Lecture 2

1/6/2010

overview of classical information theory
(descriptive fashion)

## Data Compression Example

Alice + Bob are connected by a noiseless
bit channel — $p(y|x) = \delta_{y,x}$

Alice has four symbols $\{a, b, c, d\}$

Suppose either she or someone else is
choosing them at random according to

$$Pr\{a\} = 1/2$$
$$Pr\{b\} = 1/8$$
$$Pr\{c\} = 1/4$$
$$Pr\{d\} = 1/8$$

chooses symbols independently

sps. channel accepts only bits.

Naive code is

$a \rightarrow 00$ , $b \rightarrow 01$ , $c \rightarrow 10$ , $d \rightarrow 11$

Performance measured by expected length

In this case, expected length is

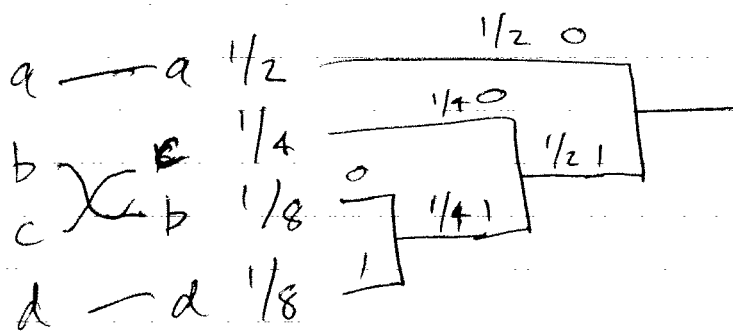$$\tfrac{1}{2} \cdot 2 + \tfrac{1}{8} \cdot 2 + \tfrac{1}{4} \cdot 2 + \tfrac{1}{8} \cdot 2 = 2$$

But they can do better b/c dist. is skewed

— use shorter codewords for more likely symbols & longer for less likely

— can do Huffman code

$a \longrightarrow a \quad 1/2$

$b \longrightarrow \quad 1/4$

$c \longrightarrow b \quad 1/8$

$d \longrightarrow d \quad 1/8$

$1/2 \quad 0$

$1/4 \quad 0$

$1/2 \quad 1$

$0 \quad 1/4 \quad 1$

$1$

$a \rightarrow 0$

$b \rightarrow 110$

$c \rightarrow 10$

$d \rightarrow 111$

Any coded sequence is uniquely decodable

E.g.)

$\underline{0}\ \underline{0}\ \underline{1\ 1\ 0}\ \underline{1\ 0}\ \underline{1\ 1\ 1}\ \underline{0\ 1}\ \underline{0\ 1}\ \underline{0}\ \underline{0}\ \underline{0\ 1\ 0}$

a  a    b    c    d   a c   c a a c

Expected length of coding scheme

$$\frac{1}{2} \cdot 1 + \frac{1}{8} \cdot 3 + \frac{1}{4} \cdot 2 + \frac{1}{8} \cdot 3 = \frac{7}{4}$$
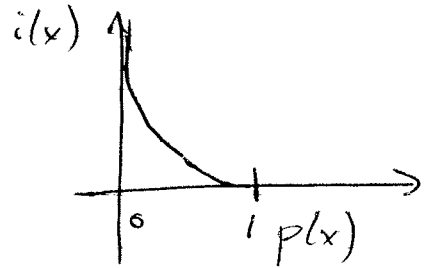
scheme suggests a measure of information

(want a measure that is higher for less likely events + lower for more likely)

one measure is

$$i(x) \equiv \log_2\left(\frac{1}{p(x)}\right)$$

called "information content"



— it so happens that the length of each symbol in Huffman code is equal to its information content
— scheme is special because we had powers of two

info. content is additive for independent events
— our info. source is memoryless

$$\Rightarrow i(x_1, x_2) = -\log\left(p(x_1, x_2)\right)$$

$$= -\log\left(p(x_1)\right) - \log\left(p(x_2)\right) = i(x_1) + i(x_2)$$

property is very important

expected information content is

$$\sum_x p(x) \, i(x) = -\sum_x p(x) \log(p(x))$$

this is the _entropy_ of the source

## Shannon's source coding theorem

What is the ultimate limit on the compressibility of information?

