

Lecture 9

Sequence Tagging

Part 1: Sequence Tagging

Machine Learning for Structured Data
Vlad Niculae · LTL, UvA · <https://vene.ro/mlsd>

Outline:

① Sequence Tagging

Definition and examples

Evaluation

② Different Scoring Models

A Simple Scoring Function

A Better Scoring Model

③ Sequence Tagging Algorithms

Dynamic Programming For Sequence Tagging

Putting It All Together

Sequence Tagging

Given a sequence of n items $\mathbf{x} = (x_1, \dots, x_n)$, assign to each of them one of K tags:

$$\mathbf{y} = (y_1, \dots, y_n) \quad \text{where each } y_i \in \{1, \dots, K\}.$$

Sequence Tagging

Given a sequence of n items $\mathbf{x} = (x_1, \dots, x_n)$, assign to each of them one of K tags:

$$\mathbf{y} = (y_1, \dots, y_n) \quad \text{where each } y_i \in \{1, \dots, K\}.$$

Example 1: **Part-of-speech (POS) tagging** in NLP

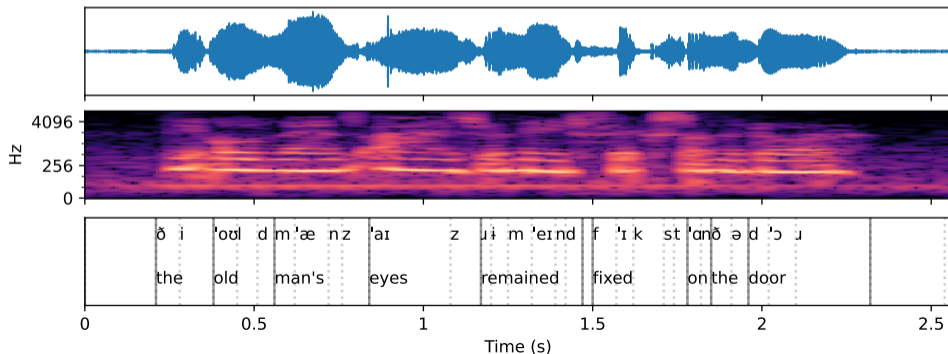
	the	old	man	the	boat
\mathbf{y}_a	det	adj	noun	det	noun
\mathbf{y}_b	det	noun	verb	det	noun

Sequence Tagging

Given a sequence of n items $\mathbf{x} = (x_1, \dots, x_n)$, assign to each of them one of K tags:

$$\mathbf{y} = (y_1, \dots, y_n) \quad \text{where each } y_i \in \{1, \dots, K\}.$$

Example 2: **Frame-level phoneme classification** (may be part of speech recognition)

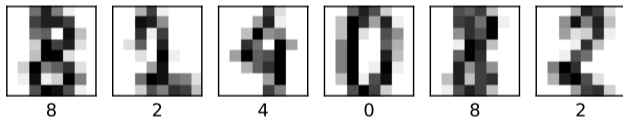


Sequence Tagging

Given a sequence of n items $\mathbf{x} = (x_1, \dots, x_n)$, assign to each of them one of K tags:

$$\mathbf{y} = (y_1, \dots, y_n) \quad \text{where each } y_i \in \{1, \dots, K\}.$$

Example 3: **Optical character recognition**



Characterizing The Output Space

Given a sequence of n items $\mathbf{x} = (x_1, \dots, x_n)$, assign to each of them one of K tags:

$$\mathbf{y} = (y_1, \dots, y_n) \quad \text{where each } y_i \in \{1, \dots, K\}.$$

Input $\mathbf{x} = (x_1, \dots, x_n)$, e.g., a sequence of words.

Output $\mathbf{y} = (y_1, \dots, y_n)$, e.g., a sequence of part-of-speech tags.

For each data point (sentence), $|\mathbf{y}| = |\mathbf{x}|$; different data points have different lengths.

Characterizing The Output Space

Given a sequence of n items $\mathbf{x} = (x_1, \dots, x_n)$, assign to each of them one of K tags:

$$\mathbf{y} = (y_1, \dots, y_n) \quad \text{where each } y_i \in \{1, \dots, K\}.$$

Input $\mathbf{x} = (x_1, \dots, x_n)$, e.g., a sequence of words.

Output $\mathbf{y} = (y_1, \dots, y_n)$, e.g., a sequence of part-of-speech tags.

For each data point (sentence), $|\mathbf{y}| = |\mathbf{x}|$; different data points have different lengths.

For fixed length n , some possible outputs:

- $(1, 1, \dots, 1, 1) \in \mathcal{Y}$
- $(1, 1, \dots, 1, 2) \in \mathcal{Y}$
- $(K, K, \dots, K, K) \in \mathcal{Y}$

How many in terms of n ?

