# Lossy Compression & The Source Coding Theorem : (§4)

Today we **fix** the number of bits but allow small **error probability** ("lossy"):

$$X \longrightarrow \boxed{C} \longrightarrow \{0,1\}^\ell \longrightarrow \boxed{D} \longrightarrow \hat{X}$$

Compressor, encoder

decompressor, decoder

**WANT:**
$$\Pr(\hat{X} \neq X) \leq \delta$$

How to achieve?

\* Take set $S \subseteq \mathcal{A}$ with $\Pr(X \notin S) \leq \delta$.

\* Then we can compress into $\ell = \lceil \log \#S \rceil$ bits with error probability $\leq \delta$. How?
Simply define $C$ by sending all $x \in S$ to distinct bitstrings. (For $x \notin S$, pick arbitrary, or fail.)

ex:

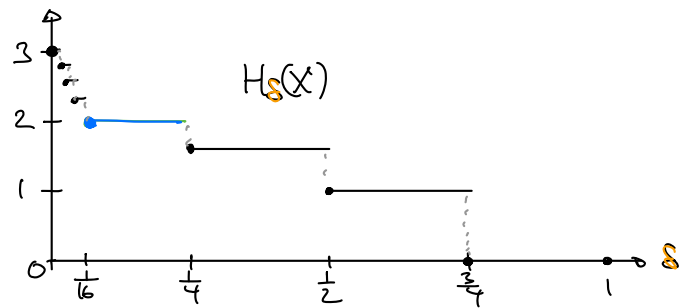| $x$ | $P(x)$ | $\delta = 0$ | $\delta = 1/16$ |
|---|---|---|---|
| a | $1/4$ | 000 | 00 |
| b | $1/4$ | 001 | 01 |
| c | $1/4$ | 010 | 10 |
| d | $3/16$ | 011 | 11 |
| e | $1/64$ | 100 | — |
| f | $1/64$ | 101 | — |
| g | $1/64$ | 110 | — |
| h | $1/64$ | 111 | — |

} arbitrary

$\ell = 2$

Define **$\delta$-essential bit content** by

$$H_\delta(X) = H_\delta(P) = \min \left\{ \log \#S \mid \Pr(X \notin S) \leq \delta \right\}$$

$$\Rightarrow \quad \boxed{\lceil H_\delta(X) \rceil \text{ is minimal } \# \text{ bits required to compress } X \text{ with error } \leq \delta}$$

if not integer, need to round up!

$H_\delta(X)$ is in general quite messy ... see D lee

Amazingly, it simplifies dramatically if we compress **blocks** of symbols.



$H_\delta(X)$

Shannon's Source Coding Theorem : Let $X_1, X_2, X_3, \ldots \overset{IID}{\sim} P$ and $0 < \delta < 1$:

$$\lim_{N \to \infty} \frac{H_\delta(X_1, \ldots, X_N)}{N} = H(P)$$

IID (memoryless) information source

optimal Compression rate for block size N and error prob $\leq \delta$

optimal asymptotic Compression rate
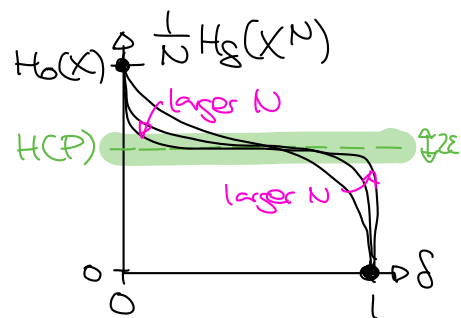
← independent of $\delta$!

(ie. $\forall \varepsilon \in (0,1), \varepsilon > 0 \; \exists N_0 \; \forall N \geq N_0 : \left| \frac{H_\delta(X_1 \cdots X_N)}{N} - H(P) \right| \leq \varepsilon$)

\* If $R > H(P)$: $\exists N_0 \, \forall N \geq N_0$:

CAN compress at rate $R$ $(\equiv$ into $\ell \leq RN$ bits$)$

\* If $R < H(P)$: $\exists N_0 \, \forall N \geq N_0$:

CANNOT compress at rate $R$



## Proof of the Source Coding Theorem

NOTATION: $x^N = x_1 \cdots x_N = (x_1, \ldots, x_N)$ for strings of length $N$.

