

Lec 11 : Universal compression and redundancy

Yanjun Han



Recall from Lec 1: If $X^n \sim P$, then \exists a uniquely decodable code $f: X^n \rightarrow \{0,1\}^*$ s.t.
 $E[\ell(f(X^n))] < H(P) + 1$ bit

- Examples: Shannon / Huffman / Arithmetic codes
- Issue: All these codes require the perfect knowledge of P

Question: Is there a universal code that:

- ① is uniquely decodable ;
- ② does not require the knowledge of P ;
- ③ achieve an expected codelength close to $H(P)$ for every P in a given class \mathcal{P} ?

A motivating example: i.i.d. Bernoulli. Suppose $X_1, \dots, X_n \stackrel{i.i.d.}{\sim} \text{Bern}(p)$, with unknown $p \in [0,1]$.
The Shannon limit is $nH(\text{Bern}(p)) = n(p \log \frac{1}{p} + (1-p) \log \frac{1}{1-p})$.

A simple universal code:

- Compute $n_1 = \sum_{i=1}^n \mathbb{1}(X_i = 1)$, the number of 1's ;
- express n_1 using $\log(n+1)$ bits ; (as $n_1 \in \{0, 1, \dots, n\}$)
- condition on n_1 , there are $\binom{n}{n_1}$ possibilities for (X_1, \dots, X_n) , so we can encode the final sequence using $\log(\binom{n}{n_1})$ bits.

Clearly this is uniquely decodable: the decoder first decodes n_1 from the $\log(n+1)$ bits, and then decodes (X_1, \dots, X_n) from n_1 and the last $\log(\binom{n}{n_1})$ bits.

If $X_1, \dots, X_n \stackrel{i.i.d.}{\sim} \text{Bern}(p)$, the expected codelength is

$$\begin{aligned} E[\ell(f(X^n))] &= \log(n+1) + E[\log(\binom{n}{n_1})] \\ &\leq \log(n+1) + nE[H(\text{Bern}(\frac{n_1}{n}))] \quad (\text{Using } H(X^n) \leq \sum H(X_i)) \\ &\leq \log(n+1) + nH(\mathbb{E} \text{Bern}(\frac{n_1}{n})) \quad (P \mapsto H(P) \text{ is concave}) \\ &= \log(n+1) + nH(\text{Bern}(p)). \end{aligned}$$

In other words, compared with the Shannon limit with the knowledge of p , there is a universal code with an extra overhead of only $O(\log n)$ bits!

General scenario: $X^n \sim P$, where $P \in \mathcal{P}$ is an unknown source distribution.

- Shannon limit with the knowledge of P : $H(P)$
- Universal code: by Kraft inequality, any uniquely decodable code f can be equivalently represented by a probability distribution Q via

$$Q(x^n) = 2^{-l(f(x^n))}.$$

Therefore, the expected codelength of the code (represented by Q) under $X^n \sim P$ is

$$\mathbb{E}_P[l(f(x^n))] = \mathbb{E}_P\left[\log \frac{1}{Q(x^n)}\right].$$

- The "overhead" of a universal code is

$$\mathbb{E}_P\left[\log \frac{1}{Q(x^n)}\right] - H(P) = \mathbb{E}_P\left[\log \frac{P(x^n)}{Q(x^n)}\right] = D_{KL}(P \parallel Q).$$

Defn (minimax redundancy) The minimax redundancy of a distribution class \mathcal{P} on X^n is defined as

$$\text{Red}(\mathcal{P}) = \inf_{Q_{X^n}} \text{Red}(Q_{X^n}; \mathcal{P}) = \inf_{Q_{X^n}} \sup_{P_{X^n} \in \mathcal{P}} D_{KL}(P_{X^n} \parallel Q_{X^n}).$$

- Q_{X^n} corresponds to a universal code that can achieve an overhead at most $\text{Red}(\mathcal{P})$ bits with respect to every $P_{X^n} \in \mathcal{P}$
- In most cases, $\text{Red}(\mathcal{P}) = o(n)$, and even $\text{Red}(\mathcal{P}) = O(\log n)$.

Bernoulli example continued. How to find good Q_{X^n} when $\mathcal{P} = \{\text{Bern}(p)^{\otimes n} : p \in [0, 1]\}$?

