

Dictionary-based Coding Techniques



National Chiao Tung University
Chun-Jen Tsai
10/16/2014

Rationale

- ❑ In previous two chapters, we looked at coding techniques that assume a source that generates a sequence of independent symbols.
 - Most data sources are correlated, thus, the coding step is generally preceded by a de-correlation step (i.e. model prediction).
- ❑ Alternatively, we can build a list of commonly occurring patterns and encode these patterns by transmitting their index in the list
→ dictionary techniques

Static vs. Adaptive Dictionary

- ❑ The dictionary holds a list of strings of symbols and it may be static or dynamic (adaptive)
- ❑ Static dictionary – permanent, sometimes allowing the addition of strings but no deletions
- ❑ Dynamic dictionary – holding strings previously found in the input stream, allowing for additions and deletions of strings as new input symbols are being read

Basic Idea of Dictionary Coding

- ❑ Given an input source, we want to
 - Identify frequent symbol patterns
 - Encode those more efficiently
 - Use a default (less efficient) encoding for the rest
 - Hopefully, the average bits per symbol gets smaller
- ❑ In general, dictionary-based techniques works well for highly correlated data (e.g. text), but less efficient for data with low correlation (e.g. i.i.d. sources)

Motivating Example

- ❑ Consider an 'English' source with 26 letters & six punctuation marks
 - Single-symbol FLC, fixed-length encoding: 5 bps
 - Four-symbol FLC, fixed-length encoding: 20 bps (32^4)
- ❑ If we assume uneven distribution of the symbols
 - Pick a dictionary which contains the 256 most-frequent patterns (probability p) and encode them with 8 bits
 - Encode the rest with 20 bits
 - Use 1-bit prefix to distinguish the two cases

then, the average rate is $9p + 21(1 - p) = 21 - 12p$.

If $p > 0.084$, $21 - 12p < 20$.

Static Dictionary

- ❑ Using a static dictionary is less complex, but the probability p of a hit highly depends on the applications
 - For student records in a university is probably ok.
- ❑ The key for success is that the most common patterns are a small subset of all possible messages
 - Out of over 100,000 English words, only less than 2,000 words are used in most writings

Digram Coding

- ❑ The dictionary is composed of
 - All letters from the alphabet
 - As many digrams (pairs of letters) as possible

- ❑ For example, if we want to encode pure ASCII text documents, we can design a dictionary of size 256 entries, and
 - Source alphabet: 95 printable ASCII symbols
 - Digrams: 161 most common pairs

Simple Digram Coding Example

- ❑ The source alphabet $A = \{a, b, c, d, r\}$
- ❑ Dictionary:

Code	Entry	Code	Entry
000	<i>a</i>	100	<i>r</i>
001	<i>b</i>	101	<i>ab</i>
010	<i>c</i>	110	<i>ac</i>
011	<i>d</i>	111	<i>ad</i>

- ❑ Try to code the sequence *abracadabra*, the output is 101100110111101100000.

Problem: Which Digrams to Use?

❑ Source 1: LaTeX documents

Pair	Count	Pair	Count
<i>e\p</i>	1128	<i>ar</i>	314
<i>\p t</i>	838	<i>at</i>	313
<i>\p \p</i>	823	<i>\p w</i>	309
<i>th</i>	817	<i>te</i>	296
<i>he</i>	712	<i>\p s</i>	295
<i>in</i>	512	<i>d\p</i>	272
<i>s\p</i>	494	<i>\p o</i>	266
<i>er</i>	433	<i>io</i>	257
<i>\p a</i>	425	<i>co</i>	256
<i>t\p</i>	401	<i>re</i>	247
<i>en</i>	392	<i>\p \$</i>	246
<i>on</i>	385	<i>r\p</i>	239
<i>n\p</i>	353	<i>di</i>	230
<i>ti</i>	322	<i>ic</i>	229
<i>\p i</i>	317	<i>ct</i>	226

