

Sequence labelling with global search

Marco Kuhlmann

Department of Computer and Information Science

High-level approaches

- **Local search**

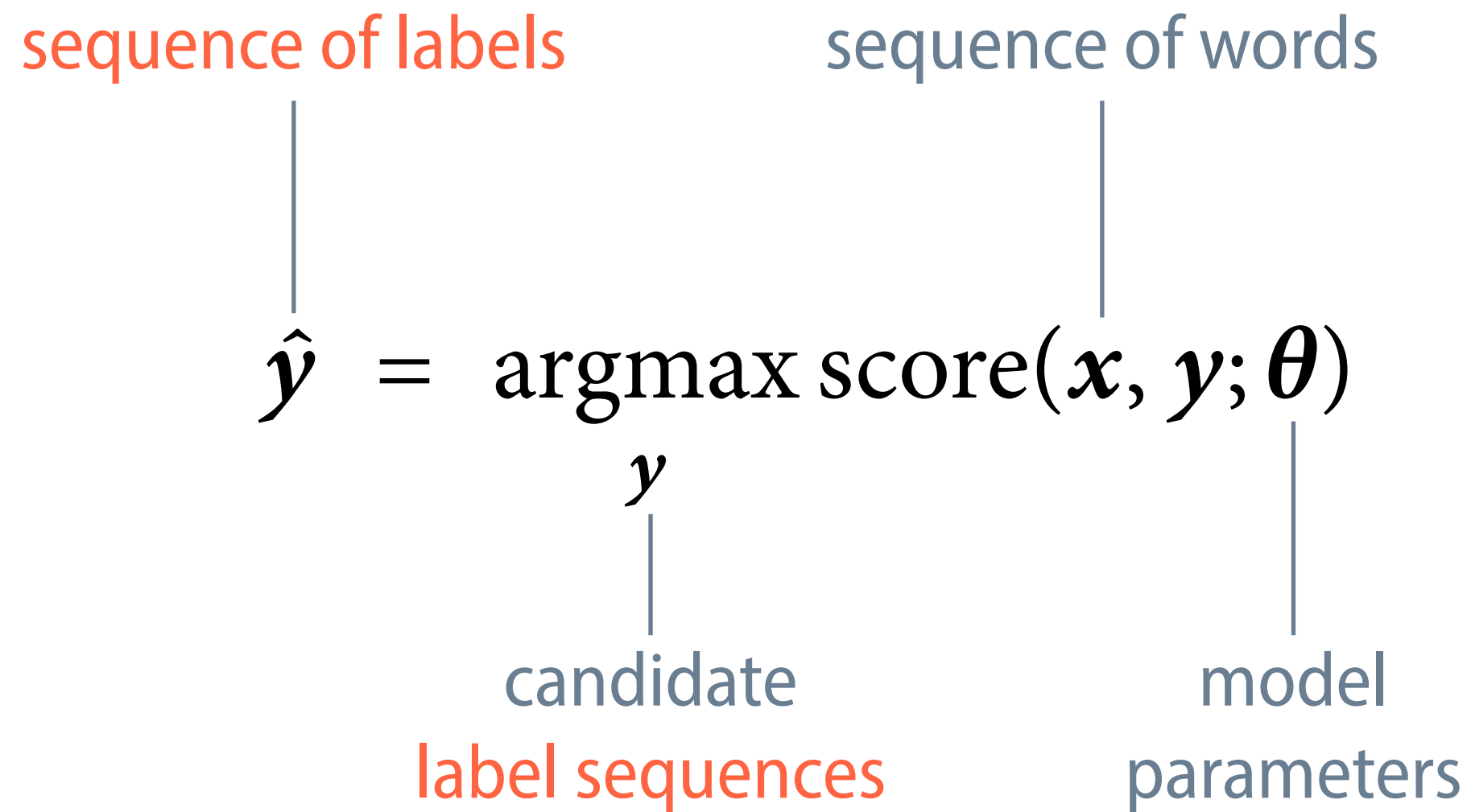
Cast sequence labelling as a sequence of classification problems: predict the label for each position in the input sequence.

- **Global search**

Cast sequence labelling as a combinatorial optimisation problem over the full set of candidate sequences.

requires specialised algorithms

Sequence labelling tasks



Ambiguity causes combinatorial explosion

I	want	to	live	in	peace
PRON	VERB	PART	VERB	ADP	NOUN
NOUN	NOUN	ADP	ADJ	ADV	VERB
		ADV	ADV	ADJ	
				NOUN	

288 possible tag sequences

Dealing with combinatorial explosion

- The number of candidate sequences is exponential in the length of the input sequence, so naive optimisation is doomed to fail.
- To make the search problem tractable, we will make assumptions about the scoring function.
- Then, in some cases, we can design special algorithms to solve the optimisation problem in polynomial time.

