

# Information theory

# What is information?

- Telecommunications is the transmission of ***information*** over distance to communicate
- So, what is information? How can we define it and measure it for engineering purposes?



Simple example: “SHE-SELLS-SEA-SHELLS”

How much information does it contain?

How many bits do we need to store that information?

# ASCII code

Dec	Hx	Oct	Char	Dec	Hx	Oct	Html	Chr	Dec	Hx	Oct	Html	Chr	Dec	Hx	Oct	Html	Chr
0	0	000	<b>NUL</b> (null)	32	20	040	&#32;	<b>Space</b>	64	40	100	&#64;	<b>@</b>	96	60	140	&#96;	<b>`</b>
1	1	001	<b>SOH</b> (start of heading)	33	21	041	&#33;	<b>!</b>	65	41	101	&#65;	<b>A</b>	97	61	141	&#97;	<b>a</b>
2	2	002	<b>STX</b> (start of text)	34	22	042	&#34;	<b>"</b>	66	42	102	&#66;	<b>B</b>	98	62	142	&#98;	<b>b</b>
3	3	003	<b>ETX</b> (end of text)	35	23	043	&#35;	<b>#</b>	67	43	103	&#67;	<b>C</b>	99	63	143	&#99;	<b>c</b>
4	4	004	<b>EOT</b> (end of transmission)	36	24	044	&#36;	<b>\$</b>	68	44	104	&#68;	<b>D</b>	100	64	144	&#100;	<b>d</b>
5	5	005	<b>ENQ</b> (enquiry)	37	25	045	&#37;	<b>%</b>	69	45	105	&#69;	<b>E</b>	101	65	145	&#101;	<b>e</b>
6	6	006	<b>ACK</b> (acknowledge)	38	26	046	&#38;	<b>&amp;</b>	70	46	106	&#70;	<b>F</b>	102	66	146	&#102;	<b>f</b>
7	7	007	<b>BEL</b> (bell)	39	27	047	&#39;	<b>'</b>	71	47	107	&#71;	<b>G</b>	103	67	147	&#103;	<b>g</b>
8	8	010	<b>BS</b> (backspace)	40	28	050	&#40;	<b>(</b>	72	48	110	&#72;	<b>H</b>	104	68	150	&#104;	<b>h</b>
9	9	011	<b>TAB</b> (horizontal tab)	41	29	051	&#41;	<b>)</b>	73	49	111	&#73;	<b>I</b>	105	69	151	&#105;	<b>i</b>
10	A	012	<b>LF</b> (NL line feed, new line)	42	2A	052	&#42;	<b>*</b>	74	4A	112	&#74;	<b>J</b>	106	6A	152	&#106;	<b>j</b>
11	B	013	<b>VT</b> (vertical tab)	43	2B	053	&#43;	<b>+</b>	75	4B	113	&#75;	<b>K</b>	107	6B	153	&#107;	<b>k</b>
12	C	014	<b>FF</b> (NP form feed, new page)	44	2C	054	&#44;	<b>,</b>	76	4C	114	&#76;	<b>L</b>	108	6C	154	&#108;	<b>l</b>
13	D	015	<b>CR</b> (carriage return)	45	2D	055	&#45;	<b>-</b>	77	4D	115	&#77;	<b>M</b>	109	6D	155	&#109;	<b>m</b>