Part-Of-Speech Tags

	Tag	Description	Example
Open Class	ADJ	Adjective: noun modifiers describing properties	<i>red, young, awesome</i>
	ADV	Adverb: verb modifiers of time, place, manner	<i>very, slowly, home, yesterday</i>
	NOUN	words for persons, places, things, etc.	<i>algorithm, cat, mango, beauty</i>
	VERB	words for actions and processes	<i>draw, provide, go</i>
	PROPN	Proper noun: name of a person, organization, place, etc..	<i>Regina, IBM, Colorado</i>
	INTJ	Interjection: exclamation, greeting, yes/no response, etc.	<i>oh, um, yes, hello</i>
Closed Class Words	ADP	Adposition (Preposition/Postposition): marks a noun's spacial, temporal, or other relation	<i>in, on, by, under</i>
	AUX	Auxiliary: helping verb marking tense, aspect, mood, etc.,	<i>can, may, should, are</i>
	CCONJ	Coordinating Conjunction: joins two phrases/clauses	<i>and, or, but</i>
	DET	Determiner: marks noun phrase properties	<i>a, an, the, this</i>
	NUM	Numeral	<i>one, two, first, second</i>
	PART	Particle: a preposition-like form used together with a verb	<i>up, down, on, off, in, out, at, by</i>
	PRON	Pronoun: a shorthand for referring to an entity or event	<i>she, who, I, others</i>
SCONJ	Subordinating Conjunction: joins a main clause with a subordinate clause such as a sentential complement	<i>that, which</i>	
Other	PUNCT	Punctuation	<i>; , ()</i>
	SYM	Symbols like \$ or emoji	<i>\$, %</i>
	X	Other	<i>asdf, qwfg</i>

Figure 8.1 The 17 parts of speech in the Universal Dependencies tagset (Nivre et al., 2016a). Features can be added to make finer-grained distinctions (with properties like number, case, definiteness, and so on).

POS Tagging Evaluation

Evaluation: sequence-level accuracy

$$\frac{\sum_{i=1}^{N_{\text{valid}}} \mathbf{y}^{(i)} = \hat{\mathbf{y}}^{(i)}}{N_{\text{valid}}}$$

or micro-averaged tag accuracy (writing $n^{(i)} = |\mathbf{y}^{(i)}|$):

$$\frac{\sum_{i=1}^{N_{\text{valid}}} \sum_{j=1}^{n^{(i)}} y_j^{(i)} = \hat{y}_j^{(i)}}{\sum_{i=1}^{N_{\text{valid}}} n^{(i)}}$$

Example:

<i>true:</i>	PRO	VERB	NUM	NOUN	ADV
<i>pred:</i>	PRO	VERB	NUM	NOUN	PRO
<i>words:</i>	there	are	70	children	there

<i>true:</i>	INTJ
<i>pred:</i>	X
<i>words:</i>	eeeeek

Lecture 9

Sequence Tagging

Part 2: Different Scoring Models

Machine Learning for Structured Data
Vlad Niculae · LTL, UvA · <https://vene.ro/mlsd>

Outline:

① Sequence Tagging

Definition and examples

Evaluation

② Different Scoring Models

A Simple Scoring Function

A Better Scoring Model

③ Sequence Tagging Algorithms

Dynamic Programming For Sequence Tagging

Putting It All Together

Designing A Simple Scorer

Writing $\mathbf{y} = (y_1, \dots, y_n)$, take

$$\text{score}(\mathbf{y}) = \sum_j a_{j,y_j}.$$

\mathbf{A} is a matrix of scores,
e.g., computed by a NN encoder.

	det	noun	adj	verb
the	5	0	0	0
old	0	1	3	0
man	0	3	0	1
the	5	0	0	0
boat	0	5	0	0

Designing A Simple Scorer

Writing $\mathbf{y} = (y_1, \dots, y_n)$, take

$$\text{score}(\mathbf{y}) = \sum_j a_{j,y_j}.$$

\mathbf{A} is a matrix of scores,
e.g., computed by a NN encoder.

	the	old	man	the	boat
\mathbf{y}_a	det	adj	noun	det	noun
\mathbf{y}_b	det	noun	verb	det	noun

	det	noun	adj	verb
the	5	0	0	0
old	0	1	3	0
man	0	3	0	1
the	5	0	0	0
boat	0	5	0	0

Designing A Simple Scorer

Writing $\mathbf{y} = (y_1, \dots, y_n)$, take

$$\text{score}(\mathbf{y}) = \sum_j a_{j,y_j}.$$

\mathbf{A} is a matrix of scores,
e.g., computed by a NN encoder.

	the	old	man	the	boat
\mathbf{y}_a	det	adj	noun	det	noun
\mathbf{y}_b	det	noun	verb	det	noun

$$\text{score}(\mathbf{y}_a) =$$

	det	noun	adj	verb
the	5	0	0	0
old	0	1	3	0
man	0	3	0	1
the	5	0	0	0
boat	0	5	0	0

Designing A Simple Scorer

Writing $\mathbf{y} = (y_1, \dots, y_n)$, take

$$\text{score}(\mathbf{y}) = \sum_j a_{j,y_j}.$$

\mathbf{A} is a matrix of scores,
e.g., computed by a NN encoder.

	the	old	man	the	boat
\mathbf{y}_a	det	adj	noun	det	noun
\mathbf{y}_b	det	noun	verb	det	noun

$$\text{score}(\mathbf{y}_a) = 21$$

	det	noun	adj	verb
the	5	0	0	0
old	0	1	3	0
man	0	3	0	1
the	5	0	0	0
boat	0	5	0	0

Designing A Simple Scorer

Writing $\mathbf{y} = (y_1, \dots, y_n)$, take

$$\text{score}(\mathbf{y}) = \sum_j a_{j,y_j}.$$

\mathbf{A} is a matrix of scores,
e.g., computed by a NN encoder.

	det	noun	adj	verb
the	5	0	0	0
old	0	1	3	0
man	0	3	0	1
the	5	0	0	0
boat	0	5	0	0

	the	old	man	the	boat
\mathbf{y}_a	det	adj	noun	det	noun
\mathbf{y}_b	det	noun	verb	det	noun

$$\text{score}(\mathbf{y}_a) = 21$$

$$\text{score}(\mathbf{y}_b) =$$

Designing A Simple Scorer

Writing $\mathbf{y} = (y_1, \dots, y_n)$, take

$$\text{score}(\mathbf{y}) = \sum_j a_{j,y_j}.$$

\mathbf{A} is a matrix of scores,
e.g., computed by a NN encoder.

