

06 – Language modelling

IA161 Advanced Techniques of Natural Language Processing

V. Baisa

NLP Centre, FI MU, Brno

November 2, 2015

- 1 Introduction to Language Modeling
- 2 N-grams
- 3 Evaluation of Language Models
- 4 Neural Networks and Language Modeling
- 5 State-of-the-art results
- 6 Practical part: generating random texts

Language models—what are they good for?

- predicting words
- generating text
- statistical machine translation
- automatic speech recognition
- optical character recognition

Predicting words

Do you speak ...

Would you be so ...

Statistical machine ...

Faculty of Informatics, Masaryk ...

WWII has ended in ...

In the town where I was ...

Lord of the ...

Generating text

Describes without errors



A person riding a motorcycle on a dirt road.

Describes with minor errors



Two dogs play in the grass.

Somewhat related to the image



A skateboarder does a trick on a ramp.

Unrelated to the image



A dog is jumping to catch a frisbee.



A group of young people playing a game of frisbee.



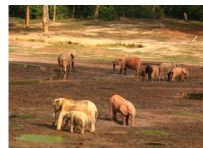
Two hockey players are fighting over the puck.



A little girl in a pink hat is blowing bubbles.



A refrigerator filled with lots of food and drinks.



A herd of elephants walking across a dry grass field.



A close up of a cat laying on a couch.

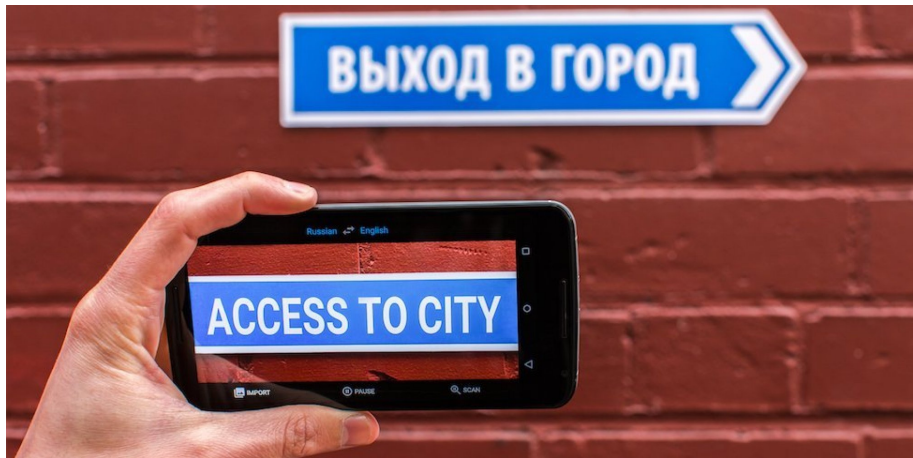


A red motorcycle parked on the side of the road.



A yellow school bus parked in a parking lot.

MT + OCR



Language models – probability of a sentence

- LM is a probability distribution over all possible word sequences.
- A sentence is a word sequence.
- What is the probability of utterance of s ?

Probability of sentence

$p_{LM}(\text{včera jsem jel do Brna})$

$p_{LM}(\text{včera jel do Brna jsem})$

$p_{LM}(\text{jel jsem včera do Brna})$

N-gram models

- intuitive idea
- classical approach
- simple implementation

Randomly generated text

“Jsi nebylo vidět vteřin přestal po schodech se dal do deníku a položili se táhl ji viděl na konci místnosti 101,” řekl důstojník.

Hungarian

A társaság kötelezettségeiért kapta a középkori temploma az volt, hogy a felhasználók az adottságai, a felhasználó azonosítása az egyesület alapszabályát.

N-gram models, naïve approach

$$W = w_1, w_2, \dots, w_n$$

$$p(W) = \prod_i p(w_i | w_1 \dots w_{i-1})$$

Markov's assumption

$$p(W) = \prod_i p(w_i | w_{i-2}, w_{i-1})$$

$$p(\textit{this is a sentence}) = p(\textit{this}) \times p(\textit{is} | \textit{this}) \times p(\textit{a} | \textit{this}, \textit{is}) \times p(\textit{sentence} | \textit{is}, \textit{a})$$

$$p(\textit{a} | \textit{this}, \textit{is}) = \frac{|\textit{this is a}|}{|\textit{this is}|}$$

Sparse data problem.

