

# CT216 Introduction to Communication Systems

## Lecture 2: Probability, Information, and Estimation

Yash M. Vasavada

Professor, DA-IICT, Gandhinagar

24th January 2024



# Overview of Today's Talk

## 1 Generation of Information



# Overview of Today's Talk

- 1 Generation of Information
- 2 A Review of Probability Theory

- Several Example Random Variables
- Fundamental Concepts



# Overview of Today's Talk

- 1 Generation of Information
- 2 A Review of Probability Theory
  - Several Example Random Variables
  - Fundamental Concepts
- 3 Bayes' Theorem



# Overview of Today's Talk

- 1 Generation of Information
- 2 A Review of Probability Theory
  - Several Example Random Variables
  - Fundamental Concepts
- 3 Bayes' Theorem
- 4 Bayes Rule and the Likelihood Ratio Test (LRT)



# Overview of Today's Talk

- 1 Generation of Information
- 2 A Review of Probability Theory
  - Several Example Random Variables
  - Fundamental Concepts
- 3 Bayes' Theorem
- 4 Bayes Rule and the Likelihood Ratio Test (LRT)
- 5 Conditional Entropy and Information Transfer



# Overview of Today's Talk

- 1 Generation of Information
- 2 A Review of Probability Theory
  - Several Example Random Variables
  - Fundamental Concepts
- 3 Bayes' Theorem
- 4 Bayes Rule and the Likelihood Ratio Test (LRT)
- 5 Conditional Entropy and Information Transfer
- 6 Summary

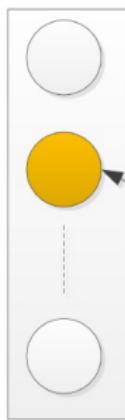




# The Model of a Communication System

Crystallized by Shannon

## 1. A Set of Messages



2. Selection of one of M messages from the set

Information Source

Selected Message

Communication Transmitter

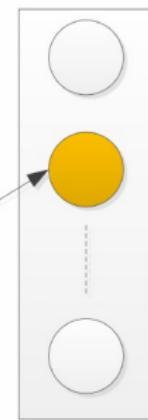
3. Signal

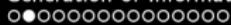
$\Sigma$

Noise

Destination

Communication Receiver





# The Definition and Measurement of Information

## The Source Alphabet

- The information source is modeled as generating symbols each of which can take  $K$  symbols from a set  $\{d_1, d_2, \dots, d_K\}$ . This set is the source alphabet.





# The Definition and Measurement of Information

## The Source Alphabet

- The information source is modeled as generating symbols each of which can take  $K$  symbols from a set  $\{d_1, d_2, \dots, d_K\}$ . This set is the source alphabet.
  - ▷ E.g., in the  $C$  disease problem, the source (i.e., the human body) is **binary**, i.e., it takes one of  $K = 2$  symbols;  $C = 0$  denotes that there is no disease,  $C = 1$  is that the body has the disease
  - ▷ When the information source (e.g., a writer) uses the English language,  $K$  has to be at least equal to the number of letters ( $= 26$ ) of the English alphabet (however, a common method for text encoding, the ASCII character set, uses  $K = 128$ )



# The Definition and Measurement of Information

## The Source Alphabet

- The information source is modeled as generating symbols each of which can take  $K$  symbols from a set  $\{d_1, d_2, \dots, d_K\}$ . This set is the source alphabet.
  - ▷ E.g., in the  $C$  disease problem, the source (i.e., the human body) is **binary**, i.e., it takes one of  $K = 2$  symbols;  $C = 0$  denotes that there is no disease,  $C = 1$  is that the body has the disease
  - ▷ When the information source (e.g., a writer) uses the English language,  $K$  has to be at least equal to the number of letters ( $= 26$ ) of the English alphabet (however, a common method for text encoding, the ASCII character set, uses  $K = 128$ )



# The Definition and Measurement of Information

## The Symbol Probabilities

- Probability  $P(D = d_k)$  of  $k^{th}$  symbol is denoted as  $p_k$ . Different values of  $p_k$  may not all be the same (i.e., some source symbols may be more likely to occur than the others). The sum of these

probabilities equals 1, i.e.,  $\sum_{k=1}^K p_k = 1$ .



# The Definition and Measurement of Information

## The Symbol Probabilities

- Probability  $P(D = d_k)$  of  $k^{th}$  symbol is denoted as  $p_k$ . Different values of  $p_k$  may not all be the same (i.e., some source symbols may be more likely to occur than the others). The sum of these

probabilities equals 1, i.e.,  $\sum_{k=1}^K p_k = 1$ .

- ▷ E.g., in the  $C$  disease problem,  $p(C = 0) = 0.99$  and  $p(C = 1) = 0.01$ .



# The Definition and Measurement of Information

## The Symbol Probabilities

- Probability  $P(D = d_k)$  of  $k^{th}$  symbol is denoted as  $p_k$ . Different values of  $p_k$  may not all be the same (i.e., some source symbols may be more likely to occur than the others). The sum of these

probabilities equals 1, i.e.,  $\sum_{k=1}^K p_k = 1$ .

- ▷ E.g., in the C disease problem,  $p(C = 0) = 0.99$  and  $p(C = 1) = 0.01$ .
- ▷ For the English alphabet, some letters such as e, a, i, t, etc. have a much higher probability of occurrence compared to the letters such as j, q, x, z.
  - See, e.g., <https://www3.nd.edu/business/handouts/cryptography/letterfrequencies.html>



# A Measure of Information Content of the Source

- What is the information  $I_k$  generated when the source emits  $k^{th}$  symbol  $d_k$  with probability  $p_k$ ? Shannon's answer:



# A Measure of Information Content of the Source

- What is the information  $I_k$  generated when the source emits  $k^{th}$  symbol  $d_k$  with probability  $p_k$ ? Shannon's answer:
  - ① Information should be inversely proportional to  $p_k$ . Why?



# A Measure of Information Content of the Source

- What is the information  $I_k$  generated when the source emits  $k^{th}$  symbol  $d_k$  with probability  $p_k$ ? Shannon's answer:
  - ① Information should be inversely proportional to  $p_k$ . Why?
    - ▷ The smaller the  $p_k$ , the greater the surprise by the occurrence of  $k^{th}$  symbol



# A Measure of Information Content of the Source

- What is the information  $I_k$  generated when the source emits  $k^{th}$  symbol  $d_k$  with probability  $p_k$ ? Shannon's answer:
  - ➊ Information should be inversely proportional to  $p_k$ . Why?
    - ▷ The smaller the  $p_k$ , the greater the surprise by the occurrence of  $k^{th}$  symbol
    - ▷ In the  $C$ -disease problem,  $p(C = 0)$  is very high (0.99), and so disease-free event is not highly surprising or **informative**. In contrast,  $C = 1$  is a low-probability event and so its occurrence conveys a greater surprise, and delivers a bigger information



# A Measure of Information Content of the Source

- What is the information  $I_k$  generated when the source emits  $k^{th}$  symbol  $d_k$  with probability  $p_k$ ? Shannon's answer:
  - ① Information should be inversely proportional to  $p_k$ . Why?
    - ▷ The smaller the  $p_k$ , the greater the surprise by the occurrence of  $k^{th}$  symbol
    - ▷ In the  $C$ -disease problem,  $p(C = 0)$  is very high (0.99), and so disease-free event is not highly surprising or **informative**. In contrast,  $C = 1$  is a low-probability event and so its occurrence conveys a greater surprise, and delivers a bigger information
    - ▷ Suppose  $C = 1$  denotes an event that “there is a snowfall in Gandhinagar”. The probability of this event is so small that in case it snows in Gandhinagar, that day would be *extremely* informative! The typical days for which  $C = 0$  are non-informative in this sense



# A Measure of Information Content of the Source

- What is the information  $I_k$  generated when the source emits  $k^{th}$  symbol  $d_k$  with probability  $p_k$ ? Shannon's answer:
  - ① Information should be inversely proportional to  $p_k$ . Why?
    - ▷ The smaller the  $p_k$ , the greater the surprise by the occurrence of  $k^{th}$  symbol
    - ▷ In the C-disease problem,  $p(C = 0)$  is very high (0.99), and so disease-free event is not highly surprising or **informative**. In contrast,  $C = 1$  is a low-probability event and so its occurrence conveys a greater surprise, and delivers a bigger information
    - ▷ Suppose  $C = 1$  denotes an event that "there is a snowfall in Gandhinagar". The probability of this event is so small that in case it snows in Gandhinagar, that day would be *extremely* informative! The typical days for which  $C = 0$  are non-informative in this sense
    - ▷ In the English language, the letters such as e carry less information compared to the letters such as q



# A Measure of Information Content of the Source

- What is the information  $I_k$  generated when the source emits  $k^{th}$  symbol  $d_k$  with probability  $p_k$ ? Shannon's answer:



# A Measure of Information Content of the Source

- What is the information  $I_k$  generated when the source emits  $k^{th}$  symbol  $d_k$  with probability  $p_k$ ? Shannon's answer:
  - ① Information should be inversely proportional to  $p_k$ .



# A Measure of Information Content of the Source

- What is the information  $I_k$  generated when the source emits  $k^{th}$  symbol  $d_k$  with probability  $p_k$ ? Shannon's answer:
  - ① Information should be inversely proportional to  $p_k$ .
  - ② Information should be a logarithmic function (the base of the logarithm would determine the unit of the information). Why?



# A Measure of Information Content of the Source

- What is the information  $I_k$  generated when the source emits  $k^{th}$  symbol  $d_k$  with probability  $p_k$ ? Shannon's answer:
  - ① Information should be inversely proportional to  $p_k$ .
  - ② Information should be a logarithmic function (the base of the logarithm would determine the unit of the information). Why?
    - ▷ When all  $p_k$  are equal, the information generated is log of  $K$



# A Measure of Information Content of the Source

- What is the information  $I_k$  generated when the source emits  $k^{th}$  symbol  $d_k$  with probability  $p_k$ ? Shannon's answer:
  - ① Information should be inversely proportional to  $p_k$ .
  - ② Information should be a logarithmic function (the base of the logarithm would determine the unit of the information). Why?
    - ▷ When all  $p_k$  are equal, the information generated is log of  $K$
    - ▷ E.g., in the  $C$ -disease problem, if  $C = 0$  and  $C = 1$  events are equally likely, they each would be assigned a single bit 0 or 1



# A Measure of Information Content of the Source

- What is the information  $I_k$  generated when the source emits  $k^{th}$  symbol  $d_k$  with probability  $p_k$ ? Shannon's answer:
  - ① Information should be inversely proportional to  $p_k$ .
  - ② Information should be a logarithmic function (the base of the logarithm would determine the unit of the information). Why?
    - ▷ When all  $p_k$  are equal, the information generated is log of  $K$
    - ▷ E.g., in the  $C$ -disease problem, if  $C = 0$  and  $C = 1$  events are equally likely, they each would be assigned a single bit 0 or 1
    - ▷ Similarly, if a source alphabet has  $K = 16$  symbols, total  $\log_2(16) = 4$  bits are required to represent all of them (0, 0, 0, 0, 0, 0, 0, 1, ..., 1, 1, 1, 0, and 1, 1, 1, 1).



# A Measure of Information Content of the Source

- What is the information  $I_k$  generated when the source emits  $k^{th}$  symbol  $d_k$  with probability  $p_k$ ? Shannon's answer:
  - ① Information should be inversely proportional to  $p_k$ .
  - ② Information should be a logarithmic function (the base of the logarithm would determine the unit of the information). Why?
    - ▷ When all  $p_k$  are equal, the information generated is log of  $K$
    - ▷ E.g., in the  $C$ -disease problem, if  $C = 0$  and  $C = 1$  events are equally likely, they each would be assigned a single bit 0 or 1
    - ▷ Similarly, if a source alphabet has  $K = 16$  symbols, total  $\log_2(16) = 4$  bits are required to represent all of them (0, 0, 0, 0, 0, 0, 0, 1, ..., 1, 1, 1, 0, and 1, 1, 1, 1).
  - Information generated by  $k^{th}$  symbol is defined as:

$$I_k = \log_2 \left( \frac{1}{p_k} \right) = -\log_2 p_k \text{ bits.}$$



# Probability Mass Functions or PMFs

- The set of the probabilities  $\{p_1, p_2, \dots, p_K\}$  of all  $K$  symbols of the source is called the Probability Mass Function or PMF



# Probability Mass Functions or PMFs

- The set of the probabilities  $\{p_1, p_2, \dots, p_K\}$  of all  $K$  symbols of the source is called the Probability Mass Function or PMF
- The PMF has the following three key properties:



# Probability Mass Functions or PMFs

- The set of the probabilities  $\{p_1, p_2, \dots, p_K\}$  of all  $K$  symbols of the source is called the Probability Mass Function or PMF
- The PMF has the following three key properties:
  - ①  $p_k \geq 0$



# Probability Mass Functions or PMFs

- The set of the probabilities  $\{p_1, p_2, \dots, p_K\}$  of all  $K$  symbols of the source is called the Probability Mass Function or PMF
- The PMF has the following three key properties:

$$\textcircled{1} \quad p_k \geq 0$$

$$\textcircled{2} \quad \sum_{k=1}^K p_k = 1$$



# Probability Mass Functions or PMFs

- The set of the probabilities  $\{p_1, p_2, \dots, p_K\}$  of all  $K$  symbols of the source is called the Probability Mass Function or PMF
- The PMF has the following three key properties:

$$\textcircled{1} \quad p_k \geq 0 \quad \underset{K}{\text{for all } k}$$

$$\textcircled{2} \quad \sum_{k=1}^K p_k = 1$$

$$\textcircled{3} \quad P(k_a \leq k \leq k_b) = \sum_{k=k_a}^{k_b} p_k$$



# Probability Mass Functions or PMFs

- The set of the probabilities  $\{p_1, p_2, \dots, p_K\}$  of all  $K$  symbols of the source is called the Probability Mass Function or PMF
- The PMF has the following three key properties:

$$\textcircled{1} \quad p_k \geq 0 \quad \underset{K}{\text{for all } k}$$

$$\textcircled{2} \quad \sum_{k=1}^K p_k = 1$$

$$\textcircled{3} \quad P(k_a \leq k \leq k_b) = \sum_{k=k_a}^{k_b} p_k$$



# The Random Variables (RVs)

Main Point: The RVs are *always* numbers (integers, fractionals, real-numbers, and sometimes even complex numbers)



# The Random Variables (RVs)

Main Point: The RVs are *always* numbers (integers, fractionals, real-numbers, and sometimes even complex numbers)

- Let the symbol  $d_k$  be assigned a unique number  $x_k$ . This number  $X = x_k$  that the source output takes is called a **random variable** or an RV in short.



