


# CRF, CTC, Word2Vec

Milan Straka

 April 3, 2023



EUROPEAN UNION  
European Structural and Investment Fund  
Operational Programme Research,  
Development and Education

Charles University in Prague  
Faculty of Mathematics and Physics  
Institute of Formal and Applied Linguistics



unless otherwise stated



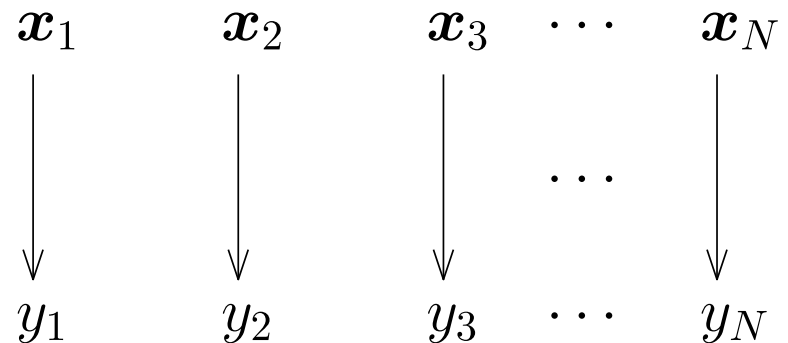
## Structured Prediction



# Structured Prediction

Consider generating a sequence of  $y_1, \dots, y_N \in Y^N$  given input  $\mathbf{x}_1, \dots, \mathbf{x}_N$ .

Predicting each sequence element independently models the distribution  $P(y_i | \mathbf{X})$ .



However, there may be dependencies among the  $y_i$  themselves, which is difficult to capture by independent element classification.

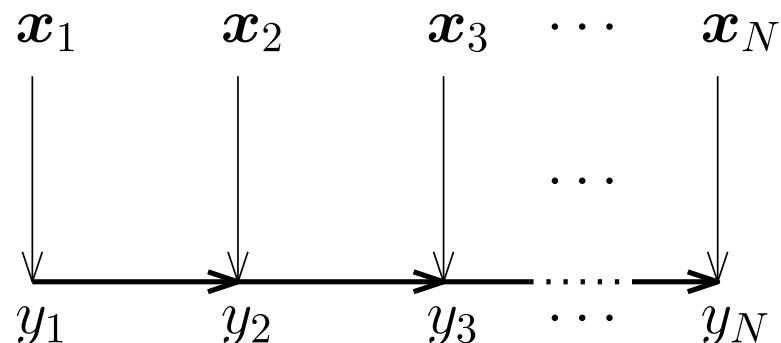


# Maximum Entropy Markov Models

We might model the dependencies by assuming that the output sequence is a Markov chain, and model it as

$$P(y_i | \mathbf{X}, y_{i-1}).$$

Each label would be predicted by a softmax from the hidden state and the *previous label*.



The decoding can be then performed by a dynamic programming algorithm.



However, MEMMs suffer from a so-called **label bias** problem. Because the probability is factorized, each  $P(y_i | \mathbf{X}, y_{i-1})$  is a distribution and **must sum to one**.

Imagine there was a label error during prediction. In the next step, the model might “realize” that the previous label has very low probability of being followed by any label – however, it cannot express this by setting the probability of all following labels to a low value, it has to “conserve the mass”.



Let  $G = (V, E)$  be a graph such that  $\mathbf{y}$  is indexed by vertices of  $G$ . Then  $(\mathbf{X}, \mathbf{y})$  is a **conditional random field**, if the random variables  $\mathbf{y}$  conditioned on  $\mathbf{X}$  obey the Markov property with respect to the graph, i.e.,

$$P(y_i | \mathbf{X}, \{y_j \mid \forall j \neq i\}) = P(y_i | \mathbf{X}, \{y_j \mid \forall j : (i, j) \in E\}).$$

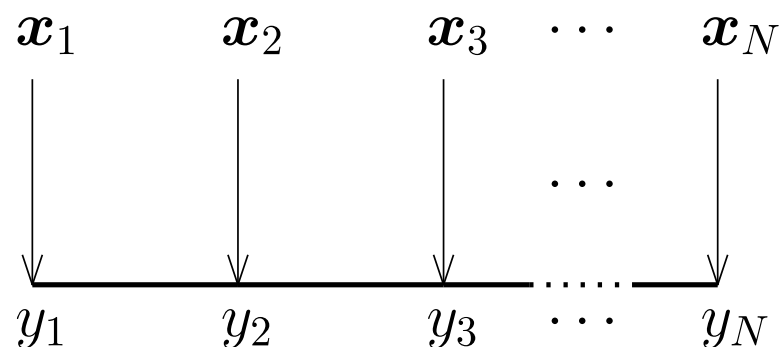
In plain speech, if you consider probabilities conditional on  $\mathbf{X}$ , all the dependencies among  $\mathbf{y}$  are captured by the edges of the graph  $G$ .