❑ Source 2: C programs

Pair	Count	Pair	Count
<i>\p \p</i>	5728	<i>st</i>	442
<i>nl\p</i>	1471	<i>le</i>	440
<i>;\p nl</i>	1133	<i>ut</i>	440
<i>in</i>	985	<i>f(</i>	416
<i>nt</i>	739	<i>ar</i>	381
<i>= \p</i>	687	<i>or</i>	374
<i>\p i</i>	662	<i>r\p</i>	373
<i>t\p</i>	615	<i>en</i>	371
<i>\p =</i>	612	<i>er</i>	358
<i>);</i>	558	<i>ri</i>	357
<i>, \p</i>	554	<i>at</i>	352
<i>nl\p nl</i>	506	<i>pr</i>	351
<i>\p f</i>	505	<i>te</i>	349
<i>e\p</i>	500	<i>an</i>	348
<i>\p *</i>	444	<i>lo</i>	347

Adaptive Dictionary Technique

- ❑ Original ideas published by Jacob Ziv and Abraham Lempel in 1977 (LZ77/LZ1) and 1978 (LZ78/LZ2)
- ❑ The most well-known dictionary-based technique, LZW, is a modification to LZ algorithms published by Terry Welch in 1984

LZ77 (1/2)

□ General approach

- Dictionary is a portion of the previously encoded sequence
- Use a sliding window for compression

□ Mechanism

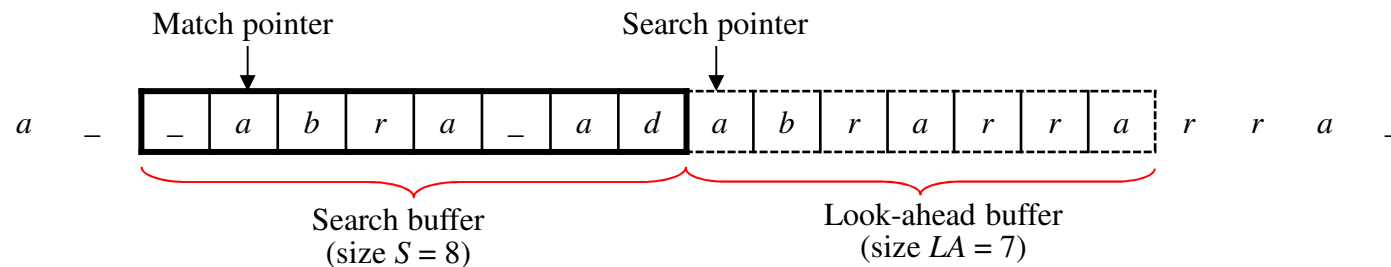
- Find the maximum length match for the string pointed to by the search pointer in the search buffer, and encode it

□ Rationale

- If patterns tend to repeat locally, we should be able to get more efficient representation

LZ77 (2/2)

- ❑ Sliding window is composed of a search buffer and a look-ahead buffer (note: window size $W = S + LA$)



- ❑ Offset = search pointer – match pointer ($o = 7$)
- ❑ Length of match = number of consecutive letters matched ($l = 4$)
- ❑ Codeword ($c = C(r)$), where $C(x)$ is the codeword for x
- ❑ Encoding triple: $\langle o, l, c \rangle = \langle 7, 4, C(r) \rangle$
 - If FLC is used and alphabet size is $|A|$, $\langle o, l, c \rangle$ can be encoded with $\lceil \log_2 S \rceil + \lceil \log_2 W \rceil + \lceil \log_2 |A| \rceil$ bits.

