S — Encoder — X — Channel — Y — Decoder — $\hat{S}$
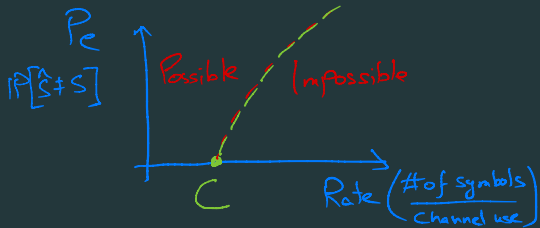
Data  
transmitted signal  
received signal  
Inferred data

# ORIE 4742 - Info Theory and Bayesian ML

Chapter 5: Channel Coding

Sid Banerjee, ORIE, Cornell

Want - $\mathbb{P}[\hat{S} \neq S] \searrow 0$

$P_e$ $\mathbb{P}[\hat{S} \neq S]$

Possible   Impossible

C   Rate $\left(\frac{\# \text{ of symbols}}{\text{Channel use}}\right)$

# dependent rv and information content

Chapter 8 from Mackay

## entropy: basic properties

rv $X$ taking values $\mathcal{X} = \{a_1, a_2, \ldots, a_k\}$, with pmf $\mathbb{P}[X = a_i] = p_i$

**Shannon's entropy function**

- outcome $X = a_i$ has *information content*: $\quad h(a_i) = \log_2\left(\frac{1}{p_i}\right)$
- random variable $X$ has *entropy*: $\quad H(X) = \mathbb{E}[h(X)] = \sum_{i=1}^{k} p_i \log_2\left(\frac{1}{p_i}\right)$

- only depends on distribution of $X$ (i.e., $H(X) = H(p_1, p_2, \ldots, p_k)$)
- $H(X) \geq 0$ for all $X$
- if $X \sim$ uniform on $\mathcal{X}$, then $H(X) = \log_2 |\mathcal{X}|$; else, $H(X) \leq \log_2 |\mathcal{X}|$
- if $X \perp\!\!\!\perp Y$, then $H(X, Y) = H(X) + H(Y)$
  where joint entropy $H(X, Y) \triangleq \sum_{(x,y)} p(x, y) \log_2 1/p(x, y)$

# mutual information

for any rvs $X, Y$:

$$I(X; Y) = H(X) - H(X|Y) = H(Y) - H(Y|X)$$

moreover, given any other conditioning rv $Z$

$$I(X; Y|Z) = H(X|Z) - H(X|Y, Z) = H(Y|Z) - H(Y|X, Z)$$

$D_{KL}(P \| Q)$

$= \sum P(x) \log \frac{P(x)}{Q(x)}$

$= 0$ if $P = Q$

$H(x, y)$

$H(x)$

$H(y|x)$

$H(x|y)$

$H(y)$

$I(x; y)$

$I(x; y) =$

$D_{KL}\left(P(x)P(y) \| P(x, y)\right)$

# conditional entropy

**conditional entropy**

for any rvs $X, Y$: $H(X|Y) = \sum_{y \in \mathcal{Y}} p(y) H(X|Y = y)$

$\qquad\qquad\qquad\quad = \sum_{y \in \mathcal{Y}} p(y) \sum_{x \in \mathcal{X}} p(x|y) \log_2(1/p(x|y))$

(Joint) $H(X,Y) = \sum_{(x,y)} p(x,y) \log_2\left(\dfrac{1}{p(x,y)}\right) = \sum_{(x,y)} p(x,y)\, h(x,y)$

(Marginals) $H(X) = \sum_{x} p(x)\, h(x)$ , $H(Y) = \sum_{y} p(y)\, h(y)$

(conditional) $\cdot \left\{ \underbrace{p(x|y)}_{} = \mathbb{P}[X=x|Y=y] \right\}_{x \in \mathcal{X}} \quad \forall\ y \in \mathcal{Y}$

$\underbrace{\qquad\qquad}_{} \; h(x|y) = \log_2\left(\dfrac{1}{p(x|y)}\right)$

$H(X|Y) = \sum_{y} p(y) \left( \sum_{x} p(x|y)\, h(x|y) \right)$

# the chain rule

## the chain rule (information content)

for any rvs $X, Y$ and realizations $x, y$:

$$h(x, y) = h(x) + h(y|x) = h(y) + h(x|y)$$

$$h(x,y) = \log_2\left(1 \middle/ P(x,y)\right), \quad h(x) = \log_2 1/P(x), \quad h(x|y) = \log_2 1/P(x|y)$$

$$\bullet \log_2 \frac{1}{P(x,y)} = \log_2\left(\frac{1}{P(x)P(y|x)}\right) = h(x) + h(y|x)$$

