

Human Language Technologies

Federico Matteoni

A.A. 2021/22

Index

- 0.1 Introduction 2
- 0.2 State of the Art 2
- 0.3 Language Modeling 3
 - 0.3.1 Evaluation and Perplexity 5
- 0.4 Representation of Words 6
 - 0.4.1 Word Embeddings 7

0.1 Introduction

Prof. Giuseppe Attardi

Prerequisites are: proficiency in Python, basic probability and statistics, calculus and linear algebra and notions of machine learning.

What will we learn Understanding of and ability to use effective modern methods for **Natural Language Processing**. From traditional methods to current advanced ones like RNN, Attentions...

Understanding the difficulties in dealing with NL and the capabilities of current technologies, with experience with **modern tools** and aiming towards the ability to build systems for some major NLP tasks: word similarities, parsing, machine translation, entity recognition, question answering, sentiment analysis, dialogue system...

Books Speech and Language Processing (Jurafsky, Martin), Deep Learning (Goodfellow, Bengio, Courville), Natural Language Processing in Python (Bird, Klein, Loper)

Exam Project (alone or team of 2-3 people) with the aim to experiment with techniques in a realistic setting using data from competitions (Kaggle, CoNLL, SemEval, Evalita...). The topic will be proposed by the team or chosen from a list of suggestions.

Experimental Approach

1. Formulate hypothesis
2. Implement technique
3. Train and test
4. Apply evaluation metric
5. If not improved:
 - Perform error analysis
 - Revise hypothesis
6. Repeat!

Motivations Language is the most distinctive feature of human intelligence, **it shapes thought**. Emulating language capabilities is a scientific challenge, a **keystone for intelligent systems** (see: Turing test)

Structured vs unstructured data The largest amount of information shared with each other is unstructured, primarily text. Information is mostly communicated by e-mails, reports, articles, conversations, media... and attempts to turn text to structured (HTML) or microformat only scratched the surface.

Problems: requires universal agreed **ontologies** and additional effort. Entity linking attempts to provide a bridge.

0.2 State of the Art

Early History During 1950s, up until AI winter.

Resurgence in the 1990s Thanks to statistical methods, novelty, to study language. Challenges arise: NIST, Netflix, DARPA Grand Challenge...

During 2010s: deep learning, neural machine translation...

Statistical Machine Learning Supervised training with **annotated** documents.

The paradigm is composed of the following:

Training set $\{x_i, y_i\}$

Representation: choose a set of features to represent data $x \mapsto \phi(x) \in R^D$

Model: choose an hypothesis function to compute $f(x) = F_{\Theta}(\phi(x))$

Evaluation: define the cost function on error with respect to examples $J(\Theta) = \sum_i (f(x_i) - y_i)^2$

Optimization: find parameters Θ that minimize $J(\Theta)$

It's a generic method, applicable to any problem.

Traditional Supervised Learning Approach Freed us from devising algorithms and rules, requiring the creation of annotated training sets and imposing the tyranny of feature engineering.

Standard approach for each new problem:

Gather as much labeled data as one can

Throw a bunch of models at it

Pick the best

Spend hours hand engineering some features or doing feature selection/dimensionality reduction

Rinse and repeat

Technological Breakthroughs Improved ML techniques but also large annotated datasets and more computing power, provided by GPUs and dedicated ML processors (like the TPU by Google).

ML exploits parallelism: stochastic gradient descent can be parallelized (asynchronous stochastic gradient descent). No need to protect shared memory access, and low (half, single) precision is enough.

Deep Learning Approach Was a big breakthrough.

Design a model architecture

Define a loss function

Run the network letting the parameters and the data representations **self-organize** as to minimize the loss

End-to-end learning: no intermediate stages nor representation

Feature representation Use a vector with each entry representing a feature of the domain element

Deep Learning represents data as vectors. Images are vectors (matrices), but words? **Word Embeddings**: transform a word into a vector of hundreds of dimensions capturing many subtle aspects of its meaning. Computed by the means of **language model**.

From a discrete to distributed representation. Words meaning are dense vectors of weights in a high dimensional space, with algebraic properties.

Background: philosophy, linguistics and statistics ML (feature vectors).

Language Model Statistical model which tells the probability that a word comes after a given word in a sentence.

Dealing with Sentences A sentence is a sequence of words: build a representation of a sequence from those of its words (compositional hypothesis). Sequence to sequence models.

