

Words: Language Modelling By N- Grams

COMP-599

Sept 12, 2016

(Fixed) Definition of FSA

A FSA consists of:

- Q finite set of states
- Σ set of input symbols
- $q_0 \in Q$ starting state
- $\delta : Q \times \Sigma \rightarrow P(Q)$

transition function from current state and input symbol to possible next states

- $P(Q)$ is the power set of Q .
- $F \subseteq Q$ set of accepting, final states

(Fixed) Definition of FST

A FST consists of:

- Q finite set of states
- Σ, Δ sets of input symbols, output symbols
- $q_0 \in Q$ starting state
- $\delta \subseteq Q \times (\Sigma \cup \{\varepsilon\}) \times (\Delta \cup \{\varepsilon\}) \times Q$

transition relation that specifies, for current state and input symbol, the possible pairs of output symbol and next state

- $F \subseteq Q$ set of accepting, final states

Identify the above components in the previous FST

Outline

Review of last class

How words are distributed: Zipf's law

Language modelling

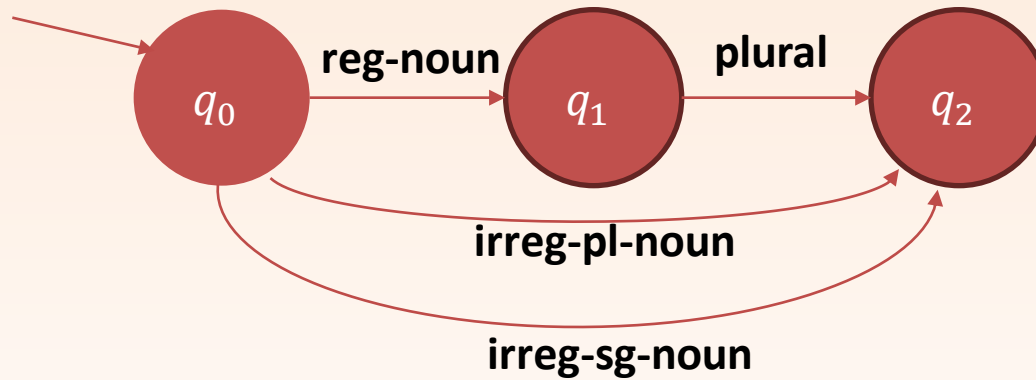
- Word sequences: N-grams

- MLE by relative frequencies

- Evaluation by cross entropy and perplexity

Last Class

FSAs and FSTs for modelling English morphology



Dealing with Regular Variations

What should be the output of the first FST that maps from the surface form to the intermediate level?

<i>jump</i>	<i>jump#</i>
<i>jumps</i>	<i>jump^s#</i>
<i>jumped</i>	<i>jump^ed#</i>
<i>jumping</i>	?
<i>chat</i>	?
<i>chats</i>	?
<i>chatted</i>	?
<i>chatting</i>	?
<i>hope</i>	?
<i>hopes</i>	?
<i>hoped</i>	?
<i>hoping</i>	?

What is a Word?

- Smallest unit that can appear in isolation

Actually not so clear cut:

Football One word, or two?

Peanut butter One word, or two?

Languages that don't separate words with spaces in writing (e.g., Chinese) – even less clear cut

- e.g., 分手信

[分][手][信] 3 words: Distribute + hand + letter ???

[分][手信] 2 words: Distribute + souvenirs

[分手][信] 2 words: Breakup + letter

[分手信] 1 word: Breakup letter

- Word segmentation a major problem in Chinese NLP

Orthographic Word, Types vs. Tokens

Convenient assumption: spaces delimit words

- Exceptions: apostrophe (e.g., 's), punctuation

Still ambiguous to ask, “How many words are there?”

e.g., the cat sat on the mat

Word **tokens**

6: *cat, mat, on, sat, the, the*

- Instances of occurrences

Word **types**

5: *cat, mat, on, sat, the*

- Kinds of words

Fuzzy Cases

Do these count as the same word type?

run, runs

happy, happily

frágment (n.), fragmént (v.)

realize, realise

We, we

srsly, seriously

Which of the above cases would be normalized by **stemming**? By **lemmatization**?

