



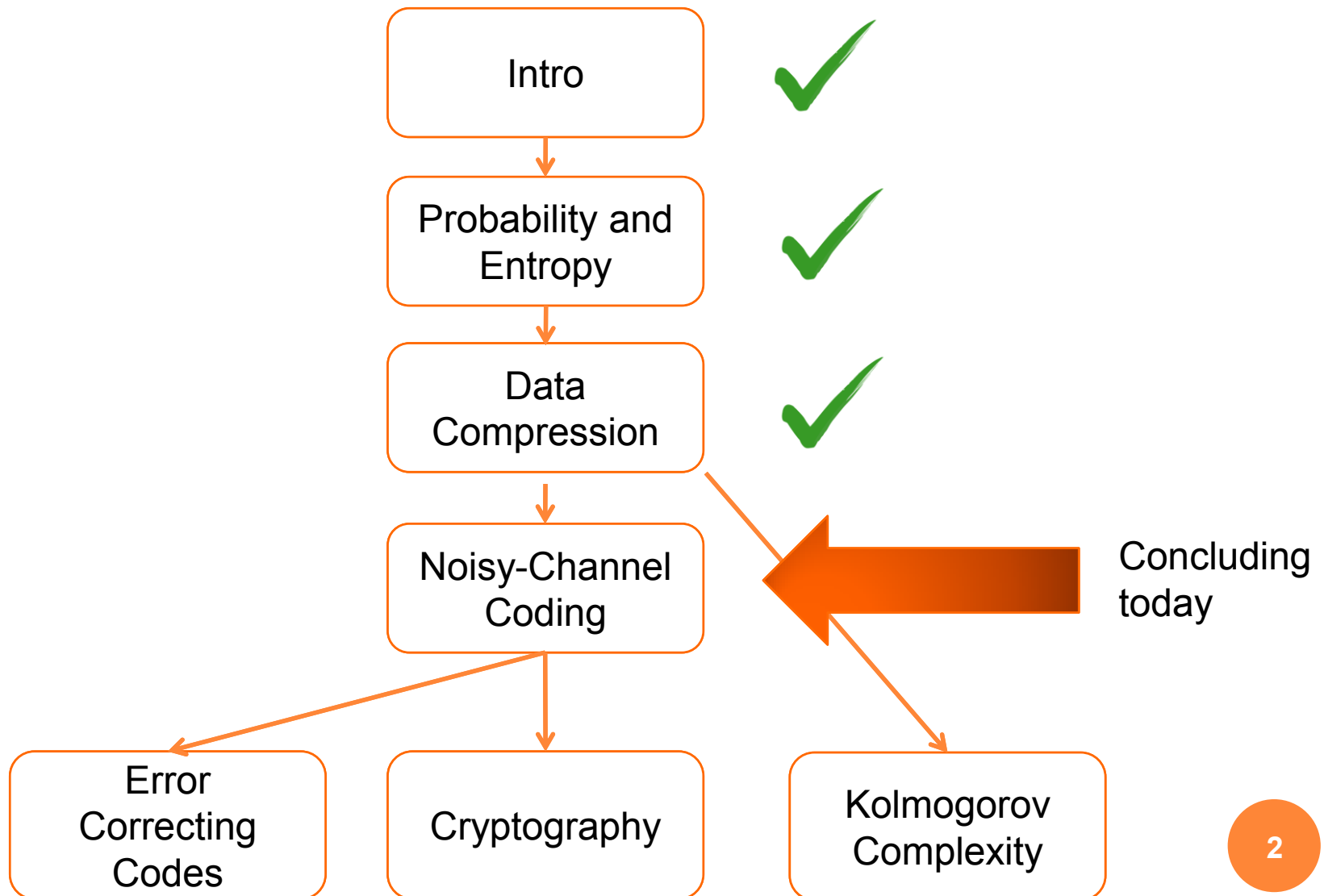
CS3236 INTRODUCTION TO INFORMATION THEORY

