

(Affiliated with Nebula Institute of Technology, Sri Lanka)

Module Name : Information Theory

Module Code : 5FTC2151

Assignment Title : CW1 – Decoding Real-World Data

Submitted by:

Student Name : S. A. L. M. Mihiliya Jayasiri

Student ID : 23113610

Degree Programme : BSc (Hons) Data Science

Department : Physics, Engineering and Computer
Science

Academic Year : Year 2 - 2025/26

Submitted to:

Lecturer : Ms. Narmada Karunaratne

Date of Submission : 30/09/2025

Table of Contents

Task 1 - Entropy Analysis	3
1) Dataset Selection	3
2) Character Frequency Calculation	3
3) Computing Shannon's Entropy	4
4) Results and Interpretation	5
Entropy Value and It's Significance	5
Predictability Analysis Using Letter Frequency Comparison	5
Impact of Non-Letter Characters	6
5) Critical Reflection on Information Content	7
Task 2: Huffman Coding	8
1) Huffman Tree Construction	8
2) Huffman Codes	10
3) Encoding the Original Text	10
4) Compression Analysis	11
5) Reflection: Comparing Huffman-Encoded Size to Theoretical Entropy	11
Task 3: Hamming Code (Error Detection and Correction)	12
1) Selected Segment and Setup	12
2) Encoding the Two Nibbles	13
Nibble 1 (1001):	13
Nibble 2 (1011):	13
3) Single-Bit Error: Detection and Correction	14
Syndrome Calculation	14
Correction	14
4) Two-Bit Error: Undetected and Miscorrected	14
Syndrome Calculation	14
Correction Attempt	15

5) Discussion on Best Suited Data Types for Huffman Coding and LZW	15
Task 4: Entropy in a Non-Latin Language	17
1) Dataset Selection	17
2) Character Frequency Calculation	18
3) Computing Shannon's Entropy	19
4) Results and Interpretation	20
5) Challenges Faced.....	20
Annex.....	22
Annex A: Huffman Tree Visualization & Huffman Encoded Text	22
References.....	23

List of Tables

Table 1: Task 1 - Character Frequencies of Selected Text	4
Table 2: Task 2 - Probability & Entropy Contribution.....	5
Table 3: Huffman Codes from highest to lowest frequency	10
Table 4: Comparison of Features is Huffman Coding & LZW	16
Table 5: Task 4 - Character Frequencies.....	18
Table 6: Task 4 - Probability & Entropy Contribution.....	19
Table 7: Sinhala & English Script Comparison	20

List of Figures

Figure 1: English Letter Frequency Comparison Line Chart.....	6
Figure 2: Task 2 - Huffman Tree - Left Subtree (0-side)	9
Figure 3: Task 2 - Huffman Tree - Right Subtree (1-side).....	9

Task 1 - Entropy Analysis

Entropy in information theory, introduced by Claude Shannon in 1948, measures the average amount of information contained per character in a dataset. This report analyzes the entropy of a short English text dataset with 517 characters using character frequency analysis. The calculations were performed in Microsoft Excel, while Python was used for visualization.

1) Dataset Selection

The selected dataset is a short excerpt from the Wikipedia article on Entropy (Information Theory) (contributors, 2025) [2].

The concept of information entropy was introduced by Claude Shannon in his 1948 paper "A Mathematical Theory of Communication", and is also referred to as Shannon entropy. Shannon's theory defines a data communication system composed of three elements: a source of data, a communication channel, and a receiver. The "fundamental problem of communication" – as expressed by Shannon – is for the receiver to be able to identify what data was generated by the source, based on the signal it receives through the channel.

This text was chosen because it is domain relevant, sufficiently long and contains a mix of letters, digits, spaces, and punctuation.

2) Character Frequency Calculation

The frequency of each character was computed using Excel and is depicted in Table 1 below. Spaces and letters such as 'e', 'n', and 'a' dominate the dataset. The total dataset length is 517 characters.

Character	Frequency	Character	Frequency
e	52	v	3
n	41	w	3
a	40	x	1
o	37	S	4
t	31	T	3
i	24	C	2
r	23	A	1
s	21	M	1
h	20	1	1
c	19	4	1
d	16	8	1

m	15	9	1
f	11	Space	83
l	10	"	4
u	10	,	4
y	9	.	3
p	8	–	2
b	7	'	1
g	3	:	1

Table 1: Task 1 - Character Frequencies of Selected Text

3) Computing Shannon's Entropy

The Shannon's Entropy of each character was computed using Excel.

Shannon's Entropy Formula:

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

Where $p(x)$ is the probability of occurrence of each symbol.