By the *fundamental theorem of random fields (the Hammersley–Clifford theorem)*, the density of a conditional random field can be factorized over the *cliques* (complete subgraphs) of the graph  $G$ :

$$P(\mathbf{y} | \mathbf{X}) = \prod_{\text{clique } C \text{ of } G} P(\mathbf{y}_C | \mathbf{X}).$$



Most often, we assume that dependencies of  $\mathbf{y}$ , conditioned on  $\mathbf{X}$ , form a chain.



Then the cliques are *nodes* and *edges*, and we can factorize the probability as:

$$P(\mathbf{y}|\mathbf{X}) \propto \exp \left( \sum_{i=1}^N \log P(y_i|\mathbf{X}) + \sum_{i=2}^N \log P(y_i, y_{i-1}|\mathbf{X}) \right).$$



Linear-chain Conditional Random Field, usually abbreviated only to CRF, acts as an output layer. It can be considered an extension of softmax – instead of a sequence of independent softmaxes, it is a sentence-level softmax, with additional weights for neighboring sequence elements.

We start by defining a score of a label sequence  $\mathbf{y}$  as

$$s(\mathbf{X}, \mathbf{y}; \boldsymbol{\theta}, \mathbf{A}) = f(y_1 | \mathbf{X}; \boldsymbol{\theta}) + \sum_{i=2}^N (\mathbf{A}_{y_{i-1}, y_i} + f(y_i | \mathbf{X}; \boldsymbol{\theta})),$$

and define the probability of a label sequence  $\mathbf{y}$  using softmax:

$$p(\mathbf{y} | \mathbf{X}) = \text{softmax}_{\mathbf{z} \in Y^N} (s(\mathbf{X}, \mathbf{z}))_{\mathbf{y}}.$$

For cross-entropy (and also to avoid underflow), we need the logarithm of the probability:

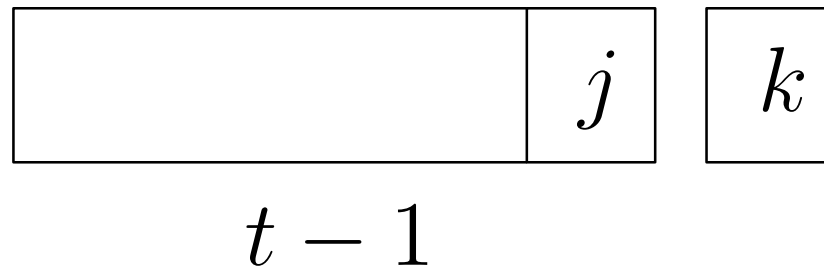
$$\log p(\mathbf{y} | \mathbf{X}) = s(\mathbf{X}, \mathbf{y}) - \text{logsumexp}_{\mathbf{z} \in Y^N} (s(\mathbf{X}, \mathbf{z})), \text{ where}$$
$$\text{logsumexp}_x (f(x)) = \log \left( \sum_x e^{f(x)} \right).$$



## Computation

We can compute  $p(\mathbf{y}|\mathbf{X})$  efficiently using dynamic programming. We denote  $\alpha_t(k)$  the logarithmic probability of all  $t$ -element sequences with the last label  $y$  being  $k$ .

The core idea is the following:



$$\alpha_t(k) = f(y_t = k|\mathbf{X}; \boldsymbol{\theta}) + \text{logsumexp}_{j \in Y} (\alpha_{t-1}(j) + \mathbf{A}_{j,k}).$$

For efficient implementation, we use the fact that

$$\begin{aligned} \log(a + b) &= \log a + \log(1 + e^{\log b - \log a}), \text{ so} \\ \text{logsumexp}_x (f(x)) &= \max_x (f(x)) + \log(\sum_x e^{f(x) - \max_x (f(x))}). \end{aligned}$$



**Inputs:** Network computing  $f(y_t = k | \mathbf{X}; \boldsymbol{\theta})$ , which is a logit (unnormalized log-probability) of output sequence label being  $k$  at time  $t$ .

**Inputs:** Transition matrix  $\mathbf{A} \in \mathbb{R}^{Y \times Y}$ .

**Inputs:** Input sequence  $\mathbf{X}$  of length  $N$ , gold labeling  $\mathbf{g} \in Y^N$ .

**Outputs:** Value of  $\log p(\mathbf{g} | \mathbf{X})$ .

**Time Complexity:**  $\mathcal{O}(N \cdot Y^2)$ .

- For  $t = 1, \dots, N$ :
  - For  $k = 1, \dots, Y$ :
    - $\alpha_t(k) \leftarrow f(y_t = k | \mathbf{X}; \boldsymbol{\theta})$
    - If  $t > 1$ :
      - $\alpha_t(k) \leftarrow \alpha_t(k) + \text{logsumexp}_{j=1}^Y (\alpha_{t-1}(j) + \mathbf{A}_{j,k})$
- Return  $\sum_{t=1}^N f(y_t = g_t | \mathbf{X}; \boldsymbol{\theta}) + \sum_{t=2}^N \mathbf{A}_{g_{t-1}, g_t} - \text{logsumexp}_{k=1}^Y (\alpha_N(k))$



# Conditional Random Fields (CRF)

## Decoding

We can perform decoding optimally, by using the same algorithm, only replacing `logsumexp` with `max`, and tracking where the maximum was attained.

