

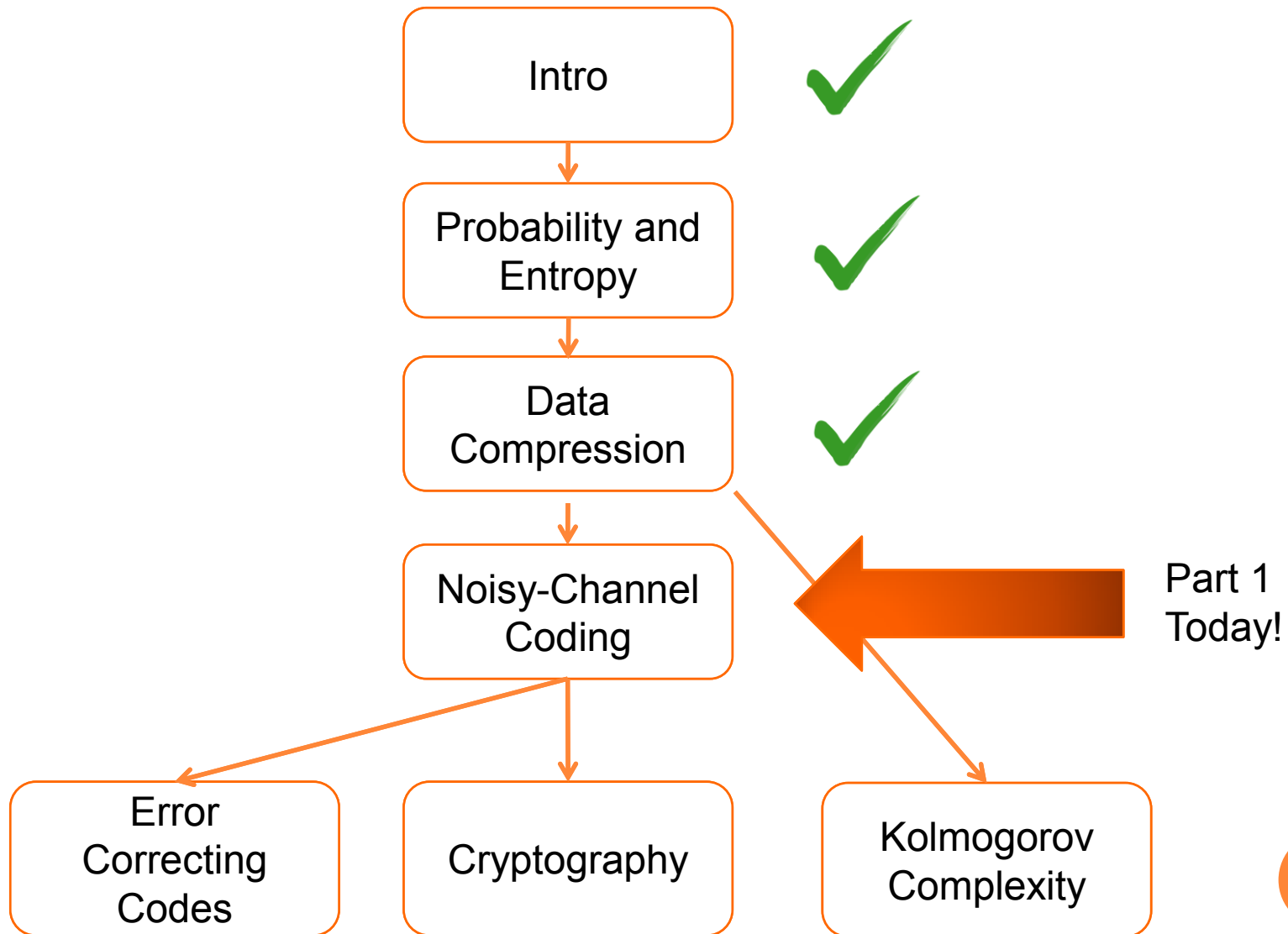
# CS3236 INTRODUCTION TO INFORMATION THEORY

## Lecture 6: Noisy Channel Coding

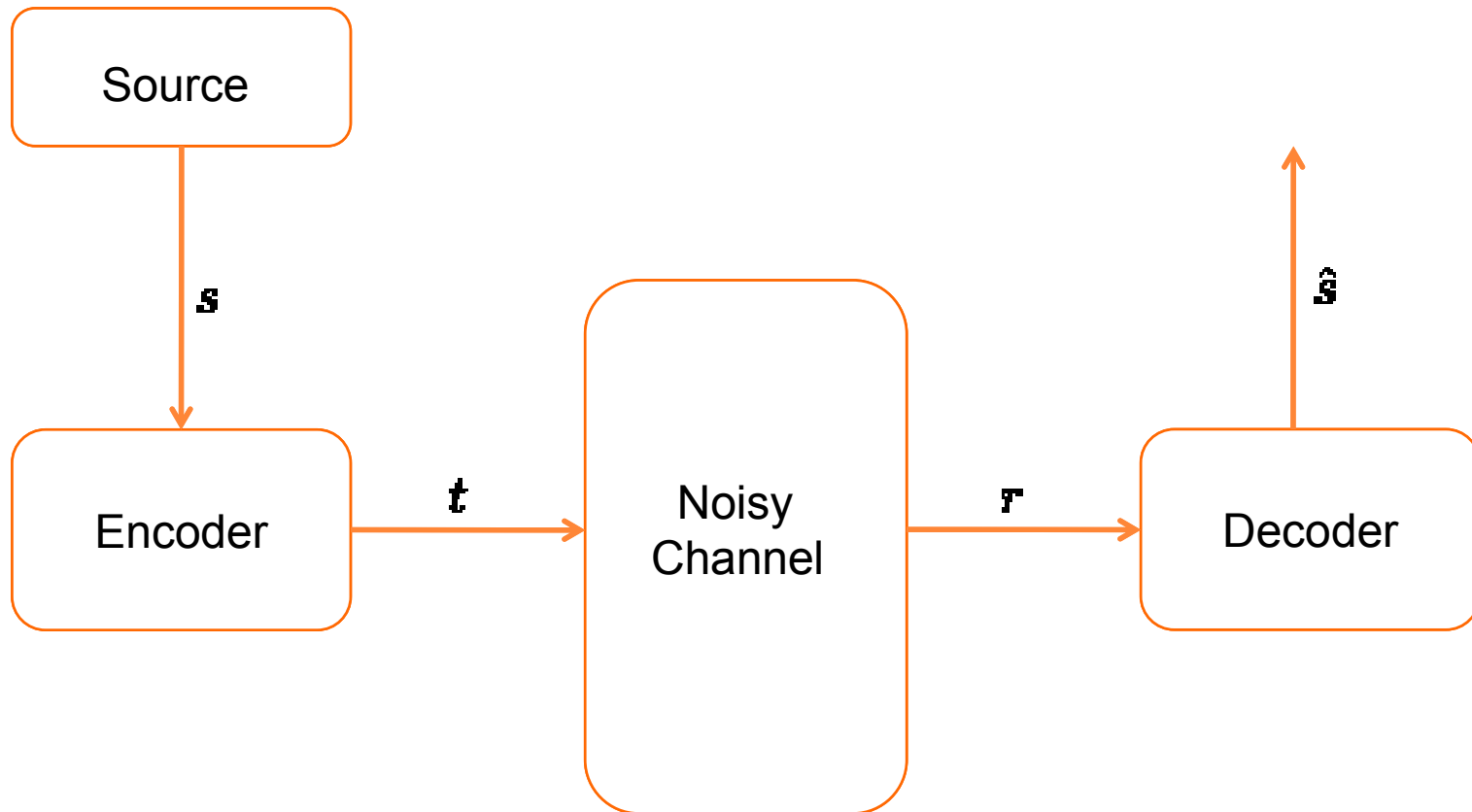
Course given by Pierre Senellart

Material by Stephanie Wehner, with additions by P. Senellart

# WHERE DO WE GO FROM HERE?



# SENDING INFORMATION WITHOUT ERROR



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

# WHAT WE'LL DO THIS TIME

- Noisy-channel coding
  - Conditional entropies
  - Mutual information
  - Capacity of a channel
  - Examples of noisy channels
  - Block error probability
- Statement of Shannon's channel coding theorem

# REMINDER: ENTROPIES

- Shannon entropy

$$H(X) = - \sum_x \text{Pr}_X(x) \log \text{Pr}_X(x)$$

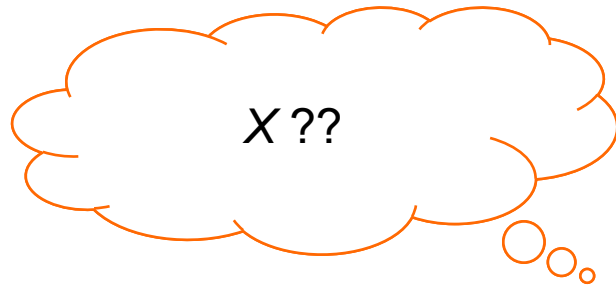
- Joint Shannon entropy

$$H(X, Y) = - \sum_{x, y} \text{Pr}_{XY}(x, y) \log \text{Pr}_{XY}(x, y)$$

- For a product distribution  $\text{Pr}_{XY}(x, y) = \text{Pr}_X(x) \text{Pr}_Y(y)$