Is there more structure in a sentence than a sequence of words? In many cases, tools forgets information when translating sentences into sequences of words, discarding much of the structure.

0.3 Language Modeling

Probabilistic Language Model The goal is to assign a probability to a sentence.

Machine Translation: $P(\text{high winds tonight}) > P(\text{large winds tonight})$

Spell Correction: $P(\text{about fifteen minutes from}) > P(\text{about fifteen minuets from})$

Speech Recognition: $P(\text{I saw a van}) > P(\text{eye saw a van})$

Language Identification: s from unknown language (italian or english) and $Lita$, $Leng$ language models for italian and english $\Rightarrow Lita(s) > Leng(s)$

Summarization, question answering...

We want to compute

$P(W) = P(w_1, w_2, \dots, w_n)$ the probability of a sequence

$P(w_4 | w_1, w_2, w_3, w_4)$ the probability of a word given some previous words

The model that computes that is called the **language model**

Markov Model and N-Grams Simplify the assumption: the probability of a word given all the previous is the same of the probability of that word given just few (one, two...) previous words. So $P(w_i | w_i-, \dots, w_1) = P(w_i | w_{i-1})$ (First order Markov chain).

With a **N-gram**: $P(w_n | w_1^{n-1}) \simeq P(w_n | w_{n-N+1}^{n-1})$

In general it's insufficient: language has **long distance dependencies**, but we can often get away with *N*-gram models. For example:

"The **man** next to the large oak tree near the grocery store on the corner **is** tall."

"The **men** next to the large oak tree near the grocery store on the corner **are** tall."

Or even semantic dependencies:

"The **bird** next to the large oak tree near the grocery store on the corner **flies** rapidly."

"The **man** next to the large oak tree near the grocery store on the corner **talks** rapidly."

So more complex models are needed to handle such dependencies.

Maximum likelihood estimate

$$P(w_n | w_{n-N+1}^{n-1}) = \frac{\text{count}(w_{n-N+1}^{n-1}, w_n)}{\text{count}(w_{n-N+1}^{n-1})}$$

Maximum because it's the one that maximize $P(\text{Training set} | \text{Model})$

Shannon Visualization Method Generate random sentences:

Choose a random bigram $(\langle s \rangle, w)$ according to its probability

Choose a random bigram (w, x) according to its probability

Repeat until we pick $\langle /s \rangle$

Shannon Game Generate random sentences by selecting the words according to the probabilities of the language models.

```
1 def gentext(cpd, word, length=20):
2     print(word, end = ' ')
3     for i in range(length):
4         word = cpd[word].generate()
5         print(word, end = ' ')
6     print("")
```

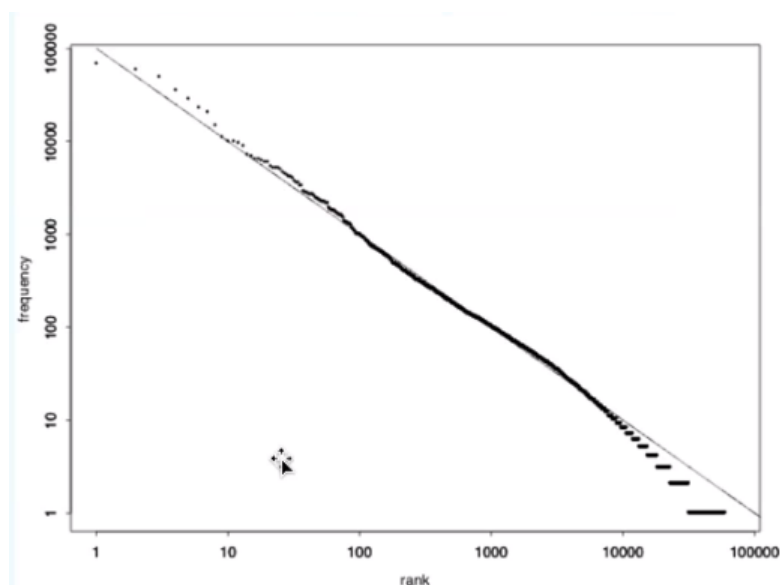
Perils of Overfitting

Zipf's Law A small number of events occur with high frequency, while a large number of events occur with small frequency. You can quickly collect statistics on the high frequency events, but you may have to wait an arbitrarily long time to get valid statistics on low frequency events.