## Applications

The CRF output layer is useful for **span labeling** tasks, like

- named entity recognition,
- dialog slot filling.

It can be also useful for image segmentation.

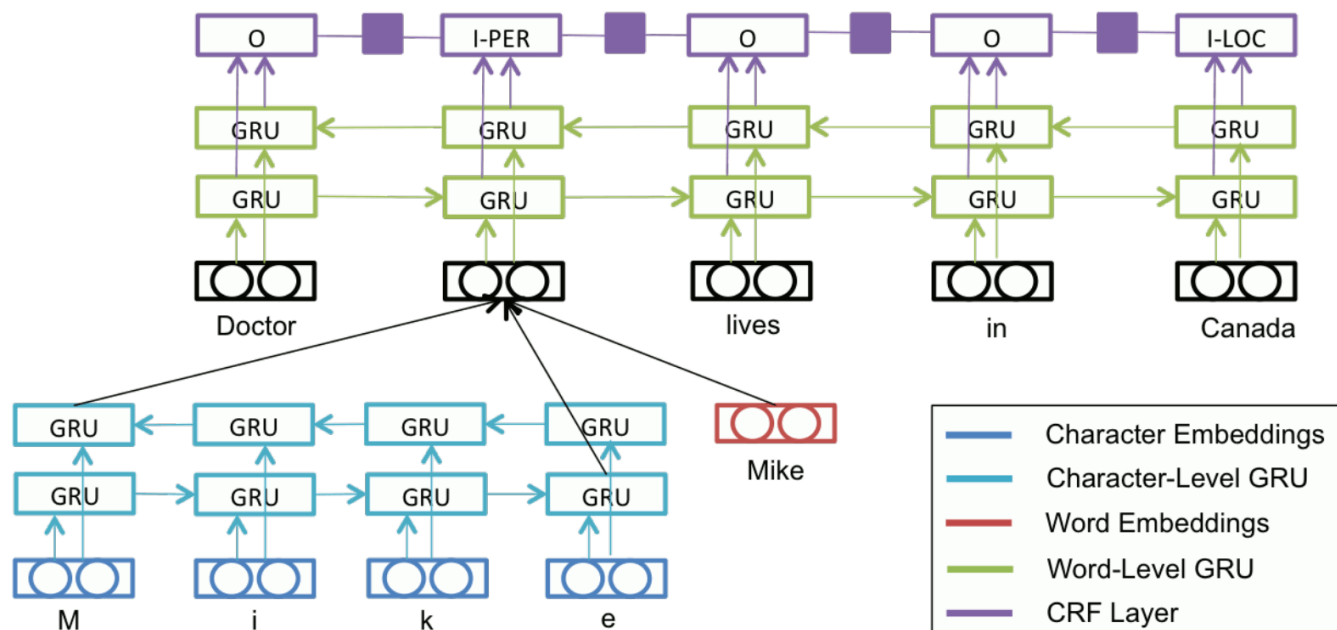


Figure 1 of "Multi-Task Cross-Lingual Sequence Tagging from Scratch", <https://arxiv.org/abs/1603.06270>



# Connectionist Temporal Classification

Let us again consider generating a sequence of  $y_1, \dots, y_M$  given input  $x_1, \dots, x_N$ , but this time  $M \leq N$ , and there is no explicit alignment of  $x$  and  $y$  in the gold data.

