

# CSC 310: Information Theory

University of Toronto, Fall 2011

Instructor: Radford M. Neal

Week 7

# Solving the Dilemma of What Order Markov Model to Use

We would like to get both:

- the advantage of fast learning of a low-order model
- the advantage of ultimately better prediction of a high-order model

We can do this by *varying* the order we use.

One scheme for this is the “prediction by partial match” (PPM) model.

# Contexts Used by PPM

PPM maintains frequencies for characters that have been seen before in *all* contexts that have occurred before, up to some maximum order.

Suppose we have so far encoded the string

`this_is_th`

If we are using contexts up to order two, then we will record frequencies for the following contexts:

Order 0: `()`

Order 1: `(t) (h) (i) (s) (_)`

Order 2: `(th) (hi) (is) (s_) (_i) (_t)`

## “Escaping” From a Context

The frequency tables maintained by PPM contain *only* the characters that have been seen before in that context.

Examples: if “x” has never occurred, none of the frequency tables will have an entry for “x”. If “x” *has* occurred before, but *not* after a “t”, the frequency table for order 1 context (t) will not contain “x”.

**The main idea:** If we need to encode a character that doesn’t appear in the context we’re using, we transmit an “escape” flag, and switch to a lower-order context.

What if we escape from every context? We end up in a special “order -1” context, in which every character has a frequency of 1.

# Frequencies in Contexts

Two details about frequencies need to be resolved.

First, what characters do we count in a context?

- We might count *every* character that appears following the characters making up the context.
- We might count a character in a context *only* when it does not appear in a higher-order context.

One could argue for either way, but we'll go for the second option.

Second, what do we use as the frequency of the “escape” symbol? There are many possibilities. We'll just give it a frequency of one.

# Basic PPM Encoding Method

Loop until end of file:

Read the next character,  $c$ . Let  $d_K, \dots, d_1$  be the preceding  $K$  characters.

Set the context size,  $k$ , to the maximum,  $K$ .

While  $(d_k, \dots, d_1)$  hasn't been seen previously:

Set  $k$  to  $k - 1$ .

While  $k \geq 0$  and  $c$  hasn't been seen in context  $(d_k, \dots, d_1)$ :

Transmit an escape flag using context  $(d_k, \dots, d_1)$ .

Set  $k$  to  $k - 1$ .

If  $k = -1$ :

Transmit  $c$  using the special “order -1” context.

Set  $k$  to 0.

Else

Transmit  $c$  using context  $(d_k, \dots, d_1)$ .

While  $k \leq K$ :

Create context  $(d_k, \dots, d_1)$  if it doesn't exist.

Increment the count for  $c$  in context  $(d_k, \dots, d_1)$ .

Set  $k$  to  $k + 1$ .

# Frequencies After Encoding `this_is_th`

Order -1: `_:1 a:1 b:1 ... z:1`

Order 0:

`() Escape:1 t:2 h:1 i:2 s:1 _:1`

Order 1:

`(t) Escape:1 h:2`

`(h) Escape:1 i:1`

`(i) Escape:1 s:2`

`(s) Escape:1 _:1`

`(_) Escape:1 i:1 t:1`

Order 2:

`(th) Escape:1 i:1`

`(hi) Escape:1 s:1`

`(is) Escape:1 _:2`

`(s_) Escape:1 i:1 t:1`

`(_i) Escape:1 s:1`

`(_t) Escape:1 h:1`

# Learning a Vocabulary

One reason PPM works well for files like English text is that it can implicitly learn the vocabulary — the dictionary of words in the language. This is because early letters of a word like “Ontario” almost completely determine the remaining letters.

A more direct approach is to store a dictionary explicitly. When a word is encountered, a short code for it is sent, rather than the letters. Or rather than store English words, we might store any string of symbols that has occurred before.

The “LZ” (for Lempel-Ziv) family of data compression algorithms build such a dictionary adaptively, based on the text seen previously. The “gzip” program is an example.



# How Well Do These Methods Work?

I applied a version of PPM (by Bill Teahan) and the gzip program to the three English text files (Latex) I previously used to test Markov models.