# The Random Variables (RVs)

Main Point: The RVs are *always* numbers (integers, fractionals, real-numbers, and sometimes even complex numbers)

Let the symbol  $d_k$  be assigned a unique number  $x_k$ . This number  $X = x_k$  that the source output takes is called a **random variable** or an RV in short.

Only for the RVs, there are two statistical concepts: the expected (or mean) value of the RV and the variance (the square of standard deviation) of the RV

# The Random Variables (RVs)

Main Point: The RVs are *always* numbers (integers, fractionals, real-numbers, and sometimes even complex numbers)

- Let the symbol  $d_k$  be assigned a unique number  $x_k$ . This number  $X = x_k$  that the source output takes is called a **random variable** or an RV in short.
- Only for the RVs, there are two statistical concepts: the expected (or mean) value of the RV and the variance (the square of standard deviation) of the RV

$$\bullet \mathbb{E}[X] = \bar{X} = \sum_{k=1}^K x_k p_k$$

$$\bullet \text{Var}[X] = \sum_{k=1}^K [x_k - \bar{X}]^2 p_k = \mathbb{E}[X^2] - (\mathbb{E}[X])^2$$



# Cumulative Mass Function (CMF)

- Defined using the PMF — but only for an RV



# Cumulative Mass Function (CMF)

- Defined using the PMF — but only for an RV
- Definition:  $F_X(x_m) = P(X \leq x_m) = \sum_{k=1}^m p_k$
- Properties:
  - ①  $F_X(x_m)$  is monotonically nondecreasing
  - ②  $F(-\infty) = 0$
  - ③  $F(\infty) = 1$
  - ④  $P(x_a < X \leq x_b) = F(x_b) - F(x_a)$
- The CMF also completely defines the probability distribution of an RV



# A Measure of Information Content of the Source

- Besides the information generated by the individual symbol of the source, Shannon showed that there is an important measure of how much the source is informative in an average sense



# A Measure of Information Content of the Source

- Besides the information generated by the individual symbol of the source, Shannon showed that there is an important measure of how much the source is informative in an average sense
- The average informative-ness of the source  $X$  is called the Entropy  $H(X)$  and it equals

$$H(X) = \mathbb{E}[I_k]$$

Note that the information  $I_k$  is a *number* and so it qualifies to be an RV whose expectation can be taken



# A Measure of Information Content of the Source

- Besides the information generated by the individual symbol of the source, Shannon showed that there is an important measure of how much the source is informative in an average sense
- The average informative-ness of the source  $X$  is called the Entropy  $H(X)$  and it equals

$$H(X) = \mathbb{E}[I_k]$$

Note that the information  $I_k$  is a *number* and so it qualifies to be an RV whose expectation can be taken

- From the definition of the expected value of an RV,  $H(X)$  can be written as

$$H(X) = \mathbb{E}[I_k] = - \sum_{k=1}^K p_k \log_2(p_k)$$



# A Binary Random Variable

Called Bernoulli( $q$ ) RV  $X$

- Outcome of the toss of a coin:  $p(X) = \begin{cases} 1 - q, & X = 0 \text{ (head),} \\ q, & X = 1 \text{ (tail)} \end{cases}$



# A Binary Random Variable

Called Bernoulli( $q$ ) RV  $X$

- Outcome of the toss of a coin:  $p(X) = \begin{cases} 1 - q, & X = 0 \text{ (head),} \\ q, & X = 1 \text{ (tail)} \end{cases}$
- Mean:  $m_X = \sum_x x p_X(x) = 0 \times (1 - q) + 1 \times q = q$



# A Binary Random Variable

Called Bernoulli( $q$ ) RV  $X$

- Outcome of the toss of a coin:  $p(X) = \begin{cases} 1 - q, & X = 0 \text{ (head),} \\ q, & X = 1 \text{ (tail)} \end{cases}$ 
  - Mean:  $m_X = \sum_x x p_X(x) = 0 \times (1 - q) + 1 \times q = q$
  - Variance:  $\sigma_X^2 = q(1 - q)$



# A Binary Random Variable

Called Bernoulli( $q$ ) RV  $X$

- Outcome of the toss of a coin:  $p(X) = \begin{cases} 1 - q, & X = 0 \text{ (head)}, \\ q, & X = 1 \text{ (tail)} \end{cases}$

$$\rightarrow \text{Mean: } m_X = \sum_x x p_X(x) = 0 \times (1 - q) + 1 \times q = q$$

$$\rightarrow \text{Variance: } \sigma_X^2 = q(1 - q)$$

$$\rightarrow \text{Information } I_k = -\log_2 q \text{ and } -\log_2(1 - q) \text{ bits}$$



# A Binary Random Variable

Called Bernoulli( $q$ ) RV  $X$

- Outcome of the toss of a coin:  $p(X) = \begin{cases} 1 - q, & X = 0 \text{ (head),} \\ q, & X = 1 \text{ (tail)} \end{cases}$

$$\rightarrow \text{Mean: } m_X = \sum_x x p_X(x) = 0 \times (1 - q) + 1 \times q = q$$

$$\rightarrow \text{Variance: } \sigma_X^2 = q(1 - q)$$

$\rightarrow$  Information  $I_k = -\log_2 q$  and  $-\log_2(1 - q)$  bits

$\rightarrow$  Average information or Entropy

$$H(X) = -q \log_2 q - (1 - q) \log_2(1 - q) \text{ bits}$$



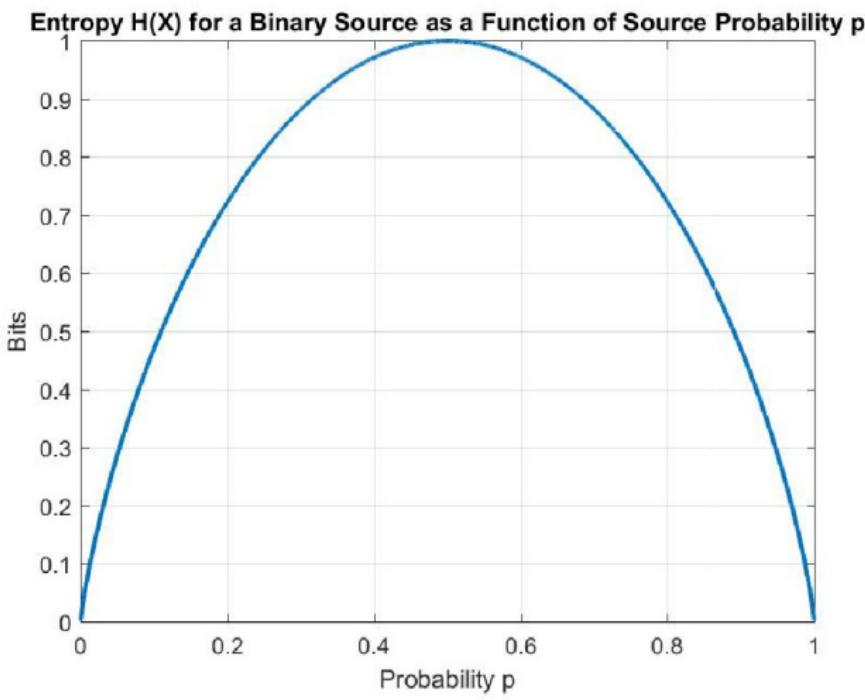
# A Binary Random Variable

Called Bernoulli( $q$ ) RV  $X$

- Outcome of the toss of a coin:  $p(X) = \begin{cases} 1 - q, & X = 0 \text{ (head),} \\ q, & X = 1 \text{ (tail)} \end{cases}$ 
  - Mean:  $m_X = \sum_x x p_X(x) = 0 \times (1 - q) + 1 \times q = q$
  - Variance:  $\sigma_X^2 = q(1 - q)$
  - Information  $I_k = -\log_2 q$  and  $-\log_2(1 - q)$  bits
  - Average information or Entropy  
 $H(X) = -q \log_2 q - (1 - q) \log_2(1 - q)$  bits
    - ▷  $H(X)$  takes the highest value of 1 bit when  $q = 0.5$ , i.e., when  $X = 0$  and  $X = 1$  both are equally likely



# Entropy of Bernoulli( $p$ ) RV



# Summary

What we have learnt so far

- We have seen a succinct definition of the information generated by a source
  - ▷ The source can be an author writing in English, Hindi, Gujarati, etc., or it can be a camera taking a picture, etc.



# Summary

What we have learnt so far

- We have seen a succinct definition of the information generated by a source
  - ▷ The source can be an author writing in English, Hindi, Gujarati, etc., or it can be a camera taking a picture, etc.
  - ▷ The source has to be discrete in nature; someone speaking is also an information source, but his speech waveform is continuous. We have not defined how to define/measure the information generated by him/her



# Summary

What we have learnt so far

- We have seen a succinct definition of the information generated by a source
  - ▷ The source can be an author writing in English, Hindi, Gujarati, etc., or it can be a camera taking a picture, etc.
  - ▷ The source has to be discrete in nature; someone speaking is also an information source, but his speech waveform is continuous. We have not defined how to define/measure the information generated by him/her
- This definition agrees with our intuitive notion about how information should be measured



# Summary

What we have learnt so far

- We have seen a succinct definition of the information generated by a source
  - ▷ The source can be an author writing in English, Hindi, Gujarati, etc., or it can be a camera taking a picture, etc.
  - ▷ The source has to be discrete in nature; someone speaking is also an information source, but his speech waveform is continuous. We have not defined how to define/measure the information generated by him/her
- This definition agrees with our intuitive notion about how information should be measured
  - the more surprising an event is, a greater amount of information  $I_k$  is defined to have been generated



# Summary

What we have learnt so far

- We have seen a succinct definition of the information generated by a source
  - ▷ The source can be an author writing in English, Hindi, Gujarati, etc., or it can be a camera taking a picture, etc.
  - ▷ The source has to be discrete in nature; someone speaking is also an information source, but his speech waveform is continuous. We have not defined how to define/measure the information generated by him/her
- This definition agrees with our intuitive notion about how information should be measured
  - the more surprising an event is, a greater amount of information  $I_k$  is defined to have been generated
- The definition also provides the units, e.g., in bits



# Summary

What we have learnt so far

- We have seen a succinct definition of the information generated by a source
  - ▷ The source can be an author writing in English, Hindi, Gujarati, etc., or it can be a camera taking a picture, etc.
  - ▷ The source has to be discrete in nature; someone speaking is also an information source, but his speech waveform is continuous. We have not defined how to define/measure the information generated by him/her
- This definition agrees with our intuitive notion about how information should be measured
  - the more surprising an event is, a greater amount of information  $I_k$  is defined to have been generated
- The definition also provides the units, e.g., in bits
  - The more informative an event is, a greater number of bits are required to represent it



# Summary

What we have learnt so far

- An information source generates, or it causes, one event to occur out of a set of multiple events
- Alternatively, it **selects** one information-bearing message out of a set of messages (these messages or events are called the symbols)
- While an information measure  $I_k$  is defined for each event/message/symbol  $k$ , the information generated on average by the source itself is measured using the Entropy  $H(X)$



# Summary

What we have learnt so far

- The Entropy  $H(X)$  captures the notion of information in a better manner compared to  $I_k$  — if  $I_k$  is large, only an individual source symbol is highly informative, whereas if  $H(X)$  is large, the source as a whole is highly informative



# Summary

What we have learnt so far

- The Entropy  $H(X)$  captures the notion of information in a better manner compared to  $I_k$  — if  $I_k$  is large, only an individual source symbol is highly informative, whereas if  $H(X)$  is large, the source as a whole is highly informative
  - ▷ An analogy: a few molecules in a sealed container filled with the gas may be vibrating a lot, but the overall temperature of the gas inside the container will be low if most of the molecules of the gas are not vibrating much;
  - ▷ In contrast, if all the molecules are vibrating about equally (i.e., if  $p_k$  is equal to  $1/K$  for all symbol indices  $k$ ), the temperature will be highest



# Summary

What we have learnt so far

- The Entropy  $H(X)$  captures the notion of information in a better manner compared to  $I_k$  — if  $I_k$  is large, only an individual source symbol is highly informative, whereas if  $H(X)$  is large, the source as a whole is highly informative
  - ▷ An analogy: a few molecules in a sealed container filled with the gas may be vibrating a lot, but the overall temperature of the gas inside the container will be low if most of the molecules of the gas are not vibrating much;
  - ▷ In contrast, if all the molecules are vibrating about equally (i.e., if  $p_k$  is equal to  $1/K$  for all symbol indices  $k$ ), the temperature will be highest
  - ▷ Thus, the entropy  $H(X)$  gives an idea of the “temperature” of the source (the “temperature”, here, indicates the informative-ness of the source)