**the chain rule (entropy)**

for any rvs $X, Y$:

$$H(X, Y) = H(X) + H(Y|X) = H(Y) + H(X|Y)$$

$$H(x,y) = \sum_{(x,y)} p(x,y) \log_2 \left( 1/p(x,y) \right)$$

$$H(x) = \sum_{x} p(x) \log_2 \left( 1/p(x) \right) = \sum_{x,y} p(x,y) \log_2 \left( 1/p(x) \right)$$

$$H(x|y) = \sum_{y} p(y) \left( \sum_{x} p(x|y) \log_2 \left( 1/p(x|y) \right) \right)$$

$$= \sum_{x,y} p(x,y) \log_2 \left( 1/p(x|y) \right)$$

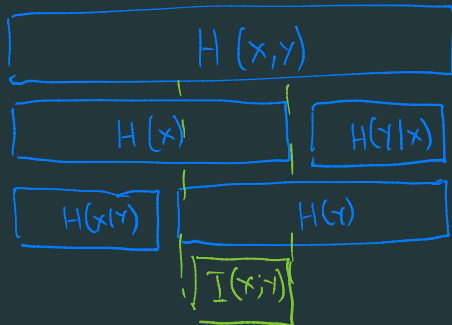# mutual information

for any rvs $X, Y$:
$$I(X; Y) = H(X) - H(X|Y) = H(Y) - H(Y|X)$$

moreover, given any other conditioning rv $Z$
$$I(X; Y|Z) = H(X|Z) - H(X|Y, Z) = H(Y|Z) - H(Y|X, Z)$$

$H(x,y)$

$H(x)$   $H(Y|x)$

$H(x|y)$   $H(y)$

$I(x;y)$

$H(x,y) = H(x) + H(y) - I(x;y)$

# example



| $P(x,y)$ | | $h(x_y)$ | $x$ | | | $P(y)$ |
|---|---|---|---|---|---|---|
| | | 1 | 2 | 3 | 4 | |
| | 1 | 1/8 ³ | 1/16 ⁴ | 1/32 ⁵ | 1/32 ⁵ | 1/4 |
| $y$ | 2 | 1/16 ⁴ | 1/8 ³ | 1/32 ⁵ | 1/32 ⁵ | 1/4 |
| | 3 | 1/16 ⁴ | 1/16 ⁴ | 1/16 ⁴ | 1/16 ⁴ | 1/4 |
| | 4 | 1/4 ² | 0 | 0 | 0 | 1/4 |
| $P(x)$ | | 1/2 | 1/4 | 1/8 | 1/8 | |
| $h(x)$ | | 1 | 2 | 3 | 3 | |

$$H(Y) = 2, \quad H(x) = \frac{1}{2} + \frac{2}{4} + \frac{3}{8} + \frac{3}{8} = \frac{7}{4}$$

$$\Rightarrow H(x,y) = \frac{1}{4}\cdot 2 + \frac{2}{8}\cdot 3 + \frac{6}{16}\cdot 4 + \frac{4}{32}\cdot 5$$

$$= \frac{4 + 6 + 12 + 5}{8} = \frac{27}{8}$$

$$H(x) + H(y) = \frac{30}{8} \geq H(x,y)$$

$P(x|y)$

| | $x$ | | | | $H(x|y=y)$ |
|---|---|---|---|---|---|
| | 1 | 2 | 3 | 4 | |
| 1 | 1/2 | 1/4 | 1/8 | 1/8 | 1 |
| $y$ 2 | 1/4 | 1/2 | 1/8 | 1/8 | 1 |
| 3 | 1/4 | 1/4 | 1/4 | 1/4 | |
| 4 | 1 | 0 | 0 | 0 | |

$P(y|x)$

| | 1 | 2 | 3 | 4 |
|---|---|---|---|---|
| 1 | 1/4 | 1/4 | 1/4 | 1/4 |
| 2 | 1/8 | 1/2 | 1/4 | 1/4 |
| $y$ 3 | 1/8 | 1/4 | 1/2 | 1/2 |
| 4 | 1/2 | 0 | 0 | 0 |

$$\Rightarrow \boxed{I(x;y) = 3/8}$$

# mutual information and KL divergence

**mutual information**

for any rvs $X, Y$:    $I(X; Y) = H(X) - H(X|Y) = H(Y) - H(Y|X)$
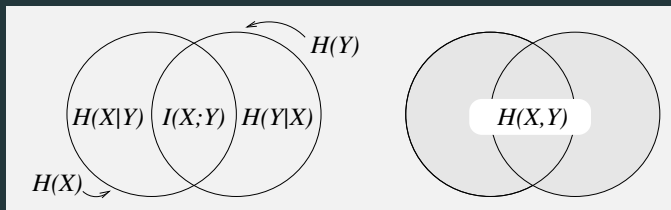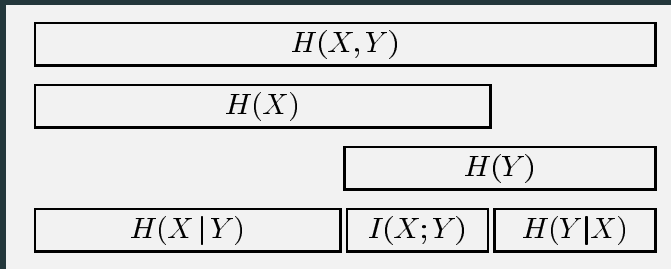
$$I(X; Y) = D_{KL}\left(P(x,y) \| P(x)P(y)\right)$$

True dist
of $(X,Y)$

Dist$^n$ of $(X,Y)$ if
$X \perp\!\!\!\perp Y$

increase in codesize if you encode
$P(x,y)$ using optimal code for $P(x)P(y)$

## visualizing mutual information

chapter 9 in Mackay

**channel coding**

# mutual information for the BSC

**Binary symmetric channel.** $\mathcal{A}_X = \{0, 1\}$. $\mathcal{A}_Y = \{0, 1\}$.

$x \overset{0 \,\overset{1-f}{\longrightarrow}\, 0}{\underset{1 \,\underset{1-f}{\longrightarrow}\, 1}{\times}} y$

$$P(y=0 \mid x=0) = 1-f; \quad P(y=0 \mid x=1) = f;$$
$$P(y=1 \mid x=0) = f; \quad P(y=1 \mid x=1) = 1-f.$$

$f = 0.1$

assume input distribution $\mathcal{P}_X = \{1-p, p\} = \{P_0, P_1\}$, $P_0 + P_1 = 1$

| X \ Y | 0 | 1 | $P(x)$ |
|---|---|---|---|
| 0 | $P_0(1-f)$ | $P_0 f$ | $P_0$ |
| 1 | $P_1 f$ | $P_1(1-f)$ | $P_1$ |

$\text{Ber}(p)$

$P(y)$: $P_0(1-f) + P_1 f$ , $P_0 f + P_1(1-f)$

$\sim \text{Ber}\left(P_1 f + P_0(1-f)\right) \sim \text{Ber}\left(P_0 f + P_1(1-f)\right)$

$\underbrace{\phantom{xxxx}}_{q} \qquad \underbrace{1-q}$

$$I(x; y) = H(y) - H(y \mid x)$$

$$= h_2(q) -$$

$P(x=1) H(y \mid x=1) + P(x=0) H(y \mid x=0)$

$$= h_2(q) - P_1 h_2(f) - P_0 h_2(f)$$

$$\boxed{= h_2(q) - h_2(f)}$$

$$\Rightarrow I(x;y) = \underbrace{h_2\left(P_1 f + P_0(1-f)\right)}_{H(Y)} - \underbrace{h_2(f)}_{H(Y|x)}$$

<u>Exercise</u> -   $I(x;y) = H(x) - H(x|y)$

$$= h_2(P_1) - \quad ?$$

# mutual information for the Z-channel

**Z channel.** $\mathcal{A}_X = \{0, 1\}$. $\mathcal{A}_Y = \{0, 1\}$.

$$x \begin{array}{c} 0 \rightarrow 0 \\ \nearrow \\ 1 \rightarrow 1 \end{array} y$$

$$P(y=0 \,|\, x=0) = 1; \quad P(y=0 \,|\, x=1) = f;$$
$$P(y=1 \,|\, x=0) = 0; \quad P(y=1 \,|\, x=1) = 1-f.$$

assume input distribution $\mathcal{P}_X = \{1-p, p\} = \{P_0, P_1\}$

$$I(x;y) = H(y) - H(y|x)$$

$$y = \begin{cases} 0 & \text{wp } P_0 + P_1 f \\ 1 & \text{wp } P_1(1-f) \end{cases} = h_2\big(P_1(1-f)\big) - P_0 \underbrace{H(y|x=0)}_{0}$$

$$- P_1 \underbrace{H(y|x=1)}_{h_2(f)}$$

$$= \boxed{h_2\big(P_1(1-f)\big) - P_1 \, h_2(f)}$$

<u>Alt-</u> $I(x;y) = H(x) - H(x|y)$

$$= h_2(P_1) - \underbrace{H(x|y=1)}_{=0}\overbrace{\mathbb{P}[y=1]}^{P_1(1-f)}$$

$$- \underbrace{H(x|y=0)}_{?}\ \overbrace{\mathbb{P}[y=0]}^{P_0 + P_1 f}$$

Bayes Thm



$P_0$ → $x=0$ → $1-f$ → $y=0$
(Start)
$P_1$ → $x=1$ → $f$ → $y=0$

$$P_{x|y=0} \sim \left\{ \underbrace{\frac{P_0(1-f)}{q_1}}_{\mathbb{P}[y=1]}, \frac{P_1 f}{q_1} \right\}$$

$$= h_2(P_1) - \left(P_0 + P_1 f\right) h_2(\theta)$$

where $\theta = \dfrac{P_1 f}{P_1 f + P_0(1-f)}$

# mutual information for the erasure channel

**Binary erasure channel.** $\mathcal{A}_X = \{0, 1\}$. $\mathcal{A}_Y = \{0, ?, 1\}$.



$$P(y=0 \,|\, x=0) = 1-f; \quad P(y=0 \,|\, x=1) = 0;$$
$$P(y=? \,|\, x=0) = f; \quad P(y=? \,|\, x=1) = f;$$
$$P(y=1 \,|\, x=0) = 0; \quad P(y=1 \,|\, x=1) = 1-f.$$

assume input distribution $\mathcal{P}_X = \{1-p, p\} = \{P_0, P_1\}$

$$I(x;y) = H(y) - H(y|x)$$

$\because H(y|x)$
$= \sum P(x) H(y|x=x)$
$= P_0 H(y|x=0) + P_1 H(y|x=1)$

However, both $H(y|x=1)$
& $H(y|x=0)$ are $h_2(f)$

$= h_2(f) + (1-f) h_2(P_1)$

$= \left[ P_0 h_2(f) - P_1 h_2(f) \right.$

$= \boxed{(1-f) h_2(P_1)}$

Note - $I(x;y)$ separates into terms

depending on $f$ and on $P_1$

$$y = \begin{cases} 0 & \text{wp } P_0(1-f) \\ ? & \text{wp } f \\ 1 & \text{wp } P_1(1-f) \end{cases}$$

$$H(y) = P_0(1-f) \log_2\left(\frac{1}{P_0(1-f)}\right)$$
$$+ P_1(1-f) \log_2\left(\frac{1}{P_1(1-f)}\right) + f \log_2 \frac{1}{f}$$

$$= h_2(f) + (1-f) h_2(P_1)$$

Let $Z = \mathbb{1}\{Y = ?\}$, then
$H(Y) = H(Z) + H(Y|Z)$
(since $Z$ if a fn of $Y$)

# capacity of a channel

the capacity of a channel $\mathcal{Q}$, with input $\mathcal{A}_{\mathcal{X}}$ and output $\mathcal{A}_{\mathcal{Y}}$, is

$$C(\mathcal{Q}) = \max_{\mathcal{P}_X} I(X; Y)$$

any $\arg\max \mathcal{P}_X^\star$ is called the optimal input distribution

**Shannon's channel coding theorem**

can communicate $\leq C$ bits of information per channel use without error!

Mackay- Ch 9 (Defines capacity, channel coding)

Ch 10 (Pf of Shannon's coding thm)

# capacity of the BSC

**Binary symmetric channel.** $\mathcal{A}_X = \{0, 1\}$. $\mathcal{A}_Y = \{0, 1\}$.

$$x \begin{array}{c} 0 \rightarrow 0 \\ \times \\ 1 \rightarrow 1 \end{array} y \qquad \begin{array}{lll} P(y=0 \mid x=0) & = & 1-f; \\ P(y=1 \mid x=0) & = & f; \end{array} \quad \begin{array}{lll} P(y=0 \mid x=1) & = & f; \\ P(y=1 \mid x=1) & = & 1-f. \end{array}$$

assume input distribution $\mathcal{P}_X = \{1 - p, p\}$

$f = 0.1$



$I(X;Y)$

$$I(x;y) = h_2(q) - h_2(f)$$

$$q = p_1 f + p_0 (1-f)$$

$$C(f) = \max I(x;y)$$
$$(p_0, p_1) \text{ s.t } p_0 + p_1 = 1$$

$$\Rightarrow P^* = \text{ set } q = \tfrac{1}{2} \Rightarrow p_1 f + (1-p_1)(1-f) = \tfrac{1}{2} \Rightarrow p_1^* = p_0^* = \tfrac{1}{2}$$

$$\Rightarrow \boxed{C(f) = 1 - h_2(f)}$$

## capacity of the Z-channel

**Z channel.** $\mathcal{A}_X = \{0, 1\}$. $\mathcal{A}_Y = \{0, 1\}$.

$$x \quad \begin{array}{c} 0 \longrightarrow 0 \\ \diagup \\ 1 \longrightarrow 1 \end{array} \quad y$$

$P(y=0 \mid x=0) = 1; \quad P(y=0 \mid x=1) = f;$
$P(y=1 \mid x=0) = 0; \quad P(y=1 \mid x=1) = 1-f.$

assume input distribution $\mathcal{P}_X = \{1 - p, p\} = \{P_0, P_1\}$



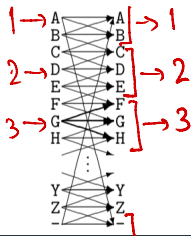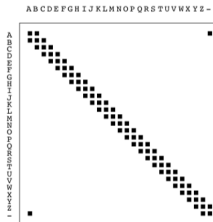$I(X;Y) = h_2\left(P_1(1-f)\right) - P_1 h_2(f)$

This is not symmetric in $P_1$ - Some what complicated to maximize (see fig)

# the noisy typewriter

**Noisy typewriter**. $\mathcal{A}_X = \mathcal{A}_Y$ = the 27 letters {A, B, ..., Z, -}. The letters are arranged in a circle, and when the typist attempts to type B, what comes out is either A, B or C, with probability $1/3$ each; when the input is C, the output is B, C or D; and so forth, with the final letter '-' adjacent to the first letter A.



$$P(y=\text{F} \mid x=\text{G}) = 1/3;$$
$$P(y=\text{G} \mid x=\text{G}) = 1/3;$$
$$P(y=\text{H} \mid x=\text{G}) = 1/3;$$

$$I(x;Y) = H(Y) - H(Y|X)$$

$$\leq \log |\mathcal{A}_Y| \qquad = \sum_x p(x) H(Y|X=x) = \log_2 3$$

$$\leq \log_2 27 \sim \log_2 3 = \log_2 9$$

Can be achieved, Eg, set $P_x = (1/27, 1/27, \cdots 1/27)$

$$C(Q) = \max_{P_X} I(x;y) = \log_2 9 \text{ bits}$$

$$\left( P_X^* \text{ can be } \{1/27, 1/27, \ldots, 1/27\} \right)$$
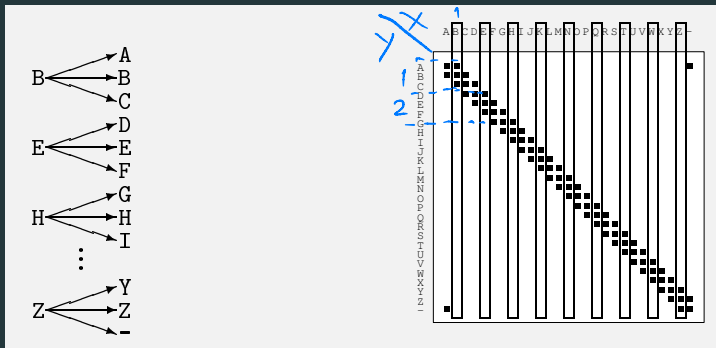
• Code for noisy typewriter: $\phi: \{1, 2, \ldots, 9\} \rightarrow \{A, B, \ldots, Z, -\}$

**Encoder** $\phi(1) = A$, $\phi(2) = D$, $\phi(3) = G$ ...

**Decoder** $\phi^{-1}(\{-, A, B\}) = 1$, $\phi^{-1}(\{C, D, E\}) = 2$, ...

Can send $\log_2 9$ bits per channel use without error

# another view of the noisy typewriter



Syndrome decoding – map set of outputs to same input

# expanded channel for the BSC

**Binary symmetric channel.** $\mathcal{A}_X = \{0, 1\}$. $\mathcal{A}_Y = \{0, 1\}$.

$$x \begin{matrix} 0 \to 0 \\ 1 \to 1 \end{matrix} y$$

$$P(y=0 \mid x=0) = 1-f; \quad P(y=0 \mid x=1) = f;$$
$$P(y=1 \mid x=0) = f; \quad \quad P(y=1 \mid x=1) = 1-f.$$

## expanded channel for the Z-channel



**Z channel.** $\mathcal{A}_X = \{0, 1\}$. $\mathcal{A}_Y = \{0, 1\}$.

$$x \begin{array}{c} 0 \to 0 \\ 1 \to 1 \end{array} y$$

$$P(y=0 \,|\, x=0) = 1; \quad P(y=0 \,|\, x=1) = f;$$
$$P(y=1 \,|\, x=0) = 0; \quad P(y=1 \,|\, x=1) = 1-f.$$

$N=1$  $N=2$  $N=4$

## typical set
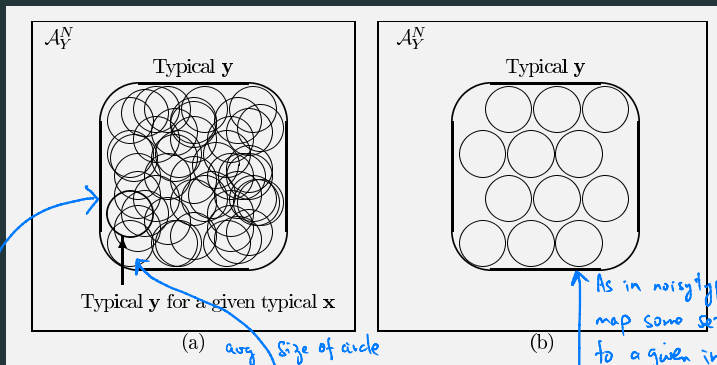
iid source produces $X^N = (X_1 X_2 \ldots X_N)$; each $X_i \in \mathcal{X}$ has entropy $H(X)$

then $X^N$ is very likely to be one of $\approx 2^{NH(X)}$ typical strings,
all of which have probability $\approx 2^{NH(X)}$

Recall- bent coin lottery

- $X_1 X_2 \ldots X_{1000} \sim Bin(1000, f)$

- Most of the time, # of 1s $= 1000f \mp \sqrt{1000f}$

# typical set and non-confusable subset



(a)

(b)

$\mathcal{A}_Y^N$

Typical $\mathbf{y}$

Typical $\mathbf{y}$ for a given typical $\mathbf{x}$

$\mathcal{A}_Y^N$

Typical $\mathbf{y}$

**Handwritten annotations:**

Union of all circles

avg size of circle

As in noisy typewriter, we map some set of outputs to a given input.

#of elements in the typical set of outputs $\approx 2^{NH(Y)}$

# of typical outputs for typical input X $\approx 2^{NH(Y|X)}$

# of non-overlapping circles $\approx 2^{NH(Y)} / 2^{NH(Y|X)}$
$= 2^{NI(X;Y)} \leq 2^{NC}$

## block codes, encoding, decoding

### block code

for channel $\mathcal{Q}$ with input $\mathcal{A}_X$, an $(N, K)$-block code is a list of $\mathcal{S} = 2^K$ codewords $\{x^{(1)}, x^{(2)}, \ldots, x^{(2^K)}\}$ with $x^{(i)} \in \mathcal{A}_X^N$ (i.e., of length $N$)

### encoder

– using $(N, K)$-block code, can encode signal $s \in \{1, 2, 3, \ldots, 2^K\}$ as $x(s)$
– the rate of the code is $R = N/K$ bits per channel use

### decoder

– mapping from each length-N string $y \in \mathcal{A}_Y^N$ of channel outputs to a codeword label $\hat{s} \in \{\varphi, 1, 2, 3, \ldots, 2^K\}$ as $x(s)$
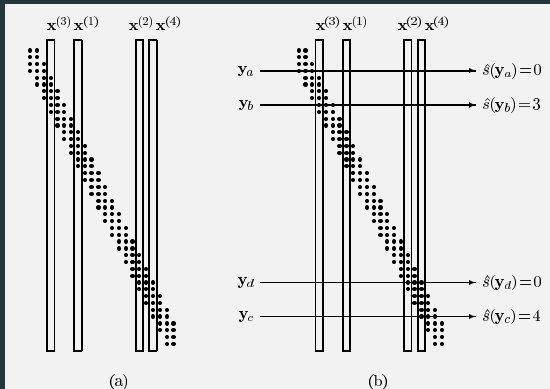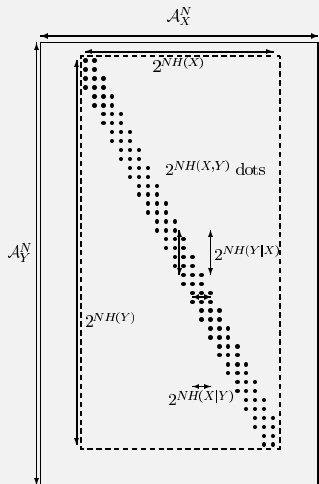– $\varphi$ indicates failure

# block codes and capacity

## block code

for channel $\mathcal{Q}$ with input $\mathcal{A}_X$, an $(N, K)$-block code is a list of $\mathcal{S} = 2^K$ codewords $\{x^{(1)}, x^{(2)}, \ldots, x^{(2^K)}\}$ with $x^{(i)} \in \mathcal{A}_X^N$ (i.e., of length $N$)
– the rate of the code is $R = N/K$ bits per channel use

## Shannon's channel coding theorem

For any $\epsilon > 0$ and $R < C$, for large enough $N$, there exists a block code of length $N$ and rate $\geq R$ such that probability of block error is $< \epsilon$.

- Only transmit $2^k$ out of $|A_x|^N$ symbols

**Binary erasure channel.** $\mathcal{A}_X = \{0, 1\}$. $\mathcal{A}_Y = \{0, ?, 1\}$.

$$x \quad \begin{matrix} 0 \to 0 \\ \searrow ? \\ 1 \to 1 \end{matrix} \quad y$$

$$\begin{aligned} P(y=0 \mid x=0) &= 1-f; & P(y=0 \mid x=1) &= 0; \\ P(y=? \mid x=0) &= f; & P(y=? \mid x=1) &= f; \\ P(y=1 \mid x=0) &= 0; & P(y=1 \mid x=1) &= 1-f. \end{aligned}$$
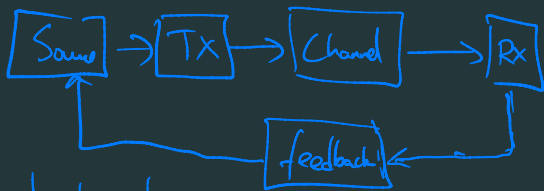
$$I(x;y) = (1-f)\, h_2(p_1) \qquad P_x = \{p_0, p_1\}$$

$$\Rightarrow \quad C = \max_{\{p_0, p_1 \mid p_0 + p_1 = 1\}} I(x;y) = 1-f \quad \text{for } P_x = \{\tfrac{1}{2}, \tfrac{1}{2}\}$$

How can we design a scheme to achieve this?

# feedback coding

<u>Idea</u> - Suppose we have feedback from the receiver



Source → TX → Channel → Rx → feedback

- <u>Code</u> - If Rx gets ?, asks for a repeat character (retransmit /ACK protocol)

- $P[\text{bit received correctly}] = 1-f$

$\Rightarrow$ # of retransmissions $\sim Geom(1-f)$ $\quad \mathbb{E}\left[\begin{smallmatrix}\#of \\ retx\end{smallmatrix}\right] = \dfrac{1}{1-f}$

(can do without feedback - fountain codes)