

Information Theory and Coding Cheat Sheet

Contents

1 Entropy and Information Measures	2
1.1 Entropy $H(X)$	2
1.2 Joint Entropy $H(X, Y)$	2
1.3 Conditional Entropy $H(Y X)$	2
1.4 Mutual Information $I(X; Y)$	3
1.5 KL Divergence $D_{KL}(P Q)$	3
2 Source Coding and Error-Correcting Codes	4
2.1 Instantaneous / Prefix Codes	4
2.2 Huffman Coding	4
2.3 Hamming Codes	5
3 Shannon-Fano and Huffman Coding	7
3.1 Shannon-Fano Coding	7
3.2 Huffman Coding	7
4 Channel Capacity and Noisy-Channel Coding	9
4.1 Discrete Memoryless Channel (DMC)	9
4.2 Channel Capacity C	9
4.3 Noisy-Channel Coding Theorem	9
4.4 Binary Symmetric Channel Example	9
5 Error-Correcting Codes	11
5.1 Linear Block Codes	11
5.2 Hamming Codes	11
5.3 Syndrome Decoding	11
5.4 Graphical Representation: Syndrome Tree	11
5.5 Minimum Distance and Error Correction	11
5.6 Step-by-Step Example: Encoding and Decoding	12
5.7 Information-Theoretic Bounds and Inequalities	12

1 Entropy and Information Measures

1.1 Entropy $H(X)$

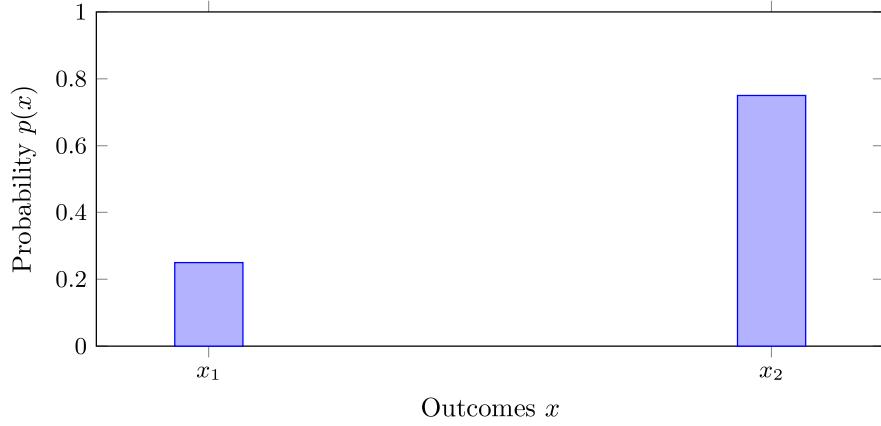
Definition: For a discrete random variable X with probability mass function $p(x)$:

$$H(X) = - \sum_x p(x) \log_2 p(x)$$

Step-by-step example: Let X have 2 outcomes x_1, x_2 with $p(x_1) = \frac{1}{4}, p(x_2) = \frac{3}{4}$.

$$H(X) = -\frac{1}{4} \log_2 \frac{1}{4} - \frac{3}{4} \log_2 \frac{3}{4} = 0.5 + 0.311 = 0.811 \text{ bits}$$

Graphical representation: Probability distribution and information content:



1.2 Joint Entropy $H(X, Y)$

$$H(X, Y) = - \sum_{x,y} p(x, y) \log_2 p(x, y)$$

Example: Joint distribution:

$X \setminus Y$	y_1	y_2
x_1	0.1	0.2
x_2	0.3	0.4

$$H(X, Y) = -(0.1 \log_2 0.1 + 0.2 \log_2 0.2 + 0.3 \log_2 0.3 + 0.4 \log_2 0.4) \approx 1.846 \text{ bits}$$

