# Source Coding/Data Compression

Claude Shannon, "A Mathematical Theory of Communication"

$c^k \in \{0,1\}^k$

$X_1 \cdots X_n$
$X_i \in \mathcal{H}$

$X^n$

ENCODER/
COMPRESSOR

$c^k$

DECODER/
DECOMPRESSOR

$\hat{X}^n$

Raw file/
Source

$n, p_{X^n}$

Compressed
file/
codeword

$n, p_{X^n}$

Decompressed file/
reconstruction/
estimate.

Universal

$p_{X^n}$ : Source distribution

Memoryless : $X^n \sim iid(p_X)$

Rate : $R = \dfrac{k}{n}$

Assumptions:

① $X^n$ is random ✓

② $X^n \sim iid(p_x)$ ✗

③ $p_x$ is known to ENC DEC ✗

$p_{x^n}$

## 1. Lossless Compression:

$$\hat{X}^n \simeq X^n$$

### Fixed-length compression:

$k \cdot$ depends only on $n, P_{X^n}$

$P_e$ = Probability of error

| Harry Potter I | 1 MB | $\dfrac{100 \, kB}{100 \, kB}$ | $R = \dfrac{k}{n}$ |
|---|---|---|---|
| II | 1 MB | | |

### Variable-length compression:

$$\underline{\hat{X}^n = X^n}$$

$\longrightarrow$ Expected length $= \mathbb{E} \, k(X^n)$

$R_{av} = \dfrac{\mathbb{E} \, k(X^n)}{n}$

$k$ can vary even if $n, P_{X^n}$ fixed

| 1 MB | 100 kB |
|---|---|
| 1 MB | 102 kB |

## 2. Lossy Compression:

$$\hat{X}^n \neq X^n$$

Distortion measure

$$\underline{\mathbb{E} \, d(X^n, \hat{X}^n) < \delta}$$

$$MSD = \mathbb{E} \| \hat{X}^n - X^n \|^2 = \mathbb{E} \sum_{i}^{n} (X_i - \hat{X}_i)^2$$

$n=2$

| | | $P(x^n)$ | | $C^k$ | | |
|---|---|---|---|---|---|---|
| 0 0 | | 0.5 | $\longrightarrow$ | 0 | $\longrightarrow$ | 0 0 |
| 0 1 | | 0.25 | $\longrightarrow$ | 1 | $\longrightarrow$ | 0 1 |
| 1 0 | | 0.125 | $\longrightarrow$ | 0 0 | $\longrightarrow$ | 1 0 |
| 1 1 | | 0.12p | $\longrightarrow$ | 0 1 | $\longrightarrow$ | 1 1 |

2 bits

Expected length: $\mathbb{E}\, k(X^n)$

$$= 0.5 \times 1 + 0.25 \times 1 + 0.125\, k2$$
$$+ 0.125\, \times L$$

$$\approx 0.5 + 0.25 + 0.5$$

$$\approx 1.25 \text{ bits} \quad < 2 \text{ bits}.$$

$R_{avg} = \dfrac{1.25}{2} < 1$

| | | | ENC | |
|---|---|---|---|---|
| 0 0 | | 0.5 | $\longrightarrow$ | 0 |
| 0 1 | | 0.25 | $\longrightarrow$ | 1 |
| 1 0 | | 0.125 | $\longrightarrow$ | 0 |
| 1 1 | | 0.12p | $\longrightarrow$ | 1 |

DEC

$0 \longrightarrow 00$   $k=1$ bit

$1 \longrightarrow 01$

$P_e$

Probability of error

$$= \Pr[\hat{X}^n \neq X^n]$$

$$= P(10) + P(11) \quad = 0.25$$

$$R_{avg} = \frac{\mathbb{E}k(X^n)}{n}$$

Assumptions: $(n, P_{X^n})$ is known

Computational power - free.

Want $\mathbb{E} \, k(X^n)$ to be as small as possible.

$$P_{X^n}(x^n) \uparrow \quad \Rightarrow \quad k(X^n) \downarrow$$

$X^n \in \mathcal{X}^n \qquad \mathcal{X} = \{0, 1, 2 \cdots a\}$

$$x_1 \text{ -- -- --- } x_n$$

$c^k$

$(a+1)^n \Bigg\{ \begin{array}{l} \\ \\ \\ \\ \\ \\ \end{array}$

$x^n(1)$    $P_{X^n}(x^n(1))$

$x^n(2)$     $\vee$

      $P_{X^n}(x^n(2))$

$\vdots$      $\vdots$

$\vdots$

$x^n((a+1)^n)$

$x^n(1)$      1 0

$x^n(2)$      0 1

$x^n(3)$      0 0

$x^n(4)$      1 1

'

'

'

0        $\phi$

1        0

00       1

01       0 0

10       01

11       10

0 0 0     11

0 0 1     0 0 0

0 1 0     !

