# Lec 3: f-divergence

Yanjun Han

<u>Defn</u>. (f-divergence, Csiszár '63)

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

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

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 $D_f(P \| Q) = \int_{q > 0} q f\left(\frac{p}{q}\right) d\mu + f'(\infty) P(q = 0)$,

with $f'(\infty) := \lim\limits_{x \to \infty} \frac{f(x)}{x}$.

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

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

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

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

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

6. $f(x) = x \log x + (x+1) \log \frac{2}{x+1}$ : $D_f(P \| Q) = JS(P, Q) = D_{KL}\left(P \| \frac{P+Q}{2}\right) + D_{KL}\left(Q \| \frac{P+Q}{2}\right)$
(Jensen-Shannon divergence)
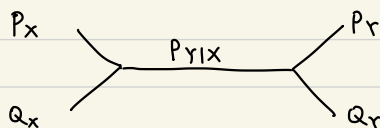
<u>Basic properties</u>. ① $D_f(P \| Q) \geq 0$

<u>Pf</u>. $D_f(P \| Q) = \mathbb{E}_Q\left[f\left(\frac{dP}{dQ}\right)\right] \geq f\left(\mathbb{E}_Q\left[\frac{dP}{dQ}\right]\right) = f(1) = 0$

☒

② $(P, Q) \longmapsto D_f(P \| Q)$ is jointly convex.

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

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

③ Data processing inequality. $D_f(P_X \| Q_X) \geq D_f(P_Y \| Q_Y)$



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: $T: X \to \{0, 1\}$
Type I error: $P(T(X) = 1)$
Type II error: $Q(T(X) = 0)$

Thm.

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

Pf. Easy to show $TV(P, Q) = \sup_A P(A) - Q(A)$
$(\leq)$ Take $T(X) = 1(X \notin A)$ for $A$ attaining the supremum;
$(\geq)$ Take $A = \{T(X) = 0\}$. 🔲

**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 <u>tensorize</u>; 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 $H^2$ : $\quad 1 - \frac{1}{2} H^2\left(\prod_i P_i, \prod_i Q_i\right) = \prod_i \left(1 - \frac{1}{2} H^2(P_i, Q_i)\right)$

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

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

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

i.e. $D_\lambda\left(\prod_i P_i \| \prod_i Q_i\right) = \sum_i D_\lambda(P_i \| Q_i)$, 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 $P \approx Q$:

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

$$\approx \mathbb{E}_Q\left[ \underbrace{f(1)}_{=0} + f'(1)\underbrace{\left(\frac{dP}{dQ} - 1\right)}_{\mathbb{E}_Q[\cdot] = 0} + \frac{f''(1)}{2}\left(\frac{dP}{dQ} - 1\right)^2 \right]$$

$$= \frac{f''(1)}{2} \chi^2(P \| Q).$$

**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 $t \approx 0$:

$$\chi^2(P_{\theta + th} \| P_\theta) = \int \frac{(f_{\theta+th} - f_\theta)^2}{f_\theta} \mu(dx) \qquad (\text{assume } \frac{dP_\theta}{d\mu} = f)$$

$$\approx t^2 h^T \int \frac{(\dot{f}_\theta)^2}{f_\theta} \mu(dx)\, h \qquad (\dot{f}_\theta(x) = \frac{\partial f}{\partial \theta}(x))$$

$$=: t^2 h^T I(\theta) h,$$

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

$$I(\theta) = \int \frac{(\dot{f}_\theta)^2}{f_\theta} d\mu = \mathbb{E}\left[ (\nabla_\theta \log f_\theta(X))(\nabla_\theta \log f_\theta(X))^T \right]$$

$$= \mathbb{E}\left[ -\nabla_\theta^2 \log f_\theta(X) \right].$$

**f-divergence as "average statistical information":**

In binary hypothesis testing, if $P(H_0) = \pi \in (0,1)$, then the Bayes error is

$$B_\pi(P, Q) = \inf_T \left( \pi\, P(T(X) = 1) + (1 - \pi)\, Q(T(X) = 0) \right)$$

$$= \int (\pi\, dP \wedge (1-\pi)\, dQ) \qquad (x \wedge y := \min\{x, y\})$$

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.