Try the following conditional distributions:

- $Q_{X^t} : \text{Unif}\{0, 1\}$
- $Q_{X_{t+1}|X^t}$: let $n_1(x^t)$, $n_0(x^t)$ be the number of 1's and 0's in x^t , respectively, set

$$Q_{X_{t+1}|X^t}(1) = \frac{n_1(x^t) + 1}{t + 2}, \quad Q_{X_{t+1}|X^t}(0) = \frac{n_0(x^t) + 1}{t + 2}.$$

This is called the add-1 estimator or the Laplace estimator.

$$\text{Then } Q_{X^n}(x^n) = \prod_{t=0}^{n-1} Q_{X_{t+1}|X^t}(x_{t+1}|x^{t+1}) = \frac{(1 \cdot 2 \cdots n_1(x^n))(1 \cdot 2 \cdots n_o(x^n))}{2 \cdot 3 \cdots (n+1)} = \frac{n_1(x^n)! n_o(x^n)!}{(n+1)!}.$$

On the other hand, for any $p \in [0, 1]$,

$$P_{X^n}(x^n) = p^{n_1(x^n)} (1-p)^{n_o(x^n)} \leq \left(\frac{n_1(x^n)}{n}\right)^{n_1(x^n)} \left(\frac{n_o(x^n)}{n}\right)^{n_o(x^n)}.$$

$$\Rightarrow \frac{P_{X^n}(x^n)}{Q_{X^n}(x^n)} = (n+1) \cdot \frac{\frac{n_1(x^n)}{n}^{n_1(x^n)}}{\frac{n_1(x^n)!}{n!}} = O(n) \quad (\text{by Stirling})$$

This means that

$$\text{Red}(Q_{X^n}; P) = \sup_{P_{X^n} \in \mathcal{P}} D_{KL}(P_{X^n} \| Q_{X^n}) = \sup_{P_{X^n}} \mathbb{E}_{P_{X^n}} \left[\log \frac{P_{X^n}}{Q_{X^n}} \right] \leq \log n + O(1).$$

Bernoulli example continued, again. Now let's use

$$Q_{X_{t+1}|X^t}(1) = \frac{n_1(x^n) + \frac{1}{2}}{t+1}, \quad Q_{X_{t+1}|X^t}(0) = \frac{n_o(x^n) + \frac{1}{2}}{t+1}$$

This is called the add- $\frac{1}{2}$ estimator, or Krichevsky-Trofimov estimator.

$$\begin{aligned} \text{In this case, } Q_{X^n}(x^n) &= \prod_{t=0}^{n-1} Q_{X_{t+1}|X^t}(x_{t+1}|x^{t+1}) \\ &= \frac{1}{n!} \left(\frac{1}{2} \cdot \frac{3}{2} \cdots (n_1(x^n) - \frac{1}{2}) \right) \left(\frac{1}{2} \cdot \frac{3}{2} \cdots (n_o(x^n) - \frac{1}{2}) \right) \\ &= \frac{(2n_1(x^n) - 1)!! (2n_o(x^n) - 1)!!}{2^n n!}, \end{aligned}$$

$$\text{and } \frac{P_{X^n}(x^n)}{Q_{X^n}(x^n)} \leq \frac{2^n n! \cdot \frac{n_1(x^n)^{n_1(x^n)} n_o(x^n)^{n_o(x^n)}}{n^n}}{(2n_1(x^n) - 1)!! (2n_o(x^n) - 1)!!} = O(\sqrt{n}). \quad (\text{Stirling})$$

$$\text{Therefore, } \text{Red}(Q_{X^n}; P) \leq \log O(\sqrt{n}) = \frac{1}{2} \log n + O(1).$$

This constant $\frac{1}{2}$ turns out to be tight: $\text{Red}(P) = \frac{1}{2} \log n + O(1)$.

In the previous examples, what we're using is actually $\text{Red}(\mathcal{P}) \leq R^*(\mathcal{P})$.

Defn (worst-case/pointwise redundancy)

$$R^*(\mathcal{P}) = \inf_{Q_{X^n}} \sup_{P_{X^n} \in \mathcal{P}} \sup_{x^n} \log \frac{P_{X^n}(x^n)}{Q_{X^n}(x^n)}.$$

- It's clear that $\text{Red}(\mathcal{P}) \leq R^*(\mathcal{P})$.
- $R^*(\mathcal{P})$ treats x^n as an individual sequence, instead of drawing from a probability distribution.
- In online learning, $R^*(\mathcal{P})$ is also the minimax regret under the log loss $\text{log}(p, x) = \log \frac{1}{p(x)}$.

$$R^*(\mathcal{P}) = \inf_{Q_{X^n}} \sup_{x^n} \left(\sum_{t=1}^n \log(Q_t(\cdot | x^{t-1}), x_t) - \inf_{P_{X^n} \in \mathcal{P}} \sum_{t=1}^n \log(P_t(\cdot | x^{t-1}), x_t) \right).$$

Unlike $\text{Red}(\mathcal{P})$ which can be hard to characterize, $R^*(\mathcal{P})$ has a combinatorial characterization.

Thm. $R^*(\mathcal{P}) = \underbrace{\log \left(\sum_{x^n} \sup_{P_{X^n} \in \mathcal{P}} P_{X^n}(x^n) \right)}$ "Shirkov sum"

Pf. (Upper bound) Let $Z = \sum_{x^n} \sup_{P_{X^n} \in \mathcal{P}} P_{X^n}(x^n)$, and

$$Q_{X^n}^*(x^n) = \frac{1}{Z} \sup_{P_{X^n} \in \mathcal{P}} P_{X^n}(x^n), \quad (\text{normalized maximum likelihood distribution})$$

Then $\sup_{P_{X^n} \in \mathcal{P}} \sup_{x^n} \log \frac{P_{X^n}(x^n)}{Q_{X^n}^*(x^n)} = \log Z \Rightarrow R^*(\mathcal{P}) \leq \log Z$.

