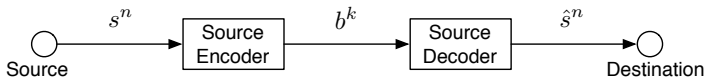# Lossy Source Coding

I-Hsiang Wang

Department of Electrical Engineering
National Taiwan University

ihwang@ntu.edu.tw

November 16, 2023

# The Fixed-to-Fixed Source Coding Problem



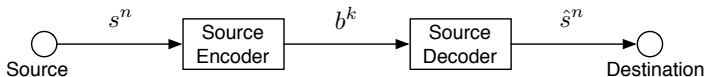Lecture 1: the fundamental limit of fixed-to-fixed lossless source coding.

- The recovery criterion is vanishing probability of error :

$$\lim_{n \to \infty} \Pr\left\{ S^n \neq \hat{S}^n \right\} = 0.$$

- The minimum compression ratio for lossless reconstruction is the entropy of the source:

$$\mathrm{R}^* = \mathrm{H}(S), \quad \text{for DMS } S \sim \mathsf{P}_S.$$

# The Fixed-to-Fixed Source Coding Problem



Turn to fixed-to-fixed lossy source coding. The setting is the same, except:

- Recovery criterion: reconstruction to within a given distortion level D :

$$\limsup_{n \to \infty} \mathsf{E}\left[ d\left( S^n, \hat{S}^n \right) \right] \le \mathsf{D}.$$

- The minimum compression ratio for recovery to within a given distortion level D is the rate-distortion function:

$$\mathrm{R}(\mathsf{D}) = \min_{\mathsf{P}_{\hat{S}|S}:\, \mathsf{E}\left[ d(S, \hat{S}) \right] \le \mathsf{D}} \mathrm{I}\left( S; \hat{S} \right), \quad \text{for DMS } S \sim \mathsf{P}_S.$$

**Why lossy source coding?**

- Sometimes it might be too expensive to reconstruct the source in a lossless way.

- Sometimes it is *impossible* to reconstruct the source losslessly.

  For example, if the source is continuous-valued, the entropy of the source is usually infinite!

- Lossy source coding has wide range of applications, including quantization/digitization of continuous-valued signals, image/video/audio compression, etc.

We first focus on discrete memoryless sources (DMS), and tell how to extend to the continuous-source case later.

# Lossless vs. Lossy source coding

The general lossy source coding problem involves quantizing all possible source sequences $s^n \in \mathcal{S}^n$ into $2^k$ reconstruction sequences $\hat{s}^n \in \hat{\mathcal{S}}^n$, which can be represented by $k$ bits.

The goal is to design the mapping from $s^n$ to $\hat{s}^n$ so that the distortion (quantization error) is below a prescribed level D.

**Lossy source coding has a couple of notable differences from lossless source coding**:

- Source alphabet $\mathcal{S}$ and the reconstruction alphabet $\hat{\mathcal{S}}$ could be **different** in general.
- Performance is determined by the chosen **distortion measure**.

**1** Lossy Source Coding Theorem for Memoryless Sources

**2** Proof of the Coding Theorem

# Distortion measures

Let us begin with the definition of the distortion measure per symbol.

## Definition 1 (Distortion Measure)

A per-symbol distortion measure is a mapping $d(s, \hat{s})$ that maps from $\mathcal{S} \times \hat{\mathcal{S}}$ to $[0, \infty)$, and it is understood as the cost of representing $s$ by $\hat{s}$.
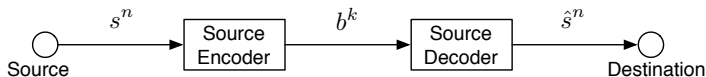
For two length $n$ sequences $s^n$ and $\hat{s}^n$, the distortion between them is defined as the average of the per-symbol distortion:

$$d(s^n, \hat{s}^n) := \frac{1}{n} \sum_{i=1}^{n} d(s_i, \hat{s}_i).$$

**Examples**: below are two widely used distortion measures:

- Hamming distortion: $\mathcal{S} = \hat{\mathcal{S}}$, $d(s, \hat{s}) := \mathbb{1}\{s \neq \hat{s}\}$.
- Squared-error distortion: $\mathcal{S} = \hat{\mathcal{S}} = \mathbb{R}$, $d(s, \hat{s}) := (s - \hat{s})^2$.

# Lossy source coding: problem setup



1. A $(n, \lfloor nR \rfloor)$ source code consists of
   - $\mathrm{enc}_n : \mathcal{S}^n \to \{0,1\}^k$ that maps each source sequence $s^n$ to a bit sequence $b^k$, $k := \lfloor nR \rfloor$.
   - $\mathrm{dec}_n : \{0,1\}^k \to \hat{\mathcal{S}}^n$ that maps each bit sequence $b^k$ to a reconstructed source sequence $\hat{s}^n$.

2. The expected distortion of the code $D^{(n)} := E[d(S^n, \hat{S}^n)]$.

3. A rate-distortion pair $(R, D)$ is achievable if $\exists$ a sequence of $(n, \lfloor nR \rfloor)$ codes such that $\limsup_{n \to \infty} D^{(n)} \le D$.

   The rate distortion function $R(D) := \inf \{R \mid (R, D) : \text{achievable}\}$.

# Landscape of the rate distortion function

$\boxed{D_{\min} := \min_{\hat{s}(s)} \mathsf{E}[d(S, \hat{s}(S))]}$
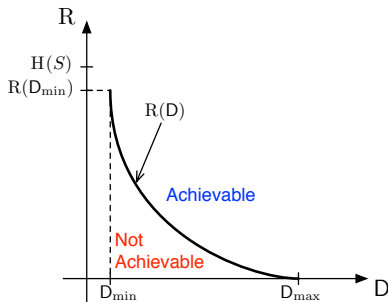
The minimum possible target distortion with finite rate.