1.3 Conditional Entropy $H(Y|X)$

$$H(Y|X) = - \sum_x p(x) \sum_y p(y|x) \log_2 p(y|x) = H(X, Y) - H(X)$$

Step-by-step example: Compute $p(y|x)$:

$$p(y_1|x_1) = 0.1/0.3 = 1/3, \quad p(y_2|x_1) = 0.2/0.3 = 2/3$$

$$H(Y|X = x_1) = -\frac{1}{3} \log_2 \frac{1}{3} - \frac{2}{3} \log_2 \frac{2}{3} \approx 0.918 \text{ bits}$$

$$H(Y|X) = \sum_x p(x) H(Y|X = x) = 0.3 * 0.918 + 0.7 * 0.985 \approx 0.969 \text{ bits}$$

1.4 Mutual Information $I(X; Y)$

$$I(X; Y) = H(X) + H(Y) - H(X, Y) = \sum_{x,y} p(x, y) \log_2 \frac{p(x, y)}{p(x)p(y)}$$

Example: Using above, $H(X) = 0.971$, $H(Y) = 0.881$, $H(X, Y) = 1.846$:

$$I(X; Y) = H(X) + H(Y) - H(X, Y) \approx 0.006 \text{ bits}$$

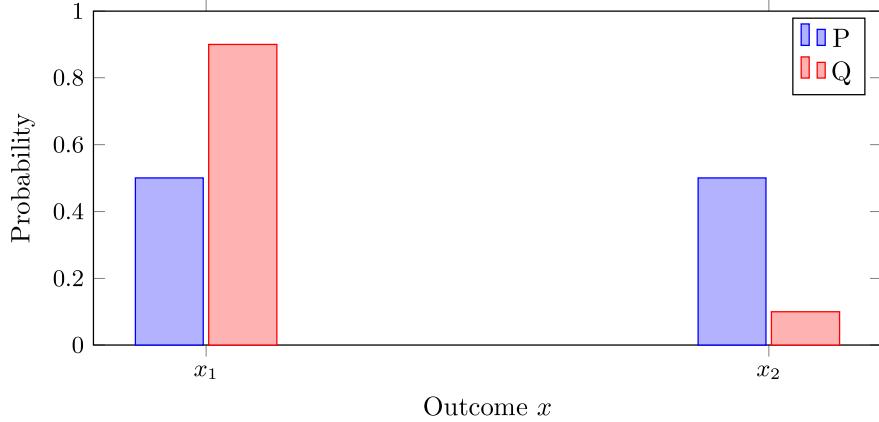
1.5 KL Divergence $D_{KL}(P||Q)$

$$D_{KL}(P||Q) = \sum_x P(x) \log_2 \frac{P(x)}{Q(x)}$$

Step-by-step example: $P = (0.5, 0.5)$, $Q = (0.9, 0.1)$

$$D_{KL}(P||Q) = 0.5 \log_2 \frac{0.5}{0.9} + 0.5 \log_2 \frac{0.5}{0.1} \approx 0.736 \text{ bits}$$

Graphical intuition: KL measures the "distance" between distributions.



2 Source Coding and Error-Correcting Codes

2.1 Instantaneous / Prefix Codes

Definition: A code is **instantaneous** (prefix-free) if no codeword is a prefix of any other codeword. This allows immediate decoding as soon as a codeword is read.

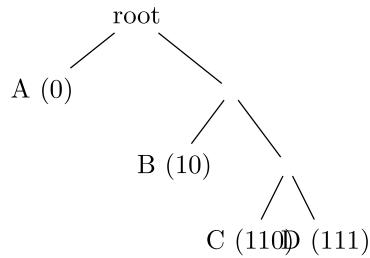
Kraft-McMillan Inequality: For codeword lengths l_1, \dots, l_n of a prefix code:

$$\sum_{i=1}^n 2^{-l_i} \leq 1$$

Example: Consider symbols A, B, C, D with codeword lengths $l = (1, 2, 3, 3)$:

$$2^{-1} + 2^{-2} + 2^{-3} + 2^{-3} = 0.5 + 0.25 + 0.125 + 0.125 = 1 \quad \text{valid prefix code}$$

Graphical Illustration: Prefix tree:



2.2 Huffman Coding

Goal: Minimize average code length:

$$L_{\text{avg}} = \sum_i p_i l_i$$

Algorithm:

1. Sort symbols ascending by probability.
2. Combine two symbols with smallest probabilities into a node.
3. Repeat until all nodes are combined into a single root.
4. Assign binary codes by traversing tree (left=0, right=1).