# Summary

What we have learnt so far

- The Entropy  $H(X)$  captures the notion of information in a better manner compared to  $I_k$  — if  $I_k$  is large, only an individual source symbol is highly informative, whereas if  $H(X)$  is large, the source as a whole is highly informative
  - ▷ An analogy: a few molecules in a sealed container filled with the gas may be vibrating a lot, but the overall temperature of the gas inside the container will be low if most of the molecules of the gas are not vibrating much;
  - ▷ In contrast, if all the molecules are vibrating about equally (i.e., if  $p_k$  is equal to  $1/K$  for all symbol indices  $k$ ), the temperature will be highest
  - ▷ Thus, the entropy  $H(X)$  gives an idea of the “temperature” of the source (the “temperature”, here, indicates the informative-ness of the source)
  - ▷ A source that is the *hottest* is the one that has the highest Entropy i.e., whose all  $K$  symbols occur with equal probability (i.e., its output is least predictable)



# Summary

What we have learnt so far

- We have obtained our answer to at least first few fundamental questions regarding the information and its generation, except



# Summary

What we have learnt so far

- We have obtained our answer to at least first few fundamental questions regarding the information and its generation, except
  - ▷ The answer is a bit ad-hoc; we have appealed to intuition, not a rigorous mathematical derivation





# Summary

What we have learnt so far

- We have obtained our answer to at least first few fundamental questions regarding the information and its generation, except
  - ▷ The answer is a bit ad-hoc; we have appealed to intuition, not a rigorous mathematical derivation
- We will now move on to the question of the information transfer over a communication channel



## Several Example Random Variables

## A Nonbinary Random Variable

## Binomial RV

- Let  $Y = \sum_{i=1}^n X_i$ , where  $\{X_i\}; i = 1, \dots, N$  are independent binary RVs with  $p(x) = \begin{cases} 1 - p, & x = 0(\text{head}), \\ p, & x = 1(\text{tail}) \end{cases}$
- In this case, RV  $Y$  follows the Binomial Distribution given as  $P_Y(y) = \binom{N}{y} p^y (1 - p)^{N-y}$
- Mean  $m_y = N \times p$
- Variance  $\sigma_y^2 = N \times p \times (1 - p)$



## Several Example Random Variables

## Example PMFs

## Binomial Distribution

- Let  $Y = \sum_{i=1}^n X_i$ , where  $\{X_i\}; i = 1, \dots, N$  are independent binary

RVs with  $p(x) = \begin{cases} 1 - p, & x = 0(\text{head}), \\ p, & x = 1(\text{tail}) \end{cases}$

- In this case, RV  $Y$  follows the Binomial Distribution given as  
 $P_Y(y) = \binom{N}{y} p^y (1 - p)^{N-y}$
- Mean  $m_y = N \times p$
- Variance  $\sigma_y^2 = N \times p \times (1 - p)$



# Probability Density Function (PDF)

- The concept of PMF generalizes to the PDF when the RV  $X$  is continuous
- Definition:  $p_X(x) = \frac{dF(x)}{dx}$
- Interpretations: PDF measures
  - how likely a random variable is to lie at a particular value
- Properties
  - $p(x) \geq 0$
  - $\int_{-\infty}^{\infty} p(x)dx = 1$
  - $P(a \leq X \leq b) = \int_a^b p(x)dx$
  - Unlike the PMFs, the value of PDF for a given value of  $X$  can exceed 1.



## Several Example Random Variables

## Cumulative Distribution Function (CDF)

- Definition:  $F_X(x) = F(x) = P(X \leq x)$
- Properties:
  - ①  $F(x)$  is monotonically nondecreasing
  - ②  $F(-\infty) = 0$
  - ③  $F(\infty) = 1$
  - ④  $P(a < X \leq b) = F(b) - F(a)$
- The CDF also completely defines the probability distribution of an RV



# Expected Values

- Sometimes the PDF is cumbersome to specify, or it may not be known
- Expected values are shorthand ways of describing the behavior of RVs
- Most important examples are:

$$\rightarrow \text{Mean: } E(x) = m_x = \int_{-\infty}^{\infty} x p(x) dx$$

$$\rightarrow \text{Variance: } E((x - m_x)^2) = \int_{-\infty}^{\infty} (x - m_x)^2 p(x) dx$$

- Expectation operator works with any function  $Y = g(X)$ .

$$\rightarrow E(Y) = E(g(X)) = \int_{-\infty}^{\infty} g(x) p(x) dx$$

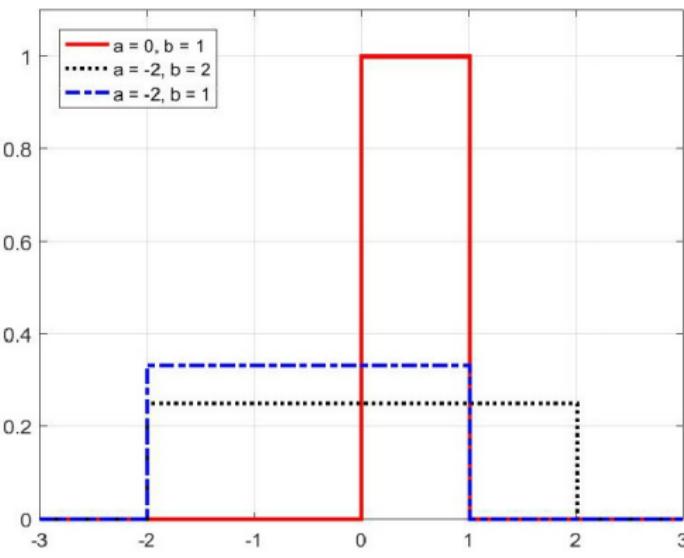


## Several Example Random Variables

## Example RVs

## Uniform PDF

$$\bullet p(x) = \begin{cases} \frac{1}{b-a}, & a \leq x \leq b \\ 0, & \text{else} \end{cases}$$

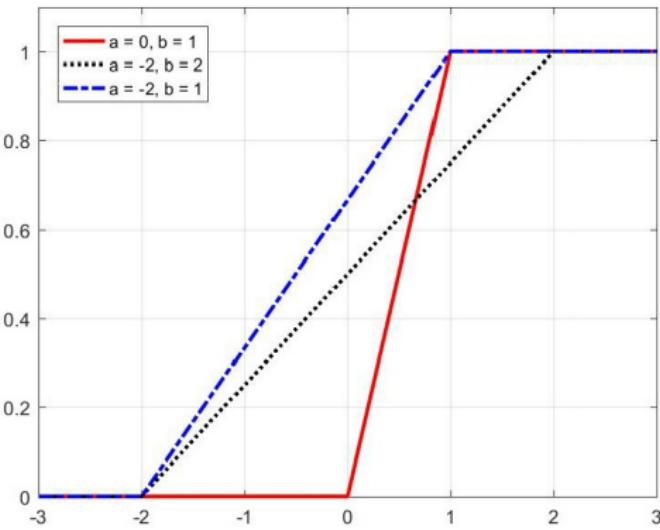


## Several Example Random Variables

## Example PDFs

## Uniform CDF

- $F(x) = \begin{cases} 0, & x \leq a \\ \frac{x-a}{b-a}, & a \leq x \leq b \\ 1, & x \geq b \end{cases}$



## Several Example Random Variables

## Example PDFs

## Uniform RV

- Mean:  $m_x = \int_a^b x p(x) dx = \frac{1}{b-a} \int_a^b x dx = \frac{a+b}{2}$
- Variance:  $\sigma_x^2 = \int_a^b (x - m_x)^2 p(x) dx = \frac{(b-a)^2}{12}$
- Probability:  
$$P(a_1 \leq x < b_1) = \int_{a_1}^{b_1} p(x) dx = \frac{b_1 - a_1}{b - a}, \quad a < a_1, b_1 < b$$

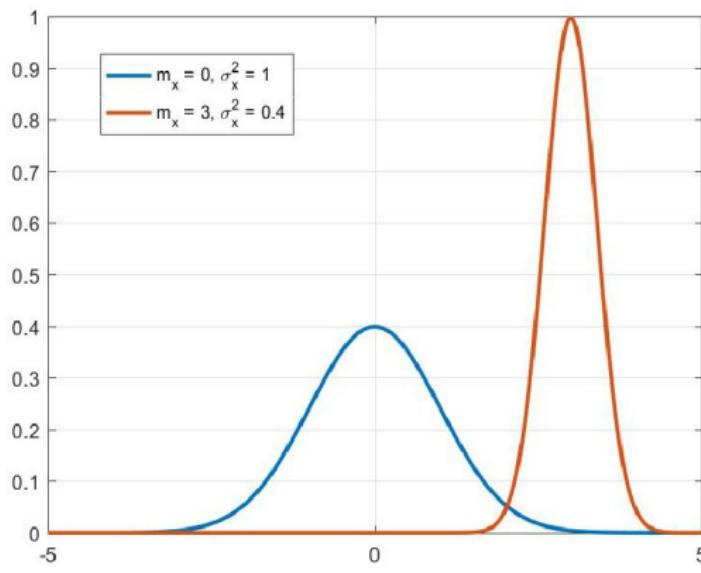


## Several Example Random Variables

## Example RVs

## Gaussian PDF

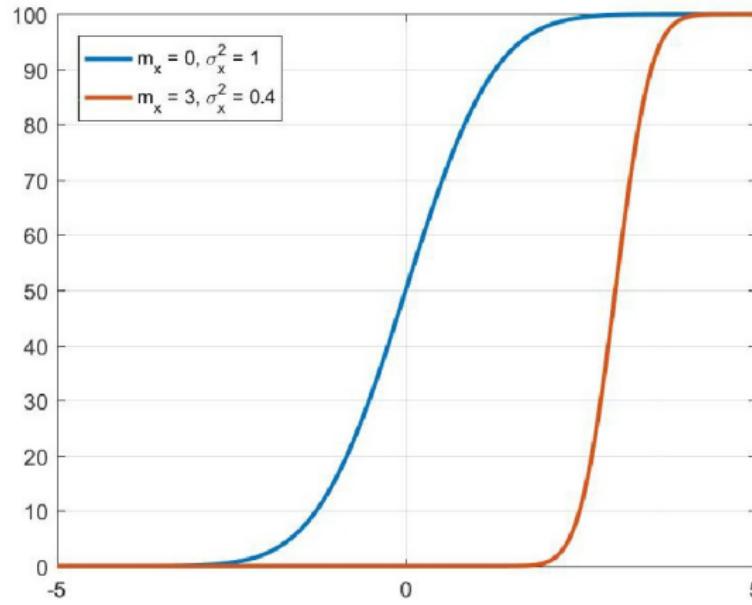
$$\bullet p(x) = \frac{1}{\sqrt{2\pi\sigma_x^2}} \exp\left(-\frac{(x - m_x)^2}{2\sigma_x^2}\right)$$



## Several Example Random Variables

## Example RVs

## Gaussian CDF



## Several Example Random Variables

## Example RVs

## Gaussian PDF

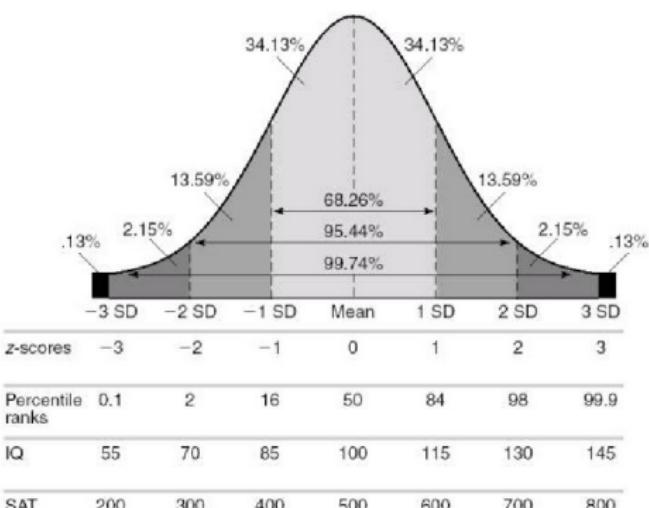


FIGURE 15.8 Percentile ranks and standard scores in relation to the normal curve.

[http://turtleinvestor888.blogspot.in/2012\\_01\\_01\\_archive.html](http://turtleinvestor888.blogspot.in/2012_01_01_archive.html)

## Fundamental Concepts

# Several Ways of Defining the Probability

Ratio of Favorable to Total Outcomes

- Ratio of favorable outcomes to the total number of outcomes



Fundamental Concepts

# Several Ways of Defining the Probability

Ratio of Favorable to Total Outcomes

- Ratio of favorable outcomes to the total number of outcomes
  - ▷ What is the probability that the roll of a six-faced die will have an even-numbered outcome?



# Several Ways of Defining the Probability

Ratio of Favorable to Total Outcomes

- Ratio of favorable outcomes to the total number of outcomes
  - ▷ What is the probability that the roll of a six-faced die will have an even-numbered outcome?
  - ▷ Favorable outcomes are half of the total outcomes, hence the probability is 0.5.



