

8: Information Theoretic Lower Bounds

①

A) Lower bounds for compression and Shannon's Source Coding Theorem

We saw that quantity $L_\varepsilon(X) = \min \{ \lceil \log |S| \rceil : S \subseteq \mathcal{X}, P(S) \geq 1 - \varepsilon \}$ is fundamental to the compression problem. We saw that

$$L_\varepsilon(X) \leq \min \{ \lambda : P(\{x : -\log p(x) \leq \lambda\}) \geq 1 - \varepsilon \}.$$

Using Chebyshev's inequality, we get

$$L_\varepsilon(X) \leq H(X) + \sqrt{\frac{\text{Var}(-\log p(X))}{\varepsilon}}$$

We now use Fano's inequality to show that this upper bound is almost tight. Consider a mapping

$f: \mathcal{X} \rightarrow \{0, 1\}^l$ and $\phi: \{0, 1\}^l \rightarrow \mathcal{X}$ s.t.

$$P(X = \phi(f(X))) \geq 1 - \varepsilon.$$

Choosing $S = \{x : x = \phi(f(x))\}$, $P(S) \geq 1 - \varepsilon$. Further, since $f(x)$ can take at most 2^l values, $|S| \leq 2^l$. Thus, $L_\varepsilon(X) \leq l$.

Also, given a set $S \subseteq \mathcal{X}$ s.t. $P(S) \geq 1 - \varepsilon$, let

$$f(x) = \begin{cases} 1, & \text{if } x \notin S \\ i, & \text{if } x = x_i \in S, \quad 1 \leq i \leq |S|. \end{cases}$$

Thus, $f: \mathcal{X} \rightarrow \{0, 1\}^{\lceil \log |S| \rceil}$, which gives

$$\begin{aligned} \tilde{L}_\varepsilon(X) &:= \min \{ l : \text{there exists } f: \mathcal{X} \rightarrow \{0, 1\}^l \text{ and } \phi: \{0, 1\}^l \rightarrow \mathcal{X}, \\ &\quad P(X = \phi(f(X))) \geq 1 - \varepsilon \} \\ &= L_\varepsilon(X). \end{aligned}$$

We derive a lower bound for $\tilde{L}_\varepsilon(x)$.

(2)

For $f: \mathcal{X} \rightarrow \{0, 1\}^k$, $\phi: \{0, 1\}^k \rightarrow \mathcal{X}$, let $M = f(X)$ and

$\hat{X} = \phi(M)$. Then, $P(X = \hat{X}) \geq 1 - \varepsilon$.

$$\begin{aligned} L &\geq H(M) \geq H(\hat{X}) = H(X, \hat{X}) - H(X|\hat{X}) \\ &\geq H(X) - H(X|\hat{X}) \\ &\geq H(X) - \varepsilon \log(|\mathcal{X}| - 1) - h(\varepsilon), \end{aligned}$$

where we used Fano's inequality in the last step.

Thus, $L_\varepsilon(x) \geq H(x) - \varepsilon \log(|\mathcal{X}| - 1) - h(\varepsilon)$, which shows $L_\varepsilon(x) \approx H(x)$.

Shannon considered an asymptotic version of this problem where the goal is to characterize the number of bits per symbol needed to store $X_1, \dots, X_n \sim \text{iid } P$. Specifically, he looked at the optimal rate R^* given by

$$R^*(P) = \lim_{\varepsilon \rightarrow 0} \lim_{n \rightarrow \infty} \frac{1}{n} L_\varepsilon(X^n).$$

→ Using our upper bound, $L_\varepsilon(X^n) \leq nH(X) + \sqrt{\frac{n}{\varepsilon} \text{Var}(-\log P(X))}$,
giving $\lim_{n \rightarrow \infty} \frac{1}{n} L_\varepsilon(X^n) \leq H(X)$.

→ Using our lower bound, $L_\varepsilon(X^n) \geq nH(X) - n\varepsilon \log|\mathcal{X}| - 1$,
giving $\lim_{n \rightarrow \infty} \frac{1}{n} L_\varepsilon(X^n) \geq H(X) - \varepsilon$.

Taking the limit $\varepsilon \rightarrow 0$, we get:

Theorem (Shannon's Source Coding Theorem)

(3)

$$R^*(P) = H(P).$$

B Lower bound for hypothesis testing and Stein's Lemma

Recall the quantity $\beta_\varepsilon(P, Q) := \min\{Q(S) : P(S) \geq 1 - \varepsilon\}$ which is fundamental for binary hypothesis testing (Neyman-Pearson formulation). We saw that

$$-\log \beta_\varepsilon(P, Q) \geq \max\{\lambda : P(\{x : \log \frac{P(x)}{Q(x)} \geq \lambda\}) \geq 1 - \varepsilon\}$$

which by Chebyshev's inequality gives,

$$-\log \beta_\varepsilon(P, Q) \geq D(P \parallel Q) - \sqrt{\frac{1}{\varepsilon} \text{Var}_P(-\log \frac{P(x)}{Q(x)})}.$$

We now show (using the data processing inequality) that this bound is almost tight. Consider a test $T: \mathcal{X} \rightarrow \{0, 1\}$.