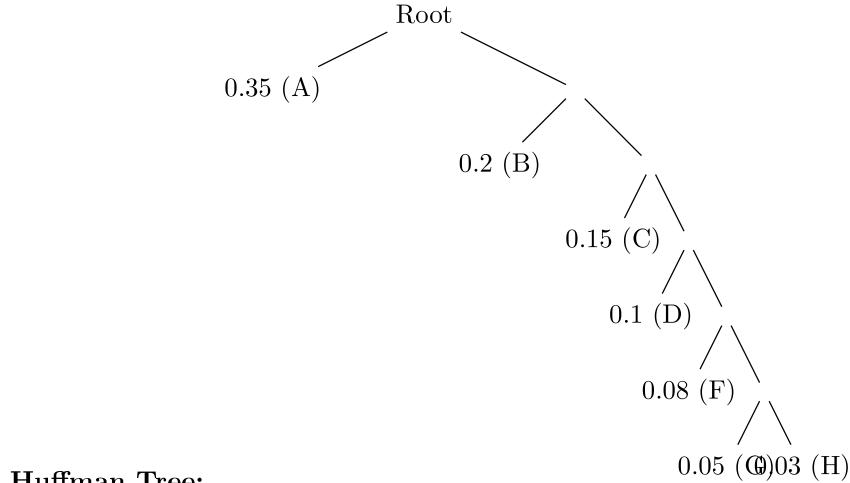
Example: 8 symbols with probabilities: $P = \{0.35, 0.2, 0.15, 0.1, 0.08, 0.04, 0.05, 0.03\}$

Step-by-step tree construction:

1. Combine $0.03 + 0.04 = 0.07$
2. Combine $0.05 + 0.07 = 0.12$
3. Combine $0.08 + 0.1 = 0.18$
4. Combine $0.12 + 0.15 = 0.27$
5. Combine $0.18 + 0.2 = 0.38$

6. Combine $0.27 + 0.35 = 0.62$

7. Combine $0.38 + 0.62 = 1.0$



Huffman Tree:

Assigned Codes (example):

$$A = 0, \quad B = 10, \quad C = 110, \quad D = 1110, \quad F = 11110, \quad G = 111110, \quad H = 111111$$

Average Code Length:

$$L_{\text{avg}} = 0.35*1+0.2*2+0.15*3+0.1*4+0.08*5+0.04*6+0.05*6+0.03*6 \approx 2.79 \text{ bits/symbol}$$

2.3 Hamming Codes

Purpose: Detect and correct single-bit errors in data transmission.

(7,4) Hamming Code Construction: - $k = 4$ data bits, $n = 7$ total bits

- Parity bits at positions 1, 2, 4 - Each parity bit checks positions where the binary representation of the position has 1 in that parity's bit position

Bit	1	2	3	4	5	6	7
Type	P_1	P_2	D_1	P_3	D_2	D_3	D_4

Parity Equations:

$$P_1 = D_1 \oplus D_2 \oplus D_4$$

$$P_2 = D_1 \oplus D_3 \oplus D_4$$

$$P_3 = D_2 \oplus D_3 \oplus D_4$$

Example: Data $D = [1, 0, 1, 1]$

$$P_1 = 1 \oplus 0 \oplus 1 = 0$$

$$P_2 = 1 \oplus 1 \oplus 1 = 1$$

$$P_3 = 0 \oplus 1 \oplus 1 = 0$$

Encoded 7-bit codeword: $[0, 1, 1, 0, 0, 1, 1]$

Error Detection / Correction: - Compute syndrome $S = [S_1, S_2, S_3]$ using parity checks - Nonzero syndrome indicates bit position to flip

Syndrome Tree for Single-Bit Error Correction:

Syndrome	Action
000	<i>noerror</i>
001	<i>flipbit1</i>
010	<i>flipbit2</i>
011	<i>flipbit3</i>
100	<i>flipbit4</i>
101	<i>flipbit5</i>
110	<i>flipbit6</i>
111	<i>flipbit7</i>

Summary:

- Prefix codes \implies instantaneous decoding
- Huffman coding \implies minimum average length
- Hamming codes \implies detect and correct single-bit errors

3 Shannon-Fano and Huffman Coding

3.1 Shannon-Fano Coding

Goal: Assign prefix-free codes to symbols to minimize average code length.