Possible Cases for Triples

- ❑ There could be three different possibilities that may be encountered during the coding process:
 - No match for the next character to be encoded in the window
 - There is a match
 - The matched string extends inside the look-ahead buffer
- ❑ For each of these cases, we have a triple to signal the case to the decoder

LZ77 Encoding Example

- ❑ Sequence
 - *cabracadabrarrarrad*
 - $W = 13, S = 7$
- ❑ |*cabraca*|*dabrarr*|*rarrad*
 - no match for *d*
 - send $\langle 0, 0, C(d) \rangle$
- ❑ |*abracad*|*abrarr*|*rarrad*
|*abracad*|*abrarr*|*rarrad*
|*abracad*|*abrarr*|*rarrad*
|*abracad*|*abrarr*|*rarrad*
 - send $\langle 7, 4, C(r) \rangle$
- ❑ |*cadabrar*|*rarrad*|
|*cadabrar*|*rarrad*|
|*cadabrar*|*rarrad*|
 - send $\langle 3, 3, C(r) \rangle$
- ❑ Could we do better?
 - send $\langle 3, 5, C(d) \rangle$ instead

LZ77 Decoding Example

- ❑ Current input: $\langle 0, 0, C(d) \rangle \langle 7, 4, C(r) \rangle \langle 3, 5, C(d) \rangle$
- ❑ Current output: *cabraca*
- ❑ Decode: $\langle 0, 0, C(d) \rangle$
 - Decode $C(d)$: *cabracad*l
- ❑ Decode: $\langle 7, 4, C(r) \rangle$
 - Start with the first 'a', copy four letters: *cabracad*abral
 - Decode $C(r)$: *cabracad*abrarl
- ❑ Decode: $\langle 3, 5, C(d) \rangle$
 - Start with the first 'r', copy three letters: *cabracadab*rrarl
 - Copy two more letters: *cabracadabr*rrarl
 - Decode $C(d)$: *cabracadabr*rrarard

LZ77 Variants

□ For LZ77, we have

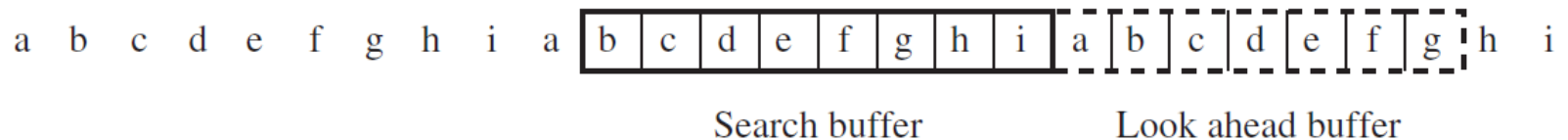
- Adaptive scheme, no prior knowledge
- Asymptotically approaches the source statistics
- Assumes that recurring patterns close to each others

□ Possible improvements

- Variable-bit encoding: PKZip, zip, gzip, ..., etc., uses a variable-length coder to encode $\langle o, l, c \rangle$.
- Variable buffer size: larger buffer requires faster searches
- Elimination of $\langle 0, 0, C(x) \rangle$
 - LZSS sends a flag bit to signal whether the next “token” is an $\langle o, l \rangle$ pair or the codeword of a symbol

Problems with LZ77

- ❑ If the recurring patterns happens with a period larger than the search window, the performance is bad
- ❑ Example:



LZ78

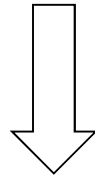
- ❑ LZ78 improvements from LZ77
 - No search buffer – explicit dictionary instead
 - Encoder/decoder must build dictionary in sync
 - Encoding: $\langle i, c \rangle$
 - i = index in the dictionary, $i = 0$ for symbols not in the dictionary
 - c = code of the following character
- ❑ Example: encode the following contents
 - *wabbabwabbabwabbabwabbabwoobwoobwoo*

LZ78 Example

- Input: *wabbawabbawabbawabbawooowooowoo*
- Dictionaries:

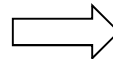
initial dictionary (empty)