Probabilities were computed as:

$$p(x) = \frac{\text{frequency of character}}{\text{Total character frequency}}$$

In the Table 2 below, $H(x)$ is the individual character entropy,

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

Character	p(x)	H(x)	Character	p(x)	H(x)
e	0.1006	0.3333	v	0.0058	0.0431
n	0.0793	0.2900	w	0.0058	0.0431
a	0.0774	0.2857	x	0.0019	0.0174
o	0.0716	0.2723	S	0.0077	0.0543
t	0.0600	0.2434	T	0.0058	0.0431
i	0.0464	0.2056	C	0.0039	0.0310
r	0.0445	0.1998	A	0.0019	0.0174
s	0.0406	0.1877	M	0.0019	0.0174
h	0.0387	0.1815	1	0.0019	0.0174
c	0.0368	0.1752	4	0.0019	0.0174
d	0.0309	0.1552	8	0.0019	0.0174

m	0.0290	0.1482	9	0.0019	0.0174
f	0.0213	0.1182	Space	0.1605	0.4237
l	0.0193	0.1101	"	0.0077	0.0543
u	0.0193	0.1101	,	0.0077	0.0543
y	0.0174	0.1017	.	0.0058	0.0431
p	0.0155	0.0931	—	0.0039	0.0310
b	0.0135	0.0840	'	0.0019	0.0174
g	0.0058	0.0431	:	0.0019	0.0174

Table 2: Task 2 - Probability & Entropy Contribution

Total Entropy $H(X)$ = 4.3159 bits per character

4) Results and Interpretation

Entropy Value and It's Significance

The Shannon entropy of the selected English text dataset was calculated to be 4.3159 bits per character. This value quantifies the average amount of information contained in each character of the text. In Information Theory, entropy reflects the degree of unpredictability or uncertainty in a dataset. A value close to 0 would indicate extreme redundancy (e.g., repeated characters), while a value approaching 5-8 bits would suggest high uncertainty, typical of encrypted or unstructured data.

An entropy of 4.3159 indicates that the text is moderately predictable, consistent with structured English writing. It contains enough variation to convey meaningful content without excessive repetition, making it suitable for efficient encoding and compression.

Predictability Analysis Using Letter Frequency Comparison

To critically interpret the entropy value, a comparison was made between the standard English letter frequencies (as published on Wikipedia (contributors, 2025) [3]) and the letter frequencies from the text dataset. Prior to analysis, all uppercase letters in the dataset were converted to lowercase to ensure consistency and alignment with the standard frequency reference, which is case-insensitive. This normalization step reduces alphabet size and stabilizes frequency distributions, allowing for a more accurate entropy calculation. The comparison was visualized using a line chart (generated in python) which is illustrated below in Figure 1.

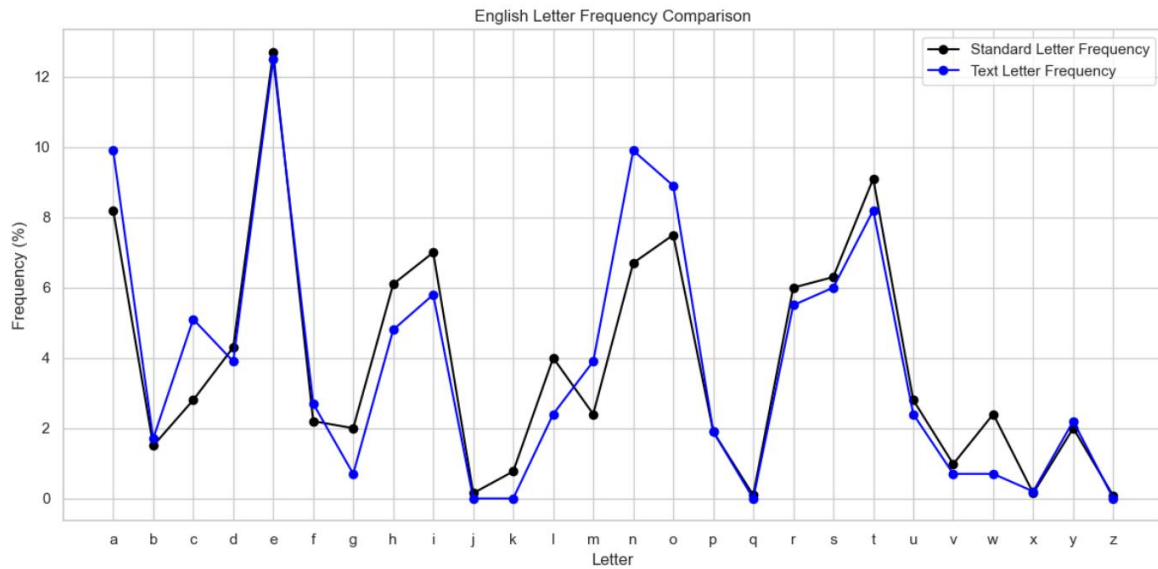


Figure 1: English Letter Frequency Comparison Line Chart

Key Observations: The letter ‘e’ has the highest frequency in both datasets, appearing slightly above 12%. This strong alignment reinforces the text’s adherence to natural English structure. Other high-frequency letters such as ‘t’, ‘a’, ‘o’, and ‘n’ also show similar proportions across both datasets, indicating that the text follows conventional linguistic patterns. Low-frequency letters like ‘q’, ‘x’, and ‘z’ remain rare in both cases, contributing minimally to the entropy.

