University of Ljubljana, Faculty of Computer and Information Science

N-gram language models



my alarm	Clark code circle shute clock	dia soil raid risk visit did	rout hot riot not must
Wake me wake me	up th up th ta th tid	ین کر کر ai m xis ha is ru er mo lo	oving aving nning orning oving

Prof Dr Marko Robnik-Šikonja

Natural Language Processing, Edition 2025

Contents

- language models
- n-grams (still used in the evaluation measures)

mostly based on Jurafsky & Martin, 3^{rd} edition, read Chapter 3.1 - 3.4

Predicting words

• The water of Walden Pond is beautifully ...

blue green clear

*refrigerator*that

Language Models

• Systems that can predict upcoming words

- Can assign a probability to each potential next word
- Can assign a probability to a whole sentence

Why word prediction?

It's a helpful part of language tasks

Grammar or spell checking
 Their are two midterms
 Their There a
 Everything has improve
 Everything
 Everything
 Everything

Their There are two midterms Everything has improve improved

Speech recognition
 I will be back soonish

I will be bassoon dish

Probabilistic Language Models

• The goal: assign a probability to a sentence

- Machine Translation:
 - P(high winds tonight) > P(large winds tonight)
- Spell Correction
- Why?
- The office is about fifteen **minuets** from my house
 - P(about fifteen minutes from) > P(about fifteen minuets from)
- Speech Recognition
 - P(I saw a van) >> P(eyes awe of an)
- + Summarization, question-answering, etc., etc.!!

Probabilistic Language Modeling

 Goal: compute the probability of a sentence or sequence of words:

 $P(W) = P(W_1, W_2, W_3, W_4, W_5...W_n)$

- Related task: probability of an upcoming word: P(w₅|w₁,w₂,w₃,w₄)
- A model that computes either of these:

P(W) or $P(w_n|w_1, w_2...w_{n-1})$ is called a **language model**.

- Another suitable name would be: the grammar model
- But language model or LM is standard

How to statistically compute P(W)

• How to compute this joint probability:

```
P(its, water, is, so, transparent, that)
```

• Intuition: let's rely on the Chain Rule of Probability

Reminder: The Chain Rule

- Recall the definition of conditional probabilities
 p(B|A) = P(A,B)/P(A) Rewriting: P(A,B) = P(A)P(B|A)
- More variables:

P(A,B,C,D) = P(A)P(B|A)P(C|A,B)P(D|A,B,C)

• The Chain Rule in General

 $P(x_1, x_2, x_3, ..., x_n) = P(x_1)P(x_2 | x_1)P(x_3 | x_1, x_2)...P(x_n | x_1, ..., x_{n-1})$

The Chain Rule applied to compute joint probability of words in sentence

$$P(w_1 w_2 \dots w_n) = \prod_i P(w_i \mid w_1 w_2 \dots w_{i-1})$$

P("its water is so transparent") = P(its) × P(water | its) × P(is | its water)

× P(so|its water is) × P(transparent|its water is so)

How to estimate these probabilities

• Could we just count and divide?

P(the | its water is so transparent that) =

Count(its water is so transparent that the)

Count(its water is so transparent that)

- No! Too many possible sentences!
- We'll never see enough data for estimating these

Markov Assumption

- The memory is short
- First order Markov assumption



Andrei Markov

 $P(\text{the} | \text{its water is so transparent that}) \gg P(\text{the} | \text{that})$