	the	old	man	the	boat
\mathbf{y}_a	det	adj	noun	det	noun
\mathbf{y}_b	det	noun	verb	det	noun

$$\text{score}(\mathbf{y}_a) = 21$$

$$\text{score}(\mathbf{y}_b) = 17$$

	det	noun	adj	verb
the	5	0	0	0
old	0	1	3	0
man	0	3	0	1
the	5	0	0	0
boat	0	5	0	0

Designing A Simple Scorer

A first attempt:

separate classifier for each position.

1. embed and encode x , eg, with a CNN.

$$(x_1, \dots, x_n) \rightarrow (z_1, \dots, z_n)$$

2. For each position j , apply a classification head with K outputs. E.g.,

$$a_j = W^T z_j + b$$

Think of A as a matrix with n rows and K columns, where $a_{j,c}$ is the score of assigning tag c at position j .

3. Writing $y = (y_1, \dots, y_n)$, take $\text{score}(y) = \sum_j a_{j,y_j}$.

```
words = [21, 79, 14] # indices
emb = Embedding(vocab_sz, dim)
clf = Linear(dim, n_tags)
```

```
# optionally add RNN, CNN, whatever
```

```
Z = emb(words) # (3 × dim)
A = clf(Z)      # (3 × n_tags)
```

```
# computing the score of a given tag sequence:
y = [2, 0, 2]
```

```
y_score = sum(A[i, yi]
               for y, yi in enumerate(y))
```

```
# or, if you want to be fancy/fast:
y_score = A[torch.arange(len(y)), y].sum()
```

Finding The Best sequence

With our $\text{score}(\mathbf{y}) = \sum_j a_{j,y_j}$, can we compute:

$$\max_{\mathbf{y} \in \mathcal{Y}} \text{score}(\mathbf{y})$$

	det	noun	adj	verb
the	5	0	0	0
old	0	1	3	0
man	0	3	0	1
the	5	0	0	0
boat	0	5	0	0

Finding The Best sequence

With our $\text{score}(\mathbf{y}) = \sum_j a_{j,y_j}$, can we compute:

$$\begin{aligned} & \max_{\mathbf{y} \in \mathcal{Y}} \text{score}(\mathbf{y}) \\ = & \max_{y_1 \in [K], \dots, y_n \in [K]} \text{score}([y_1, \dots, y_n]) \end{aligned}$$

	det	noun	adj	verb
the	5	0	0	0
old	0	1	3	0
man	0	3	0	1
the	5	0	0	0
boat	0	5	0	0

Finding The Best sequence

With our $\text{score}(\mathbf{y}) = \sum_j a_{j,y_j}$, can we compute:

$$\begin{aligned} & \max_{\mathbf{y} \in \mathcal{Y}} \text{score}(\mathbf{y}) \\ = & \max_{y_1 \in [K], \dots, y_n \in [K]} \text{score}([y_1, \dots, y_n]) \\ = & \max_{y_1 \in [K], \dots, y_n \in [K]} \sum_j a_{j,y_j} \end{aligned}$$

	det	noun	adj	verb
the	5	0	0	0
old	0	1	3	0
man	0	3	0	1
the	5	0	0	0
boat	0	5	0	0

Finding The Best sequence

With our $\text{score}(\mathbf{y}) = \sum_j a_{j,y_j}$, can we compute:

$$\begin{aligned} & \max_{\mathbf{y} \in \mathcal{Y}} \text{score}(\mathbf{y}) \\ &= \max_{y_1 \in [K], \dots, y_n \in [K]} \text{score}([y_1, \dots, y_n]) \\ &= \max_{y_1 \in [K], \dots, y_n \in [K]} \sum_j a_{j,y_j} \\ &= \sum_j \max_{y_j \in [K]} a_{j,y_j} \end{aligned}$$

	det	noun	adj	verb
the	5	0	0	0
old	0	1	3	0
man	0	3	0	1
the	5	0	0	0
boat	0	5	0	0

Finding The Best sequence

With our $\text{score}(\mathbf{y}) = \sum_j a_{j,y_j}$, can we compute:

$$\begin{aligned} & \max_{\mathbf{y} \in \mathcal{Y}} \text{score}(\mathbf{y}) \\ &= \max_{y_1 \in [K], \dots, y_n \in [K]} \text{score}([y_1, \dots, y_n]) \\ &= \max_{y_1 \in [K], \dots, y_n \in [K]} \sum_j a_{j,y_j} \\ &= \sum_j \max_{y_j \in [K]} a_{j,y_j} \end{aligned}$$

	det	noun	adj	verb
the	5	0	0	0
old	0	1	3	0
man	0	3	0	1
the	5	0	0	0
boat	0	5	0	0

$\mathbf{A} =$

So, $\arg \max_{\mathbf{y}} \text{score}(\mathbf{y})$ is made up of the tags selected independently at each position.