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

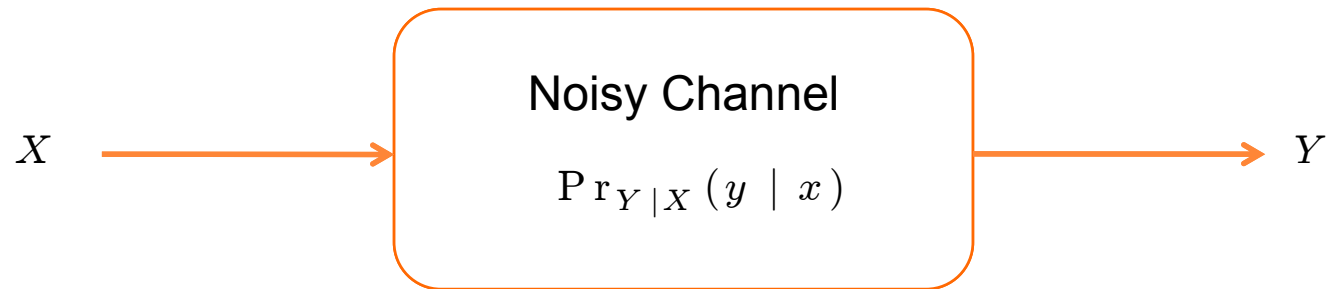
# CONDITIONAL ENTROPY $H(X|Y)$



Y

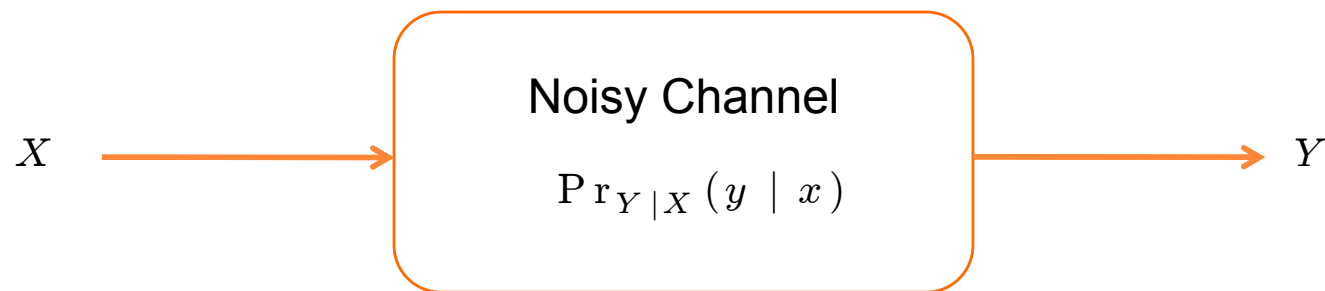
# MOTIVATION

- How can we quantify the entropy of  $X$  **given**  $Y$ ?



# MOTIVATION

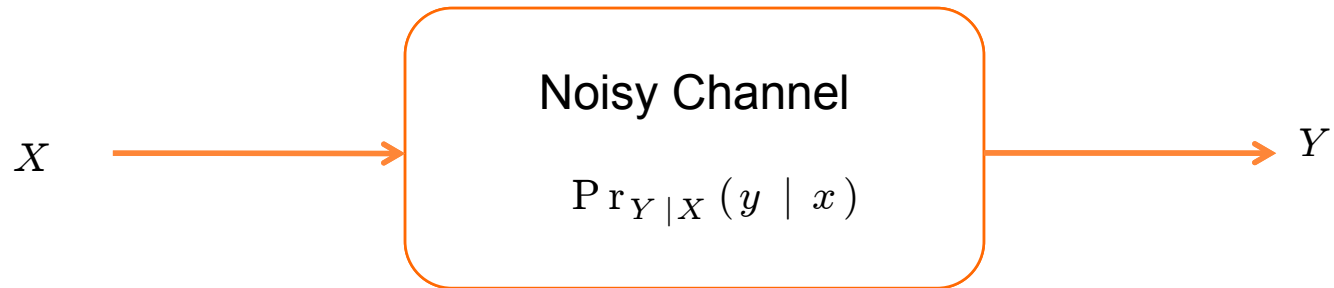
- How can we quantify the entropy of  $X$  **given**  $Y$ ?



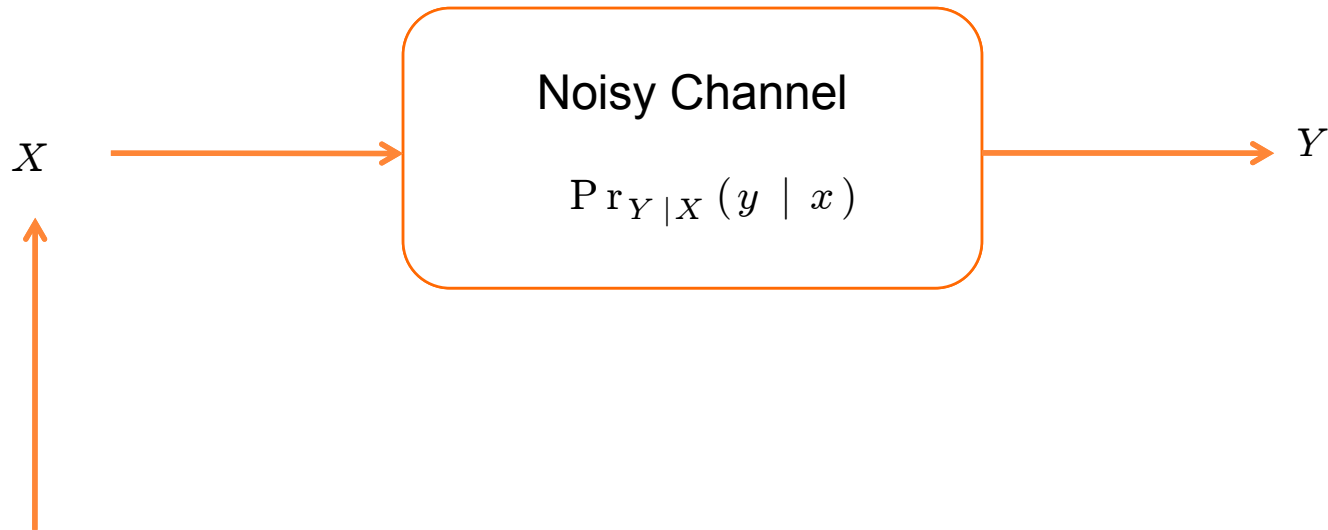
Example: Binary symmetric channel

$$\Pr_{Y|X}(y|x) = \begin{cases} 1 - f & \text{if } y = x \in \{0, 1\} \\ f & \text{otherwise} \end{cases}$$