Typical set: 
$$T_{N,\varepsilon}(P) = \left\{ x^N \in \mathcal{A}_X^N : \left| \frac{1}{N} \log \frac{1}{P(x^N)} - H(P) \right| \leq \varepsilon \right\}$$

$$\stackrel{\text{iid}}{=} \left\{ x^N \in \mathcal{A}_X^N : \left| \frac{1}{N} \sum_{k=1}^{N} \log \frac{1}{P(x_k)} - H(P) \right| \leq \varepsilon \right\}$$

Properties:

⓪ $2^{-N(H(P)+\varepsilon)} \leq P(x^N) \leq 2^{-N(H(P)-\varepsilon)}$ (by definition)

① $\#T_{N,\varepsilon} \leq 2^{N(H(P)+\varepsilon)}$

Pf: $1 \geq \Pr(X^N \in T_{N,\varepsilon}) = \sum_{x^N \in T_{N,\varepsilon}} P(x^N) \geq \#T_{N,\varepsilon} \cdot 2^{-N(H(P)+\varepsilon)}$. ☐

② $\Pr(X^N \notin T_{N,\varepsilon}) \leq \frac{\sigma^2}{N\varepsilon^2} \longrightarrow 0$, where $\sigma^2 = \text{Var}\left( \log \frac{1}{P(X_k)} \right)$.

Pf: Let $L_k = \log \frac{1}{P(X_k)}$ and $\mu := E[L_k] = H(X_k) = H(P)$. Then:

$$\text{LHS} = \Pr\left( \left| \frac{1}{N} \sum_{k=1}^{N} L_k - \mu \right| > \varepsilon \right) \leq \frac{\text{Var}(L_k)}{N\varepsilon^2}. \quad ☐$$

$\underbrace{\hspace{10cm}}$

"Asymptotic Equipartition Property" (AEP)

"For large $N$... ...typical probabilities are $2^{-N(H(P) \pm \varepsilon)}$."

Proof of the theorem: Let $\delta \in (0,1)$ and $\varepsilon > 0$ be arbitrary.

Ⓐ $\Pr(X^N \in T_{N,\varepsilon}) \stackrel{②}{\geq} 1 - \frac{\sigma^2}{N\varepsilon^2} \geq 1 - \delta$ if $N$ large enough

$\implies \frac{H_\delta(X^N)}{N} \leq \frac{\log \#T_{N,\varepsilon}}{N} \stackrel{①}{\leq} H(P) + \varepsilon$ for large $N$. $\overset{\infty}{\smile}$

(B) Want to prove that $\frac{H_\delta(X^N)}{N} \geq H(P) - \varepsilon$ for $N$ large.

If not: $\exists$ sets $S_N$ for $N \to \infty$ s.th.

$$\Pr(X^N \in S_N) \geq 1 - \delta \quad \text{and} \quad \#S_N < 2^{N(H(P)-\varepsilon)}.$$

$$\Longrightarrow 1-\delta \leq \Pr(X^N \in S_N) = \Pr(X^N \in S_N \cap T_{N,\varepsilon/2}) + \Pr(X^N \in S_N \setminus T_{N,\varepsilon/2})$$

$$\leq \underbrace{\Pr(X^N \in S_N \cap T_{N,\varepsilon/2})}_{\substack{⑥ \\ \leq \#S_N \cdot 2^{-N(H(P)-\frac{\varepsilon}{2})}}} + \underbrace{\Pr(X^N \notin T_{N,\varepsilon/2})}_{\to 0 \text{ by } ②} \longrightarrow 0 \; \lightning$$

$$\leq 2^{-N\varepsilon/2} \longrightarrow 0 \qquad\qquad \square$$

Remark: $T_{N,\varepsilon}$ is usually NOT the smallest set $S_N$ w/ $\Pr(X^N \in S_N) \geq 1-\delta$...

... but small enough and easy to handle as $N \to \infty$! $\longrightarrow$ EX CLASS

## How to use this in practice?

SCENARIO: Want to compress IID (memoryless) data source $P$
(we know $P$, but NOT which string will be emitted)

FIX: * block size $N$
* parameter $\varepsilon > 0$
* a way to order the typical set $T_{N,\varepsilon}$

| index | element |
|-------|---------|
| 0 | - - - - |
| 1 | - - - - |
| ⋮ | - - - - |
| $\#T_{N,\varepsilon}-1$ | - - - - |

COMPRESSOR: Input: A string $x^N = x_1 \cdots x_N$

* If $x^N \notin T_{N,\varepsilon}^{(P)}$: FAIL
* Determine index $p$ of $x^N$ in $T_{N,\varepsilon}$.
* Return $p$ in binary.