Normalizing Constant (log-sum-exp)

With our $\text{score}(\mathbf{y}) = \sum_j a_{j,y_j}$, can we compute:

$$\log \sum_{\mathbf{y} \in \mathcal{Y}} \exp(\text{score}(\mathbf{y}))$$

	det	noun	adj	verb
the	5	0	0	0
old	0	1	3	0
man	0	3	0	1
the	5	0	0	0
boat	0	5	0	0

Normalizing Constant (log-sum-exp)

With our score(\mathbf{y}) = $\sum_j a_{j,y_j}$, can we compute:

$$\begin{aligned} & \log \sum_{\mathbf{y} \in \mathcal{Y}} \exp(\text{score}(\mathbf{y})) \\ &= \log \sum_{y_1=1}^K \dots \sum_{y_n=1}^K \exp \sum_{j=1}^n a_{j,y_j} \end{aligned}$$

	det	noun	adj	verb
the	5	0	0	0
old	0	1	3	0
man	0	3	0	1
the	5	0	0	0
boat	0	5	0	0

$\mathbf{A} =$

Normalizing Constant (log-sum-exp)

With our $\text{score}(\mathbf{y}) = \sum_j a_{j,y_j}$, can we compute:

$$\begin{aligned} & \log \sum_{\mathbf{y} \in \mathcal{Y}} \exp(\text{score}(\mathbf{y})) \\ &= \log \sum_{y_1=1}^K \dots \sum_{y_n=1}^K \exp \sum_{j=1}^n a_{j,y_j} \\ &= \log \sum_{y_1=1}^K \dots \sum_{y_n=1}^K \prod_{j=1}^n \exp a_{j,y_j} \end{aligned}$$

	det	noun	adj	verb
the	5	0	0	0
old	0	1	3	0
man	0	3	0	1
the	5	0	0	0
boat	0	5	0	0

$\mathbf{A} =$

Normalizing Constant (log-sum-exp)

With our $\text{score}(\mathbf{y}) = \sum_j a_{j,y_j}$, can we compute:

$$\begin{aligned} & \log \sum_{\mathbf{y} \in \mathcal{Y}} \exp(\text{score}(\mathbf{y})) \\ &= \log \sum_{y_1=1}^K \dots \sum_{y_n=1}^K \exp \sum_{j=1}^n a_{j,y_j} \\ &= \log \sum_{y_1=1}^K \dots \sum_{y_n=1}^K \prod_{j=1}^n \exp a_{j,y_j} \\ &= \log \prod_{j=1}^n \sum_{y_j=1}^K \exp a_{j,y_j} \end{aligned}$$

	det	noun	adj	verb
the	5	0	0	0
old	0	1	3	0
man	0	3	0	1
the	5	0	0	0
boat	0	5	0	0

$\mathbf{A} =$

Normalizing Constant (log-sum-exp)

With our $\text{score}(\mathbf{y}) = \sum_j a_{j,y_j}$, can we compute:

$$\begin{aligned} & \log \sum_{\mathbf{y} \in \mathcal{Y}} \exp(\text{score}(\mathbf{y})) \\ &= \log \sum_{y_1=1}^K \dots \sum_{y_n=1}^K \exp \sum_{j=1}^n a_{j,y_j} \\ &= \log \sum_{y_1=1}^K \dots \sum_{y_n=1}^K \prod_{j=1}^n \exp a_{j,y_j} \\ &= \log \prod_{j=1}^n \sum_{y_j=1}^K \exp a_{j,y_j} \\ &= \sum_{j=1}^n \log \sum_{y_j=1}^K \exp a_{j,y_j} \end{aligned}$$

	det	noun	adj	verb
the	5	0	0	0
old	0	1	3	0
man	0	3	0	1
the	5	0	0	0
boat	0	5	0	0

$\mathbf{A} =$

Normalizing Constant (log-sum-exp)

With our $\text{score}(\mathbf{y}) = \sum_j a_{j,y_j}$, can we compute:

$$\begin{aligned} & \log \sum_{\mathbf{y} \in \mathcal{Y}} \exp(\text{score}(\mathbf{y})) \\ &= \log \sum_{y_1=1}^K \dots \sum_{y_n=1}^K \exp \sum_{j=1}^n a_{j,y_j} \\ &= \log \sum_{y_1=1}^K \dots \sum_{y_n=1}^K \prod_{j=1}^n \exp a_{j,y_j} \\ &= \log \prod_{j=1}^n \sum_{y_j=1}^K \exp a_{j,y_j} \\ &= \sum_{j=1}^n \log \sum_{y_j=1}^K \exp a_{j,y_j} \end{aligned}$$

	det	noun	adj	verb
the	5	0	0	0
old	0	1	3	0
man	0	3	0	1
the	5	0	0	0
boat	0	5	0	0

$\mathbf{A} =$

Probabilistic interpretation: independence

$$\begin{aligned} \log \Pr(\mathbf{y}) &= \text{score}(\mathbf{y}) - \log \sum_{\mathbf{y}' \in \mathcal{Y}} \exp \text{score}(\mathbf{y}') \\ &= \sum_j \underbrace{\left(a_{j,y_j} - \log \sum_{k \in [K]} \exp a_{j,k} \right)}_{\log \Pr(y_j)} \end{aligned}$$

Fully-Local vs. Fully-Global

For sequence tagging, the separable (fully-local) score

$$\text{score}(\mathbf{y}) = \sum_j a_{j,y_j}$$

amounts to applying a probabilistic classifier to each of the n positions separately!
(any “magic” comes from the feature representation / neural net encoder.)

Can we design a richer score(\mathbf{y}) taking into account the sequential structure of \mathbf{y} ?

Fully-Local vs. Fully-Global

Entirely global model: like classification, where *each possible sequence* is a class.