Lecture 7: Noisy-channel coding theorem

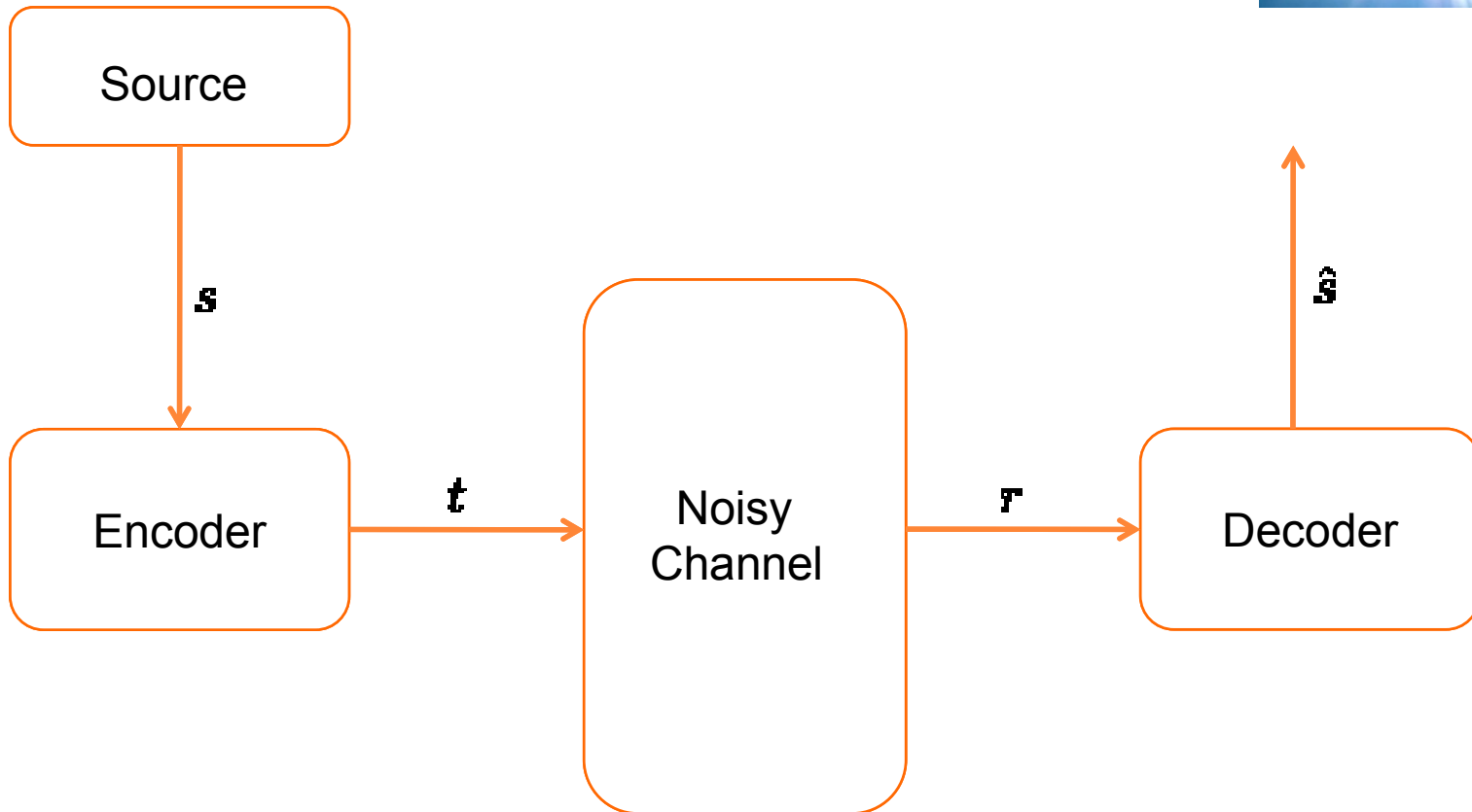
Course given by Pierre Senellart

Material by Stephanie Wehner, with additions by P. Senellart

WHERE DO WE GO FROM HERE?

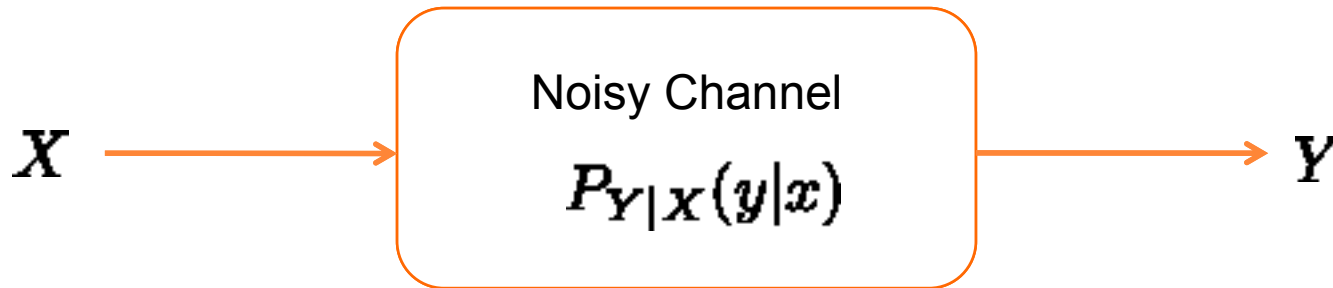


SENDING INFORMATION



Goal: Construct an encoder and decoder such that we can recover the information $\hat{s} = s$ with very low error

HOW MUCH INFORMATION DOES THE OUTPUT GIVE ABOUT THE INPUT?



- Intuition: mutual information measures how much our uncertainty about X is reduced by learning Y

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

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

CAPACITY OF A CHANNEL

- Shannon proved that the capacity of a channel \mathcal{N} given by

$$C(\mathcal{N}) = \max_{P_X} I(X; Y)$$

- Determines how well we can send information.

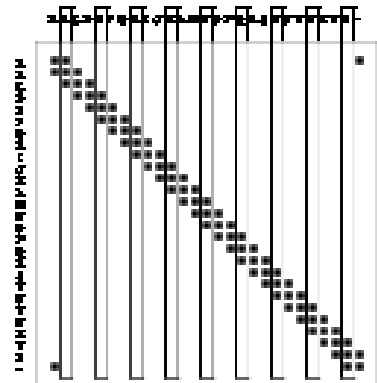
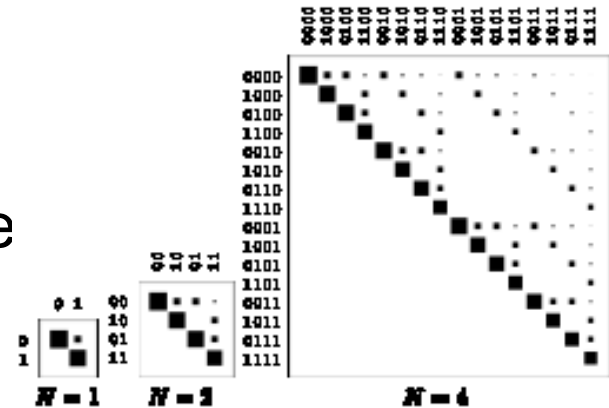
SHANNON'S NOISY CHANNEL CODING THEOREM (PART 1: ACHIEVABILITY)