Fluctuations and Their Impacts: Slight increases in the frequency of letters such as ‘c’, ‘m’, ‘n’, and ‘o’ suggest the presence of domain-specific vocabulary, likely due to technical nature of the text (e.g., terms like “communication”, “entropy”, and “Shannon”). These fluctuations introduce mild unpredictability, enriching the information content without disrupting the overall structure.

Interpretation: The close alignment between the two frequency curves confirms that text is linguistically typical, supporting the entropy value of 4.3159. The minor fluctuations may assume informational depth, but not uncertainty which resulting a text that is predictable enough for efficient encoding yet varied enough to avoid redundancy.

Impact of Non-Letter Characters

The inclusion of digits and punctuation marks though less frequent adds to the character diversity. These elements introduce less predictable symbols, which slightly increase the entropy.

5) Critical Reflection on Information Content

The calculated entropy value of 4.3159 bits per character suggests that the text is efficiently encoded, with a good balance between redundancy and novelty. It is neither overly repetitive nor excessively random.

Furthermore, the entropy reflects the clarity and coherence of the text. It is informative without being noisy, making it suitable for both human reading and machine processing.

By extension, this observation suggests that Wikipedia tends to exhibit moderate predictability and rich information content. Its editorial standards and structured language contribute to a character distribution that is neither overly uniform nor chaotic, aligning with the entropy characteristics observed in the text dataset. This reinforces Wikipedia's role as a reliable and computationally tractable source for tasks involving compression, encoding, and information analysis.

Task 2: Huffman Coding

Huffman coding is a fundamental lossless data compression technique that assigns variable-length codes to characters based on their frequencies. Characters that occur more frequently receive shorter codes, while less frequent characters are assigned longer codes.

In this task, the character frequency data obtained in Task 1 was used to construct a Huffman tree, generate Huffman codes, and encode the selected text. The efficiency of Huffman coding is then evaluated by comparing the encoded text length with the theoretical entropy and fixed-length ASCII encoding.

1) Huffman Tree Construction

The Huffman tree was constructed using the character frequencies derived from Task 1, which analyzed a 517-character English text excerpt. Python was used to implement the Huffman coding algorithm, ensuring that each character was assigned a binary code based on its relative frequency.

To enhance clarity and visual interpretation, the tree was separated into two distinct subtrees:

- **Left Subtree:** Represents all codes beginning with 0.
- **Right Subtree:** Represents all codes beginning with 1.

This bifurcation was not algorithmically necessary but was intentionally applied to improve the readability of the tree structure and to clearly demonstrate the hierarchical placement of each node.

The sum of the highest- frequency nodes from both subtrees equals 517, which corresponds to the total number of characters in the dataset. This confirms that the tree accurately encapsulates the full character distribution and preserves the integrity of the original dataset.

Frequent characters such as the space and 'e' appear closer to the root, resulting in shorter Huffman codes, while rarer characters are placed deeper in the tree, receiving longer codes. This structure ensures optimal compression efficiency, aligning with the principles of prefix coding.

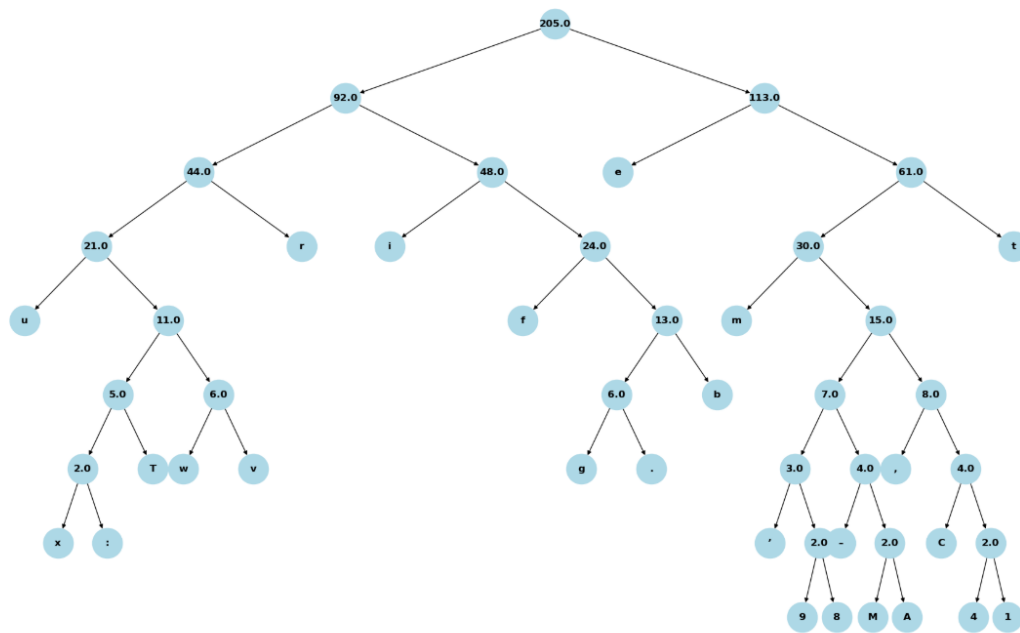


Figure 2: Task 2 - Huffman Tree - Left Subtree (0-side)