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

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

$$f(t) = \int_1^t (t-x) f''(dx)$$
$$\overset{\text{check}}{=} \int_0^1 (x - t\wedge x) f''(dx) + \int_1^\infty (t - t\wedge x) f''(dx).$$

(For $f \in C^2$, $f''(dx) = f''(x) dx$; in general, any convex function gives rise to a "measure" $f''(dx)$)

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

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

On the other hand,

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

so

$$\int_0^\infty (1+x) I_{\frac{1}{1+x}}(P, Q) f''(dx) = \mathbb{E}_Q\left[\int_0^\infty (1+x) f_{\frac{1}{1+x}}(\tfrac{dP}{dQ}) f''(dx)\right]$$
$$= \mathbb{E}_Q[\tilde{f}(\tfrac{dP}{dQ})] = \mathbb{E}_Q[f(\tfrac{dP}{dQ})] = D_f(P\|Q),$$

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

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

# Different guarantees on contiguity

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

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

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

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

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

$$\frac{(P_n(A_n) - Q_n(A_n))^2}{Q_n(A_n)(1 - Q_n(A_n))} \overset{PPI}{\leq} \chi^2(P_n \| Q_n) \leq C$$

$$\Rightarrow P_n(A_n) \leq Q_n(A_n) + \sqrt{C \cdot Q_n(A_n)}.$$

Therefore, different $f$-divergences have different powers in establishing contiguity
results, due to different growth of $f(t)$ as $t \to \infty$. In this context, a popular
choice is to upper bound $\chi^2(P_n \| Q_n)$, 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

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

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

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

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

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

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

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

$$D_{KL}(P \| Q) = \sup_g \sup_{a \in \mathbb{R}} \mathbb{E}_P[g + a] - \mathbb{E}_Q\, e^{g+a-1}$$

$$= \sup_g \left( \mathbb{E}_P[g] - \underbrace{\inf_{a \in \mathbb{R}} \left( \mathbb{E}_Q\, e^{g+a-1} - a \right)}_{= \log \mathbb{E}_Q\, e^g, \text{ by taking } a = 1 - \log \mathbb{E}_Q\, e^g.} \right)$$

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

$$\chi^2(P \| Q) = \sup_g \mathbb{E}_P[g] - \mathbb{E}_Q\left[g + \frac{g^2}{4}\right]$$

$$= \sup_g \sup_{\lambda, c \in \mathbb{R}} \mathbb{E}_P[\lambda(g+c)] - \mathbb{E}_Q\left[\lambda(g+c) + \frac{\lambda^2(g+c)^2}{4}\right]$$

$$= \sup_g \frac{(\mathbb{E}_P g - \mathbb{E}_Q g)^2}{\operatorname{Var}_Q g}.$$

<u>Corollary</u> (Hammersley-Chapman-Robbins (HCR) lower bound)

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

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

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

$$\operatorname{Var}_\theta(\hat{\theta}) \geqslant \frac{1}{I(\theta)}.$$

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

$$JS(P, Q) = \sup_{g \leqslant \log 2} \mathbb{E}_P g - \mathbb{E}_Q[\log(2 - e^g)]$$

$$\overset{h = \frac{e^g}{2}}{=} \sup_{0 < h < 1} \mathbb{E}_P[\log h] + \mathbb{E}_Q[\log(1-h)] + \log 2.$$

So generative adversarial networks (GAN) aim to minimize

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

generator  data  noise  discriminator
distribution

<u>Joint range</u> : given two f-divergences, how to prove inequalities between them?

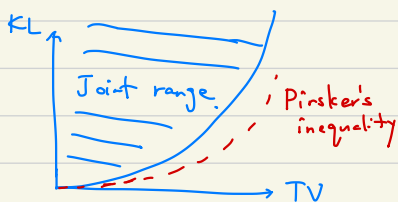<span style="color:red">(For example, is there a general paradigm to prove Pinsker's inequality</span>

<span style="color:red">$$2\,TV(P, Q)^2 \leq D_{KL}(P \| Q)? \quad )$$</span>

<u>Def (Joint range)</u> . Fix two f-divergences $D_f(P \| Q)$ and $D_g(P \| Q)$.

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

$R_k = \{ (D_f(P \| Q), D_g(P \| Q)) : P, Q \text{ prob. measures on } [k] \}$.

<span style="color:blue"><u>Example (TV vs. KL)</u> :</span>



Joint range

Pinsker's inequality

KL

TV

---

<u>Thm (Harremoës – Vajda '11)</u>     $R = \text{conv}(R_2) = R_4$.