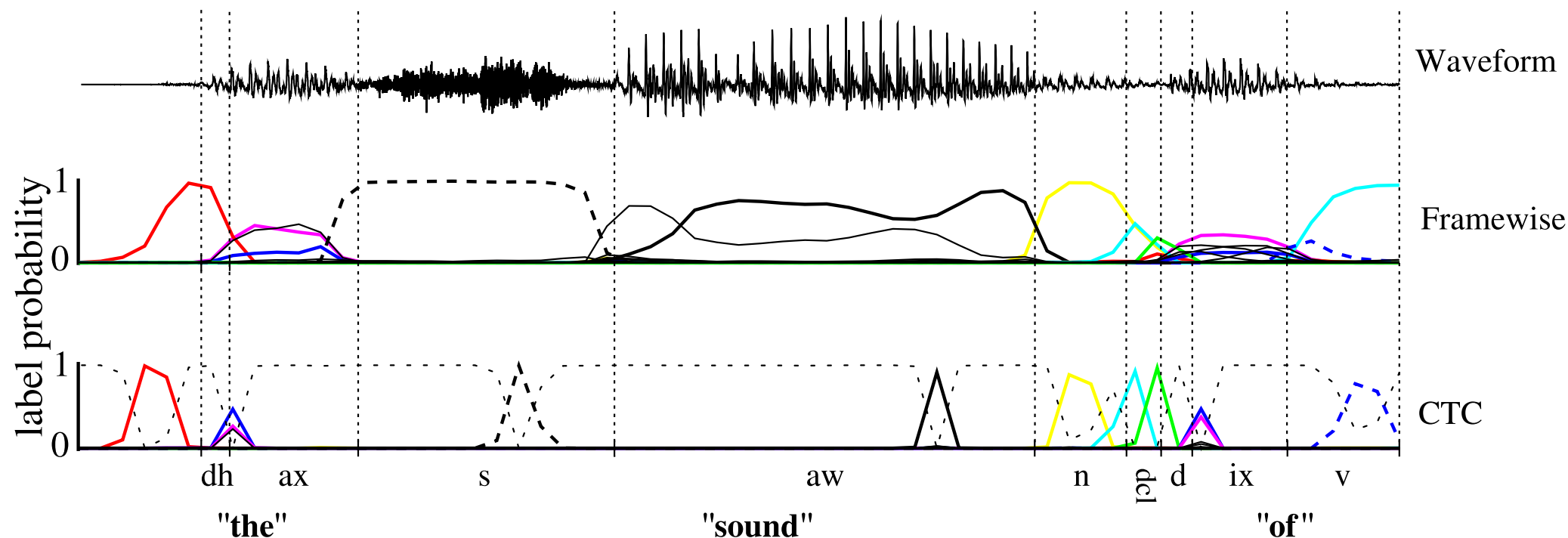


Figure 7.1 of "Supervised Sequence Labelling with Recurrent Neural Networks" dissertation by Alex Graves



# Connectionist Temporal Classification

We enlarge the set of the output labels by a – (**blank**), and perform a classification for every input element to produce an **extended labeling** (in contrast to the original **regular labeling**). We then post-process it by the following rules (denoted as  $\mathcal{B}$ ):

1. We collapse multiple neighboring occurrences of the same symbol into one.
2. We remove the blank –.

Because the explicit alignment of inputs and labels is not known, we consider *all possible* alignments.

Denoting the probability of label  $l$  at time  $t$  as  $p_l^t$ , we define

$$\alpha^t(s) \stackrel{\text{def}}{=} \sum_{\substack{\text{extended} \\ \text{labelings } \boldsymbol{\pi}: \\ \mathcal{B}(\boldsymbol{\pi}_{1:t}) = \mathbf{y}_{1:s}}} \prod_{i=1}^t p_{\pi_i}^i.$$



In CRF, we normalize the whole sentences, therefore we need to compute unnormalized probabilities for all the (exponentially many) sentences. Decoding can be performed optimally.

In CTC, we normalize per each label. However, because we do not have explicit alignment, we compute probability of a labeling by summing probabilities of (generally exponentially many) extended labelings.



## Computation

When aligning an extended labeling to a regular one, we need to consider whether the extended labeling ends by a *blank* or not. We therefore define

$$\alpha_{-}^t(s) \stackrel{\text{def}}{=} \sum_{\substack{\text{extended} \\ \text{labelings } \boldsymbol{\pi}: \\ \mathcal{B}(\boldsymbol{\pi}_{1:t}) = \mathbf{y}_{1:s}, \pi_t = -}} \prod_{i=1}^t p_{\pi_i}^i$$

$$\alpha_{*}^t(s) \stackrel{\text{def}}{=} \sum_{\substack{\text{extended} \\ \text{labelings } \boldsymbol{\pi}: \\ \mathcal{B}(\boldsymbol{\pi}_{1:t}) = \mathbf{y}_{1:s}, \pi_t \neq -}} \prod_{i=1}^t p_{\pi_i}^i$$

and compute  $\alpha^t(s)$  as  $\alpha_{-}^t(s) + \alpha_{*}^t(s)$ .



## Computation – Initialization

We initialize  $\alpha^1$  as follows:

- $\alpha_{-}^1(0) \leftarrow p_{-}^1$
- $\alpha_{*}^1(1) \leftarrow p_{y_1}^1$
- all other  $\alpha^1$  to zeros

## Computation – Induction Step

We then proceed recurrently according to:

- $\alpha_{-}^t(s) \leftarrow p_{-}^t(\alpha_{*}^{t-1}(s) + \alpha_{-}^{t-1}(s))$
- $\alpha_{*}^t(s) \leftarrow \begin{cases} p_{y_s}^t(\alpha_{*}^{t-1}(s) + \alpha_{-}^{t-1}(s-1) + \alpha_{*}^{t-1}(s-1)), & \text{if } y_s \neq y_{s-1} \\ p_{y_s}^t(\alpha_{*}^{t-1}(s) + \alpha_{-}^{t-1}(s-1) + \alpha_{*}^{t-1}(s-1)), & \text{if } y_s = y_{s-1} \end{cases}$

We can write the update as  $p_{y_s}^t(\alpha_{*}^{t-1}(s) + \alpha_{-}^{t-1}(s-1) + [y_s \neq y_{s-1}] \cdot \alpha_{*}^{t-1}(s-1))$ .