Word Frequencies

First thing we can do with words? Count them!

Term frequency:

$$TF(w, S) = \#w \text{ in corpus } S$$

- e.g., $TF(\text{cat}, \text{the cat sat on the mat}) = 1$

Relative frequency:

$$RF(w, S) = \frac{TF(w, S)}{|S|}$$

- e.g., $RF(\text{cat}, \text{the cat sat on the mat}) = \frac{1}{6}$

Corpus (n. sing.)

We need a **corpus** (pl.: **corpora**) of text to count.

Some well-known English text corpora:

Brown corpus

British National Corpus (BNC)

Wall Street Journal corpus

English Gigaword

Zipf's Law

When counting word frequencies in corpora, this is one striking effect that you'll notice:

$$f \propto \frac{1}{r}$$

Frequency of word type

Rank of word type (by frequency)

Some Empirical Counts

Rank	Word	Frequency
1	<i>the</i>	228,257,001
2	<i>to</i>	96,247,620
3	<i>of</i>	93,917,643
10	<i>for</i>	34,180,099
100	<i>most</i>	3,499,587
1,000	<i>work</i>	1,999,899
10,000	<i>planning</i>	299,996

Word counts from the English Gigaword corpus

Zipf's Law is (very) roughly true

Zipf-Mandelbrot Law

To get a better fit to the word counts we see, we can add parameters to the equation:

$$f \propto \frac{1}{r} \quad \text{means} \quad f = \frac{P}{r} \quad \text{for some } P$$

Add additional parameters ρ, B :

$$f = \frac{P}{(r + \rho)^B}$$

Or equivalently:

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

“The Long Tail”

Practical implications:

- Most word (types) are very rare!
- A small number of word (types) make up the majority of word (tokens) that you see in any corpus.
- These issues will cause problems for us in terms of designing models and evaluating their performance, as we will see.

Cross-linguistically Speaking

The parameters in the Zipf-Mandelbrot equation will differ by language

English: top handful of word types will account for most tokens. ~40% of words appear once in a corpus.

Hungarian: same number of word types account for fewer tokens

Inuktitut: ~80% of words appear only once (Langlais and Patry, 2006)

Why the disparity?

Why Count Words?

Word frequencies turn out to be very useful:

- Text classification (for genre, sentiment, authorship, ...)
- Information retrieval
- Many, many, other applications

Task we will be considering: **language modelling**

Language Modelling

Predict the next word given some context

Mary had a little _____

- *lamb* GOOD
- *accident* GOOD?
- *very* BAD
- *up* BAD

Viewed Probabilistically

Learn a probability distribution

- $P(W = w | C)$

Random variable
 W takes on a value
which is a word in the
lexicon

w represents that
value

C is the context that
we are conditioning
on

e.g.,

$$P(W = \text{"lamb"} | C = \text{"Mary had a little"}) = 0.6$$

People are often lazy:

$$P(\text{"lamb"} | \text{"Mary had a little"})$$

If any of this notation is not obvious, go review basics of probability theory now! (As in, right after class.)

Equivalently

Learn probability distribution over sequences of words

Let the context be all of the previous words. Then,

$$\begin{aligned} P(w_1 w_2 \dots w_k) \\ &= P(w_k | w_1 \dots w_{k-1}) P(w_1 \dots w_{k-1}) \quad \text{By the chain rule} \\ &= P(w_k | w_1 \dots w_{k-1}) P(w_{k-1} | w_1 \dots w_{k-2}) P(w_1 \dots w_{k-2}) \end{aligned}$$

Keep decomposing further...

$$= P(w_k | w_1 \dots w_{k-1}) \dots P(w_2 | w_1) P(w_1)$$

Example

A good language model should assign:

- higher probability to a grammatical string of English

You are wearing a fancy hat.

- lower probability to ungrammatical strings

Fancy you are hat a wearing.

Your waring a fency haat.

Note

The absolute probability from a language model isn't a good indicator of grammaticality.

- e.g., $P(\textit{artichokes intimidate zippers})$
- Likely low probability, but grammatical

Also, the length of the sentence and the rarity of the words in the sentences affect the probability

- e.g., $P(\textit{I ate the}) > P(\textit{I ate the cake})$ in most language models, but the former is clearly not a well formed sentence!