Computing, LM probabilities estimation

Trigram model uses 2 preceding words for probability learning. Using **maximum-likelihood estimation**:

$$p(w_3|w_1, w_2) = \frac{\text{count}(w_1, w_2, w_3)}{\sum_w \text{count}(w_1, w_2, w)}$$

trigram: (*the, green, w*) (1748)

<i>w</i>	count	$p(w)$
paper	801	0.458
group	640	0.367
light	110	0.063
party	27	0.015
ecu	21	0.012

Large LM – n-gram counts

How many unique n-grams are in a corpus?

order	unique	singletons
unigram	86 700	33 447 (38,6 %)
bigram	1 948 935	1 132 844 (58,1 %)
trigram	8 092 798	6 022 286 (74,4 %)
4-gram	15 303 847	13 081 621 (85,5 %)
5-gram	19 882 175	18 324 577 (92,2 %)

Taken from Europarl with 30 mil. tokens.

Language models smoothing

Issue: if an n-gram is missing in the data but is in $w \rightarrow p(w) = 0$.

We need to distinguish p also for *unseen data*. This must hold:

$$\forall w. p(w) > 0$$

The issue is more serious for high-order models.

Smoothing: attempt to amend real counts of n-grams to expected counts in any data (different corpora).

Add-one smoothing (Laplace)

Maximum likelihood estimation assigns p based on

$$p = \frac{c}{n}$$

Add-one smoothing uses

$$p = \frac{c + 1}{n + v}$$

where v is amount of all possible n-grams. That is quite inaccurate since all permutations might outnumber real (possible) n-grams by several magnitudes.

Europarl has (unique) 139,000 words, so 19 billion possible bigrams. In fact it has only 53 mil. tokens, so maximally 53 mil. bigrams.

This smoothing overvalues unseen n-grams.

Add- α smoothing

We won't add 1, but α . This can be estimated for the smoothing to be the most just and balanced.

$$p = \frac{c + \alpha}{n + \alpha v}$$

α can be obtained experimentally: we can try several different values and find the best one.

Usually it is very small (0.0001).

Deleted estimation

We can find unseen n-grams in another corpus. N-grams contained in one of them and not in the other help us to estimate general amount of unseen n-grams.

E.g. bigrams not occurring in a training corpus but present in the other corpus million times (given the amount of all possible bigrams equals 7.5 billions) will occur approx.

$$\frac{10^6}{7.5 \times 10^9} = 0.00013 \times$$

Good–Turing smoothing

We need to amend occurrence counts (frequencies) of n-grams in a corpus in such a way they correspond to general occurrence in texts. We use *frequency of frequencies*: number of various n-grams which occur $n \times$.

We use frequency of hapax legomena (singletons in data) to estimate unseen data.

$$r^* = (r + 1) \frac{N_{r+1}}{N_r}$$

Especially for n-grams not in our corpus we have

$$r_0^* = (0 + 1) \frac{N_1}{N_0} = 0.00015$$

where $N_1 = 1.1 \times 10^6$ a $N_0 = 7.5 \times 10^9$ (Europarl).

Example of Good–Turing smoothing (Europarl)

r	FF	r^*
0	7,514,941,065	0.00015
1	1,132,844	0.46539
2	263,611	1.40679
3	123,615	2.38767
4	73,788	3.33753
5	49,254	4.36967
6	35,869	5.32929
8	21,693	7.43798
10	14,880	9.31304
20	4,546	19.54487

Comparing smoothing methods (Europarl)

method	perplexity
add-one	382.2
add- α	113.2
deleted est.	113.4
Good-Turing	112.9

Interpolation and back-off

Previous methods treated all unseen n-grams the same. Consider trigrams

nádherná červená řepa

nádherná červená mrkev

Despite we don't have any of these in our training data, the former trigram should be probably more probable.

We will use probability of lower order models, for which we have necessary data:

červená řepa

červená mrkev

nádherná červená

Interpolation

$$p_I(w_3|w_1w_2) = \lambda_1 p(w_3) \times \lambda_2 p(w_3|w_2) \times \lambda_3 p(w_3|w_1w_2)$$

If we have enough data we can trust higher order models more and assign a higher significance to corresponding n-grams.

p_I is probability distribution, thus this must hold:

$$\begin{aligned}\forall \lambda_n : 0 \leq \lambda_n \leq 1 \\ \sum_n \lambda_n = 1\end{aligned}$$

Quality and comparison of LMs

We need to compare quality of various LM (various orders, various data, smoothing techniques etc.)