# CHANGE OF ENTROPY



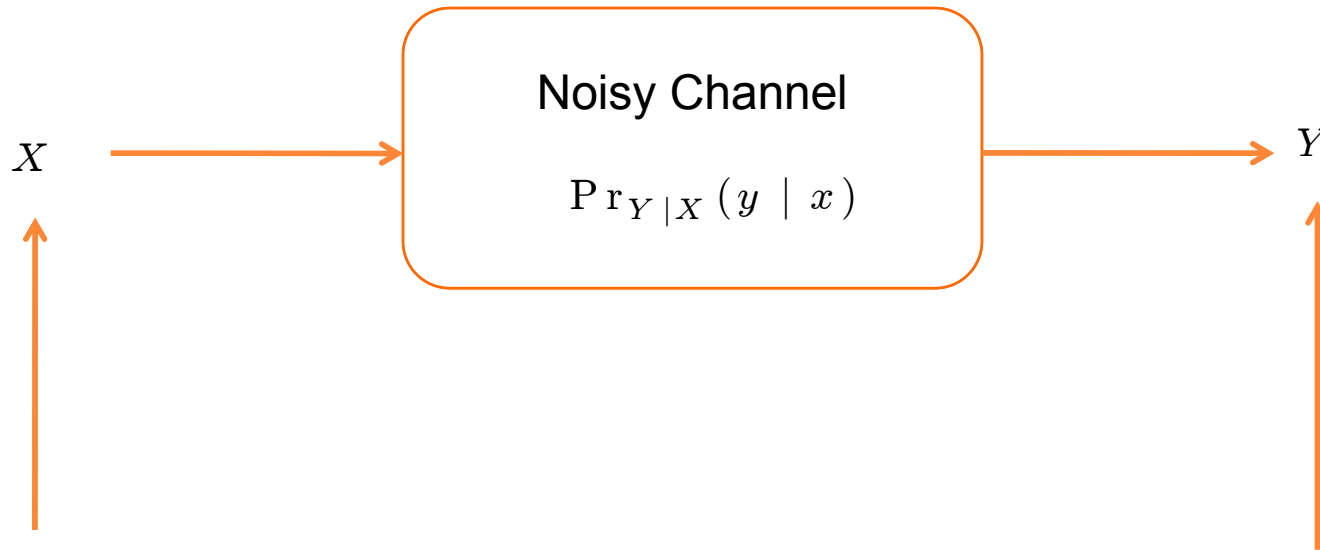
# CHANGE OF ENTROPY



Initial uncertainty  
(without any observation!)  
given by the entropy

$$H(X) = - \sum_x \Pr_X(x) \log \Pr_X(x)$$

# CHANGE OF ENTROPY



Initial uncertainty  
(without any observation!)  
given by the entropy

$$H(X) = - \sum_x \Pr_X(x) \log \Pr_X(x)$$

Now we observe  $Y$   
which could decrease  
our uncertainty

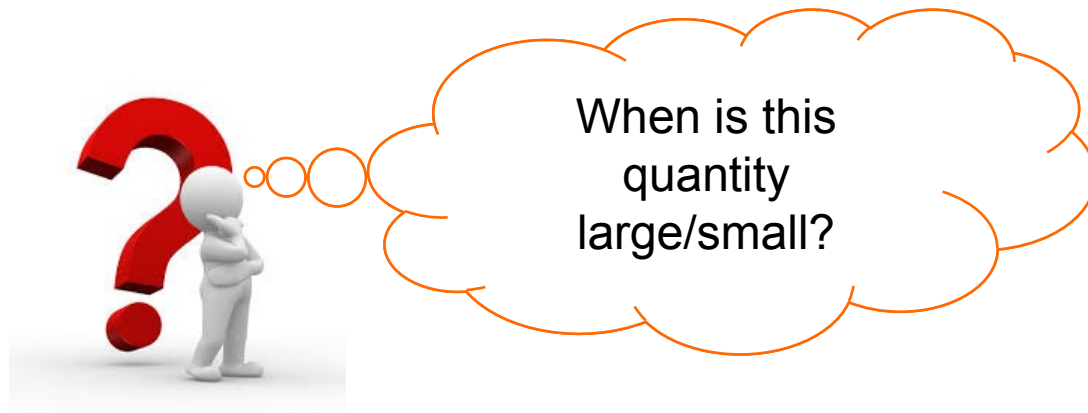
$$H(X | Y)$$

$$H(X | Y) = \sum_y \Pr_Y(y) H(X | Y = y)$$

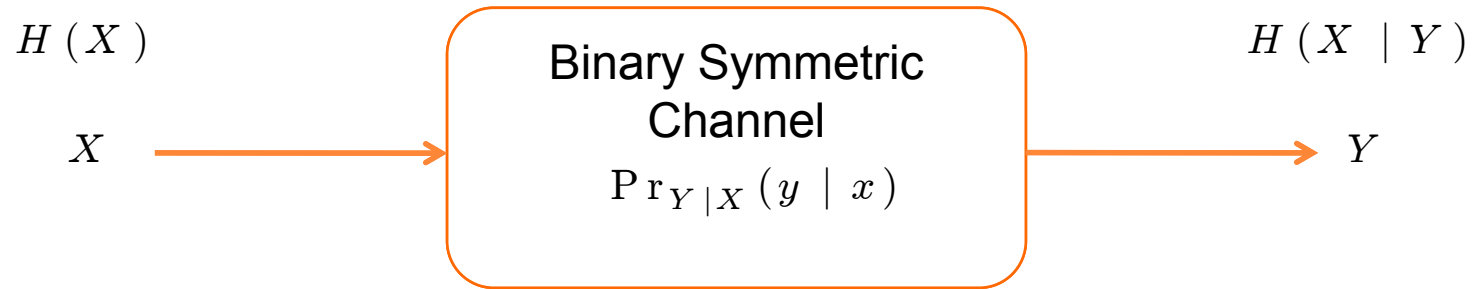
# CONDITIONAL ENTROPY

$$H(X | Y) = \sum_y \Pr_Y(y) H(X | Y = y)$$

$$H(X | Y = y) = - \sum_x \Pr_{X|Y}(x | y) \log \Pr_{X|Y}(x | y)$$



# WHAT DO YOU EXPECT?



$$\Pr_{Y|X}(y|x) = \begin{cases} 1 - f & \text{if } y = x \in \{0, 1\} \\ f & \text{otherwise} \end{cases}$$

# CONDITIONAL ENTROPY

- Entropy of  $X$  given  $Y$

$$H(X | Y) = - \sum_{x,y} \text{Pr}_Y(y) \text{Pr}_{X|Y}(x | y) \log \text{Pr}_{X|Y}(x | y) \geq 0$$



“Side information”

- Other examples

- Compression:  $Y = C(X)$  is the compressed string
- Cryptography:  $Y$  is some information that the adversary has gathered

# CHAIN RULE FOR INFORMATION CONTENT

- We know that

$$\Pr_{XY}(x, y) = \Pr_{X|Y}(x | y) \Pr_Y(y) = \Pr_{Y|X}(y | x) \Pr_X(x)$$

- *Chain rule* for information content

$$\log \frac{1}{\Pr_{XY}(x, y)} = \log \frac{1}{\Pr_X(x)} + \log \frac{1}{\Pr_{Y|X}(y | x)}$$

# A FEW WAYS TO WRITE THE CONDITIONAL ENTROPY: CHAIN RULE FOR ENTROPY

- Chain rule for information content

$$\log \frac{1}{\text{Pr}_{XY}(x, y)} = \log \frac{1}{\text{Pr}_X(x)} + \log \frac{1}{\text{Pr}_{Y|X}(y | x)}$$

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

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

# A FEW WAYS TO WRITE THE CONDITIONAL ENTROPY: CHAIN RULE FOR ENTROPY

- Chain rule for information content

$$\log \frac{1}{\text{Pr}_{XY}(x, y)} = \log \frac{1}{\text{Pr}_X(x)} + \log \frac{1}{\text{Pr}_{Y|X}(y | x)}$$

Translated into entropies, this gives us the so-called *chain rule*

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

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

# A FEW WAYS TO WRITE THE CONDITIONAL ENTROPY: CHAIN RULE FOR ENTROPY

- Chain rule for information content

$$\log \frac{1}{\text{Pr}_{XY}(x, y)} = \log \frac{1}{\text{Pr}_X(x)} + \log \frac{1}{\text{Pr}_{Y|X}(y | x)}$$

Translated into entropies, this gives us the so-called *chain rule*

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

- Solving for  $H(X | Y) = H(X, Y) - H(Y)$

# A FEW WAYS TO WRITE THE CONDITIONAL ENTROPY: CHAIN RULE FOR ENTROPY

- Chain rule for information content

$$\log \frac{1}{\text{Pr}_{XY}(x, y)} = \log \frac{1}{\text{Pr}_X(x)} + \log \frac{1}{\text{Pr}_{Y|X}(y | x)}$$

Translated into entropies, this gives us the so-called *chain rule*

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

- Solving for  $H(X | Y) = H(X, Y) - H(Y)$

# A FEW WAYS TO WRITE THE CONDITIONAL ENTROPY: CHAIN RULE FOR ENTROPY

- Chain rule for information content

$$\log \frac{1}{\text{Pr}_{XY}(x, y)} = \log \frac{1}{\text{Pr}_X(x)} + \log \frac{1}{\text{Pr}_{Y|X}(y | x)}$$

Translated into entropies, this gives us the so-called *chain rule*

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

- S

olving for  $H(X | Y) = H(X, Y) - H(Y)$

Intuitively, total entropy of  $X$  and  $Y$ , minus what we know, namely  $Y$ , gives  $X$

# HOW ARE THEY RELATED?

- Would think that the entropy, i.e., uncertainty becomes less if we have  $Y$  then if we had nothing....
- Conditioning reduces entropy

$$H(X | Y) \leq H(X)$$

# HOW SMALL OR LARGE IS THIS QUANTITY?

- Conditional entropy  $H(X | Y) = H(X, Y) - H(Y)$

# HOW SMALL OR LARGE IS THIS QUANTITY?

- Conditional entropy  $H(X | Y) = H(X, Y) - H(Y)$
- Upper bound  $H(X | Y) \leq H(X) \leq \log |X|$

When is this attained?

# HOW SMALL OR LARGE IS THIS QUANTITY?

○ Conditional entropy  $H(X | Y) = H(X, Y) - H(Y)$

○ Upper bound  $H(X | Y) \leq H(X) \leq \log |X|$

When is this attained?