(Lower bound) For any Q_{X^n} ,

$$\begin{aligned} \sup_{P_{X^n} \in \mathcal{P}} \sup_{x^n} \log \frac{P_{X^n}(x^n)}{Q_{X^n}(x^n)} &= \sup_{P_{X^n} \in \mathcal{P}} \sup_{x^n} \left(\log \frac{P_{X^n}(x^n)}{Q_{X^n}^*(x^n)} + \log \frac{Q_{X^n}^*(x^n)}{Q_{X^n}(x^n)} \right) \\ &= \log Z + \underbrace{\sup_{x^n} \log \frac{Q_{X^n}^*(x^n)}{Q_{X^n}(x^n)}}_{\geq \sum_{x^n} Q_{X^n}^*(x^n) \log \frac{Q_{X^n}^*(x^n)}{Q_{X^n}(x^n)}} \geq \log Z. \end{aligned}$$

This combinatorial nature of $R^*(P)$ makes it easy to upper bound $\text{Red}(P)$ for non-i.i.d. families P .

Example (Markov chain) Let $P = \{P_{X^n} = p(x_t) \prod_{t=1}^{n-1} M(x_{t+1}|x_t)\}$ be the class of all time-homogeneous (first-order) Markov chains on state space $[k]$.

Claim: $\text{Red}(P) \leq \frac{k(k-1)}{2} \log n + O_k(1)$.

Pf. Apply add- $\frac{1}{2}$ estimator to all transition probabilities:

$$Q_{X_{t+1}|X^t}(i) = \frac{n_j \rightarrow i(x^t) + \frac{1}{2}}{n_j(x^t) + \frac{k}{2}} \quad \text{if } x_t = j,$$

where $n_j \rightarrow i(x^t) = \sum_{s=1}^{t-1} \mathbb{1}(x_s = j, x_{s+1} = i)$,
 $n_j(x^t) = \sum_{s=1}^{t-1} \mathbb{1}(x_s = j)$.

Then for any $x^n \in [k]^n$,

$$\begin{aligned} \frac{P_{X^n}(x^n)}{Q_{X^n}(x^n)} &= \frac{p(x_1)}{1/k} \cdot \frac{k}{\prod_{\substack{t \in [n-1]: \\ x_t = j}} \underbrace{\frac{M(x_{t+1}|j)}{(Q_{X_{t+1}})_{x^t}(x_{t+1}|x^t)}}_{= O(\sqrt{n}) \text{ by the analysis in the i.i.d. model}}} \\ &\leq k \cdot (\sqrt{n})^{k(k-1)} \end{aligned}$$

$$\Rightarrow \text{Red}(P) \leq R^*(P) \leq \log(k \cdot (\sqrt{n})^{k(k-1)}) = \frac{k(k-1)}{2} \log n + O(k^2).$$

□

The same programs can be extended to other processes such as the hidden Markov models; we refer to the book [Gassiat, 2018] for more examples.

Redundancy bounds for i.i.d. families.

Entropic upper bound. By the global Fano proof in Lec 9, we have

$$\text{Thm. } \text{Red}(P^{\otimes n}) \leq \inf_{\varepsilon > 0} (n\varepsilon^2 + \log N_{KL}(P, \varepsilon)).$$

Example. When $P = (P_\theta)_{\theta \in \mathbb{R}^d}$ is a parametric family with d parameters, usually $\log N_{KL}(P, \varepsilon) \sim d \log \frac{1}{\varepsilon}$. Choosing $\varepsilon \sim \sqrt{\frac{d}{n}}$ gives the upper bound
 $\text{Red}(P^{\otimes n}) \leq \frac{d}{2} \log \frac{n}{d} + O(d)$.

Lower bound by Rissanen. We begin with a variational representation of $\text{Red}(P)$.

Redundancy-capacity Thm. For $P = (P_\theta)_{\theta \in \Theta}$,

$$\text{Red}(P) = \sup_{P \in \Delta(\Theta)} I(\theta; X) \text{ where } \theta \sim p, X|\theta \sim P_\theta.$$

(The quantity $\sup_p I(\theta; X)$ is the capacity of the channel $P = (P_\theta)_{\theta \in \Theta}$.)