# Several Ways of Defining the Probability

Ratio of Favorable to Total Outcomes

- Ratio of favorable outcomes to the total number of outcomes
  - ▷ What is the probability that the roll of a six-faced die will have an even-numbered outcome?
  - ▷ Favorable outcomes are half of the total outcomes, hence the probability is 0.5.
- Problems with this definition:



# Several Ways of Defining the Probability

Ratio of Favorable to Total Outcomes

- Ratio of favorable outcomes to the total number of outcomes
  - ▷ What is the probability that the roll of a six-faced die will have an even-numbered outcome?
  - ▷ Favorable outcomes are half of the total outcomes, hence the probability is 0.5.
- Problems with this definition:
  - ▷ Cannot be used when the number of outcomes is infinite



## Fundamental Concepts

# Several Ways of Defining the Probability

Ratio of Favorable to Total Outcomes

- Ratio of favorable outcomes to the total number of outcomes
  - ▷ What is the probability that the roll of a six-faced die will have an even-numbered outcome?
  - ▷ Favorable outcomes are half of the total outcomes, hence the probability is 0.5.
- Problems with this definition:
  - ▷ Cannot be used when the number of outcomes is infinite
  - ▷ Does not work when the outcomes are not equally likely (e.g., roll of a biased coin)



## Fundamental Concepts

# Several Ways of Defining the Probability

## Relative Frequency Definition

- Run the probabilistic experiment a large number of times (say,  $N$  times) and count the number of times (say,  $N_T$ ) a particular event  $T$  occurs. The relative frequency is the ratio  $N_T/N$



## Fundamental Concepts

# Several Ways of Defining the Probability

## Relative Frequency Definition

- Run the probabilistic experiment a large number of times (say,  $N$  times) and count the number of times (say,  $N_T$ ) a particular event  $T$  occurs. The relative frequency is the ratio  $N_T/N$ 
  - What is the probability that the toss of coin will come up Tails?



## Fundamental Concepts

# Several Ways of Defining the Probability

## Relative Frequency Definition

- Run the probabilistic experiment a large number of times (say,  $N$  times) and count the number of times (say,  $N_T$ ) a particular event  $T$  occurs. The relative frequency is the ratio  $N_T/N$ 
  - What is the probability that the toss of coin will come up Tails?
  - Toss the coin  $N$  times and measure the number of times the event  $T$ (ails) occurs.



## Fundamental Concepts

# Several Ways of Defining the Probability

## Relative Frequency Definition

- Run the probabilistic experiment a large number of times (say,  $N$  times) and count the number of times (say,  $N_T$ ) a particular event  $T$  occurs. The relative frequency is the ratio  $N_T/N$ 
  - What is the probability that the toss of coin will come up Tails?
  - Toss the coin  $N$  times and measure the number of times the event  $T$ (ails) occurs.
- The probability of occurrence of  $T$  is defined as the limit or asymptotic value of the relative frequency as  $N \rightarrow \infty$ , i.e., as  
 $P_T = \lim_{N \rightarrow \infty} N_T/N.$



## Fundamental Concepts

# Several Ways of Defining the Probability

## Relative Frequency Definition

- Run the probabilistic experiment a large number of times (say,  $N$  times) and count the number of times (say,  $N_T$ ) a particular event  $T$  occurs. The relative frequency is the ratio  $N_T/N$ 
  - What is the probability that the toss of coin will come up Tails?
  - Toss the coin  $N$  times and measure the number of times the event  $T$ (ails) occurs.
- The probability of occurrence of  $T$  is defined as the limit or asymptotic value of the relative frequency as  $N \rightarrow \infty$ , i.e., as
$$P_T = \lim_{N \rightarrow \infty} N_T/N.$$
- Problems with this definition:



## Fundamental Concepts

# Several Ways of Defining the Probability

## Relative Frequency Definition

- Run the probabilistic experiment a large number of times (say,  $N$  times) and count the number of times (say,  $N_T$ ) a particular event  $T$  occurs. The relative frequency is the ratio  $N_T/N$ 
  - What is the probability that the toss of coin will come up Tails?
  - Toss the coin  $N$  times and measure the number of times the event  $T$ (ails) occurs.
- The probability of occurrence of  $T$  is defined as the limit or asymptotic value of the relative frequency as  $N \rightarrow \infty$ , i.e., as  
 $P_T = \lim_{N \rightarrow \infty} N_T/N$ .
- Problems with this definition:
  - No one has time to run the experiment infinite number of times



## Fundamental Concepts

# Several Ways of Defining the Probability

## Relative Frequency Definition

- Run the probabilistic experiment a large number of times (say,  $N$  times) and count the number of times (say,  $N_T$ ) a particular event  $T$  occurs. The relative frequency is the ratio  $N_T/N$ 
  - What is the probability that the toss of coin will come up Tails?
  - Toss the coin  $N$  times and measure the number of times the event  $T$ (ails) occurs.
- The probability of occurrence of  $T$  is defined as the limit or asymptotic value of the relative frequency as  $N \rightarrow \infty$ , i.e., as  
$$P_T = \lim_{N \rightarrow \infty} N_T/N.$$
- Problems with this definition:
  - No one has time to run the experiment infinite number of times
  - When  $N$  is not large enough, the relative frequency  $N_T/N$  will itself be a random number.



## Fundamental Concepts

# Several Ways of Defining the Probability

## Relative Frequency Definition

- Run the probabilistic experiment a large number of times (say,  $N$  times) and count the number of times (say,  $N_T$ ) a particular event  $T$  occurs. The relative frequency is the ratio  $N_T/N$ 
  - ▷ What is the probability that the toss of coin will come up Tails?
  - ▷ Toss the coin  $N$  times and measure the number of times the event  $T$ (ails) occurs.
- The probability of occurrence of  $T$  is defined as the limit or asymptotic value of the relative frequency as  $N \rightarrow \infty$ , i.e., as
$$P_T = \lim_{N \rightarrow \infty} N_T/N.$$
- Problems with this definition:
  - ▷ No one has time to run the experiment infinite number of times
  - ▷ When  $N$  is not large enough, the relative frequency  $N_T/N$  will itself be a random number.
    - If you run the sequence of  $N$  experiments today, you will get one answer; when you run the same sequence of  $N$  experiments tomorrow, you may get (at least a slightly) different answer
  - ▷ Cannot define the probability using a ratio whose outcome is itself probabilistic!



## Fundamental Concepts

# Several Ways of Defining the Probability

## An Axiomatic Set-Theoretic Definition

- The mathematicians struggled with the slippery-slope definitions of the probability
- Eventually, the following systematic approach was proposed by Kolmogorov, which is now widely accepted
  - ▷ This approach does not try to answer the question of “how do we know that the probability of the tails coming up on the toss of a coin is 0.5?”
  - ▷ Instead it defines a consistent and methodical way of treating the problems related to the probability



## Fundamental Concepts

# Several Ways of Defining the Probability

## An Axiomatic Set-Theoretic Definition

- The mathematicians struggled with the slippery-slope definitions of the probability
- Eventually, the following systematic approach was proposed by Kolmogorov, which is now widely accepted
  - ▷ This approach does not try to answer the question of “how do we know that the probability of the tails coming up on the toss of a coin is 0.5?”
  - ▷ Instead it defines a consistent and methodical way of treating the problems related to the probability
- This approach considers atomic events and a universal set made up of all the atomic events.
- Each set member, and all the possible subsets, are assigned a number between 0 to 1. These numbers are called the probabilities.
- The approach outlined by Kolmogorov ensured that the assignment of the probability numbers to the set members and all the subsets are such that there are no contradictions or inconsistencies



# Several Ways of Defining the Probability

## An Axiomatic Set-Theoretic Definition

- Examples of the atomic events and the universal sets
  - Roll a six-sided die:
    - ▷ Set of atomic events:  $S = \{1, 2, \dots, 6\}$
  - Toss a coin:
    - ▷ Set of atomic events:  $S = \{\text{Head}, \text{Tail}\}$
- Set  $S$  of all atomic events is a certain event. Its probability is 1.
- Set  $\emptyset$  is a null event. Its probability is zero.
- Any subset of the set  $S$  denotes a probable event. Its probability is a variable between 0 and 1.



# Several Ways of Defining the Probability

## An Axiomatic Set-Theoretic Definition

- An event  $A$  can be defined by its two complementary outputs:  $X = 1$  denotes that  $A$  has occurred and  $X = 0$  denotes the complementary event
  - Let  $A$  denote the even-numbered outcome when a die is rolled;
  - $X = 1$  denotes that either 0, 2, 4 or 6 have come up;
  - $X = 0$  denotes that 1, 3 or 5 are rolled
- Thus, for  $X = i$ ,  $i = 1$  denotes that the event  $A$  has occurred, and  $i = 0$  denotes the complementary event (i.e., that the event  $A$  has not occurred)



# Several Ways of Defining the Probability

## An Axiomatic Set-Theoretic Definition

- An event  $A$  can be defined by its two complementary outputs:  $X = 1$  denotes that  $A$  has occurred and  $X = 0$  denotes the complementary event
  - Let  $A$  denote the even-numbered outcome when a die is rolled;
  - $X = 1$  denotes that either 0, 2, 4 or 6 have come up;
  - $X = 0$  denotes that 1, 3 or 5 are rolled
- Thus, for  $X = i$ ,  $i = 1$  denotes that the event  $A$  has occurred, and  $i = 0$  denotes the complementary event (i.e., that the event  $A$  has not occurred)
  - The certain event  $S$  has a complementary event  $\emptyset$
- We will drop the event  $A$  of the actual world and just focus on the binary variable  $X$



## Fundamental Concepts

# Axioms of Probability

- Probability  $P(X = i)$  is a number which measures the likelihood of event  $X = i$ .
- Following are three axioms of probability:
  - ①  $P(X = i) \geq 0$  (i.e., no event has a probability less than zero).
  - ②  $P(X = i) \leq 1$ , and  $P(X = i) = 1$  only if  $(X = i)$  equals  $S$ , i.e., if  $X = i$  is a certain event.
  - ③ Let  $X = i$  and  $Y = j$  be two events (note that  $j$  is also either 0 or 1). The probability of the union of these events is given as
$$P((X = i) \cup (Y = j)) = P(X = i) + P(Y = j) - P((X = i) \cap (Y = j)).$$
- The entire theory of probability follows from these three axioms.
  - ▷ E.g., the probabilities of mutually exclusive events add; for some combination of  $i$  and  $j$ ,  $(X = i) \cap (Y = j) = \emptyset$ , i.e.,
$$P((X = i) \cap (Y = j)) = 0.$$



## Fundamental Concepts

# A Highly Useful Rule of Probability

## Marginalization of the Joint Probability

- Given the joint probability  $P((X = i) \cap (Y = j))$  (i.e., the probability of the joint event  $(X = i) \cap (Y = j)$ ) often abbreviated as  $P(X = i, Y = j)$ , it is possible to obtain the marginal probability  $P(X = i)$  by summing or integrating over both the possibilities of  $Y$ .

$$\rightarrow P(X = i) = \sum_{j=0}^1 P(X = i, Y_j)$$



# Joint Probability

- The joint probability  $P(X = i, Y = j) = P(X = i \cap Y = j)$  is the probability that both  $X = i$  and  $Y = j$  occur



# Joint Probability

- The joint probability  $P(X = i, Y = j) = P(X = i \cap Y = j)$  is the probability that both  $X = i$  and  $Y = j$  occur
  - Suppose  $X$  and  $Y$  are both Bernoulli sources, i.e., they take two possibilities each (either 0 or 1). Therefore, the joint event  $(X, Y)$  takes four possibilities



# Joint Probability

- The joint probability  $P(X = i, Y = j) = P(X = i \cap Y = j)$  is the probability that both  $X = i$  and  $Y = j$  occur
  - Suppose  $X$  and  $Y$  are both Bernoulli sources, i.e., they take two possibilities each (either 0 or 1). Therefore, the joint event  $(X, Y)$  takes four possibilities
  - For simplicity in the notations, the dependence on  $i, j$  is dropped; though it is to be noted that  $(X, Y)$  is an event with four possibilities



# Joint Probability

- The joint probability  $P(X = i, Y = j) = P(X = i \cap Y = j)$  is the probability that both  $X = i$  and  $Y = j$  occur
  - Suppose  $X$  and  $Y$  are both Bernoulli sources, i.e., they take two possibilities each (either 0 or 1). Therefore, the joint event  $(X, Y)$  takes four possibilities
  - For simplicity in the notations, the dependence on  $i, j$  is dropped; though it is to be noted that  $(X, Y)$  is an event with four possibilities



Fundamental Concepts

# An Example

- Let a student from your batch be picked at random and asked whether he/she wishes to be a part of DA-IICT Dance Club, Sports Club, both or none. Let



# An Example

- Let a student from your batch be picked at random and asked whether he/she wishes to be a part of DA-IICT Dance Club, Sports Club, both or none. Let
  - ▷  $X$  represent the preference regarding the Sports Club



# An Example

- Let a student from your batch be picked at random and asked whether he/she wishes to be a part of DA-IICT Dance Club, Sports Club, both or none. Let
  - ▷  $X$  represent the preference regarding the Sports Club
    - $X = 1 \Rightarrow$  the respondent wishes to take part in the Sports Club
    - $X = 0 \Rightarrow$  the respondent does not wish to be a sports person



## Fundamental Concepts

## An Example

- Let a student from your batch be picked at random and asked whether he/she wishes to be a part of DA-IICT Dance Club, Sports Club, both or none. Let
  - ▷  $X$  represent the preference regarding the Sports Club
    - $X = 1 \Rightarrow$  the respondent wishes to take part in the Sports Club
    - $X = 0 \Rightarrow$  the respondent does not wish to be a sports person
  - ▷  $Y$  represent the choice regarding the Dance Club



## Fundamental Concepts

## An Example

- Let a student from your batch be picked at random and asked whether he/she wishes to be a part of DA-IICT Dance Club, Sports Club, both or none. Let
  - ▷  $X$  represent the preference regarding the Sports Club
    - $X = 1 \Rightarrow$  the respondent wishes to take part in the Sports Club
    - $X = 0 \Rightarrow$  the respondent does not wish to be a sports person
  - ▷  $Y$  represent the choice regarding the Dance Club
    - $Y = 1 \Rightarrow$  the respondent wishes to take part in the Dance Club
    - $Y = 0 \Rightarrow$  the respondent does not wish to be part of the Dance person
- What are the four possibilities? How to express  $X$  and  $Y$  in terms of them?