Index	Entry



dictionary after encoding *w, a, b*

Encoder Output	Index	Entry
$\langle 0, C(w) \rangle$	01	<i>w</i>
$\langle 0, C(a) \rangle$	02	<i>a</i>
$\langle 0, C(b) \rangle$	03	<i>b</i>



final dictionary

Encoder Output	Index	Dictionary Entry
$\langle 0, C(w) \rangle$	01	<i>w</i>
$\langle 0, C(a) \rangle$	02	<i>a</i>
$\langle 0, C(b) \rangle$	03	<i>b</i>
$\langle 3, C(a) \rangle$	04	<i>ba</i>
$\langle 0, C(\emptyset) \rangle$	05	\emptyset
$\langle 1, C(a) \rangle$	06	<i>wa</i>
$\langle 3, C(b) \rangle$	07	<i>bb</i>
$\langle 2, C(\emptyset) \rangle$	08	<i>a\emptyset</i>
$\langle 6, C(b) \rangle$	09	<i>wab</i>
$\langle 4, C(\emptyset) \rangle$	10	<i>ba\emptyset</i>
$\langle 9, C(b) \rangle$	11	<i>wabb</i>
$\langle 8, C(w) \rangle$	12	<i>a\emptyset w</i>
$\langle 0, C(o) \rangle$	13	<i>o</i>
$\langle 13, C(\emptyset) \rangle$	14	<i>o\emptyset</i>
$\langle 1, C(o) \rangle$	15	<i>wo</i>
$\langle 14, C(w) \rangle$	16	<i>o\emptyset w</i>
$\langle 13, C(o) \rangle$	17	<i>oo</i>

Remarks on LZ78

- ❑ Observation
 - If we keep on encoding, the dictionary will keep on growing
- ❑ Possible solutions
 - Stop growing the dictionary
 - Effectively switch to a static dictionary
 - Prune it
 - Based on usage statistics
 - Reset it
 - Start all over again
- ❑ The best solution depends on the knowledge of the source

LZ78 Variants: LZW

- ❑ Invented by Terry Welch in 1984

- ❑ Idea

- Instead of $\langle i, c \rangle$, encode i only

- ❑ Algorithm

- Initial dictionary contains all alphabet letters, $p = \text{null}$

```
while (!done)
  read next symbol into  $a$ 
  if ( $p*a$ ) is in the dictionary // Note: '*' stands for concatenation
     $p = p*a$ 
  else
    send out index of  $p$ 
    add  $p*a$  to the dictionary
     $p = a$ 
end
```

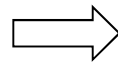
Example: LZW Encoding

❑ Input: *wabbawabbawabbawabbawooowooowoo*

❑ Dictionaries:

initial dictionary (source alphabet)

Index	Entry
1	<i>ϕ</i>
2	<i>a</i>
3	<i>b</i>
4	<i>o</i>
5	<i>w</i>



final dictionary

Index	Entry	Index	Entry
01	<i>ϕ</i>	14	<i>aϕ w</i>
02	<i>a</i>	15	<i>wabb</i>
03	<i>b</i>	16	<i>baϕ</i>
04	<i>o</i>	17	<i>ϕ wa</i>
05	<i>w</i>	18	<i>abb</i>
06	<i>wa</i>	19	<i>baϕ w</i>
07	<i>ab</i>	20	<i>wo</i>
08	<i>bb</i>	21	<i>oo</i>
09	<i>ba</i>	22	<i>oϕ</i>
10	<i>aϕ</i>	23	<i>ϕ wo</i>
11	<i>ϕ w</i>	24	<i>ooϕ</i>
12	<i>wab</i>	25	<i>ϕ woo</i>
13	<i>bba</i>		

❑ Output: 5 2 3 3 2 1 6 8 10 12 9 11 7 16 5 4 4 11 21 23 4

Problems with LZW Decoding

- ❑ Decoding of LZW is simple, in general
 - Output symbols from the dictionary as indexed by the inputs
 - Construct the dictionary on-the-fly as the encoder does