Pf. The "golden formula" for mutual info (Lec 7) says

$$I(\theta; X) = \inf_{Q_X} \mathbb{E}_{\theta \sim p} [D_{KL}(P_\theta \| Q_X)].$$

$$\text{Then } \sup_p I(\theta; X) = \sup_p \inf_{Q_X} \mathbb{E}_{\theta \sim p} [D_{KL}(P_\theta \| Q_X)]$$

$$= \inf_{Q_X} \sup_p \mathbb{E}_{\theta \sim p} [D_{KL}(P_\theta \| Q_X)] \text{ (minimax thm.)}$$

$$= \inf_{Q_X} \sup_{\theta} D_{KL}(P_\theta \| Q_X) = \text{Red}(P).$$

④

Implication: to lower bound $\text{Red}(P)$, can find a proper prior distribution p such that $I(\theta; X)$ is large when $\theta \sim p$.

Rissanen's program: Find an estimator $\hat{\theta}(x^*)$ s.t. $\sup_{\theta \in \Theta} \mathbb{E}_{\theta} [\|\theta - \hat{\theta}(x^*)\|^2] \leq \varepsilon_n^2$.

Ihm. Let $\Theta \subseteq \mathbb{R}^d$ have a non-empty interior. Then

$$\text{Red}(P^{\Theta}) \geq \log \text{Vol}_d(\Theta) - \frac{d}{2} \log \left(\frac{2\pi e \varepsilon_n^2}{d} \right).$$

Pf. Let $\theta \sim \rho = \text{Unif}(\Theta)$, and $h(\cdot)$ denote the differential entropy on \mathbb{R}^d .

Then

$$\begin{aligned} I(\theta; X^*) &= h(\theta) - h(\theta | X^*) \\ &= \log \text{Vol}_d(\Theta) - h(\theta | X^*), \end{aligned}$$

and

$$\begin{aligned} h(\theta | X^*) &= h(\theta - \hat{\theta}(X^*) | X^*) \\ &\leq h(\theta - \hat{\theta}(X^*)) \quad (\text{conditioning reduces entropy}) \\ &\leq \frac{d}{2} \log(2\pi e) + \frac{1}{2} \log \det \left(\mathbb{E}[(\theta - \hat{\theta}(X^*))(\theta - \hat{\theta}(X^*))^T] \right) \\ &\leq \frac{d}{2} \log(2\pi e) + \frac{d}{2} \log \left(\frac{\mathbb{E}\|\theta - \hat{\theta}(X^*)\|^2}{d} \right) \quad (\det(A) = \prod_i \lambda_i \leq \left(\frac{\sum_i \lambda_i}{d} \right)^d) \\ &= \frac{d}{2} \log \left(\frac{2\pi e \varepsilon_n^2}{d} \right). \end{aligned}$$

\square

Example. In parametric families, usually $\text{Vol}_d(\Theta) = \Omega(\frac{1}{n})^{d/2}$ and $\varepsilon_n^2 = O(\frac{d}{n})$.

Therefore, Rissanen's lower bound gives $\text{Red}(P^{\Theta}) \geq \frac{d}{2} \log \frac{n}{d} - O(d)$.

Lower bounds by Haussler & Opper.

The argument of Haussler & Opper (1997) chooses ρ to be a uniform mixture

$$\rho = \frac{1}{M} \sum_{i=1}^M \delta_{\theta_i}.$$

Lemma. For $X|\theta \sim P_\theta$ and $0 < \lambda \leq 1$:

$$I(\theta; X) \geq -\mathbb{E}_{\theta, X} \left[\log \mathbb{E}_{\theta'} \left(\frac{P_{\theta'}(X)}{P_\theta(X)} \right)^\lambda \right],$$

where θ' is an independent copy of θ . ($\theta' \perp\!\!\!\perp (\theta, X)$)

Pf. Let $f(\lambda) = \text{RHS}$, then $f(1) = I(\theta; X)$. Since CGF is convex, f is concave.

Finally,

$$f'(\lambda) = \mathbb{E}_{\theta, X} [\log P_\theta(X)] - \mathbb{E}_{\theta, X} \underbrace{\frac{\mathbb{E}_{\theta'} [P_{\theta'}(X)^\lambda \log P_{\theta'}(X)]}{\mathbb{E}_{\theta'} [P_{\theta'}(X)^\lambda]}}_{\text{when } \lambda=1: (\cdot) = \int_X \int_{\theta} p(\theta|x) \int_{\theta'} p(\theta') P_\theta(x) \log P_{\theta'}(x) d\theta d\theta' dx}.$$

$$= \int_X \int_{\theta'} p(\theta') P_\theta(x) \log P_{\theta'}(x) d\theta' dx$$

$$= \mathbb{E}_{\theta, X} [\log P_\theta(X)].$$

Therefore, $f'(1) = 0 \Rightarrow f'(\lambda) \geq 0$ by concavity of f
 $\Rightarrow f(\lambda) \leq f(1) = I(\theta; X)$

④

Thm. $\text{Red}(P^{\otimes n}) \geq \sup_{\varepsilon \geq 0} \min \left\{ \frac{n\varepsilon^2}{2}, \log M_H(P, \varepsilon) \right\} - \log 2.$

Pf. Let $P_{\theta_1}, \dots, P_{\theta_M}$ be an ε -packing of P under Hellinger, and $p = \frac{1}{M} \sum_{i=1}^M \delta_{\theta_i}$.