– need more general techniques & setting to answer this question in a satisfying way

– given random variable $X$, the information content of it is itself a random variable

$$i(X) = -\log\left(p_X(x)\right).$$

seemingly self-referential, but OK

on to source coding,

could associate a binary codeword for each $x$ as before, but this ~~scheme~~ may lose some efficiency ~~if~~ if alphabet or probabilities are not a power of two.
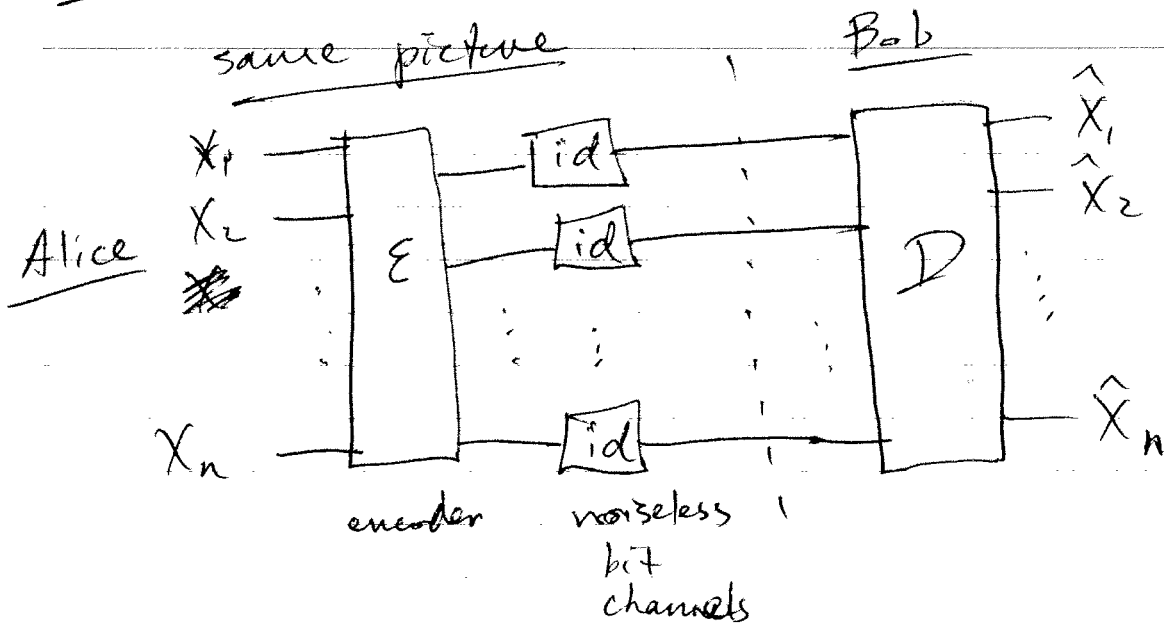
Shannon's idea: 1) block coding

Source emits large number of realizations & then code data as emitted block

2) allow for slight error, but show that it vanishes asymptotically

same picture

Bob

Alice

$X_1$
$X_2$

$X_n$

$\mathcal{E}$

id

id

id

$D$

$\hat{X}_1$
$\hat{X}_2$

$\hat{X}_n$

encoder    noiseless
bit
channels

Consider the distribution of the source
under the IID assumption

$$P_{X^n}(x^n) = P_{X_1, \ldots, X_n}(x_1, \ldots, x_n)$$

$$= P_{X_1}(x_1) \cdots P_{X_n}(x_n)$$

$$= P_X(x_1) \cdots P_X(x_n)$$

$$= \prod_{i=1}^{n} P_X(x_i)$$

Let $a_1, \ldots, a_{|\mathcal{X}|}$ denote letters in $\mathcal{X}$
+ Let $N(a_i | x^n)$ denote number of occurrences of
$a_i$ in $x^n$
then

$$P_{X^n}(x^n) = \prod_{i=1}^{n} P_X(x_i) = \prod_{j=1}^{|\mathcal{X}|} P_X(a_j)^{N(a_j | x^n)}$$

probability for a
particular sequence

Consider information content of

~~the the~~ a <u>random</u> sequence:

$$\frac{i(X^n)}{n} = -\frac{1}{n} \log\left(P_{X^n}(X^n)\right) \quad \left(\begin{array}{c} \text{AKA sample} \\ \text{entropy} \end{array}\right)$$