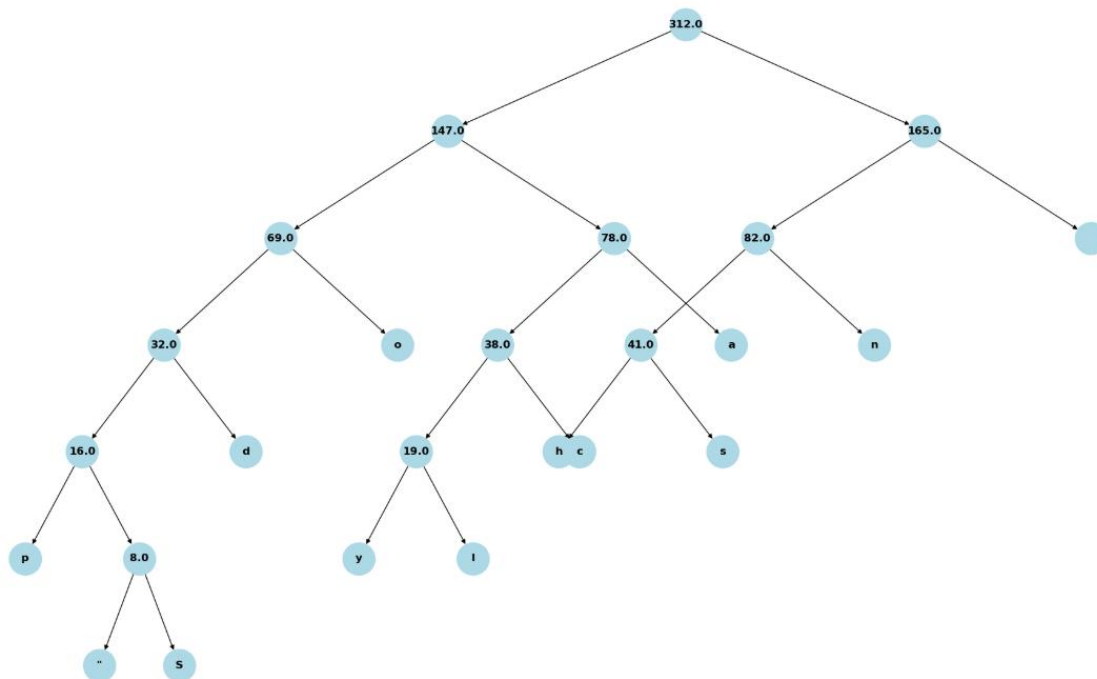


Figure 3: Task 2 - Huffman Tree - Right Subtree (1-side)

For better visualization, a clear Huffman tree was constructed using draw.io. The complete tree is provided in Annex A (Figure A1) for reference.

2) Huffman Codes

The Huffman coding algorithm generated the following codes for the dataset characters as in Table 3 using python.

Character	Frequency	Huffman Code	Character	Frequency	Huffman Code
Space	83	111	S	4	1000011
e	52	010	"	4	1000010
n	41	1101	,	4	0110110
a	40	1011	g	3	0011100
o	37	1001	v	3	0000111
t	31	0111	w	3	0000110
i	24	0010	T	3	0000101
r	23	0001	.	3	0011101
s	21	11001	C	2	01101110
h	20	11000	–	2	01101010
c	19	10101	x	1	0000100
d	16	10001	A	1	011010111
m	15	01100	M	1	011010110
f	11	00110	1	1	011011111
l	10	101001	4	1	011011110
u	10	00000	8	1	011010011
y	9	101000	9	1	011010010
p	8	100000	,	1	01101000
b	7	001111	:	1	00001001

Table 3: Huffman Codes from highest to lowest frequency

3) Encoding the Original Text

The original 517-character text excerpt was encoded using the Huffman codes using python.

- **Original text length (ASCII, fixed length):** $517 \times 8 = 4136$ bits
- **Encoded text length (Huffman):** 2252 bits

The resulting encoded text is a binary string composed of Huffman codewords for each character; fully encoded output of the original text is in Annex A (Figure A2) for reference.

4) Compression Analysis

Calculations were done using python.

- **Total bits in the encoded message:** 2252 bits
- **Average bits per character:** 4.3559 bits per character
- **Compression ratio compared to fixed-length ASCII encoding (assuming 8 bits per character):**

$$\text{Compression Ratio} = 4136 / 2252 = 1.8366$$

This means the Huffman-encoded message is almost twice as compact as ASCII encoding.

5) Reflection: Comparing Huffman-Encoded Size to Theoretical Entropy

The Shannon entropy calculated in Task 1 was 4.3159 bits per character, representing the theoretical lower bound for any lossless compression scheme based on the character distribution of the given text. After applying Huffman coding, the encoded text resulted in a total of 2252 bits, with an average of 4.3559 bits per character.