14	E	016	<b>SO</b> (shift out)	46	2E	056	&#46;	<b>.</b>	78	4E	116	&#78;	<b>N</b>	110	6E	156	&#110;	<b>n</b>
15	F	017	<b>SI</b> (shift in)	47	2F	057	&#47;	<b>/</b>	79	4F	117	&#79;	<b>O</b>	111	6F	157	&#111;	<b>o</b>
16	10	020	<b>DLE</b> (data link escape)	48	30	060	&#48;	<b>0</b>	80	50	120	&#80;	<b>P</b>	112	70	160	&#112;	<b>p</b>
17	11	021	<b>DC1</b> (device control 1)	49	31	061	&#49;	<b>1</b>	81	51	121	&#81;	<b>Q</b>	113	71	161	&#113;	<b>q</b>
18	12	022	<b>DC2</b> (device control 2)	50	32	062	&#50;	<b>2</b>	82	52	122	&#82;	<b>R</b>	114	72	162	&#114;	<b>r</b>
19	13	023	<b>DC3</b> (device control 3)	51	33	063	&#51;	<b>3</b>	83	53	123	&#83;	<b>S</b>	115	73	163	&#115;	<b>s</b>
20	14	024	<b>DC4</b> (device control 4)	52	34	064	&#52;	<b>4</b>	84	54	124	&#84;	<b>T</b>	116	74	164	&#116;	<b>t</b>
21	15	025	<b>NAK</b> (negative acknowledge)	53	35	065	&#53;	<b>5</b>	85	55	125	&#85;	<b>U</b>	117	75	165	&#117;	<b>u</b>
22	16	026	<b>SYN</b> (synchronous idle)	54	36	066	&#54;	<b>6</b>	86	56	126	&#86;	<b>V</b>	118	76	166	&#118;	<b>v</b>
23	17	027	<b>ETB</b> (end of trans. block)	55	37	067	&#55;	<b>7</b>	87	57	127	&#87;	<b>W</b>	119	77	167	&#119;	<b>w</b>
24	18	030	<b>CAN</b> (cancel)	56	38	070	&#56;	<b>8</b>	88	58	130	&#88;	<b>X</b>	120	78	170	&#120;	<b>x</b>
25	19	031	<b>EM</b> (end of medium)	57	39	071	&#57;	<b>9</b>	89	59	131	&#89;	<b>Y</b>	121	79	171	&#121;	<b>y</b>
26	1A	032	<b>SUB</b> (substitute)	58	3A	072	&#58;	<b>:</b>	90	5A	132	&#90;	<b>Z</b>	122	7A	172	&#122;	<b>z</b>
27	1B	033	<b>ESC</b> (escape)	59	3B	073	&#59;	<b>;</b>	91	5B	133	&#91;	<b>[</b>	123	7B	173	&#123;	<b>{</b>
28	1C	034	<b>FS</b> (file separator)	60	3C	074	&#60;	<b>&lt;</b>	92	5C	134	&#92;	<b>\</b>	124	7C	174	&#124;	<b> </b>
29	1D	035	<b>GS</b> (group separator)	61	3D	075	&#61;	<b>=</b>	93	5D	135	&#93;	<b>]</b>	125	7D	175	&#125;	<b>}</b>
30	1E	036	<b>RS</b> (record separator)	62	3E	076	&#62;	<b>&gt;</b>	94	5E	136	&#94;	<b>^</b>	126	7E	176	&#126;	<b>~</b>
31	1F	037	<b>US</b> (unit separator)	63	3F	077	&#63;	<b>?</b>	95	5F	137	&#95;	<b>_</b>	127	7F	177	&#127;	<b>DEL</b>