- Associated with every discrete memoryless channel, there is a non-negative number C (the capacity) such that
 - For any error $\epsilon > 0$ and $R \leq C$ for large enough N , there exists a block code of length N and rate R , and a decoding algorithm, such that the maximum probability of block error is $< \epsilon$



ESSENTIAL IDEA

- For large N , the diagrams for the extended channel start to look a lot like the noisy-typewriter!
- Restricting to typical inputs and outputs allows us to do the same general channels



LET'S FIRST GATHER AN ESSENTIAL TOOL..
JOINT TYPICALITY!



TYPICAL INPUTS: X TYPICAL FOR P(X)

- Remember the lecture on the source coding theorem!
- We say that some iid input $\mathbf{x}_1, \dots, \mathbf{x}_N$ with probability $P(\mathbf{x}) = P(\mathbf{x}_1) \cdot \dots \cdot P(\mathbf{x}_N)$ is typical if

$$\left| \frac{1}{N} \log \frac{1}{P(\mathbf{x})} - H(X) \right| < \beta$$

- Size of the typical set was roughly $2^{NH(X)}$
- All elements of the typical set occur with roughly the same probability

TYPICAL OUTPUTS: Y TYPICAL FOR P(Y)

- For the outputs, we can also define the typical set!
- For iid input the output distribution is also iid

$$\left| \frac{1}{N} \log \frac{1}{P(\mathbf{y})} - H(Y) \right| < \beta$$

- Size of the typical set is roughly $2^{NH(Y)}$
- All elements of the typical set occur with roughly the same probability

JOINT TYPICALITY: (X,Y) TYPICAL FOR P(X,Y)

- Same idea, but now for the joint distribution

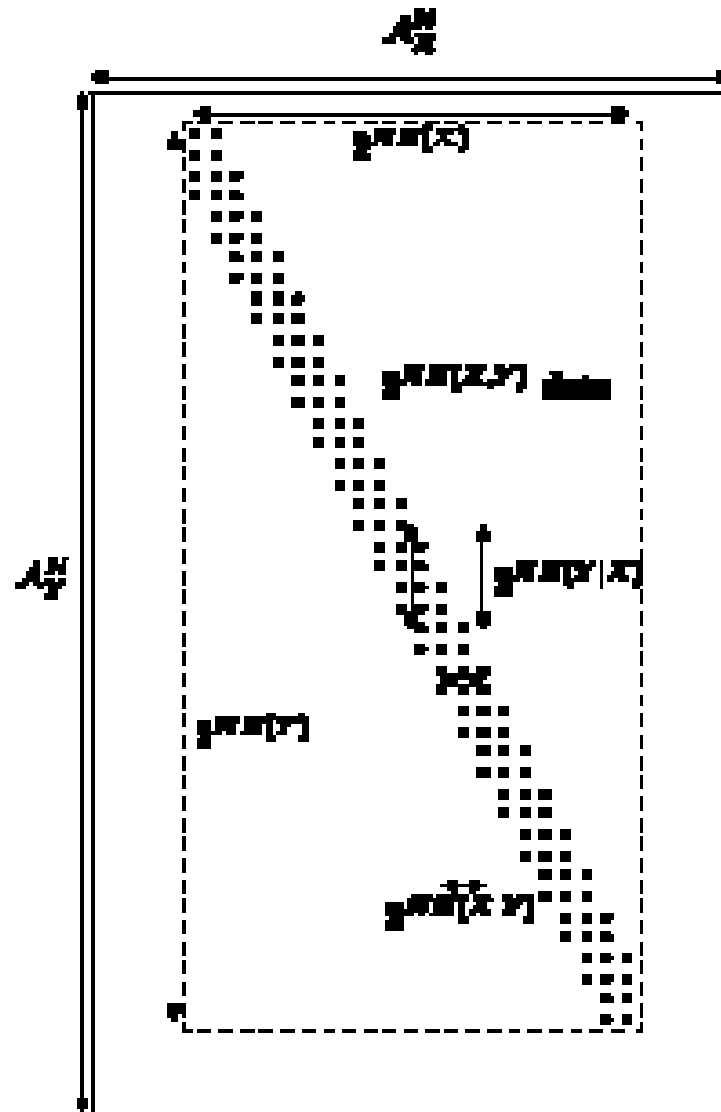
$$P(\mathbf{x}, \mathbf{y}) = P(x_1, y_1) \dots p(x_N, y_N)$$

- A sequence \mathbf{x}, \mathbf{y} is jointly typical if
 - \mathbf{x} is typical for $P(x)$
 - \mathbf{y} is typical of $P(y)$
 - (\mathbf{x}, \mathbf{y}) is typical for $P(\mathbf{x}, \mathbf{y})$

$$\left| \frac{1}{N} \log \frac{1}{P(\mathbf{x}, \mathbf{y})} - H(X, Y) \right| < \beta$$

