

analysis-entropy-of-text

November 1, 2020

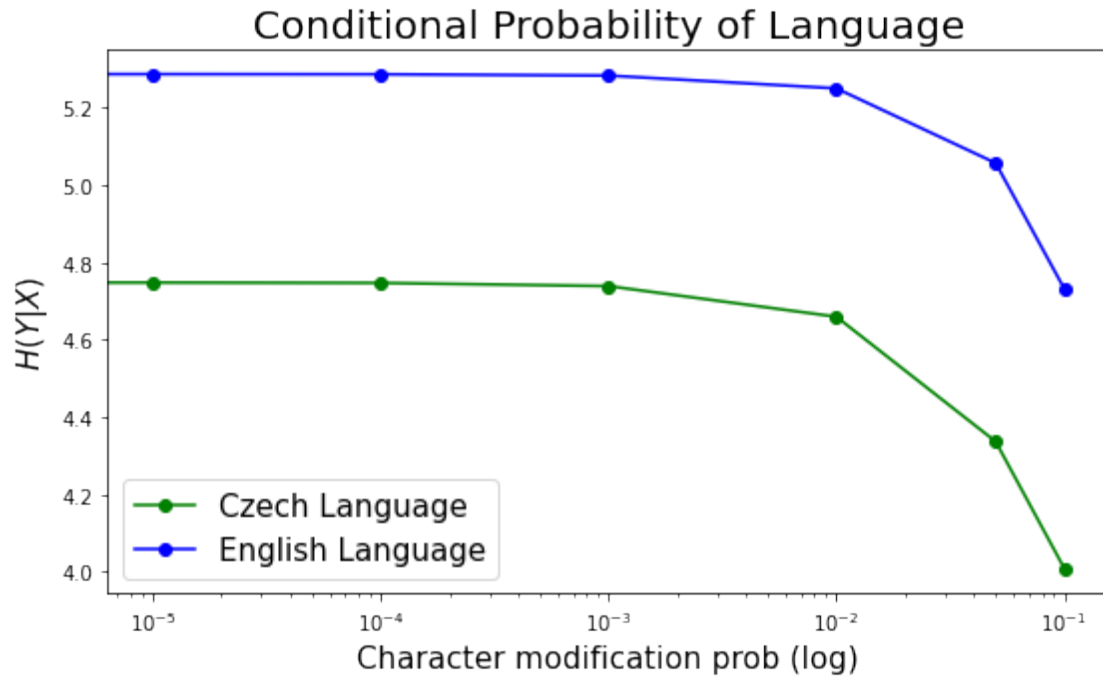
1 Conditional Entropy $H(Y|X)$

- All results are computed using the 10 repeat average.
- Note that `mess_prob=0.0` represents original (unmodified) dataset...

1.1 Character modification

- In the table below (and also in the plot) we can see, that modifying the characters leads to a linear decrease of the conditional entropy.
- When we modify characters of words, we are creating new words and thus making the vocabulary much larger.
- Moreover we destroy the relationship between the words, because for example if in the original text we had “can of coke” 100 times, now we will probably end up each time with a different word, thus each combination of words w_{i-1}, w_i will be more and more unique, because each word w_i will start becoming more unique.
- This means that $p(w_{i-1}w_i)$ will go towards $1/|V|$ and $p(w_i|w_{i-1})$ will go to 1, thus $\log_2(p(w_i|w_{i-1}))$ will more often be 0 and the conditional entropy will be lower.