# A Highly Useful Rule of Probability

## Marginalization of the Joint Probability

- Given the joint probability  $P((X = i) \cap (Y = j))$  (i.e., the probability of the joint event  $(X = i) \cap (Y = j)$ ) often abbreviated as  $P(X = i, Y = j)$ , it is possible to obtain the marginal probability  $P(X = i)$  by summing or integrating over both the possibilities of  $Y$ .

$$\rightarrow P(X = i) = \sum_{j=0}^1 P(\underline{X = i}, \underline{Y_j})$$



# Conditional Probability

- Conditional probability  $P(X | Y)$  is defined as the ratio  $\frac{P(X, Y)}{P(Y)}$



# Conditional Probability

- **Conditional probability**  $P(X | Y)$  is defined as the ratio  $\frac{P(X, Y)}{P(Y)}$
- It denotes the probability that  $X$  will occur given  $Y$  has occurred



# Statistical Independence versus Dependence

- **Statistical independence:** events  $X$  and  $Y$  are statistically independent if  $P(X, Y) = P(X) \times P(Y)$  (for all four possibilities when  $X$  and  $Y$  are Bernoulli RVs)
  - If  $X$  and  $Y$  are statistically independent,  $\underline{P(X | Y) = P(X)}$  and  $\underline{P(Y | X) = P(Y)}$



# Statistical Independence versus Dependence

- **Statistical independence:** events  $X$  and  $Y$  are statistically independent if  $P(X, Y) = P(X) \times P(Y)$  (for all four possibilities when  $X$  and  $Y$  are Bernoulli RVs)
  - If  $X$  and  $Y$  are statistically independent,  $P(X | Y) = P(X)$  and  $P(Y | X) = P(Y)$
- **Statistical Dependence:** the joint probabilities cannot be written as the products of the marginal probabilities; the conditional probabilities are not the same as the marginal probabilities



# Statistical Independence versus Dependence

- **Statistical independence:** events  $X$  and  $Y$  are statistically independent if  $P(X, Y) = P(X) \times P(Y)$  (for all four possibilities when  $X$  and  $Y$  are Bernoulli RVs)
  - If  $X$  and  $Y$  are statistically independent,  $P(X | Y) = P(X)$  and  $P(Y | X) = P(Y)$
- **Statistical Dependence:** the joint probabilities cannot be written as the products of the marginal probabilities; the conditional probabilities are not the same as the marginal probabilities
- Although the words and the math symbolism may seem abstract and unfamiliar, these are some of the concepts that all of us understand **very well!**



# Statistical Independence versus Dependence

- **Statistical independence:** events  $X$  and  $Y$  are statistically independent if  $P(X, Y) = P(X) \times P(Y)$  (for all four possibilities when  $X$  and  $Y$  are Bernoulli RVs)
  - If  $X$  and  $Y$  are statistically independent,  $P(X | Y) = P(X)$  and  $P(Y | X) = P(Y)$
- **Statistical Dependence:** the joint probabilities cannot be written as the products of the marginal probabilities; the conditional probabilities are not the same as the marginal probabilities
- Although the words and the math symbolism may seem abstract and unfamiliar, these are some of the concepts that all of us understand **very well!**
  - They relate to the notions of influence, causation, correlation, etc.



$X = \text{Teacher}$

$Y = \text{Student}$

0.99

$X = 0$  (Physics)

0.01

$X = 1$  (math)

0.5

0.9

0.05

0.05

0.9

0.05

0.3

0.05

$Y = 0$

$Y = 1$

$Y = 2$

$Y = 3$

student 1

## Fundamental Concepts

# Statistical Independence versus Dependence

- **Statistical independence:** events  $X$  and  $Y$  are statistically independent if  $P(X, Y) = P(X) \times P(Y)$  (for all four possibilities when  $X$  and  $Y$  are Bernoulli RVs)
  - If  $X$  and  $Y$  are statistically independent,  $P(X | Y) = P(X)$  and  $P(Y | X) = P(Y)$
- **Statistical Dependence:** the joint probabilities cannot be written as the products of the marginal probabilities; the conditional probabilities are not the same as the marginal probabilities
- Although the words and the math symbolism may seem abstract and unfamiliar, these are some of the concepts that all of us understand **very well!**
  - They relate to the notions of influence, causation, correlation, etc.
  - Probably the best way to come to grips with the concept of the joint and the conditional probabilities is through a framework in which Bayes' Theorem applies



# Bayes Rule

## The Causes and the Effects

- The Bayes' Rule deals with the events of the real world which is governed by underlying causes (that are often not observable) and the observed outcomes or the effects



# Bayes Rule

## The Causes and the Effects

- The Bayes' Rule deals with the events of the real world which is governed by underlying causes (that are often not observable) and the observed outcomes or the effects
- The Bayes' Rule strikes a chord in many ways possibly because it comes close to **us**, i.e., what we do all the time — our brains probably are programmed to implement the Bayes' Rule without our active awareness



# Bayes Rule

## The Causes and the Effects

- The Bayes' Rule deals with the events of the real world which is governed by underlying causes (that are often not observable) and the observed outcomes or the effects
- The Bayes' Rule strikes a chord in many ways possibly because it comes close to **us**, i.e., what we do all the time — our brains probably are programmed to implement the Bayes' Rule without our active awareness
- Due to the Bayes' Rule, the probability theory applies to, e.g., the intensely absorbing episodes of, say, Shakespeare's dramas!



# Bayes Rule

## A Visualization

### The World of Causes

- We never have a direct contact with this world

$$p(C = 0) = 1 - q$$

$$C = 0$$

$$p(C = 1) = q$$

$$C = 1$$

### The "Communication" Channel

- The invisible causes exert the influence on the real world through a "communication channel", which can be noisy, i.e., it can have  $p_{Error} > 0$

$$p_{Error,0}$$

$$p_{Error,1}$$

$$1 - p_{Error,0}$$

$$1 - p_{Error,1}$$

### The World of Effects

- We only can observe the effects or the outcomes influenced by the underlying causes

$$T = 0$$

$$p(T = 0)$$

$$T = 1$$

$$p(T = 1)$$

$p(C)$  is called the **prior** (defines the probability distribution of the causes)

$p_{Error,0}$  and  $p_{Error,1}$  together define the **likelihood** for a binary source and a binary sink. For a binary symmetric channel,  $p_{Error,0} = p_{Error,1}$

$p(T)$  is called the **evidence** (defines the probability distribution of the effects)



# Bayes Rule

## Terminology

- Let  $x_i$  be the cause and  $y_j$  be the effect.
  - $i \in \{0, 1\}$  and  $j \in \{0, 1\}$
  - i.e., the cause can be one of two possibilities, and the observed effect can be one of two possibilities
  - We will generalize the number of possibilities later



# Bayes Rule

## Terminology

- Let  $x_i$  be the cause and  $y_j$  be the effect.
  - $i \in \{0, 1\}$  and  $j \in \{0, 1\}$
  - i.e., the cause can be one of two possibilities, and the observed effect can be one of two possibilities
  - We will generalize the number of possibilities later
- It is often the case that the conditional probability  $p(y_j | x_i)$  is available.
  - This is called **the likelihood function** or **"forward probability"**



# Bayes Rule

## Terminology

- Let  $x_i$  be the cause and  $y_j$  be the effect.
  - $i \in \{0, 1\}$  and  $j \in \{0, 1\}$
  - i.e., the cause can be one of two possibilities, and the observed effect can be one of two possibilities
  - We will generalize the number of possibilities later
- It is often the case that the conditional probability  $p(y_j | x_i)$  is available.
  - This is called **the likelihood** function
- The probabilities of the causes  $p(x_i)$  may be available also — this is called **the prior**; and the probabilities of the effect  $p(y_j)$  can be computed — this is called **the evidence**



# Bayes Rule

## Terminology

- Let  $x_i$  be the cause and  $y_j$  be the effect.
  - $i \in \{0, 1\}$  and  $j \in \{0, 1\}$
  - i.e., the cause can be one of two possibilities, and the observed effect can be one of two possibilities
  - We will generalize the number of possibilities later
- It is often the case that the conditional probability  $p(y_j | x_i)$  is available.
  - This is called **the likelihood** function
- The probabilities of the causes  $p(x_i)$  may be available also — this is called **the prior**; and the probabilities of the effect  $p(y_j)$  can be computed — this is called **the evidence**
- We are oftentimes interested in knowing or *inferring*, given the observation  $y_j$  that we have, the most likely  $x_i$  that could have caused it. This is called the **posterior**
- The Bayes' Theorem relates the the likelihood, the prior and the evidence to the posterior



# Bayes Rule

The Math Remains the Same; the Problem Description Changes

- Cause  $x$ :
  - $x = 1$ : patient has (some) disease
  - $x = 0$ : patient does not have the disease
- Effect  $y$ :
  - $y = 1$ : the test that the patient takes is positive
  - $y = 0$ : test result is negative
- A person takes the test and  $y = 1$ , does she/he have the disease, i.e., is  $x = 1$ ?
  - If  $p$  and  $q$  are known, the posterior on  $x$  can be calculated using the Bayes' rule



# Toothache and Cavity

The Math Remains the Same; the Problem Description Changes

- Cause  $x$ :
  - $x = 1$ : your friend has cavity
  - $x = 0$ : the friend does not have the cavity
- Effect  $y$ :
  - $y = 1$ : your friend complains of toothache
  - $y = 0$ : she/he does not
- Your friend has the toothache  $y = 1$ , does she/he have the cavity, i.e., is  $x = 1$ ?
  - If  $p$  and  $q$  are known, the posterior on  $x$  can be calculated using the Bayes' rule



# Burglary and Car Alarm

The Math Remains the Same; the Problem Description Changes

- Cause  $x$ :
  - $x = 1$ : burglars are attempting to steal your car
  - $x = 0$ : the car is safe
- Effect  $y$ :
  - $y = 1$ : you hear the car alarm at 3 am
  - $y = 0$ : the car alarm is silent
- You hear the car alarm  $y = 1$ , are the burglars at the work, i.e., is  $x = 1$ ?
  - Again, the same math can solve the problem provided  $p$  and  $q$  are known



# A Day in the Court

The Math Remains the Same; the Problem Description Changes

- Suppose you're witness to a night-time hit and run accident involving a taxi in Gandhinagar. Of all taxis in Gandhinagar, 1% are light gray, and the rest are white colored. Extensive testing has earlier shown that under dim light conditions, discrimination between the light gray and white is 95% reliable.
- You swear, under oath, that the taxi you saw was gray. Should all the suspects that own the white taxi be released from the jail?



# Is It Raining Outside?

The Math Remains the Same; the Problem Description Changes

- Suppose you're a security guard at an underground nuclear facility. Only observation of the outside world you have is when you see the Director coming in. There is a 1% that it rains on any given day. You see the Director coming with an umbrella. When it is raining, there is a 95% chance that the Director has the umbrella with him.
- You see one day Director arriving with his umbrella. Is it raining outside?



# Who is the favorite player?

The Math Remains the Same; the Problem Description Changes

- Suppose who gets to play in the Indian cricket team is influenced to a large extent either by the captain of the team  $X = 0$  or the coach  $X = 1$ . A particular player  $Y = 1$  is favored highly (with 95% weight) by the captain, whereas the coach of the team is strongly against including this player in the team, i.e., the coach favors the outcome  $Y = 0$  with 95% weight.
- You see that the team declared on the day of the match has the player included in the team, i.e.,  $Y = 1$ . Your immediate reaction is that “oh, this happened only because of ...”.

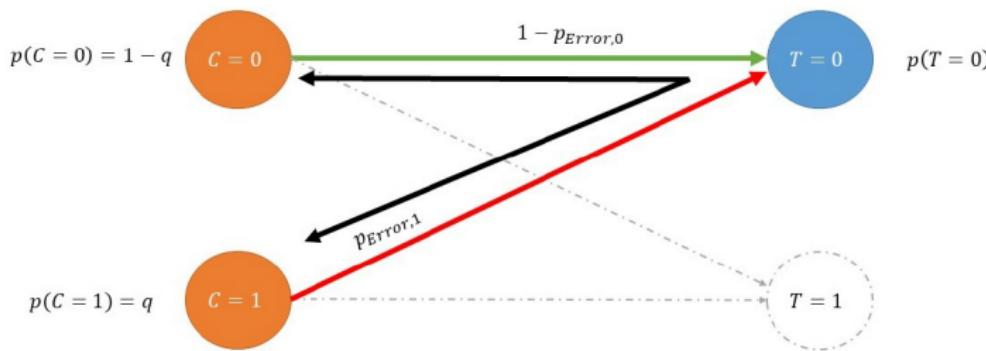


# Bayes Rule

## A Visualization

Upon the observation of an effect, say,  $T = 0$ , the task is to determine the likely cause that could have generated this observation, i.e., whether  $C = 1$  is more likely compared to  $C = 0$  given  $T = 0$ .