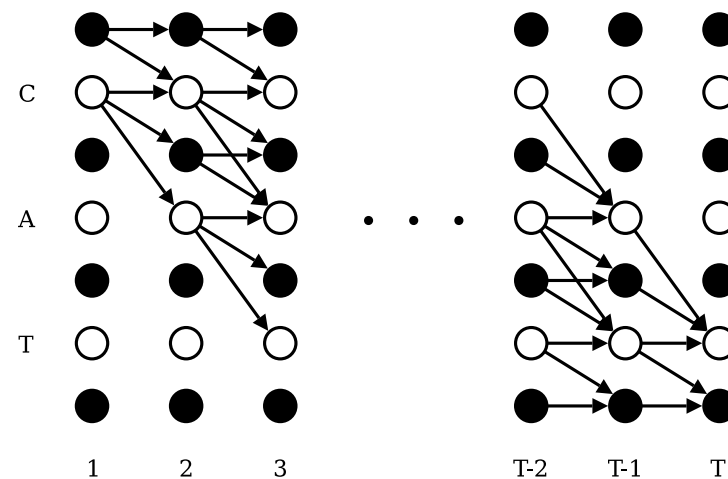


Figure 7.3 of "Supervised Sequence Labelling with Recurrent Neural Networks" dissertation by Alex Graves



# CTC Decoding

Unlike CRF, nobody knows how to perform decoding optimally in polynomial time.

The key observation is that while an optimal extended labeling can be extended into an optimal labeling of a greater length, the same does not apply to a regular labeling. The problem is that regular labeling corresponds to many extended labelings, which are modified each in a different way during an extension of the regular labeling.

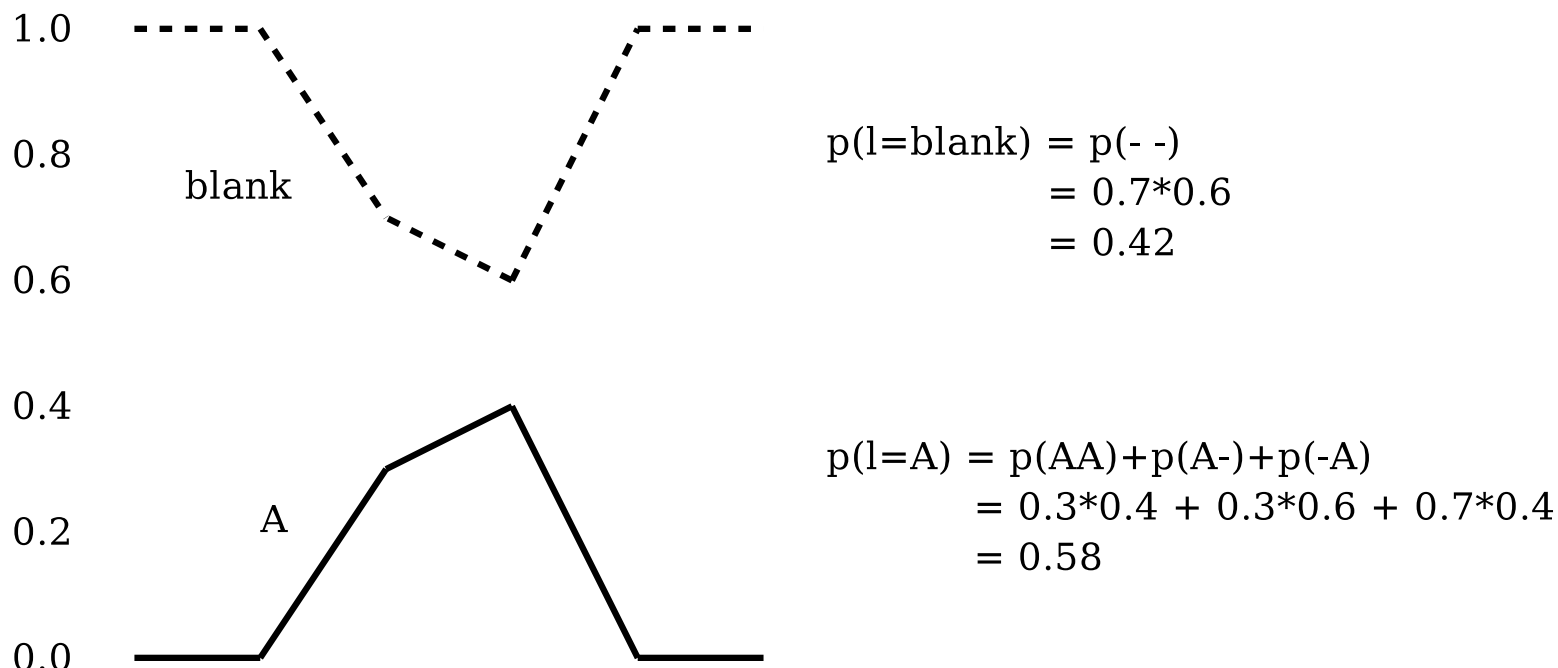


Figure 7.5 of "Supervised Sequence Labelling with Recurrent Neural Networks" dissertation by Alex Graves



## Beam Search

To perform a beam search, we keep  $k$  best **regular** (non-extended) labelings. Specifically, for each regular labeling  $\mathbf{y}$  we keep both  $\alpha_{-}^t(\mathbf{y})$  and  $\alpha_{*}^t(\mathbf{y})$ , which are probabilities of all (modulo beam search) extended labelings of length  $t$  which produce the regular labeling  $\mathbf{y}$ ; we therefore keep  $k$  regular labelings with the highest  $\alpha_{-}^t(\mathbf{y}) + \alpha_{*}^t(\mathbf{y})$ .