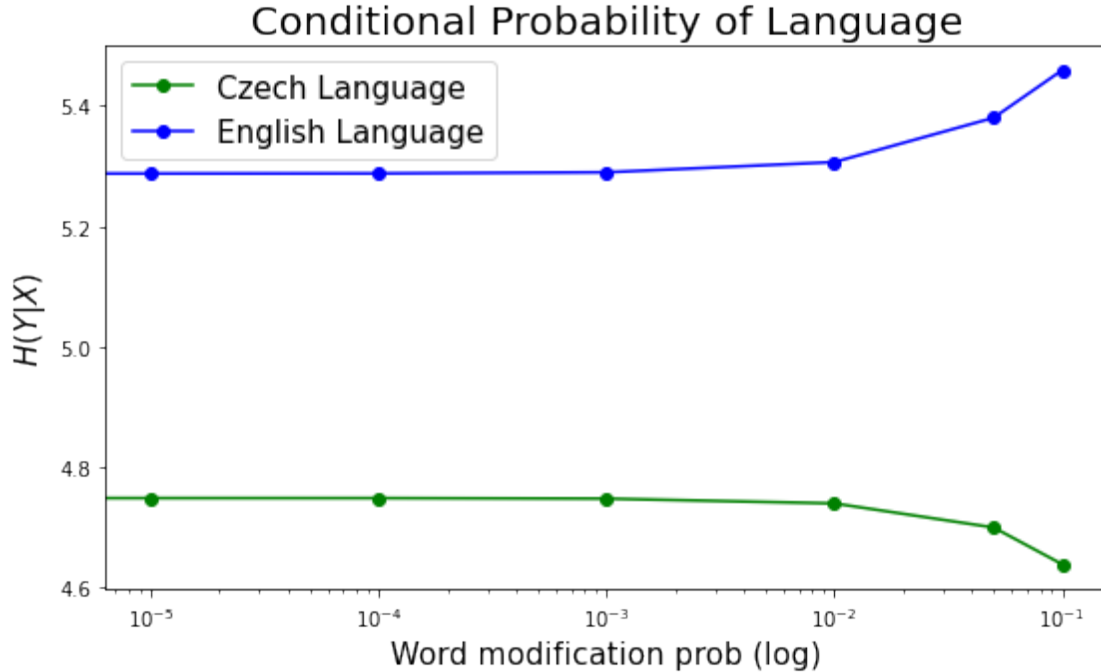
dataset	cz		en	cz_perplexity		en_perplexity
mess_prob						
0.0	4.747801	5.287398		26.867694		39.053984
1e-05	4.747718	5.287392		26.866157		39.053815
0.0001	4.746771	5.287000		26.848525		39.043208
0.001	4.738600	5.283650		26.696888		38.952674
0.01	4.659225	5.250462		25.267749		38.066804
0.05	4.336144	5.057127		20.198049		33.292550
0.1	4.007994	4.731907		16.088900		26.573320



1.2 Word modification

- Contrary to character modification, when we are substituting the words, we are not making the vocabulary larger. But we are changing the distribution of words. Each substitution leads to the fact that this word is in this position with uniform probability.
- Let's see the results first

dataset	cz	en
mess_prob		
0.0	4.747801	5.287398
1e-05	4.747787	5.287406
0.0001	4.747751	5.287613
0.001	4.746902	5.289483
0.01	4.738953	5.306460
0.05	4.698919	5.380170
0.1	4.637736	5.459387



- I think that $p(w_{i-1}w_i)$ will decrease for the most common word combinations.. because there is a high probability of modification of one of those two words (and result will uniformly random word).
- I though $p(w_i|w_{i-1})$ that will decrease, because now we will start having more and more random history for this word... and i=that it will tend to go to $1/|V|$... so that the Entropy will be increasing ... and the distribution of the words will go towards the maximal entropy...
- But we can see, that this is the case only for English. For the Czech, it quite surprisingly decreases the entropy. My explanation is that this is happening because of the very large Czech dictionary. Let us see the stats about the languages... In the table, we have the size of the dataset, vocabulary size, number of different characters... and quantiles with counts of most common words...

	dataset	T	V	C	max	q25	q50	q75	q95
0	cz	222412	42826	117	13788	1	1	2	11
1	en	221098	9607	74	14721	1	2	7	56