Applications

- Text prediction for mobile devices
- Automatic speech recognition (ASR)
- Machine translation

Typically, find the solution that maximizes a combination of:

1. Task-specific quality
 - ASR: acoustic model quality
 - MT: word/phrase alignment probability
2. Language model probability

Building Models

Given lots of data from the real world, we can build a **model**, which is a set of **parameters** that describes the data, and can be used to **predict** or **infer** future or unseen data.

e.g.,

Task: language modelling

Model: a probability distribution, $P(W = w | C)$

Parameters: the parameters to this probability distribution

Application: tell us how likely it is to observe w_N given its context

Steps

1. Gather a large, representative training corpus
2. Learn the parameters from the corpus to build the model
3. Once the model is fixed, use the model to evaluate on testing data

Steps

1. Gather a large, representative training corpus
2. **Learn the parameters from the corpus to build the model**
3. Once the model is fixed, use the model to evaluate on testing data

Learning the Model

How do we actually learn the parameters to $P(W = w | C)$ given training data?

Need to:

- Specify exactly what the context of a word is
- Use corpus counts to derive the parameter values

N-grams

Make a **conditional independence assumption** to make the job of learning the probability distribution easier.

- Context = the previous N-1 words

Common choices: N is between 1 and 3

Unigram model

$$P(w_N|C) = P(w_N)$$

Bigram model

$$P(w_N|C) = P(w_N|w_{N-1})$$

Trigram model

$$P(w_N|C) = P(w_N|w_{N-1}, w_{N-2})$$

Deriving Parameters from Counts

Simplest method: count N-gram frequencies, then divide by the total count

e.g.,

Unigram: $P(\textit{cats}) = \text{Count}(\textit{cats}) / \text{Count}(\text{all words in corpus})$

Bigram: $P(\textit{the cats}) = \text{Count}(\textit{the cats}) / \text{Count}(\textit{the})$

Trigram: $P(\textit{feed the cats}) = ?$

These are the **maximum likelihood estimates (MLE)**.

Exercise

Come up with the MLE estimate of a unigram and a bigram language model using the following sentence as training data:

A statistical language model is a probability distribution over sequences of words.

N-grams as Linguistic Knowledge

N-grams can crudely capture some linguistic knowledge and even facts about the world

- e.g., $P(\text{English} | \text{want}) = 0.0011$
 $P(\text{Chinese} | \text{want}) = 0.0065$

World knowledge:
culinary preferences?

$$P(\text{to} | \text{want}) = 0.66$$

$$P(\text{eat} | \text{to}) = 0.28$$

$$P(\text{food} | \text{to}) = 0$$

Syntax

$$P(I | \langle \text{start-of-sentence} \rangle) = 0.25$$

Discourse

Steps

1. Gather a large, representative training corpus
2. Learn the parameters from the corpus to build the model
3. **Once the model is fixed, use the model to evaluate on testing data**

Training and Testing Data

After training a model, we need to evaluate it on unseen data that the model has not been exposed to.

- We are testing the model's ability to generalize.
- More on this topic next class

Given a corpus, how is the data usually split?

Training data: often 60-90% of the available data

Testing data: often 10-20% of the available data

There is often also a **development** or **validation** data set, for deciding between different versions of a model.

Cross Validation

k-fold cross validation: splitting data into k partitions or folds; iteratively test on each after training on the rest

e.g., 3-fold CV: dataset1 dataset2 dataset3

train on {2,3}, test on 1

train on {1,3}, test on 2

train on {1,2}, test on 3

Average results from above folds

- CV is often used if the corpus is small

Evaluation Measures

Likelihood of generating the test corpus

i.e., $P(\text{test_corpus}; \theta)$, where θ represents the parameters learned by training our LM on the training data

Intuition: a good language model should give a high probability of generating some new, valid English text.

Absolute number is not very meaningful—this can only be used to compare the quality of different language models!

Unwieldy because of small values, so not actually used in the literature. Alternatives to likelihood:

Cross-entropy

Perplexity

Basic Information Theory

Consider some random variable X , distributed according to some probability distribution.

We can define information in terms of how much certainty we gain from knowing the value of X .