The result is that our estimates are sparse. No counts at all for the vast bulk of things we want to estimate. Some of the zeroes in the table are really zeroes, but others are simply low frequency events you haven't seen yet: after all, anything can happen. How to address? By estimating the likelihood of unseen N-grams.

Word	Freq. (f)	Rank (r)	$f \cdot r$	Word	Freq. (f)	Rank (r)	$f \cdot r$
the	3332	1	3332	turned	51	200	10200
and	2972	2	5944	you'll	30	300	9000
a	1775	3	5235	name	21	400	8400
he	877	10	8770	comes	16	500	8000
but	410	20	8400	group	13	600	7800
be	294	30	8820	lead	11	700	7700
there	222	40	8880	friends	10	800	8000
one	172	50	8600	begin	9	900	8100
about	158	60	9480	family	8	1000	8000
more	138	70	9660	brushed	4	2000	8000
never	124	80	9920	sins	2	3000	6000
Oh	116	90	10440	Could	2	4000	8000
two	104	100	10400	Applausive	1	8000	8000

Result: f is proportional to $\frac{1}{r}$, meaning that there is a constant $k \mid f \cdot r = k$



0.3.1 Evaluation and Perplexity

Evaluation Train parameters of our model on a training set. For the evaluation, we apply the model on new data and look at the model's performance: **test set**. We need a metric which tells us how well the model is performing: one of such is **perplexity**.

Extrinsic Evaluation Evaluating N-gram models. Put model A in a task (language identification, speech recognizer, machine translation system...)

Language Identification Task Build a model for each language and compute probability that text is of such language.

$$lang = \arg \max_l P_l(text)$$

Difficulty of Extrinsic Evaluation Extrinsic evaluation is time-consuming, can take days to run an experiment. Recent benchmarks like GLUE have become popular due to the effectiveness.

Intrinsic evaluation uses an approximation called **perplexity**: it's a poor approximation unless the test data looks just like the training data.

Perplexity The intuition is the notion of surprise: how surprised is the language model when it sees the test set? Where surprise is a measure of "I didn't see that coming". The more surprised is, the lower the probability it assigned to the test set, and the higher the probability the less surprised it is.

So perplexity measures how well a model "fits" the test data. It uses the probability that the model assigns to the test corpus and normalizes for the number of words in the test corpus, taking the inverse.

$$PP(w) = \sqrt[N]{\frac{1}{P(w_1 \dots w_N)}}$$

Measures the weighted average branching factor in predicting the next word (lower is better).
In the chain rule

$$PP(w) = \sqrt[N]{\prod_{i=1}^N \frac{1}{P(w_i | w_1 \dots w_{i-1})}}$$

And for bigrams

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

Lower perplexity \Rightarrow better model

0.4 Representation of Words

Word Meaning Definition of **meaning**:

Idea represented by a word, phrase...

Idea that a person wants to express by using words, signs...

Idea expressed in a work of writing, art...

Signifier (symbol) \Leftrightarrow **Signified** (idea or concept) = **denotation**

Linguistic Solution Dictionary or lexical resources provide word definitions as well as synonyms, hypernyms, antonyms... example: **Wordnet**.

Problems with lexical resources

Imperfect, sometimes there are errors

Missing new meanings of words, hard to keep up-to-date

Subjective

Require human labor to create and adapt

Hard to compute accurate word similarity

Vector Space Model The **VSM** is a representation for text used in information retrieval. A document is represented by a vector in a n -dimensional space

$$v(d_1) = [t_1, \dots, t_n]$$

Each dimension corresponds to a separate term.
The VSM considers words as discrete symbols.

One Hot Representation A word occurrence is represented by one-hot vectors, with the vector dimension being the number of words in a vocabulary and 1 in correspondence of the position of the represented word. Useful for representing a document by the sum of the word occurrence vector, and document similarity given by **cosine distance** (angle between vectors).

tf*idf Measure Term frequency: frequency of the word.
Inverse document frequency:

Classical VSM Vector is all zeroes with idf_t instead of 1. The vector of a document is the sum of the vectors of its terms occurrence vectors.
Doesn't capture similarity between terms.

Problems with Discrete Symbols All vectors are orthogonal, with no notion of similarity. Search engines try to address the issue using WordNet synonyms but it's better to encode the similarity in the vectors themselves.