	\mathbf{y}	score(\mathbf{y})
	det det det det det	-1000
	det det det det noun	-940
	det det det det verb	-800
	...	
	det noun verb det noun	400
	...	
	verb verb verb verb verb	-1100

As expressive as possible: score is any function of the sequence.

Fully-Local vs. Fully-Global

Entirely global model: like classification, where *each possible sequence* is a class.

	\mathbf{y}	score(\mathbf{y})
	det det det det det	-1000
	det det det det noun	-940
	det det det det verb	-800
	...	
	det noun verb det noun	400
	...	
	verb verb verb verb verb	-1100

As expressive as possible: score is any function of the sequence.

But completely intractable: $O(K^n)$ time and space.

Fully-Local vs. Fully-Global

Entirely global model: like classification, where *each possible sequence* is a class.

	\mathbf{y}	score(\mathbf{y})
	det det det det det	-1000
	det det det det noun	-940
	det det det det verb	-800
	...	
	det noun verb det noun	400
	...	
	verb verb verb verb verb	-1100

As expressive as possible: score is any function of the sequence.

But completely intractable: $O(K^n)$ time and space.

Structure output prediction is about the space in between these two extremes.

Idea: scoring transitions between adjacent tags

$$\text{score}(\mathbf{y}) = \sum_{j=1}^n a_{j,y_j} + \sum_{j=2}^n t_{y_{j-1},y_j}$$

For example, $\text{score}([\text{NOUN}, \text{DET}, \text{VERB}]) = +a_{2,\text{DET}} a_{1,\text{NOUN}} + a_{3,\text{VERB}} + t_{\text{NOUN},\text{DET}} + t_{\text{DET},\text{VERB}}$

Scoring Transitions Between Tags

A rich scorer that takes into account the sequential nature of \mathbf{y} while still allowing efficient computation:

scoring transitions between adjacent tags

$$\text{score}(\mathbf{y}) = \sum_{j=1}^n a_{j,y_j} + \sum_{j=2}^n t_{y_{j-1},y_j}$$

For example, $\text{score}([\text{NOUN}, \text{DET}, \text{VERB}]) = a_{1,\text{NOUN}} + a_{2,\text{DET}} + a_{3,\text{VERB}} + t_{\text{NOUN},\text{DET}} + t_{\text{DET},\text{VERB}}$

Sequence Modeling With Transition Scores

$$\text{score}(\mathbf{y}) = \sum_{j=1}^n a_{j,y_j} + \sum_{j=2}^n t_{y_{j-1},y_j}$$

The tag scores $\mathbf{A} \in \mathbb{R}^{n \times K}$ can be computed as before (e.g., with a convnet.)

The transition scores $\mathbf{T} \in \mathbb{R}^{K \times K}$:

- could be a learned parameter. (size does not depend on n)
- could be predicted by the neural net as a function of \mathbf{x} .

Unlike in the separable case, with transition scores, we no longer get n parallel classifiers: the different tags impact one another. (This makes the model more expressive and more interesting.)

Lecture 9

Sequence Tagging

Part 3: Sequence Tagging Algorithms

Machine Learning for Structured Data
Vlad Niculae · LTL, UvA · <https://vene.ro/mlsd>

Outline:

① Sequence Tagging

Definition and examples

Evaluation

② Different Scoring Models

A Simple Scoring Function

A Better Scoring Model

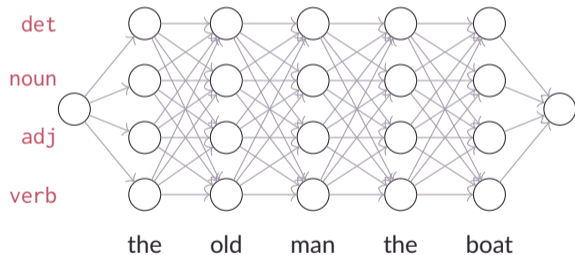
③ Sequence Tagging Algorithms

Dynamic Programming For Sequence Tagging

Putting It All Together

Sequence Tagging As A DAG

$$\text{score}(\mathbf{y}) = \sum_{j=1}^n a_{j,y_j} + \sum_{j=2}^n t_{y_{j-1},y_j}$$



$G = (V, E, w)$ where:

$$V = \{(j, c) : j \in [n], c \in [K]\} \\ \cup \{s, t\}$$

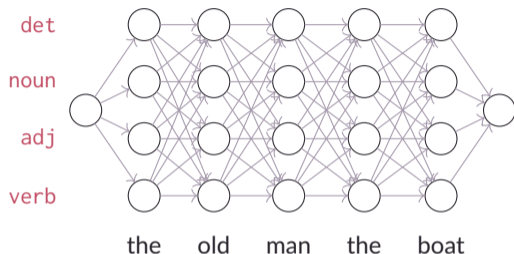
$$E = \{(j-1, c') \rightarrow (j, c) : j \in [2, n], c, c' \in [K]\} \\ \cup \{s \rightarrow (1, c) : c \in [K]\} \\ \cup \{(n, c) \rightarrow t : c \in [K]\}$$

$$w((j-1, c') \rightarrow (j, c)) = a_{j,c} + t_{c',c} \\ w(s \rightarrow (1, c)) = a_{1,c} \\ w((n, c) \rightarrow t) = 0$$

$$|V| \in \Theta(nK); \quad |E| \in \Theta(nK^2)$$

Topological ordering?