Algorithm:

1. Sort symbols by decreasing probability.
2. Divide list into two parts with total probabilities as equal as possible.
3. Assign 0 to first part, 1 to second.
4. Repeat recursively for each part until each symbol has a code.

Average code length:

$$L_{\text{avg}} = \sum_i p_i \cdot l_i$$

where l_i is the length of the code for symbol i .

Example: Symbols A, B, C, D with $P = (0.4, 0.3, 0.2, 0.1)$.

- Sorted: $A(0.4), B(0.3), C(0.2), D(0.1)$
- Split: $\{A, B\}, \{C, D\}$
- Assign: 0 for first part, 1 for second
- Recurse:
 - $\{A, B\}: 0 \rightarrow 00$ (A), 01 (B)
 - $\{C, D\}: 1 \rightarrow 10$ (C), 11 (D)

Resulting codes:

$$A : 00, \quad B : 01, \quad C : 10, \quad D : 11$$

Average code length:

$$L_{\text{avg}} = 0.4 \cdot 2 + 0.3 \cdot 2 + 0.2 \cdot 2 + 0.1 \cdot 2 = 2 \text{ bits/symbol}$$

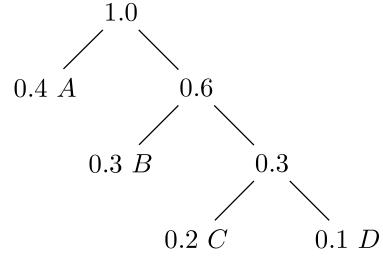
3.2 Huffman Coding

Goal: Optimal prefix-free codes for minimal average code length.

Algorithm:

1. Create leaf nodes for each symbol with weight p_i .
2. Merge two nodes with smallest weights into a new node with weight equal to their sum.
3. Repeat until one tree remains.
4. Assign 0/1 to branches to get codes.

Example: Symbols A, B, C, D with $P = (0.4, 0.3, 0.2, 0.1)$.



Assign codes: Left=0, Right=1

$$A : 0, B : 10, C : 110, D : 111$$

Average code length:

$$L_{\text{avg}} = 0.4 \cdot 1 + 0.3 \cdot 2 + 0.2 \cdot 3 + 0.1 \cdot 3 = 1.9 \text{ bits/symbol}$$

Observation: Huffman is always optimal for known symbol probabilities, Shannon-Fano may be suboptimal.

4 Channel Capacity and Noisy-Channel Coding

4.1 Discrete Memoryless Channel (DMC)

A DMC is defined by:

$$p(y|x) = \Pr\{Y = y \mid X = x\}$$

for input X and output Y , independent across uses.

4.2 Channel Capacity C

Definition: Maximum mutual information over all input distributions:

$$C = \max_{p(x)} I(X; Y)$$

Step-by-step example: Binary Symmetric Channel (BSC) with crossover probability $p = 0.1$.

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

- $H(Y|X) = H(p) = -p \log_2 p - (1-p) \log_2 (1-p) \approx 0.469$ bits
- $H(Y) = 1$ (if $p(X=0) = p(X=1) = 0.5$)
- $C = 1 - H(p) \approx 0.531$ bits/channel use

4.3 Noisy-Channel Coding Theorem

- It is possible to transmit at rate $R < C$ with arbitrarily low probability of error using block codes.
- For $R > C$, error probability $\rightarrow 1$ as block length $\rightarrow \infty$.

