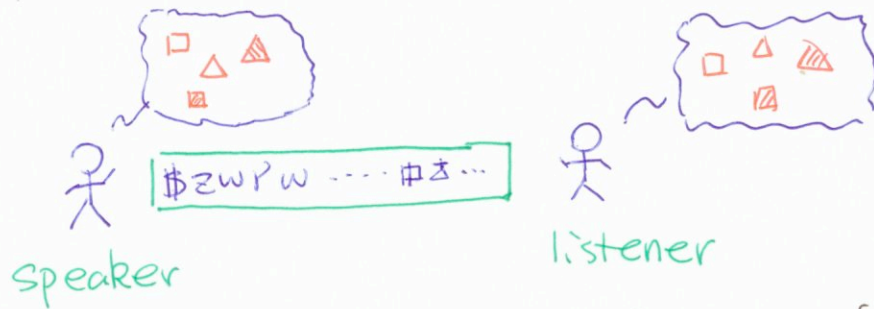


1. Natural Language processing Basics

1.1. what is language for? - Communication



A person has a representational space of semantics, by using sounds, words or even gestures to tell others about what you are thinking.

1.2 Two approaches to language "Universal Grammar"

Rationalist (Poverty of the stimulus - Chomsky)

Empiricist (emphasis on learning rather than hardwired principals of human brains)

1.3 Grammaticality and Conventionality

Colorless green ideas sleep furiously

grammatically correct!

Sheep year China February begin

unconventional but
Semantically understandable.

* | A good novelist should be with an imaginable mind as well as a technique of writing or communicating with words.

1.4 why Statistical?

- ✓ human cognition is probabilistic. So as the language.
- ✓ Computing probability of sentence from a corpus of utterances would assign the same low probability to all untested sentences, grammatical or ungrammatical, which is not true in linguistics (Chomsky)

1.5. Statistical analysis of language - Some examples,

✓ example of Tom Sawyer by Mark Twain.

8018 words for the novel (11,100 for news)

✓ How to Count words - Stemming

✓ Zipf's Law $f \propto \frac{1}{r}$

frequency

rank

✓ Mandelbrot's formule

$f = P(r+p)^{-B}$ or by taking logarithm:

$$\log f = \log P - B \log (r+p)$$

parameters to make Zipf's law more practical based on the texts given.

1.6 Can you tell an author's style from statistical analysis of their works?

It will be great to try! our first homework!

1.7. Mathematical Foundation for NLP!

✓ Essential probability theory

✓ Conditional probability and Bayes theorem

✓ Statistical distributions.

✓ estimation

✓ Entropy (Information Theory)

✓ Noisy Channel Model

1.8 Zipf's law in Chinese language Research Question!

(Character-based, or word-based)

Zipf-Mandelbrot Law extension.

1. Language Model

For some NLP tasks as the following

A. Machine Translation

$P(\text{我是一个学生} | \text{I am a student})$

B. Pinyin (Spelling)

$P(\text{我是一个学生} | \text{wo shi yige xue sheng})$

C. Speech Recognition

$P(\text{我是一个学生} | \text{wuhunwuhunwuhunwuhun})$

The signals sounds like

D. Word Segmentation

$P(\text{我是}|\text{一个学生})$ $P(\text{我是}|\text{一个学生})$

$P(\text{我是}|\text{一个学生})$ $P(\text{我是}|\text{一个学生})$

which is bigger??

All these problems are related to one central question.
What are the probabilities of particular sentences?

$$P(e) = ?? \quad e = \{w_1, w_2, \dots, w_N\}$$

E.g. in machine translation:

$$P(\underline{c} | e) \propto P(e | \underline{c}) P(\underline{c}) \quad (\text{Bayesian Rule})$$

English Chinese

$P(e) \sim$ likelihood of an English sentence

$P(\underline{c} | e) \sim$ probability of a Chinese sentence \underline{c}
given an English sentence e .

2. It is always too hard to estimate the probability of a sentence (why?). We then go for the probability of substrings of words. E.g.

$$P(e) = P(w_1, w_2, w_3, \dots, w_N) \approx \prod_{i=1}^N P(w_i)$$

A good model??

$$2.1 \quad p(e) = p(w_1, w_2 \dots w_N)$$

Independence assumption of words (A Bag-of-words)

$$\begin{aligned} P(w_1, w_2 \dots w_N) &= P(w_2, w_3 \dots w_N, w_1) \quad \text{unigram} \\ &= P(w_3, \dots, w_{N-1}, w_N, w_1, w_2) \\ &= P(\text{All possible permutations}) \end{aligned}$$

am Student
a I
Beihang of

The word order does not matter (really?)
how about the letter order in a word?