To compute the best regular labelings for a longer prefix of extended labelings, for each regular labeling in the beam we consider the following cases:

- adding a *blank* symbol, i.e., contributing to  $\alpha_{-}^{t+1}(\mathbf{y})$  both from  $\alpha_{-}^t(\mathbf{y})$  and  $\alpha_{*}^t(\mathbf{y})$ ;
- adding a non-blank symbol, i.e., contributing to  $\alpha_{*}^{t+1}(\bullet)$  from  $\alpha_{-}^t(\mathbf{y})$  and contributing to a possibly different  $\alpha_{*}^{t+1}(\bullet)$  from  $\alpha_{*}^t(\mathbf{y})$ .

Finally, we merge the resulting candidates according to their regular labeling, and keep only the  $k$  best.



The embeddings can be trained for each task separately.

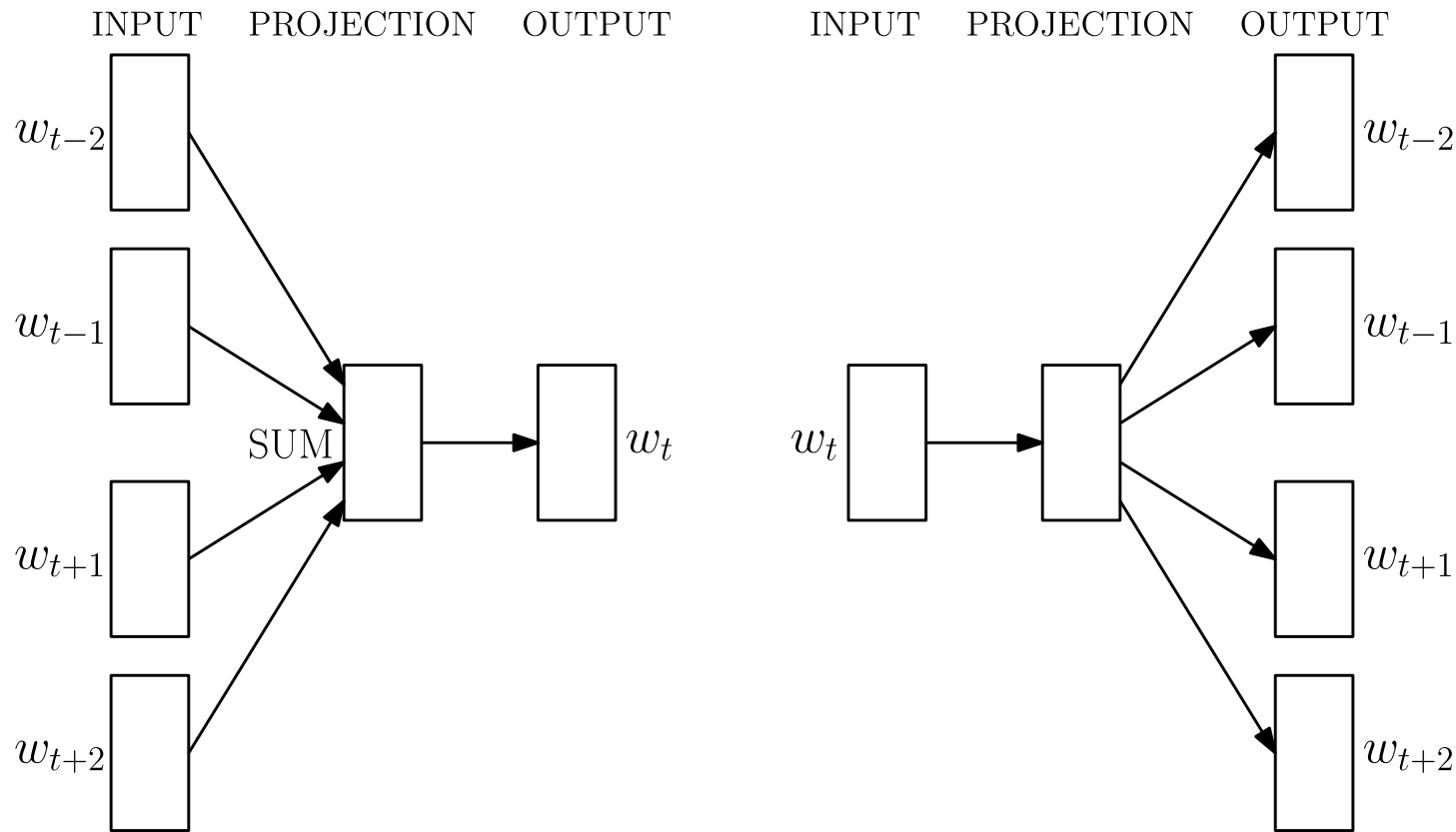
However, a method of precomputing word embeddings have been proposed, based on *distributional hypothesis*:

**Words that are used in the same contexts tend to have similar meanings.**

The distributional hypothesis is usually attributed to Firth (1957):

*You shall know a word by a company it keeps.*





CBOW (Continuous Bag Of Words)

Skip-gram

Mikolov et al. (2013) proposed two very simple architectures for precomputing word embeddings, together with a C multi-threaded implementation `word2vec`.