---

<span style="color:red"><u>Implication</u> : to establish inequalities between $D_f$ and $D_g$, suffices to</span>

<span style="color:red">prove them for $P = (p, 1-p)$ and $Q = (q, 1-q)$ !</span>

<span style="color:blue"><u>Pf</u> (of a simpler case $P \ll Q$)</span>

<span style="color:blue">① $R \subseteq \text{conv}(R_2)$: Fix any point $(D_f(P \| Q), D_g(P \| Q)) \in R$.</span>

<span style="color:blue">Then $L = \dfrac{dP}{dQ}$ is a RV in $[0, \infty)$ with $\mathbb{E}_Q[L] = 1$, and</span>

<span style="color:blue">$$(D_f(P \| Q), D_g(P \| Q)) = (\mathbb{E}_Q[f(L)], \mathbb{E}_Q[g(L)]).$$</span>

<span style="color:blue">Next consider the set $C$ of all prob. measures on $[0, \infty)$ with mean 1.</span>

<span style="color:blue">For $\mu \in C$, we associate a point $(\mathbb{E}_\mu f(L), \mathbb{E}_\mu g(L)) \in \mathbb{R}^2$.</span>

<span style="color:blue">Clearly $C$ is convex, and</span>

<span style="color:blue">extremal points of $C$ = {distributions with mean 1 and support size $\leq 2$}.</span>

<span style="color:green">(i.e. all points $x$ that cannot</span>
<span style="color:green">be expressed as $x = \lambda y + (1-\lambda) z$</span>
<span style="color:green">with $y, z \in C$, $\lambda \in (0,1)$ )</span>

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 $\text{conv}(R_2) \subseteq \mathbb{R}^2$ (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 $S' = \{x_1, \cdots, x_k\}$ 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.

Examples of inequalities:
  ① TV vs. $H^2$:    $\frac{H^2}{2} \leq TV \leq \sqrt{H^2(1-\frac{H^2}{4})}$    (also the joint range)
  ② TV vs. KL:    $TV^2 \leq \frac{1}{2} KL$
                     $TV \leq 1 - \frac{1}{2} \exp(-KL)$
  ③ KL vs. $\chi^2$:    $KL \leq \log(1+\chi^2)$    (also the joint range)

Special topic: chain rule for $H^2$

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

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

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 $P_{2^k} = P_{X^n} \otimes P_o^{2^k - n}$, $Q_{2^k} = Q_{X^n} \otimes P_o^{2^k - n}$. The proof uses several properties of $H^2$.

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

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

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

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

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

$$= \frac{1}{m} \sum_{i,j=1}^{m} \| P_i - P_o - (P_j - P_o) \|^2$$

$$= \frac{1}{m} \sum_{i,j=1}^{m} \left( \| P_i - P_o \|^2 + \| P_j - P_o \|^2 - 2 \langle P_i - P_o, P_j - P_o \rangle \right)$$

$$= 2 \cdot RHS - \frac{2}{m} \| \sum_{i=1}^{m} (P_i - P_o) \|^2 \leq 2 \cdot RHS.$$

Finally, note that

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

is indeed an $L^2$ distance.

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

$$P^A = \prod_{i=1}^{n} \left( P_{x_i \mid x^{i-1}} \right)^{\mathbb{1}(i \notin A)} \left( Q_{x_i \mid x^{i-1}} \right)^{\mathbb{1}(i \in A)} .$$

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

<u>Lemma 2 (cut-paste property)</u> 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 $H^2(P^A, P^B) = H^2(P^C, P^D)$.

<u>Pf.</u> $\quad H^2(P^A, P^B) = 2 - 2 \int \sqrt{P^A P^B}$

$$= 2 - 2 \int \sqrt{\prod_{i=1}^{n} P_{x_i \mid x^{i-1}}^{1-a_i + 1 - b_i} Q_{x_i \mid x^{i-1}}^{a_i + b_i}}$$

$$= 2 - 2 \int \sqrt{\prod_{i=1}^{n} P_{x_i \mid x^{i-1}}^{1-c_i + 1 - d_i} Q_{x_i \mid x^{i-1}}^{c_i + d_i}} = H^2(P^C, P^D) \qquad \boxed{\triangle}$$

<u>Lemma 3 (1-factorization of cliques)</u> 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 $m = 0, 1, \cdots, k$ that for any partition $A_1, \cdots, A_{2^m}$ of $[n]$ (each of size $2^{k-m}$):

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

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

Base $m = 0$: trivial.

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

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

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

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

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

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

**Conclusion:** choosing $m = k$ yields

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

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