

**HKUSTx:** ELEC1200.2x A System View of Communications: From Signals to...

- Pre-course Materials
- ► Topic 1: Course Overview
- Topic 2: LosslessSource Coding: HammingCodes
- ► Topic 3: The Frequency Domain
- ► Topic 4: Lossy Source Coding
- Topic 5: Filters and the Frequency Response
- Topic 6: The Discrete Fourier Transform
- ➤ Topic 7: Signal Transmission -Modulation
- Topic 8: Signal Transmission -Demodulation
- ▶ Topic 9: IQ

### LOSSLESS SOURCE CODING/ENTROPY

### SECTION 1 QUESTION 1 BACKGROUND

Suppose these 7 symbols appear with the following probabilities:

Symbol	Probability
A	0.25
Е	0.02
Н	0.1
K	0.2
S	0.05
T	0.03
U	0.35

## SECTION 1 QUESTION 1 PART A (2/2 points)

What is the entropy of the symbols, assuming the probabilities given above? Give your answer to two significant digits (e.g. 1.00).

2.31

You have used 1 of 1 submissions

# SECTION 1 QUESTION 1 PART B (2/2 points)

Find a Huffman code for these symbols.

How long is the codeword for each symbol? Please select the correct answer.

ullet T 5 E 5 S 4 H 3 A 2 K 2 U 2 ullet

#### Modulation

- Topic 10: Summary and Review
- **▼** Final Exam

#### **Final Exam**

Final Exam due Dec 07, 2015 at 16:00 UTC

- MATLAB download and tutorials
- MATLAB Sandbox
- Post Course Survey

- lacksquare  $E \ 6 \ T \ 5 \ H \ 4 \ S \ 4 \ K \ 3 \ A \ 2 \ U \ 1$
- lacksquare T 5 E 5 S 4 H 3 A 2 K 2 U 1
- lacksquare  $E \ 6 \ T \ 5 \ S \ 4 \ H \ 3 \ K \ 3 \ A \ 2 \ U \ 2$

You have used 1 of 1 submissions

## SECTION 1 QUESTION 1 PART C (2/2 points)

Find the average code length assuming the Huffman code is used to encode these symbols.

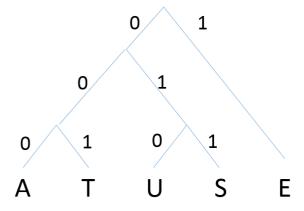
2.35



2.35

You have used 1 of 1 submissions

### SECTION 1 QUESTION 2 (2/2 points)



Suppose that the 5 symbols are encoded using the above tree, find the symbol sequence corresponding to this bit stream: 01111010000001010011001

Please input the corresponding characters below:

SEEUATUST



You have used 1 of 1 submissions

### **SECTION 1 MATLAB QUESTION** (3 points possible)

The goal of this task is to encode an English text using the Huffman code. The text, which is stored inside the variable **text\_code**, is represented using the 8-bit ASCII code, where each symbol is described by 8 bits. For example, the letter "a" is represented by the binary number "01100001", which corresponds to the decimal number 97. In an English text some characters are more likely to be used than others, hence the message can be efficiently encoded using the Huffman code. Your task is to analyze the text, find the most used symbols and create the dictionary. As in Lab 1, instead of encoding every character with the Huffman code, we will use a hybrid code: most likely characters will be represented with the Huffman code and the others by an escape code followed by their ASCII code.

The initial MATLAB code loads the input text from the file "G4jHRdGpVaY.txt" and stores it inside the variable **text\_code**. This variable is a vector of decimal integers, each lying between 0 and 255 and representing the ASCII code of the corresponding character in the input text.

The code then encodes the text using the dictionary specified by the variables **dict** and **code**. The variable **code** contains the decimal values of the ASCII codes of the N most common symbols in the text. The cell array **dict** contains N+1 elements. The first N are the binary vectors for the codewords of the corresponding symbols in **code**. The last one is for the escape code. Initially, N = 0, so the dictionary contains only the escape code. All characters in the text are encoded by the escape code [0] followed by the 8 bit ASCII code of that character. The resulting bitstream is stored in the variable **huffman**.

The code then computes the length of the ASCII and the length of the encoding using the dictionary. The length of the ASCII code is 79,632, since the number of characters in the text is 9,954. The length of the bitstream **huffman** is 89,586, since each character is encoded with nine bits (0 plus the 8 bit ASCII). Finally, the code decodes the bitstream in **huffman** using the dictionary and compares it with the input text. If the dictionary is valid, the two texts will be equal.

Your first task is to find the probabilities of each of the symbols (decimal numbers) in the vector **text\_code** (e.g. using the MATLAB function histogram). You can find the most common symbols by sorting this histogram. Find the smallest value of N, such that the N most common symbols account for more than 70% of symbols in **text\_code**. Store the N most common symbols inside the variable **code**.

Your next task is to find the Huffman dictionary that encodes these N symbols as well as the escape sequence. Store this dictionary in the cell array **dict**. Note that the probability of the escape sequence is the remainder of the probability not accounted for by the N symbols (i.e. the probability that one of the N symbols does not occur). The total probability of the N symbols and the escape sequence should sum to one.

If your dictionary is correct, then the length of the Huffman encoding should be 54,648, which is less than that of the ASCII encoding. You are graded only on the length of the dictionary (i.e. whether you find N correctly) and the dictionary itself (i.e. whether the lengths of the codewords are correct and the dictionary can be used to encode and decode the text correctly). You are free to use any of the functions used in the previous labs.

Modify the code between the lines

% % % Revise the following code % % % %

and

% % % % Do not change the code below % % % %

Please, do not change other parts of the code.



Unanswered

Size of ASCII encoding: 79632 Size of Huffman encoding: 89586

The sent and the received messages are the same.

**Run Code** 

You have used 0 of 5 submissions

© All Rights Reserved



© edX Inc. All rights reserved except where noted. EdX, Open edX and the edX and Open EdX logos are registered trademarks or trademarks of edX Inc.

















