

1. a. Uniform distribution, meaning that all outcomes are equally probable.

Remark: You can think of this as a 0-gram model! Consider that

b. 4. Usually the answer would depend on the text, but in this case, the model assigns exactly $-\log(1/16) = 4$ bits to each letter. (In information theory, logs are base 2 by convention.)

c. $2^4 = 16$. Yes, this is related to the $1/16$ probability.

d. infinity.

That's because $p(\text{an English text}) = 0$ under your model. Your model doesn't allow the letter "x", for instance.

(We say that $-\log(0) = \text{infinity}$ because infinity is the limit of $-\log(x)$ as $x \rightarrow 0$.)

e. infinity.

2. False.

But almost true! It's true of cross-entropy, not perplexity.

$\text{crossentropy/word} = 6 * \text{crossentropy/letter}$, so
 $\text{perplexity/word} = (2^{(\text{crossentropy/word})})$
 $= 2^{(6 * (\text{crossentropy/letter}))}$
 $= (2^{(\text{crossentropy/letter})})^6$
 $= (\text{perplexity/letter})^6$

For example, if you are 10-ways perplexed about what each letter is, then you will be a million-ways perplexed (10^6) about what each word is.

But wait, I wrote above:

$\text{crossentropy/word} = 6 * (\text{crossentropy/letter})$.

Why is that true?

Well, suppose there are N words in the corpus.

Then there must be $6N$ letters (so that each word is 6 letters on average).

So total cross-entropy of the corpus can be expressed in 2 different ways:

$(\text{crossentropy/word}) * N \text{ words} = (\text{crossentropy/letter}) * 6N \text{ letters}$

Dividing both sides by N gives

$(\text{crossentropy/word}) = 6 * (\text{crossentropy/letter})?$

More intuitively, this is for the same reason that when you take a bus trip,

$\text{cost/week} = 7 * \text{cost/day}$

If you travel by bus for 100 weeks and it costs you \$840 total, then the cost was \$8.40/week or \$1.20/day. (Maybe some days were cheaper and some days were more expensive; this is just the AVERAGE cost.)

In the same way, if the overall log-probability of your test corpus is 2^{-840} , giving a cost of 840 bits to encode the whole corpus, and the corpus consisted of 100 words = 600 letters, then the cross-entropy was 8.4 bits/word or 1.4 bits/letter.

Think of the log-probability as the cost of driving through the corpus using your model (your bus). You have to pay to be surprised.

Some portions of the journey cost more bits than others, because

they are longer, or more surprising in context (driving the bus uphill requires more gas!).

But to get the average cost per word or per letter, you look at the total cost and divide by the total number of words or letters. This is the cross-entropy of the model.

3. (a)

	HHHH	HHHH	THHH	HTHH	THHT	HTHT
unsmoothed	1	1	3/4	3/4	1/2	1/2
add-1 smoothed	5/6	5/6	2/3	2/3	1/2	1/2

- (b) 3/4 (just average the first line above)
 2/3 (just average the second line above)

(c) Less bias: unsmoothed.

The unsmoothed technique has no bias at all. That means that ON AVERAGE, it will get the right answer. In fact, the average answer above happened to be exactly correct (my coin does in fact have $p(H) = 3/4$). So the bias is 0 (formally, the bias is the average of (estimated value - correct value)).

By contrast, the smoothed technique will be biased toward giving answers that are closer to 1/2 than the truth, because adding 1 to the H and T counts always makes the estimate closer to 1/2 than what was observed.

Less variance: smoothed.

You will note that in the 6 columns above (alternate universes), the smoothed estimates are closer to one another than the unsmoothed estimates. So our estimate varies less from universe to universe.

(Formally, the variance is the average of (estimated value - mean estimated value)². Over these 6 universes, that average is 0.022, but the true variance of the smoothed estimator considers what would happen if you repeated the experiment infinitely many times, not just 6.)

In the unsmoothed case, the AVERAGE estimate (over all universes) will actually match the truth. That is what it means for the unsmoothed estimator to be unbiased (0 bias).

But who cares that the unsmoothed estimate is right on AVERAGE? The SINGLE estimate in YOUR universe could be wildly wrong. For example, in the first two columns, it estimates $p(H)=1$ and $p(T)=0$, which would put you in a world of hurt when you eventually see a T. You may be better off with the smoothed estimate, which goes a little bit wrong on average (bias > 0) but is more consistent (low variance) and so rarely goes too far wrong.

- (d) Each of the H flips costs $-\log(0.8)$ bits under your model, and each of the T flips costs $-\log(0.2)$ bits. All logs are base 2 as usual.

So the total cost per flip is
 $(-4*\log(0.8) - 2*\log(0.2))/6$
 or equivalently
 $-(2/3)*\log(0.8) - (1/3)*\log(0.2)$

Note: If you want to get really fancy, you can notice that
 $\log(0.8) = \log(0.2*4) = \log(0.2)+\log(4) = \log(0.2)+2$
 so the answer becomes

$-\log(0.2) - 4/3$
 or equivalently
 $\log(5) - 4/3$
 where $\log(5)$ is a little more than $\log(4)=2$; so the answer is about 1 bit.
 (With a calculator, the exact number turns out to be about 0.99 bit.)