# Example

- How many bits do we need to encode “SHE-SELLS-SEA-SHELLS”
  - ASCII coding gives 8 bits per character, i.e. 1 byte
  - 20 characters = 20 bytes =  $20 \times 8 = 160$  bits
- Can we say the information carried by that string is 160 bits?

# Compressing information

- Unfortunately is not as easy, for example I could compress that string:
- Define new code based on frequency of letters:

Character	Occurrences	Relative frequencies	Coding
A	1	0.05	0000
H	2	0.10	0001
-	3	0.15	001
E	4	0.20	10
L	4	0.20	11
S	6	0.30	01

i.e. we use shorter codewords for more frequent characters

➔ We can use 49 bits instead of 160...

# Why is it interesting?

- This theory is the foundation for:
  - All modern digital communication (obviously)
- But also:
  - Cryptography and cryptanalysis (that's what won World War II)
  - Data compression
  - It has implication in many other fields: physics, linguistics, neurobiology

# Definition of information

- So, how can we define and quantify information? – this is a hard question!!
- Refer to Claude Shannon (1916-2001), the father of information theory.

Information is a measure of a reduction of the entropy of a random variable

entropy is a measure of the uncertainty associated with a random variable



# What???

- Weaver gives an explanation of Shannon's information:

Information is a measure of one's freedom of choice in selecting a message. The greater this freedom of choice, the greater the information, the greater is the uncertainty that the message actually selected is some particular one. Greater freedom of choice, greater uncertainty greater information go hand in hand.<sup>1</sup>

- Notice that Shannon does not consider the information carried out by the meaning of a message:

“...Frequently the messages have meaning; that is they refer to or are correlated according to some system with certain physical or conceptual entities. These semantic aspects of communication are irrelevant to the engineering problem. The significant aspect is that the actual message is one selected from a set of possible messages...”

<sup>1</sup> C.E.Shannon, W.Weaver, “The Mathematical Theory of Communication”



# Shannon entropy

- *Entropy is a measure of the uncertainty associated with a random variable*
- A random variable is used to describe a message in a communication system, it indicates the choice of a message over every possible message



Alice

If Alice sends a message to Bob, she could choose any expression in the English language

Before receiving the message Bob doesn't know what the message is about, it could be anything. At this stage a potential message can be described with a random variable.



Bob

# Random variable

- A random variable is a variable whose value is unknown.
- A random variable will follow a certain probability distribution.
- A discrete random variable will have a discrete number of outcomes:
  - flipping a coin has 2 outcomes
  - throwing one die has 6 outcomes
- A continuous random variable will have an infinite number of possible outcomes within a given range:
  - Computer time required to process a certain program
  - The amount of rain that falls in a certain location

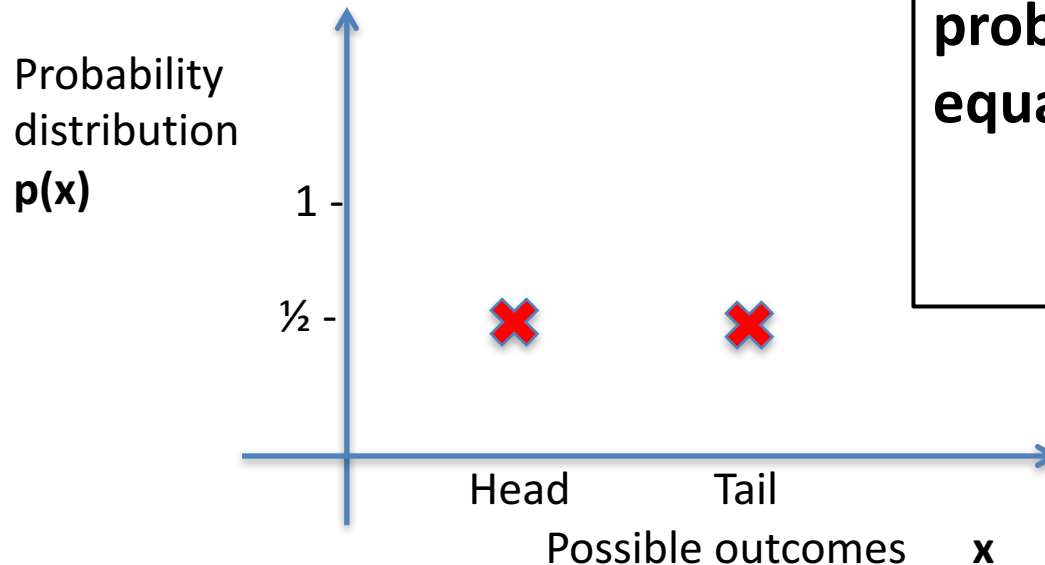
# Probability distribution

- It defines how likely is that any of the possible values of the random variable would come out.
- Examples:
  - Flip a coin, what is the probability distribution?  
What is the probability to get head or tail?  
➔ There are two possible outcomes, equally likely... so the probability of each one is  $1/2 = 0.5$ .
  - Throw a die, what is the probability distribution?  
➔ 6 possible outcomes, equally likely... so the probability of each outcome is  $1/6 = 0.166666..$

# Random variable example I

Flipping a coin:

- possible outcomes are head or tail
- Since they are equally likely, the probability distribution is uniform