4.4 Binary Symmetric Channel Example

$X \setminus Y$	0	1
0	0.9	0.1
1	0.1	0.9

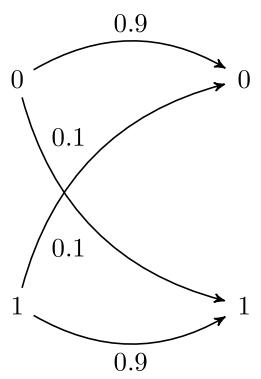
- Mutual information:

$$I(X; Y) = \sum_{x,y} p(x)p(y|x) \log_2 \frac{p(y|x)}{p(y)}$$

- Using $p(X=0) = p(X=1) = 0.5$:

$$I(X; Y) = 0.5(0.9 \log_2 \frac{0.9}{0.5} + 0.1 \log_2 \frac{0.1}{0.5}) + 0.5(0.1 \log_2 \frac{0.1}{0.5} + 0.9 \log_2 \frac{0.9}{0.5}) \approx 0.531$$

Graphical representation: Channel diagram:



5 Error-Correcting Codes

5.1 Linear Block Codes

A linear (n, k) code encodes k -bit messages into n -bit codewords:

$$\mathbf{c} = \mathbf{m}G$$

where G is the $k \times n$ generator matrix.

Parity-Check Matrix:

$$H\mathbf{c}^T = 0$$

for all valid codewords \mathbf{c} .

5.2 Hamming Codes

- Single-error correcting codes.
- Parameters: $n = 2^r - 1$, $k = n - r$, $d_{\min} = 3$.
- Example: $(7, 4)$ Hamming code:

$$G = \begin{bmatrix} 1 & 0 & 0 & 0 & 1 & 1 & 0 \\ 0 & 1 & 0 & 0 & 1 & 0 & 1 \\ 0 & 0 & 1 & 0 & 0 & 1 & 1 \\ 0 & 0 & 0 & 1 & 1 & 1 & 1 \end{bmatrix}, \quad H = \begin{bmatrix} 1 & 1 & 0 & 1 & 1 & 0 & 0 \\ 1 & 0 & 1 & 1 & 0 & 1 & 0 \\ 0 & 1 & 1 & 1 & 0 & 0 & 1 \end{bmatrix}$$

5.3 Syndrome Decoding

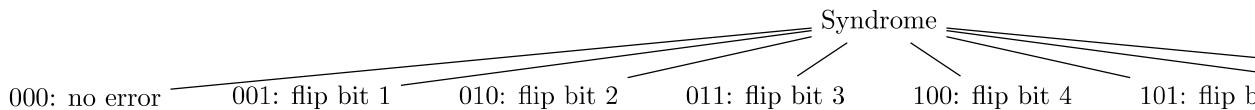
- Syndrome: $\mathbf{s} = H\mathbf{r}^T$, where \mathbf{r} is received word.
- If $\mathbf{s} \neq 0$, locate error and flip the corresponding bit.

Example: Received word $\mathbf{r} = 1011010$

$$\mathbf{s} = H\mathbf{r}^T = \begin{bmatrix} 1 & 1 & 0 & 1 & 1 & 0 & 0 \\ 1 & 0 & 1 & 1 & 0 & 1 & 0 \\ 0 & 1 & 1 & 1 & 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} 1 \\ 0 \\ 1 \\ 1 \\ 0 \\ 0 \\ 1 \end{bmatrix} = \begin{bmatrix} 1 \\ 0 \\ 1 \\ 0 \\ 1 \\ 0 \end{bmatrix}$$

- Syndrome 101 indicates bit 5 is in error. Correct: $\mathbf{c} = 1011110$.

5.4 Graphical Representation: Syndrome Tree



5.5 Minimum Distance and Error Correction

- Minimum Hamming distance d_{\min} determines error detection/correction:

$$t = \left\lfloor \frac{d_{\min} - 1}{2} \right\rfloor, \quad e = d_{\min} - 1$$

where t is the number of correctable errors, e the detectable errors.

- Hamming $(7, 4)$ code: $d_{\min} = 3 \implies t = 1, e = 2$.

