# Lecture 2

Review:   *  $\underline{P \text{ vs } Q}$:   $P_e^*(\frac{1}{2}, \frac{1}{2}) = \frac{1}{2}(1 - d(P, Q))$

Thus, we need to take as many samples $n$
as needed to make $d(P^n, Q^n)$ constant

*  $\boxed{d(P^n, Q^n) \leq n\, d(P, Q)}$   (#)

Agenda:   *  Kullback-Leibler divergence $D(P \| Q)$

- Data processing inequality

- Pinsker's inequality and improvement over (#)

- Fano's inequality  <u>Bonus: A new proof !!</u>

.

[A]  <u>Kullback-Leibler Divergence</u>

$D(P \| Q) = \begin{cases} \sum\limits_{x} P(x) \log \dfrac{P(x)}{Q(x)} & , \text{ iff } \mathrm{supp}(P) \subseteq \mathrm{supp}(Q) \\ \infty & , \qquad\qquad\qquad o.w. \end{cases}$

$= \begin{cases} \displaystyle\int f(x) \log \dfrac{f(x)}{g(x)} \mu(dx), & \text{ iff } \mathrm{supp}(f) \subseteq \mathrm{supp}(g) \\ \infty, & o.w. \end{cases}$

$= \begin{cases} \mathbb{E}_Q\left[\dfrac{dP}{dQ} \log \dfrac{dP}{dQ}\right], & \text{if } P \ll Q \\ \infty, & o.w. \end{cases}$

(a) <u>Data Processing Inequality</u>

Let $W: \mathcal{X} \to \mathcal{Y}$ be a fixed channel.

Denote by $PW$ the distribution $\sum_x P(x) W(y|x)$.

$\longrightarrow d(PW, QW) \leq d(P, Q)$ <span style="color:red">(triangle inequality)</span>

$\longrightarrow D(PW, QW) \leq D(P, Q)$

<u>Pf.</u> Follows from the <u>log-sum inequality</u>

$$\sum a_i \log \frac{a_i}{b_i} \geq \sum a_i \log \frac{\sum a_i}{\sum b_i} \quad \text{for } a_i, b_i \geq 0.$$

(b) <u>Chain rules</u>

$\longrightarrow d\left(P_{X_1 \dots X_n}, Q_{Y_1 \dots Y_n}\right) \leq \sum_{i=1}^{n} d\left(P_{X_i}, P_{X^{i-1}} Q_{Y_i | X^{i-1}}\right)$

$\longrightarrow D\left(P_{X_1 \dots X_n}, Q_{Y_1 \dots Y_n}\right) = \mathbb{E}\left[\log \frac{P_{X_1 \dots X_n}(x^n)}{Q_{Y_1 \dots Y_n}(x^n)}\right]$

$$= \mathbb{E}\left[\sum_{i=1}^{n} \log \frac{P_{X_i | X^{i-1}}(X_i | x^{i-1})}{Q_{Y_i | Y^{i-1}}(X_i | X^{i-1})}\right]$$

$$= \sum_{i=1}^{n} \mathbb{E}_{x^{i-1} \sim P_{X^{i-1}}}\underbrace{\left[D\left(P_{X_i | X^{i-1} = x^{i-1}} \| Q_{Y_i | Y^{i-1} = x^{i-1}}\right)\right]}$$

<span style="color:red">$$=: D\left(P_{X_i | X^{i-1}} \| Q_{Y_i | Y^{i-1}} | P_{X^{i-1}}\right)$$</span>

$$= \sum_{i=1}^{n} D\left(P_{X_i} \| Q_{Y_i | Y^{i-1}} P_{X^{i-1}}\right)$$

(c) <u>Pinsker's inequality</u>

$$\boxed{d^2(P, Q) \leq \frac{1}{2\ln 2} D(P \| Q)}$$

## How this improves over (#)

$$d^2(P^n, Q^n) \leq \frac{1}{2\ln 2} \cdot D(P^n \| Q^n)$$

$$= \frac{n \, D(P \| Q)}{2\ln 2}$$

$$\Rightarrow \boxed{d(P^n, Q^n) \leq \sqrt{\frac{n \, D(P \| Q)}{2\ln 2}}} \qquad (\#\#)$$

$$\left[ \text{If } D(P \| Q) \text{ is of the same order as } d^2(P, Q), \text{ then} \atop (\#\#) \text{ is an improvement over } (\#). \right]$$

## Proof of Pinsker's inequality

**Step 1.** For any $A \subseteq \mathcal{Z}$.

$$D(P \| Q) \geq D(P(A) \| Q(A)) \qquad \longrightarrow \quad P(A) \log \frac{P(A)}{Q(A)} + (1-P(A)) \log \frac{1-P(A)}{1-Q(A)}$$

by the data processing inequality.

**Step 2.** $D(p \| q) = p \log \frac{p}{q} + (1-p) \log \frac{1-p}{1-q}$

Suffices: $p \ln \frac{p}{q} + (1-p) \ln \frac{1-p}{1-q} \geq 2(p-q)^2$

**Proof:** $f(p,q) = p \ln \frac{p}{q} + (1-p) \ln \frac{1-p}{1-q} - 2(p-q)^2$

$$\frac{df}{dq} = -\frac{p}{q} + \frac{(1-p)}{1-q} + 4(p-q)$$

$$= (p-q) \left[ 4 - \underbrace{\frac{1}{q(1-q)}}_{\geq 0} \right]$$

$$\Rightarrow \frac{df}{dq} \geq 0 \quad \text{iff} \quad p \geq q \Rightarrow f(p,q) \geq f(q,q) = 0. \qquad \blacksquare$$

(d) Underline{Convexity of $\underline{D(P\|Q)}$}

$\quad$ $D(P\|Q)$ is convex in $(P,Q)$. $\left(\begin{array}{l}\text{Proof uses only log-sum}\\ \text{inequality}\end{array}\right)$

[B] $\underline{\text{Fano's inequality}}$

Recall that

$$P_e^*\left(\frac{1}{2},\frac{1}{2}\right) \geq \frac{1}{2}\left(1- d(P,Q)\right)$$

$$\geq \frac{1}{2}\left(1- \sqrt{\frac{1}{2\ln 2} D(P\|Q)}\right)$$

This bound allows us to quantize the difficulty of hypothesis

testing in terms of "distance" $D(P\|Q)$.

The next result provides a similar bound for M-ary

hypothesis testing.

$\underline{\text{Problem.}}$ $\quad H_m : X \sim P_m, \quad m = 1,\ldots, M$

$\qquad\quad d: X \rightarrow \{1,\ldots,M\}$ be a randomized map.

$\qquad\quad P_e^*(\text{unif}) = \underset{d}{\inf} \ \frac{1}{M}\sum_{m=1}^{M} P_m\left(d(X) \neq m\right)$

$\underline{\text{Theorem}}$ (Fano's inequality)

$$P_e^*(\text{unif}) \geq 1- \frac{\frac{1}{M}\sum_{m=1}^{M} D\left(P_m\|\frac{1}{M}\sum_{m'=1}^{M} P_{m'}\right) + 1}{\log M}$$

$\underline{\text{Remark.}}$ Think of $U \sim \text{unif}[M]$ as input to a channel which

then produces the output $X \sim P_u$. The quantity

$$\frac{1}{M} \sum_{m=1}^{M} D\left(P_m \,\middle\|\, \frac{1}{M} \sum_{m'=1}^{M} P_{m'}\right)$$

is then called the <u>mutual information</u> between $U$ and $X$,

denoted $I(U \wedge X)$. Note the following properties of

$I(U \wedge X)$:

$\rightarrow I(U \wedge X) = D(P_{UX} \,\|\, P_U P_X)$

$\left(\begin{array}{l} \text{for discrete} \\ U, X \end{array}\right) \quad = H(X) - H(X|U)$

$$\underbrace{\phantom{H(X|U)}}_{\color{red}{\mathbb{E}_{(X,U) \sim P_{XU}}\left[\log \frac{1}{P_{X|U}(X|U)}\right]}}$$

$\rightarrow I(U \wedge X) \leq \max_{u, u'} D(P_U \,\|\, P_{U'})$,

using the convexity of $D(P \| Q)$.

$\rightarrow I(U \wedge XY) = I(U \wedge X) + \underbrace{I(U \wedge Y | X)}$

$$\qquad\qquad\qquad \underbrace{\sum_{x} P_X(x) \; \underbrace{I(U \wedge Y | X = x)}_{\color{red}{I \text{ under } P_{UY|X=x}}}}$$

This follows from the chain rule of KL divergence.

<u>Proof of Fano's inequality</u>. Restrict first to a deterministic $d$.

$U \sim \text{unif } \{1, \ldots, M\}$; Let $Q_{UX} = P_U P_X$

Let $B = \mathbb{1}(U = d(X))$. Then,

$I(U \wedge X) = D(P_{UX} \| Q_{UX}) \geq D(P_B \| Q_B)$ $\color{red}{\left(\begin{array}{l}\text{by data processing} \\ \text{inequality}\end{array}\right)}$

Denote $\quad \mathbb{E}_P[B] = p, \quad \mathbb{E}_Q[B] = q$.

Then, the right-side above equal

$$p \log p + (1-p) \log (1-p) + p \log \frac{1}{q} + (1-p) \log \frac{1}{1-q}$$

$$\geq p \log \frac{1}{q} - h(p)$$

Note that $p = P(U = d(X))$ and

$$q = Q(U = d(X)) = \frac{1}{M} \sum_{m=1}^{M} Q_X(D_m)$$

$$\leq \frac{1}{M}$$

$\Rightarrow \quad I(U \wedge X) \geq (1 - P_e(d)) \log M - h(p)$

$\Longleftrightarrow \quad P_e(d) \geq 1 - \dfrac{I(U \wedge X) + h(p)}{\log M}$

$$\geq 1 - \frac{I(U \wedge X) + 1}{\log M},$$

which completes the proof for deterministic $d$.

For a randomized $d$, note that $\exists \, v$ s.t.

$$P_e(d) = \underset{\downarrow}{\mathbb{E}_V} \left[ P_e(d_v) \right] \geq P_e(d_v)$$

$$\text{randomization using } V \qquad \hookrightarrow \text{deterministic rule.} \qquad \blacksquare$$