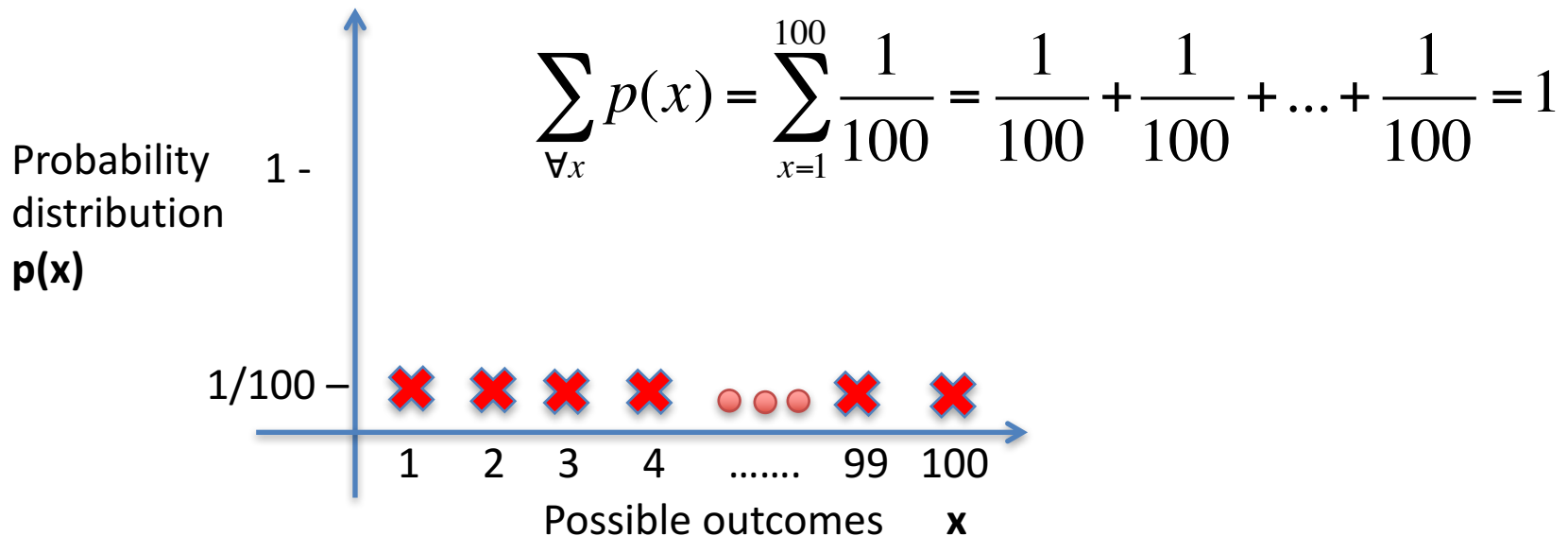
**Notice that the sum of all probabilities needs to be equal to 1:**

$$\sum_{\forall x} p(x) = 1$$

# Random variable example II

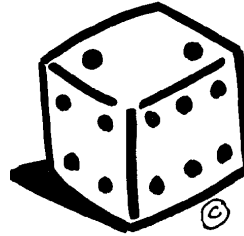
Lottery draw:

- possible outcomes are numbers, say from 1 to 100
- Since they are equally likely, the probability distribution is uniform



# Random variable example III

- Throwing one die:
  - Equal outcome for 6 values → uniform distribution with every outcome has probability  $1/6$