- We can see that Czech vocabulary has 42826 different words. Also half of the words are only single time in the dataset and even 75% of words from the vocabulary is at most two times in the dataset.
- So if we randomly change 10% of the words, we are changing about 22,241 of the words. Most probably we will choose again and again some common words.
- These words will be changed uniformly to some of 42,826 different words. So there is a high chance that often we sample a very common word and substitute it with a word that was in the

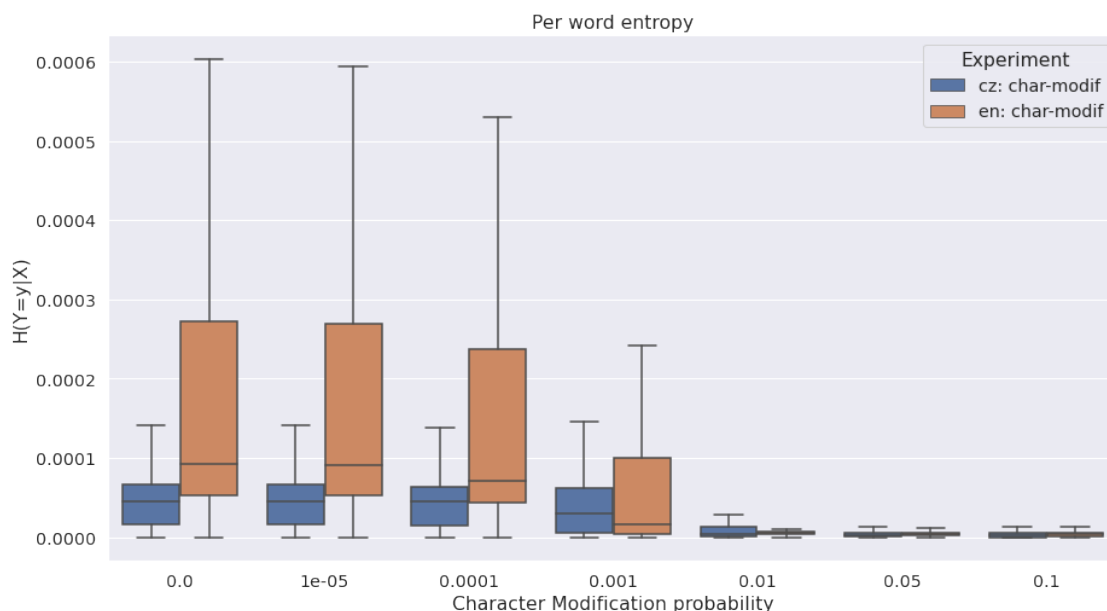
dataset only once or twice ... If this new changed ‘rare’ word r is in the history e.g. $p(w_i|r_{i-1})$ then the probability will be very high 1 or 0.5 .. thus this will lower the entropy.

- This is not the case for English, where is still a higher chance that a new word was at least multiple times in the dataset.