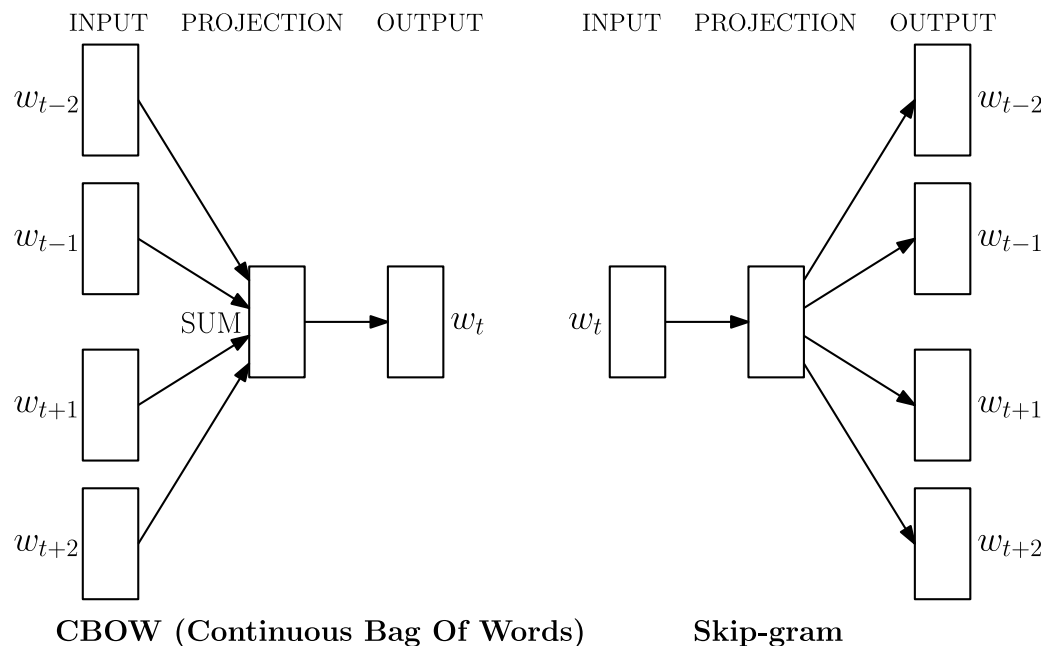


Table 8: *Examples of the word pair relationships, using the best word vectors from Table 4 (Skip-gram model trained on 783M words with 300 dimensionality).*

Relationship	Example 1	Example 2	Example 3
France - Paris	Italy: Rome	Japan: Tokyo	Florida: Tallahassee
big - bigger	small: larger	cold: colder	quick: quicker
Miami - Florida	Baltimore: Maryland	Dallas: Texas	Kona: Hawaii
Einstein - scientist	Messi: midfielder	Mozart: violinist	Picasso: painter
Sarkozy - France	Berlusconi: Italy	Merkel: Germany	Koizumi: Japan
copper - Cu	zinc: Zn	gold: Au	uranium: plutonium
Berlusconi - Silvio	Sarkozy: Nicolas	Putin: Medvedev	Obama: Barack
Microsoft - Windows	Google: Android	IBM: Linux	Apple: iPhone
Microsoft - Ballmer	Google: Yahoo	IBM: McNealy	Apple: Jobs
Japan - sushi	Germany: bratwurst	France: tapas	USA: pizza

Table 8 of "Efficient Estimation of Word Representations in Vector Space", <https://arxiv.org/abs/1301.3781>





Considering input word  $w_i$  and output  $w_o$ , the Skip-gram model defines

$$p(w_o|w_i) \stackrel{\text{def}}{=} \frac{e^{\mathbf{V}_{w_i}^\top \mathbf{W}_{w_o}}}{\sum_w e^{\mathbf{V}_{w_i}^\top \mathbf{W}_w}}.$$

After training, the final embeddings are the rows of the  $\mathbf{V}$  matrix.



Instead of a large softmax, we construct a binary tree over the words, with a sigmoid classifier for each node.

If word  $w$  corresponds to a path  $n_1, n_2, \dots, n_L$ , we define

$$p_{\text{HS}}(w|w_i) \stackrel{\text{def}}{=} \prod_{j=1}^{L-1} \sigma([+1 \text{ if } n_{j+1} \text{ is right child else } -1] \cdot \mathbf{V}_{w_i}^\top \mathbf{W}_{n_j}).$$



Instead of a large softmax, we could train individual sigmoids for all words.

We could also only sample several *negative examples*. This gives rise to the following *negative sampling* objective (instead of just summing all the sigmoidal losses):

$$l_{\text{NEG}}(w_o, w_i) \stackrel{\text{def}}{=} -\log \sigma(\mathbf{V}_{w_i}^\top \mathbf{W}_{w_o}) - \sum_{j=1}^k \mathbb{E}_{w_j \sim P(w)} \log (1 - \sigma(\mathbf{V}_{w_i}^\top \mathbf{W}_{w_j})).$$

The usual value of negative samples  $k$  is 5, but it can be even 2 for extremely large corpora.

Each expectation in the loss is estimated using a single sample.