Viterbi For Sequence Tagging



General Viterbi (reminder sketch)

initialize $m_1 \leftarrow 0$

for $i = 2, \dots, n$ do

$$m_i \leftarrow \max_{j \in P_i} (m_j + w(ji))$$

$$\pi_i \leftarrow \arg \max_{j \in P_i} (m_j + w(ji))$$

follow backpointers to get best path

Viterbi for sequence tagging

input: Unary scores \mathbf{A} ($n \times K$ array)

Transition scores \mathbf{T} ($K \times K$ array)

Forward: compute scores recursively

$$m_{1c} = a_{1c} \quad \text{for all } c \in [K]$$

for $j = 2$ to n do

for $c = 1$ to K do

$$m_{j,c} \leftarrow \max_{c' \in [K]} (m_{j-1,c'} + a_{j,c} + t_{c',c})$$

$$\pi_{j,c} \leftarrow \arg \max_{c' \in [K]} (m_{j-1,c'} + a_{j,c} + t_{c',c})$$

$$f^* = \max_{c' \in [K]} m_{n,c'}$$

Backward: follow backpointers

$$y_n = \arg \max_{c'} m_n(c')$$

for $j = n - 1$ down to 1 do

$$y_j = \pi_{j+1,y_{j+1}}$$

output: f^* and $\mathbf{y}^* = [y_1, \dots, y_n]$

Viterbi For Sequence Tagging: Example

$m_{j,c}$ is stored as a matrix \mathbf{M} , same shape as \mathbf{A} .

Apply $m_{1,c} = a_{1,c}$ to get the first row: (copied from \mathbf{A})

Then iteratively: $m_{j,c} = \max_{c' \in [K]} m_{j-1,c'} + a_{j,c} + t_{c',c}$

At the end, take the maximum over the last row.

	det	noun	adj	verb
$\mathbf{M} =$	the			
	old			
	man			
	the			
	boat			

unary and transition scores:

	det	noun	adj	verb	
$\mathbf{A} =$	the	5	0	0	0
	old	0	1	3	0
	man	0	3	0	1
	the	5	0	0	0
	boat	0	5	0	0

	det	noun	adj	verb	
$\mathbf{T} =$	det	-4	3	2	-1
	noun	-3	-2	-1	2
	adj	-2	2	1	1
	verb	1	-1	0	0

Viterbi For Sequence Tagging: Example

$m_{j,c}$ is stored as a matrix M , same shape as A .

Apply $m_{1,c} = a_{1,c}$ to get the first row: (copied from A)

Then iteratively: $m_{j,c} = \max_{c' \in [K]} m_{j-1,c'} + a_{j,c} + t_{c',c}$

At the end, take the maximum over the last row.

	det	noun	adj	verb
the	5	0	0	0
old				
man				
the				
boat				

unary and transition scores:

	det	noun	adj	verb
the	5	0	0	0
old	0	1	3	0
man	0	3	0	1
the	5	0	0	0
boat	0	5	0	0

	det	noun	adj	verb
det	-4	3	2	-1
noun	-3	-2	-1	2
adj	-2	2	1	1
verb	1	-1	0	0

Viterbi For Sequence Tagging: Example

$m_{j,c}$ is stored as a matrix M , same shape as A .

Apply $m_{1,c} = a_{1,c}$ to get the first row: (copied from A)

Then iteratively: $m_{j,c} = \max_{c' \in [K]} m_{j-1,c'} + a_{j,c} + t_{c',c}$

At the end, take the maximum over the last row.

	det	noun	adj	verb
the	5	0	0	0
old	↓	↙	↘	↘
man				
the				
boat				

unary and transition scores:

	det	noun	adj	verb
the	5	0	0	0
old	0	1	3	0
man	0	3	0	1
the	5	0	0	0
boat	0	5	0	0

	det	noun	adj	verb
det	-4	3	2	-1
noun	-3	-2	-1	2
adj	-2	2	1	1
verb	1	-1	0	0

Viterbi For Sequence Tagging: Example

$m_{j,c}$ is stored as a matrix \mathbf{M} , same shape as \mathbf{A} .

Apply $m_{1,c} = a_{1,c}$ to get the first row: (copied from \mathbf{A})

Then iteratively: $m_{j,c} = \max_{c' \in [K]} m_{j-1,c'} + a_{j,c} + t_{c',c}$

At the end, take the maximum over the last row.

	det	noun	adj	verb
the	5	0	0	0
old	1			
man				
the				
boat				

unary and transition scores:

	det	noun	adj	verb
the	5	0	0	0
old	0	1	3	0
man	0	3	0	1
the	5	0	0	0
boat	0	5	0	0

	det	noun	adj	verb
det	-4	3	2	-1
noun	-3	-2	-1	2
adj	-2	2	1	1
verb	1	-1	0	0

Viterbi For Sequence Tagging: Example

$m_{j,c}$ is stored as a matrix M , same shape as A .

Apply $m_{1,c} = a_{1,c}$ to get the first row: (copied from A)

Then iteratively: $m_{j,c} = \max_{c' \in [K]} m_{j-1,c'} + a_{j,c} + t_{c',c}$

At the end, take the maximum over the last row.

	det	noun	adj	verb
the	5	0	0	0
old	1			
man				
the				
boat				

unary and transition scores:

	det	noun	adj	verb
the	5	0	0	0
old	0	1	3	0
man	0	3	0	1
the	5	0	0	0
boat	0	5	0	0

	det	noun	adj	verb
det	-4	3	2	-1
noun	-3	-2	-1	2
adj	-2	2	1	1
verb	1	-1	0	0

Viterbi For Sequence Tagging: Example

$m_{j,c}$ is stored as a matrix M , same shape as A .