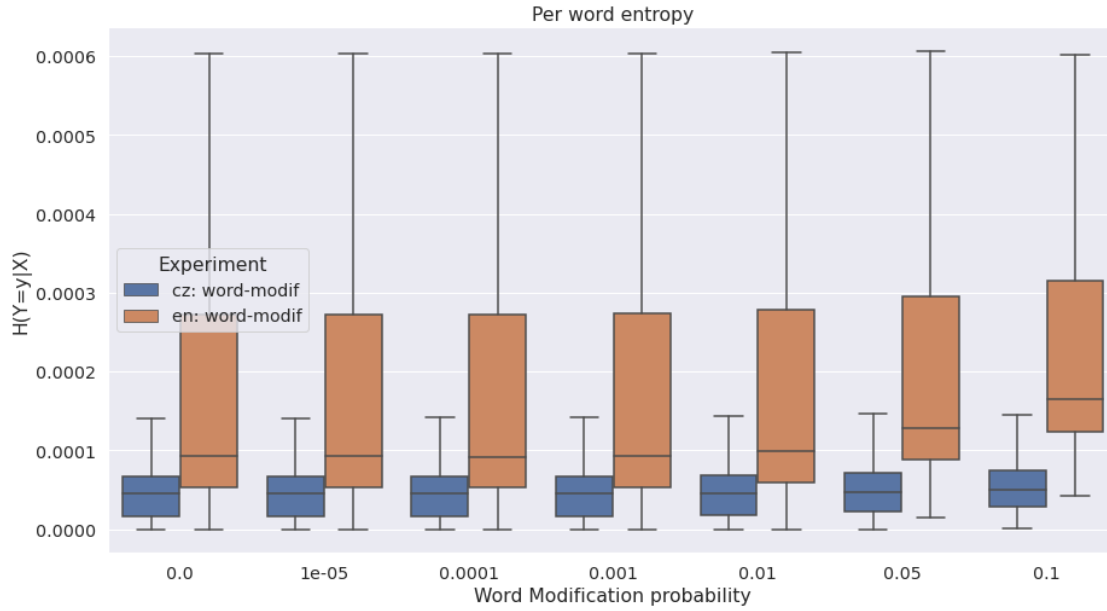
2 Per-word contribution to Conditional Entropy $H(Y = y|X), y \in \mathcal{Y}$

- I thought that it would also be interesting to see how individual words contribute to the conditional entropy.
- In the following table, we will find maximum, minimum, average, and quantiles for $H(Y = y|X)$ for different modification probabilities.
- First table and plot are for character modification. The second table and plot are for word modification.
- For the character modification we can see, that all words reduce their contribution to the conditional entropy quite linearly. We can also see that some words contribute a lot (common words) but most words contribute a little.
- For the word modification, we can see, that for English the skew of the distribution is changing and that the more we change words uniformly randomly to other words, the mean per-word contribution to the conditional entropy is getting closer to 0.5 quantile. For Czech, this is happening also, but too little.

<pandas.io.formats.style.Styler at 0x7fcdcd8f280>



<pandas.io.formats.style.Styler at 0x7fcdcd8f280>



We can also show those words, which contribute to conditional entropy the highest. Result is not suprsing ...

English

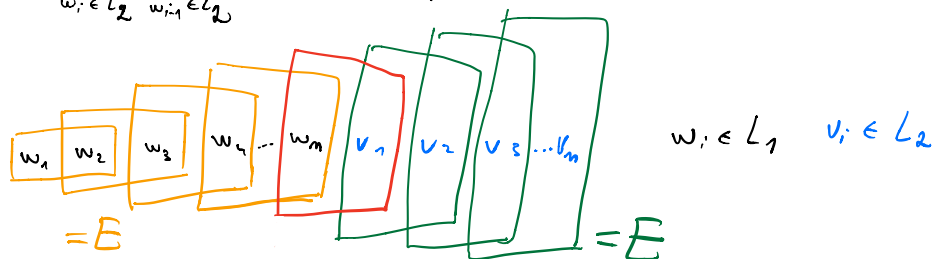
	word	avg_1
1713947	,	0.146136
1713956	the	0.143909
1713958	of	0.083818
1713955	in	0.080007
1713964	and	0.073680
1713941	.	0.068793
1714039	a	0.061798
1713968	to	0.050796

Czech

	word	avg_1
15	,	0.107800
1	.	0.101425
120	v	0.056312
53	a	0.047288
142	"	0.040737
46	se	0.040549
80	na	0.035754
18	-	0.025175

$$H(L_1) = \sum_{w_i \in L_1} \sum_{w_{i-1} \in L_1} p(w_{i-1}, w_i) \log_2(w_i | w_{i-1}) = E$$

$$H(L_2) = \sum_{v_i \in L_2} \sum_{v_{i-1} \in L_2} p(v_{i-1}, v_i) \log_2(v_i | v_{i-1}) = E$$