- Size of the typical set roughly $2^{NH(X,Y)}$

JOINTLY TYPICAL SETS - DIAGRAM



DRAWING X AND Y INDEPENDENTLY

- Marginals should be the same

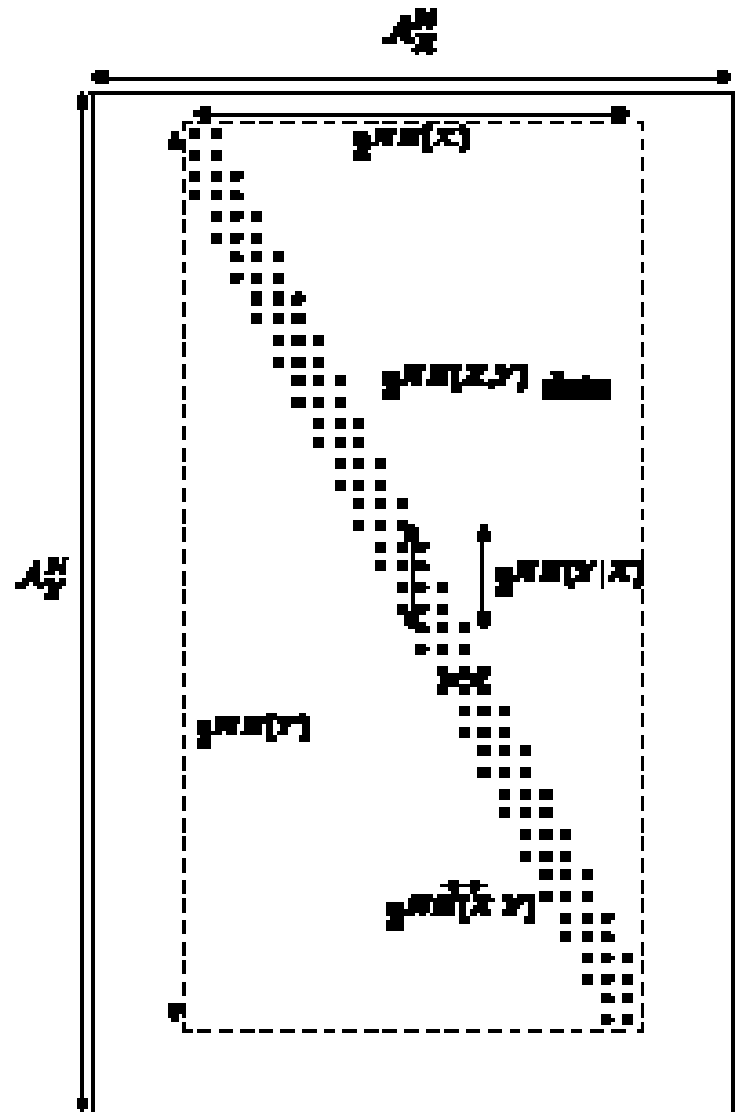
$$P(x') := \sum_y P(x', y)$$

$$P(y') := \sum_x P(x, y')$$

- Draw x' and y' from

$$P(x', y') := P(x')P(y')$$

- Are x' and y' likely to be in the original jointly typical set?



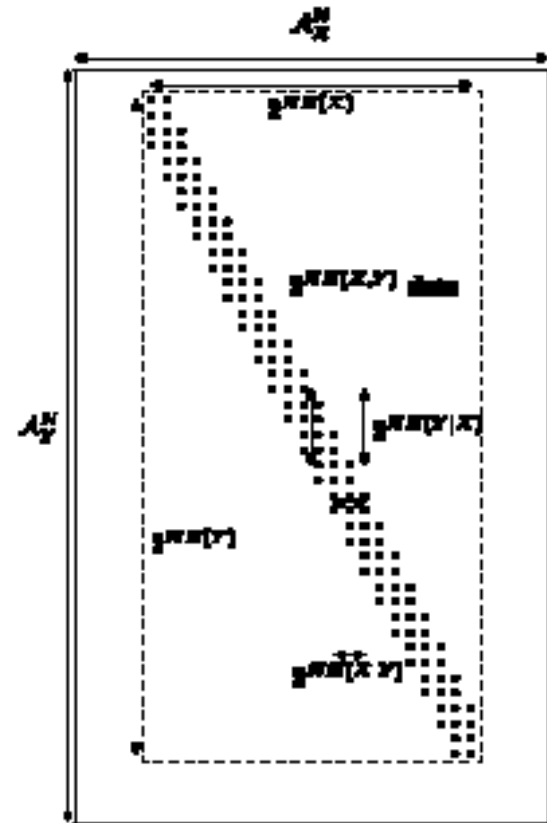
EMBEDDING IN THE JOINTLY TYPICAL SET

- Original jointly typical set

$$J_{N\beta} := \left\{ (x, y) \mid \begin{aligned} & \left| \frac{1}{N} \log \frac{1}{P(x)} - H(X) \right| \\ & \left| \frac{1}{N} \log \frac{1}{P(y)} - H(Y) \right| \\ & \left| \frac{1}{N} \log \frac{1}{P(x,y)} - H(X, Y) \right| \end{aligned} \right\}$$

