

Lecture 3 - Compression

PORLAND STATE
UNIVERSITY

Administrative

- Programming assignment on web site.
 - ◊ Due Monday, October 19, 2015
 - ◊ Start early
- HW 1 – Due on Wednesday
- QUIZ #1 – Next Wednesday
 - ◊ October 14, 2015

PORLAND STATE
UNIVERSITY

Huffman encoding

Optimal encoding one symbol \rightarrow one codeword

Huffman codebook required

- optimized for single instance
 -) have to include in compressed file
- pick generic codebook (e.g. English dictionary)
 - +) no codebook in file Huffman table
 -) not quite as small

PORLAND STATE
UNIVERSITY

Run Length Encoding

Sequence of data + encode it as

$\langle \text{symbol}, \text{run length} \rangle$

- applicable to repeating data

$\underbrace{1111}_{300} \quad \underbrace{222}_{20} \quad \underbrace{66}_{100}$ } Need to be careful w/ how to encode

$\langle 1, 300 \rangle \langle 2, 20 \rangle \langle 6, 100 \rangle$

- Effective for graphical images
- Limited use in text files
- JPEG uses this as part of its lossless encoding \rightarrow entropy encoding process

PORLAND STATE
UNIVERSITY

Compression Fundamental

- CF3 – Escape sequences are necessary for control and disambiguation of data

AAAAA...A B ...3 C ...C
9 12 7
A9B12C7

Create special characters or sequences that
never occur in regular data

Have a special meaning

Example printf

printf("\\\\");
" "

printf("\\\\");
" "

PORLAND STATE
UNIVERSITY

Run Length Encoding Variations

Variation 1

AABB CC DD EEEE FFFF

A2B2C2D2E4F5

- wait for two symbols, then use RLE

AA1 BB1 CC1 DD1 EE3 FF4

Variation 2

- Use escape sequence to indicate run

AB EEEE DDC BBBB# escape char
A B E*4 D*2 C B*6

Variation 3

Using runs to skip a single common symbol

8 0002 4 005 000003

8 ([4], 2) ([1], 4) ([3], 5) ([1], 3)
skipped places

JPEG
uses
this form

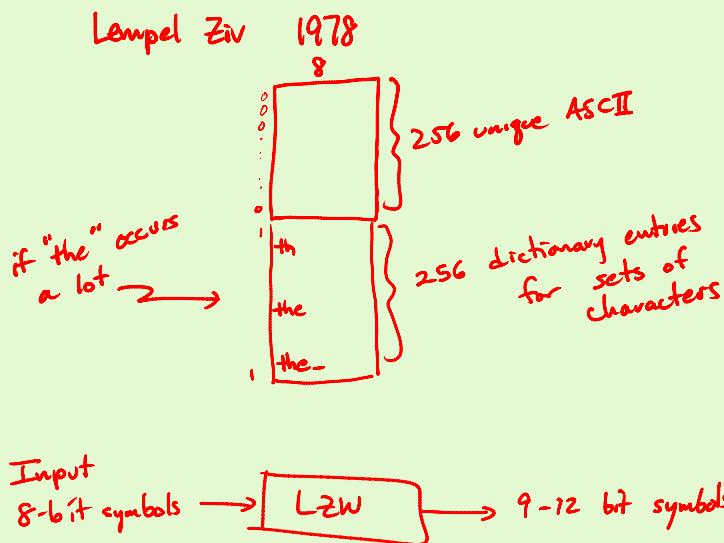
PORLAND STATE
UNIVERSITY

Adaptive Compression

- Huffman compression takes one symbol and generates variable length codeword
 - ◊ Decompression complicated because of variable bit processing
- Alternative: Learn about common sequences that occur in the input stream and generate one codeword
 - common sequences in English, for example
the
and

PORLAND STATE
UNIVERSITY

LZW Compression



PORLAND STATE
UNIVERSITY

LZW Compression – summary

(accidentally skipped but was discussed)

- Lempel-Ziv Compression
 - ◊ convert variable length strings into fixed length codes
 - ◊ builds “table” during both compression and decompression
 - +) No table
 -) Needs conditioning
- Lempel-Ziv-Welch (LZW) Compression - is essentially a fast implementation of the LZ compression method
 - ◊ Uses:
 - UNIX compress and gzip use LZ and LZW
 - v.42bis network data compression
 - gif and tiff image compression format

PORLAND STATE
UNIVERSITY

LZW Compression

LZW Compression

$w = nil$

while (read char k)

 if wk exists in dictionary

$w = wk$

 else

 add wk to dictionary

 output code for w

$w = k$

PORLAND STATE
UNIVERSITY

Example: ABCDABCABBABBDABC

w	k	dict	output	
A	A	① AB <256>	A ①	
B	C	BC <257>	B	
C	D	CD <258>	C	
D	A	DA <259>	D	
A	B	② ③ dict		
AB	C	ABC <260>	<256>	
C	A	CA <261>	C	
A	B	④ dict		
AB	B	ABB <262>	<256>	
B	A	BA <263>	B	
A	B	⑤		
AB	B	⑥		
ABD	D	ABBD <264>	<262>	

8 slots

 slot 1: A

 slot 2: B

 slot 3: C

 slot 4: D

 slots 5-8: dict

9-bit output symbols

 ABC...D padded w/ 0 @ front

PORLAND STATE
UNIVERSITY

LZW Decompression

LZW Compression

Read k

Output k

$w = k$

while (read char k)

 entry = dictionary entry for k

 output entry

 add $w + entry[0]$ to dictionary

$w = entry$

PORLAND STATE
UNIVERSITY

A B C <257> D <258> E <257>

ENTRY

W	k	k-dict	dictionary	output
A	A ①	B ⑤	AB <256> ⑦	A ③
B	B ②	C ⑥	BC <257>	B ④
C	C ③	BC ⑦	CB <258> ⑧	C
BC	D ④	D ⑨	BCD <259>	BC
D	D ⑤	CB ⑩	DC <260>	D
CB	E ⑥	E ⑪	CBE <261>	CB
E	E ⑦	EB ⑫	EB <262>	E

PORLAND STATE
UNIVERSITY

KwKwK String

K is a single symbol
w is a dictionary entry or single symbol

Kw exists in dictionary before KwKwK appears

w + entry(o) → dictionary

w + w(o) → dictionary

PORLAND STATE
UNIVERSITY