The for $\theta \sim p$,

$$\begin{aligned} I(\theta; X^n) &\geq -\frac{1}{M} \sum_{i=1}^M \mathbb{E}_{P_{\theta_i}^{\otimes n}} \left[\log \left(\frac{1}{M} \sum_{j=1}^M \sqrt{\frac{P_{\theta_i}(X)}{P_{\theta_j}(X)}} \right) \right] \\ &\geq -\frac{1}{M} \sum_{i=1}^M \log \left(\frac{1}{M} \sum_{j=1}^M \mathbb{E}_{P_{\theta_i}^{\otimes n}} \sqrt{\frac{P_{\theta_i}(X)}{P_{\theta_j}(X)}} \right) \quad (x \mapsto -\log x \text{ is convex}) \\ &= -\frac{1}{M} \sum_{i=1}^M \log \left(\frac{1}{M} \sum_{j=1}^M \left(1 - \frac{1}{2} H^2(P_{\theta_i}, P_{\theta_j}) \right) \right) \quad (\text{tensorisation of } H^2) \\ &\geq -\frac{1}{M} \sum_{i=1}^M \log \left(\frac{1}{M} + e^{-\frac{n}{2}\varepsilon^2} \right) \quad (H^2(P_{\theta_i}, P_{\theta_j}) \geq \varepsilon^2 \text{ for } i \neq j) \\ &\geq \min \left\{ M, \frac{n}{2}\varepsilon^2 \right\} - \log 2 \quad \left(\frac{1}{a} + \frac{1}{b} \leq \frac{2}{\min\{a, b\}} \right). \end{aligned}$$

④

Example. In parametric families, usually $\log M_H(P, \varepsilon) \sim d \log \frac{1}{\varepsilon}$, so the Haussler - Upper lower bound is also $\text{Red}_n(P^{\otimes n}) \geq \frac{d}{2} \log \frac{n}{Cd \log n}$.

Relationship between redundancy & prediction risk.

Defn (prediction risk)

$$\text{Risk}_n(P) = \inf_{Q_{X_{n+1}|X^n}} \sup_{P_{X_{n+1}} \in P} \mathbb{E}_{P_{X^n}} [D_{KL}(P_{X_{n+1}|X^n} \| Q_{X_{n+1}|X^n})]$$

("KL risk" for next-symbol prediction)

Mutual info representation. If $P = (P_\theta)_{\theta \in \Theta}$,

$$\text{Risk}_n(P) = \sup_{P \in \Delta(\Theta)} I(\theta; X_{n+1} | X^n).$$

$$\begin{aligned} \text{Pf. } \text{Risk}_n(P) &= \inf_{Q_{X_{n+1}|X^n}} \sup_P \mathbb{E}_{P \sim P} [D_{KL}(P_{X_{n+1}|X^n} \| Q_{X_{n+1}|X^n})] \\ &= \sup_{P} \inf_{Q_{X_{n+1}|X^n}} \mathbb{E}_{P \sim P} [D_{KL}(P_{X_{n+1}|X^n} \| Q_{X_{n+1}|X^n})] \quad (\text{minimax thm.}) \\ &= \sup_P I(\theta; X_{n+1} | X^n). \end{aligned}$$

□

Redundancy - risk inequality: $\text{Red}_n(P) \leq \sum_{t=0}^{n-1} \text{Risk}_t(P)$.

$$\text{Pf. Chain rule: } I(\theta; X^n) = \sum_{t=0}^{n-1} I(\theta; X_n | X^{n-1})$$

$$\Rightarrow \sup_P I(\theta; X^n) \leq \sum_{t=0}^{n-1} \sup_P I(\theta; X_n | X^{n-1}).$$

□

Tightness: For i.i.d. $P^{\otimes n}$ with $P = (P_\theta)_{\theta \in \Theta} \subseteq \mathbb{R}^d$, the MLE $\hat{\theta}_t$ based on X^t asymptotically achieves $\mathbb{E}_{\theta} [D_{KL}(P_\theta \| P_{\hat{\theta}})] \sim \frac{d}{2t}$ by Wilk's Thm.
So $\text{Risk}_t \sim \frac{d}{2t}$, and $\text{Red}_n \sim \frac{d}{2} \log n \sim \sum_{t=1}^{n-1} \text{Risk}_t$.

"Online-to-batch" conversion: if each $P_{X^{n+1}} \in \mathcal{P}$ is stationary, i.e. $P_{X_{t_1}, \dots, X_{t_n}} = P_{X_{t_1+1}, \dots, X_{t_n+1}}$, then

$$\text{Risk}_n(\mathcal{P}) \leq \frac{1}{n} \text{Red}(\mathcal{P}) + \text{Mem}(\mathcal{P}).$$

where the "memory term" is

$$\text{Mem}(\mathcal{P}) = \sup_{P_{X^{n+1}} \in \mathcal{P}} \frac{1}{n} \sum_{t=1}^n I(X_{nt}; X^{n-t} | X_{n-t+1}^n).$$

Pf. Let $Q_{X^{n+1}} = \prod_{t=1}^{n+1} Q_{X_t | X^{t-1}}$ attain the minimax redundancy $\text{Red}(\mathcal{P})$.