This outcome aligns closely with theoretical expectations. Huffman coding is known to produce an average code length that satisfies the inequality:

$$H \leq L \leq H + 1$$

Where H is the Shannon entropy and L is the average number of bits per character in the encoded text (Ayfer Ozgur, 2024) [1].

$$4.3159 \leq 4.3559 \leq 5.3159$$

In this case, the Huffman average is only 0.04 bits per character above the entropy, which is a minimal overhead. This slight increase is due to Huffman's use of integer length codewords and the finite sample size of the dataset. Rare characters such as punctuations typically receive longer codewords, contributing to the marginal difference.

Overall, the Huffman encoding demonstrates near optimal performance for the given character distribution. The result confirms that Huffman coding is highly efficient for symbol-by-symbol compression and closely tracks the theoretical entropy, with only minor deviation due to structural constraints of prefix coding.

Task 3: Hamming Code (Error Detection and Correction)

Hamming codes, introduced by Richard Hamming in 1950, are a family of linear error-correcting codes that add redundant parity bits to transmitted data. Their main purpose is to detect and correct single-bit errors and, in some cases, detect two-bit errors.

In this task, the Hamming (7, 4) scheme is applied to a short binary segment derived from the Huffman-encoded text in Task 2. The report demonstrates how single-bit errors are corrected, why two-bit errors cause problems, and concludes with a comparison between Huffman coding and Run-Length Encoding (RLE).

1) Selected Segment and Setup

- **Selected 8-bit segment:** 10011011 (Random)
- **Nibble split:** 1001 and 1011
- **Encoding scheme:** Hamming (7, 4) with even parity
- **Bit layout:** positions 1, 2, 4 are parity; data bits at positions 3 = d1, 5 = d2, 6 = d3, 7 = d4
- **Parity coverage:**

Bit position	7	6	5	4	3	2	1
Bit Role	d4	d3	d2	p4	d1	p2	p1
Bit position in binary	0111	0110	0101	0100	0011	0010	0001

- p1: positions {1, 3, 5, 7}
 - p2: positions {2, 3, 6, 7}
 - p4: positions {4, 5, 6, 7}
- **Syndrome bits on receive (s1, s2, s4) are computed under even parity:**
 - s1 = parity (r1, r3, r5, r7)
 - s2 = parity (r2, r3, r6, r7)
 - s4 = parity (r4, r5, r6, r7)

and error position is $(s4s2s1)_2$ with 0 meaning “no error”

2) Encoding the Two Nibbles

Nibble 1 (1001):

Bit Position	7	6	5	4	3	2	1
Bit Role	d4	d3	d2	p4	d1	p2	p1
Parity Bits	1	0	0	1	1	0	0
Parities in the Positions	p1, p2, p3	p2, p3	p1, p3	p3	p1, p2	p2	p1

- p1: positions 1, 3, 5, 7 → _, 1, 0, 1 → even → p1 = 0
- p2: positions 2, 3, 6, 7 → _, 1, 0, 1 → even → p2 = 0
- p4: positions 4, 5, 6, 7 → _, 0, 0, 1 → odd → p3 = 1

Codeword: 1001100

Nibble 2 (1011):

Bit Position	7	6	5	4	3	2	1
Bit Role	d4	d3	d2	p4	d1	p2	p1
Parity Bits	1	0	1	0	1	0	1
Parities in the Positions	p1, p2, p3	p2, p3	p1, p3	p3	p1, p2	p2	p1

- p1: positions 1, 3, 5, 7 → _, 1, 1, 1 → odd → p1 = 1
- p2: positions 2, 3, 6, 7 → _, 1, 0, 1 → even → p2 = 0
- p4: positions 4, 5, 6, 7 → _, 1, 0, 1 → even → p3 = 0

Codeword: 1010101

Encoded 14-bit segment: 10011001010101

3) Single-Bit Error: Detection and Correction

- **Error introduced:** bit flipped at position 5 in first codeword.
 - **Original:** 1001100
 - **Received:** 1011100

Syndrome Calculation

7	6	5	4	3	2	1
1	0	1	1	1	0	0

- s1: parity of positions 1, 3, 5, 7 $\rightarrow 0, 1, 1, 1 \rightarrow \text{odd} \rightarrow s1 = 1$
- s2: parity of positions 2, 3, 6, 7 $\rightarrow 0, 1, 0, 1 \rightarrow \text{even} \rightarrow s2 = 0$
- s3: parity of positions 4, 5, 6, 7 $\rightarrow 1, 1, 0, 1 \rightarrow \text{odd} \rightarrow s3 = 1$

Syndrome: 101 \rightarrow binary 5

Correction

Bit at position 5 flipped back.

Corrected codeword: 1001100

In conclusion, single-bit errors are reliably detected and corrected using the syndrome.