Y says nothing about X: product distribution

$$\Pr_{XY}(x, y) = \Pr_X(x) \Pr_Y(y)$$

# HOW SMALL OR LARGE IS THIS QUANTITY?

○ Conditional entropy  $H(X | Y) = H(X, Y) - H(Y)$

○ Upper bound  $H(X | Y) \leq H(X) \leq \log |X|$

When is this attained?

Y says nothing about X: product distribution

$$\Pr_{XY}(x, y) = \Pr_X(x) \Pr_Y(y)$$

Lower bound  $0 \leq H(X | Y)$

When is this attained?

# HOW SMALL OR LARGE IS THIS QUANTITY?

○ Conditional entropy  $H(X | Y) = H(X, Y) - H(Y)$

○ Upper bound  $H(X | Y) \leq H(X) \leq \log |X|$

When is this attained?

Y says nothing about X: product distribution

$$\Pr_{XY}(x, y) = \Pr_X(x) \Pr_Y(y)$$

Lower bound  $0 \leq H(X | Y)$

When is this attained?

Y says everything about X

# A WORD OF CAUTION

- Conditioning only reduces entropy *on average!*

$$H(X | Y) \leq H(X)$$

# A WORD OF CAUTION

- Conditioning only reduces entropy *on average!*

$$H(X | Y) \leq H(X)$$

But it is not always the case that for all  $y$

$$H(X | Y = y) = - \sum_x \Pr_{X|Y}(x | y) \log \Pr_{X|Y}(x | y) \leq H(X)$$

# EXAMPLE

- Conditioning reduces entropy but only on average

$P(x   y)$		$x$				$H(X   y)/\text{bits}$
		1	2	3	4	
$y$	1	$1/2$	$1/4$	$1/8$	$1/8$	$7/4$
	2	$1/4$	$1/2$	$1/8$	$1/8$	$7/4$
	3	$1/4$	$1/4$	$1/4$	$1/4$	2
	4	1	0	0	0	0

$$H(X) = 7/4 = 14/8 \quad H(X | Y) = 11/8$$

# EXAMPLE

- Conditioning reduces entropy but only on average

$P(x   y)$		$x$				$H(X   y)/\text{bits}$
		1	2	3	4	
$y$	1	$1/2$	$1/4$	$1/8$	$1/8$	$7/4$
	2	$1/4$	$1/2$	$1/8$	$1/8$	$7/4$
	3	$1/4$	$1/4$	$1/4$	$1/4$	$2$
	4	$1$	$0$	$0$	$0$	$0$

$$H(X) = 7/4 = 14/8 \quad H(X | Y) = 11/8$$

$$H(X | Y) < H(X)$$

# EXAMPLE

- Conditioning reduces entropy but only on average

$P(x   y)$		$x$				$H(X   y)/\text{bits}$
		1	2	3	4	
$y$	1	$1/2$	$1/4$	$1/8$	$1/8$	$7/4$
	2	$1/4$	$1/2$	$1/8$	$1/8$	$7/4$
	3	$1/4$	$1/4$	$1/4$	$1/4$	2
	4	1	0	0	0	0

←  $< H(X)$

$$H(X) = 7/4 = 14/8 \quad H(X | Y) = 11/8$$

$$H(X | Y) < H(X)$$

# EXAMPLE

- Conditioning reduces entropy but only on average

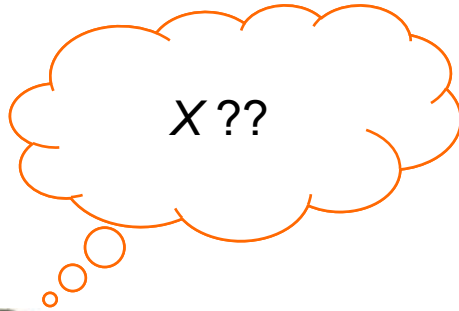
$P(x   y)$		$x$				$H(X   y)/\text{bits}$
		1	2	3	4	
$y$	1	1/2	1/4	1/8	1/8	7/4
	2	1/4	1/2	1/8	1/8	7/4
	3	1/4	1/4	1/4	1/4	2 ← $> H(X)$
	4	1	0	0	0	0 ← $< H(X)$

$$H(X) = 7/4 = 14/8 \quad H(X | Y) = 11/8$$

$$H(X | Y) < H(X)$$

# MUTUAL INFORMATION $I(X; Y) = H(X) - H(X|Y)$

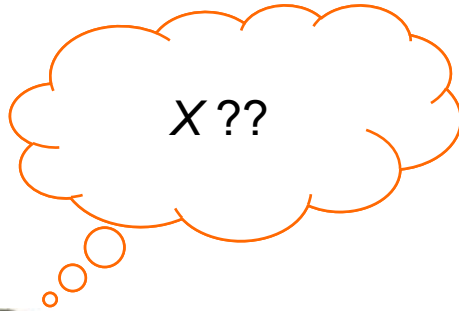
$H(X)$



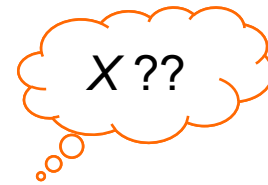
$Y$

# MUTUAL INFORMATION $I(X; Y) = H(X) - H(X|Y)$

$H(X)$



$H(X|Y)$



Y



# MUTUAL INFORMATION

- Measures the reduction of uncertainty about  $X$  after we learn  $Y$ , i.e., the average amount of information that  $Y$  gives us about  $X$

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

# MUTUAL INFORMATION

- Measures the reduction of uncertainty about  $X$  after we learn  $Y$ , i.e., the average amount of information that  $Y$  gives us about  $X$

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

Is symmetric

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

$$I(X; Y) = I(Y; X)$$

# MUTUAL INFORMATION

- Measures the reduction of uncertainty about  $X$  after we learn  $Y$ , i.e., the average amount of information that  $Y$  gives us about  $X$

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

Is symmetric

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

$$I(X; Y) = I(Y; X)$$

Conditioning reduces entropy implies that

$$I(X; Y) \geq 0$$

# CONDITIONAL MUTUAL INFORMATION

- What if we are already given some extra information,  $Z$ ?

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



Simply add the side information  $Z$  to all conditionings

# HOW LARGE OR SMALL CAN THIS GET?

- Mutual information

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

# HOW LARGE OR SMALL CAN THIS GET?

- Mutual information

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

Lower bound  $I(X; Y) \geq 0$

When is this attained?

Product distribution  $\Pr_{X Y}(x, y) = \Pr_X(x) \Pr_Y(y)$

# HOW LARGE OR SMALL CAN THIS GET?

- Mutual information

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

Lower bound  $I(X; Y) \geq 0$

When is this attained?

Product distribution  $\Pr_{X Y}(x, y) = \Pr_X(x) \Pr_Y(y)$

Upper bound  $I(X; Y) \leq H(X) \leq \log |X|$

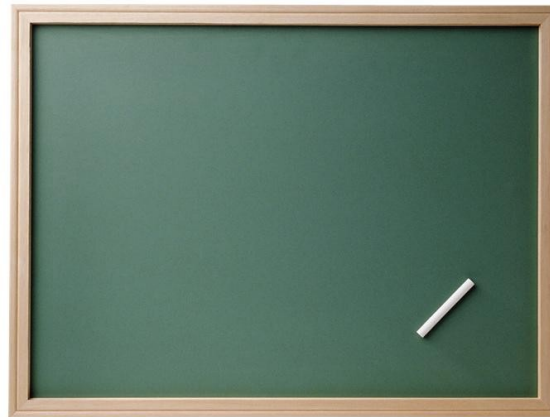
- When is this attained?
- Maximally far from product,  $Y$  tells us everything about  $X$

# CHAIN RULE FOR THE MUTUAL INFORMATION

$$\begin{aligned} I(X ; Y Z) &= I(X ; Y) + I(X ; Z | Y) \\ &= I(X ; Z) + I(X ; Y | Z) \end{aligned}$$

# CHAIN RULE FOR THE MUTUAL INFORMATION

$$\begin{aligned} I(X; YZ) &= I(X; Y) + I(X; Z | Y) \\ &= I(X; Z) + I(X; Y | Z) \end{aligned}$$