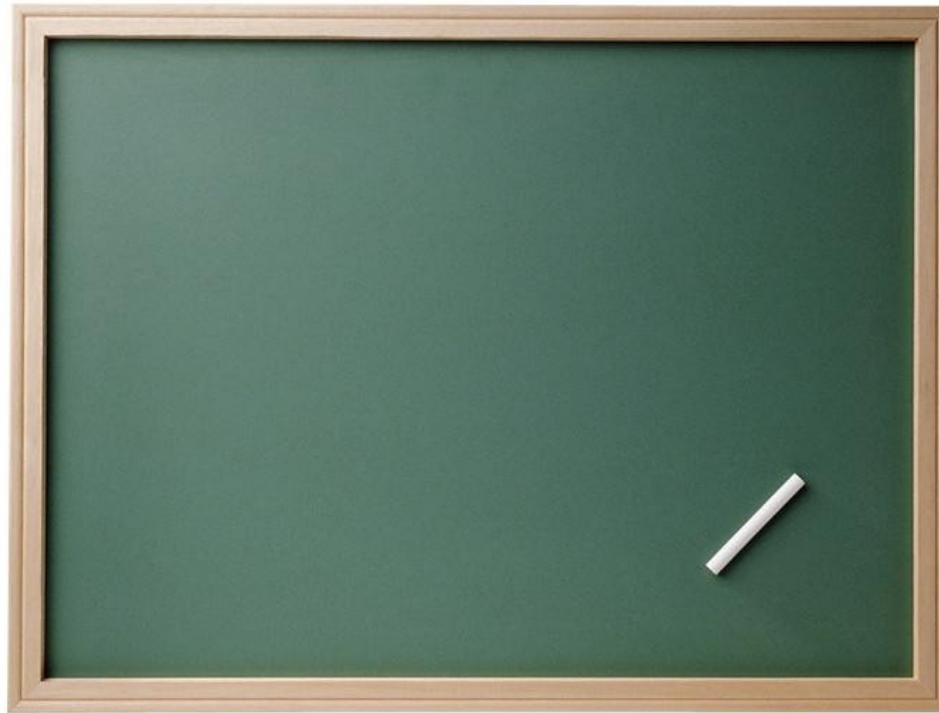
- From our picture we expect

$$\begin{aligned} P((x', y') \in J_{N\beta}) &\approx \frac{2^{NH(X,Y)}}{2^{N(H(X)+H(Y))}} \\ &= 2^{-NI(X;Y)} \end{aligned}$$



COMPUTATION

○ Let's verify:
$$P((x', y') \in J_{N\beta}) \approx \frac{2^{NH(X,Y)}}{2^{N(H(X)+H(Y))}} = 2^{-NI(X;Y)}$$



SUMMARY: JOINT TYPICALITY THEOREM

- Consider strings x and y with
$$P(x, y) = \prod_{j=1}^N P(x_j, y_j)$$

- 1. The probability that x and y are in the typical set goes to 1 for large N
- 2. The number of jointly typical sequences obeys

$$|J_{N\beta}| \leq 2^{N(H(X,Y)+\beta)}$$

- Proofs same as in source coding theorem
- 3. For x' and y' drawn independently according to the marginal distributions

$$P((x', y') \in J_{N\beta}) \leq 2^{-N(I(X,Y)-3\beta)}$$

- Just proved this!

ONE MORE TRICK....



A USEFUL TRICK: AVERAGING IS SOMETIMES EASIER

- Imagine we have 100 durians, and we want to determine whether at least one of them weighs 1.5 kg
- Question: in the best case, how many times do we have to use a (large) weighing scale?



A USEFUL TRICK: AVERAGING IS SOMETIMES EASIER

- If we are lucky, just once!
- Weigh them all together, and compute average weight $\frac{1}{100}W$

$$\frac{1}{100}W \geq 1.5kg$$

There is at least one durian of weight 1.5kg


$$\frac{1}{100}W < 1.5kg$$

Hmm, who knows? There might be a big durian if there are many small ones too.

HOW IS THIS USEFUL?

- We want to show that there exists a code with rate achieving the capacity
- Average over all possible codes!