First-order factorized scoring function

candidate
output sequence

sequence
length

$$\text{score}(\mathbf{x}, \mathbf{y}; \boldsymbol{\theta}) = \sum_{i=1}^L \text{score}^e(\mathbf{x}, y_i; \boldsymbol{\theta}) + \sum_{i=1}^L \text{score}^t(\mathbf{x}, y_{i-1}, y_i; \boldsymbol{\theta})$$

input
sequence

score for a
single label

score for a
pair of labels

emission score

transition score

Maximum Entropy Markov Model (MEMM)

sequence of labels

sequence of words

$$P(\mathbf{y} | \mathbf{x}) = \prod_{i=1}^L \frac{\exp(\text{score}(\mathbf{x}, y_{i-1}, y_i))}{\sum_{y'} \exp(\text{score}(\mathbf{x}, y_{i-1}, y'))}$$

candidate
label

Algorithmic problems

- **Decoding:** For a trained model, we want to find the most probable label sequence \mathbf{y} , given the input sequence \mathbf{x} .

involves the search over exponentially many candidate sequences \rightarrow Viterbi

- **Training:** To train a model, we want to minimise the negative log likelihood on gold-standard examples $((\mathbf{x}, y_{i-1}), y_i)$.

standard softmax regression problem

- Search during decoding is global, search during training is local.

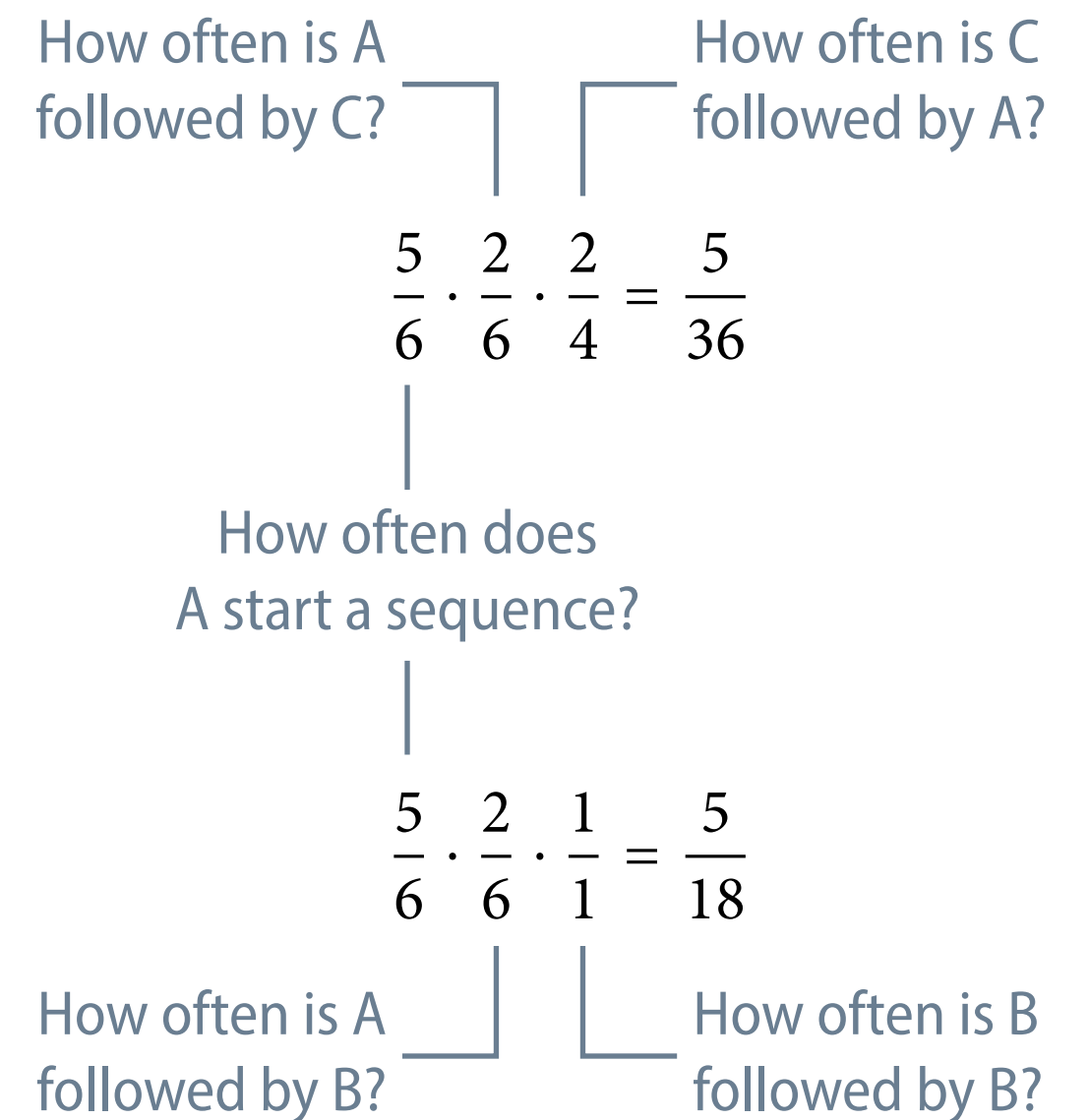
The label bias problem

Zhang and Teng (2021)