Now choose the Yang-Barron type predictor:

$$\tilde{Q}_{X_{n+1}|X^n} = \frac{1}{n} \sum_{t=1}^n Q_{X_{t+1}|X^t} (\cdot | X_{n-t+1}^n).$$

$$\begin{aligned} & \text{Then } \mathbb{E}_{P_{X^n}} [D_{KL}(P_{X_{n+1}}|X^n || \tilde{Q}_{X_{n+1}|X^n})] \\ & \leq \frac{1}{n} \sum_{t=1}^n \mathbb{E}_{P_{X^{n+1}}} \left[\log \frac{P_{X_{n+1}}|X^n(X_{n+1}|X^n)}{Q_{X_{t+1}|X^t}(X_{n+1}|X_{n-t+1}^n)} \right] \quad (\text{convexity of KL}) \\ & = \frac{1}{n} \sum_{t=1}^n \mathbb{E}_{P_{X^{n+1}}} \left[\log \frac{P_{X_{n+1}}|X_{n-t+1}^n(X_{n+1}|X_{n-t+1}^n)}{Q_{X_{t+1}|X^t}(X_{n+1}|X_{n-t+1}^n)} + \log \frac{P_{X_{n+1}}|X^n(X_{n+1}|X^n)}{P_{X_{n+1}}|X_{n-t+1}^n(X_{n+1}|X_{n-t+1}^n)} \right] \\ & = \frac{1}{n} \sum_{t=1}^n \left(\mathbb{E}_{P_{X^{n+1}}} \left[\log \frac{P_{X_{n+1}}|X^t(X_{n+1}|X^t)}{Q_{X_{t+1}}|X^t(X_{n+1}|X^t)} \right] + I(X_{n+1}; X^{n-t} | X_{n-t+1}^n) \right) \quad (\text{stationarity}) \\ & \leq \frac{1}{n} \sum_{t=1}^n \mathbb{E}_{P_{X^t}} [D_{KL}(P_{X_{n+1}}|X^t || Q_{X_{n+1}}|X^t)] + \text{Mem}(\mathcal{P}) \\ & \leq \frac{1}{n} D_{KL}(P_{X^{n+1}} || Q_{X^{n+1}}) + \text{Mem}(\mathcal{P}) \leq \frac{1}{n} \text{Red}(\mathcal{P}) + \text{Mem}(\mathcal{P}) \end{aligned}$$

②

Example (Markov chain prediction) $\mathcal{P} = \{ \text{stationary Markov chains on } [k] \text{ of length } n+1 \}$

$$\text{Then } \text{Red}(\mathcal{P}) = O(k^2 \log n)$$

$$\text{Mem}(\mathcal{P}) = \sup_{P_{X^{n+1}}} \frac{1}{n} I(X_{n+1}; X^n) \leq \frac{\log k}{n}$$

$$\Rightarrow \text{Risk}_n(\mathcal{P}) = O\left(\frac{k^2 \log n}{n}\right).$$

- The main surprising feature is that this upper bound on $\text{Risk}_k(P)$ does not depend on the mixing property of the Markov chain
- A pure statistical proof of this upper bound is unknown without mixing conditions
- This bound is tight for $3 \leq k \ll \sqrt{n}$.

Special topic: characterization of R^* in Gaussian models (Mourtada. 2023)

Gaussian family: $\mathcal{P}_A = \{N(\theta, I_n) : \theta \in A\}$, with $A \subseteq \mathbb{R}^n$.

By the entropic upper bound and Haussler-Dupper lower bound, with

$$D_{KL}(N(\theta, I_n) \| N(\theta', I_n)) = \frac{1}{2} \|\theta - \theta'\|^2$$

$$\int \sqrt{dN(\theta, I_n) dN(\theta', I_n)} = \exp(-\frac{1}{8} \|\theta - \theta'\|^2),$$

we get the following characterization of $\text{Red}(P_A)$:

$$\text{Red}(P_A) \asymp \inf_{r>0} \log N(A, \|\cdot\|_2, r) + r^2$$

The main result of this section is a similar characterization of $R^*(P_A)$:

$$R^*(P_A) \asymp \inf_{r>0} \log N(A, \|\cdot\|_2, r) + w_A(r)$$

where $w_A(r)$ is the local Gaussian width:

$$w_A(r) = \sup_{\theta \in A} w(A \cap B(\theta, r)) = \sup_{\theta \in A} \mathbb{E} \left[\sup_{w \in A \cap B(\theta, r)} \langle w, z \rangle \right], \quad z \sim N(0, I_n)$$