- This is called the **posterior** on  $C$  given  $T = 0$  and it is the conditional probability  $p(C | T = 0)$
- This requires the knowledge of  $p_{Error,0}$  and  $p_{Error,1}$  and the prior  $p(C)$

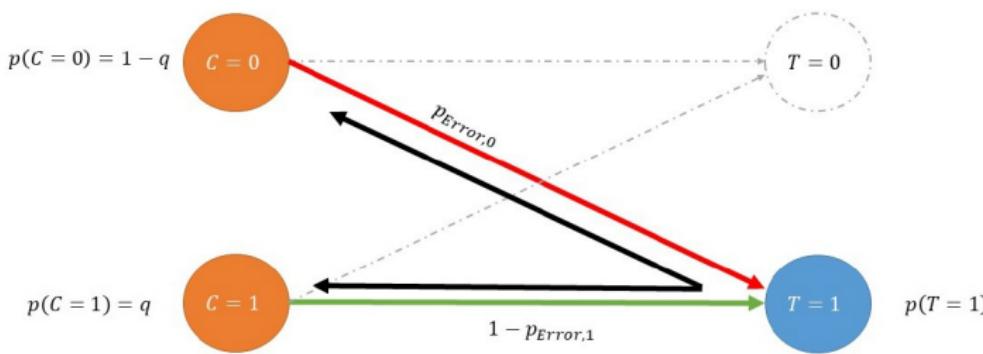


# Bayes Rule

## A Visualization

If  $T = 1$  is observed instead, the posterior on  $C$  given  $T = 1$ , i.e.,  $p(C | T = 1)$ , is required to be computed

- This case requires the knowledge of  $p_{Error,1}$  and  $p_{Error,0}$ , and the prior  $p(C)$



# Bayes Rule

## The Main Formula

- Since  $p(y_j, x_i) = p(y_j | x_i) \times p(x_i) = p(x_i | y_j) \times p(y_j)$ , we have:

### Theorem (Bayes')

$$p(x_i | y_j) = \frac{p(y_j | x_i) \times p(x_i)}{p(y_j)}$$

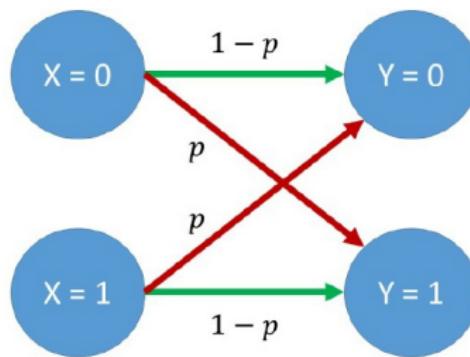
$$\text{Posterior} = \frac{\text{Likelihood} \times \text{Prior}}{\text{Evidence}}$$



# Binary Symmetric Channel or BSC

- The likelihood function  $p(Y = 1 | X = 0)$  is identical to  $p(Y = 0 | X = 1)$ . Both equal  $p$
- The prior on the Bernouilli variable  $X$  is denoted as  $q$

$$1 - q = 0.5$$



$$q = 0.5$$

- (1)  $P(Y)$
- (2)  $P(X|Y)$
- (3)  $H(X)$
- (4)  $H(X|Y)$
- (5)  $I(X;Y)$

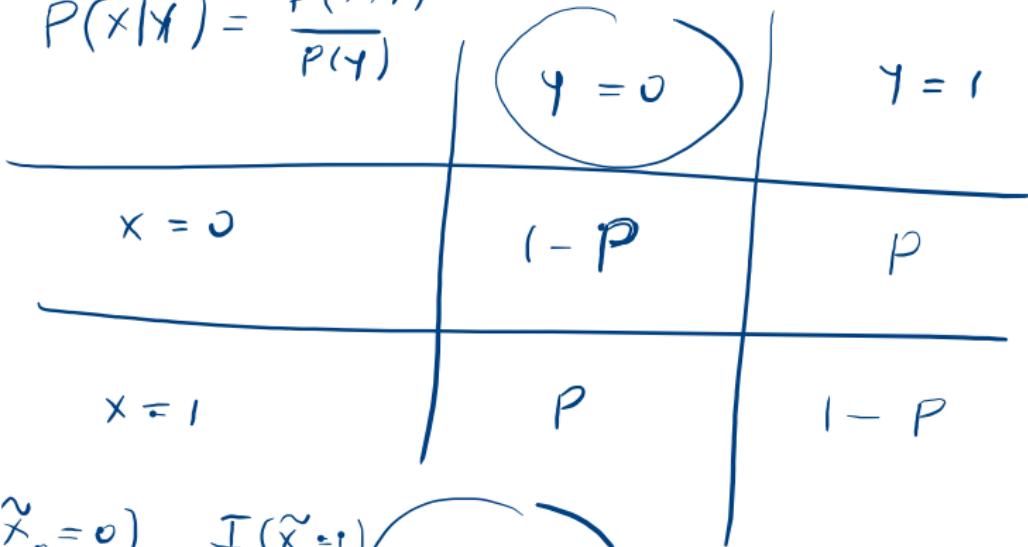


|                    |                          | $P(Y/x)$   |       |
|--------------------|--------------------------|------------|-------|
|                    |                          | $y=0$      | $y=1$ |
| $x=0$              | $0.5$                    | $1-P$      | $P$   |
|                    | $0.5$                    | $P$        | $1-P$ |
|                    | $\uparrow$               | $\uparrow$ |       |
| $P_{\text{prior}}$ | $Likelihood$<br>$P(Y/x)$ |            |       |

|       |       | $P(x,y)$   |            |
|-------|-------|------------|------------|
|       |       | $y=0$      | $y=1$      |
| $x=0$ | $x=0$ | $(1-P)0.5$ | $P(0.5)$   |
|       | $x=1$ | $P(0.5)$   | $(1-P)0.5$ |
|       |       | $0.5$      | $0.5$      |

$$P(x,y) = P(Y/x)P(x)$$

$$P(x|y) = \frac{P(x,y)}{P(y)}$$

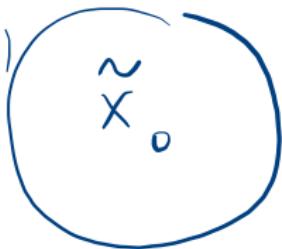


$$\mathbb{I}(\tilde{x}_0=0)$$

$$\mathbb{I}(\tilde{x}_0=1)$$

$$-\log_2(1-P)$$

$$-\log_2 P$$



$\tilde{x}_0$

$y=1$

$P$

$1 - P$

$\tilde{x}_0$

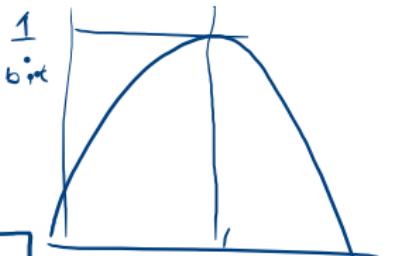
$x=1$

$P$

$$P(x|y) = \frac{P(x,y)}{P(y)}$$

(3)

$$H(X) = 1 \text{ bit}$$



$$H_2(x) = -x \log_2(x) - (1-x) \log_2(1-x)$$

Binary Entropy  
Function

(4)

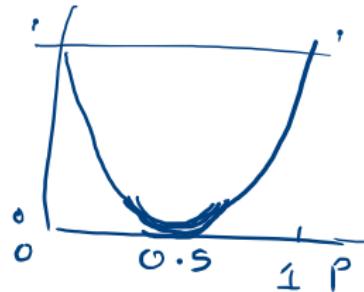
$$H(X|Y)$$

$$(4.1) H(\tilde{X}_0) = H(X|Y=0) = H_2(P)$$

$$(4.2) H(\tilde{X}_1) = H(X|Y=1) = H_2(1-P) = H_2(P)$$

$$\begin{aligned}
 \textcircled{4} \quad H(x|y) &= P(y=0) H(x|y=0) + P(y=1) H(x|y=1) \\
 &= 0.5 H_2(p) + 0.5 H_2(p) \\
 &= H_2(p) \quad I(x;y)
 \end{aligned}$$

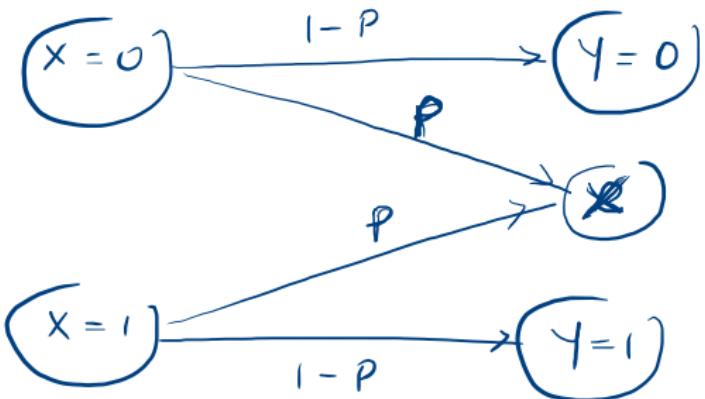
$$\begin{aligned}
 \textcircled{5} \quad I(x;y) &= H(x) - H(x|y) \\
 &= 1 - H_2(p)
 \end{aligned}$$



(BEC ( $p$ ))

# Binary Erasure channel

$$1-q = 0.5$$



$$q = 0.5$$

- (1)  $P(Y)$
- (2)  $P(X|Y)$
- (3)  $H(X)$
- (4)  $H(X|Y)$
- (5)  $I(X; Y)$

(Plot This function  
 $I(X; Y)$ )

# Bayes Rule

## The MAP Decision Rule

- Sometimes we don't care about the exact probabilities; all that we want to know, given  $Y = 1$ , is whether  $X = 1$  is more likely than  $X = 0$



# Bayes Rule

## The MAP Decision Rule

- Sometimes we don't care about the exact probabilities; all that we want to know, given  $Y = 1$ , is whether  $X = 1$  is more likely than  $X = 0$ 
  - ▷ E.g., the only allowed answers to the patient's query "do I have the disease?" may be "yes" or "no"; the former implies that the patient has to be hospitalized, the latter means the hospitalization is not needed. Note that the hospitalization is a binary decision, the patient cannot be partially hospitalized!



# Bayes Rule

## The MAP Decision Rule

- Sometimes we don't care about the exact probabilities; all that we want to know, given  $Y = 1$ , is whether  $X = 1$  is more likely than  $X = 0$ 
  - ▷ E.g., the only allowed answers to the patient's query "do I have the disease?" may be "yes" or "no"; the former implies that the patient has to be hospitalized, the latter means the hospitalization is not needed. Note that the hospitalization is a binary decision, the patient cannot be partially hospitalized!
- Given this question, the obvious strategy, when  $Y = 1$  is observed, is to say "yes" if  $p(X = 1 | Y = 1) > p(X = 0 | Y = 1)$  and say "no" otherwise.



# Bayes Rule

## The MAP Decision Rule

- Sometimes we don't care about the exact probabilities; all that we want to know, given  $Y = 1$ , is whether  $X = 1$  is more likely than  $X = 0$ 
  - ▷ E.g., the only allowed answers to the patient's query "do I have the disease?" may be "yes" or "no"; the former implies that the patient has to be hospitalized, the latter means the hospitalization is not needed. Note that the hospitalization is a binary decision, the patient cannot be partially hospitalized!
- Given this question, the obvious strategy, when  $Y = 1$  is observed, is to say "yes" if  $p(X = 1 | Y = 1) > p(X = 0 | Y = 1)$  and say "no" otherwise.
- This is called the Maximum A-Posteriori (MAP) decision rule



# Bayes Rule

## The MAP Decision Rule

- The MAP rule:



# Bayes Rule

## The MAP Decision Rule

- The MAP rule:

→ Given an observation  $Y_j$ , evaluate  $p(X_i | Y_j)$  for all  $i$ , and select that cause  $i$  for which  $p(X_i | Y_j)$  is the maximum



# Bayes Rule

## The MAP Decision Rule

- The MAP rule:

- Given an observation  $Y_j$ , evaluate  $p(X_i | Y_j)$  for all  $i$ , and select that cause  $i$  for which  $p(X_i | Y_j)$  is the maximum
- When  $X$  is binary, this amounts to evaluating the ratio of the posteriors when  $Y = 1$

$$\lambda_1 = \frac{p(X=1|Y=1)}{p(X=0|Y=1)} = \frac{p(Y=1|X=1)p(X=1)}{p(Y=0|X=1)p(X=0)}$$

Prior

*Likelihood function*

$$= \frac{(1-p)q}{p(1-q)}$$

This  $\lambda_1$  is called the “odds in favor of  $X = 1$ ”.



# Bayes Rule

## The MAP Decision Rule

- The MAP rule:

- Given an observation  $Y_j$ , evaluate  $p(X_i | Y_j)$  for all  $i$ , and select that cause  $i$  for which  $p(X_i | Y_j)$  is the maximum
- When  $X$  is binary, this amounts to evaluating the ratio of the posteriors when  $Y = 1$

$$\lambda_1 = \frac{p(X = 1 | Y = 1)}{p(X = 0 | Y = 1)} = \frac{p(Y = 1 | X = 1)p(X = 1)}{p(Y = 0 | X = 1)p(X = 0)}$$

Prior on  $X$

$$= \frac{(1 - p)q}{p(1 - q)}$$

This  $\lambda_1$  is called the “odds in favor of  $X = 1$ ”.

