

Data Compression - Report

lrk99

Version 2021.01.25

1 Lempel-Ziv-Welch (LZW)

Lempel-Ziv-Welch (LZW) is appropriate for a large tex file since there are likely to be repeated occurrences of both sequences of letters (words) and sequences of words (sentences). This is where LZW excels, since it builds up longer and longer sequences of characters in its dictionary which can then be encoded as a vastly shorter sequence, possibly just 2 bytes. My idea for how to use LZW was to encode each token as the smallest possible number of bytes, each byte can encode a number between 0 and 255, this means with 2 bytes we can store any number between 0 and 65535. Therefore if our maximum token is 65535, we can encode every token as 2 bytes, supporting the worst case scenario for a file of size 65KB. However if the file is large enough that we require more than 65535, we can increase the number of bytes used to encode each token to 4, accommodating the worst case scenario for files up to 4GB. The change is made by switching the mode used to encode the tokens (`struct.pack()`) from short unsigned integer (H) to long unsigned integer (L). Since the tex file that my program will be tested on is around 200 pages, which should be around 1MB in size, this upper limit of 4 bytes is clearly sufficient. It is also clear that only using 2 bytes would likely be insufficient. The number of bytes required to encode each token is included in the encoded file (as the first byte), which means the decoder knows how to decode the tokens and build the dictionary properly. The program performs well on the lecture notes files (1.3 ratio) and achieves good results on many large text benchmarks, scoring a ratio of 1.84 on the Bible, 1.99 on As You Like It, 2.17 on Alice in Wonderland, and 1.71 on The Project Gutenberg World Factbook. I chose to use LZW over LZ variants such as LZ77 and ROLZ since I think that for this assignment LZW is most appropriate. There may be huge gaps between repeated instances of sequences in the large tex file, which would hurt the efficiency of LZ77, and ROLZ is not necessary here since we only have to compress one tex file, so there is no context variation.

2 Prediction by Partial Matching (PPM)

Prediction by Partial Matching (PPM) is appropriate for a large English tex file since it is possible to predict the next character given the previous character(s). As such, a statistical model of english text is appropriate to use. This statistical model can be 'trained' on lots of English text such as Alice in Wonderland, and other tex documents such as the lecture notes for this module. It works by using a statistical model, stored as an external file to be used by both the encoder and decoder, which denotes how likely a given character is to appear next, given the context of the previous n characters, where n is the length of the context. This model is 'trained' on other English tex files, and stored in a json format. The encoder uses the model to encode each byte of input data from the tex file as a shorter sequence of bits, thus achieving compression. In my implementation I take the entire bit stream from the encoder and convert it to bytes by splitting every 8 bits up and writing the byte representation to the lz file. The decoder then reads the binary file as input and converts it back to a stream of bits. Using the same external json file it can decode this bit stream back into the original tex file, since it follows the same algorithm as the encoder but in reverse. With regards to using Method A, B, or C to assign frequencies to the escape character, I chose to use Method C. This is because it is the most suitable for this implementation since it takes into account the fact that some contexts can be followed by virtually any other character by giving the escape symbol an appropriate count, whilst not reducing

the count of other symbols. I chose to use a maximum context length of 5, but given more time I would further investigate how increasing this maximum up to 10 affects the tradeoff between time taken for the encoding and decoding processes and the level of compression achieved. Also, I chose to use Arithmetic Coding rather than Huffman Coding, since its codewords encode messages rather than being limited to symbols, and it's more efficient. PPM performs very well on text compression benchmarks such as enwik8 and enwik9 [5]. I was unable to get a working version of PPM within the time constraints, so I have left the code of my attempt in the encoder, but commented out, as a demonstration of my efforts.