(Alternative representation: let $r_N = \sup \{r > 0 : \log N(A, r, \|\cdot\|_2) \geq r^2\}$
 $r_w = \sup \{r > 0 : w_A(r) \geq r^2\}$

then

$$\text{Red}(P_A) \asymp r_N^2$$

$$R^*(P_A) \asymp r_N^2 + r_w^2.$$

)

Example (ellipsoids) If $A = \{\theta \in \mathbb{R}^n : \sum_{i=1}^n \frac{\theta_i^2}{a_i^2} \leq 1\}$, then

$$Rd(P_A) \approx \inf_{r>0} \left(\sum_{i: a_i > 2r} \log\left(\frac{a_i}{r}\right) + r^2 \right),$$

$$R^*(P_A) \approx \inf_{r>0} \left(\sum_{i=1}^n \log\left(1 + \frac{a_i^2}{r^2}\right) + r^2 \right).$$

The key result in the proof lies in the following lemma.

Lemma. $w(A) - \sup_{\theta \in A} \frac{\|\theta\|^2}{2} \leq R^*(P_A) \leq w(A).$

Pf of thm assuming the lemma,

(Upper bound of R^*) First, observe a simple inequality:

$$R^*\left(\bigcup_{i=1}^N P_i\right) \leq \max_{i \in [N]} R^*(P_i) + \log N.$$

The proof is easy: let Q_i attain $R^*(P_i)$, then $\bar{Q} = \frac{1}{N} \sum_{i=1}^N Q_i$ attains the claimed upper bound of $R^*\left(\bigcup_{i=1}^N P_i\right)$.

To apply this inequality, consider a r -covering $\theta_1, \dots, \theta_N$ of A under $\|\cdot\|_2$.

Then

$$\begin{aligned} R^*(P_A) &\leq \max_{i \in [N]} R^*(P_{A \cap B(\theta_i, r)}) + \log N \\ &\leq \max_{i \in [N]} w(A \cap B(\theta_i, r)) + \log N \quad (\text{by Lemma}) \\ &\leq w_A(r) + \log N. \end{aligned}$$

(Lower bound of R^*) First, $R^*(P_A) \geq Rd(P_A) \gtrsim r_N^2$ by Hansler-Opper.

Second, for $r = r_w$ and $\theta \in A$,

$$R^*(P_A) \geq R^*(P_{A \cap B(\theta, r)})$$

$$\geq w(A \cap B(\theta, r)) - \frac{r^2}{2} \quad (\text{by Lemma \& translation invariance})$$

$$\Rightarrow R^*(P_A) \geq w_A(r) - \frac{r^2}{2} \geq \frac{r^2}{2} \quad \text{for } r = r_w \text{ and defn. of } r_w. \quad \square$$

Next we prove the lemma. First we write out the Shtarkov sum.

$$\text{Lemma. } R^*(P_A) = \log \int_{\mathbb{R}^n} \frac{1}{(2\pi)^{n/2}} e^{-\frac{1}{2}d(x, A)^2} dx,$$

$$\text{where } d(x, A) = \inf_{y \in A} \|x - y\|_2.$$

Pf. Trivially follow from Shtarkov.

Using auxiliary $\mathbf{z} \sim N(\mathbf{0}, I_n)$,

$$\begin{aligned} \int_{\mathbb{R}^n} \frac{1}{(2\pi)^{n/2}} e^{-\frac{1}{2}d(x, A)^2} dx &= \mathbb{E}_z [e^{\frac{1}{2}(z^2 - \text{dist}(z, A)^2)}] \\ &= \mathbb{E}_z [\exp(\sup_{w \in A} (\frac{z^2}{2} - \frac{(z-w)^2}{2}))] \\ &= \mathbb{E}_z [\exp(\sup_{w \in A} \langle w, z \rangle - \frac{\|w\|^2}{2})]. \end{aligned}$$

(Pf of lower bound)

$$\begin{aligned} \log \mathbb{E}_z [\exp(\sup_{w \in A} \langle w, z \rangle - \frac{\|w\|^2}{2})] &\geq \log \mathbb{E}_z [\exp(\sup_{w \in A} \langle w, z \rangle)] - \frac{1}{2} \sup_{w \in A} \|w\|^2 \\ &\geq \mathbb{E}_z [\sup_{w \in A} \langle w, z \rangle] - \frac{1}{2} \sup_{w \in A} \|w\|^2 \quad (\text{Jensen}) \\ &= w(A) - \frac{1}{2} \sup_{w \in A} \|w\|^2. \end{aligned}$$

(Pf of upper bound) Let $v = N(\mathbf{0}, I_n)$, and $f(z) = \sup_{w \in A} \langle w, z \rangle - \frac{\|w\|^2}{2}$.