2 approaches: extrinsic (WER, MT, ASR, OCR) and intrinsic (perplexity) evaluation.

A good LM should assign a higher probability to a good (looking) text than to an incorrect text. For a fixed testing text we can compare various LMs.

Cross-entropy

$$\begin{aligned} H(p_{LM}) &= -\frac{1}{n} \log p_{LM}(w_1, w_2, \dots w_n) \\ &= -\frac{1}{n} \sum_{i=1}^n \log p_{LM}(w_i | w_1, \dots w_{i-1}) \end{aligned}$$

Cross-entropy is average value of negative logarithms of words probabilities in testing text. It corresponds to a measure of uncertainty of a probability distribution. **The lower the better.**

A good LM should reach entropy close to real entropy of language. That can not be measured but quite reliable estimates do exist, e.g. Shannon's game. For English, entropy is estimated to approx. 1.3 bit per letter.

Perplexity

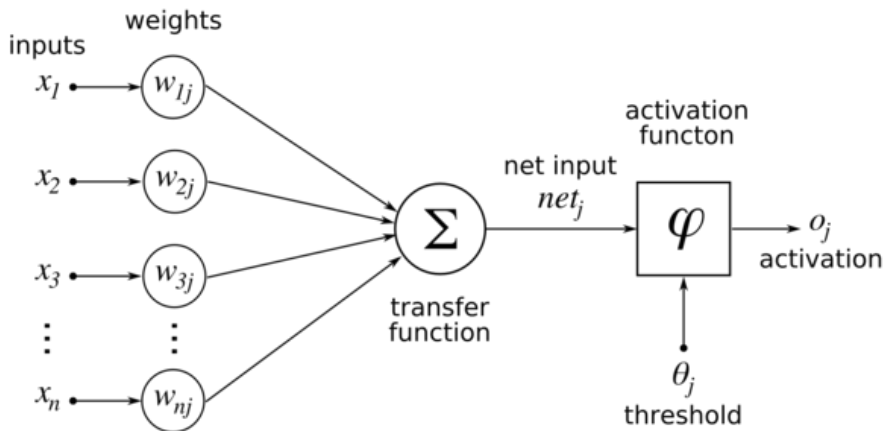
$$PP = 2^{H(p_{LM})}$$

Perplexity is simple transformation of cross-entropy.

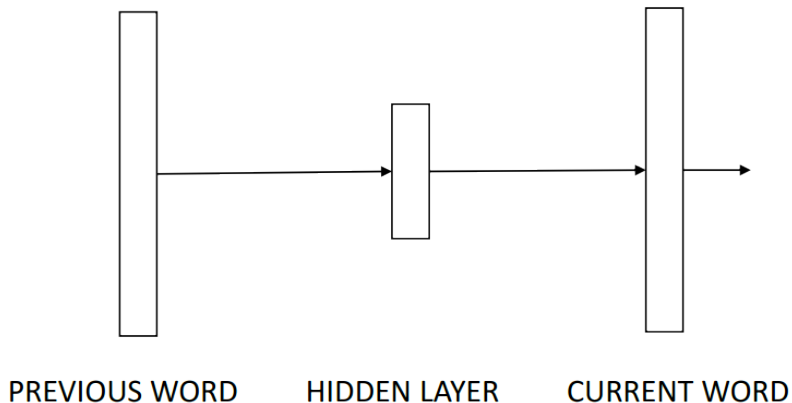
A good LM should not waste p for improbable phenomena.

The lower entropy, the better \rightarrow the lower perplexity, the better.

Neuron in artificial neural network



Basic NN

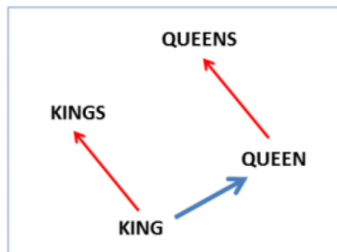
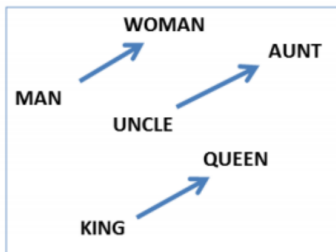


Bigram neural language model.