Intuition Model the meaning of a word by "embedding" it in a vector space. The meaning of a word is a vectors of numbers (vector models are also called **embeddings**), so a **dense vector space**. Contrast: word meaning is represented in many computational

Word Vectors/Embeddings Build dense vector for each word chosen so that it is similar to vectors of words that appear in similar contexts.
Four kinds:

Sparse Vector Representations

1. Mutual-information weighted word co-occurrence matrices

Dense Vector Representations

- 2.
- 3.
- 4.

Distributional Hypothesis A word's meaning is given by the words that frequently appear close-by. Old idea but successful just with modern statistical NLP.

When a word w appears in a text, its context is the set of words that appear nearby (within a fixed-size window). Use of the many contexts of w to build a representation of w .

Word Context Matrix Is a $|V| \times |V|$ matrix X that counts the frequencies of co-occurrence of words in a collection of contexts (i.e. text spans of a given length).

Co-Occurrence Matrix Words \equiv Context words. Rows of X capture similarity yet X is still high dimensional and sparse. One row per word with counts of occurrences of any other word.
We can compute the distance vectors between words, but neighboring words are not semantically related.

0.4.1 Word Embeddings

Dense Representations Project word vectors $v(t)$ into a low dimensional space R^k with $k \ll |V|$ of continuous space word representations (a.k.a. **embeddings**)

$$Embed : R^{|V|} \rightarrow R^k$$

$$Embed(v(t)) = e(t)$$

Desired properties: remaining dimensions

Collobert Build embeddings and estimate whether the word is in the proper context using a neural network. Positive examples from text, and negative examples made replacing center word with random one. The loss for the training is

$$Loss(\Theta) = \sum_{x \in X} \sum_{w \in W} \max(0, 1 - f_{\Theta}(x) + f_{\Theta}(x^{(w)}))$$

with $x^{(w)}$ obtained by replacing the central word of x with a random word.

Word2Vec Framework for learning words, much faster to train. Idea:

Collect a large corpus of text

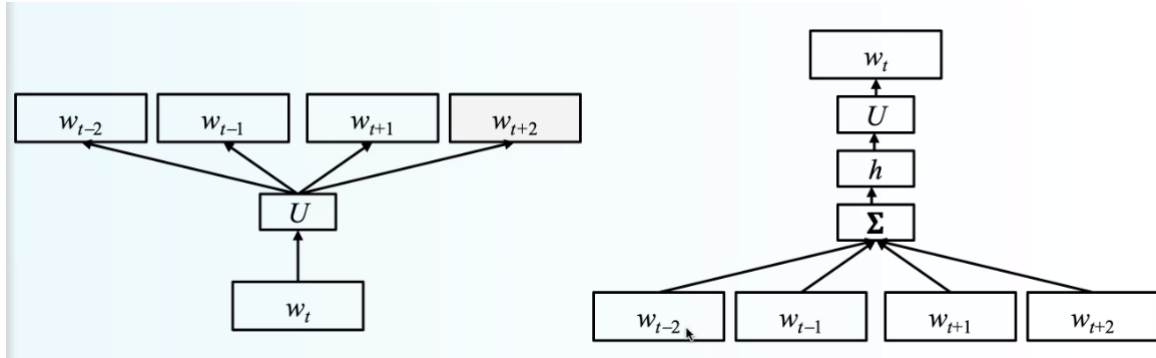
Every word in a fixed vocabulary is represented by a vector

Go through each position t in the text, which has a center word c and context ("outside") words o

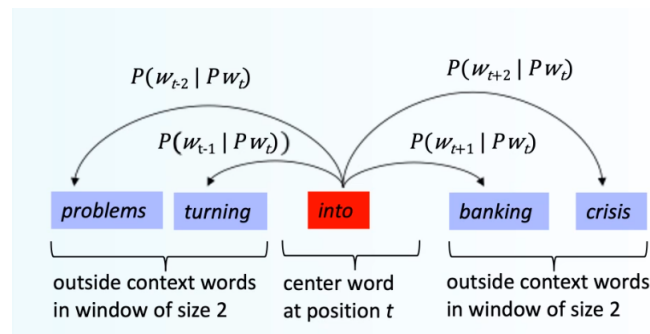
Use the similarity of the word vectors for c and o to calculate the probability of o given c

Keep adjusting word vectors to maximize this probability

Two kinds:



Skip-gram: predict context words within window of size m given the center word w_t



CBoW: predict center word w_t given context words within window of size m

Embeddings are a by-product of the word prediction task. Even though it's a prediction task, the network can be trained on any text (no need for human-labeled data!).