We can even allow a randomized test which can be viewed as a channel. We can interpret the output 0 as P and 1 as Q . Then, $\sum_x P(x) T(0|x) \geq 1 - \varepsilon$. Then, by DPI,

$$D(T_0 P \parallel T_0 Q) \leq D(P \parallel Q). \text{ Further,}$$

$$D(T_0 P \parallel T_0 Q) \geq (T_0 P)(0) \log \frac{1}{(T_0 Q)(0)} - h((T_0 P)(0))$$

$$\geq (1 - \varepsilon) \log \frac{1}{(T_0 Q)(0)} - 1.$$

$$\text{Thus, } -\log (T \circ Q)(0) \leq \frac{D(P \parallel Q) + 1}{1 - \epsilon} \quad (4)$$

Note that $(T \circ Q)(0) = \sum_x Q(x) T(0|x)$ is the prob of missed detection and $(T \circ P)(1) = \sum_x P(x) T(1|x)$ is the prob. of false alarm.

$$\text{Thus, } -\log \beta_\epsilon(P, Q) \leq \frac{D(P \parallel Q) + 1}{1 - \epsilon}$$

That is, the best prob. of missed detection $\approx 2^{-D(P \parallel Q)}$

Stein's Lemma

$$\lim_{\epsilon \rightarrow 0} \lim_{n \rightarrow \infty} \frac{1}{n} -\log \beta_\epsilon(P^n \parallel Q^n) = D(P \parallel Q).$$

(i.e., $\beta_\epsilon(P^n \parallel Q^n) \approx 2^{-n D(P \parallel Q)}$ for iid observations.

To make this small, we must have $n \gtrsim \frac{1}{D(P \parallel Q)}$.

C Lower bound for randomness generation

Consider the uniformity generation problem. Suppose a mapping $f: \mathcal{X} \rightarrow \{0, 1\}^l$ satisfies

$$d(P_{f(x)}, \text{unif}(\{0, 1\}^l)) \leq \epsilon,$$

where $x \sim P$, then by the continuity bound for entropy,

$$\begin{aligned} H(f(x)) &\geq H(\text{unif}(\{0, 1\}^l)) - \epsilon \log |\mathcal{X}| - h(\epsilon) \\ &= l - \epsilon \log |\mathcal{X}| - 1 \end{aligned}$$

$$\Rightarrow l \leq H(x) + \epsilon \log |\mathcal{X}| + 1. \quad (\text{no more } \approx H(x) \text{ random bits can be generated})$$

D) Strong Converse for Shannon's Source Coding Theorem ⁽⁵⁾

We saw earlier that

$$\lim_{n \rightarrow \infty} \frac{1}{n} L_{\varepsilon}(x^n) \leq H(X)$$

and

$$\lim_{n \rightarrow \infty} \frac{1}{n} L_{\varepsilon}(x^n) \geq H(X) - \varepsilon \log |\mathcal{X}|.$$

In fact, it can be shown that

$$\lim_{n \rightarrow \infty} \frac{1}{n} L_{\varepsilon}(x^n) = H(X), \quad \text{for every } 0 < \varepsilon < 1.$$

Thus $\lim_{n \rightarrow \infty} \frac{1}{n} L_{\varepsilon}(x^n) = R^* \quad \forall 0 < \varepsilon < 1$. Such

results are called strong converse results. We show this now.

Consider a λ s.t. $P(\{x: -\log p(x) \geq \lambda\}) \geq 1 - \delta$
(we used $-\log p(x) \leq \lambda$ for the upper bound)

Denoting $B_{\lambda} = \{x: -\log p(x) \geq \lambda\}$, for any $S \subseteq \mathcal{X}$ s.t.

$P(S) \geq 1 - \varepsilon$, we have $P(B \cap S) \geq P(B) + P(S) - 1 \geq 1 - \varepsilon - \delta$.

Thus, $1 - \varepsilon - \delta \leq \sum_{x \in B \cap S} p(x) \leq \sum_{x \in B \cap S} 2^{-\lambda} \leq 2^{-\lambda} |B \cap S|$,

since $p(x) \leq 2^{-\lambda}$ for every $x \in B$. We have,

$$\log |S| \geq \log |B \cap S| \geq \lambda - \log \frac{1}{1 - \varepsilon - \delta}.$$

⑥

We have shown:

$$L_{\varepsilon}(X) \geq \max \left\{ \lambda : P(\{x: -\log p(x) \geq \lambda\}) \geq \frac{1-\varepsilon}{2} \right\} \\ - \log \frac{2}{1-\varepsilon}.$$

By using this bound with Chebyshev's inequality,

$$\frac{1}{n} L_{\varepsilon}(X^n) \geq H(X) - \sqrt{\frac{1}{n\varepsilon} \text{Var}(-\log P(X))} \\ - \frac{1}{n} \log \frac{2}{1-\varepsilon}.$$

$$\Rightarrow \lim_{n \rightarrow \infty} \frac{1}{n} L_{\varepsilon}(X^n) \geq H(X).$$

E Minmax lower bound for estimation

The final lower bound we present is very useful in statistics.

We present it using an example. Consider the problem

of estimating the mean of a Gaussian random variable.

Let $X_1, \dots, X_n \sim \text{iid } N(\mu, I)$. The empirical mean

$$\hat{\mu} = \frac{1}{n} \sum_{t=1}^n X_t \text{ satisfies}$$

$$\mathbb{E}[\|\mu - \hat{\mu}\|_2^2] = \sum_{i=1}^d \mathbb{E}[(\mu_i - \hat{\mu}_i)^2] \\ = \sum_{i=1}^d \frac{1}{n^2} \sum_{t=1}^n \text{Var}(X_t) = \frac{d}{n}.$$

Thus, we found an estimator with MSE $\frac{d}{n}$.

Can we do better? Not in the worst-case.

(7)

For any estimator $\hat{\mu} : \mathcal{X}^n \rightarrow \mathbb{R}^d$, we have

$$\begin{aligned} \max_{\mu} \mathbb{E}_{\rho_{\mu}}[\|\mu - \hat{\mu}\|_2^2] &\geq \max_{\mu \in \{\mu_1, \dots, \mu_M\}} \mathbb{E}_{\rho_{\mu}}[\|\mu - \hat{\mu}\|_2^2] \\ &\geq \frac{1}{M} \sum_{i=1}^M \mathbb{E}_{\rho_{\mu_i}}[\|\mu_i - \hat{\mu}\|_2^2] \\ &\geq \delta^2 \cdot \left(\frac{1}{M} \sum_{i=1}^M P_{\mu_i}(\|\mu_i - \hat{\mu}\|_2^2 > \delta^2) \right) \end{aligned}$$

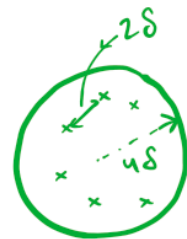
(This is a different application of Markov's inequality)

This bound holds for any selection of μ_1, \dots, μ_M . The expression on the right-side is probability of error for an M -ary hypothesis testing problem. We choose μ_1, \dots, μ_M to get a very difficult hypothesis testing problem.

Specifically, we choose μ_1, \dots, μ_M s.t.

$$(i) \|\mu_i\|_2 \leq 4\delta, \quad 1 \leq i \leq M;$$

$$(ii) \|\mu_i - \mu_j\|_2 \geq 2\delta, \quad i \neq j.$$



We can use our estimator $\hat{\mu}$ to form a test for our M -ary hypothesis testing problem. Specifically, we find the unique i s.t. $\|\hat{\mu} - \mu_i\|_2 < \delta$. Note that by (ii) we can only find one such i . Let $J \sim \text{unif}\{1, \dots, M\}$.

We generate $X_1, \dots, X_n \sim \text{iid } \mathcal{N}(\mu_J, I)$ and find $\hat{J} = i$ s.t. $\|\hat{\mu} - \mu_i\|_2 < \delta$, we get

$$\frac{1}{M} \sum_{i=1}^M P_{\mu_i}(\|\mu_i - \hat{\mu}\|_2 > \delta) \geq P(J \neq \hat{J}) \quad (8)$$

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

where the final bound is by Fano's inequality.

Note that $I(J \wedge X^n) \leq \max_{P_J} I(J \wedge X^n)$

(using the information-radius interpretation of capacity)

$$\begin{aligned} &\leq \max_{1 \leq i \neq j \leq M} D(P_{\mu_i}^n \| P_{\mu_j}^n) \\ &= n \max_{1 \leq i \neq j \leq M} \frac{\|\mu_i - \mu_j\|_2^2}{2} \\ &\leq 8n\delta^2 \end{aligned}$$

Thus, the overall bound gives

$$\max_{\mu} \mathbb{E}_{P_{\mu}}[\|\mu - \hat{\mu}\|_2^2] \geq \delta^2 \cdot \left(1 - \frac{8n\delta^2 + 1}{\log M}\right)$$

We want M to be as large as possible. We can simply choose these points in a greedy manner by selecting a point and removing the ball of radius 2δ around it.

Note that $\text{vol}(B_d(r)) = c_d r^d$. Thus, we can find at least $M \geq \frac{c_d (4\delta)^d}{c_d (2\delta)^d}$ points, which gives

$\log M \geq d \log 2 = d$. With this choice, we have

$$\max_{\mu} \mathbb{E}_{\rho_{\mu}} (\|\mu - \hat{\mu}\|_2^2) \geq \sigma^2 \left(1 - \frac{8n\sigma^2}{d} - \frac{1}{d}\right) \quad (9)$$

Choosing $\sigma^2 = \frac{d}{16n}$, the right-side is bounded by

$$\geq \frac{d}{16n} \left(\frac{1}{2} - \frac{1}{d}\right) \geq \frac{d}{(6 \times 16)n} \text{ if } d \geq 3.$$

$$\Rightarrow \boxed{\text{MSE} \geq \frac{d}{n}}$$