Sequence id	Sequence
1	ACA
2	ACA
3	AA
4	AAB
5	ABB
6	CCC

Global probability

First-order probability



$\frac{2}{6}$

$\frac{1}{6}$

Conditional random field

Lafferty et al. (2001)

sequence of labels

sequence of words

$$P(\mathbf{y} | \mathbf{x}) = \frac{\exp(\text{score}(\mathbf{x}, \mathbf{y}))}{\sum_{\mathbf{y}'} \exp(\text{score}(\mathbf{x}, \mathbf{y}'))}$$

partition
function

candidate
label sequence

Algorithmic problems

- **Decoding:** For a trained model, we want to find the most probable label sequence \mathbf{y} , given the input sequence \mathbf{x} .

involves the search over exponentially many candidate sequences \rightarrow Viterbi

- **Training:** To train a model, we want to minimise the negative log likelihood on gold-standard examples (\mathbf{x}, \mathbf{y}) .

involves the search over exponentially many candidate sequences \rightarrow Viterbi

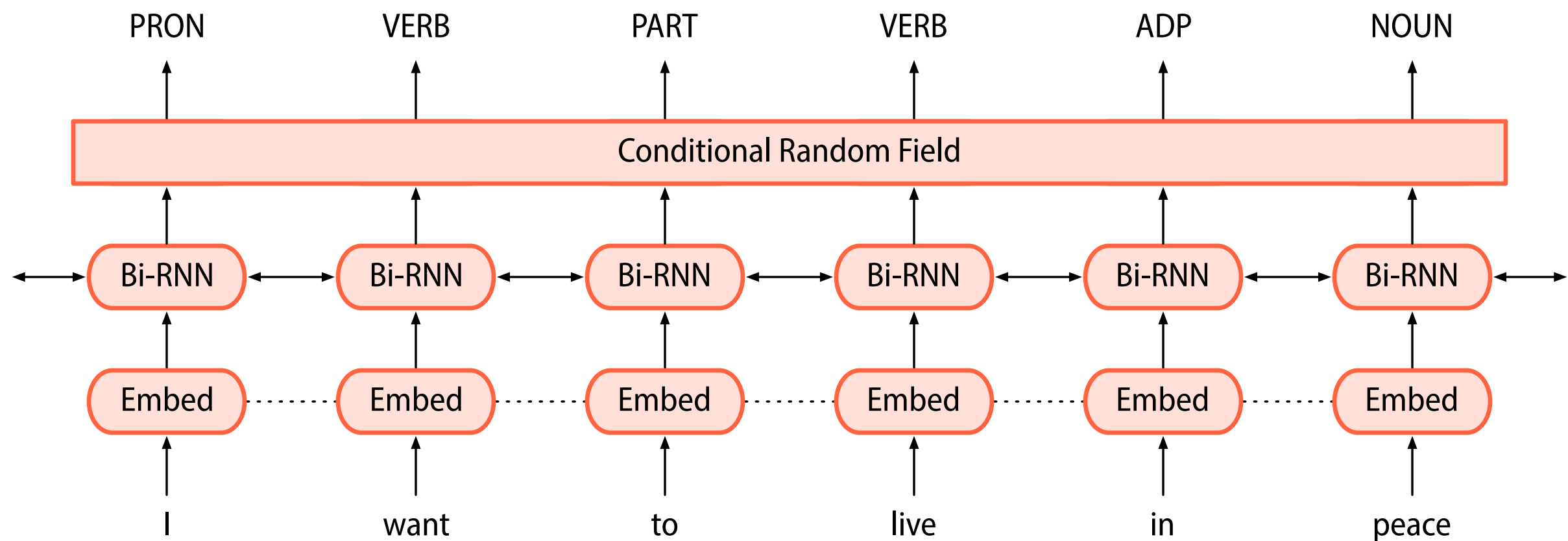
- Search is global during both decoding and training.

Finding the highest-probability sequence

$$\begin{aligned}\arg \max_y P(\mathbf{y} | \mathbf{x}) &= \arg \max_y \frac{\exp(\text{score}(\mathbf{x}, \mathbf{y}))}{\sum_{\mathbf{y}'} \exp(\text{score}(\mathbf{x}, \mathbf{y}'))} \\ &= \arg \max_y \exp(\text{score}(\mathbf{x}, \mathbf{y})) \\ &= \arg \max_y \text{score}(\mathbf{x}, \mathbf{y})\end{aligned}$$

We do not need the partition function to find the highest-scoring output.

Adding a CRF to a Bi-RNN model



PyTorch-compatible implementation available in [AllenNLP](#)