(e) The problem is that 6 flips is not really enough to see the true behavior of the coin. The easiest solution is to flip (say) 1000 times instead of 6 times. A 1000-flip sequence will give a lower-variance estimate of cross-entropy, because one 1000-flip sequence of this coin is much like another.

(Since the coin has $p(H)=3/4$, very close to $3/4$ of the 1000 flips will be heads. By contrast, we got a 6-flip sequence HHHHTT in which only $2/3$ of the flips were heads, and this is far from unusual.)

4. The correct answers are a, c, d.

(a) should look familiar from an exercise on homework 2. It can be obtained by backing off from this use of the chain rule:

$$\begin{aligned}
 & p(w_0, w_1, w_2, w_3, w_4, w_5, w_6) \\
 &= p(w_6) \\
 & \quad * p(w_5 \mid w_6) \\
 & \quad * p(w_4 \mid w_5, w_6) \\
 & \quad * p(w_3 \mid w_4, w_5, w_6) \\
 & \quad * p(w_2 \mid w_3, w_4, w_5, w_6) \\
 & \quad * p(w_1 \mid w_2, w_3, w_4, w_5, w_6) \\
 & \quad * p(w_0 \mid w_1, w_2, w_3, w_4, w_5, w_6)
 \end{aligned}$$

(c) can be obtained by backing off from this use of the chain rule:

$$\begin{aligned}
 & p(w_0, w_1, w_2, w_3, w_4, w_5, w_6) \\
 &= p(w_0) \\
 & \quad * p(w_1 \mid w_0) \\
 & \quad * p(w_2 \mid w_0, w_1) \\
 & \quad * p(w_3 \mid w_0, w_1, w_2) \\
 & \quad * p(w_4 \mid w_0, w_1, w_2, w_3) \\
 & \quad * p(w_5 \mid w_0, w_1, w_2, w_3, w_4) \\
 & \quad * p(w_6 \mid w_0, w_1, w_2, w_3, w_4, w_5)
 \end{aligned}$$

The same is true of (d). If you're not convinced, notice that to back off, we always generalize the condition to the right of the bar. Let's do this formally. Let W_0, W_1, W_2, \dots be random variables denoting the words, and U_i be a random variable denoting the number of uncommon words in (W_0, W_1, \dots, W_i) . The notation

$$\begin{aligned}
 & p(w_3 \mid w_0, w_1, w_2) \\
 \text{is really just shorthand for} \\
 & p(W_3 = w_3 \mid W_0 = w_0, W_1 = w_1, W_2 = w_2) \\
 &= p(W_3 = \text{song} \mid W_0 = \text{BOS}, W_1 = \text{Sing}, W_2 = a) \\
 &= p(W_3 = \text{song} \mid W_0 = \text{BOS}, W_1 = \text{Sing}, W_2 = a, U_2 = 1)
 \end{aligned}$$

We can back off from this to

$$= p(W_3 = \text{song} \mid U_2 = 1)$$

as claimed.

However, (b) is NOT valid. The problem is that, for example, W_2 is generated dependent on W_3 and also vice-versa. This is circular.

In any valid use of the chain rule, you introduce variables one at a time, in some fixed order, with each variable's value being conditioned on only the values already chosen for the previous variables -- or, with backoff, on some subset of those values.

But (b) cannot have been obtained this way, because of the circularity noted above. (W_2 must have been introduced after W_3 , and also vice-versa, which is impossible.)

5a. Answer 1: Because we have no training data for Synglish.

(Although it is possible to create such training data. Just simulate the plagiarism process by sampling sentences from C and randomly replacing their words with synonyms. This will give you a sample of Synglish, at least if our model of the plagiarism process is correct. Now you can indeed train a log-linear model to distinguish between English and Synglish.)

Answer 2: Even if we have training data, it is not clear what features to use. Whether the log-linear model performs well depends on whether it has adequate features. The method in the rest of the problem doesn't require any feature design -- it simply draws the appropriate conclusions from the language models.