4) Two-Bit Error: Undetected and Miscorrected

- **Error introduced:** bit flipped at positions 2 and 6 in second codeword.
 - **Original:** 1010101
 - **Received:** 1110111

Syndrome Calculation

7	6	5	4	3	2	1
1	1	1	0	1	1	1

- s1: parity of positions 1, 3, 5, 7 $\rightarrow 1, 1, 1, 1 \rightarrow \text{even} \rightarrow s1 = 0$
- s2: parity of positions 2, 3, 6, 7 $\rightarrow 1, 1, 1, 1 \rightarrow \text{even} \rightarrow s2 = 0$
- s3: parity of positions 4, 5, 6, 7 $\rightarrow 0, 1, 1, 1 \rightarrow \text{odd} \rightarrow s3 = 1$

Syndrome: 001 → binary 1

This incorrectly suggests an error at position 1, even though the actual errors are at positions 2 and 6.

Correction Attempt

Decoder flips bit at position 1.

Miscorrection: 1110110 \neq 1010101

In conclusion, two-bit errors produce misleading syndromes that point to a third bit, resulting in incorrect correction, where the original data is corrupted.

5) Discussion on Best Suited Data Types for Huffman Coding and LZW

Huffman Coding and Lempel-Ziv-Welch (LZW) are both lossless compression techniques, yet they differ fundamentally in their approach to redundancy and data structure. Understanding these differences is essential when selecting the appropriate method for a given dataset.

Huffman coding is a statistical compression algorithm that assigns variable-length codes to input symbols based on their frequency. As (Sayood, 2018) [6], it is optimal when the symbol probabilities are known and data source is memoryless (no memory of past outputs). The technique constructs a binary tree where more frequent symbols receive shorter codes, minimizing the average code length.

This method is particularly effective for datasets with skewed symbol distribution, such as English text, grayscale images, or sensor logs. Its performance closely aligns Shannon entropy, making it ideal for entropy-aware encoding tasks. However, Huffman coding requires a preprocessing step to calculate character frequencies, which can be a limitation in streaming or real-time applications.

LZW, in contrast, is a dictionary-based algorithm that builds a codebook dynamically during encoding. It does not require prior knowledge of symbol frequencies and instead compresses data by referencing previously seen substrings. (Sayood, 2018) [6] highlights LZW's adaptability, noting its strength in handling data with repeated patterns such as source code, markup languages, and structured logs.

LZW is particularly well suited for medium or large datasets where patterns emerge over time. Its single pass encoding makes it efficient for streaming scenarios, and its use in formats like GIF and TIFF underscores its practical relevance. However, LZW may be less effective on short or highly random data, where dictionary growth is limited.

Feature	Huffman Coding	LZW Compression
Compression Type	Statistical	Dictionary-based
Optimization Target	Symbol frequency	Pattern repetition
Preprocessing Required	Yes	No
Adaptability	Static (unless extended)	Highly adaptive
Best Data Types	Skewed text, grayscale images	Source code, markup, structured logs
Entropy Alignment	Direct (Shannon entropy)	Indirect (pattern entropy)

Table 4: Comparison of Features is Huffman Coding & LZW

Both techniques offer distinct advantages depending on the structure and entropy profile of the input data. Huffman coding excels in symbol-level optimization, while LZW thrives in pattern-rich environments. The choice between them should be guided by the dataset's characteristics and operational constraints such as whether preprocessing is feasible or streaming efficiency is required.

Task 4: Entropy in a Non-Latin Language

Information entropy can be applied across different languages and scripts to measure the unpredictability and information density of text. While English uses the Latin alphabet with relatively simple character composition, non-Latin scripts such as Sinhala introduce unique challenges due to complex grapheme clusters, diacritics, and joining mechanisms.

This task analyses the entropy of a Sinhala text excerpt using Shannon's entropy formula, compares the results with the English dataset from Task 1, and reflects on the challenges of working with non-Latin scripts.

1) Dataset Selection

The selected dataset is a short excerpt from the Wikipedia article “Sri Lanka” (contributors, 2025) [5].

ශ්‍රී ලංකාවේ විශාලත්වය වර්ග කි.මී 65,610 (වර්ග සැතපුම් 25,330) වේ. එය ලොව 123 වැනි විශාලතම රට වන අතර විශාලතම දූපත් අතරින් 25 වන ස්ථානයෙහි සිටියි. එහි භූමි ප්‍රදේශය පළාත් නවයකට සහ දිස්ත්‍රික්ක විසිහතරකට වෙන් කර ඇත. රට පළාත් 9 කින්, දිස්ත්‍රික්ක 25 හා ප්‍රදේශය ලේකම් කොට්ඨාශය 331 හා ග්‍රාම නිලධාරී කොට්ඨාශ 14,022 කින් සමන්විත වේ. අගනුවර ශ්‍රී ජයවර්ධනපුර කොට්ටේ වන අතර විශාලතම නගරය කොළඹ වේ. කොළඹ ජාතියේ ආර්ථික, දේශපාලන කේන්ද්‍රස්ථානයයි. ශ්‍රී ලංකාව ලෝක උරුම ස්ථාන අටක්, සැතපුම් ගණනක් පුරා විහිදුනු රන්වන් වැලි සහිත වෙරළවල්, වැසි වනාන්තර සහ කඳුකර තේ වතු සහිත ඉතා සුන්දර දූපතකි. නිවර්තන වනාන්තර, වෙරළ සහ භූ දර්ශනවල ස්වභාවික සුන්දරත්වය, ජෛව විවිධත්වය මෙන්ම පොහොසත් සංස්කෘතික උරුමයන් සඳහා ප්‍රසිද්ධය, එය ලෝක ප්‍රසිද්ධ සංචාරක ගමනාන්තයක් බවට පත් කළේය.