:

1 1 1

0 0 0 0

'

'

'

$\mathbb{E} k(X^n)$

$= \sum_{i=1}^{(a+1)^n} k(x^n(i))$

$P_{X^n}(x^n(i))$

$n(i)$

$$k(n(1)) = 0$$

$$k(n(2)) = 1$$

$$k(n(3)) = 1$$

$$4 = 2$$

$$5 = 2$$

$$6 = 2$$

$$7 = 2$$

$$8 = 3$$

$$\vdots$$

$$15 = 3$$

$$16 = 4$$

$$\vdots$$

$$31 = 4$$

$$\vdots$$

$$k(n(i)) = \lfloor \log_2 i \rfloor$$

$$k(x^n(i)) = \lfloor \log_2 i \rfloor$$

$$P_{X^n}(x^n(1)) \geq P_{X^n}(x^n(2)) \geq \cdots$$

**Claim:** $\quad P_{X^n}(x^n(i)) \leq \dfrac{1}{i}$

$$P_{X^n}(1) \leq 1 \qquad P_{X^n}(2) \leq \tfrac{1}{2}$$

$$P_{X^n}(3) \leq \tfrac{1}{3}$$

$\underline{\text{Assume:}}\ P_{X^n}(2) > \tfrac{1}{2} \ \Rightarrow\ P_{X^n}(1) \geq P_{X^n}(2) > \tfrac{1}{2}$

$\qquad\qquad \Rightarrow\ P_{X^n}(1) + P_{X^n}(2) > \tfrac{1}{2} + \tfrac{1}{2} = 1$

$($

---

## EXPECTED LENGTH

$$\mathbb{E}\, k(X^n) = \sum_{i=1}^{(a+1)^n} k(x^n(i))\, P_{X^n}(x^n(i))$$

$$= \sum_{i=1}^{(a+1)^n} P_{X^n}(x^n(i)) \ \boxed{\lfloor \log_2 i \rfloor}$$

$$\leq \sum_{i=1}^{(a+1)^n} P_{X^n}(x^n(i)) \ \left\lfloor \log_2 \frac{1}{P_{X^n}(x^n(i))} \right\rfloor$$

$$\leq \sum_{i=1}^{(a+1)^n} P_{X^n}(x^n(i))\, \log_2 \frac{1}{P_{X^n}(x^n(i))}$$

$$= \underbrace{\sum_{x^n} P_{X^n}(x^n)\, \log \frac{1}{P_{X^n}(x^n)}}_{H(X^n)}$$

ENTROPY OF $P_{X^n}$

In general, if $P_{X^n}(i) > \frac{1}{i}$ for some $i$

$$\Rightarrow \quad P_{X^n}(1) \geq P_{X^n}(2) \geq \cdots \cdots \geq P_{X^n}(i) > \frac{1}{i}$$