• The second order Markov assumption *P*(the | its water is so transparent that) » *P*(the | transparent that) Using Markov assumption of order k

$$P(w_1w_2...w_n) \approx \prod_i P(w_i \mid w_{i-k} \dots w_{i-1})$$

 In other words, we approximate each component in the product

$$P(w_i \mid w_1 w_2 ... w_{i-1}) \approx P(w_i \mid w_{i-k} ... w_{i-1})$$

Simplest case: Unigram model

$$P(w_1 w_2 \dots w_n) \approx \prod_i P(w_i)$$

Some automatically generated sentences from a unigram model

fifth, an, of, futures, the, an, incorporated, a, a, the, inflation, most, dollars, quarter, in, is, mass

thrift, did, eighty, said, hard, 'm, july, bullish

that, or, limited, the

Bigram model

Condition on the previous word:

$$P(w_i | w_1 w_2 ... w_{i-1}) \approx P(w_i | w_{i-1})$$

texaco, rose, one, in, this, issue, is, pursuing, growth, in, a, boiler, house, said, mr., gurria, mexico, 's, motion, control, proposal, without, permission, from, five, hundred, fifty, five, yen

outside, new, car, parking, lot, of, the, agreement, reached this, would, be, a, record, november

Estimating bigram probabilities

• The Maximum Likelihood Estimate

$$P(w_i | w_{i-1}) = \frac{count(w_{i-1}, w_i)}{count(w_{i-1})}$$

$$P(w_i | w_{i-1}) = \frac{C(w_{i-1}, w_i)}{C(w_{i-1})}$$

An example

<s> I am Sam </s>

<s> Sam I am </s>

<s> I do not like green eggs and ham </s>

$$P(w_i | w_{i-1}) = \frac{c(w_{i-1}, w_i)}{c(w_{i-1})}$$

$$P(I | < s >) = \frac{2}{3} = .67 \qquad P(Sam | < s >) = \frac{1}{3} = .33 \qquad P(am | I) = \frac{2}{3} = .67 P(| Sam) = \frac{1}{2} = 0.5 \qquad P(Sam | am) = \frac{1}{2} = .5 \qquad P(do | I) = \frac{1}{3} = .33$$

N-gram models

- We can extend to trigrams, 4-grams, 5-grams
- In general this is an insufficient model of language
 - because language has long-distance dependencies:

"The computer(s) which I had just put into the machine room on the fifth floor is (are) crashing."

• N-gram models are better in English than in Slovene and many other languages. Why?

Example: Restaurant sentences

- can you tell me about any good cantonese restaurants close by
- mid priced thai food is what i'm looking for
- tell me about chez panisse
- can you give me a listing of the kinds of food that are available
- i'm looking for a good place to eat breakfast
- when is caffe venezia open during the day

Raw bigram counts

• Out of 9222 sentences

	i	want	to	eat	chinese	food	lunch	spend
i	5	827	0	9	0	0	0	2
want	2	0	608	1	6	6	5	1
to	2	0	4	686	2	0	6	211
eat	0	0	2	0	16	2	42	0
chinese	1	0	0	0	0	82	1	0
food	15	0	15	0	1	4	0	0
lunch	2	0	0	0	0	1	0	0
spend	1	0	1	0	0	0	0	0

Raw bigram probabilities

• Normalize by unigrams:

• Result:

i	want	to	eat	chinese	food	lunch	spend
2533	927	2417	746	158	1093	341	278

	i	want	to	eat	chinese	food	lunch	spend
i	0.002	0.33	0	0.0036	0	0	0	0.00079
want	0.0022	0	0.66	0.0011	0.0065	0.0065	0.0054	0.0011
to	0.00083	0	0.0017	0.28	0.00083	0	0.0025	0.087
eat	0	0	0.0027	0	0.021	0.0027	0.056	0
chinese	0.0063	0	0	0	0	0.52	0.0063	0
food	0.014	0	0.014	0	0.00092	0.0037	0	0
lunch	0.0059	0	0	0	0	0.0029	0	0
spend	0.0036	0	0.0036	0	0	0	0	0

Bigram estimates of sentence probabilities

- P(<s> I want english food </s>) = P(I|<s>)
 - × P(want|I)
 - × P(english|want)
 - × P(food|english)
 - \times P(</s>|food)
 - = .000031

What kinds of knowledge bigram LM contains?