Even the decoder knows $s^n$ and finds a best representative $\hat{s}^n$, the expected distortion is still $D_{\min}$.

$\boxed{D_{\max} := \min_{\hat{s}} \mathsf{E}[d(S, \hat{s})]}$



Let $\hat{s}^* := \arg\min_{\hat{s}} \mathsf{E}[d(S, \hat{s})]$. For target distortion $D \geq D_{\max}$, we can use a single representative $\hat{\boldsymbol{s}}^* := (\hat{s}^*, \hat{s}^*, \ldots, \hat{s}^*)$ to reconstruct all $s^n \in \mathcal{S}^n$ (rate is $0$!), and the expected distortion will be

$$\mathsf{E}\left[d(S^n, \hat{\boldsymbol{s}}^*)\right] = \frac{1}{n} \sum_{i=1}^n \mathsf{E}\left[d(S_i, \hat{s}^*)\right] = D_{\max} \leq D.$$

Hence, $R(D) = 0$ for all $D \geq D_{\max}$.

# Lossy source coding theorem



---

**Theorem 1 (A Lossy Source Coding Theorem for DMS)**

*For a DMS $S$, let $\mathrm{R}^{\mathrm{I}}(\mathsf{D}) := \min_{\mathsf{P}_{\hat{S}|S} : \, \mathsf{E}[d(S,\hat{S})] \leq \mathsf{D}} \mathrm{I}\left(S; \hat{S}\right)$. Then,*

$$\mathrm{R}(\mathsf{D}) = \mathrm{R}^{\mathrm{I}}(\mathsf{D}). \tag{1}$$

**Interpretation**:

$$\mathrm{H}(S) \quad - \quad \mathrm{H}\left(S \big| \hat{S}\right) \quad = \quad \mathrm{I}\left(S; \hat{S}\right)$$

Uncertainty of source $S$ — Uncertainty of $S$ after learning $\hat{S}$ = The rate used in compressing $S$ to $\hat{S}$

# Properties of $R^I(D)$



$R^I(D)$ satisfies the following properties:

1. Nonnegative.

2. Non-increasing in D.

3. Convex in D.

4. Continuous in D.

5. $R^I(D_{\min}) \leq H(S)$.

6. $R^I(D) = 0$ if $D \geq D_{\max}$.

These properties are all quite intuitive. Nonnegativity and monotonicity are clear from the definition. Below we sketch the proof of other properties.

**Convexity**     Goal: show that $D_1, D_2 \geq D_{\min}$ and $\lambda \in [0,1], \overline{\lambda} := 1 - \lambda$,

$$R^I(\lambda D_1 + \overline{\lambda} D_2) \leq \lambda R^I(D_1) + \overline{\lambda} R^I(D_2).$$

- Let $Q_i(\hat{s}|s) := \arg\min_{P_{\hat{S}|S}: \, E[d(S,\hat{S})] \leq D_i} I\left(S; \hat{S}\right)$, the optimizing conditional distribution that achieves distortion $D_i$, $i = 1, 2$. Let $Q_\lambda := \lambda Q_1 + \overline{\lambda} Q_2$.

- Under $Q_\lambda(\hat{s}|s)$, the expected distortion between $S$ and $\hat{S}$ is not greater than $\lambda D_1 + \overline{\lambda} D_2$, because

$$E_{P_S Q_\lambda(\hat{s}|s)}\left[d(S,\hat{S})\right] = \sum_s \sum_{\hat{s}} P_S(s) \left[\lambda Q_1(\hat{s}|s) + \overline{\lambda} Q_2(\hat{s}|s)\right] d(s,\hat{s}).$$

- Proof is complete since $I\left(S; \hat{S}\right)$ is convex in $P_{S|\hat{S}}$ with a fixed $P_S$:

$$R^I\left(\lambda D_1 + \overline{\lambda} D_2\right) \leq I\left(S; \hat{S}\right)_{Q_\lambda} \leq \lambda \underbrace{I\left(S; \hat{S}\right)_{Q_1}}_{R^I(D_1)} + \overline{\lambda} \underbrace{I\left(S; \hat{S}\right)_{Q_2}}_{R^I(D_2)}.$$

**Continuity**    Note that convexity within an open interval implies continuity within that open interval.



Hence, it suffices to check the continuity at the boundary $D = D_{\min}$.

Since it is at the boundary, it suffices to show its *right continuity* at $D = D_{\min}$. The proof is technical. Details are in the next slide.

At the boundary $D = D_{\min}$, due to the monotonicity and convexity, it is obvious that $\lim_{D \downarrow D_{\min}} R^I(D)$ exists and $\lim_{D \downarrow D_{\min}} R^I(D) \leq R^I(D_{\min})$.

It remains to show that $\lim_{D \downarrow D_{\min}} R^I(D) \geq R^I(D_{\min})$.

To do so, let us consider a sequence $\{D_m\}$ that converges to $D_{\min}$ from above. Let $Q_m(\hat{s}|s)$ be the conditional PMF that achieves $R^I(D_m)$. Since the set of conditional PMF's is compact, there exists a subsequence $\{m_\ell\}$ such that $\lim_{\ell \to \infty} Q_{m_\ell} := \tilde{Q}$ exists.

Hence, by the continuity of expectation with respect to the distribution, the expected distortion achieved by $Q_{m_\ell}$, which is not greater than $D_{m_\ell}$, converges to that achieved by $\tilde{Q}$, as $\ell \to \infty$. As a result, the expected distortion achieved by $\tilde{Q}$ is not greater than $\lim_{\ell \to \infty} D_{m_\ell} = D_{\min}$.