The non-Latin dataset (Sinhala script) was encoded using UTF-8, a variable-length character encoding capable of representing all Unicode characters. UTF-8 ensures compatibility across platforms and preserves the integrity of Sinhala graphemes, including conjunct forms and diacritics. Sinhala text differs from English in that many visible characters (syllables) are formed by combining multiple Unicode code points (e.g., consonant + virama + ZWJ + another consonant). This encoding was explicitly set during file reading and verified using python's chardet and manual inspection.

It consists of letters, spaces, punctuation, digits and Zero Width Joiners (ZWJ) (contributors, 2025) [4].

2) Character Frequency Calculation

The frequency of each unique character was calculated using Excel, except for ZWJ, which was calculated using python. The total length of the Sinhala dataset is 741 characters.

Character	Frequency	Character	Frequency
ව	43	ඊ	61
ර	40	උ	41
න	36	ඌ	27
ක	33	ඹ	13
ක	33	ඹේ	13
ස	26	ඹො	8
ය	20	ඹු	5
ප	15	ඹේ	5
ම	15	ඹ	4
ල	15	ඹු	4
ද	14	ඹො	4
ට	13	ඹේ	2
ශ	13	ඹා	1
හ	12	ඹේ	1
ග	7		
ඒ	7	0	3
අ	5	1	4
ධ	5	2	6
ථ	4	3	5
ඵ	3	4	1
ඡ	3	5	4
භ	3	6	2
උ	2	9	1
ඳ	2		
ඹ	2	(1
ආ	1)	1
ඉ	1	,	10
ඇ	1	.	9
ච	1		
ඡ	1	Space	117
බ	1	ZWJ	11

Table 5: Task 4 - Character Frequencies

3) Computing Shannon's Entropy

The Shannon's Entropy of each character was computed using Excel, following the same method as in Task 1.

Character	p(x)	H(x)	Character	p(x)	H(x)
ଏ	0.0580	0.2383	ୠ	0.0823	0.2966
ଋ	0.0540	0.2273	ୡ	0.0553	0.2310
ଅ	0.0486	0.2120	ୢ	0.0364	0.1741
ଇ	0.0445	0.1999	ୣ	0.0175	0.1023
ଈ	0.0445	0.1999	୤	0.0175	0.1023
ଊ	0.0351	0.1696	୥	0.0108	0.0705
ଋ	0.0270	0.1407	୦	0.0067	0.0487
ୠ	0.0202	0.1139	୧	0.0067	0.0487
ୡ	0.0202	0.1139	୨	0.0054	0.0407
ୢ	0.0202	0.1139	୩	0.0054	0.0407
ୣ	0.0189	0.1082	୪	0.0054	0.0407
୤	0.0175	0.1023	୫	0.0027	0.0230
୥	0.0175	0.1023	୬	0.0013	0.0129
୦	0.0162	0.0963	୭	0.0013	0.0129
୧	0.0094	0.0635			
୨	0.0094	0.0635	୦	0.0040	0.0322
୩	0.0067	0.0487	୧	0.0054	0.0407
୪	0.0067	0.0487	୨	0.0081	0.0563
୫	0.0054	0.0407	୩	0.0067	0.0487
୬	0.0040	0.0322	୪	0.0013	0.0129
୭	0.0040	0.0322	୫	0.0054	0.0407
୮	0.0040	0.0322	୬	0.0027	0.0230
୯	0.0027	0.0230	୭	0.0013	0.0129
୦	0.0027	0.0230			
୧	0.0027	0.0230	(0.0013	0.0129
୨	0.0013	0.0129)	0.0013	0.0129
୩	0.0013	0.0129	,	0.0135	0.0838
୪	0.0013	0.0129	.	0.0121	0.0773
୫	0.0013	0.0129			
୬	0.0013	0.0129	Space	0.1579	0.4205
୭	0.0013	0.0129	ZWJ	0.0148	0.0902

Table 6: Task 4 - Probability & Entropy Contribution

Total Entropy $H(X)$ = 4.8561 bits per character

4) Results and Interpretation

Entropy Value: The calculated entropy of the Sinhala dataset is 4.8561 bits per character which is 0.5402 bits higher than the English dataset analyzed in Task 1 (4.3159 bits per character). This increase reflects the greater unpredictability and diversity of symbols in Sinhala text.

Feature	Sinhala	English
Alphabet Size	60+ distinct grapheme components (base letters, vowel signs, diacritics, ZWJ)	26 base letters; 52 with case sensitivity
Character Distribution	More balanced across characters; includes rare graphemes and modifiers	Skewed toward frequent letters like 'e', 't', 'a'
Script Complexity	Dense encoding; syllables formed by combining multiple components	Simple encoding: characters are independent units

Table 7: Sinhala & English Script Comparison

The broader and more balanced distribution in Sinhala contributes to a higher entropy value, as each character carries more information on average. Although, the visual length of Sinhala text may appear shorter, the underlying encoding is denser and more complex.