$$\sum_{j=1}^{i} P_{X^n}(j) = P_{X^n}(1) + P_{X^n}(2) + \cdots + P_{X^n}(i)$$

$$> \frac{1}{i} + \frac{1}{i} + \frac{1}{i} + \cdots + \frac{1}{i}$$

$$\geq i \times P_{X^n}(i) > i \times \frac{1}{i} = 1 \quad \Rightarrow \quad \text{NOT POSSIBLE}$$
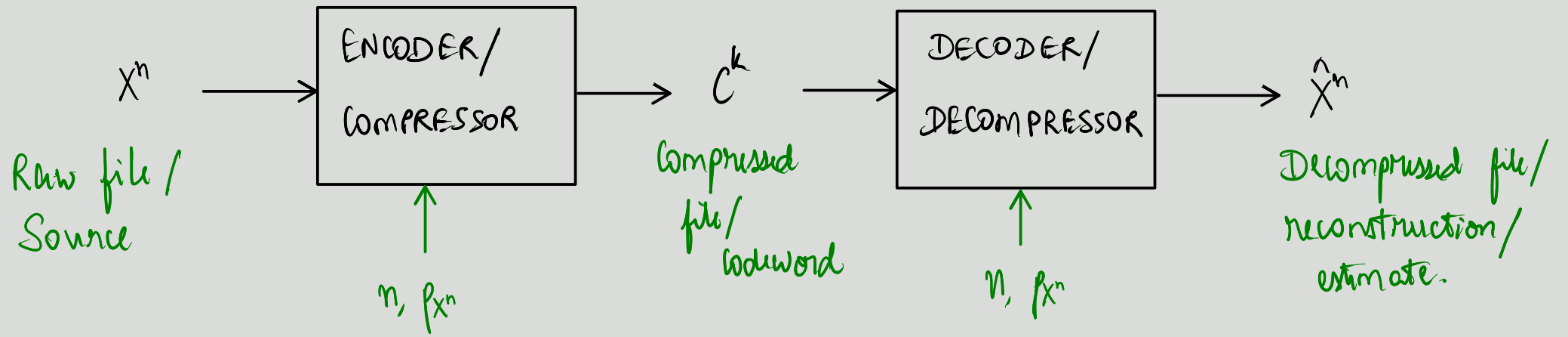$$\text{Contradiction!}$$

$$P_{X^n}(i) \leq \frac{1}{i} \quad \Rightarrow \quad i \leq \frac{1}{P_{X^n}(i)} \quad = \quad \frac{1}{P_{X^n}(X^n(i))}$$

$$R_{avg} = \lim_{n \to \infty} \frac{\mathbb{E} k(X^n)}{n} = \lim_{n \to \infty} \frac{H(X^n)}{n} \quad \longrightarrow \quad \underline{\text{ENTROPY RATE}}$$

Can construct compressor for which $R_{avg} = $ ENTROPY RATE.

For iid $X^n$,

$$\lim_{n \to \infty} \frac{H(X^n)}{n} = H(X) = \sum_{x \in \mathcal{X}} P_X(x) \log_2 \frac{1}{P_X(x)}$$

$X^n$

Raw file /
Source

ENCODER/
COMPRESSOR

$n, p_{X^n}$

$C^k$

Compressed
file /
codeword

DECODER/
DECOMPRESSOR

$n, p_{X^n}$

$\hat{X}^n$

Decompressed file/
reconstruction/
estimate.

*Lossless source coding theorem:*

① $X^n \sim iid(p_X)$

Achievability : We can construct (ENC, DEC) $\sigma$r

for every $\varepsilon > 0$, (fixed length compressor)

$$R = H(X) + \varepsilon$$

$$P_e \to 0 \quad \text{as} \quad n \to \infty$$

find
length
compression

Converse : for every (ENC, DEC) with $R < H(X)$

$$P_e \to 1 \quad \text{as} \quad n \to \infty$$

$$X^n \sim iid(p_X) \implies H(X^n) = n H(X) \implies \frac{H(X^n)}{n} = H(X)$$

$$H(X^n) = \sum_{x} p_{X^n}(x^n) \log_2 \frac{1}{p_{X^n}(x^n)}$$

$$= \sum_{x^n} \left( \prod_{i=1}^{n} p_X(x_i) \right) \log_2 \frac{1}{\left( \prod_{i=1}^{n} p_X(x_i) \right)} \qquad (iid)$$

$$= \sum_{x^n} \prod_{i=1}^{n} p_X(x_i) \sum_{j=1}^{n} \log \frac{1}{p_X(x_j)}$$

$$= \sum_{j=1}^{n} \sum_{x^n} \left( \prod_{i=1}^{n} p_X(x_i) \right) \left( \log \frac{1}{p_X(x_j)} \right)$$

$$= \sum_{j=1}^{n} \sum_{x_j} p_X(x_j) \log_2 \frac{1}{p_X(x_j)} = n \sum_{x} p_X(x) \log_2 \frac{1}{p_X(x)}$$

$$H(X^n) = H(X_1) + H(X_2) + \cdots \cdot H(X_n) \quad \text{if } X^n \text{ is iid}$$

**Entropy is non-negative**

$$X \sim p_X$$

$$H(X) = \sum_{x \in \mathcal{X}} p_X(x) \log \frac{1}{p_X(x)}$$

By convention,

$$p_X(x) = 0$$

$$p_X(x) \log_2 \frac{1}{p_X(x)} = 0$$

$$0 < p_X(x) \leq 1$$

$$\lim_{\substack{p \to 0 \\ p \downarrow 0}} p \log_2 \frac{1}{p} = 0$$

$$\Rightarrow \quad 1 \leq \frac{1}{p_X(x)}$$

$$H(X) = \sum_x p_X(x) \log_2 \frac{1}{p_X(x)} \geq 0$$

Entropy: $H(X)$ $\qquad$ $H(p_X) \longrightarrow$ Better notation for entropy

$\qquad\qquad\qquad\quad \downarrow$

actually a function of $p_X$ & not $x$.

# Entropy is invariant to relabeling

$\mathcal{X} = \{0, 1, 2\}$.  $P_X(0) = \frac{1}{2}$

$P_X(1) = \frac{1}{4}$

$P_X(2) = \frac{1}{4}$

$H(X) = 1.5 \text{ bits}$

Entropy is a function only of the probability multiset $\{\frac{1}{2}, \frac{1}{4}, \frac{1}{4}\}$

$q_X(0) = \frac{1}{4}$

$q_X(1) = \frac{1}{2}$

$H(q_X) = 1.5 \text{ bits}$

$q_X(2) = \frac{1}{4}$

$$H(X) = \sum_{x \in \mathcal{X}} P_X(x) \log_2 \frac{1}{P_X(x)} \qquad \underline{\underline{\text{bits}}}$$

$$(\ ) \log_e (\ ) \qquad \underline{\underline{\text{nats}}}$$

① $X \sim \text{Unif}(\mathcal{X})$ $\qquad |\mathcal{X}| = m$

$$p_X(x) = \frac{1}{m} \quad \forall x \in \mathcal{X}$$

$$H(X) = \sum_x \frac{1}{m} \log_2 \frac{1}{(1/m)}$$

$$= \sum_x \frac{1}{m} \log_2 m$$

$$= \log_2 m$$

$$H(X) = \log_2 |\mathcal{X}|$$

② $\mathcal{X} = \{1, 2, 3, \text{----}\} \qquad \mathcal{Y} = \mathbb{Z}_{>0}$

Geometric rv: $\qquad p_X(x) = \left(\frac{1}{2}\right)^x \qquad \sum_{x=1}^{\infty} \left(\frac{1}{2}\right)^x = 1$

$$H(X) = \sum_{x=1}^{\infty} \frac{1}{2^x} \log_2 \frac{1}{1/2^x} = \sum_{x=1}^{\infty} \frac{x}{2^x} = 2 \text{ bits.}$$

$$\sum_{x_n} \frac{x}{2^n} \qquad\qquad \sum_{n=1}^{\infty} p^n = \frac{p}{1-p}$$

$$\sum_{n=1}^{\infty} n p^n = p \sum_{n=1}^{\infty} n p^{n-1} = p \frac{d}{dp}\left(\sum_{n=1}^{\infty} p^n\right)$$

$$= p \frac{d}{dp}\left(\frac{p}{1-p}\right) = p\left(\frac{1}{1-p} + \frac{p}{(1-p)^2}(\times 1)\right)$$

$$= p \times \frac{1-p+p}{(1-p)^2}$$

$$= \frac{p}{(1-p)^2}$$

$$X_i \in \mathbb{Z} \qquad \mathbb{E}\, k(X^n) \approx H(X^n) \qquad\qquad |\mathbb{Z}| = m$$

$$X^n \text{ compressd } \underline{\underline{C^k}} \quad \underline{\underline{\mathbb{E}\, k}} \approx H(X^n) = n H(X) \leq n \log|\mathbb{Z}| \qquad \underline{n\lceil \log_2 m\rceil}$$
$$\text{(iid)}$$

$$\text{iid} \qquad \text{Geometric}(\tfrac{1}{2}) \text{ seq} \rightarrow \underline{\underline{2n}} \text{ bits on } \underline{\text{avg}}$$

$$
\begin{array}{cc}
1 & 0\ 0\ 0 \\
2 & 0\ 0\ 1 \\
3 & 0\ 1\ 0 \\
4 & 0\ 1\ 1 \\
5 & 1\ 0\ 0 \\
\end{array}
$$

$$\lceil \log_2 m \rceil$$

# Fixed-length compression

$$k \to fn \ of \ n, \ P_{X^n}$$

$$X_1 \cdots X_n$$
$$X_i \in \mathcal{X}$$

$$c^k \in \{0,1\}^k$$

$$X^n \longrightarrow \boxed{\begin{array}{c} \text{ENCODER/} \\ \text{COMPRESSOR} \end{array}} \longrightarrow C^k \longrightarrow \boxed{\begin{array}{c} \text{DECODER/} \\ \text{DECOMPRESSOR} \end{array}} \longrightarrow \hat{X}^n$$

Raw file / Source

$\uparrow$ $n, \ P_{X^n}$

Compressed file / Codeword

$\uparrow$ $n, \ P_{X^n}$

Decompressed file / reconstruction / estimate.

$$P_e = Pr[\hat{X}^n \neq X^n]$$

$$k - fixed \qquad depends \ only \ on \ n, P_X$$

$$X^n \sim iid \ (P_X) \qquad\qquad X_i \sim P_X$$

$$R = \frac{k}{n} \qquad , \qquad P_e \qquad small.$$

$$n \gg 1 \qquad n \to \infty \qquad \lim_{n \to \infty} \frac{k}{n} = R \quad small, \quad P_e \to 0 \quad as \atop n \to \infty$$

# Compressing Bernoulli(p) sequences

$$P_X(1) = p, \qquad P_X(0) = 1-p.$$

$$X^n = 0\ 1\ 0\ 1\ 1\ 0\ 1\ 1\ 1\ 1\ 1\ 0\ 0\ 0\ 0\ 0$$

$$\downarrow$$

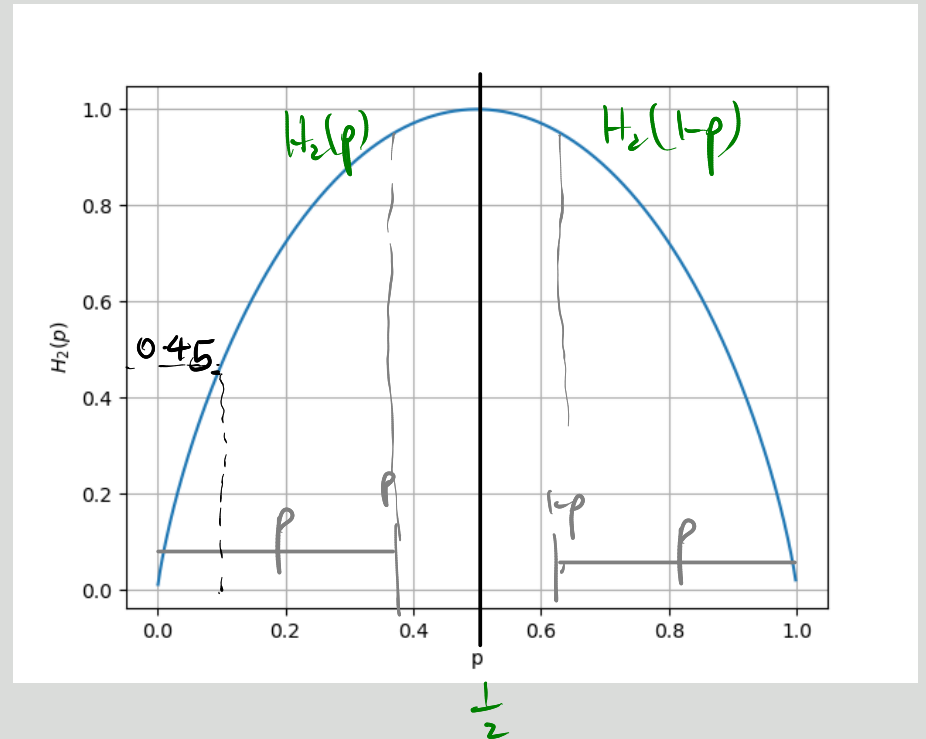$$C^n \qquad 0\ 0\ 0\ 0\ 1\ 0\ 0\ 1\ 0$$

$$R_{opt} = H(X)$$

$$= p \log_2 \frac{1}{p} + (1-p) \log_2 \frac{1}{1-p}$$

$$R_{opt}(p) = H_2(p) \longrightarrow \text{Binary entropy } (p)$$

$$\text{as } p \to 0 \qquad p \log_2 \frac{1}{p} \to 0 \quad , \quad (1-p) \log_2 \frac{1}{1-p} \to 0 \quad \Rightarrow \quad H_2(p) \to 0$$

$$p \to 1 \qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad H(p) \to 0$$

**GOAL:** Construct ENC, DEC st

① $\dfrac{k}{n} \longrightarrow H(X)$ as $n \to \infty$

② $P_e = \Pr\left[\hat{X}^n \neq X^n\right] \to 0$ as $n \to \infty$.

**Structure of ANY fixed-length compressor**

$X^n \sim$ iid $\mathrm{Ber}(0.1)$ $\qquad H_2(0.1) = 0.45$

$n \longrightarrow 0.45n$

$X^n \longrightarrow C^{0.45n}$ $\qquad\qquad\qquad\qquad 10^{-6} \qquad\qquad 10^7$

$\underline{X}^{1000} \longrightarrow C^{450} \longrightarrow \underline{\hat{X}}^{1000} \qquad\qquad \Pr\left[\hat{X}^{1000} \neq X^{1000}\right] \leq 10^{-2}$

EXTREME: $p = 0 \qquad n$ $\qquad\qquad\qquad 0\,0\,0\,0 \,\cdots\, 0 \qquad k = 0$
$\overset{\longleftarrow n \longrightarrow}{}$
$\qquad\qquad\quad p = 1 \qquad\qquad\qquad\qquad\qquad 11\cdots\! - \qquad 1 \qquad k = 0$

$X^n$ : ENC, DEC, know $\#1's = 1$, $n$

$n = 100,000$

ENC : location of the 1.

$$k = \lceil \log_2 100,000 \rceil = \lceil \log_2 n \rceil \ll n$$

ASSUME: $Pr[\#1's \text{ in } X^n > 1] \leq \frac{1}{n}$

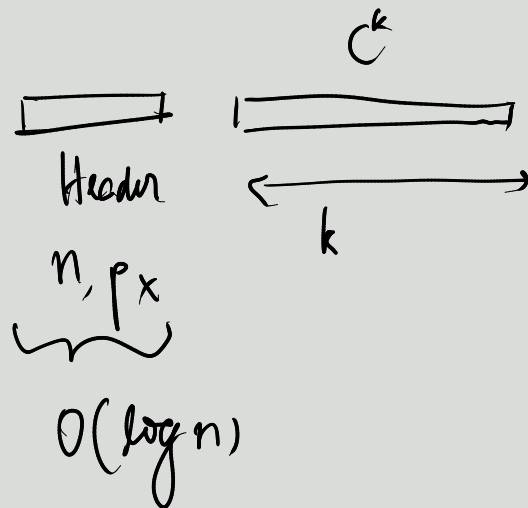$\frac{k}{n} = \frac{\lceil \log_2 n \rceil}{n}$      $P_e = \leq \frac{1}{n} \rightarrow 0$    as $n \rightarrow \infty$

$\xrightarrow[n \rightarrow \infty]{} 0$

$p = 0 \qquad k = 0 \qquad p_c = 0$

Assumptions: $(n, p)$ known

Compressed file:



Header

$\underbrace{n, p_x}$

$O(\log n)$

Can be ignored if $n \gg 1$.

$a \rightarrow 0$

$b \rightarrow 10$

$c \rightarrow 110$

$d \rightarrow 1100$

$e \rightarrow 1110$

$\lceil \log_2 5 \rceil \simeq 3 \text{ bit}$

$2^8$

. $-$

$X^n \sim$ iid Ber$(p)$

#1's $\approx np$

$\Pr\left[ |\#1's \text{ in } X^n - np| > n\varepsilon \right] \to 0$     as $n \to \infty$

$\Pr\left[ |(\text{fraction of 1's in } X^n) - p| > \varepsilon \right] \to 0$     as $n \to \infty$

(WLLN)

$2^n - 2^{nH_2(p)}$

$S_{\varepsilon,p}^n = \left\{ x^n \in \{0,1\}^n : np(1-\varepsilon) \le M_1(x^n) \le np(1+\varepsilon) \right\}$     $|S| \approx 2^{nH_2(p)}$

$\ll 2^n$

$\Pr\left[ X^n \in S \right] \to 1$     as $n \to \infty$

$x^n(1)$          $0\,0\,0 - 0$

$\vdots$            $0\,0 \cdots - 1$

                    $k = \lceil \log_2 |S| \rceil$

$\vdots$

$\approx nH(x) = \underline{n\,H_2(p)}$

$x^n(|S|)$        $1\,1 \cdots 1$

$x^h(1) \longrightarrow 0\,0 - 0 \longrightarrow x^h(1)$
$\quad\rightharpoonup x^h(2)$

$\vdots$

$x^h(|S|) \longrightarrow 1\,1 - - - 1 \longrightarrow x^h(|S|)$

$x^h(|S|+1) \longrightarrow$

$\vdots$

$x^h(2^r)$

$\widehat{q_c}$

$an$

$as$

$\widehat{q_2}$