$$2.2 \quad P(\text{I am a student}) = P(I | \text{start-of-sentence}) * P(\text{am} | I) * P(\text{a} | \text{am}) * P(\text{student} | \text{a}) * P(\text{end-of-sentence} | \text{student})$$

$$\begin{aligned} P(w_1, w_2 \dots w_N) &= \prod_{i=1}^{N-1} P(w_{i+1} | w_i) \quad \text{Bi-gram} \\ P(w_1, w_2 \dots w_N) &= \prod_{i=1}^{N-2} P(w_{i+2} | w_{i+1}, w_i) \quad \text{tri-gram} \end{aligned}$$

Markov Property

2.3 Smoothing in N-gram Model

$$P(w_{i+2} | w_{i+1}, w_i) \rightarrow P(z | x, y) = \frac{\#(x, y, z)}{\#(x, y)}$$

#(.): Counting occurrences

By Smoothing:

$$\begin{aligned} P(z | x, y) &= 0.95 * \frac{\#(x, y, z)}{\#(x, y)} + 0.04 * \frac{\#(y, z)}{\#(y)} \\ &\quad + 0.008 * \frac{\#(z)}{\#(\text{all words})} + 0.002 \end{aligned}$$

$P(z | x, y) > 0.002$, these smoothing coefficients are manually assigned.

$$\blacksquare P(\text{I want to eat British food}) = P(\text{I} | \text{start} >) P(\text{want} | \text{I})$$

$$\frac{P(\text{to} | \text{want})}{*} \frac{P(\text{eat} | \text{to})}{\diamond} \frac{P(\text{British} | \text{eat})}{\triangle} \frac{P(\text{food} | \text{British})}{\square} \frac{P(\text{British})}{\star}$$

$$= .25 * .32 * .65 * .26 * .001 * .60 = .000080$$

Eat on	.16	Eat Thai	.03
Eat some	.06	Eat	.03
Eat lunch	.06	Eat in	.02
Eat dinner	.05	Eat	.02
Eat at	.04	Eat	.02
Eat a	.04	Eat	.01
Eat Indian	.04	Eat	.007
Eat today	.03	Eat British	.001 ∇

<start> I	.25 \triangle	Want some	.04
<start> I'd	.06	Want Thai	.01
<start> Tell	.04	To eat	.26 \diamond
<start> I'm	.02	To have	.14
I want	.32 \square	To spend	.09
I would	.29	To be	.02
I don't	.08	British food	.60 \star
I have	.04	British restaurant	.15
Want to	.65 \star	British cuisine	.01
Want a	.05	British lunch	.01

2.4 Perplexity

Based on what we have learnt from last a few sections on language model. A model often consists of a generative "story" (people produce words based on the last one (two or more) words they said).

uni, bigram or trigram

$$P(\text{model} | \text{test-data}) \propto P(\text{model}) * P(\text{test-data} | \text{model})$$

uniform unless we have more information.

$$P(\text{test-data} | \text{model}) = P(e) = \prod_{i=1}^N P(w_{i+1} | w_i)$$

product of small numbers

$$\text{Perplexity} = 2^{-\log(P(e))/N}$$

$$\prod_j a_j \rightarrow \log \sum_j a_j$$

In information theory, perplexity is used to measure how well a probability distribution or probability model predicts a sample, used to compare probability models.

(First used by Frederick Jelinek, whose legacy in NLP can be remembered by his famous saying: "Everytime I fire a linguist, the performance of the speech recognizer goes up")

✓ lower perplexity implies better language models.

✓ perplexity of a probability distribution is:

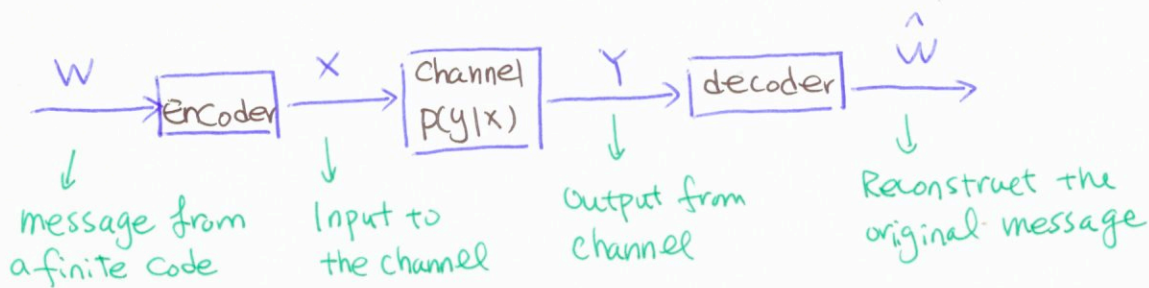
$$2^{H(P)} = 2^{-\sum_x P(x) \log_2 P(x)} \quad H(P) \sim \text{entropy}$$

✓ How to estimate the entropy of a language?