By Gibbs' variational principle,

$$\begin{aligned} \log \mathbb{E}_z [\exp(-f(z))] &= \sup_{\mu} \mathbb{E}_{z \sim \mu} [f(z)] - D_{KL}(\mu || v) \\ &\leq \sup_{\mu} \mathbb{E}_{z \sim \mu} [f(z)] - \frac{1}{2} W_2^2(\mu, v) \quad (\text{Talagrand T}_2 \text{ ineq.}) \\ &= \sup_{\mu} \mathbb{E}_{z \sim \mu} [f(z)] - \sup_g (\mathbb{E}_{\mu} g + \mathbb{E}_v g^c) \\ &\quad (g^c(z) = \inf_x \frac{\|x - z\|^2}{2} - g(x)) \\ &\leq \mathbb{E}_{z \sim \mu} \left[\sup_x f(x) - \frac{\|x - z\|^2}{2} \right]. \end{aligned}$$

On the other hand,

$$\begin{aligned} \sup_x \left(f(x) - \frac{\|x-z\|^2}{2} \right) &= \sup_{w \in A} \sup_x \left(\langle w, x \rangle - \frac{\|w\|^2}{2} - \frac{\|x-z\|^2}{2} \right) \\ &= \sup_{w \in A} \langle w, z \rangle \\ \Rightarrow \log \mathbb{E}_z e^{\frac{f(z)}{2}} &\leq \mathbb{E}_z [\sup_{w \in A} \langle w, z \rangle] = w(A). \end{aligned}$$

(An alternative proof is (Mourtada, 2023), using convex geometry)

Defn (mixed volume) Let K_1, \dots, K_r be convex bodies in \mathbb{R}^n . Write

$$\text{Vol}_n(\lambda_1 K_1 + \dots + \lambda_r K_r) = \sum_{j_1, \dots, j_n=1}^r V(K_{j_1}, \dots, K_{j_n}) \lambda_{j_1} \dots \lambda_{j_n},$$

the quantity $V(K_{j_1}, \dots, K_{j_n})$ is called the mixed volume.

Defn (intrinsic volume)

$$V_j(K) = \binom{n}{j} \frac{\underbrace{V(K, \dots, K, B, \dots, B)}_{K^{n-j}}}{K^{n-j}},$$

where $K_j = \frac{\pi^{j/2}}{\Gamma(\frac{j}{2}+1)}$ is the volume of the unit ball in \mathbb{R}^j .

$$\text{Thm (Steiner formula)} \quad \text{Vol}_n(K+tB) = \sum_{j=0}^n V_{n-j}(K) K_j t^j.$$

Thm (Alexandrov-Fenchel)

$$V(K_1, K_2, \dots, K_n)^2 \geq V(K_1, K_1, K_3, \dots, K_n) V(K_2, K_2, K_3, \dots, K_n).$$

Corollary. By choosing $(K_1, K_2, \dots, K_n) = (K, B, K, \dots, K, B, \dots, B)$, we get

$$j V_j(K)^2 \geq (j+1) V_{j+1}(K) V_{j-1}(K).$$

In particular,

$$V_j(K) \leq \frac{V_1(K)^j}{j!}.$$

Back to the proof of upper bound: since $R^*(P_A) \leq R^*(P_{\text{conv}(A)})$ and $w(A) = w(\text{conv}(A))$,

WLOG we may assume $A = K$ is convex. Then

$$\begin{aligned} \int_{\mathbb{R}^n} \exp\left(-\frac{1}{2}d(x, K)^2\right) dx &= \int_0^\infty \text{Vol}_n(\{x \in \mathbb{R}^n : e^{-\frac{1}{2}d(x, K)^2} \geq t\}) dt \\ &= \int_0^\infty \text{Vol}_n(\{x \in \mathbb{R}^n : d(x, K) \leq r\}) r e^{-r^2/2} dr \\ &= \int_0^\infty \text{Vol}_n(K + rB) r e^{-r^2/2} dr \\ &= \int_0^\infty \sum_{j=0}^n V_{n-j}(K) k_j r^j \cdot r e^{-r^2/2} dr \\ &= \sum_{j=0}^n V_{n-j}(K) (2\pi)^{j/2} \cdot \left(\int_0^\infty r^{j+1} e^{-r^2/2} dr = 2^{j/2} \Gamma(\frac{j}{2} + 1) \right) \end{aligned}$$

$$\Rightarrow R^*(P_K) = \log \sum_{j=0}^n V_j(K) (2\pi)^{-j/2} = \log \underbrace{\sum_{j=0}^n V_j\left(\frac{K}{\sqrt{2\pi}}\right)}_{\text{called the Wills functional}}$$

Using the corollary,

$$\begin{aligned} R^*(P_K) &\leq \log \sum_{j=0}^n \frac{V_i\left(\frac{K}{\sqrt{2\pi}}\right)^j}{j!} < \log \exp\left(V_i\left(\frac{K}{\sqrt{2\pi}}\right)\right) = V_i\left(\frac{K}{\sqrt{2\pi}}\right) \\ &= w(K). \end{aligned}$$