Usual context size is 5 words before and after. Features can be multi-word expressions. Longer windows can capture more semantics and less syntax.

A typical size for h is 200-300.

Skip-Gram

Inputs are one-hot representation of the word

$w \in R^{|Vocabulary|}$ are high-dimensional

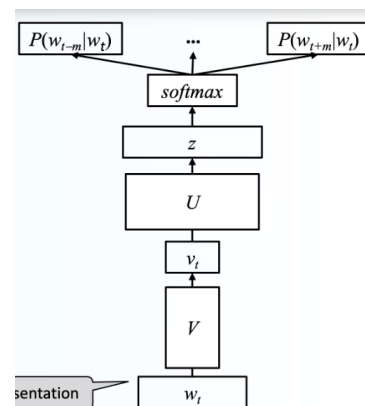
$v \in R^d$ is low dimensional: size of the embedding space d

$V \in R^{|Voc| \times d}$ input word matrix

row t of V is the input vector, representation for **center** word w_t

$U \in R^{d \times |Voc|}$ output matrix

column o of U is the output vector, representation for **context** word w_o



$v_t = w_t V$ embedding of word w_t
 $z = v_t U$ z_i is the similarity of w_t with w_i
 Softmax converts z to a probability distribution p_i

$$p_i = \frac{e^{z_i}}{\sum_{j \in V} e^{z_j}}$$

Procedure

Lookup the embedded word vector for the center word $v_c = V[c] \in R^n$

Generate score vector $z = v_c U$

Turn the score vector into probabilities $\hat{y} = \text{softmax}(z)$

\hat{y}_{c-m}, \dots Those are the estimates

So the training iterates through the whole corpus predicting surrounding words given the center word.

Objective Function For each position $t = 1, \dots, T$ predict context words within a windows of fixed size m given each center word w_t

$$\text{Likelihood} = L(\Theta) = \prod_{t=1}^T \prod_{-m \leq j \leq m, j \neq 0} P(w_{t+j} | w_t, \Theta)$$

While the objective function $J(\Theta)$ is the average negative log likelihood

$$J(\Theta) = -\frac{1}{T} \log L(\Theta) = -\frac{1}{T} \sum_{t=1}^T \sum_{-m \leq j \leq m, j \neq 0} \log(P(w_{t+j} | w_t, \Theta)) =$$

To compute $P(o | c, \Theta)$ we use two vectors per word w : v_w when w center word and u_w when w context word. For a center word c and a context word o :

$$P(o | c) = \frac{e^{u_o v_c}}{\sum_{w \in V} e^{u_w v_c}}$$

With the dot product $u \cdot v = \sum_i u_i v_i$: larger product \Rightarrow larger probability. Normalize over entire vocabulary to give probability distribution.

$J(\Theta)$ is a function of all windows in the corpus, potentially billions: too expensive to compute. The solution is the stochastic gradient descent, sampling windows randomly and update after each one.

Softmax

Soft because still assign some probability to smaller x_i

Max because amplifies the probability of largest x_i

Can we really capture the concept represented by a word? Philosophical debate.

Negative Sampling $\log \sum_{j \in F} e^{u_j}$ has lots of terms, costly to compute. The solution is to compute it only on a small sample of negative samples, i.e. $\log \sum_{j \in E} e^{u_j}$ where words E are a few (e.g. 5)

CBoW Continuous Bag of Words

Mirror of the skip-gram, where context words are used to predict target words.

h is computed from the average of the embeddings of the input context, z_i is the similarity of h with the words embedding of w_i from U .

Which Embeddings V and U both define embeddings, which to use? Usually just V . Sometimes average pairs of vectors from V and U into a single one or append one embedding vector after the other, doubling the length.

GloVe Global Vectors for Word Representation. Insight: the ratio of conditional probabilities may capture meaning.

$$J = \sum_{i,j=1}^V f(X_{ij}) \dots$$

fastText Similar to CBoW, word embeddings averaged to obtain good sentence representation. Pretrained models.