<b>PPM:</b>	Uncompressed file size	Compressed file size	Compression factor	Bits per character
	2344	1042	2.25	3.56
	20192	5903	3.42	2.34
	235215	51323	4.58	1.75

  

<b>GZIP:</b>	Uncompressed file size	Compressed file size	Compression factor	Bits per character
	2344	1160	2.02	3.96
	20192	7019	2.88	2.78
	235215	70030	3.36	2.38

One other difference: On the long file, PPM took 2.2s to encode and 2.3s to decode; gzip needed only 0.06s to encode, and an unmeasurably small time to decode.

# Merits of Probabilistic Models

$N$ -th order Markov models and PPM models cleanly separate the *model* for symbol probabilities from the *coding* based on those probabilities.

Such models have several advantages:

- Coding can be nearly optimal (eg, using arithmetic coding).
- It's easy to try out various modeling ideas.
- You can get very good compression, if you use a good model.

The big disadvantage:

- The coding and decoding involves operations for every symbol and every bit, plus possibly expensive model updates, which limits how fast these methods can be.

# Merits of Dictionary Methods

Compression using adaptive dictionaries may be less elegant, but has its own advantages:

- Dictionary methods can be quite fast (especially at decoding), since whole sequences of symbols are specified at once.
- The idea that the data contain many repeated strings fits many sources quite well — eg, English text, machine-language programs, files of names and addresses.

The main disadvantage is that compression may not be as good as a model based method:

- Dictionaries are inappropriate for some sources — eg, noisy images.
- Even when dictionaries work well, a good model-based method may do better — and can't do worse, if it uses the same modeling ideas as the dictionary method.

# The LZ77 Scheme

This scheme was devised by Ziv and Lempel in 1977. There are many variants, including the method used by `gzip`.

The idea of LZ77 is to use the past text as the dictionary — avoiding the need to transmit a dictionary separately. We need a buffer of size  $W$  that contains the previous  $S$  characters plus the following  $W - S$  characters.

We encode up to  $W - S$  characters at once by sending the following:

- A pointer to a past character in the buffer (an integer from 1 to  $S$ ).
- The number of characters to take from the buffer (an integer from 0 to  $W - S - 1$ , or maybe more).
- The single character that follows the string taken from the buffer.

# An Example of LZ77 Coding

Suppose we look at the past 16 characters, and look ahead at the next 8 characters.

After encoding the first 16 characters of the following string, we would proceed as follows:

W a y _ o v e r _ t h e r e _ i s _ w h e r e _	i t _ i s
---	-----------

No match with string in window.

Transmit (-,0,s)

W a y _ o v e r _ t h e r e _ i s _ w h e r e _ i	t _ i s
---	---------

Match 3 back with \_

Transmit (3,1,w)

W a y _ o v e r _ t h e r e _ i s _ w h e r e _ i t _	i s
---	-----

Match with 9 back with here\_i

Transmit (9,6,t)

# Encoding the Pointers

If we look back  $S$  characters, we can encode a pointer back in  $\lceil \log_2(S) \rceil$  bits.

If we look forward  $W - S$  characters, we can encode the length of the match in  $\lceil \log_2(W - S) \rceil$  bits.

The character after the match can be encoded in  $\lceil \log_2(I) \rceil$  bits, if we have  $I$  symbols.

If these lengths are multiples of 8, we can quickly output these codes as one or more bytes.

An alternative: Use Huffman or arithmetic coding. This will give better compression, but won't be as fast.

# LZ77 Encoding and Decoding Speed

Even if writing the codes for the match is fast, *finding* the longest match may be slow.

Techniques such as hashing can speed this up, however. The `gzip` program builds a hash table for all strings of length three, then searches within the hash bucket for the next three characters to find the longest match.

Decoding can be very fast. Reading the codes is very quick if they take up fixed numbers of bytes. Even if we use Huffman codes, table look up on the next few bits (as in `gzip`) can be pretty fast. Once we have the codes, we just copy text from the buffer.

# The LZ78 Scheme

Ziv and Lempel introduced another scheme in 1978, in which the dictionary is kept explicitly, and contains phrases from the entire past text.

In the LZW variant, due to Welch, we start with a dictionary containing just the alphabet. We then proceed as follows:

- Find the longest match of following characters with a dictionary item.
- Transmit the index of that dictionary item.
- Add the matched phrase plus the character following it to the dictionary.
- Continue coding with the character following the matched phrase.

Codes for dictionary indexes will have to get longer as we go, but at a fairly slow rate.