Apply $m_{1,c} = a_{1,c}$ to get the first row: (copied from A)

Then iteratively: $m_{j,c} = \max_{c' \in [K]} m_{j-1,c'} + a_{j,c} + t_{c',c}$

At the end, take the maximum over the last row.

	det	noun	adj	verb
the	5	0	0	0
old	1	9		
man				
the				
boat				

unary and transition scores:

	det	noun	adj	verb
the	5	0	0	0
old	0	1	3	0
man	0	3	0	1
the	5	0	0	0
boat	0	5	0	0

	det	noun	adj	verb
det	-4	3	2	-1
noun	-3	-2	-1	2
adj	-2	2	1	1
verb	1	-1	0	0

Viterbi For Sequence Tagging: Example

$m_{j,c}$ is stored as a matrix \mathbf{M} , same shape as \mathbf{A} .

Apply $m_{1,c} = a_{1,c}$ to get the first row: (copied from \mathbf{A})

Then iteratively: $m_{j,c} = \max_{c' \in [K]} m_{j-1,c'} + a_{j,c} + t_{c',c}$

At the end, take the maximum over the last row.

	det	noun	adj	verb
the	5	0	0	0
old	1	9	10	4
man				
the				
boat				

unary and transition scores:

	det	noun	adj	verb
the	5	0	0	0
old	0	1	3	0
man	0	3	0	1
the	5	0	0	0
boat	0	5	0	0

	det	noun	adj	verb
det	-4	3	2	-1
noun	-3	-2	-1	2
adj	-2	2	1	1
verb	1	-1	0	0

Viterbi For Sequence Tagging: Example

$m_{j,c}$ is stored as a matrix M , same shape as A .

Apply $m_{1,c} = a_{1,c}$ to get the first row: (copied from A)

Then iteratively: $m_{j,c} = \max_{c' \in [K]} m_{j-1,c'} + a_{j,c} + t_{c',c}$

At the end, take the maximum over the last row.

	det	noun	adj	verb
the	5	0	0	0
old	1	9	10	4
man	8	15	11	12
the	18	13	14	17
boat	18	26	20	17

unary and transition scores:

	det	noun	adj	verb
the	5	0	0	0
old	0	1	3	0
man	0	3	0	1
the	5	0	0	0
boat	0	5	0	0

	det	noun	adj	verb
det	-4	3	2	-1
noun	-3	-2	-1	2
adj	-2	2	1	1
verb	1	-1	0	0

Viterbi For Sequence Tagging: Example

$m_{j,c}$ is stored as a matrix M , same shape as A .

Apply $m_{1,c} = a_{1,c}$ to get the first row: (copied from A)

Then iteratively: $m_{j,c} = \max_{c' \in [K]} m_{j-1,c'} + a_{j,c} + t_{c',c}$

At the end, take the maximum over the last row.

	det	noun	adj	verb
$M =$				
the	5	0	0	0
old	1	9	10	4
man	8	15	11	12
the	18	13	14	17
boat	18	26	20	17

unary and transition scores:

	det	noun	adj	verb
$A =$				
the	5	0	0	0
old	0	1	3	0
man	0	3	0	1
the	5	0	0	0
boat	0	5	0	0

	det	noun	adj	verb
$T =$				
det	-4	3	2	-1
noun	-3	-2	-1	2
adj	-2	2	1	1
verb	1	-1	0	0

To find the best tag sequence \mathbf{y}^* , keep track of the path.

Viterbi For Sequence Tagging: Example

$m_{j,c}$ is stored as a matrix M , same shape as A .

Apply $m_{1,c} = a_{1,c}$ to get the first row: (copied from A)

Then iteratively: $m_{j,c} = \max_{c' \in [K]} m_{j-1,c'} + a_{j,c} + t_{c',c}$

At the end, take the maximum over the last row.

	det	noun	adj	verb
the	5	0	0	0
old	1	9	10	4
man	8	15	11	12
the	18	13	14	17
boat	18	26	20	17

Diagram showing the path of maximum values in the matrix M:

- From (1,1) to (2,2)
- From (2,2) to (3,4)
- From (3,4) to (4,1)
- From (4,1) to (5,2)

To find the best tag sequence \mathbf{y}^* , keep track of the path.

unary and transition scores:

	det	noun	adj	verb
the	5	0	0	0
old	0	1	3	0
man	0	3	0	1
the	5	0	0	0
boat	0	5	0	0

$A =$

	det	noun	adj	verb
det	-4	3	2	-1
noun	-3	-2	-1	2
adj	-2	2	1	1
verb	1	-1	0	0

$T =$

The Two Main Recurrences Of Sequence Tagging:

(Dynamic programming applied to the sequence tagging DAG)

$$m_{j,c} = \max_{c' \in [K]} (m_{j-1,c'} + a_{jc} + t_{c'c}),$$
$$q_{j,c} = \log \sum_{c' \in [K]} \exp (q_{j-1,c'} + a_{jc} + t_{c'c}).$$

The Forward Algorithm

Forward algorithm for sequence tagging

input: Unary scores \mathbf{A} ($n \times K$ array)

Transition scores \mathbf{T} ($K \times K$ array)

Forward: compute scores recursively

$q_{1,c} = a_{1,c}$ for all $c \in [K]$

for $j = 2$ **to** n **do**

for $c = 1$ **to** K **do**

$$q_{j,c} = \log \sum_{c' \in [K]} \exp(q_{j-1,c'} + a_{j,c} + t_{c',c})$$

return $\log Z = \log \sum_{c' \in [K]} \exp(q_{n,c'})$

	the	old	man	the	boat	
y_a	det	adj	noun	det	noun	score(y_a) = 25
y_b	det	noun	verb	det	noun	score(y_b) = 26
y_c	noun	noun	noun	noun	noun	score(y_c) = 1