Rank the following in terms of how much information we gain by knowing its value:

- Fair coin flip

- An unfair coin flip where we get tails $\frac{3}{4}$ of the time

- A very unfair coin that always comes up heads

Likely vs Unlikely Outcomes

Observing a likely outcome – less information gained

Intuition: you kinda knew it would happen anyway

- e.g., observing the word *the*

Observing a rare outcome: more information gained!

Intuition: it's a bit surprising to see something unusual!

- e.g., observing the word *armadillo*

Formal definition of information in bits:

$$I(x) = \log_2\left(\frac{1}{P(x)}\right)$$

Minimum number of bits needed to communicate some outcome x

Entropy

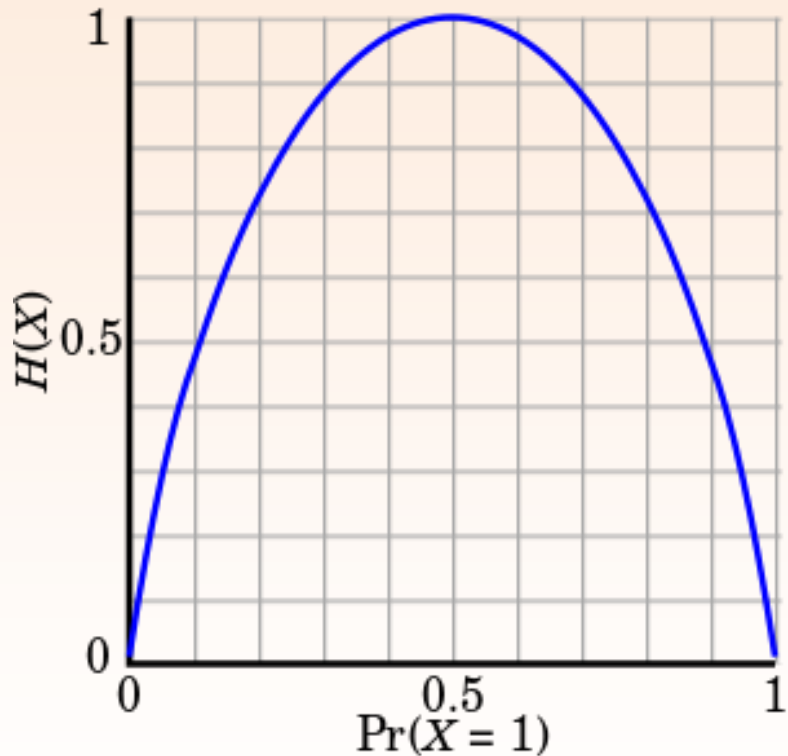
The expected amount of information we get from observing a random variable.

Let a discrete random variable be drawn from distribution p take on one of k possible values with probabilities $p_1 \dots p_k$

$$\begin{aligned} H(p) &= \sum_{i=1}^k p_i I(x_i) \\ &= \sum_{i=1}^k p_i \log_2 \frac{1}{p_i} \\ &= - \sum_{i=1}^k p_i \log_2 p_i \end{aligned}$$

Entropy Example

Plot of entropy vs. coin toss “fairness”



Maximum fairness =
maximum entropy

Completely biased =
minimum entropy

Image source: Wikipedia, by Brona and Alessio Damato

Cross Entropy

Entropy is the minimum number of bits needed to communicate some message, *if we know what probability distribution the message is drawn from.*

Cross entropy is for when we don't know.

e.g., language is drawn from some true distribution, the language model we train is an approximation of it

$$H(p, q) = - \sum_{i=1}^k p_i \log_2 q_i$$

p : “true” distribution

q : model distribution

Estimating Cross Entropy

When evaluating our LM, we assume the test data is a good representative of language drawn from p .

So, we estimate cross entropy to be:

$$H(p, q) = -\frac{1}{N} \log_2 q(w_1 \dots w_N)$$

True language distribution, which we don't have access to.

Language model under evaluation

Size of test corpus in number of tokens

The words in the test corpus

Perplexity

Cross entropy gives us a number in bits, which is sometimes hard to read. Perplexity makes this easier.

$$\text{Perplexity}(p, q) = 2^{H(p, q)}$$