↑ random variable

we can use formula from before &

random quantity $N(a_j | X^n)$

$$-\frac{1}{n} \log\left(P_{X^n}(X^n)\right) = -\frac{1}{n} \log\left(\prod_{j=1}^{|X|} P_X(a_j)^{N(a_j|X^n)}\right)$$

$$= -\frac{1}{n} \sum_{j=1}^{|X|} \log\left(P_X(a_j)^{N(a_j|X^n)}\right)$$

$$= \sum_{j=1}^{|X|} -\left(\frac{1}{n} N(a_j|X^n) \log\left(P_X(a_j)\right)\right)$$

<u>state law of large numbers</u>

$$\lim_{n\to\infty} Pr\left\{|\overline{X^n} = \mu\xi\right\} = 1$$

where $\overline{X^n}$

implies $\dfrac{N(a_j|X^n)}{n} \longrightarrow P_X(x)$

$$\Rightarrow \lim_{n\to\infty} Pr\left\{\left|\frac{i(X^n)}{n} - H(x)\right| < \epsilon\right\} = 1$$

"It is highly likely that the source emits a sequence w/ sample entropy close to the true entropy."

Leads to the notion of "high probability set"

or "typical set"

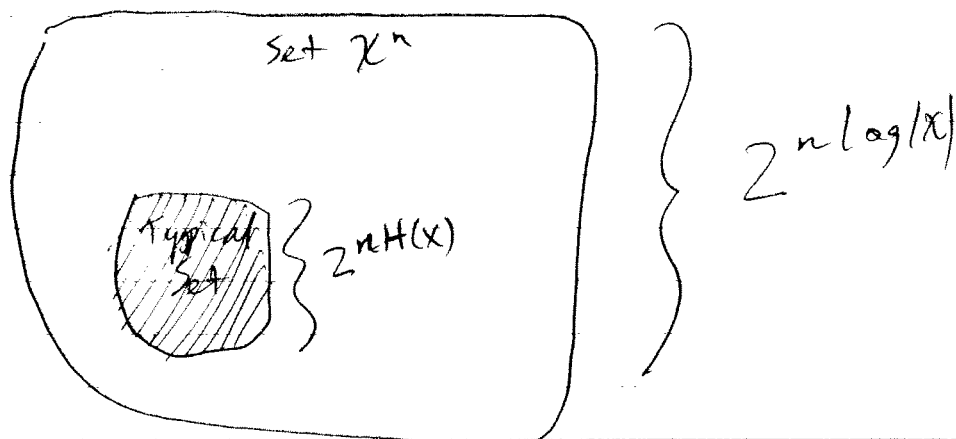set of sequences for which the sample entropy is close to the true entropy

asymptotically, this set has all of the probability.

Since we're concerned w/ probability of error in communication, it seems reasonable to give attention only to the high probability set of sequences

- wonderful thing about typical set is that its size is $2^{nH(X)}$ whereas set of all sequences is $|X|^n = 2^{n \log |X|}$

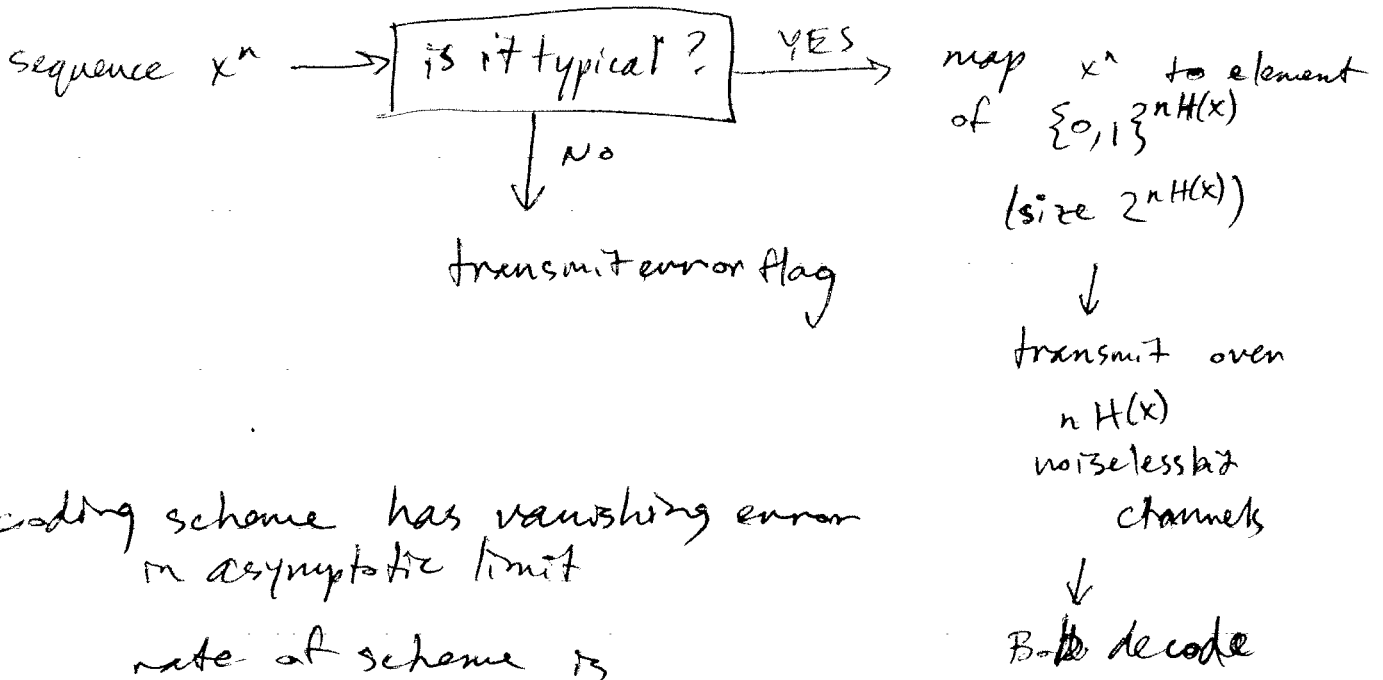- probability is concentrating on an exponentially small set