5) Challenges Faced

Issued Identified: During initial character frequency analysis in Excel, a discrepancy was observed between the total Unicode code points and the visible character count. Specifically,

- Total string length: 741
- Characters tallied: 730

This mismatch raised concerns about the accuracy of entropy calculations.

Root Cause: The Sinhala text uses Zero Width Joiner (ZWJ, U+200D) to form request conjunct formation and half-forms (ligatures) (e.g., rakaransaya and yansaya). These joiners are:

- Invisible in rendered text
- Counted as individual code points
- Essential for accurate grapheme formation

Thus, the discrepancy stemmed from treating grapheme clusters as Sinhala characters, while the encoding counted each constituent code point.

Resolution Strategy: To resolve the issue:

- Developed a Python script to explicitly detect ZWJ (U+200D) and other special Unicode marks.
- Identified 11 ZWJ characters, reconciling the count between code points and visible graphemes.

This confirmed that grapheme clusters \neq code points, and that entropy analysis must be based on code point-level frequency, not visual character count.

Insight Gained: This challenge highlighted a critical nuance in non-Latin text analysis:

Unicode-aware processing is essential for accurate entropy computation in languages like Sinhala, where visual characters are often composites of multiple code points.

It also reinforced the importance of distinguishing between rendered text and encoded data, especially when working with compression, entropy, or linguistic diversity.

Annex A: Huffman Tree Visualization & Huffman Encoded Text



Encoded binary string: 000010111000010111101011001110101010100000011111100100110111001011010
01101001000101100101101110010100111011110101101011100010011000001010001110000110101111001111001
011010111000110011000100000101010100011100111101000111011011101001101100000100010101111000
011110001011110111010011101111001011011110000101100111101110101001001101110011010011
110000010111000001000011110000100110101111101010101010111100001001100101110010101010
111010011110000101110000101001000110100011100100110111011011010010110001100000001101001010101
0110111001010011101100001001101101110111011000111001011001111011101001110011001110001010001
1001000010001010100011101111001111011100111100001110001011101110110111011101011010111000
11001100000101000011011110000111000010111011101101101101101000100011101111000010001000011
0100011110001010001100010101010111101111100010110111011110101100101100011000000001010
0101010110110111001010011101111100110100011001011101001100111101011001011001000010011100101010
00111110010011101111100000010100101110101010010100110001011010111100100001001111011111100
100100000000110101011110010011011100011001101110110110111101110101100101100011000000011
0100101010110111001010011101111010110001011101101010100101101111011101100011110111
11000101010101010001000001101000010011011100001011100001011100001000100000011010001011011
00010101011101110100111100000000110010011110100101001100111100100110111101011001011000110000
00011010010101011011100101001110110000101101101010111011110011110100001000100000001010110
011100101010001111001111000011100001111000101110111011001101110101011100101100111100110
1001000111101111000010111000101010101000100001101000011101110011100111100111101101111
010010101110111100111001010001010110010001100100011100001011000100110111111000110110111
10111110000110101111001111001110001011010100001101101110101000111100111110100011101111000010111
110011001000000001010101001101011100111101111001010100011110011101111101111100001011111001001
00011100110110110100111100100111110001010101000100000111010100111101111100000011001000000
01110011000110111111000010111010110001011101101010100100011011101111100000011001000000

22

References

- [1] Ayfer Ozgur, S. U., 2024. *Lecture 4: Entropy and Lossless Compression*. [Online]
Available at: <https://web.stanford.edu/class/engr76/lectures/lecture4.pdf>
[Accessed 30 08 2025].
- [2] contributors, W., 2025. *Entropy (information theory)*. [Online]
Available at:
[https://en.wikipedia.org/w/index.php?title=Entropy_\(information_theory\)&oldid=1300594028](https://en.wikipedia.org/w/index.php?title=Entropy_(information_theory)&oldid=1300594028)
[Accessed 30 08 2025].
- [3] contributors, W., 2025. *Letter frequency*. [Online]
Available at:
https://en.wikipedia.org/w/index.php?title=Letter_frequency&oldid=1311139883
[Accessed 30 08 2025].
- [4] contributors, W., 2025. *Zero-width joiner*. [Online]
Available at: https://en.wikipedia.org/w/index.php?title=Zero-width_joiner&oldid=1267965011
[Accessed 01 09 2025].
- [5] contributors, ව., 2025. *ශ්‍රී ලංකාව*. [Online]
Available at:
https://si.wikipedia.org/w/index.php?title=%E0%B7%81%E0%B7%8A%E2%80%8D%E0%B6%BB%E0%B7%93_%E0%B6%BD%E0%B6%82%E0%B6%9A%E0%B7%8F%E0%B7%80&oldid=739740
[Accessed 01 09 2025].
- [6] Sayood, K., 2018. *Introduction to Data Compression*. 5th ed. Cambridge: Morgan Kaufmann Publishers.