Average block error
for code

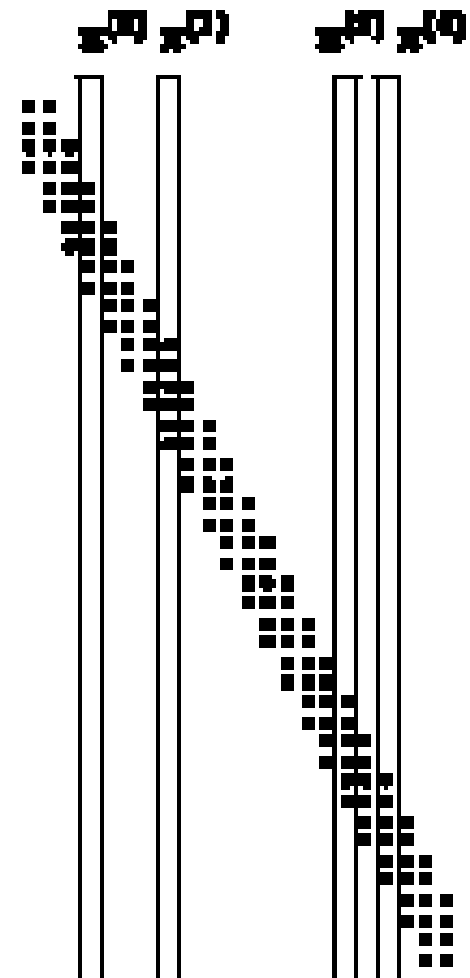
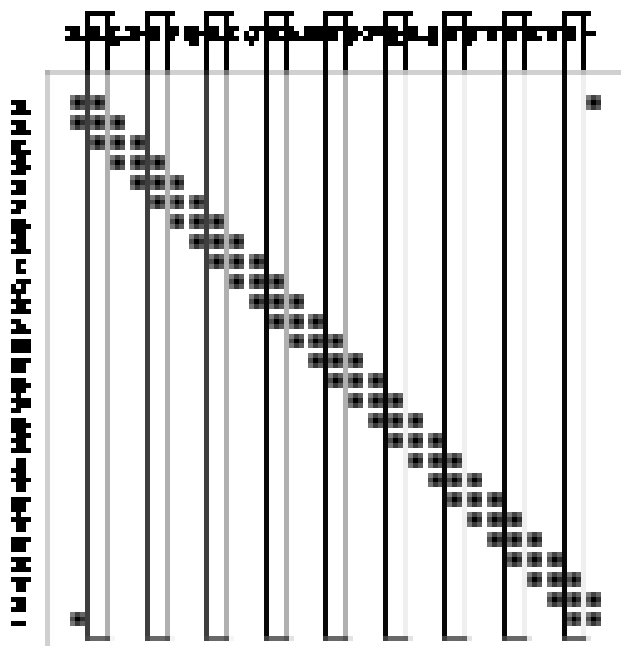
$$\mathbb{E}(P_{\text{err,Block}}) = \sum_{\text{code } C} P(C) P(\hat{s} \neq s | C)$$


- If we can show that this average is small for codes with a rate at capacity then we know there exists at least one such code!
- ... later worry about maximum block error

WHAT DOES IT MEAN TO CHOOSE CODES RANDOMLY?

- Generate $2^{NR'}$ codewords for a $(N, NR') = (N, K)$ code

Compare noisy-typewriter



FOR A RANDOM CODE

- Sender and receiver both know the code!
- A source symbol $s \in \{1, \dots, 2^{NR'}\}$ is encoded as $\mathbf{x}^{(s)}$
- The received signal is \mathbf{y} with probability

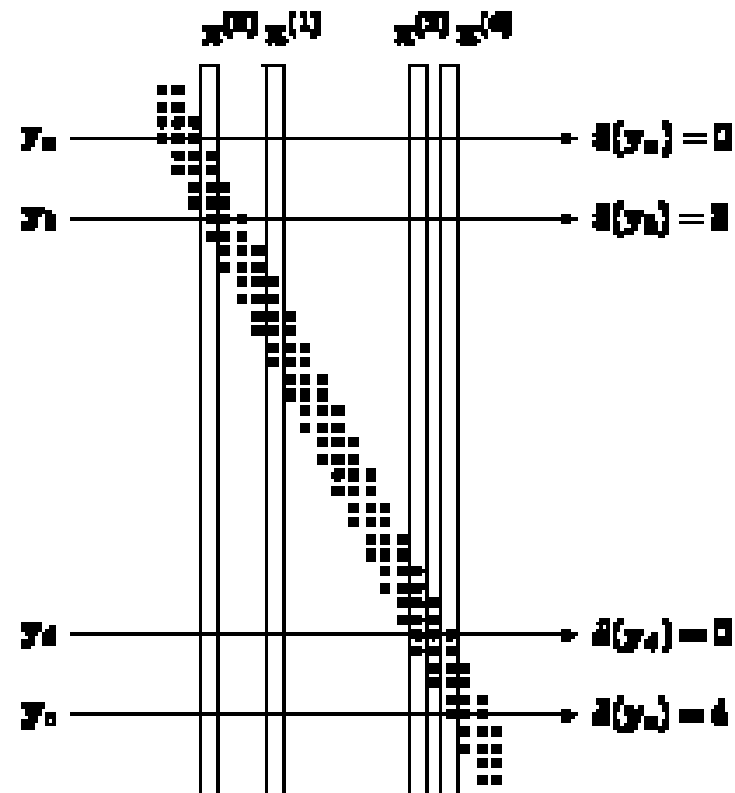
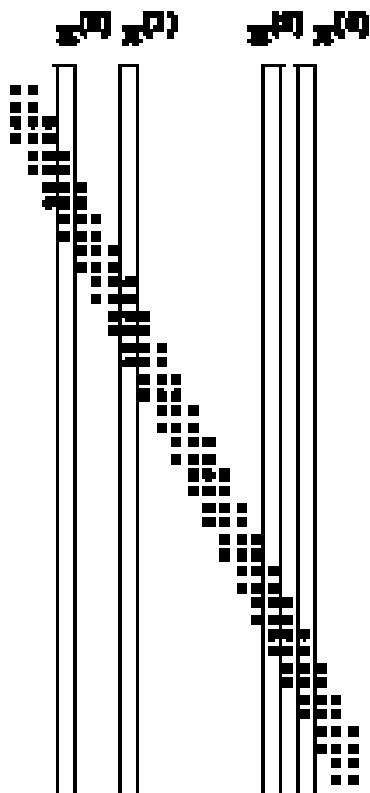
$$P(\mathbf{y}|\mathbf{x}^{(s)}) = \prod_{j=1}^N P(y_j|x_j^{(s)})$$

ONTO TO THE PROOF!