# Coding Strategy

"Keep the typical sequences + throw away
the rest"
(b/c the others' probability is negligible)

sequence $x^n$ → | is it typical ? | → YES → map $x^n$ to element
of $\{0,1\}^{nH(x)}$

(size $2^{nH(x)}$)

↓ No

transmit error flag

↓

transmit over
$nH(x)$
noiseless bit
channels

↓

Bob decode

coding scheme has vanishing error
in asymptotic limit

rate of scheme is

$$\frac{\# \text{ channel bits}}{\# \text{ source symbols}} = \frac{nH(x)}{n} = H(x)$$

any rate $R > H(x)$ is achievable in the sense that
$\forall \epsilon > 0$ + sufficiently large n  there exists
an $(n, R, \epsilon)$ compression code

this is called direct part of coding theorem

demonstrate an achievable strategy

If rate > entropy, there exists achievable scheme

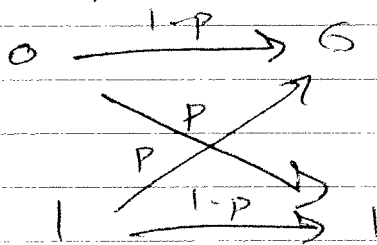Converse Part (optimality)

If ∃ achievable scheme, then
rate > entropy.

Simple Example of Error Correction

binary symmetric channel

$$0 \xrightarrow{\quad 1-p \quad} 0$$

$$1 \xrightarrow{\quad 1-p \quad} 1$$

(with crossing $p$ paths)

transmit data as is, then error probability
is
$$p(e) = p(e|0)\,p(0) + p(e|1)\,p(1)$$
$$= p$$

would like to supress errors

can try to engineer better channel, but we
would like a "systems engineering" solution

use a code

rate is $1/3$

$$0 \rightarrow 0\,0\,0$$

$$1 \rightarrow 1\,1\,1$$

take a majority vote at output

.if Alice sends 0, then

| output | Prob | Sum | |
|--------|------|-----|---|
| 000 | $(1-p)^3$ | $(1-p)^3$ | $\}$ correct |
| 001, 010, 100. | $p(1-p)^2$ | $3p(1-p)^2$ | |
| 110, 101, 011 | $p^2(1-p)$ | $3p^2(1-p)$ | $\}$ error |
| 111 | $p^3$ | $p^3$ | $\}$ |

$$p(e) = 3p^2(1-p) + p^3$$

$$= 3p^2 - 2p^3$$

coding is helping if

$$3p^2 - 2p^3 < p$$

or equivalently

$$0 < p(2p-1)(p-1)$$

only values of $p$ satisfying are

$$0 < p < 1/2$$

error prob. goes like $O(p^2)$

How to reduce even more?
concatenate...

$$0 \to 000 \to 000 \ \ 000 \ \ 000$$

$$1 \to 111 \to 111 \ \ 111 \ \ 111$$

rate $1/9$

error prob. goes like

$$3p^2(e) - 2p^3(e) = O(p^4)$$

we would like error free communication...

keep on concatenating

rate is

$$\frac{1}{3^n}$$    for error suppression $O(p^{2n})$

rate goes to zero to make

error prob go to zero?