By the definition of $R^I(\cdot)$,

$$R^I(D_{\min}) \leq I\left(\hat{S}; S\right)_{\tilde{Q}} = \lim_{\ell \to \infty} I\left(\hat{S}; S\right)_{Q_{m_\ell}} = \lim_{\ell \to \infty} R^I(D_{m_\ell}) = \lim_{D \downarrow D_{\min}} R^I(D).$$

This completes the proof.

# Example: Bernoulli source with Hamming distortion

Let us give a first example about how to compute the rate distortion function, that is, how to solve (2) in the lossy source coding theorem.

**Source** (binary) $\qquad\qquad\qquad S_i \in \mathcal{S} = \{0, 1\}$, and $S_i \overset{\text{i.i.d.}}{\sim} \text{Ber}(p) \; \forall \, i$.

**Distortion** (Hamming) $\qquad\qquad d(s, \hat{s}) = \mathbb{1}\{s \neq \hat{s}\}$.

---

### Example 1

Show that the rate distortion function of the Bernoulli $p$ source with Hamming distortion is

$$
\mathrm{R}(\mathrm{D}) = \begin{cases} \mathsf{H}_{\mathsf{b}}(p) - \mathsf{H}_{\mathsf{b}}(\mathrm{D}), & 0 \leq \mathrm{D} \leq \min\{p, 1-p\} \\ 0, & \mathrm{D} > \min\{p, 1-p\} \end{cases}.
$$

---

**sol.** **The first step is to identify** $D_{\min}$ **and** $D_{\max}$.

$$D_{\min} = 0 \qquad \text{because one can choose } \hat{s}(s) = s.$$

$$D_{\max} = \min\{p, 1-p\} \qquad \text{because one can choose } \hat{s} = \begin{cases} 0 & p \leq \frac{1}{2} \\ 1 & p \geq \frac{1}{2} \end{cases}.$$

**The next step is to lower bound** $\mathrm{I}\left(S; \hat{S}\right) = \mathrm{H}(S) - \mathrm{H}\left(S \middle| \hat{S}\right)$. **It is equivalent to upper bounding** $\mathrm{H}\left(S \middle| \hat{S}\right)$:

$$\mathrm{H}\left(S \middle| \hat{S}\right) = \mathrm{H}\left(S \oplus \hat{S} \middle| \hat{S}\right) \leq \mathrm{H}\left(S \oplus \hat{S}\right) = \mathsf{H}_{\mathsf{b}}(q),$$

where we assume $S \oplus \hat{S} \sim \mathrm{Ber}(q)$ for some $q \in [0,1]$. Note: $d(S, \hat{S}) \equiv S \oplus \hat{S}$. Hence, $\mathsf{E}[d(S, \hat{S})] \leq \mathsf{D} \implies q \leq \mathsf{D}$. Since $\mathsf{D} \leq \mathsf{D}_{\max} \leq \frac{1}{2}$, we see that $\mathsf{H}_{\mathsf{b}}(q)$ is maximized when $q = \mathsf{D}$.

**Hence,** $\mathrm{I}\left(S; \hat{S}\right) \geq \mathsf{H}_{\mathsf{b}}(p) - \mathsf{H}_{\mathsf{b}}(\mathsf{D})$.

**Final step: show that the lower bound** $H_b(p) - H_b(D)$ **can be attained.**

Try to find a conditional PMF $P_{\hat{S}|S}$ such that the following holds:

$$\hat{S} \perp\!\!\!\perp S \oplus \hat{S} \quad \text{so that } H\left(S \oplus \hat{S}\middle|\hat{S}\right) = H\left(S \oplus \hat{S}\right) \quad \text{and} \quad \Pr\{S \oplus \hat{S} = 1\} = D.$$

**The difficulty can be resolved via an auxiliary reverse channel.**

Consider a channel with input $\hat{S}$, output $S$, additive noise $Z \sim \text{Ber}(D) \perp\!\!\!\perp \hat{S}$. $S = \hat{S} \oplus Z \iff Z = S \oplus \hat{S}$. The reverse channel specifies the joint distribution $P_{S,\hat{S}}$ and hence $P_{\hat{S}|S}$!



$$p = (1 - \alpha)D + \alpha(1 - D) \implies \alpha = \frac{p - D}{1 - 2D}$$

# Example: source alphabet $\neq$ reconstruction alphabet

**Source** (ternary)      $S_i \in \mathcal{S} = \{0, *, 1\}$, and $S_i \overset{\text{i.i.d.}}{\sim} \mathsf{P}_S \; \forall i$, where $\mathsf{P}_S(0) = \mathsf{P}_S(1) = \varepsilon \leq \frac{1}{2}$.

**Reconstruction** (binary)     $\hat{\mathcal{S}} = \{0, 1\}$.

**Distortion**            $d(s, \hat{s}) = \begin{cases} 1 & \text{if } s \neq * \text{ and } s \neq \hat{s} \\ 0 & \text{if } s = * \text{ or } \hat{s} \end{cases}$.

In other words, there is a <span style="color:green">don't-care</span> symbol $*$, and $\mathcal{S} \neq \hat{\mathcal{S}}$.

---

### Exercise 1

Derive the rate distortion function and show that it is given by

$$\mathrm{R}(\mathsf{D}) = \begin{cases} 2\varepsilon \left(1 - \mathsf{H}_\mathsf{b}\left(\frac{\mathsf{D}}{2\varepsilon}\right)\right), & 0 \leq \mathsf{D} \leq \varepsilon \\ 0, & \mathsf{D} > \varepsilon \end{cases}.$$

# Lossy source coding for continuous-valued sources



## Theorem 2 (A Lossy Source Coding Theorem for CMS)

*For a continuous memoryless source $S$ with PDF $f_S$,*

$$\mathrm{R}(\mathsf{D}) = \inf_{f_{\hat{S}|S}:\; \mathsf{E}\left[d(S,\hat{S})\right] \leq \mathsf{D}} \mathrm{I}\left(S;\hat{S}\right). \tag{2}$$

**Remark**: One can use weak typicality or the discretization method used in channel coding to extend the lossy source coding theorem from discrete memoryless sources to continuous.