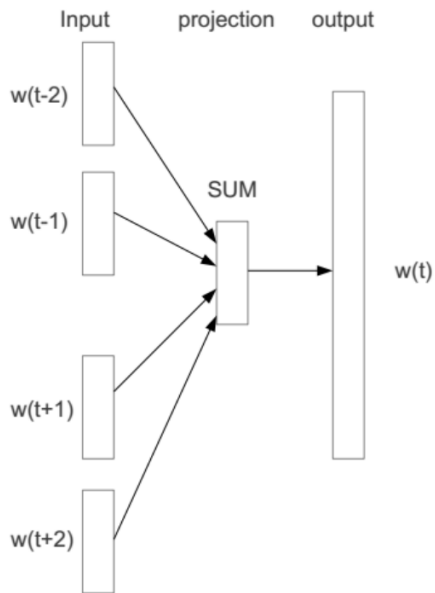
One-hot representation of words: [0 0 0 0 0 0 0 1 0 0 0 0]

Distributional Representation of Words

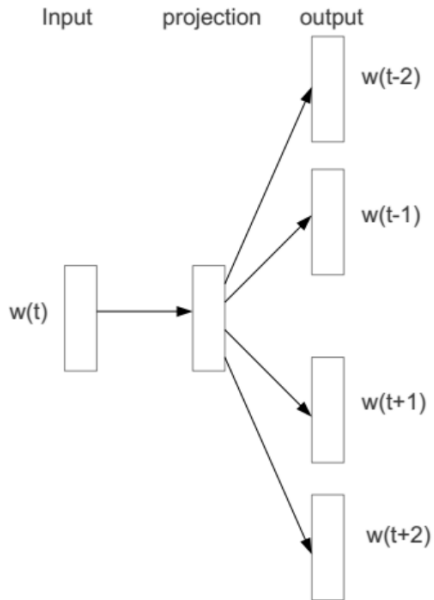
- goal: more compact representation of vectors
- limited dimensionality (500–1000)
- [Mikolov et al., 2013b]
- word vectors capture many linguistic properties (gender, tense, plurality, even semantic concepts like “capital city of”)