# INBETWEEN EXTREMES

- A measure how far we are from the product distribution
- Can write the mutual information in terms of the relative entropy

$$I(X; Y) = D_{\text{KL}}(\text{Pr}_{XY} \parallel \text{Pr}_X \text{Pr}_Y)$$

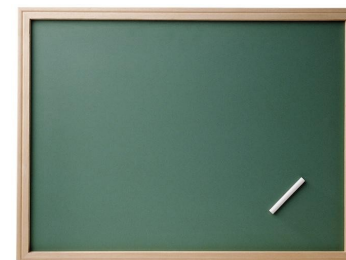


Intuitively, measures how far we are from the product distribution

# INBETWEEN EXTREMES

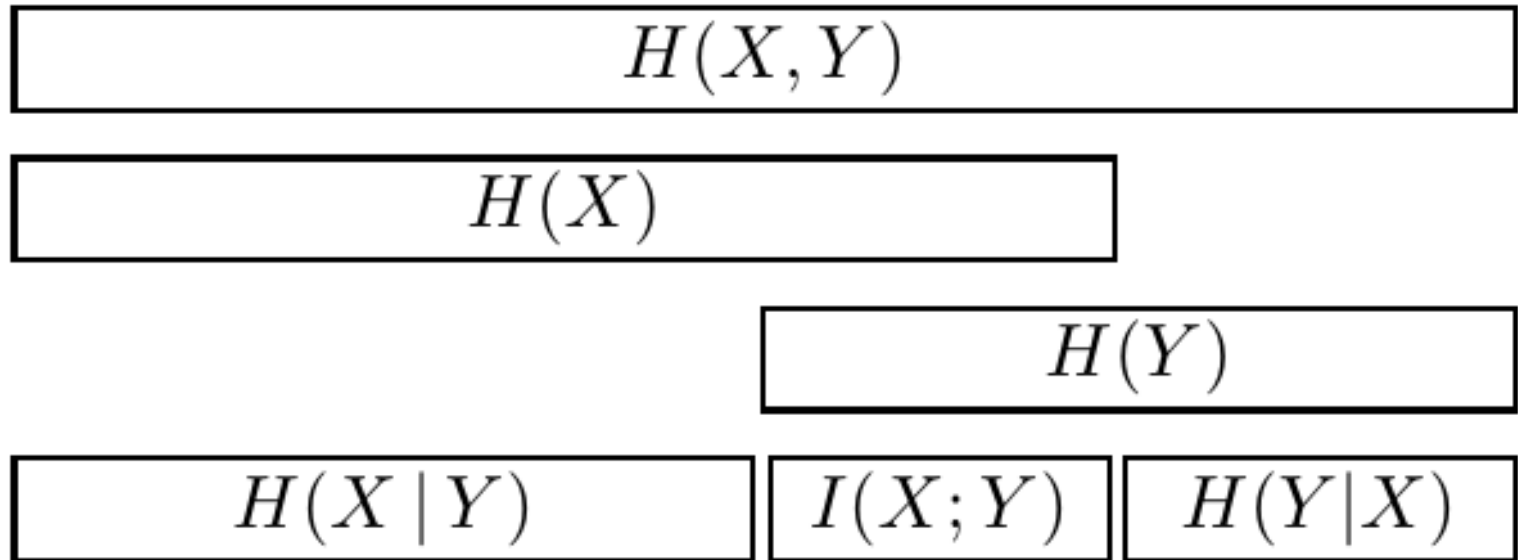
- A measure how far we are from the product distribution
- Can write the mutual information in terms of the relative entropy

$$I(X; Y) = D_{\text{KL}}(\text{Pr}_{XY} \parallel \text{Pr}_X \text{Pr}_Y)$$



Intuitively, measures how far we are from the product distribution

# RELATION AMONG ENTROPIES



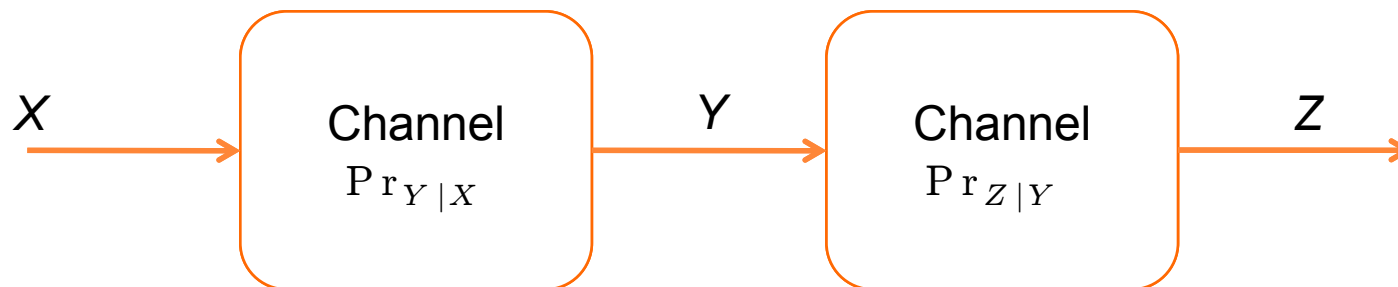
# MARKOV CHAIN

- Random variables form a Markov Chain

$$X \rightarrow Y \rightarrow Z$$

if  $Z$  only depends on  $Y$  (but is conditionally independent of  $X$ )

$$\Pr_{XYZ}(x, y, z) = \Pr_X(x) \Pr_{Y|X}(y | x) \Pr_{Z|Y}(z | y)$$



# PROPERTIES OF A MARKOV CHAIN

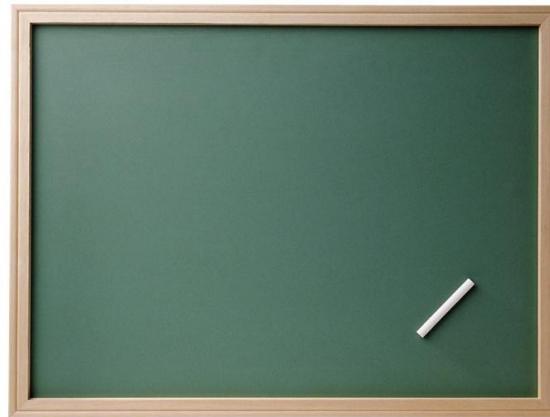
- X and Z are conditionally independent

$$\Pr_{X Z | Y} (x, z | y) = \Pr_{X | Y} (x | y) \Pr_{Z | Y} (z | y)$$

# PROPERTIES OF A MARKOV CHAIN

- X and Z are conditionally independent

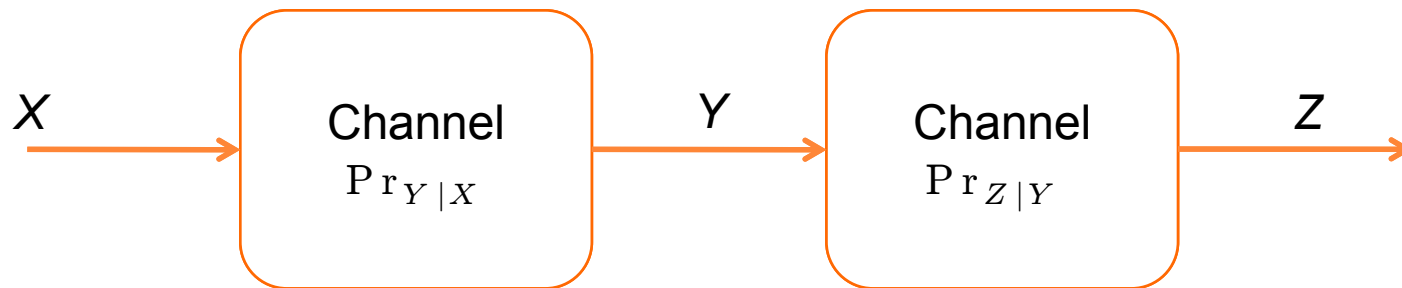
$$\Pr_{X Z | Y} (x, z | y) = \Pr_{X | Y} (x | y) \Pr_{Z | Y} (z | y)$$



# PROPERTIES OF A MARKOV CHAIN

Symmetric: If  $X \rightarrow Y \rightarrow Z$  then also  $Z \rightarrow Y \rightarrow X$

If  $Z = f(Y)$  then  $X \rightarrow Y \rightarrow Z$



# DATA PROCESSING

- Suppose we have a Markov chain  $X \rightarrow Y \rightarrow Z$
- What can we say about  $I(X; Y)$  vs.  $I(X; Z)$ ?

# DATA PROCESSING

- Suppose we have a Markov chain  $X \rightarrow Y \rightarrow Z$
- What can we say about  $I(X; Y)$  vs.  $I(X; Z)$ ?

Data processing

If  $X$ ,  $Y$ , and  $Z$  form a Markov Chain, then

$$I(X; Y) \geq I(X; Z)$$

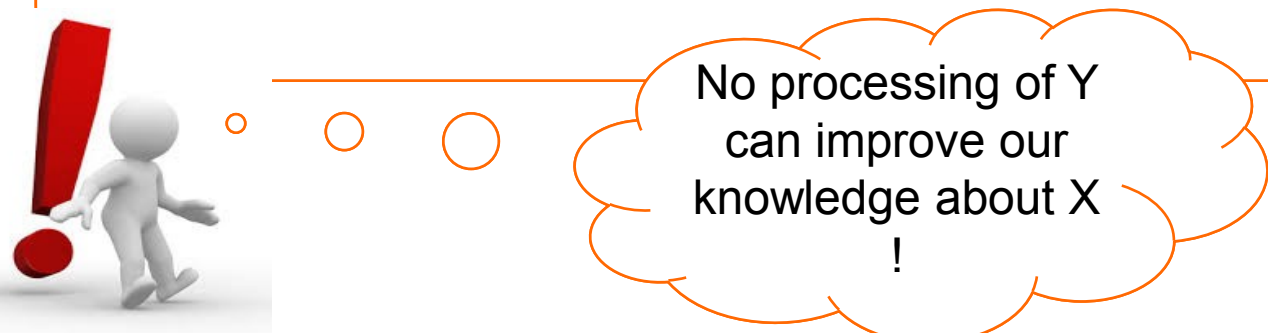
# DATA PROCESSING

- Suppose we have a Markov chain  $X \rightarrow Y \rightarrow Z$
- What can we say about  $I(X; Y)$  vs.  $I(X; Z)$ ?