# Gaussian source with squared error distortion

**Source** (Gaussian)               $\mathcal{S} = \mathbb{R}$, and $S \sim \mathrm{N}(\mu, \sigma^2)$.

**Distortion** (Squared Error)      $d(s, \hat{s}) = |s - \hat{s}|^2$.

---

**Example 2 (Rate distortion function of Gaussian source)**

For Gaussian source with squared error distortion (as defined above), the rate distortion function is

$$\mathrm{R}(\mathsf{D}) = \begin{cases} \frac{1}{2} \log \left( \frac{\sigma^2}{\mathsf{D}} \right), & 0 \leq \mathsf{D} \leq \sigma^2 \\ 0, & \mathsf{D} > \sigma^2 \end{cases}.$$

---

**Note**: In particular, note that $\mathrm{R}(0) = \infty$, which is quite intuitive!

**sol**: First step: identify $D_{\min}$ and $D_{\max}$.

$$D_{\min} = 0 \qquad \text{because one can choose } \hat{s}(s) = s.$$
$$D_{\max} = \sigma^2 \qquad \text{because one can choose } \hat{s} = \mu, \text{ the mean of } S.$$

Next step: lower bound $\mathrm{I}\left(S; \hat{S}\right) = \mathrm{h}(S) - \mathrm{h}\left(S \big| \hat{S}\right)$.

It is equivalent to upper bounding $\mathrm{h}\left(S \big| \hat{S}\right)$:

$$\mathrm{h}\left(S \big| \hat{S}\right) = \mathrm{h}\left(S - \hat{S} \big| \hat{S}\right) \leq \mathrm{h}\left(S - \hat{S}\right) \leq \frac{1}{2} \log\left(2\pi e\, \mathsf{D}\right),$$

where the last inequality holds since $\mathsf{Var}\left[S - \hat{S}\right] \leq \mathsf{E}\left[\left|S - \hat{S}\right|^2\right] \leq \mathsf{D}$.

Hence, $\mathrm{I}\left(S; \hat{S}\right) \geq \frac{1}{2} \log\left(2\pi e\, \sigma^2\right) - \frac{1}{2} \log\left(2\pi e\, \mathsf{D}\right) = \frac{1}{2} \log\left(\frac{\sigma^2}{\mathsf{D}}\right)$.

Final step: show that the lower bound $\frac{1}{2} \log \left( \frac{\sigma^2}{D} \right)$ can be attained. The goal is to find a conditional density $f_{\hat{S}|S}$ such that

$$\hat{S} \perp\!\!\!\perp \left( S - \hat{S} \right) \text{ so that } h\left( S - \hat{S} \big| \hat{S} \right) = h\left( S - \hat{S} \right) \text{ and } \left( S - \hat{S} \right) \sim N(0, D).$$

Again, this can be done via an auxiliary reverse channel.

Consider a channel with input $\hat{S}$, output $S$, additive noise $Z \sim N\left( 0, D \right) \perp\!\!\!\perp \hat{S}$.

$$S = \hat{S} + Z \implies Z = S - \hat{S}.$$

The reverse channel specifies the joint distribution $f_{S, \hat{S}}$ and hence $f_{\hat{S}|S}$!

# Gaussian source is the hardest to compress

**Theorem 3 (Gaussian is the Worst-Case Source to Compress)**

*Consider a zero-mean CMS $S$ with density $f_S$ and variance $\sigma^2$. Then, the rate distortion function with squared error distortion is maximized when $S \sim N(0, \sigma^2)$, and*

$$R(D) \leq R^G(D) := \max\left\{0, \tfrac{1}{2}\log\left(\tfrac{\sigma^2}{D}\right)\right\}.$$

**pf:** Note that $R(D) = \min_{f_{\hat{S}|S}: \, E[d(S,\hat{S})]\leq D} I\left(S; \hat{S}\right)$.

Hence, to obtain an upper bound, we simply need to choose some $f_{\hat{S}|S}$ that yields the desired

$$I\left(S; \hat{S}\right) = h\left(\hat{S}\right) - h\left(\hat{S}\Big|S\right) \leq R^G(D).$$

Further note that

- $\mathrm{h}\left(\hat{S}\big|S\right)$ can be computed if $\hat{S} = aS + bZ^{\mathrm{G}}$ for $Z^{\mathrm{G}} \sim \mathrm{N}(0,1) \perp\!\!\!\perp S$.
- $\mathrm{h}\left(\hat{S}\right)$ is upper bounded by $\frac{1}{2}\log(2\pi e\, \mathsf{Var}[\hat{S}]) = \frac{1}{2}\log(2\pi e(a^2\sigma^2 + b^2))$.

How to find the coefficients $a$ and $b$? Reverse engineering:

1. The distortion should be D: $\mathsf{E}[(S - \hat{S})^2] = (1-a)^2\sigma^2 + b^2 \overset{!}{=} \mathsf{D}$.

2. The induced mutual information is upper bounded by $\mathrm{R}^{\mathrm{G}}(\mathsf{D})$:

$$\mathrm{I}\left(S;\hat{S}\right) = \mathrm{h}\left(\hat{S}\right) - \mathrm{h}\left(\hat{S}\big|S\right) \leq \tfrac{1}{2}\log\left(2\pi e(a^2\sigma^2 + b^2)\right) - \tfrac{1}{2}\log(2\pi e\, b^2)$$
$$= \tfrac{1}{2}\log\frac{a^2\sigma^2 + b^2}{b^2} \overset{!}{=} \tfrac{1}{2}\log\frac{\sigma^2}{\mathsf{D}}$$