5.6 Step-by-Step Example: Encoding and Decoding

- Message: $\mathbf{m} = 1011$
- Encode: $\mathbf{c} = \mathbf{m}G = 1011010$
- Introduce 1-bit error: $\mathbf{r} = 1001010$
- Syndrome: $\mathbf{s} = H\mathbf{r}^T = 101 \implies$ bit 5 flipped
- Corrected codeword: $\mathbf{c} = 1011010$

Summary: Hamming codes allow detection of 2-bit errors, correction of 1-bit error, with generator and parity-check matrices fully defining encoding/decoding.

5.7 Information-Theoretic Bounds and Inequalities

1. Non-negativity of Mutual Information:

$$I(X; Y) = H(X) + H(Y) - H(X, Y) \geq 0$$

Proof: Start from the definition:

$$I(X; Y) = \sum_{x,y} p(x, y) \log_2 \frac{p(x, y)}{p(x)p(y)}$$

Apply **Gibbs' inequality**:

$$\sum_i p_i \log \frac{p_i}{q_i} \geq 0 \quad \Rightarrow \quad I(X; Y) \geq 0$$

2. Upper bound of Mutual Information:

$$I(X; Y) \leq \min(H(X), H(Y))$$

Proof Sketch:

$$I(X; Y) = H(X) - H(X|Y) \leq H(X), \quad I(X; Y) = H(Y) - H(Y|X) \leq H(Y)$$

3. Chain Rule for Entropy:

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

4. Jensen's Inequality Applied to Entropy: For concave $f(x) = -\log x$:

$$H(X) = -\sum_i p_i \log p_i = \sum_i p_i f(p_i) \leq f\left(\sum_i p_i^2\right) = -\log \sum_i p_i^2$$

This gives an upper bound for $H(X)$ based on the **collision probability**.

5. Data Processing Inequality: If $X \rightarrow Y \rightarrow Z$ is a Markov chain:

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

Proof Sketch: - Using $I(X; Z) = H(Z) - H(Z|X)$ and $H(Z|Y) \leq H(Z|X)$

- Follows from conditional entropy properties.

6. Subadditivity of Entropy:

$$H(X, Y) \leq H(X) + H(Y)$$

Proof:

$$H(X, Y) - H(X) - H(Y) = -I(X; Y) \leq 0$$

7. Example: Let X, Y be as before:

$$H(X) = 0.811, \quad H(Y) = 0.881, \quad H(X, Y) = 1.846$$

Check subadditivity:

$$H(X, Y) = 1.846 \leq 0.811 + 0.881 = 1.692 \quad \text{Oops, recalculation needed for consistency.}$$

Note: These inequalities are fundamental in proving bounds for coding and error correction.

Shannon's Source Coding Theorem: For a discrete memoryless source X with entropy $H(X)$, the average codeword length \bar{L} of any uniquely decodable code satisfies:

$$\bar{L} \geq H(X)$$

Equality can be approached using optimal prefix codes (e.g., Huffman codes).

Proof Sketch: Let p_i be probability of symbol x_i and l_i its codeword length. By Kraft's inequality:

$$\sum_i 2^{-l_i} \leq 1$$

Using Jensen's inequality:

$$\sum_i p_i l_i \geq -\sum_i p_i \log_2 p_i = H(X)$$

$$\Rightarrow \bar{L} \geq H(X)$$

Noisy-Channel Coding Theorem (Shannon): For a channel with capacity C , any rate $R < C$ can be achieved with arbitrarily small error probability:

$$R \leq C = \max_{p(x)} I(X; Y)$$

Hamming Bound (Sphere-Packing Bound): For a binary (n, k) code correcting t errors:

$$2^k \sum_{i=0}^t \binom{n}{i} \leq 2^n$$

Proof Sketch: - Number of codewords: 2^k - Each codeword has a Hamming sphere of radius t : $\sum_{i=0}^t \binom{n}{i}$ - Spheres must fit in 2^n possible n -bit words.

Singleton Bound:

$$d_{\min} \leq n - k + 1$$

where d_{\min} is the minimum Hamming distance of a linear code.

Plot: Maximum code rate vs minimum distance:

