

程序代写代做 CS编程辅导



COMP9519 Web Data Compression and Search

WeChat: cstutorcs
Assignment Project Exam Help

Recap for Compression;
Email: tutores@163.com

Preview on Search;
QQ: 749389476

Q&A for a1
<https://tutorcs.com>

Agenda for today



Where we are?

- Recap for H & AC
- LZW, Adaptive Huffman & BWT overview
- Roadmap: Compression -> Search

Assignment Project Exam Help

Other course-related matters

Email: tutorcs@163.com

- Reference papers on WebCMS3
- Q&As in Ed Forum & Consultations
- Regular exercises (<https://tutorcs.com> started this week)
- Assignment 1 spec (how to start / Q&A)

Compression



- Minimize amount of information to be stored / transferred

WeChat: cstutorcs

- Transform a sequence of characters into a new bit sequence
 - same information content (for lossless)
 - as short as possible

<https://tutorcs.com>

Run-length coding

- Run-length (encoding) is a very widely used simple compression technique
 - does not assume a memoryless source
 - replace runs of symbols (possibly of length one) with pairs of (symbol, run-length)

WeChat: cstutorcs

Assignment Project Exam Help

Email: tutorcs@163.com

QQ: 749389476

<https://tutorcs.com>



Uniquely decodable



- Uniquely decodable if no codeword is a prefix of any other

WeChat: cstutorcs

- For example $\{1, 100000, 00\}$ is uniquely decodable, but is not a prefix code
 - consider the codeword $\{\dots1000000001\dots\}$
- In practice, we prefer prefix code (why?)

Static codes



- Mapping is done before transmission
 - E.g., Huffman coding

WeChat: cstutorcs

- probabilities known in advance
 - Assignment Project Exam Help

Email: tutorcs@163.com

QQ: 749389476

<https://tutorcs.com>

Dynamic codes



- Mapping changes over time
 - i.e. adaptive coding
- Attempts to exploit locality of reference
 - periodic, frequent occurrences of messages
 - e.g., dynamic Huffman

Email: tutorcs@163.com

QQ: 749389476

<https://tutorcs.com>

Variable length coding



- Also known as entropy coding
 - The number of bits used to code symbols in the alphabet is variable
 - E.g. Huffman coding, Arithmetic coding

Assignment Project Exam Help

Email: tutorcs@163.com

QQ: 749389476

<https://tutorcs.com>

程序代写代做CS编程辅导



- What is the minimum number of bits per symbol?
- Answer: Shannon's result – theoretical minimum average number of bits per code word is known as Entropy (H)
WeChat: estutores
Assignment Project Exam Help
Email: tutorcs@163.com

QQ: 749389476

$$\sum_{i=1}^n p(s_i) \log_2 p(s_i)$$

Huffman coding algorithm



1. Take the two least probable symbols in the alphabet
(longest code words, equal length, differing in last digit)
2. Combine these two symbols into a single symbol
3. Repeat

Assignment Project Exam Help

Email: tutorcs@163.com

QQ: 749389476

<https://tutorcs.com>

Example

程序代写代做CS编程辅导

S	Freq	Huff
a	30	00
b	30	01
c	20	10
d	10	110
e	10	111

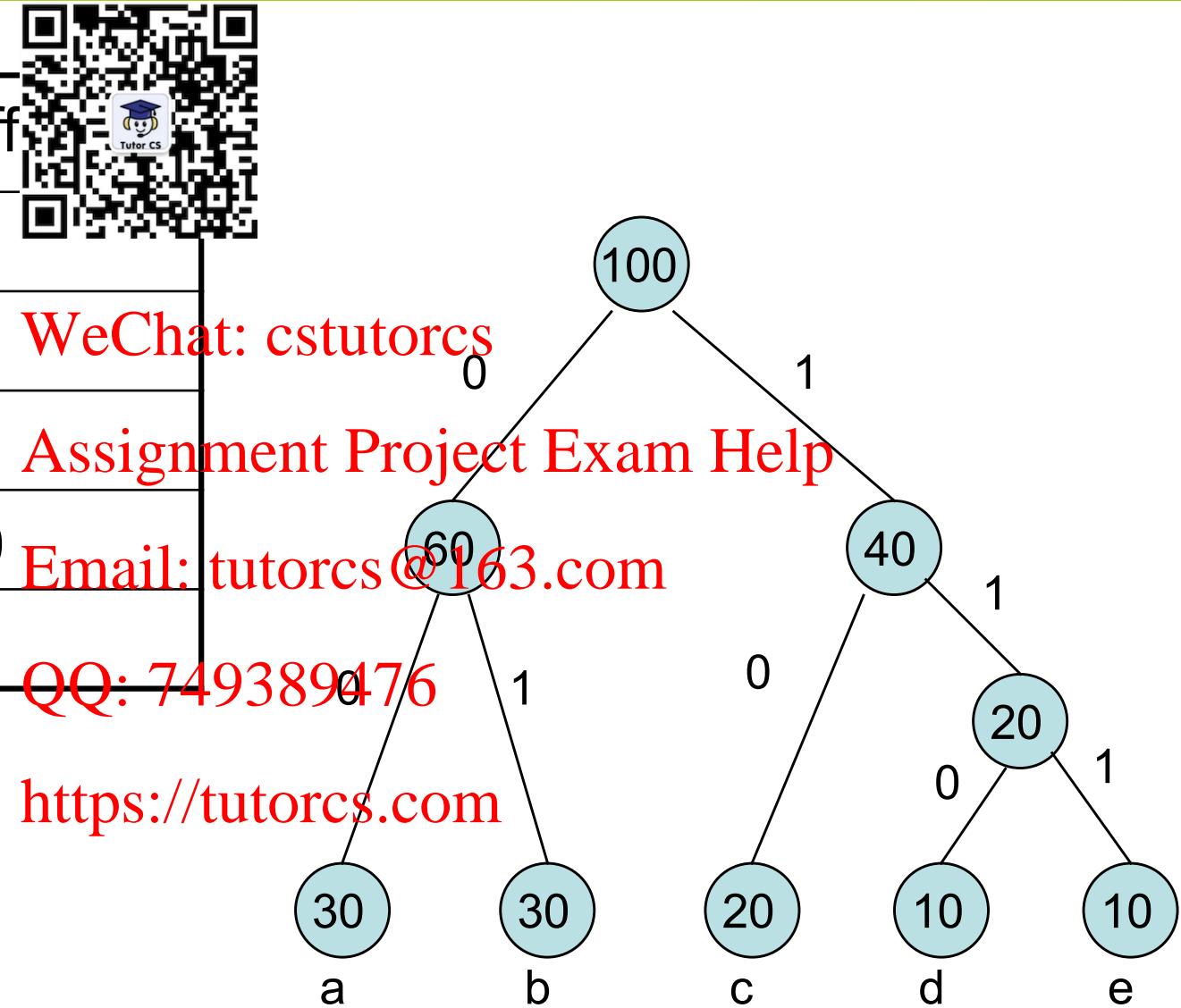


WeChat: cstutorcs
Assignment Project Exam Help

Email: tutorcs@163.com

QQ: 749389476

<https://tutorcs.com>



Another example

程序代写代做CS编程辅导

- $S = \{a, b, c, d\}$ with  $\{2, 1, 1\}$
- $H = 4/8 \log_2 2 + 2/8 \log_2 1 + 1/8 \log_2 8 + 1/8 \log_2 8$

WeChat: cstutorcs

- $H = 1/2 + 1/2 + 3/8 + 3/8 = 1.75$

Assignment Project Exam Help

- $a \Rightarrow 0 \quad b \Rightarrow 10 \quad c \Rightarrow 110 \quad d \Rightarrow 111$
- Message: {abcdabaa}  110 111 0 10 0 0}

<https://tutorcs.com>

- Average length $L = 14 \text{ bits} / 8 \text{ chars} = 1.75$
- If equal probability, i.e. fixed length, need $\log_2 4 = 2 \text{ bits}$

Problems of Huffman coding



- Huffman codes always have an integral # of bits.
 - E.g., $\log_2(3) = 1.5$ while Huffman may need 2 bits
- Noticeable non-optimality when prob of a symbol is high.