- P(english | want) = .0011
- P(chinese | want) = .0065
- P(to | want) = .66
- P(eat | to) = .28
- P(food | to) = 0
- P(want | spend) = 0
- P (i | <s>) = .25

Dealing with scale in large n-grams

- •LM probabilities are stored and computed in log format, i.e. **log probabilities**
- •This avoids underflow from multiplying many small numbers

$$\log(p_1 f_2 f_3 f_2) = \log p_1 + \log p_2 + \log p_3 + \log p_4$$

If we need probabilities we can do one exp at the end

$$p_1 \times p_2 \times p_3 \times p_4 = \exp(\log p_1 + \log p_2 + \log p_3 + \log p_4)$$

Larger ngrams

- 4-grams, 5-grams
- Large datasets of large n-grams have been released
 - N-grams from Corpus of Contemporary American English (COCA) 1 billion words (Davies 2020)
 - Google Web 5-grams (Franz and Brants 2006) 1 trillion words)
 - Efficiency: quantize probabilities to 4-8 bits instead of 8-byte float

Newest model: infini-grams (∞-grams) (Liu et al 2024)

• No precomputing! Instead, store 5 trillion words of web text in **suffix arrays**. Can compute n-gram probabilities with any n!

N-gram LM Toolkits

- •SRILM
 - •http://www.speech.sri.com/projects/srilm/
- •KenLM
 - <u>https://kheafield.com/code/kenlm/</u>

Language Modeling Tools

- are ngram language models still useful?
- yes, e.g., in speech processing
- mostly replaced by neural LMs
- many variants of adapted neural LMs exist, e.g., word2vec, fastText, BERT, GPT

Evaluation: How good is our model?

- Does our language model prefer good sentences to bad ones?
 - Assign higher probability to "real" or "frequently observed" sentences
 - Than "ungrammatical" or "rarely observed" sentences?
- We train parameters of our model on a training set.
- We test the model's performance on data we haven't seen.
 - A **test set** is an unseen dataset that is different from our training set, totally unused.
 - An evaluation metric tells us how well our model does on the test set.
- Two types of evaluation
 - intrinsic (internal)
 - extrinsic (external, on a downstream task)

Extrinsic evaluation of N-gram models

- Best evaluation for comparing models A and B
 - Use each model in a task
 - spelling corrector, speech recognizer, MT system
 - Run the task, get an accuracy for A and for B
 - How many misspelled words corrected properly
 - How many words translated correctly
 - Compare accuracy for A and B

Intrinsic (in-vitro) evaluation

- Extrinsic evaluation not always possible
 - Expensive, time-consuming
 - Doesn't always generalize to other applications
- Intrinsic evaluation: perplexity
 - Directly measures language model performance at predicting words.
 - Doesn't necessarily correspond with real application performance
 - But gives us a single general metric for language models
 - Useful for large language models (LLMs) as well as n-grams

Training sets and test sets

We train parameters of our model on a training set.

We test the model's performance on data we haven't seen.

- A test set is an unseen dataset; different from training set.
 - Intuition: we want to measure generalization to unseen data
- An evaluation metric (like perplexity) tells us how well our model does on the test set.

Choosing training and test sets

- If we're building an LM for a specific task
 - The test set should reflect the task language we want to use the model for
- If we're building a general-purpose model
 - We'll need lots of different kinds of training data
 - We don't want the training set or the test set to be just from one domain or author or language.

Training on the test set

We can't allow test sentences into the training set

- Or else the LM will assign that sentence an artificially high probability when we see it in the test set
- And hence assign the whole test set a falsely high probability.
- Making the LM look better than it really is

This is called "Training on the test set"

Bad science, bad practice!

Dev sets

- If we test on the test set many times we might implicitly tune to its characteristics
 - Noticing which changes make the model better.
- So we run on the test set only once, or a few times
- That means we need a third dataset:
 - A development test set or, devset.
 - We test our LM on the devset until the very end
 - And then test our LM on the **test set** once

Intuition of Perplexity

- The Shannon Game:
 - How well can we predict the next word?