CBOW: Continuous Bag of Words



Skip gram



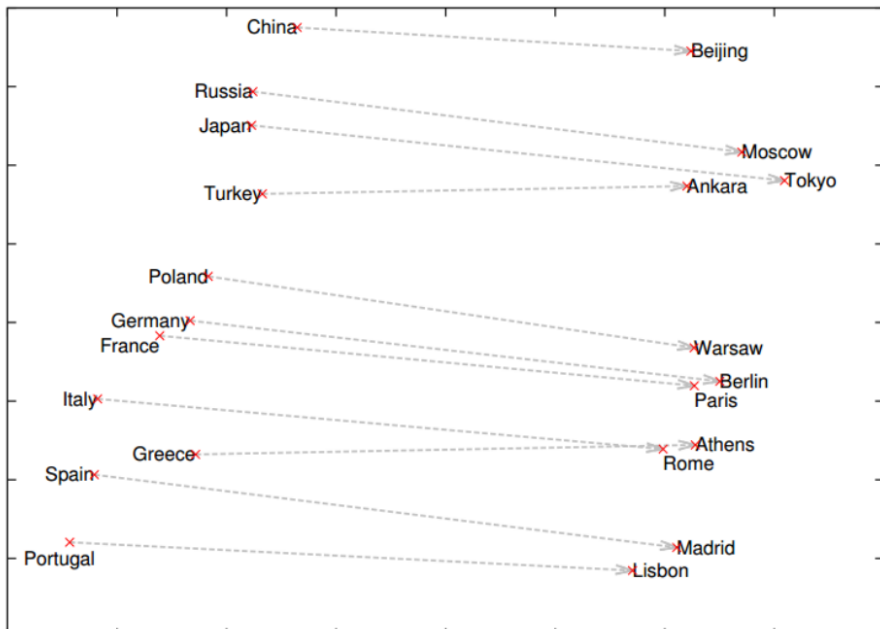
Features: nearest neighbors

	Redmond	Havel	graffiti	capitulate
Collobert NNLM	conyers lubbock keene	plauen dzerzhinsky osterreich	cheesecake gossip dioramas	abdicate accede rearm
Turian NNLM	McCarthy Alston Cousins	Jewell Arzu Ovitz	gunfire emotion impunity	- - -
Mnih NNLM	Podhurst Harlang Agarwal	Pontiff Pinochet Rodionov	anaesthetics monkeys Jews	Mavericks planning hesitated
Skip-gram (phrases)	Redmond Wash. Redmond Washington Microsoft	Vaclav Havel president Vaclav Havel Velvet Revolution	spray paint grafitti taggers	capitulation capitulated capitulating

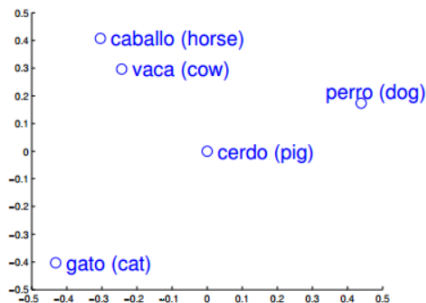
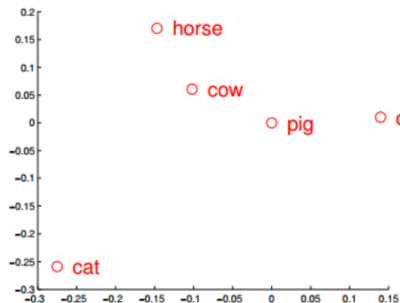
Features: vector arithmetics I

<i>Expression</i>	<i>Nearest token</i>
Paris - France + Italy	Rome
bigger - big + cold	colder
sushi - Japan + Germany	bratwurst
Cu - copper + gold	Au
Windows - Microsoft + Google	Android
Montreal Canadiens - Montreal + Toronto	Toronto Maple Leafs

Features: vector arithmetics II



Features: vector arithmetics III



Best models

Model	Num. Params [billions]	Training Time		Perplexity
		[hours]	[CPUs]	
Interpolated KN 5-gram, 1.1B n-grams (KN)	1.76	3	100	67.6
Katz 5-gram, 1.1B n-grams	1.74	2	100	79.9
Stupid Backoff 5-gram (SBO)	1.13	0.4	200	87.9
Interpolated KN 5-gram, 15M n-grams	0.03	3	100	243.2
Katz 5-gram, 15M n-grams	0.03	2	100	127.5
Binary MaxEnt 5-gram (n-gram features)	1.13	1	5000	115.4
Binary MaxEnt 5-gram (n-gram + skip-1 features)	1.8	1.25	5000	107.1
Hierarchical Softmax MaxEnt 4-gram (HME)	6	3	1	101.3
Recurrent NN-256 + MaxEnt 9-gram	20	60	24	58.3
Recurrent NN-512 + MaxEnt 9-gram	20	120	24	54.5
Recurrent NN-1024 + MaxEnt 9-gram	20	240	24	51.3

References I



Bahl, L. R., Brown, P. F., De Souza, P. V., and Mercer, R. L. (1989).
A tree-based statistical language model for natural language speech
recognition.

Acoustics, Speech and Signal Processing, IEEE Transactions on,
37(7):1001–1008.



Bengio, Y., Ducharme, R., Vincent, P., and Janvin, C. (2003).
A neural probabilistic language model.

The Journal of Machine Learning Research, 3:1137–1155.



Chelba, C., Mikolov, T., Schuster, M., Ge, Q., Brants, T., Koehn, P.,
and Robinson, T. (2013).

One Billion Word Benchmark for Measuring Progress in Statistical
Language Modeling.

ArXiv e-prints.

References II



Heafield, K., Pouzyrevsky, I., Clark, J. H., and Koehn, P. (2013). Scalable modified kneser-ney language model estimation. In *Proceedings of the 51st Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)*, pages 690–696, Sofia, Bulgaria. Association for Computational Linguistics.



Mikolov, T. (2012). Statistical language models based on neural networks. *Presentation at Google, Mountain View, 2nd April*.



Mikolov, T., Chen, K., Corrado, G., and Dean, J. (2013a). Efficient estimation of word representations in vector space. *arXiv preprint arXiv:1301.3781*.