For  $P(w)$ , both uniform and unigram distribution  $U(w)$  work, but

$$U(w)^{3/4}$$

outperforms them significantly (this fact has been reported in several papers by different authors).



<i>increased</i>	<i>John</i>	<i>Noahshire</i>	<i>phding</i>
reduced	Richard	Nottinghamshire	mixing
improved	George	Bucharest	modelling
expected	James	Saxony	styling
decreased	Robert	Johannesburg	blaming
targeted	Edward	Gloucestershire	christening

Table 2: Most-similar in-vocabular words under the C2W model; the two query words on the left are in the training vocabulary, those on the right are nonce (invented) words.

Table 2 of "Finding Function in Form: Compositional Character Models for Open Vocabulary Word Representation", <https://arxiv.org/abs/1508.02096>



	In Vocabulary					Out-of-Vocabulary		
	<i>while</i>	<i>his</i>	<i>you</i>	<i>richard</i>	<i>trading</i>	<i>computer-aided</i>	<i>misinformed</i>	<i>loooooook</i>
LSTM-Word	<i>although</i>	<i>your</i>	<i>conservatives</i>	<i>jonathan</i>	<i>advertised</i>	–	–	–
	<i>letting</i>	<i>her</i>	<i>we</i>	<i>robert</i>	<i>advertising</i>	–	–	–
	<i>though</i>	<i>my</i>	<i>guys</i>	<i>neil</i>	<i>turnover</i>	–	–	–
	<i>minute</i>	<i>their</i>	<i>i</i>	<i>nancy</i>	<i>turnover</i>	–	–	–
LSTM-Char (before highway)	<i>chile</i>	<i>this</i>	<i>your</i>	<i>hard</i>	<i>heading</i>	<i>computer-guided</i>	<i>informed</i>	<i>look</i>
	<i>whole</i>	<i>hhs</i>	<i>young</i>	<i>rich</i>	<i>training</i>	<i>computerized</i>	<i>performed</i>	<i>cook</i>
	<i>meanwhile</i>	<i>is</i>	<i>four</i>	<i>richer</i>	<i>reading</i>	<i>disk-drive</i>	<i>transformed</i>	<i>looks</i>
	<i>white</i>	<i>has</i>	<i>youth</i>	<i>richter</i>	<i>leading</i>	<i>computer</i>	<i>inform</i>	<i>shook</i>
LSTM-Char (after highway)	<i>meanwhile</i>	<i>hhs</i>	<i>we</i>	<i>eduard</i>	<i>trade</i>	<i>computer-guided</i>	<i>informed</i>	<i>look</i>
	<i>whole</i>	<i>this</i>	<i>your</i>	<i>gerard</i>	<i>training</i>	<i>computer-driven</i>	<i>performed</i>	<i>looks</i>
	<i>though</i>	<i>their</i>	<i>doug</i>	<i>edward</i>	<i>traded</i>	<i>computerized</i>	<i>outperformed</i>	<i>looked</i>
	<i>nevertheless</i>	<i>your</i>	<i>i</i>	<i>carl</i>	<i>trader</i>	<i>computer</i>	<i>transformed</i>	<i>looking</i>

**Table 6:** Nearest neighbor words (based on cosine similarity) of word representations from the large word-level and character-level (before and after highway layers) models trained on the PTB. Last three words are OOV words, and therefore they do not have representations in the word-level model.

Table 6 of "Character-Aware Neural Language Models", <https://arxiv.org/abs/1508.06615>



Another simple idea appeared simultaneously in three nearly simultaneous publications as [Charagram](#), [Subword Information](#) or [SubGram](#).

A word embedding is a sum of the word embedding plus embeddings of its character  $n$ -grams. Such embedding can be pretrained using same algorithms as `word2vec`.

The implementation can be

- dictionary based: only some number of frequent character  $n$ -grams is kept;
- hash-based: character  $n$ -grams are hashed into  $K$  buckets (usually  $K \sim 10^6$  is used).

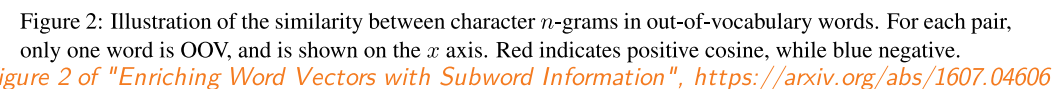


query	tiling	tech-rich	english-born	micromanaging	eateries	dendritic
sisg	tile flooring	tech-dominated tech-heavy	british-born polish-born	micromanage micromanaged	restaurants eaterie	dendrite dendrites
sg	bookcases built-ins	technology-heavy .ixic	most-capped ex-scotland	defang internalise	restaurants delis	epithelial p53

Table 7: Nearest neighbors of rare words using our representations and skipgram. These hand picked examples are for illustration.

Table 7 of "Enriching Word Vectors with Subword Information", <https://arxiv.org/abs/1607.04606>







The word2vec enriched with subword embeddings is implemented in publicly available fastText library <https://fasttext.cc/>.

Pre-trained embeddings for 157 languages (including Czech) trained on Wikipedia and CommonCrawl are also available at <https://fasttext.cc/docs/en/crawl-vectors.html>.