WeChat: cstutorcs

Assignment Project Exam Help

Email: tutorcs@163.com

QQ: 749389476

=> Arithmetic coding

<https://tutorcs.com>

Arithmetic coding

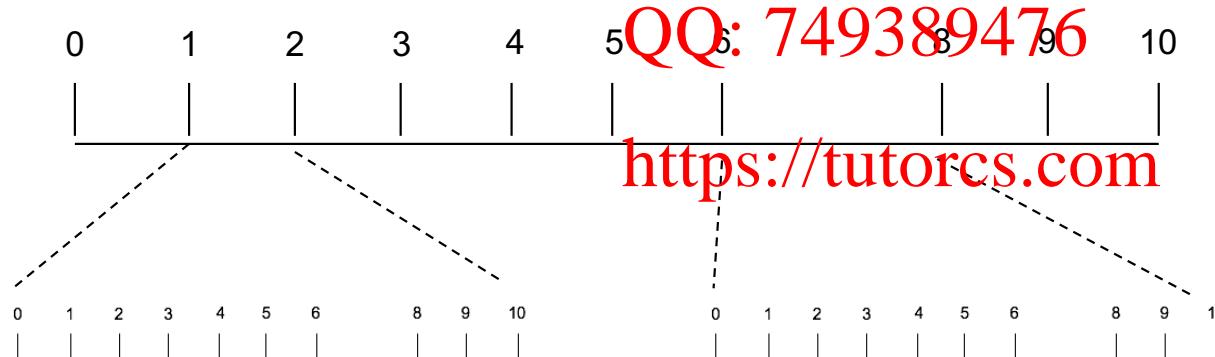
程序代写代做 CS编程辅导

Character	Probability
SPACE	1/10
A	1/10
B	1/10
E	1/10
G	1/10
I	1/10
L	2/10
S	1/10
T	1/10



WeChat: cstutorcs
Assignment Project Exam Help

Email: tutorcs@163.com



Arithmetic coding

程序代写代做 CS编程辅导

New Character	new value	High Value
B	0.2	0.3
I	0.25	0.26
L	0.256	0.258
L	0.2572	0.2576
SPACE	WeChat: cstutorcs Assignment Project Exam Help Email: tutorcs@163.com	0.25724
G	QQ: 749389476	0.257220
A	0.2572164	0.2572168
T	https://tutorcs.com	0.25721676
E	0.257216772	0.257216776
S	0.2572167752	0.2572167756



程序代写代做 CS编程辅导



COMP9519 Web Data Compression and Search

WeChat: cstutorcs
Assignment Project Exam Help

Email: tutorcs@163.com

Adaptive Huffman
QQ: 749389476

<https://tutorcs.com>

Dictionary coding



- Patterns: compressions between part of the data
- Idea: replace recurring patterns with references to dictionary
 - Assignment Project Exam Help
- LZ algorithms are adaptive:
 - Universal coding (the prob. distr. of a symbol is unknown)
 - Single pass (dictionary created on the fly)
 - No need to transmit/store dictionary

Lempel-Ziv-Welch (LZW) Algorithm



- Most popular application to LZ78
- Very common, e.g., Unix compress, TIFF, GIF, PDF (until recently)
- Read [Assignment Project Exam Help](http://en.wikipedia.org/wiki/LZW) regarding its patents
Email: tutorcs@163.com
- Fixed-length references (12bit 4096 entries)
QQ: 749389476
- Static after max entries reached
<https://tutorcs.com>

Problems of Huffman coding



Need statistics: e.g., single pass over the data to collect stat & stat unchanged during encoding

WeChat: cstutorcs

To decode, the stat table need to be transmitted. Table size can be significant for small msg.

Email: tutorcs@163.com

QQ: 749389476

=> Adaptive compression e.g., adaptive huffman

<https://tutorcs.com>

Adaptive Huffman Coding (dummy)

Encoder

Reset the stat



Repeat for each input char

(

 Encode char

 Update the stat

 Rebuild huffman tree

)

Decoder

Reset the stat

Repeat for each input char

WeChat: cstutorcs

 Decode char

 Update the stat

 Rebuild huffman tree

QQ: 749389476

<https://tutorcs.com>
This works but too slow!