- If the odds in favor of  $X = 1$  exceeds 1,  $X = 1$  is more likely than  $X = 0$ .
- According to the MAP rule, the receiver decides  $\hat{X} = 1$ ; else it sets  $\hat{X} = 0$ . (here,  $\hat{X}$  is the decision formed by the receiver about the transmitted bit  $X$ )



# Bayes Rule

## The MAP and the ML Decision Rules

- The MAP rule simplifies to Maximum Likelihood (ML) rule when all the priors are equally likely.
  - Given an observation  $Y_j$ , select that  $i$  for which  $p(Y_j | X_i)$  is maximum
  - Note that Bayes rule is not explicitly invoked
  - When  $X$  is binary, the ML rule amounts to evaluating the ratio of the likelihoods:

$$\lambda_1 = \frac{1-p}{p}$$

Thus, the ratio of the posteriors becomes simply the likelihood ratio (LR):  $(1-p)/p$



# The MAP and the ML Rules for the BSC

The Ratio of the Posteriors when  $Y = 0$  is received

- The ratio of the posteriors when  $Y = 0$  is received:



# The MAP and the ML Rules for the BSC

The Ratio of the Posteriors when  $Y = 0$  is received

- The ratio of the posteriors when  $Y = 0$  is received:

$$\begin{aligned}\lambda_0 &= \frac{p(X = 0 | Y = 0)}{p(X = 1 | Y = 0)} = \frac{p(Y = 0 | X = 0)p(X = 0)}{p(Y = 0 | X = 1)p(X = 1)} \\ &= \frac{(1 - p)(1 - q)}{p q}\end{aligned}$$



# The MAP and the ML Rules for the BSC

The Ratio of the Posteriors when  $Y = 0$  is received

- The ratio of the posteriors when  $Y = 0$  is received:

$$\begin{aligned}\lambda_0 &= \frac{p(X = 0 | Y = 0)}{p(X = 1 | Y = 0)} = \frac{p(Y = 0 | X = 0)p(X = 0)}{p(Y = 0 | X = 1)p(X = 1)} \\ &= \frac{(1 - p)(1 - q)}{p q}\end{aligned}$$

- The MAP receiver decides  $\hat{X} = 0$  if  $\lambda_0 \geq 1$ ; else it decides  $\hat{X} = 1$



# The MAP and the ML Rules for the BSC

The Ratio of the Posteriors when  $Y = 0$  is received

- The ratio of the posteriors when  $Y = 0$  is received:

$$\begin{aligned}\lambda_0 &= \frac{p(X = 0 | Y = 0)}{p(X = 1 | Y = 0)} = \frac{p(Y = 0 | X = 0)p(X = 0)}{p(Y = 0 | X = 1)p(X = 1)} \\ &= \frac{(1 - p)(1 - q)}{p q}\end{aligned}$$

- The MAP receiver decides  $\hat{X} = 0$  if  $\lambda_0 \geq 1$ ; else it decides  $\hat{X} = 1$
- If  $q = 0.5$ , the above degenerates to the ML rule:



# The MAP and the ML Rules for the BSC

The Ratio of the Posteriors when  $Y = 0$  is received

- The ratio of the posteriors when  $Y = 0$  is received:

$$\begin{aligned}\lambda_0 &= \frac{p(X = 0 | Y = 0)}{p(X = 1 | Y = 0)} = \frac{p(Y = 0 | X = 0)p(X = 0)}{p(Y = 0 | X = 1)p(X = 1)} \\ &= \frac{(1 - p)(1 - q)}{p q}\end{aligned}$$

- The MAP receiver decides  $\hat{X} = 0$  if  $\lambda_0 \geq 1$ ; else it decides  $\hat{X} = 1$
- If  $q = 0.5$ , the above degenerates to the ML rule:

$$\lambda_0 = \frac{1 - p}{p}$$



# The MAP and the ML Rules for the BSC

The Ratio of the Posteriors when  $Y = 0$  is received

- The ratio of the posteriors when  $Y = 0$  is received:

$$\begin{aligned}\lambda_0 &= \frac{p(X = 0 | Y = 0)}{p(X = 1 | Y = 0)} = \frac{p(Y = 0 | X = 0)p(X = 0)}{p(Y = 0 | X = 1)p(X = 1)} \\ &= \frac{(1 - p)(1 - q)}{p q}\end{aligned}$$

- The MAP receiver decides  $\hat{X} = 0$  if  $\lambda_0 \geq 1$ ; else it decides  $\hat{X} = 1$
- If  $q = 0.5$ , the above degenerates to the ML rule:

$$\lambda_0 = \frac{1 - p}{p}$$



# Observations

- In typical cases, the source (i.e., its PMF) appears vastly different depending on whether  $Y = 0$  is received versus  $Y = 1$  is received



# Observations

- In typical cases, the source (i.e., its PMF) appears vastly different depending on whether  $Y = 0$  is received versus  $Y = 1$  is received
  - ▷ If  $Y = 1$  is received,  $\lambda_1$  typically exceeds 1, i.e., the source PMF becomes *biased* toward the possibility that  $X = 1$



# Observations

- In typical cases, the source (i.e., its PMF) appears vastly different depending on whether  $Y = 0$  is received versus  $Y = 1$  is received
  - ▷ If  $Y = 1$  is received,  $\lambda_1$  typically exceeds 1, i.e., the source PMF becomes *biased* toward the possibility that  $X = 1$
  - ▷ Similarly, if  $Y = 0$  is received,  $\lambda_0$  typically exceeds 1, i.e., the source PMF favors the possibility that  $X = 0$



# Observations

- In typical cases, the source (i.e., its PMF) appears vastly different depending on whether  $Y = 0$  is received versus  $Y = 1$  is received
  - ▷ If  $Y = 1$  is received,  $\lambda_1$  typically exceeds 1, i.e., the source PMF becomes *biased* toward the possibility that  $X = 1$
  - ▷ Similarly, if  $Y = 0$  is received,  $\lambda_0$  typically exceeds 1, i.e., the source PMF favors the possibility that  $X = 0$
  - ▷ Note that the receiver never compares  $\lambda_1$  and  $\lambda_0$  with each other — after a given reception, the receiver either calculates  $\lambda_1$  (if  $Y = 1$  is received) or  $\lambda_0$  (if  $Y = 0$ ). The receiver compares either  $\lambda_1$  or  $\lambda_0$  against 1



# Observations

- In typical cases, the source (i.e., its PMF) appears vastly different depending on whether  $Y = 0$  is received versus  $Y = 1$  is received
  - ▷ If  $Y = 1$  is received,  $\lambda_1$  typically exceeds 1, i.e., the source PMF becomes *biased* toward the possibility that  $X = 1$
  - ▷ Similarly, if  $Y = 0$  is received,  $\lambda_0$  typically exceeds 1, i.e., the source PMF favors the possibility that  $X = 0$
  - ▷ Note that the receiver never compares  $\lambda_1$  and  $\lambda_0$  with each other — after a given reception, the receiver either calculates  $\lambda_1$  (if  $Y = 1$  is received) or  $\lambda_0$  (if  $Y = 0$ ). The receiver compares either  $\lambda_1$  or  $\lambda_0$  against 1
- We have here assumed that the receiver knows  $p$  and  $q$ . Is this a fair assumption? Maybe. In many cases, the exact knowledge of  $p$  and  $q$  may not be available and so they may be estimated



# Observations

- In typical cases, the source (i.e., its PMF) appears vastly different depending on whether  $Y = 0$  is received versus  $Y = 1$  is received
  - ▷ If  $Y = 1$  is received,  $\lambda_1$  typically exceeds 1, i.e., the source PMF becomes *biased* toward the possibility that  $X = 1$
  - ▷ Similarly, if  $Y = 0$  is received,  $\lambda_0$  typically exceeds 1, i.e., the source PMF favors the possibility that  $X = 0$
  - ▷ Note that the receiver never compares  $\lambda_1$  and  $\lambda_0$  with each other — after a given reception, the receiver either calculates  $\lambda_1$  (if  $Y = 1$  is received) or  $\lambda_0$  (if  $Y = 0$ ). The receiver compares either  $\lambda_1$  or  $\lambda_0$  against 1
- We have here assumed that the receiver knows  $p$  and  $q$ . Is this a fair assumption? Maybe. In many cases, the exact knowledge of  $p$  and  $q$  may not be available and so they may be estimated
- The values of  $\lambda_1$  or  $\lambda_0$  indicate the belief — a value of 1 implies both the causes are equally likely, values either much greater or much less than 1 imply a strong belief about the likely cause



# Observations

- In typical cases, the source (i.e., its PMF) appears vastly different depending on whether  $Y = 0$  is received versus  $Y = 1$  is received
  - ▷ If  $Y = 1$  is received,  $\lambda_1$  typically exceeds 1, i.e., the source PMF becomes *biased* toward the possibility that  $X = 1$
  - ▷ Similarly, if  $Y = 0$  is received,  $\lambda_0$  typically exceeds 1, i.e., the source PMF favors the possibility that  $X = 0$
  - ▷ Note that the receiver never compares  $\lambda_1$  and  $\lambda_0$  with each other — after a given reception, the receiver either calculates  $\lambda_1$  (if  $Y = 1$  is received) or  $\lambda_0$  (if  $Y = 0$ ). The receiver compares either  $\lambda_1$  or  $\lambda_0$  against 1
- We have here assumed that the receiver knows  $p$  and  $q$ . Is this a fair assumption? Maybe. In many cases, the exact knowledge of  $p$  and  $q$  may not be available and so they may be estimated
- The values of  $\lambda_1$  or  $\lambda_0$  indicate the belief — a value of 1 implies both the causes are equally likely, values either much greater or much less than 1 imply a strong belief about the likely cause
- The prior acts as a “dead weight” that prevents these beliefs to move by much



# An Example Calculation

The Ratio of the Posteriors = LR × the Ratio of the Priors

- For the C-disease problem,



# An Example Calculation

The Ratio of the Posteriors = LR  $\times$  the Ratio of the Priors

- For the C-disease problem,  
 $\rightarrow p = 0.05$  and  $q = 0.01$



# An Example Calculation

The Ratio of the Posteriors = LR × the Ratio of the Priors

- For the C-disease problem,  
→  $p = 0.05$  and  $q = 0.01$
- Therefore, if  $Y = 1$ , the posterior ratio

$$\lambda_1 = \frac{(1-p)q}{p(1-q)} \approx \frac{q}{p} = 0.01/0.05 = 0.2$$



# An Example Calculation

The Ratio of the Posteriors = LR  $\times$  the Ratio of the Priors

- For the C-disease problem,  
 $\rightarrow p = 0.05$  and  $q = 0.01$
- Therefore, if  $Y = 1$ , the posterior ratio

$$\lambda_1 = \frac{(1-p)q}{p(1-q)} \approx \frac{q}{p} = 0.01/0.05 = 0.2$$

- Let  $p(X = 1 | Y = 1)$  be denote as  $q'$ . Thus,

$$\lambda_1 = \frac{q'}{1 - q'} = 0.2 \Rightarrow q' = 0.16$$

- The person detected positive is quite **unlikely** to have the C-disease



# An Example Calculation

When the Ratio of the Priors is Unity, The Ratio of the Posteriors = LR

- For the C-disease problem, suppose



# An Example Calculation

When the Ratio of the Priors is Unity, The Ratio of the Posteriors = LR

- For the C-disease problem, suppose  
→  $p = 0.05$  and  $q = 0.5$



# An Example Calculation

When the Ratio of the Priors is Unity, The Ratio of the Posteriors = LR

- For the C-disease problem, suppose  
→  $p = 0.05$  and  $q = 0.5$
- Now, the posterior ratio  $\lambda_1 = q'/(1 - q')$  becomes identical to the likelihood ratio  $(1 - p)/p = 0.95/0.05 = 19$



# An Example Calculation

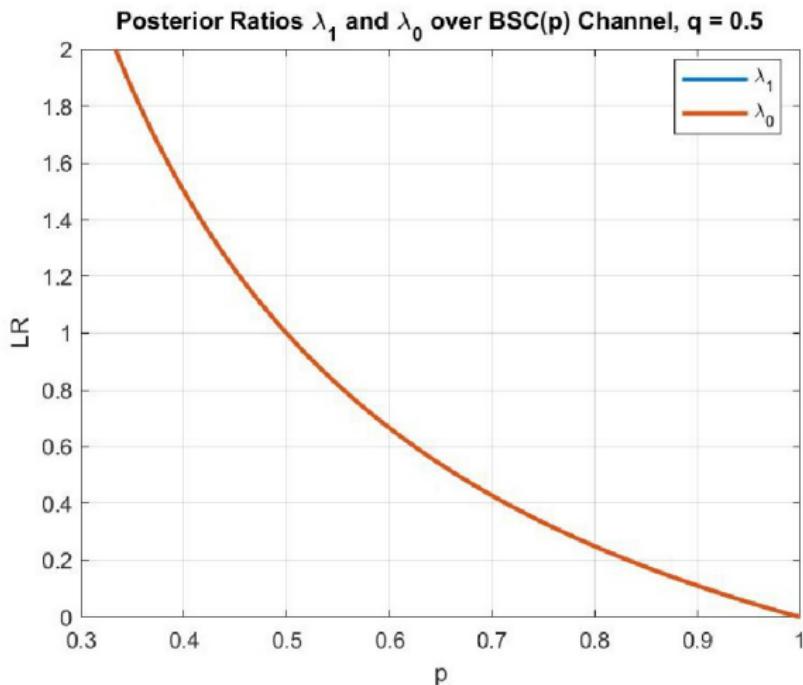
When the Ratio of the Priors is Unity, The Ratio of the Posteriors = LR

- For the C-disease problem, suppose  
→  $p = 0.05$  and  $q = 0.5$
- Now, the posterior ratio  $\lambda_1 = q'/(1 - q')$  becomes identical to the likelihood ratio  $(1 - p)/p = 0.95/0.05 = 19$
- This shows that  $q' = 0.95$ . The person detected positive is now **highly likely** to have the C-disease