Data processing

If  $X$ ,  $Y$ , and  $Z$  form a Markov Chain, then

$$I(X; Y) \geq I(X; Z)$$

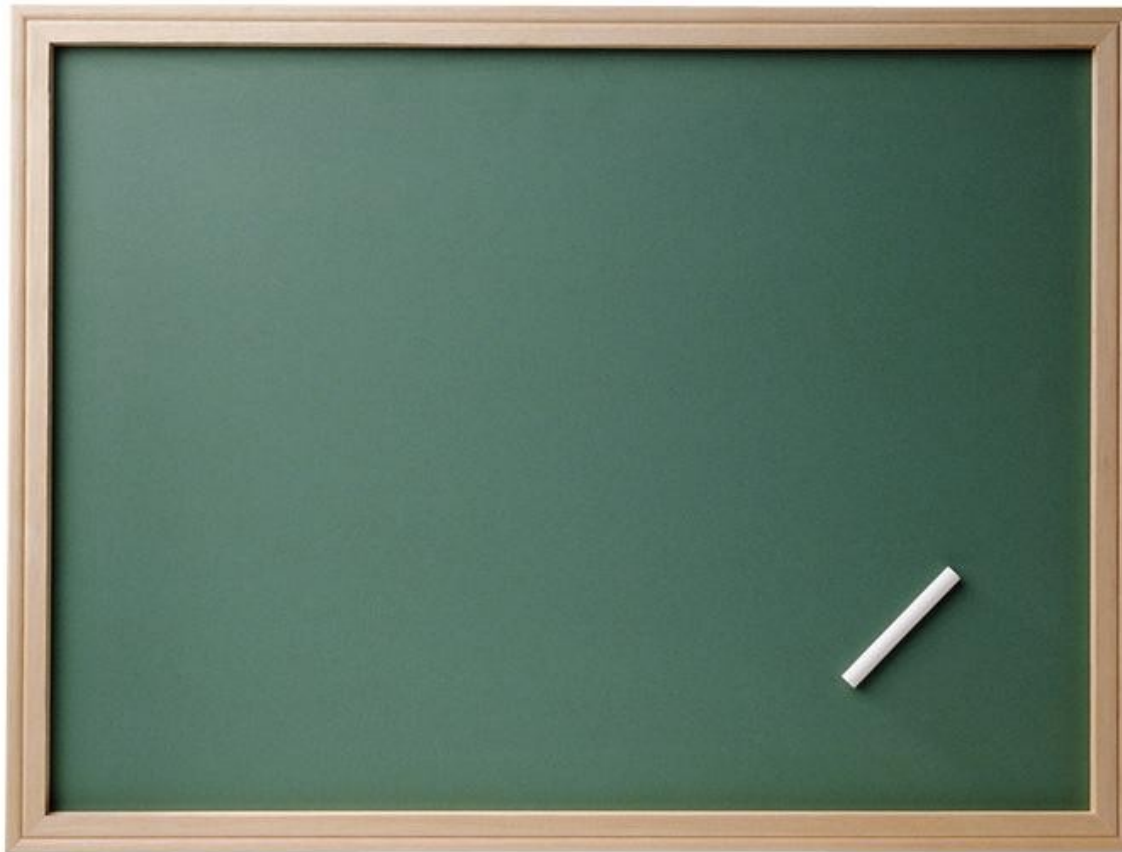


No processing of  $Y$   
can improve our  
knowledge about  $X$   
!

# PROOF

If  $X$ ,  $Y$  and  $Z$  form a Markov Chain, then

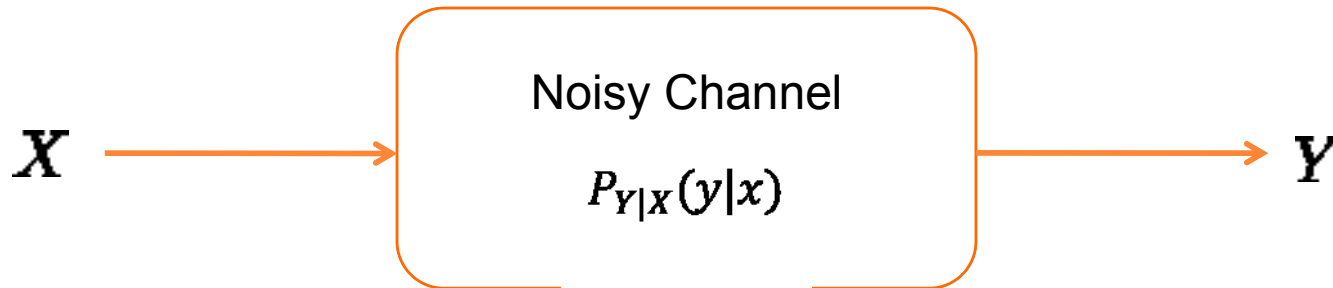
$$I(X : Y) \geq I(X : Z)$$



# SUMMARY

- Conditional entropy
  - Measures the uncertainty about  $X$  given  $Y$
- Mutual information
  - Measures how much knowing  $Y$  reduces the uncertainty about  $X$  (on average)

# HOW MUCH INFORMATION DOES THE OUTPUT GIVE ABOUT THE INPUT?

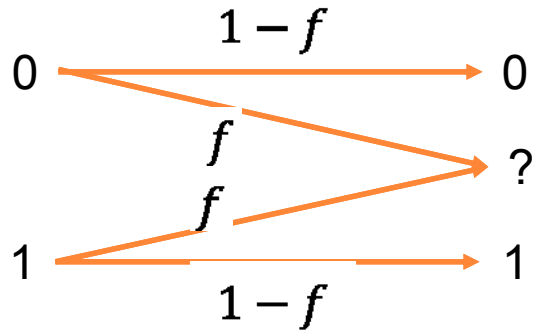


- Intuition: mutual information measures how much our uncertainty about  $X$  is reduced by learning  $Y$

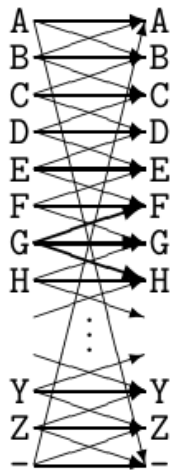
$$I(X; Y) = H(X) - H(X|Y)$$

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

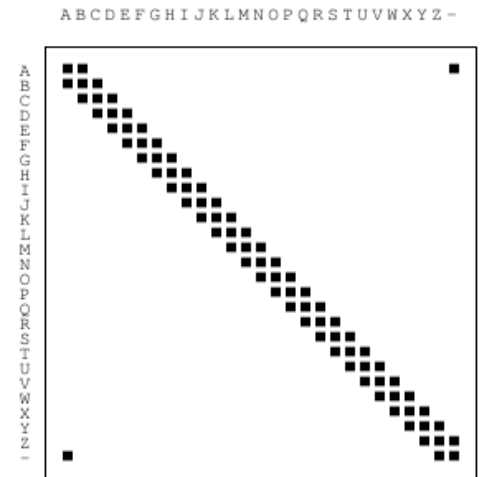
# EXAMPLE OF CHANNELS: BINARY ERASURE CHANNEL



# EXAMPLE OF CHANNELS: NOISY-TYPEWRITER

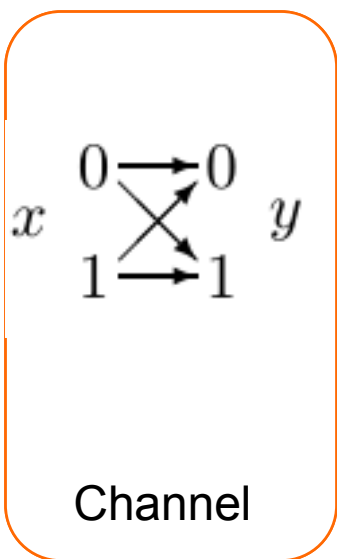


$$\begin{aligned} & \vdots \\ P(y = F | x = G) &= 1/3; \\ P(y = G | x = G) &= 1/3; \\ P(y = H | x = G) &= 1/3; \\ & \vdots \end{aligned}$$



# EXAMPLES OF CHANNELS: BSC

- Binary symmetric channel



$$\Pr_{Y|X}(y | x) = \begin{cases} 1 - f & \text{if } y = x \in \{0, 1\} \\ f & \text{otherwise} \end{cases}$$