Terminology (Types)



- Block-block
 - source message: fixed; codeword: fixed length
 - e.g., ASCII
- Block-variable
 - source message: fixed; codeword: variable
 - e.g., Huffman Coding
- Variable-block
 - source message: variable; codeword: fixed
 - e.g., LZW
- Variable-variable
 - source message and codeword: variable
 - e.g., Arithmetic coding

Summarised schedule



-
- 0. Information presentation (today)
 - 1. Compression
 - 2. Search
 - WeChat: cstutorcs
 - Assignment Project Exam Help
 - 3. Compression + Search on plain text
Email: tutorcs@163.com
 - 4. “Compression + Search” on Web text
QQ: [749389476](https://tutorcs.com)
 - 5. Selected advanced topics (if time allows)
<https://tutorcs.com>

程序代写代做 CS编程辅导



COMP 9519 Web Data Compression and Search

Assignment Project Exam Help

Basic BWT

Email: tutorcs@163.com

QQ: 749389476

<https://tutorcs.com>

程序代写代做 CS编程辅导



Basic BWT
WeChat: cstutorcs

(to be discussed more detailed
Assignment Project Exam Help

Email: tutorcs@163.com
next week)

QQ: 749389476

<https://tutorcs.com>

Recall from Lecture 1's RLE and 程序代写代做 CS 编程辅导 BWT example

rabcabcababa\$ bcabca**WeChat: cs_tutorcs**bcababaaa\$



aabbccacc**WeChat: cs_tutorcs**cbaaaaaaaaaabbba\$

Assignment Project Exam Help

aab4ccac3rcb**Email: tutorcs@163.com**a10b5a\$

QQ: 749389476

<https://tutorcs.com>

A simple example

程序代写与代做CS编程辅导

Input:

#BANANAS



WeChat: cstutorcs

Assignment Project Exam Help

Email: tutorcs@163.com

QQ: 749389476

<https://tutorcs.com>

All rotations



NANAS
BANANA

WeChat: cstutoros
AS#BANAN

NAS#BANA
Assignment Project Exam Help

Email: tutoros@163.com
ANAS#BAN

QQ: 749389476
NANAS#BA

<https://tutoros.com>
ANANAS#B

BANANAS#

Sort the rows
程序代写代做 CS 编程辅导



NANAS
NAS#B
ANAS#BAN

Assignment Project Exam Help

Email: tutorcs@163.com

QQ: 749389476

NANAS#BA

NAS#BANA

S#BANANA

程序代写代做CS编程辅导



NANAS
NAS#B

WeChat: cstutoros
ANAS#BAN

Assignment Project Exam Help
AS#BANAN

Email: tutorcs@163.com
BANANAS#

QQ: 749389476
NANAS#BA

<https://tutorcs.com>
NAS#BANA

S#BANANA

Exercise: you can try this example

程序代写与代做CS编程辅导

rabcabcababa\$bcabca



aabbccaccorbaaaaaaaaaaabbbba\$

WeChat: cstutorcs
Assignment Project Exam Help

Email: tutorcs@163.com

QQ: 749389476

<https://tutorcs.com>

Now the inverse, for decoding... 程序代写代做 CS 编程辅导

Input:

S
B
N
N

A
A
A



WeChat: cstutorcs

Assignment Project Exam Help

Email: tutorcs@163.com

QQ: 749389476

<https://tutorcs.com>

程序代写代做CS编程辅导
First add



WeChat: cstutorcs

S
B
N

Assignment Project Exam Help

N

A
A
A

Email: tutorcs@163.com

QQ: 749389476

<https://tutorcs.com>

程序代写代做CS编程辅导



WeChat: cstutorcs

A
A

Assignment Project Exam Help

A
B
N
N
S

Email: tutorcs@163.com

QQ: 749389476

<https://tutorcs.com>

Add again

程序代写代做CS编程辅导



WeChat: cstutorcs

S#

BA

NA

Assignment Project Exam Help

NA

Email: tutorcs@163.com

#B

QQ: 749389476

AN

<https://tutorcs.com>

AN

AS

程序代写代做CS编程辅导



WeChat: cstutorcs

#B

AN

AN

Assignment Project Exam Help

AS

Email: tutorcs@163.com

BA

QQ: 749389476

NA

<https://tutorcs.com>

NA

S#

程序代写代做CS编程辅导



WeChat: cstutorcs

Assignment Project Exam Help

Email: tutorcs@163.com

QQ: 749389476

<https://tutorcs.com>

S#B

BAN

NAN

NAS

#BA

ANA

ANA

AS#

程序代写代做CS编程辅导



WeChat: cstutorcs

Assignment Project Exam Help

Email: tutorcs@163.com

QQ: 749389476

<https://tutorcs.com>

#BA

ANA

ANA

AS#

BAN

NAN

NAS

S#B

程序代写与代做CS编程辅导



WeChat: cstutorcs

Assignment Project Exam Help

Email: tutorcs@163.com

QQ: 749389476

<https://tutorcs.com>

S#BA

BANA

NANA

NAS#

#BAN

ANAN

ANAS

AS#B

程序代写代做CS编程辅导



WeChat: cstutorcs

Assignment Project Exam Help

Email: tutorcs@163.com

QQ: 749389476

<https://tutorcs.com>

#BAN

ANAN

ANAS

AS#B

BANA

NANA

NAS#

S#BA

程序代写与代做CS编程辅导



WeChat: cstutorcs **NANAS**
Assignment Project Exam Help
Email: tutorcs@163.com
QQ: 749389476 **#BANA**
<https://tutorcs.com> **ANANA**
ANAS#
AS#BA

程序代写代做CS编程辅导



WeChat: cstutorcs **#BANA**
Assignment Project Exam Help **ANANA**
Email: tutorcs@163.com **ANAS#**
QQ: 749389476 **AS#BA**
<https://tutorcs.com> **BANAN**
NANAS **NAS#B**
S#BAN

程序代写代做CS编程辅导
Then add



S#BANA
BANANA

WeChat: cstutorcs
Assignment Project Exam Help
Email: tutorcs@163.com

QQ: 749389476
<https://tutorcs.com>

#BANAN
ANANAS
ANAS#B
AS#BAN

程序代写代做CS编程辅导



#BANAN
ANANAS
ANAS#B

WeChat: cstutorcs
Assignment Project Exam Help

Email: tutorcs@163.com
QQ: 749389476
https://tutorcs.com

BANANA
NANAS#
NAS#BA
S#BANA

程序代写代做CS编程辅导
Then add



S#BANAN
BANANAS
NANAS#B
WeChat: cstutorcs
NAS#BAN
Assignment Project Exam Help
#BANANA
Email: tutorcs@163.com
ANANAS#
QQ: 749389476
ANAS#BA
<https://tutorcs.com>
AS#BANA

程序代写代做CS编程辅导



#BANANA
ANANAS#

WeChat: cstutorcs
Assignment Project Exam Help
AS#BANA

BANANAS
Email: tutorcs@163.com
NANAS#B
QQ: 749389476
NAS#BAN
<https://tutorcs.com>
S#BANAN

程序代写代做CS编程辅导
Then add



S#BANANA
BANANAS#
NANAS#BA
NAS#BANA
Assignment Project Exam Help
#BANANAS
Email: tutorcs@163.com
ANANAS#B
QQ: 749389476
ANAS#BAN
<https://tutorcs.com>
AS#BANAN

The n sort (???)



#BANANAS
ANANAS#B

WeChat: cstutors
Assignment Project Exam Help
Email: tutorcs@163.com
QQ: 749389476
<https://tutorcs.com>

Exercise: you can try this example



rabcabcababa\$bcabca

aabbccaccorbaaaaaaaaaabbba\$

WeChat: cstutorcs
Assignment Project Exam Help

Email: tutorcs@163.com

QQ: 749389476

<https://tutorcs.com>

Reference Papers on WebCMS3

1098



TINGS OF THE I.R.E.

September

A Method for the Construction of Minimum Redundancy Codes*

DAVID A. HUFFMAN*, ASSOCIATE, IRE

Summary—An optimum method of coding an ensemble of messages consisting of a finite number of members is developed. A minimum-redundancy code is one constructed in such a way that the average number of coding digits per message is minimized.

INTRODUCTION

ONE IMPORTANT METHOD of transmitting messages is to transmit in their place sequences of symbols. If there are more messages which might be sent than there are kinds of symbols available, then some of the messages must use more than one symbol. If it is assumed that each symbol requires the same time for transmission, then the time for transmission (length) of a message is directly proportional to the number of symbols associated with it. In this paper, the symbol or sequence of symbols associated with a given message will be called the "message code." The entire number of messages which might be transmitted will be

will be defined here as an ensemble code which, for a message ensemble consisting of a finite number of members, N , and for a given number of coding digits, D , yields the lowest possible average message length. In order to avoid the use of the lengthy term "minimum-redundancy," this term will be replaced here by "optimum." It will be understood then that, in this paper, "optimum code" means "minimum-redundancy code."