---

## Shannon's Channel Coding Theorem

Alice chooses message from a set

$$[M] \equiv \{1, ..., M\}$$

(chooses uniformly at random)

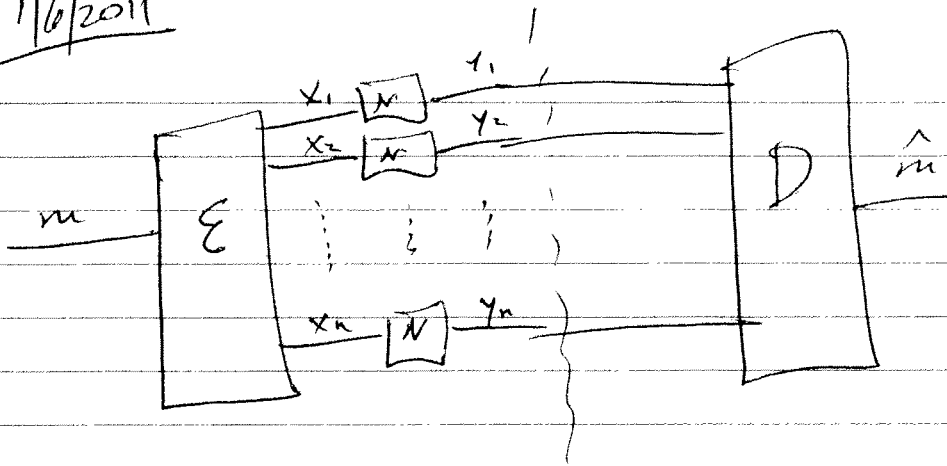set requires $\log_2 M$ bits to represent it

one "layer of randomness"

---

noisy channel $N \equiv P_{Y|X}(y|x)$

IID

$$P_{Y^n|X^n}(y^n|x^n) \equiv \prod_{i=1}^{n} P_{Y|X}(y_i|x_i)$$

for every message $m$, there is some
codeword $x^n(m)$ & Bob makes
his best estimate of $x^n(m)$
from received sequence $y$

$$\text{rate} = \frac{\#\text{ message bits}}{\#\text{ channel uses}} = \frac{\log M}{n}$$

capacity = highest rate of reliable comm.

$$C \equiv \{x^n(m)\}_{m \in [M]}$$

$p_e(m, c) = $ prob. of error for message $m$
under code $c$

$$\overline{p_e}(c) = \frac{1}{M} \sum_{m=1}^{M} p_e(m, c) = \text{avg. prob. of error}$$

$$p_e^*(c) = \max_m \; p_e(m, c) = \text{maximal prob. of error}$$

These error criteria are difficult to analyze.

Shannon's idea: choose the code randomly & analyze the <u>expectation</u> of the average error prob. rather

use some distribution $P_X(x)$

for message 1, choose $x_1(1)$ according to $p(x)$

$\qquad\qquad\qquad\qquad x_2(1)$ " " $p(x)$

$\qquad\qquad\qquad\qquad\vdots$

$\qquad\qquad\qquad\qquad x_n(1)$ " " $p(x)$

$$x^n(1) = x_1(1) \, x_2(1) \cdots x_n(1)$$

same for $x^n(2)$ (message 2)

$$x^n(2) = x_1(2) \, x_2(2) \cdots x_n(2)$$

$$\vdots$$

$$x^n(M) = x_1(M) \, x_2(M) \cdots x_n(M)$$

every codeword chosen independently of the message $m$

probability for a particular code $C_0$ is

$$P_C(C_0) = \prod_{m=1}^{M} \prod_{i=1}^{n} P_X(x_i(m))$$

expectation of the average error probability

$$\mathbb{E}_C \left\{ \overline{P_e}(c) \right\} = \mathbb{E}_C \left\{ \frac{1}{M} \sum_{m=1}^{M} P_e(m,c) \right\}$$

$$= \frac{1}{M} \sum_{m=1}^{M} \underbrace{\mathbb{E}_C \left\{ P_e(m,c) \right\}}$$

expectation does not depend on particular message m

$$= \frac{1}{M} \sum_{m=1}^{M} \mathbb{E}_C \left\{ P_e(1,c) \right\}$$

$$= \underbrace{\mathbb{E}_C \left\{ P_e(1,c) \right\}}$$

much easier to analyze

show that

$$\mathbb{E}_C \left\{ P_e(1,c) \right\} \leq \epsilon$$

$\Rightarrow$ there exists some code $c_0$ such that

$$\overline{P_e}(c) \leq \epsilon$$

can use this to bound $p_e^*(c) \leq 2\epsilon$

---

What about size of code?