Solving the above two, we get $a = \frac{\sigma^2 - \mathsf{D}}{\sigma^2}$ and $b = a\sqrt{\frac{\sigma^2\mathsf{D}}{\sigma^2 - \mathsf{D}}} = \sqrt{(\sigma^2 - \mathsf{D})\frac{\mathsf{D}}{\sigma^2}}$, for $\mathsf{D} < \sigma^2$. For $\mathsf{D} \geq \sigma^2$, it is obvious that $\mathrm{R}(\mathsf{D}) = 0$. Proof complete. $\qquad\square$

# Proof of the converse part of Theorem 1

**Goal**: Prove that for any sequence of $(n, \lfloor nR \rfloor)$ source codes with

$$\limsup_{n \to \infty} D^{(n)} \leq D,$$

the rate $R$ must satisfy $R \geq R^I(D)$ (defined in (2)). We begin with similar steps as in lossless source coding.

**pf**: Note that $B^k$ is a r.v. because it is generated by another r.v, $S^n$.

$$nR \geq k \geq H\left(B^k\right) \geq I\left(B^k; \hat{S}^n\right) \overset{(a)}{\geq} I\left(S^n; \hat{S}^n\right)$$

$$\overset{(b)}{=} \sum_{i=1}^n I\left(S_i; \hat{S}^n \Big| S^{i-1}\right) \overset{(c)}{=} \sum_{i=1}^n I\left(S_i; \hat{S}^n, S^{i-1}\right) \geq \sum_{i=1}^n I\left(S_i; \hat{S}_i\right)$$

- (a) is due to $S^n - B^k - \hat{S}^n$ and the data processing inequality.
- (b) is due to Chain Rule. (c) is due to $S_i \perp\!\!\!\perp S^{i-1}$ (memoryless source).

So far, we have not used the condition on distortion yet.

Further working on the inequality:

$$
\begin{aligned}
n\mathrm{R} &\geq \sum_{i=1}^{n} \mathrm{I}\left(S_i; \hat{S}_i\right) \\
&\overset{(d)}{\geq} \sum_{i=1}^{n} \mathrm{R}^{\mathrm{I}}\left(\mathsf{E}[d(S_i, \hat{S}_i)]\right) = n \sum_{i=1}^{n} \tfrac{1}{n}\mathrm{R}^{\mathrm{I}}\left(\mathsf{E}[d(S_i, \hat{S}_i)]\right) \\
&\overset{(e)}{\geq} n\mathrm{R}^{\mathrm{I}}\left(\sum_{i=1}^{n} \tfrac{1}{n}\mathsf{E}[d(S_i, \hat{S}_i)]\right) = n\mathrm{R}^{\mathrm{I}}\left(\mathsf{E}\left[\tfrac{1}{n}\sum_{i=1}^{n} d(S_i, \hat{S}_i)\right]\right) \\
&= n\mathrm{R}^{\mathrm{I}}\left(\mathsf{E}[d(S^n, \hat{S}^n)]\right) = n\mathrm{R}^{\mathrm{I}}(\mathsf{D}^{(n)}).
\end{aligned}
$$

- $(d)$ is due to the definition of $\mathrm{R}^{\mathrm{I}}(\mathsf{D})$ in (2).
- $(e)$ is due to the convexity of $\mathrm{R}^{\mathrm{I}}(\mathsf{D})$ and Jensen's inequality.

Hence, $\mathrm{R} \geq \limsup\limits_{n \to \infty} \mathrm{R}^{\mathrm{I}}\left(\mathsf{D}^{(n)}\right) \overset{(f)}{\geq} \mathrm{R}^{\mathrm{I}}\left(\limsup\limits_{n \to \infty} \mathsf{D}^{(n)}\right) \overset{(g)}{\geq} \mathrm{R}^{\mathrm{I}}(\mathsf{D})$.

Here $(f)$ is due to continuity of $\mathrm{R}^{\mathrm{I}}(\mathsf{D})$.

$(g)$ is due to $\limsup_{n \to \infty} \mathsf{D}^{(n)} \leq \mathsf{D}$ and $\mathrm{R}^{\mathrm{I}}(\mathsf{D})$ is non-increasing. $\qquad\square$

# Remarks

You might note that in the proof of converse in the previous slide, we do not make use of lower bounds on error probability such as Fano's inequality.

This is because in our formulation of the lossy source coding problem, the reconstruction criterion is laid on the **expected distortion**.

Instead of the criterion $\limsup_{n\to\infty} \mathsf{D}^{(n)} \leq \mathsf{D}$ ($\mathsf{D}^{(n)} := \mathsf{E}[d(S^n, \hat{S}^n)]$), we could use a stronger one:

$$\mathsf{P}_{\mathsf{e}}^{(n,\delta)} := \mathsf{Pr}\left\{d(S^n, \hat{S}^n) > \mathsf{D} + \delta\right\}, \ \delta > 0 \qquad \text{(Probability of Error)}$$

$$\lim_{n\to\infty} \mathsf{P}_{\mathsf{e}}^{(n,\delta)} = 0, \ \forall \, \delta > 0 \qquad\qquad \text{(Reconstruction Criterion)}$$

We can then re-define the rate distortion function and re-formulate the lossy source coding problem accordingly.

It turns out that Theorem 1 remains the same!

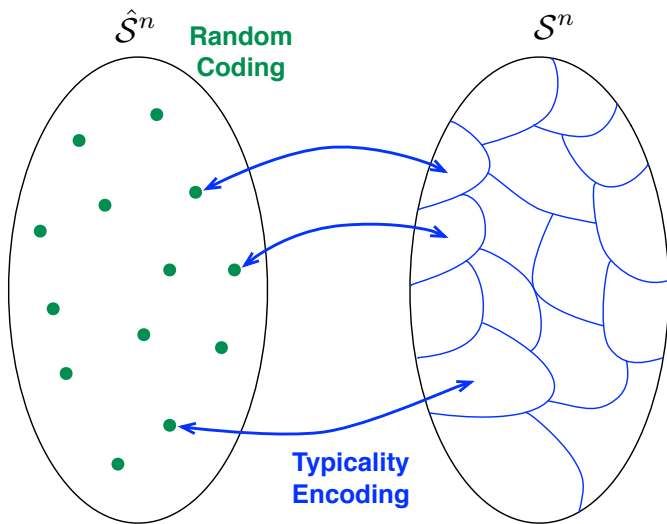# Ideas of constructing a good source code

Key in source coding:

1. Find a good set of **representatives** (quantization codewords).
2. For each source sequence, determine which codeword to be used.

Main tools we use so far in developing achievability of coding theorems:

1. **Random coding**: construct the codebook randomly and show that at least one realization can achieve the desired target performance.
2. **Typicality**: help give bounds in performance analysis.

In the following, we prove the achievability part of Theorem 1 by

1. Random coding – show existence of good quantization codebook.
2. Typicality encoding – determine which codeword to be used.

$\hat{\mathcal{S}}^n$ **Random Coding** $\mathcal{S}^n$

**Typicality Encoding**

# Proof program

**1** **Random Codebook Generation**: Generate a random ensemble of quantization codebooks. Each contains $2^k$ codewords.

**2** **Analysis of Expected Distortion**: Our goal is to show that

$$\limsup_{n\to\infty} \mathsf{E}_{\mathbf{C},S^n}\Big[d(S^n, \hat{S}^n)\Big] \leq \mathsf{D},$$

and conclude that there must exist a codebook $\mathbf{c}$ such that the expected distortion satisfies $\limsup_{n\to\infty} \mathsf{E}_{S^n}\Big[d(S^n, \hat{S}^n)\Big] \leq \mathsf{D}$.

Note that for a source sequence $s^n$, the optimal encoder chooses an index $w \in \{1, ..., 2^k\}$, that is, a codeword $\hat{s}^n(w)$ in the codebook, so that $d(s^n, \hat{s}^n(w))$ is minimized.

However, similar to ML decoding in channel coding, such optimal encoder is hard to analyze. To simplify analysis, we shall introduce a suboptimal encoder based on **typicality**.

# Random codebook generation

Fix the conditional PMF that attains $\mathrm{R}^{\mathrm{I}}\left(\frac{\mathrm{D}}{1+\varepsilon}\right)$:

$$\mathrm{Q}_{\hat{S}|S} = \arg\min_{\mathsf{P}_{\hat{S}|S}: \, \mathsf{E}[d(S,\hat{S})] \leq \frac{\mathrm{D}}{1+\varepsilon}} \mathrm{I}\left(S; \hat{S}\right) \tag{3}$$

Based on the chosen $\mathrm{Q}_{\hat{S}|S}$ and the source distribution $\mathsf{P}_S$, calculate $\mathsf{P}_{\hat{S}}$, the marginal distribution of the reconstruction $\hat{S}$.

Generate $2^k$ codewords $\left\{\hat{s}^n\left(w\right) \middle| w = 1, 2, ..., 2^k\right\}$, i.i.d. according to $\mathsf{P}(\hat{s}^n) = \prod_{i=1}^n \mathsf{P}_{\hat{S}}(\hat{s}_i)$.

In other words, if we think of the quantization codebook as a $2^k \times n$ matrix $\mathbf{C}$, the elements of $\mathbf{C}$ will be i.i.d. distributed according to $\mathsf{P}_{\hat{S}}$.

**Remark**: observe the resemblance with the channel coding achievability.

# Encoding and decoding

**Encoding**: compared to channel coding, the encoder here is more involved.

We use typically encoding: (resembling typically decoding in channel coding)

- Given a source sequence $s^n$, find an index $w \in \{1, 2, ..., 2^k\}$ (need not be unique) such that

$$(s^n, \hat{s}^n(w)) \in \mathcal{T}_{\varepsilon}^{(n)}\left(\mathsf{P}_{S,\hat{S}}\right).$$

  Recall the joint distribution $\mathsf{P}_{S,\hat{S}} = \mathsf{P}_S\,\mathsf{Q}_{\hat{S}|S}$, $\mathsf{Q}_{\hat{S}|S}$ as defined in (3).

- If there is no such index, give up and uniformly randomly pick a $w \in \{1, ..., 2^k\}$.

- Send out the bit sequence that represent the chosen $w$.

**Decoding**: Upon receiving the bit sequence representing $w$, generate the reconstructed $\hat{s}^n(w)$ by looking up the quantization codebook.

# Analysis of expected distortion

Why typicality encoder? Typical average lemma:

> For any nonnegative function $g(x)$ on $\mathcal{X}$, if $x^n \in \mathcal{T}_\varepsilon^{(n)}(\mathsf{P}_X)$, then
> $$(1 - \varepsilon)\, \mathsf{E}_{\mathsf{P}_X}[g(X)] \leq \tfrac{1}{n} \sum_{i=1}^n g(x_i) \leq (1 + \varepsilon)\, \mathsf{E}_{\mathsf{P}_X}[g(X)]\,.$$

To analyze $\mathsf{E}_{\mathbf{C},S^n}[d(S^n, \hat{S}^n)]$, let us distinguish into two cases: $\mathcal{E}$ and $\mathcal{E}^c$, where $\mathcal{E} := \{(S^n, \hat{S}^n) \notin \mathcal{T}_\varepsilon^{(n)}(\mathsf{P}_{S,\hat{S}})\}$, denoting the event that typicality encoding fails.

$$
\begin{aligned}
& \mathsf{E}_{\mathbf{C},S^n}\left[d(S^n, \hat{S}^n)\right] \\
&= \Pr\{\mathcal{E}\}\, \mathsf{E}_{\mathbf{C},S^n}\left[d(S^n, \hat{S}^n)\,\Big|\,\mathcal{E}\right] + \Pr\{\mathcal{E}^c\}\, \mathsf{E}_{\mathbf{C},S^n}\left[d(S^n, \hat{S}^n)\,\Big|\,\mathcal{E}^c\right] \\
&\leq \Pr\{\mathcal{E}\}\, \max_{s,\hat{s}} d(s,\hat{s}) + \Pr\{\mathcal{E}^c\}\,(1+\varepsilon)\tfrac{\mathsf{D}}{1+\varepsilon} \\
&\leq \Pr\{\mathcal{E}\}\, \max_{s,\hat{s}} d(s,\hat{s}) + \mathsf{D}.
\end{aligned}
$$

Hence, as long as $\Pr\{\mathcal{E}\}$ vanishes as $n \to \infty$, we are done.

# Analysis of expected distortion $\rightarrow$ Analysis of $\Pr\{\mathcal{E}\}$

With typicality encoding, the analysis of expected distortion is made easy – **just need to control** $\Pr\{\mathcal{E}\}$!

Look at event $\mathcal{E}$: it is the event that the reconstructed $\hat{S}^n$ is not jointly typical with $S^n$, which is equivalent to the event that *none of the quantization codewords in the codebook is jointly typical with $S^n$*.

Hence, $\mathcal{E} = \bigcap_{w=1}^{2^k} \mathcal{A}_w^c$, where $\mathcal{A}_w := \left\{ \left( S^n, \hat{S}^n(w) \right) \in \mathcal{T}_\varepsilon^{(n)} \right\}$. Therefore,

$$\Pr\{\mathcal{E}\} = \Pr\left\{ \bigcap_{w=1}^{2^k} \mathcal{A}_w^c \right\}.$$

Unfortunately, events $\left\{ \mathcal{A}_w^c \,\middle|\, w = 1, ..., 2^k \right\}$ may not be mutually independent, because they all involve a common random sequence $S^n$.

Nevertheless, for a *fixed* $s^n$, the events $\mathcal{A}_w^c(s^n) := \left\{ \left( s^n, \hat{S}^n(w) \right) \notin \mathcal{T}_\varepsilon^{(n)} \right\}$, $w = 1, ..., 2^k$, are indeed mutually independent!

# Analysis of $\Pr\{\mathcal{E}\}$, $\mathcal{E} := \{(S^n, \hat{S}^n) \notin \mathcal{T}_\varepsilon^{(n)}\}$

Motivated by the above observation, we give an alternative upper bound:

$$\Pr\{\mathcal{E}\} = \sum_{s^n \in \mathcal{S}^n} \mathsf{P}(s^n)\Pr\left\{\bigcap_{w=1}^{2^k} \mathcal{A}_w^c(s^n)\right\}$$

$$= \sum_{s^n \in \mathcal{S}^n} \mathsf{P}(s^n) \prod_{w=1}^{2^k} \Pr\{\mathcal{A}_w^c(s^n)\}$$

$$= \sum_{s^n \in \mathcal{S}^n} \mathsf{P}(s^n) \prod_{w=1}^{2^k} \left(1 - \Pr\{\mathcal{A}_w(s^n)\}\right)$$

**Question**: Is there a way to *lower bound*

$$\Pr\{\mathcal{A}_w(s^n)\} := \Pr\left\{\left(s^n, \hat{S}^n(w)\right) \in \mathcal{T}_\varepsilon^{(n)}(P_{S,\hat{S}})\right\} ?$$

**Yes!** As long as $s^n \in \mathcal{T}_{\varepsilon'}^{(n)}(\mathsf{P}_S)$ for some $\varepsilon' < \varepsilon$, Lemma 1 (next slide) guarantees

$$\Pr\{\mathcal{A}_w(s^n)\} \geq 2^{-n\left(\mathrm{I}(S;\hat{S})+\delta(\varepsilon)\right)}$$

for sufficiently large $n$, where $\lim_{\varepsilon\to 0} \delta(\varepsilon) = 0$.

# Joint typicality lemma

The following lemma formally states the bounds.

## Lemma 1 (Joint Typicality Lemma)

*Consider a joint PMF $P_{X,Y} = P_X P_{Y|X} = P_Y P_{X|Y}$. Then, there exists $\delta(\varepsilon) > 0$ with $\lim_{\varepsilon \to 0} \delta(\varepsilon) = 0$ such that:*

**1** *For an arbitrary sequence $x^n$ and random $Y^n \sim \prod_{i=1}^{n} P_Y(y_i)$,*

$$\Pr\left\{(x^n, Y^n) \in \mathcal{T}_\varepsilon^{(n)}(P_{X,Y})\right\} \leq 2^{-n(\mathrm{I}(X;Y) - \delta(\varepsilon))}.$$

**2** *For an $\varepsilon'$-typical sequence $x^n \in \mathcal{T}_{\varepsilon'}^{(n)}(P_X)$ with $\varepsilon' < \varepsilon$, and random $Y^n \sim \prod_{i=1}^{n} P_Y(y_i)$, for sufficiently large $n$,*

$$\Pr\left\{(x^n, Y^n) \in \mathcal{T}_\varepsilon^{(n)}(P_{X,Y})\right\} \geq 2^{-n(\mathrm{I}(X;Y) + \delta(\varepsilon))}.$$