The following basic restrictions will be imposed on an ensemble code:

- No two messages will consist of identical arrangements of coding digits.
- The message codes will be constructed in such a way that no additional indication is necessary to specify where a message code begins and ends once the starting point of a sequence of messages is known.

QQ: 749389476

<https://tutorcs.com>

Q&As in Ed Forum

Hey all,



I just finished the arithmetic coding [1] video and was wondering why we need to worry about using different probabilities for encoding?

E.g. if we knew we were just encoding English alphabet letters and spaces, is there a strong downside to just using a $\frac{1}{27}$ split of the $[0, 1]$ range and running the algorithm? I assume it has something to do with switching from pure real numbers to binary representations that the probabilities come in to play, but just wanted to ask in case there is a different explanation as well?

WeChat: cstutorcs
Assignment Project Exam Help

Thanks!

Email: tutorcs@163.com

This is my somewhat limited understanding from reading parts of the AC paper:

QQ: 749389476

If you allocate a *smaller* range for a symbol, then you must transmit *more* bits in order to encode that symbol. This is because a smaller (i.e. narrower) range requires more decimal places to represent a number that lies in that range. More decimal places \Rightarrow longer binary representation.

For the English language, you would be using unnecessary extra bits to encode those characters that appear more frequently than others (vowels, for example).

Q&As from Consultations



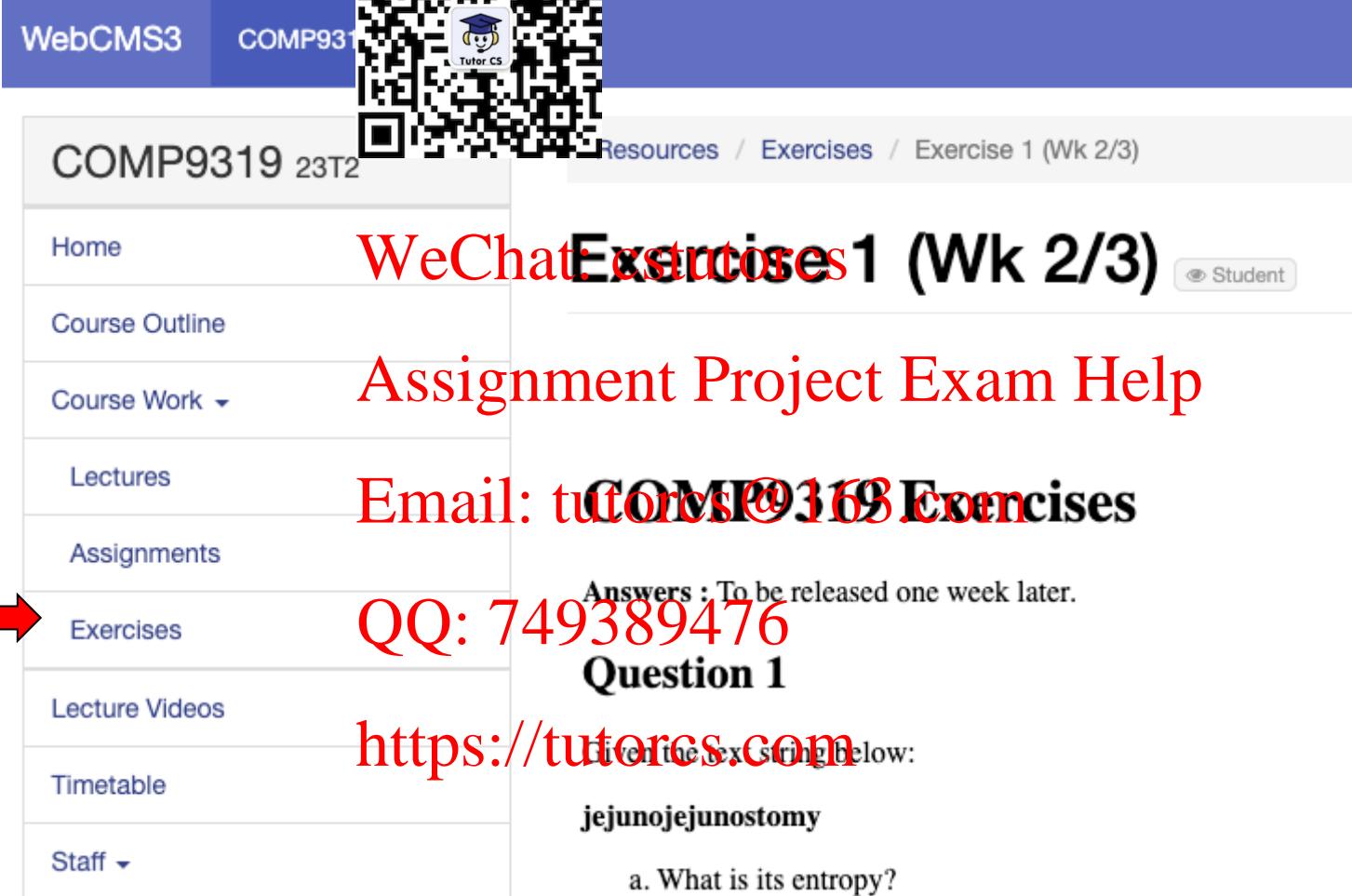
For Huffman, AC, LZW we covered so far, what if
we consider source messages in UTF8 instead of
ASCII? WeChat: cstutorcs Assignment Project Exam Help