5b. Bayes' Theorem says that the two probabilities in the odds ratio are respectively proportional to $p(\text{Synglish}) * p(W | \text{Synglish})$ and $p(\text{English}) * p(W | \text{English})$. Therefore, their ratio is

$$\frac{p(\text{Synglish}) p(W | \text{Synglish})}{p(\text{English}) p(W | \text{English})}$$

Note that $p(\text{Synglish})$ is our prior probability that a student decides to plagiarize. $p(W | \text{Synglish})$ is given by a model of how text is generated in this case. The problem says that we have such a model, and it is called p_{Synglish} . So the final answer is

$$\frac{p(\text{Synglish}) p_{\text{Synglish}}(W)}{p(\text{English}) p_{\text{English}}(W)}$$

Remark: If you are feeling insecure about the odds ratio above, here is a more explicit use of Bayes' Theorem to get it:

$$\frac{p(\text{Synglish} | W)}{p(\text{English} | W)} = \frac{p(\text{Synglish}) * p(W | \text{Synglish}) / p(W)}{p(\text{English}) * p(W | \text{English}) / p(W)}$$

When we take the odds ratio of these two, the identical denominators cancel out. This is very convenient because the denominator has the slightly complicated form

$$p(W) = p(\text{Synglish}) * p(W | \text{Synglish}) + p(\text{English}) * p(W | \text{English})$$

In general, when using Bayes' Theorem, it is sometimes easier to think about RELATIVE probabilities of different outcomes (odds ratios) without bothering to convert them to absolute probabilities.

5c. Here is one reasonable answer:

$$p_{\text{Synglish}}(W) = \sum_V p(V) * p(w_1 | v_1) * \dots * p(w_n | v_n)$$

This assumes that for each word, the student makes an INDEPENDENT decision about whether and how to replace it.

Notice that this formula sums over the many ways of getting

W = In the senior sunshines, this was easy.

We don't know whether V was

V = In the old days, this was easy.

V = Inside the old days, this was simple.

V = In the super sunshines, this was simple.

...

So we have to sum over all the possibilities. What we have really done is to build a joint model of $p(V,W)$, and then we obtain $p(W) = \sum_V p(V,W)$.

Now, we have to fill in the details. First of all, what is $p(V)$? This is the English sentence that the cheating student started with. As noted in question 4d, there are at least two reasonable models that we could use:

$$(A) \quad \begin{aligned} p(V) &= 1/|C| \text{ if } V \text{ is in } C \\ p(V) &= 0 \quad \text{otherwise} \end{aligned}$$

$$(B) \quad p(V) = p_{\text{English}}(V)$$

Second, what is $p(w_i | v_i)$? Well, the student doesn't replace all the words. A simple model is that the student has a 90% chance of leaving v_i alone (so $w_i = v_i$), and a 10% chance of replacing it with an element of $\text{syn}(v_i)$, where all elements including v_i itself are equally likely to be chosen. Under this model,

$$\begin{aligned} p(w_i | v_i) &= 0.9 + 0.1 / \text{syn}(v_i) \text{ if } w_i = v_i \\ p(w_i | v_i) &= 0.1 / \text{syn}(v_i) \quad \text{if } w_i \text{ is in } \text{syn}(v_i) \text{ but } w_i \neq v_i \\ p(w_i | v_i) &= 0 \quad \text{otherwise} \end{aligned}$$

(Remark: A more complex model might say that if v_i is a function word, then the student has a 98% chance of leaving it alone, but if v_i is a content word, then the student has only an 80% chance of leaving it alone.)

(These numbers such as 0.9, 0.98, 0.8 don't have to be hard-coded. They could be learned parameters of the model. But we would need English data to learn them from.)

Note that are other possible models. For example, another model is that the student first chooses k words from the n word sentence to replace, and then replaces those. This means that replacing one word makes the student less likely to replace other words. That is different from making an independent decision for each word.

5d. Model (A) is like a 100-gram model that is trained on C , and just memorizes C .

Model (B) is a 3-gram model that is trained on C , so it may generalize to sentences outside of C .

In general, reducing the value of n in an n -gram model will increase the bias (bad) but reduce the variance (good).

i. Model (B) suffers from greater bias because it pretends that English is a trigram model. Model (B) cannot ever learn properties of English like "a sentence almost always has exactly one main verb."

Model (A) is not biased at all, because ON AVERAGE it will give you the correct distribution of English. (If you average over many random training sets C .)

ii. Model (A) suffers from greater variance because the estimates are highly sensitive to the particular training set C . In particular, model (A) estimates sentences that are in C to

have probability $1/|C|$, but estimates sentences that are not in C to have probability 0, which is worse by a factor of infinity!

Model (B) smooths the training data somewhat simply by using a trigram model, so that it can assign positive probability to sentences outside C . But model (A) doesn't smooth the training data at all.

- iii. There are a lot of options between a 3-gram model and a 100-gram model. How about a 4-gram model? Better yet, how about a 100-gram model estimated with backoff smoothing?

Another possibility is to take a linear interpolation of the two models:

$$p(V) = \lambda * p_A(V) + (1-\lambda) * p_B(V)$$

But this is probably not as clever as using backoff smoothing, which factors $p(V)$ into a product of contextual probabilities, each of which can be interpolated to a *different* degree depending on how well its context was observed in C .

- 5e. i. We want V that maximizes $p(V) * p(W | V)$.
(Note that in question 4c, we summed over the possible source sentences V , but here we are looking for the single most likely one.)

- ii. Decoding.

(We have a noisy channel model. V is the source sentence, and synonym substitution is the noisy channel that messes up V into W . We observe W and want to "decode" this message to recover V .)