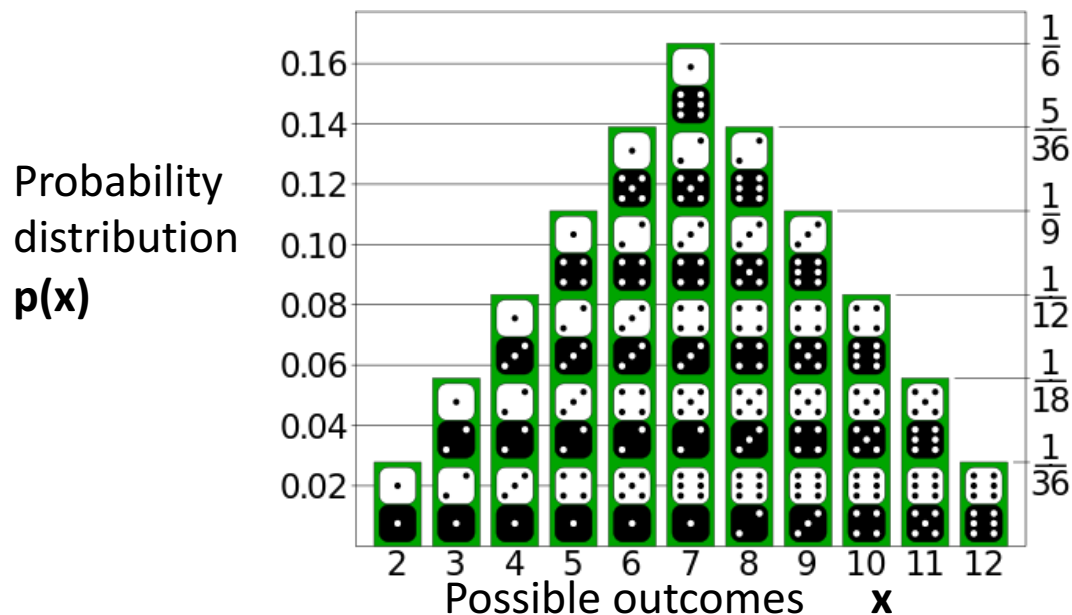


- Throwing two dice?
  - This is the sum of two independent outcomes:
  - What's the probability distribution?

# Probability distribution of two dice

- Let's see all possible outcomes, sum them and count the occurrences
- We have 36 possible outcomes
- We check how many times the sum is 1, how many times is 2...

Dice 1	Dice 2	Sum
1	1	2
1	2	3
1	3	4
...	...	...
2	1	3
2	2	4
...	...	...
6	6	12



# Shannon entropy

$$H(x) = - \sum_{\forall x} p(x) \log_2 p(x)$$

- $H(x)$  is a measure of the information carried by the random variable and is measured in bits
- *“Greater freedom of choice, greater uncertainty greater information go hand in hand”*
- This means that the higher the uncertainty of the random variable the greater the information it carries
- What carries more information, a random variable with less or more possible outcomes?
- What carries more information, a random variable with uniform or non-uniform distribution?





# Information of flipping a coin

$$\begin{aligned} H(x) &= -\sum_{\forall x} p(x) \log_2 p(x) = -\left(\frac{1}{2} \log_2 \frac{1}{2} + \frac{1}{2} \log_2 \frac{1}{2}\right) = \\ &= -\left(\frac{1}{2}(-1) + \frac{1}{2}(-1)\right) = -(-1) = 1 \end{aligned}$$

- Flipping a coin carries 1 bit of information

# Information of flipping an unfair coin

- Say the probability of heads are 0.9 and that of tails 0.1

$$\begin{aligned} H(x) &= - \sum_{\forall x} p(x) \log_2 p(x) = -(0.9 \log_2 0.9 + 0.1 \log_2 0.1) = \\ &= -(0.9(-0.152) + 0.1(-3.322)) = -(-0.469) = 0.469 \end{aligned}$$

**Flipping an unfair coin carries less information than one that is totally unpredictable.**

- This is understandable: if you know that heads are much more likely to occur, then by telling you that the outcome of the coin was head, you get less information, as you already expected it would probably be head

# Information of throwing one die

$$\begin{aligned} H(x) &= -\sum_{\forall x} p(x) \log_2 p(x) = -\left(\frac{1}{6} \log_2 \frac{1}{6} + \dots + \frac{1}{6} \log_2 \frac{1}{6}\right) = \\ &= -6 \cdot \left(\frac{1}{6} \log_2 \frac{1}{6}\right) = -\log_2 \frac{1}{6} = 2.585 \end{aligned}$$

- Throwing one die carries 2.585 bits of information

**A random variable with more possible outcomes carries more information**

# Uniform distributions carry more information

- Information of throwing two dice

$$H(x) = - \sum_{\forall x} p(x) \log_2 p(x) = - \left( \frac{1}{36} \log_2 \frac{1}{36} + \frac{1}{18} \log_2 \frac{1}{18} + \dots \right) = 3.2744$$

- Information of a random distribution between 2 and 12

$$H(x) = - \sum_{\forall x} p(x) \log_2 p(x) = - \sum_{x=2}^{12} \frac{1}{11} \log_2 \frac{1}{11} = - \log_2 \frac{1}{11} = 3.4594$$

**A random variable with uniform distribution carries the maximum amount of information**