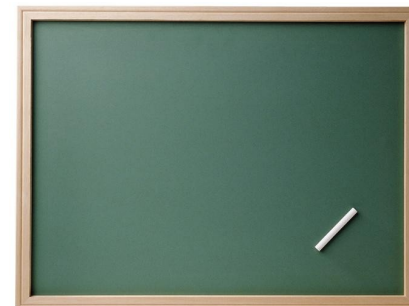
Example

$$f = 0.15$$

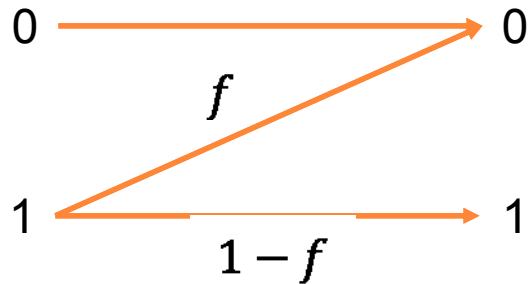
$$\Pr_X(0) = 0.9, \Pr_X(1) = 0.1$$

- Mutual information for

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



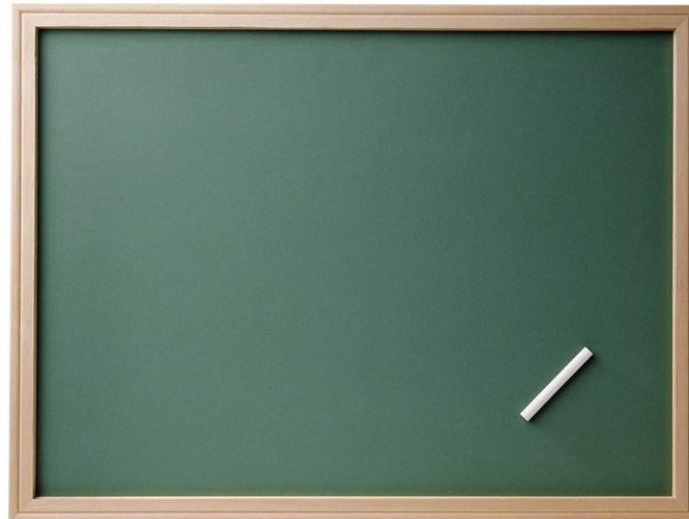
# EXAMPLE OF CHANNELS: Z-CHANNEL



$$f = 0.15$$

$$\Pr_X(0) = 0.9, \Pr_X(1) = 0.1$$

Mutual information?



# INTUITION

- How can we ensure we learn as much as possible from the output over the input?
- Remember: so far we don't "encode" we simply choose inputs  $X$  with some probability
- Mutual information measures how much the output reduces the uncertainty about the input

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

- Idea: Want to maximize the mutual information!

# CAPACITY OF A CHANNEL

- We define the capacity of a channel  $\mathcal{N}$  by:

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

- Determines how well we can send information.



## FOR THE BINARY SYMMETRIC CHANNEL

- What's  $C(\mathcal{N}) = \max_{P_X} I(X; Y)$  ?

- Remember our example

- More generally

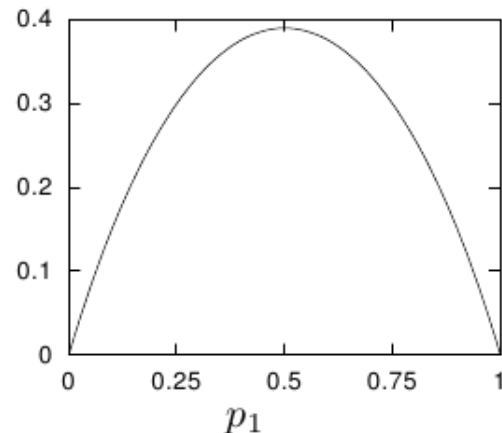
$$I(X; Y) = H_2((1 - f)p_1 + f(1 - p_1)) - H_2(f)$$

$I(X; Y)$

- Maximum at  $p_1 = 1/2$

- Capacity

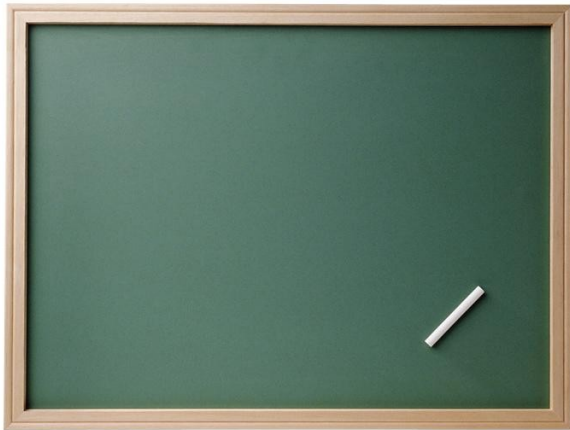
$$C(BSC) = 1 - H_2(f)$$



# HOW ABOUT THE Z-CHANNEL?

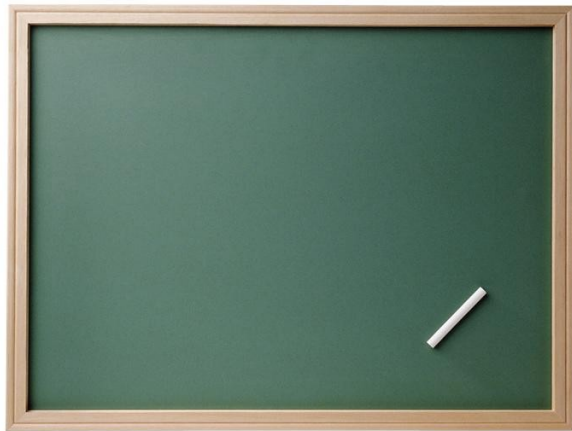
# HOW ABOUT THE Z-CHANNEL?

- Computing the mutual information



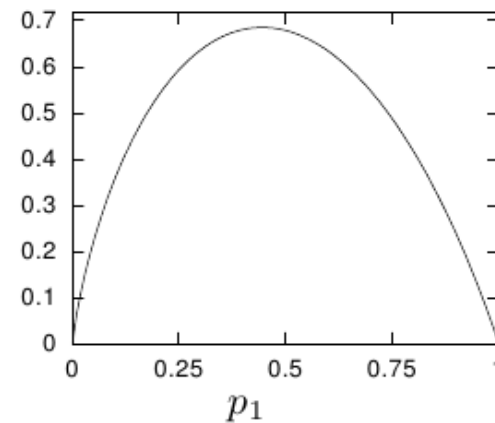
# HOW ABOUT THE Z-CHANNEL?