Email: tutorcs@163.com

QQ: 749389476

<https://tutorcs.com>

Exercises on WebCMS3



WebCMS3 COMP9319

COMP9319 23T2 Resources / Exercises / Exercise 1 (Wk 2/3)

Home Course Outline Course Work Lectures Assignments Exercises Lecture Videos Timetable Staff

WeChat: cstutors Student

Exercise 1 (Wk 2/3)

Assignment Project Exam Help

Email: **COMP9319.Exercises**

Answers : To be released one week later.

QQ: **749389476**

Question 1

Given the hex string below:

jejunojejunostomy

a. What is its entropy?

A red arrow points to the 'Exercises' link in the sidebar menu.

Assignment¹

程序代写代做CS编程辅导

COMP9319 2023T2 Assignment 1: LZW Encoding and Decoding



Your task in this assignment is to implement an LZW encoder and its decoder with 15-bit 32768 dictionary entries (excluding those entries for the individual ASCII characters), called `lencode` and `ldecode`, in C or C++. After the dictionary is full, no new entries can be added. You may assume the source file may contain any valid ASCII characters.

WeChat: cstutorc

```
%grieg> lencode ~cs9319/a1/test1.txt test1.encoded  
%grieg> ldecode test1.encoded test1.decoded  
%grieg> diff ~cs9319/a1/test1.txt test1.decoded  
%grieg>
```

Assignment Project Exam Help

Email: tutorcs@163.com

test1.lzw using xxd:

```
cs9319@grieg:~/a1$ xxd -b test1.txt  
00000000: 01011110 01010111 01000101 01000100 01011110 01010111 ^WED^W  
00000006: 0000101 0001101 0101111 01000101 01000101 01011110 E^WEE^  
0000000c: 01010111 01000101 01000010 01011110 01010111 01000101 WEB^WE  
00000012: 01010100 T  
cs9319@grieg:~/a1$ xxd -b test1.lzw  
00000000: 01011110 01010111 01000101 01000100 01011110 01010111 ^WED^W  
00000006: 01000101 00000000 00000100 01000101 01011110 01010111 E..E^W  
0000000c: 01000101 01000010 10000000 00000100 01010100 EB..T  
cs9319@grieg:~/a1$
```

Hint: Setting MSB to 1



Unsigned_T val ;

val = ... WeChat: cstutorcs

Assignment Project Exam Help

msb = ((*Unsigned_T*) >> 1) + 1;

QQ: 749389476

val |= msb;

Important: if you use your PC to

程序代写代做 CS编程辅导

code a1



Make sure you reserve time to:

- port & compile
- test & debug

Assignment Project Exam Help
on grieg.cse.unsw.edu.au before you submit.

Email: tutorcs@163.com

QQ: 749389476

<https://tutorcs.com>