# Finalizing the proof

From Lemma 1, the additional condition that $s^n \in \mathcal{T}_{\varepsilon'}^{(n)}(\mathsf{P}_S)$ for some $\varepsilon' < \varepsilon$ motivates us to split the upper bound on $\Pr\{\mathcal{E}\}$ as follows:

$$\Pr\{\mathcal{E}\}$$

$$= \sum_{s^n \in \mathcal{S}^n} \mathsf{P}(s^n) \prod_{w=1}^{2^k} \left(1 - \Pr\{\mathcal{A}_w(s^n)\}\right)$$

$$\leq \sum_{s^n \notin \mathcal{T}_{\varepsilon'}^{(n)}(\mathsf{P}_S)} \mathsf{P}(s^n) + \sum_{s^n \in \mathcal{T}_{\varepsilon'}^{(n)}(\mathsf{P}_S)} \mathsf{P}(s^n) \prod_{w=1}^{2^k} \left(1 - \Pr\{\mathcal{A}_w(s^n)\}\right)$$

$$\leq \Pr\left\{S^n \notin \mathcal{T}_{\varepsilon'}^{(n)}(\mathsf{P}_S)\right\} + \sum_{s^n \in \mathcal{T}_{\varepsilon'}^{(n)}(\mathsf{P}_S)} \mathsf{P}(s^n) \left(1 - 2^{-n\left(\mathrm{I}\left(S;\hat{S}\right)+\delta(\varepsilon)\right)}\right)^{2^k}$$

$$\leq \Pr\left\{S^n \notin \mathcal{T}_{\varepsilon'}^{(n)}(\mathsf{P}_S)\right\} + \left(1 - 2^{-n\left(\mathrm{I}\left(S;\hat{S}\right)+\delta(\varepsilon)\right)}\right)^{2^k}$$

We then obtain a nice upper bound by using the fact that $(1-x)^r \le e^{-rx}$ for $x \in [0,1]$ and $r \ge 0$:

$$
\begin{aligned}
\Pr\{\mathcal{E}\} &\le \Pr\left\{ S^n \notin \mathcal{T}_{\varepsilon'}^{(n)}(\mathsf{P}_S) \right\} + \left( 1 - 2^{-n\left( \mathrm{I}(S;\hat{S}) + \delta(\varepsilon) \right)} \right)^{2^k} \\
&\le \Pr\left\{ S^n \notin \mathcal{T}_{\varepsilon'}^{(n)}(\mathsf{P}_S) \right\} + \exp\left( -2^k \, 2^{-n\left( \mathrm{I}(S;\hat{S}) + \delta(\varepsilon) \right)} \right).
\end{aligned}
$$

The first term vanishes as $n \to \infty$ due to AEP.

The second term vanishes as $n \to \infty$ if

$$
\mathrm{R} > \mathrm{I}\left( S; \hat{S} \right) + \delta(\varepsilon) = \mathrm{R}^{\mathrm{I}}\left( \tfrac{\mathrm{D}}{1+\varepsilon} \right) + \delta(\varepsilon).
$$

Hence, $\forall\, \mathrm{R} > \mathrm{R}^{\mathrm{I}}\left( \tfrac{\mathrm{D}}{1+\varepsilon} \right) + \delta(\varepsilon)$, the scheme achieves average distortion $\le \mathrm{D}$.

Finally, due to continuity of the rate-distortion function, we take $\varepsilon \to 0$ and complete the proof. $\qquad\square$

# Extensions to stationary ergodic sources

As in lossless source coding, the lossy source coding theorem can be extended to stationary ergodic sources, where the rate distortion function is characterized as follows: $\mathrm{R}(\mathsf{D}) = \mathcal{R}^{\mathrm{I}}(\mathsf{D})$, where

$$\mathcal{R}^{\mathrm{I}}(\mathsf{D}) := \lim_{n \to \infty} \min_{\substack{\mathsf{P}_{\hat{S}^n \mid S^n}: \\ \mathsf{E}\left[d(S^n, \hat{S}^n)\right] \leq \mathsf{D}}} \frac{1}{n} \mathrm{I}\left(S^n; \hat{S}^n\right).$$

Details can be found in Gallager [Gal68].

# Summary

- For a discrete memoryless source $S \sim \mathsf{P}_S$,

$$\mathrm{R}(\mathsf{D}) = \min_{\mathsf{P}_{\hat{S}|S}:\, \mathsf{E}\left[d(S,\hat{s})\right] \leq \mathsf{D}} \mathrm{I}\left(S; \hat{S}\right).$$

  For a continuous memoryless source $S$ with density $\mathsf{f}_S$,

$$\mathrm{R}(\mathsf{D}) = \inf_{\mathsf{f}_{\hat{S}|S}:\, \mathsf{E}\left[d(S,\hat{s})\right] \leq \mathsf{D}} \mathrm{I}\left(S; \hat{S}\right).$$

- Properties of a rate distortion function $\mathrm{R}(\mathsf{D})$:
  1. Nonnegative; Non-increasing in $\mathsf{D}$; Convex in $\mathsf{D}$; Continuous in $\mathsf{D}$.
  2. $\mathrm{R}(\mathsf{D}_{\min}) \leq \mathrm{H}(S)$; $\mathrm{R}(\mathsf{D}) = 0$ if $\mathsf{D} \geq \mathsf{D}_{\max}$.

- Weak converse: Data processing, DMS assumption, functional properties of a rate distortion function.

- Achievability: random coding, typicality encoder, typical average lemma, joint typicality lemma.

- Bernoulli with Hamming Distortion:

$$\mathrm{R(D)} = \max\{\mathsf{H_b}(p) - \mathsf{H_b}(\mathsf{D}), 0\}.$$

- Gaussian with squared error distortion:

$$\mathrm{R(D)} = \max\left\{\frac{1}{2}\log\left(\frac{\sigma^2}{\mathsf{D}}\right), 0\right\}.$$

- Gaussian source is the worst source under a second moment constraint.

# References

[Gal68] Robert G. Gallager.
*Information Theory and Reliable Communication*.
John Wiley & Sons Inc., USA, 1968.