rate is $R = \frac{\log M}{n}$

so think of message set size as

$$M = 2^{nR}$$

Alice is using distribution $p(x)$
to generate sequences $x^n$.

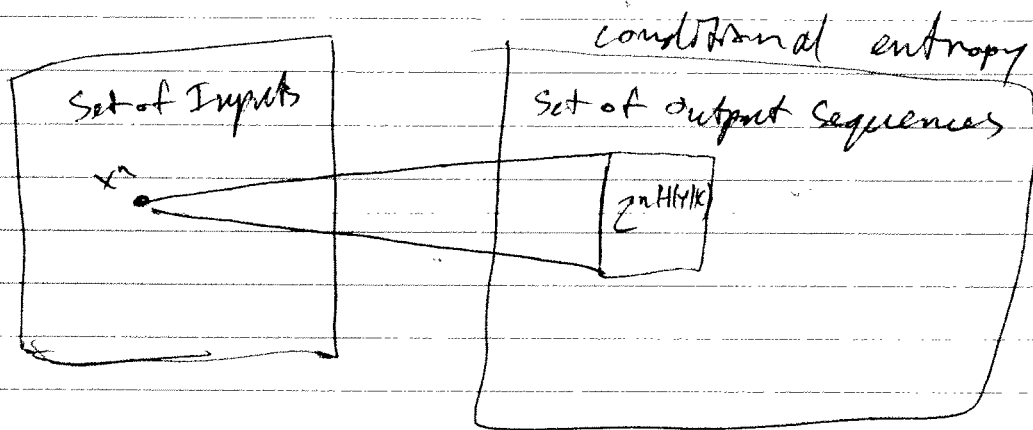Highly likely that $x^n$ is typical

$$x^n \longrightarrow P_{Y^n|X^n}(y^n|x^n) \longrightarrow Y^n$$

number of output sequences <u>likely</u> to

correspond to $x^n$

<u>tool</u>: conditional typicality

set of size $2^{nH(Y|X)}$

                      conditional entropy



Set of Inputs     Set of output sequences

$x^n$         $2^{nH(Y|X)}$

From Bob's perspective, distribution
he "sees"    is $p(y) = \sum\limits_{x} p(y|x) p(x)$

holds that $H(Y) \geq H(Y|X)$

Alice transmits codeword $x^n(m)$

Bob is ignorant of $m$ so from his perspective, sequences generated according to $p(y)$

1) Determine if $y^n$ lies in typical set
for $y$ of size $2^{n H(Y)}$

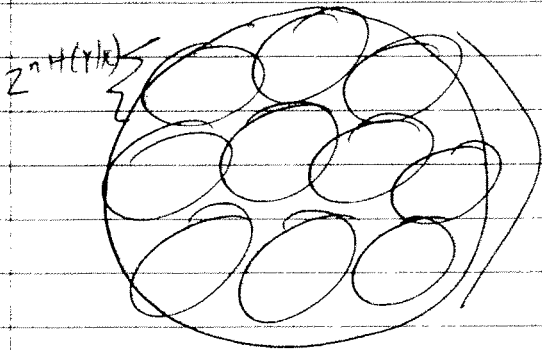If not, declare error. If so, proceed

2) Determine (w/ knowledge of code)
the conditionally typical set from
which $y^n$ would come.

If $y^n$ lies in wrong set, error occurs.

If not, done.

How many message can we decode?

Same question as how many ~~message~~ conditionally typical sets we can fit into typical set



$2^{n H(Y|X)}$

$2^{n H(Y)}$

$$2^{nR} \approx \frac{2^{n H(Y)}}{2^{n H(Y|X)}}$$

$$= 2^{n(H(Y) - H(Y|X))}$$

$$= 2^{n I(X;Y)}$$

choose best dist. for generating code
$R = \max_x I(X;Y)$