3 Burrows-Wheeler Transform (BWT)

The Burrows-Wheeler Transform (BWT) rearranges a character string into runs of similar characters [2]. The BWT is appropriate for this assignment because it is possible to scan and manipulate the entire file before having to encode it. This is a fundamental difference between BWT and PPM, since PPM does not require access to the entire file at once. BWT works by taking an input string with little redundancy or predictability and outputting a string with a large amount of repeating characters and thus high redundancy. The crucial reason this works is that the 'sorting' of the characters done by Burrows-Wheeler Forward in the encoding process can be reversed by Burrows-Wheeler Back in the decoding process. To transform a string in the forward section of the algorithm, we insert a marker character at the start of the string, and then generate all of the cyclic shifts of this new string. Next, we sort these cyclic shifts in ascending order and take the last character of each shifted string. This is our output, it's a partially sorted version of our original string, with a marker character added at some index. The important property of this partially sorted string is that it has more redundancy than the original, meaning it lends itself much better to compression algorithms. The BWT is a pre-processing step which is done before using a compression algorithm such as LZW or even just run-length encoding. This algorithm improves the efficiency of the compression algorithm used afterwards, without storing any extra data other than the position of the first original character. This means that BWT is a virtually free way to get better compression results since it only increases computation requirements. Many of the best text-based compression algorithms use BWT [5]. Within the time constraints I was unable to get a working version of the BWT.

4 Move-to-front transform (MTF)

The move-to-front (MTF) transform, which I first read about in [6], is a very powerful transform which significantly improves the compression performance of various algorithms such as Huffman coding and LZW. It works by moving the last symbol seen to the front, such that its index in the table becomes 0. This is similar to what the Run-Length-Transform (RLT) does, however, the MTF transform is a big improvement over RLT since it transforms not just runs of bytes but also other kinds of string patterns. This is possible since it actually outputs a character's position in the symbol table, rather than the character itself, allowing it to output the same sequence of position codes but this sequence might denote different strings at different states of the process. This idea of the same token representing different strings at different times is possible because the symbol table is transformed in a defined way at both the encoding and decoded stages, allowing the decoder to recreate the table exactly as it was when the encoder reached a given character. While this transform works well on highly-redundant files such as images, it does not perform well on user-created text-based files since they are often not highly-redundant and don't contain many runs of distinct bytes and similar-context strings. For the purposes of this assignment I need to use a transform which can accept seemingly-random input and produce an output of redundant bytes, and this can be done much more effectively by the BWT. As such I will not use the MTF transform here.

5 Context Mixing (CM)

Context mixing is based on PPM in that it is comprised of a predictor and a coder, but it is an improvement over PPM in that the predictor is actually a combination of multiple models which have been trained on different contexts. The most useful contexts for compressing tex files are the previous n bytes before the symbol, (the context used in PPM), and the previous n words. There are other contexts which can be used for different file types to achieve better results. The weighted combination of these two models often provides a more accurate prediction than either model on its own. The disadvantage of using context mixing is that it takes longer and uses more memory, since we have to store two models and read both of them to calculate our more accurate prediction. This is a tradeoff worth making in this assignment since time taken and memory usage are not assessed, but compression ratio is, and combining statistical models will yield more accurate predictions of the next character and thus better compression results. Context mixing produces the best lossless text compression results outside of using a neural network [5], and is also used in the PAQ series of lossless data compression programs [4]. Within the time constraints I was unable to produce a working context mixing encoder and decoder.

6 Long short-term memory (LSTM)

Long short-term memory (LSTM) is appropriate for this assignment since it achieves state-of-the-art results on numerous text compression benchmarks [3, 5]. LSTM is similar to context mixing in that it uses a combination of multiple models to predict the next character. The difference is that these models are built using recurrent neural networks, instead of simply the previous n words or n bytes, providing more accurate predictions and better compression results. LSTM is used in cmix [3], which produces the best results for enwik8 and up until 2021 [1] produced the best compression on enwik9 [5]. I was unable to implement a working LSTM system within the time constraints.

References

- [1] Fabrice Bellard. Nncp v2: Lossless data compression with transformer, 2021. https://bellard.org/nncp/nncp_v2.pdf.
- [2] Michael Burrows and David J Wheeler. A block-sorting lossless data compression algorithm. 1994.
- [3] Byron Knoll. cmix homepage, 2019. <https://www.byronknoll.com/cmix.html>, Last accessed on 2021-01-19.
- [4] Matt Mahoney. Data compression programs. <http://mattmahoney.net/dc>, Last accessed on 2021-01-19.
- [5] Matt Mahoney. Large text compression benchmark, 2021. <http://mattmahoney.net/dc/text.html>, Last accessed on 2021-01-19.
- [6] Gerald Tamayo. The data compression guide, 2020. <https://sites.google.com/site/datacompressionguide>, Last accessed on 2021-01-19.