DECOMPRESSOR:

Input: A binary string $s$
* Interpret $s$ as integer $p$
* Return $p$-th element of $T_{N,\varepsilon}$.

This is a lossy compression protocol:

AEP

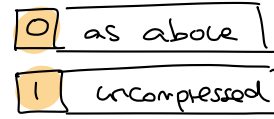* Error probability: $\Pr(X^N \notin T_{N,\varepsilon}) \leq \frac{\sigma^2}{N\varepsilon^2} \longrightarrow 0$ as $N \to \infty$

* Rate $R = \dfrac{\#\text{bits required to represent } p}{N}$

AEP

$$\leq \frac{\log \#T_{N,\varepsilon} + 1}{N} \leq H(P) + \varepsilon + \frac{1}{N} \longrightarrow 0$$

# Variations

Ⓐ How to make it LOSSLESS?

When $x^N \notin T_{N,\varepsilon}$, send uncompressed

using $N \cdot \lceil \log \#A_X \rceil$ bits.

| 0 | as above |
| 1 | uncompressed |

} prefix code!

"flag" bit

$\implies$ average rate $\bar{R} \leq \frac{1}{N} + \underbrace{Pr(X^N \in T_{N,\varepsilon})}_{\to 1}\underbrace{(H(P) + \varepsilon + \frac{1}{N})}_{\text{from above}}$

$+ \underbrace{Pr(X^N \notin T_{N,\varepsilon})}_{\to 0} \cdot \lceil \log \#A_X \rceil$

$\approx H(P) + \varepsilon$ for large $N$

← agrees with symbol code discussion

Ⓑ How to also make it UNIVERSAL? (IID, but we do NOT know P)

For simplicity: assume $A = \{0,1\}$ i.e. data source of bits.

FIX: * block size N

* a way to order the sets

$B(N,k) := \{x^N$ with $k$ ones and $N-k$ zeros$\}$

B(3,2)

| index | string |
|-------|--------|
| 0 | 011 |
| 1 | 101 |
| 2 | 110 |

→ ex class, HW

COMPRESSOR: Input: A bitstring $x^N = x_1 \cdots x_N$

* Compute $k :=$ #ones in $x^N$

* Determine index $p$ of $x^N$ in $B(N,k)$

* Return $k$ and $p$ in binary.

$\underbrace{\phantom{xx}}_{\approx \log(N) + 1}$ $\underbrace{\phantom{xx}}_{\approx \log \#B(N,k) + 1}$ bits

Key idea: $B(N,k)$ can be MUCH SMALLER than $\{0,1\}^N$

(e.g. imagine k=1)

DECOMPRESSOR

clear !? :

Not used in protocol, only in the analysis !!!

Average rate $\bar{R}$? Assume that $X_1, \ldots, X_N \overset{IID}{\sim} P$. Then:

depends on $P$, but only used in analysis !

$x^N \in T_{N,\varepsilon} \implies B(N,k) \subseteq T_{N,\varepsilon} \implies \#B(N,k) \leq \#T_{N,\varepsilon}$ ⊛

Since typicality only depends on # zeros and ones in $x^N$ !

Thus we can argue as above:

$$\overline{R} = \frac{\#\text{bits required to represent } k \ + \ \#\text{bits required to represent } P}{N}$$

$$\leq \frac{\log(N)}{N} + \frac{\log \#B(n,k)}{N}$$

$\frac{\log(N)}{N} \to 0$, so can ignore

use ⊛ to obtain the following bound:

$$\leq \Pr(X^N \in T_{N,\varepsilon}) \cdot \frac{\log \#T_{N,\varepsilon}}{N} + \Pr(X^N \notin T_{N,\varepsilon}) \frac{\log 2^N}{N}$$

$\Pr(X^N \in T_{N,\varepsilon}) \to 1$

$\frac{\log \#T_{N,\varepsilon}}{N} \leq H(P) + \varepsilon$

$\Pr(X^N \notin T_{N,\varepsilon}) \to 0$, as before

$\frac{\log 2^N}{N} = 1$

dropping some $\frac{1}{N}$ terms

$$\approx H(P) + \varepsilon \quad \text{for large } N!$$

HW: Program this protocol & compress the donkey!

Discussion: Many disadvantages!

* Have to look at entire $x^N$ to compress. Can we compress by looking at a few symbols at a time?

* Assume IID distribution... what if P changes? Or if we have local correlations?

Q → U  frequent
Q → R  rare

↳ Wednesday 😎