n]$ , define interpolations

$$p^A = \prod_{i=1}^n (P_{x_i | x^{i-1}})^{1(i \notin A)} (Q_{x_i | x^{i-1}})^{1(i \in A)}.$$

Then  $p^\emptyset = P_{x^n}$ ,  $p^{[n]} = Q_{x^n}$ .

Lemma 2 (cut-paste property) Let  $a, b, c, d \in \{0, 1\}^n$  be the indicators of sets  $A, B, C, D \subseteq [n]$ . If  $a + b = c + d$ , then  $H^2(p^A, p^B) = H^2(p^C, p^D)$ .

Pf. 
$$\begin{aligned} H^2(p^A, p^B) &= 2 - 2 \int \sqrt{p^A p^B} \\ &= 2 - 2 \int \sqrt{\prod_{i=1}^n P_{x_i | x^{i-1}}^{1-a_i + (1-b_i)} Q_{x_i | x^{i-1}}^{a_i + b_i}} \\ &= 2 - 2 \int \sqrt{\prod_{i=1}^n P_{x_i | x^{i-1}}^{1-c_i + (1-d_i)} Q_{x_i | x^{i-1}}^{c_i + d_i}} = H^2(p^C, p^D) \quad \square \end{aligned}$$

Lemma 3 (1-factorization of cliques) For even  $n$ , the complete graph  $K_n$  can be decomposed into  $(n-1)$  edge-disjoint perfect matchings.  
(i.e. round-robin tournaments)

Example of  $n=4$ .



A geometric construction.



Put node 1 in the center of a regular polygon with  $(n-1)$  vertices. Use color  $i$  for  $(1, i)$  and all edges perpendicular to  $(1, i)$ .

Completing the proof. For  $n = 2^k$ , prove by induction on  $m = 0, 1, \dots, k$  that for any partition  $A_1, \dots, A_{2^m}$  of  $[n]$  (each of size  $2^{k-m}$ ),

$$\sum_{i=1}^{2^m} H^2(P^{A_i}, P^\emptyset) \geq c_m \cdot H^2(P^{[n]}, P^\emptyset),$$

with  $c_m = \prod_{i=1}^m (1 - 2^{-i})$ .

Base  $m = 0$ : trivial.

Induction from  $m-1$  to  $m$ :

$$\sum_{i=1}^{2^m} H^2(P^{A_i}, P^\emptyset) \stackrel{\text{Lemma 1}}{\geq} \frac{1}{2^m} \sum_{1 \leq s < t \leq 2^m} H^2(P^{A_s}, P^{A_t})$$

$$\stackrel{\text{Lemma 2}}{=} \frac{1}{2^m} \sum_{1 \leq s < t \leq 2^m} H^2(P^{A_s \cup A_t}, P^\emptyset)$$

$$\stackrel{\text{Lemma 3}}{=} \frac{1}{2^m} \sum_{a=1}^{2^m-1} \sum_{(s,t) \in E_a} H^2(P^{A_s \cup A_t}, P^\emptyset),$$

where each  $E_a$  is a perfect matching of  $K_{2^m}$ . By induction hypothesis,

$$\sum_{i=1}^{2^m} H^2(P^{A_i}, P^\emptyset) \geq \frac{2^m - 1}{2^m} c_{m-1} H^2(P^{[n]}, P^\emptyset) = c_m H^2(P^{[n]}, P^\emptyset).$$

Conclusion: choosing  $m = k$  yields

$$H^2(P^{[n]}, P^\emptyset) \leq \frac{1}{c_k} \sum_{i=1}^n H^2(P^{[i]}, P^\emptyset)$$

$$= \frac{1}{c_k} \sum_{i=1}^n \mathbb{E}_P [H^2(P_{X_{:i}} | X_{:i-1}, Q_{X_{:i}} | X_{:i-1})].$$