• For $w_i \in L_1, v_j \in L_2$ is $p(w_i | v_j) = \frac{c_2(v_j, w_i)}{c_1(v_j)} = \frac{0}{c_1(v_j)} = 0$ [v_j never before w_i]

$p(v_j | w_i) = \frac{c_2(w_i, v_j)}{|L_1| + |L_2| - 1} = 0$

• For $w_i \in L_1, v_j \in L_2$ is $p(v_j | w_i) = \frac{c_2(w_i, v_j)}{c_1(w_i)} = \begin{cases} \frac{1}{c_1(w_i)} & \text{FOR } w_i = w_m \\ 0 & \text{ELSE} \end{cases}$

$p(w_i | v_j) = \frac{c_2(w_i, v_j)}{|L_1| + |L_2| - 1} = \begin{cases} \frac{1}{|L_1| + |L_2| - 1} & \text{FOR } w_i = w_m \\ 0 & \text{ELSE} \end{cases}$

• For $w_i \in L_1, w_j \in L_1$ $p(w_i | w_j) = \frac{c_2(w_j, w_i)}{c_1(w_j)} = \text{SAME AS BEFORE}$

• For $w_i \in L_1, w_j \in L_1$ $p(w_i, w_j) = \frac{c_2(w_i, w_j)}{|L_1| + |L_2| - 1}$ much smaller...

• For $v_i \in L_2, v_j \in L_2$ $p(v_i | v_j) = \frac{c_2(v_j, v_i)}{c_1(v_j)} = \text{SAME AS BEFORE}$

• For $v_i \in L_2, v_j \in L_2$ $p(v_i, v_j) = \frac{c_2(v_i, v_j)}{|L_1| + |L_2| - 1}$ much smaller...

So... Conditional Entropy will be

$$H(L_1, L_2) = -1 \cdot \left[\sum_{w_i, w_j \in L_1} \frac{c_2(w_i, w_j)}{|L_1| + |L_2| - 1} \cdot \log_2 \left[p(w_j | w_i) \right] + \sum_{v_i, v_j \in L_2} \frac{c_2(v_i, v_j)}{|L_1| + |L_2| - 1} \cdot \log_2 \left[p(v_j | v_i) \right] + \frac{1}{|L_1| + |L_2| - 1} \log_2 \left[\frac{1}{c_1(w_m)} \right] \right]$$

$$\frac{c_2(w_i, w_j)}{|L_1| + |L_2| - 1} = \frac{c_2(w_i, w_j)}{|L_1| - 1} \cdot \frac{|L_1| - 1}{|L_1| + |L_2| - 1}$$

SAME NEW
DOES NOT DEPEND
ON i, j

$$\frac{c_2(v_i, v_j)}{|L_1| + |L_2| - 1} = \frac{c_2(v_i, v_j)}{|L_2| - 1} \cdot \frac{|L_2| - 1}{|L_1| + |L_2| - 1}$$

SAME DOES NOT DEPEND
ON i, j

$$= -1 \cdot \left[\frac{|L_1| - 1}{|L_1| + |L_2| - 1} \cdot (-E) + \frac{|L_2| - 1}{|L_1| + |L_2| - 1} \cdot (-E) + \frac{1}{|L_1| + |L_2| - 1} \log_2 \left[\frac{1}{c_1(w_m)} \right] \right] = E \left(\frac{|L_1| + |L_2| - 2}{|L_1| + |L_2| - 1} \right) - \frac{1}{|L_1| + |L_2| - 1} \log_2 \left[\frac{1}{c_1(w_m)} \right]$$

$c < 1$ d

$= cE + d$ where $0 < c < 1$ c will be close to 1
 $0 \leq d < 1$ d will be very small

SO I THINK IT WILL PROBABLY
BE A BIT LOWER THAN E