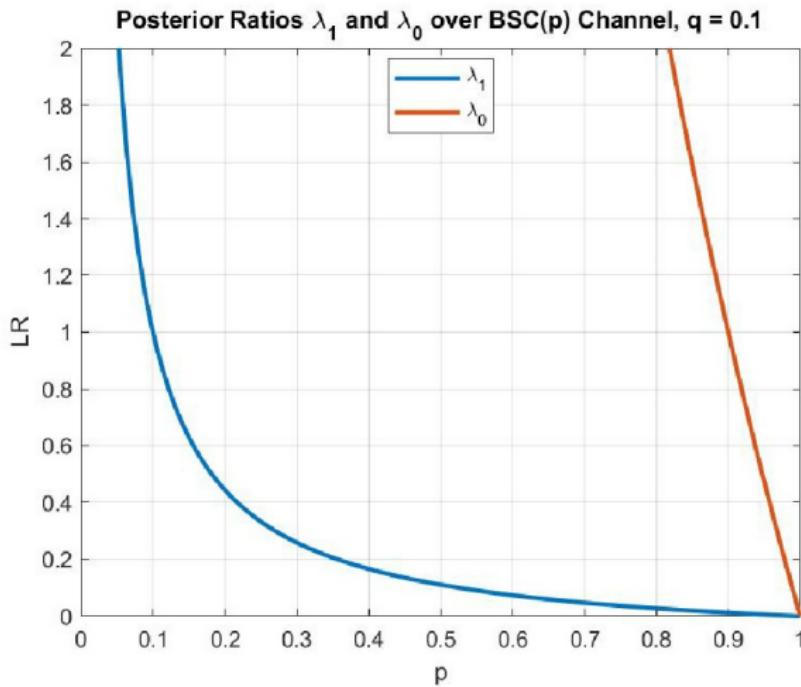
# Posterior Ratios $\lambda_1 = \lambda_0 = \text{Likelihood Ratio } (1-p)/p$

over BSC( $p$ ) Channel,  $q = 0.5$



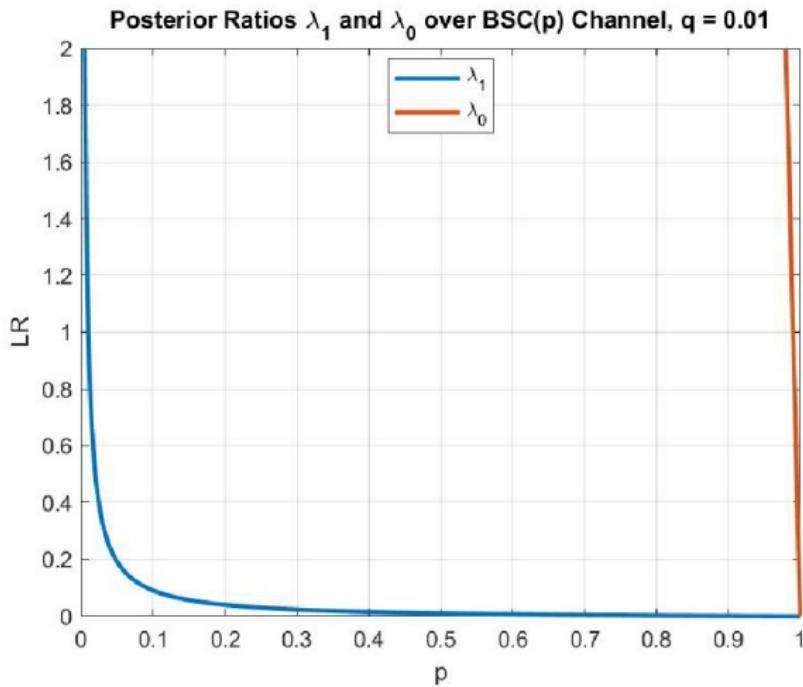
# Posterior Ratios $\lambda_1$ , and $\lambda_0$

over BSC( $p$ ) Channel,  $q = 0.1$



# Posterior Ratios $\lambda_1$ , and $\lambda_0$

over BSC( $p$ ) Channel,  $q = 0.01$



## Repeated Tests

- Suppose your friend tests positive for C-disease. For further confirmation, she decides to take another test in which she again tests positive.



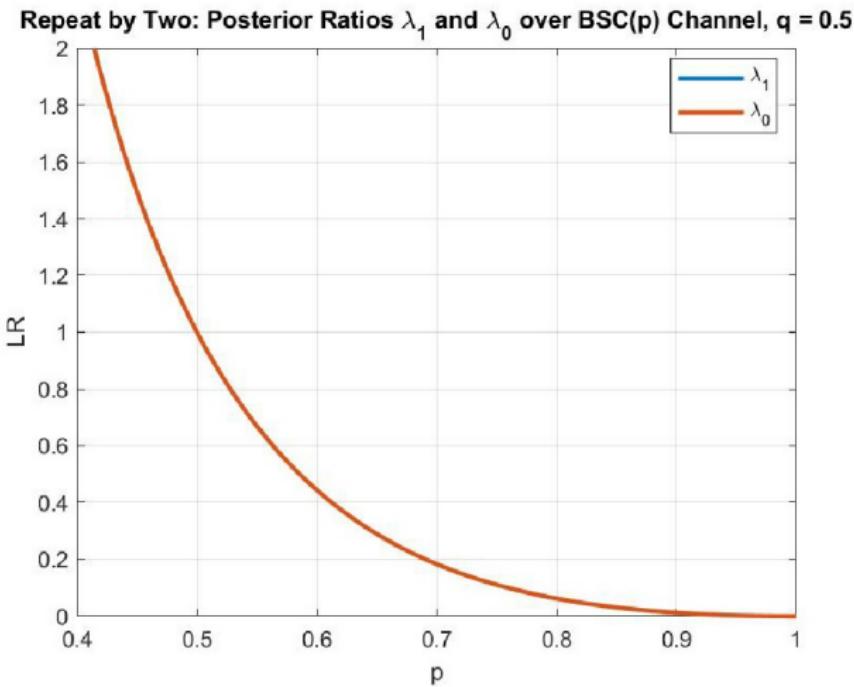
## Repeated Tests

- Suppose your friend tests positive for C-disease. For further confirmation, she decides to take another test in which she again tests positive.
- What is the probability that she has the disease? (Take  $p = 0.05$ . Use  $q = 0.5, 0.1$ , and  $0.01$ ).



# Posterior Ratios $\lambda_1 = \lambda_0$

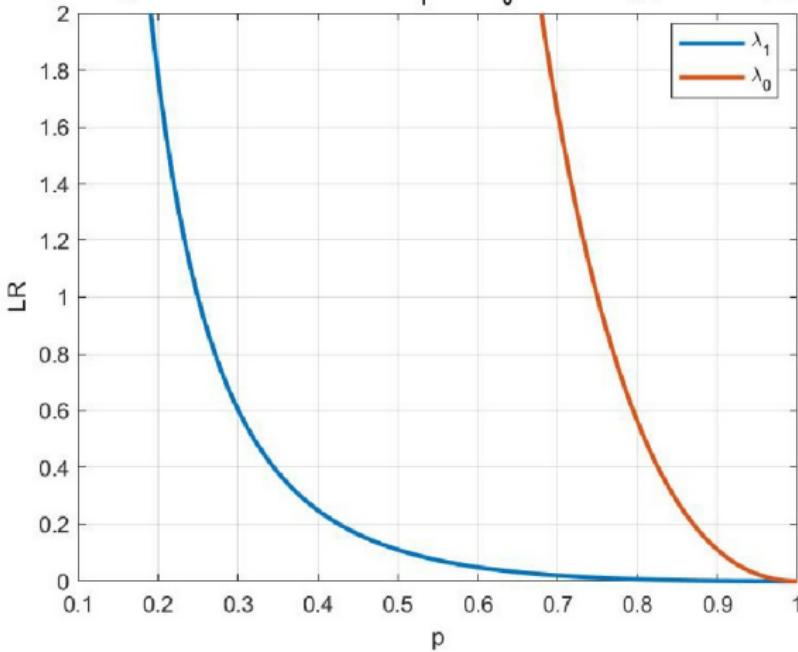
over BSC( $p$ ) Channel,  $q = 0.5$ . Repeated Transmission.



# Posterior Ratios $\lambda_1$ , and $\lambda_0$

over  $\text{BSC}(p)$  Channel,  $q = 0.1$ . Repeated Transmission.

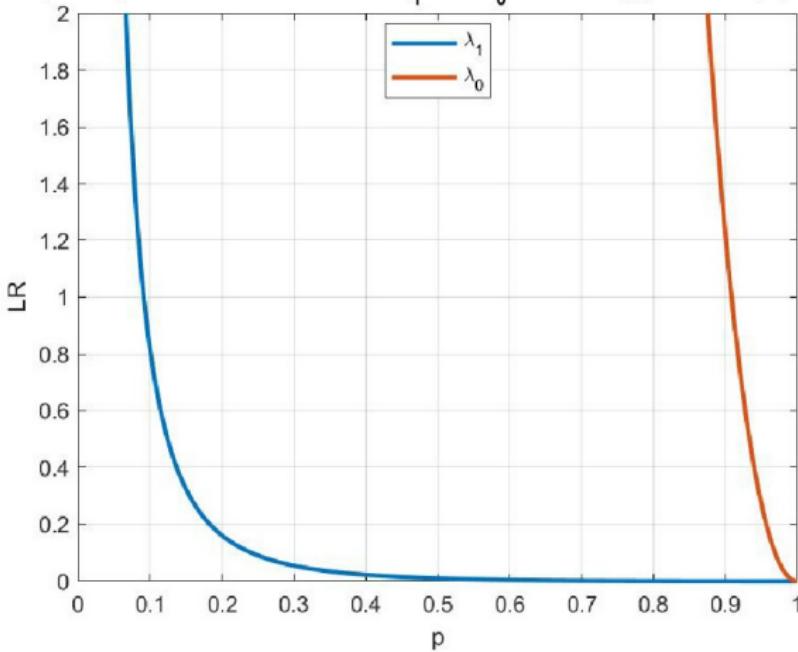
Repeat by Two: Posterior Ratios  $\lambda_1$  and  $\lambda_0$  over  $\text{BSC}(p)$  Channel,  $q = 0.1$



# Posterior Ratios $\lambda_1$ , and $\lambda_0$

over  $\text{BSC}(p)$  Channel,  $q = 0.01$ . Repeated Transmission.

Repeat by Two: Posterior Ratios  $\lambda_1$  and  $\lambda_0$  over  $\text{BSC}(p)$  Channel,  $q = 0.01$



# The Conditional Entropy

- We have earlier seen that the Entropy  $H(X)$  of the source can be thought of as its **temperature**
- Bayesian framework tells us that we never observe a source  $X$  directly; we observe it through the evidence  $Y$  we obtain



# The Conditional Entropy

- We have earlier seen that the Entropy  $H(X)$  of the source can be thought of as its **temperature**
- Bayesian framework tells us that we never observe a source  $X$  directly; we observe it through the evidence  $Y$  we obtain
- A question: what is the Entropy of the source after we obtain the evidence  $Y$ ?



# The Conditional Entropy

- We have earlier seen that the Entropy  $H(X)$  of the source can be thought of as its **temperature**
- Bayesian framework tells us that we never observe a source  $X$  directly; we observe it through the evidence  $Y$  we obtain
- A question: what is the Entropy of the source after we obtain the evidence  $Y$ ?
  - ▷ This is called the Conditional Entropy  $H(X | Y)$ , called the entropy of  $X$  given  $Y$ .



# The Conditional Entropy

- We have earlier seen that the Entropy  $H(X)$  of the source can be thought of as its **temperature**
- Bayesian framework tells us that we never observe a source  $X$  directly; we observe it through the evidence  $Y$  we obtain
- A question: what is the Entropy of the source after we obtain the evidence  $Y$ ?
  - ▷ This is called the Conditional Entropy  $H(X | Y)$ , called the entropy of  $X$  given  $Y$ .
- Hopefully, the source temperature, after we receive the information-bearing message  $Y$ , should reduce. If not, the information transfer cannot be said to have occurred



# The Conditional Entropy

- The conditional entropy when the observed effect  $Y$  takes a particular value  $j$  is given as

$$H(X | Y = j) = \mathbb{E} [I_k | Y = j] = - \sum_{k=1}^K p(X_k | Y = j) \log_2 (p(X_k | Y = j))$$



# The Conditional Entropy

- The conditional entropy when the observed effect  $Y$  takes a particular value  $j$  is given as

$$H(X | Y = j) = \mathbb{E}[I_k | Y = j] = - \sum_{k=1}^K p(X_k | Y = j) \log_2(p(X_k | Y = j))$$

- The conditional entropy is obtained by averaging over all realizations of the variable  $Y$ , i.e., as

$$H(X | Y) = \mathbb{E}[\mathbb{E}[I_k | Y = j]]$$

$$= - \sum_{j=1}^J p(Y = j) \sum_{k=1}^K p(X_k | Y = j) \log_2(p(X_k | Y = j))$$



# The Information Transfer

- Having defined the entropies  $H(X)$  and  $H(X | Y)$  of the source *before* and *after* the reception of the message  $Y$ , we can now define the information that is transferred to the receiver after it receives the message



# The Information Transfer

- Having defined the entropies  $H(X)$  and  $H(X | Y)$  of the source *before* and *after* the reception of the message  $Y$ , we can now define the information that is transferred to the receiver after it receives the message
  - In our analogy, the information transferred equals the amount by which the *temperature* of the source **reduces** because it has received the message



# The Information Transfer

- Having defined the entropies  $H(X)$  and  $H(X | Y)$  of the source *before* and *after* the reception of the message  $Y$ , we can now define the information that is transferred to the receiver after it receives the message
  - In our analogy, the information transferred equals the amount by which the *temperature* of the source **reduces** because it has received the message
- The information transfer achieved is defined as the difference between  $H(X)$  and  $H(X | Y)$ ; specifically,

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

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



## How can the bits take fractional values?

- We have seen that the concepts of information is obtained in terms of fraction of bits. However, how can the bits be in fractional units?



## How can the bits take fractional values?

- We have seen that the concepts of information is obtained in terms of fraction of bits. However, how can the bits be in fractional units?
- Shannon's answer: one bit at a time cannot be in fraction. However, the best way to think about the information generation and its transfer is to consider a large sequence of bits