OUTLINE

- Choose a random code, and random inputs s
- Typical set decoding: Given y output \hat{s} such that
 - $(x^{(\hat{s})}, y)$ are jointly typical
 - There is no other s' such that $(x^{(s')}, y)$ are jointly typical
 - Otherwise output failure $\hat{s} = 0$



HOW WELL CAN THIS WORK?

- Two ways our typical set decoder can fail
 - There's no \hat{s} such that $(x^{(\hat{s})}, y)$ is jointly typical
 - There is another s' such that $(x^{(s')}, y)$ is jointly typical
- Since all codewords play the same role, it is sufficient to show what happens for some fixed input $s = 1$
- Want to show that the probability both failure events happen is small

PART 1

- Want to show that the probability that there is no jointly typical element can be made arbitrarily small
- Joint typicality theorem, part 1:

$$P((\mathbf{x}^{(1)}, \mathbf{y}) \in J_{N\beta}) \geq 1 - \delta$$

- But then there exists at least one jointly typical element (with high probability)

PART 2

- Want to show that there is no other s' such that $(x^{(s')}, y)$ is jointly typical

- By the joint typicality theorem for any fixed $s' \neq 1$

$$P((x^{(s')}, y) \in J_{N\beta}) \leq 2^{-N(I(X;Y)-3\beta)}$$

- Since there are $2^{NR'} - 1$ values for $s' \neq 1$

$$\begin{aligned} \mathbb{E}(p_{\text{err,Block}}) &\leq \delta + \sum_{s'=2}^{2^{NR'}} 2^{-N(I(X;Y)-3\beta)} \\ &\leq \delta + 2^{-N(I(X;Y)-R'-3\beta)} \end{aligned}$$

WRAPPING IT UP

- We know the average block error obeys

$$\mathbb{E}(p_{\text{err,Block}}) \leq \delta + 2^{-N(I(X;Y) - R' - 3\beta)}$$

- Hence for $R' < I(X;Y) - 3\beta$ we have for large enough N

$$\mathbb{E}(p_{\text{err,Block}}) \leq 2\delta$$

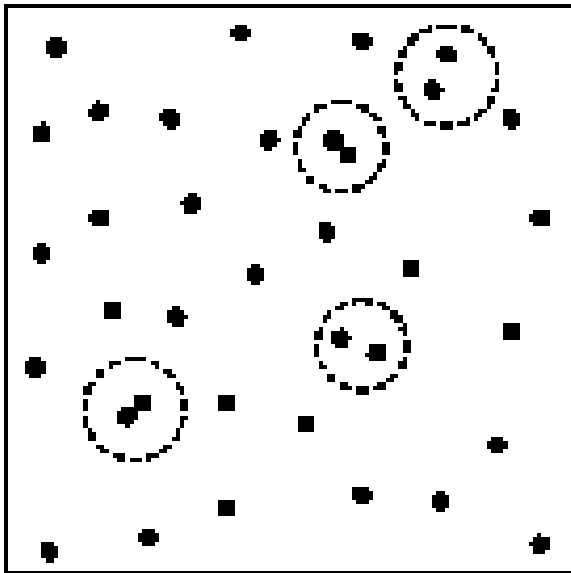
- To find the largest rate, we maximize $I(X;Y)$

$$R' < \max_{P_X} I(X;Y) - 3\delta = C - 3\delta$$

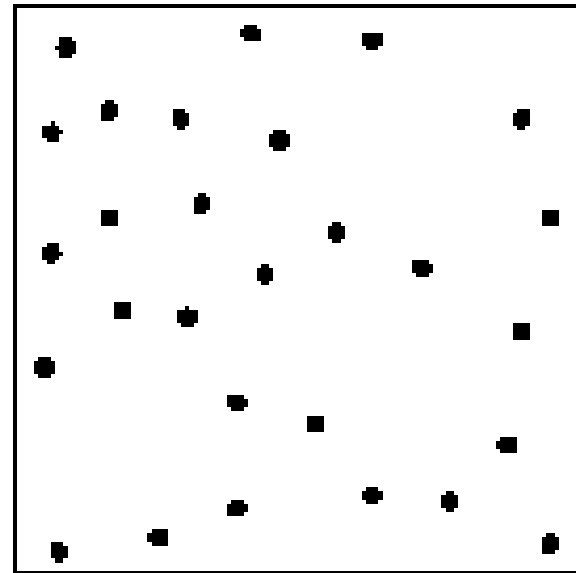
- Since the average over codes is small, there exists a code with small average block error

AVERAGE ERROR VS MAXIMUM ERROR

- To find a bound on maximum error we discard some codewords which are easily confused