- Computing the mutual information



For our example  $f = 0.15$

$I(X; Y)$



# HOW ABOUT THE NOISY-TYPEWRITER?

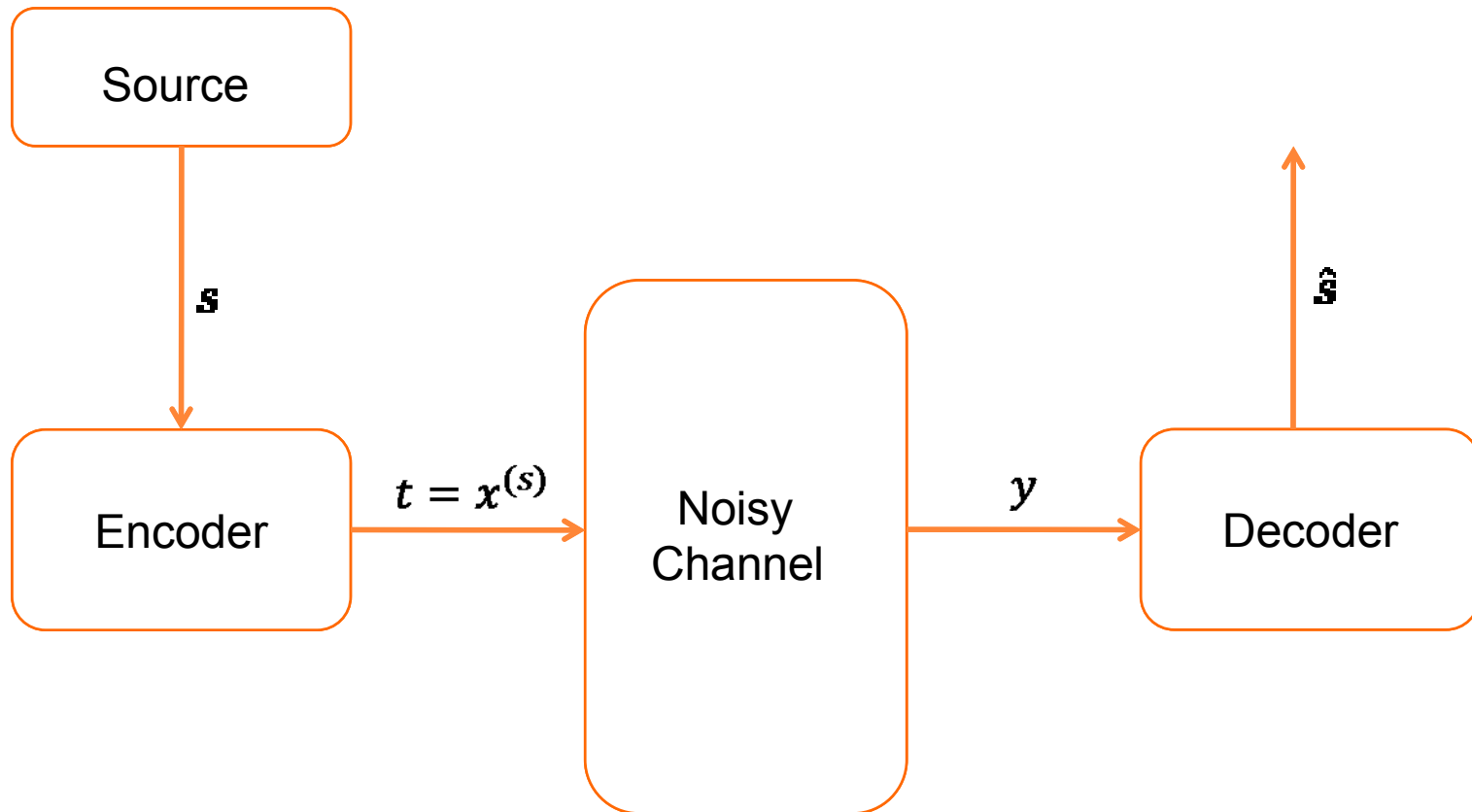
# TERMINOLOGY: BLOCK CODES

- A  $(N,K)$  block code
  - Source symbols  $s \in \{1, \dots, 2^K\}$                    ie K bits!
  - Codewords  $x^{(s)}$  of length N bits!
  
- List of possible codewords  $\{x^{(1)}, \dots, x^{(2^K)}\}$
  
- Can you name an example of block code? 😊

## TERMINOLOGY: RATE OF THE CODE

- ③ For an  $(N,K)$  block code has rate  $R = \frac{K}{N}$
- What was the rate of the  $(7,4)$  Hamming code? 😊
- Rate measures the price of transmission: to send  $K$  bits of source signal we need to send  $N = \frac{K}{R}$  bits
  - Lower rate means higher price
  - Higher rate means lower price

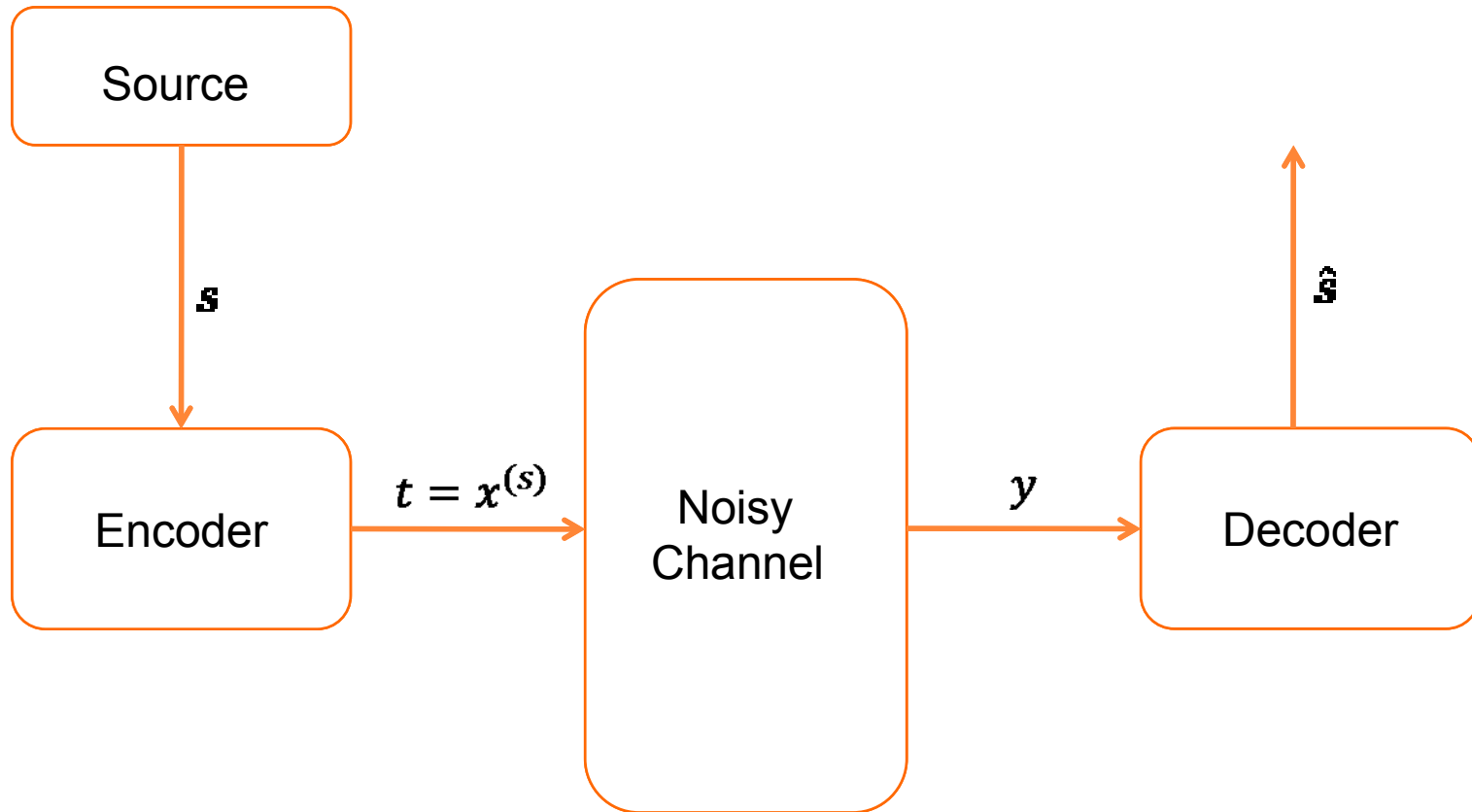
# TERMINOLOGY: ENCODER



Encoder: mapping from source symbols to channel inputs

Caution: channel inputs  $x$  are not the same as source symbols!

# TERMINOLOGY : DECODER



Decoder: mapping from channel outputs to possible source symbols

# TERMINOLOGY: OPTIMAL DECODER

- Most likely input symbol:  $\hat{s} = \arg \max_s \Pr(s | y)$

# TERMINOLOGY: PROBABILITY OF BLOCK ERROR

- Average probability of observing a wrong block

$$P_{\text{Block}} = \sum_s \text{Pr}(s) \text{Pr}(\hat{s} \neq s \mid s)$$

- A block error occurs if a single input bit is decoded incorrectly!
- Do you remember what is a bit error in comparison?

# TERMINOLOGY: MAXIMUM PROBABILITY OF BLOCK ERROR

- Worst case error

$$P_{\max} = \max_s \Pr(\hat{s} \neq s \mid s)$$

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

- Associated with every discrete memoryless (IID) channel, there is a non-negative number  $C$  (the capacity) such that
  - For any error  $\epsilon > 0$  and  $R < C$  for large enough  $N$ , there exists a block code of length  $N$  and rate  $\geq R$ , and a decoding algorithm, such that the maximum probability of block error is  $< \epsilon$
  - If a probability of bit error  $p_b$  is acceptable, rates up to  $R(p_b)$  are achievable, where

$$R(p_b) = \frac{C}{1 - H_2(p_b)}$$

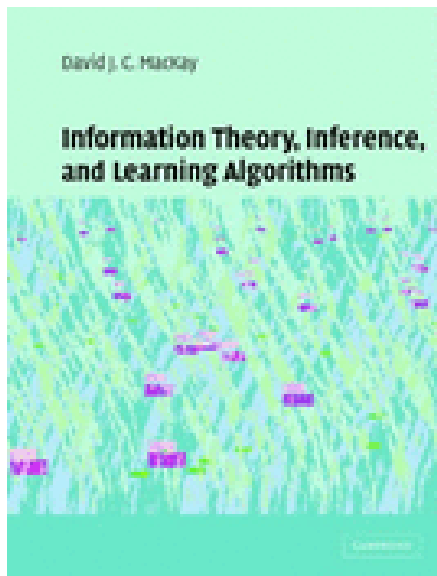
- For any  $p_b$ , rates greater than  $R(p_b)$  are not achievable

## NEXT TIME

- Explanation and proof of noisy channel coding theorem!

# READING FOR THIS LECTURE

- Chapters 8, 9.1 to 9.5 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