

CS222: Compression Problems

1. Compare the performance of `bzip2`, `gzip`, and `compress` (or comparable programs for your system— if you have to, download these utilities off the Web) on various file types and sizes. Make charts that detail your comparison. You should note encoding and decoding times as well; you may wish to experiment with different settings of parameters. (For example, `gzip` can try to compress more at the expense of speed.) (Write approximately 1 page.)

2. A fair coin is flipped until the first head occurs. Let X denote the number of flips required. Find the entropy $H(X)$ in bits. Suppose your friend flips a fair coin until the first head is flipped to generate a value for X , and now you want to ask a series of yes-no questions (of the form “is X contained in the following set?”) to determine the value of X generated. Describe what questions you ask, determine the expected number of questions you ask, and compare your result with $H(X)$.

3. Let X and Y be random variables, where X takes on values x_1, x_2, \dots, x_r and Y takes on values y_1, y_2, \dots, y_r . Let $Z = X + Y$.

- Show that $H(Z|X) = H(Y|X)$.
- Show that if X and Y are independent, $H(Z) \geq H(X)$.
- Give an example where $H(X) > H(Z)$. (Note X, Y must be dependent!)
- Under what conditions does $H(Z) = H(X) + H(Y)$?

4. Consider an n -sided die, where the i -th face comes up with probability p_i . Show that the entropy of a die roll is maximized when each face comes up with equal probability $1/n$. (Hint: you may wish to show more generally that if two faces have probability p_i and p_j with $p_i < p_j$, the entropy increases if you change the probabilities of these two faces to $p_i + \epsilon$ and $p_j - \epsilon$ for some suitable ϵ .)

5. Explain in reasonably pedantic detail the steps of the Burrows-Wheeler algorithm (both compression and decompression) on the Sesame Street phrase “wabbawabbawoo”.

6. Consider arithmetic coding for the string *aacbca* when the probability of an a is 0.2, a b is 0.3, and a c is 0.5. Show each step for idealized arithmetic coding, with real number arithmetic. Now suppose someone gave you the real number 0.63215699. Decode a sequence of length 10 corresponding to the above model.

7. Give the grammar SEQUITUR would derive for the string

aactgaacatgagacatagagacag.

Give enough detail so that we can check your work. Do the same for the reverse of the above string. (Hint: be careful— it’s easy to make a mistake!)

8. Consider the following eight by eight table, corresponding to pixel values for an image.

124	125	122	120	122	119	117	118
120	120	120	119	119	120	120	120
125	124	123	122	121	120	119	118
125	124	123	122	121	120	119	118
130	131	132	133	134	130	126	122
140	137	137	133	133	137	135	130
150	147	150	150	150	150	150	150
160	160	162	164	168	170	172	175

Go through the steps of compressing and decompressing this eight by eight block with JPEG. (We will skip the entropy coding step, and just do the quantization.)

Recall the appropriate compression steps:

- Translate so that the original values 0 to 255 instead vary between -128 and 127 .
- Apply the discrete cosine transform.
- Round according to the quantization table, given below.

And to decompress:

- Take the compressed form, and reverse the quantization.
- Apply the inverse of discrete cosine transform.
- Round to integers and necessary.

16	11	10	16	24	40	51	61
12	12	14	19	26	58	60	55
14	13	16	24	40	57	69	56
14	17	22	29	51	87	80	62
18	22	37	56	68	109	103	77
24	35	55	64	81	104	113	92
49	64	78	87	103	121	120	101
72	92	95	98	112	100	103	99

In one sentence, how well do you think JPEG does?

9. Compress the string “abracadabraabadacarba” using the LZ77 approach (sliding windows), with a window size of six and a lookahead buffer of size six. Similarly, compress it using the LZ78 and LZW approaches. For the LZW approach, assume the letters a, b, c, d, and r start in the dictionary in alphabetical order.

10. Delta encoding has been suggested as a means of compressing sorted lists of URLs and in other aspects of Web compression. We’ll examine a variation here for sorted lists of numbers. You will need to find a good random number generator for your platform. You may also use any sorting routine available on your platform.

Generate a list of 10,000 random numbers on the range $[0, 2^{30})$. If we were not using any compression at all, we would need to use 300,000 bits to send all the numbers.

Now suppose we sort the numbers, so they form a sorted sequence $a_0, a_1, a_2, \dots, a_{9999}$. Now we calculate all the differences $a_1 - a_0, a_2 - a_1, a_3 - a_2$, and so on. Take the largest difference D , and calculate $k = \lceil \log_2 D \rceil$. We know that representing each of the differences requires only k bits.

To transmit the numbers, we first transmit a_0 using 30 bits and k using 5 bits, just to be safe (since $k < 32$, we only need five bits to represent it). We then transmit all the differences $a_1 - a_0, a_2 - a_1, a_3 - a_2$, and so on, using only k bits for each. Argue that this is enough information to reconstruct the numbers. How many bits does it take to send this information? Repeat the experiment 10 times, and give the results for each trial.

11. *For your edification, bonus, do not turn anything in.* Recall that in arithmetic coding, a sequence X ends up corresponding to an interval $(F(X) - p(X), F(X)]$. Show that using the first $\lceil \log 1/p(X) \rceil$ bits of $F(X)$ could give you a codeword from which you can decode correctly. Give an example to show that this choice does not lead to a prefix-free code. Then show that using $F(X) - p(X)/2$ rounded to $\lceil \log 1/p(X) \rceil + 1$ bits gives a prefix-free code.