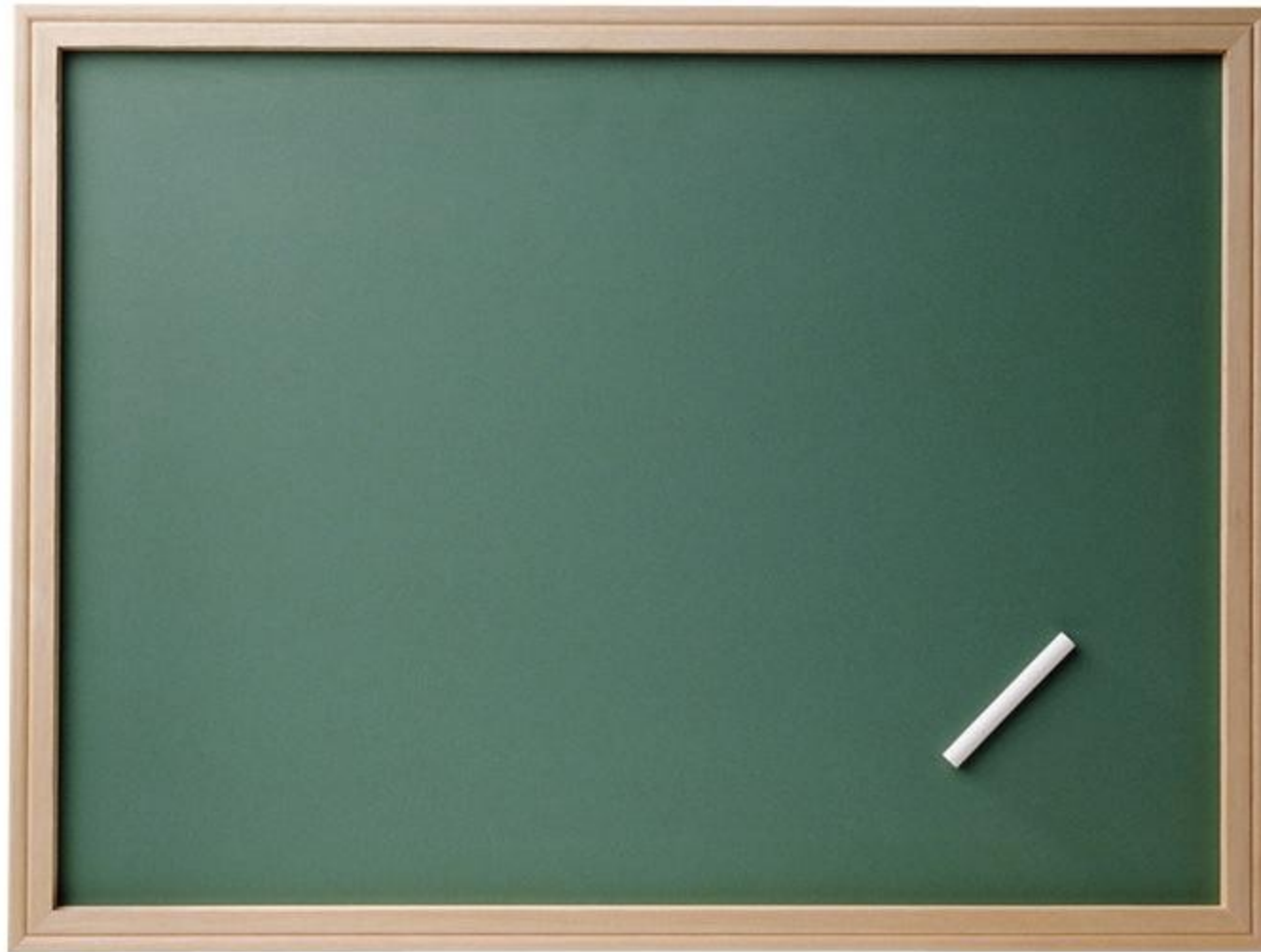


(a) A random code ...



(b) pruned

BOUNDING THE MAXIMUM ERROR



SHANNON'S NOISY CHANNEL CODING THEOREM (PART 1: ACHIEVABILITY)

- Associated with every discrete memoryless channel, there is a non-negative number C (the capacity) such that
 - For any error $\epsilon > 0$ and $R \leq C$ for large enough N , there exists a block code of length N and rate R , and a decoding algorithm, such that the maximum probability of block error is $< \epsilon$



CONVERSE TO SHANNON'S CODING THEOREM

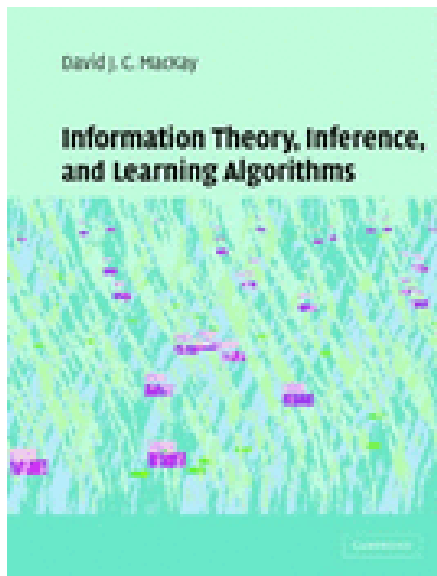
- For rates above the capacity, the error cannot be made small
- Proof during tutorial session using Fano's inequality
- Capacity forms sharp threshold for information transmission!

NEXT TIME

- Back to error correcting codes: what did we learn from Shannon's theory?

READING FOR THIS LECTURE

- Chapters 9.6 to 9.8, 10.1 to 10.3 in the book



Information Theory, Inference and
Learning Algorithms
by David J. C. MacKay
Cambridge University Press, 2003

- Homework due by Monday 2pm two weeks from now (Hari Raya Haji next Monday, no lecture, but tutorial sessions held)