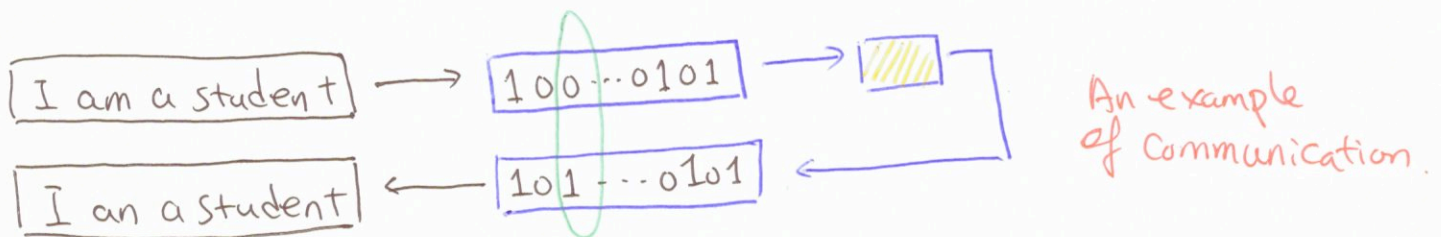
Obviously, the entropy is related to the language model.

✓ Based on research of Brown et al (1992), the upper bound of entropy of characters in printed English is 1.75 bits per character. How?

2. Information Theory for NLP.



Classical noisy channel model was proposed by Shannon as the foundation of the Information Theory.



2.1 why the noisy channel model is important for NLP?

Machine Translation

Input: L_1 word sequence	我是学生
Output: L_2 word sequence	I am a student
$P(i)$: $P(L_1)$ language model	$P(\text{我})P(\text{是})P(\text{学})P(\text{生})$
$P(o i)$: translation model	$P(I \text{我}) \dots P(\text{is} \text{是})$

Speech recognition

Input: word sequence
Output: speech signal.
$P(i)$: probability of word sequence
$P(o i)$: acoustic model (HMM)

From the above, we can ~~see~~ easily see how the information theory is related to machine learning. Most of the learning cases, we ~~are~~ just need the right model of $P(i)$ and $P(o|i)$, where "i" is examples and "o" is classes.

Any discussions or comments??

2.2. Entropy

Information measure $I(x)$

$$\begin{aligned} \text{Def: } H(X) &= E\left(\log \frac{1}{P(X)}\right) = \sum_i P(x_i) \log \frac{1}{P(x_i)} \\ &= - \sum_i P(x_i) \log P(x_i) \\ &= - \sum_i P_i \log P_i \end{aligned}$$

where $I(x)$ satisfy the following:

$$\begin{cases} I(X_1, X_2) = I(X_1) + I(X_2) \\ I(X) \propto \frac{1}{P(X)} \end{cases}$$

that is why we use this equation to define the content of information

2.3 Joint entropy

$$\begin{aligned} H(X, Y) &= \sum_x \sum_y P(x, y) \log \frac{1}{P(x, y)} \\ &= - \sum_x \sum_y P(x, y) \log P(x, y) \end{aligned}$$

2.4. Conditional entropy

$$\begin{aligned} H(Y|X) &= \sum_{x_i} P(x_i) H(Y|x_i) \\ &= \sum_x P(x) \left[- \sum_y P(y|x) \log P(y|x) \right] \\ &= - \sum_x \sum_y P(x, y) \log P(y|x) \end{aligned}$$

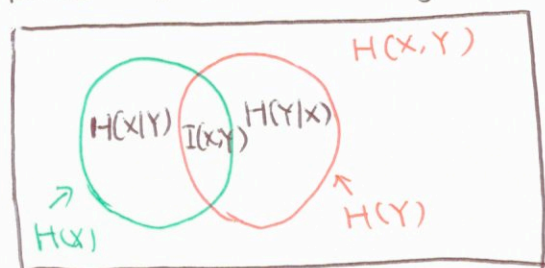
So that: $H(X, Y) = H(X) + H(Y|X)$

2.5 Mutual information

$$\text{Since } H(X, Y) = \underline{H(X) + H(Y|X)} = \underline{H(Y) + H(X|Y)}$$

(remember $p(x, y) = p(x)p(y|x) = p(y)p(x|y)$?)

then: $H(X) - H(X|Y) = H(Y) - H(Y|X) = I(X; Y)$
is named the mutual information.



Information Relations

According to a research at Cambridge

University, it doesn't matter in what order the

letters in a word are, the only important thing is

that the first and last letter be at the right place.

The rest can be a total mess and you can still read

it without problem. This is because the human

mind does not read every letter by itself, but the

word as a whole.

According to a researcher at Cambridge University, it doesn't matter in what order the letters in a word are, the only important thing is that the first and last letter be at the right place. The rest can be a total mess and you can still read it without problem. This is because the human mind does not read every letter by itself but the word as a whole.