I always order pizza with cheese and

The 33rd President of the US was

I saw a _

- Unigrams are terrible at this game. (Why?)
- A better model of a text
 - is one which assigns a higher probability to the word that actually occurs

mushrooms 0.1
pepperoni 0.1
anchovies 0.01
....
fried rice 0.0001
....
and 1e-100

Perplexity

The best language model is one that best predicts an unseen test set

• Gives the highest P(sentence)

Perplexity is the inverse probability of the test set, normalized by the number of words:

$$PP(W) = P(w_1 w_2 ... w_N)^{-\frac{1}{N}}$$

$$= \sqrt[N]{\frac{1}{P(w_1w_2...w_N)}}$$

Chain rule:

$$PP(W) = \sqrt[N]{\prod_{i=1}^{N} \frac{1}{P(w_i|w_1...w_{i-1})}}$$
For bigrams:

$$PP(W) = \sqrt[N]{\prod_{i=1}^{N} \frac{1}{P(w_i|w_{i-1})}}$$

Minimizing perplexity is the same as maximizing probability

Perplexity example

- Let us suppose a sentence consisting of random digits
- What is the perplexity of this sentence according to a model that assign P=1/10 to each digit?

$$PP(W) = P(w_1 w_2 \dots w_N)^{-\frac{1}{N}} \\ = (\frac{1}{10}^N)^{-\frac{1}{N}} \\ = \frac{1}{10}^{-1} \\ = 10$$

Lower perplexity = better model

• Training 38 million words, test 1.5 million words, WSJ

N-gram Order	Unigram	Bigram	Trigram
Perplexity	962	170	109

The Shannon Visualization Method

• Choose a random bigram

(<s>, w) according to its
probability

- Now choose a random bigram (w, x) according to its probability
- And so on until we choose </s>
- Then string the words together I want to eat Chinese food

 $\langle s \rangle I$

```
I want
want to
to eat
eat Chinese
Chinese food
food </s>
```

Approximating Shakespeare

1 gram	 To him swallowed confess hear both. Which. Of save on trail for are ay device and rote life have Hill he late speaks; or! a more to leg less first you enter
2 gram	Why dost stand forth thy canopy, forsooth; he is this palpable hit the King Henry. Live king. Follow.What means, sir. I confess she? then all sorts, he is trim, captain.
3 gram	–Fly, and will rid me these news of price. Therefore the sadness of parting, as they say, 'tis done.–This shall forbid it should be branded, if renown made it empty.
4 gram	 –King Henry. What! I will go seek the traitor Gloucester. Exeunt some of the watch. A great banquet serv'd in; –It cannot be but so.

Shakespeare as corpus

- N=884,647 tokens, |V|=29,066
- Shakespeare produced 300,000 bigram types out of |V|²= 844 million possible bigrams.
 - So 99.96% of the possible bigrams were never seen (have zero entries in the table)
- Quadrigrams are even worse: What's coming out looks like Shakespeare because it *is* Shakespeare

The Wall Street Journal

Months the my and issue of year foreign new exchange's september were recession exchange new endorsed a acquire to six executives
 Last December through the way to preserve the Hudson corporation N.

B. E. C. Taylor would seem to complete the major central planners one point five percent of U. S. E. has already old M. X. corporation of living on information such as more frequently fishing to keep her

3 gram They also point to ninety nine point six billion dollars from two hundred four oh six three percent of the rates of interest stores as Mexico and Brazil on market conditions

What is the source of these random 3-gram sentences?

- They also point to ninety nine point six billion dollars from two hundred four oh six three percent of the rates of interest stores as Mexico and gram Brazil on market conditions
- This shall forbid it should be branded, if renown made it empty.
- "You are uniformly charming!" cried he, with a smile of associating and now and then I bowed and they perceived a chaise and four to wish for.

The perils of overfitting

- N-grams only work well for word prediction if the test corpus looks like the training corpus
- In real life, it often doesn't
- We need to train robust models that generalize!
- One kind of (outdated) generalization: Zeros!
 - Things that don't ever occur in the training set
 - But occur in the test set
- In practice, we use neural language models