- ❑ However, if we have a message pattern $cScS \dots$, where c is a character, S is a string, we may run into a situation that the indexed entry is in partial construction

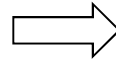
- ❑ Solution: the current dictionary entry under construction is in p , we should allow reading partial data out of p during decoding

Example: Special Case in Decoding

- ❑ Alphabet $A = \{a, b\}$, input is $abababab$, encoder output is 1235
- ❑ Decoding dictionaries:

initial dictionary

Index	Entry
1	a
2	b



intermediate dictionary

Index	Entry
1	a
2	b
3	ab
4	ba
5	$a \dots$

when we reach decoding of 5, $p = ab???$, we do not have the complete output!

Application: Compress

- ❑ An early implementation of LZW
- ❑ Adaptive dictionary, starts with 2^9 entries
- ❑ User can configure max codeword length $b_{max} = 9\sim 16$
- ❑ Dictionary grows up to double in size
 - When dictionary reaches $2^{b_{max}}$ entries, it becomes a static dictionary encoder
- ❑ If compression ratio falls below a threshold, dictionary is reset

Application: GIF Images

- ❑ LZW scheme, similar to compress:
- ❑ Clear code is used to reset the encoder/decoder. For b bits/pixel images, 2^b is used as the clear code
- ❑ Dictionary size is initially 2^{b+1}
- ❑ Dictionary size can grows up to 4096 entries
- ❑ Format:
 - Codewords stored in blocks of 8-bit characters
 - Each block begins with a header with a size count up to 255, and ends with a block terminator symbol (8 zero bits)
 - The last block has a end-of-information code, $2^b + 1$, before the block terminator

GIF Performance

□ GIF vs. arithmetic coding

Image	GIF	Arithmetic Coding of Pixel Values	Arithmetic Coding of Pixel Differences
Sena	51,085	53,431	31,847
Sensin	60,649	58,306	37,126
Earth	34,276	38,248	32,137
Omaha	61,580	56,061	51,393

Application: PNG Images

- ❑ Based on LZ77, patent-free alternative to GIF
- ❑ Designed specifically for lossless image compression
- ❑ Modes: true color, grayscale, 8-bit palette
- ❑ Two autonomous compression components
 - Deflate (RFC 1951) — LZ77-style dictionary compression algorithm plus Huffman coding
 - Filtering — lossless transformations of byte-level image data

PNG – Deflate

- ❑ Deflate = LZ77 + Huffman
- ❑ Three types of data blocks
 - Uncompressed, LZ77 + fixed Huffman, LZ77 + adaptive Huffman
- ❑ Match length is between 3 and 258 bytes
 - A sliding window of at least 3-byte long is examined
 - If match is not found, encode the first byte and slide window
 - At each step, LZ77 either outputs a codeword for a literal or a paired value of <match_length, offset>
 - Match length is encoded by index code (257~285) and a selector code (0~5 bits)
 - Offset (1~32768) is encoded using Huffman code

PNG – Filtering

- ❑ Filters are applied on a scanline-by-scanline basis
- ❑ All algorithms applied to bytes (not pixels)
- ❑ Filter types:
 - None: unmodified value
 - Sub: difference from previous byte value (mod 256)
 - Up: difference from the byte value above
 - Average: subtract average of the left and the above bytes
 - Paeth:
 - Compute initial estimate by $\text{left} + \text{above} - \text{upper_left}$
 - The value of left, above, or upper_left that is closest to the initial estimate is used as the estimate

PNG: Performance

□ PNG vs. GIF vs. arithmetic coding

Image	PNG	GIF	Arithmetic Coding of Pixel Values	Arithmetic Coding of Pixel Differences
Sena	31,577	51,085	53,431	31,847
Sensin	34,488	60,649	58,306	37,126
Earth	26,995	34,276	38,248	32,137
Omaha	50,185	61,580	56,061	51,393