Mikolov, T., Yih, W.-t., and Zweig, G. (2013b). Linguistic regularities in continuous space word representations. In *HLT-NAACL*, pages 746–751.

References III



Vinyals, O., Toshev, A., Bengio, S., and Erhan, D. (2014).
Show and Tell: A Neural Image Caption Generator.
ArXiv e-prints.

Character-based Language Model

The goal of the CBLM is:

- get rid of tokenization and any other rule-based processing
- limit language units not by length but by frequency
- do not work with word units, use bytes as they are universal for all languages represented in computers

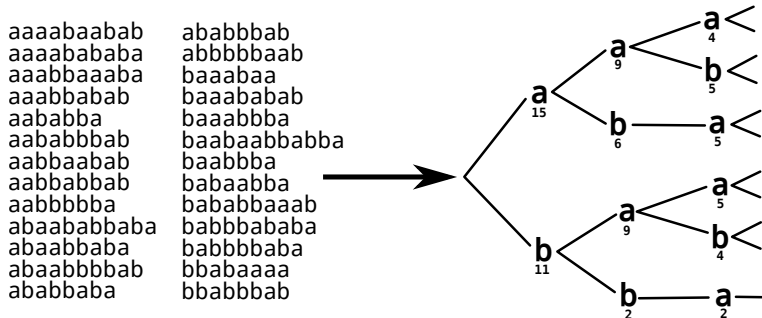
Suffix tree example

Input is any plain text data, one sentence per line.

I	suffix	SA	sorted suffix	LCP
0	popocatepetl	5	atepetl	0
1	opocatepetl	4	catepetl	0
2	pocatepetl	7	epetl	0
3	ocatepetl	9	etl	1
4	catepetl	11	l	0
5	atepetl	3	ocatepetl	0
6	tepetl	1	opocatepetl	1
7	epetl	8	petl	0
8	petl	2	pocatepetl	1
9	etl	0	popocatepetl	2
10	tl	6	tepetl	0
11	l	10	tl	1

LCP up to 255 (longer sequences not stored).

From SA to trie



Parameter N : all sequences occurring $> N \times$ are put to trie.

Trie as stored on disk

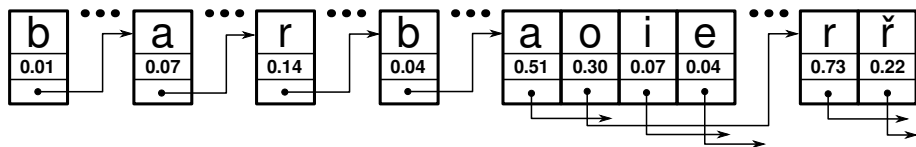


Figure : Example of a Czech model, prefix *barb*

Language model

Assign probability to any byte sequence:

```
P = ROOT          // pointer to a node in trie
index = 0          // position in input sequence
probability = 1.0
current_path = [] // path from ROOT to P
while index < length(input):
    find input[index] among children(P)
    if not found:
        shorten from left and translate the current_path to ROOT
        until it can be prolonged with input[index] byte
        // current_path may be emptied
        P = current_path[-1] or ROOT
    else:
        P = found children position
    probability *= probability of P
    current_path.append(P)
    index++
return probability
```

Example random sentences

English First there is the fact that he was listening to the sound of the shot and killed in the end a precise answer to the control of the common ancestor of the modern city of Katherine Street, and when the final result may be the structure of conservative politics; and they were standing in the corner of the room.

Czech Pornoherečka Sharon Stone se nachází v blízkosti lesa. ¶ Máme malý byt, tak jsem tu zase. ¶ Změna je život a tak by nás nevolili. ¶ Petrovi se to začalo projevovat na veřejnosti. ¶ Vojáci byli po zásluze odměněni pohledem na tvorbu mléka. ¶ Graf znázorňuje utrpení Kristovo, jež mělo splňovat následující kritéria.

Hungarian Az egyesület székhelye: 100 m-es uszonyos gyorsúszásban a következő években is részt vettek a díjat az égre nézve szójaszármazékot. ¶ Az oldal az első lépés a tengeri akvarisztikával foglalkozó szakemberek számára is ideális szállás költsége a vevőt terhelik.