Applying the Forward algorithm yields

		det	noun	adj	verb
$Q =$	the	5.00	0.00	0.00	0.00
	old	1.73	9.00	10.00	4.19
	man	8.18	15.01	11.05	12.70
	the	18.88	13.92	14.37	17.03
	boat	18.08	26.88	20.90	18.38

unary and transition scores:

		det	noun	adj	verb
$A =$	the	5	0	0	0
	old	0	1	3	0
	man	0	3	0	1
	the	5	0	0	0
	boat	0	5	0	0
		det	noun	adj	verb
$T =$	det	-4	3	2	-1
	noun	-3	-2	-1	2
	adj	-2	2	1	1
	verb	1	-1	0	0

unary and transition scores:

	the	old	man	the	boat	
y_a	det	adj	noun	det	noun	score(y_a) = 25
y_b	det	noun	verb	det	noun	score(y_b) = 26
y_c	noun	noun	noun	noun	noun	score(y_c) = 1

Applying the Forward algorithm yields

		det	noun	adj	verb
$Q =$	the	5.00	0.00	0.00	0.00
	old	1.73	9.00	10.00	4.19
	man	8.18	15.01	11.05	12.70
	the	18.88	13.92	14.37	17.03
	boat	18.08	26.88	20.90	18.38

$$\log Z \approx 26.885$$

		det	noun	adj	verb
$A =$	the	5	0	0	0
	old	0	1	3	0
	man	0	3	0	1
	the	5	0	0	0
	boat	0	5	0	0
		det	noun	adj	verb
$T =$	det	-4	3	2	-1
	noun	-3	-2	-1	2
	adj	-2	2	1	1
	verb	1	-1	0	0

	the	old	man	the	boat	
y_a	det	adj	noun	det	noun	score(y_a) = 25
y_b	det	noun	verb	det	noun	score(y_b) = 26
y_c	noun	noun	noun	noun	noun	score(y_c) = 1

Applying the Forward algorithm yields

		det	noun	adj	verb
$Q =$	the	5.00	0.00	0.00	0.00
	old	1.73	9.00	10.00	4.19
	man	8.18	15.01	11.05	12.70
	the	18.88	13.92	14.37	17.03
	boat	18.08	26.88	20.90	18.38

$$\log Z \approx 26.885$$

$$\log P(y_a) = \text{score}(y_a) - \log Z = 25 - 26.885 = -1.885$$

unary and transition scores:

		det	noun	adj	verb
$A =$	the	5	0	0	0
	old	0	1	3	0
	man	0	3	0	1
	the	5	0	0	0
	boat	0	5	0	0
		det	noun	adj	verb
$T =$	det	-4	3	2	-1
	noun	-3	-2	-1	2
	adj	-2	2	1	1
	verb	1	-1	0	0

	the	old	man	the	boat	
y_a	det	adj	noun	det	noun	score(y_a) = 25
y_b	det	noun	verb	det	noun	score(y_b) = 26
y_c	noun	noun	noun	noun	noun	score(y_c) = 1

Applying the Forward algorithm yields

		det	noun	adj	verb
$Q =$	the	5.00	0.00	0.00	0.00
	old	1.73	9.00	10.00	4.19
	man	8.18	15.01	11.05	12.70
	the	18.88	13.92	14.37	17.03
	boat	18.08	26.88	20.90	18.38

$$\log Z \approx 26.885$$

$$\log P(y_a) = \text{score}(y_a) - \log Z = 25 - 26.885 = -1.885$$

$$\log P(y_b) = \text{score}(y_b) - \log Z = 26 - 26.885 = -0.885$$

unary and transition scores:

		det	noun	adj	verb
$A =$	the	5	0	0	0
	old	0	1	3	0
	man	0	3	0	1
	the	5	0	0	0
	boat	0	5	0	0
		det	noun	adj	verb
$T =$	det	-4	3	2	-1
	noun	-3	-2	-1	2
	adj	-2	2	1	1
	verb	1	-1	0	0

	the	old	man	the	boat	
y_a	det	adj	noun	det	noun	score(y_a) = 25
y_b	det	noun	verb	det	noun	score(y_b) = 26
y_c	noun	noun	noun	noun	noun	score(y_c) = 1

unary and transition scores:

		det	noun	adj	verb
$A =$	the	5	0	0	0
	old	0	1	3	0
	man	0	3	0	1
	the	5	0	0	0
	boat	0	5	0	0

Applying the Forward algorithm yields

		det	noun	adj	verb
$Q =$	the	5.00	0.00	0.00	0.00
	old	1.73	9.00	10.00	4.19
	man	8.18	15.01	11.05	12.70
	the	18.88	13.92	14.37	17.03
	boat	18.08	26.88	20.90	18.38

		det	noun	adj	verb
$T =$	det	-4	3	2	-1
	noun	-3	-2	-1	2
	adj	-2	2	1	1
	verb	1	-1	0	0

$$\log Z \approx 26.885$$

$$\log P(y_a) = \text{score}(y_a) - \log Z = 25 - 26.885 = -1.885$$

$$\log P(y_b) = \text{score}(y_b) - \log Z = 26 - 26.885 = -0.885$$

$$\log P(y_c) = \text{score}(y_c) - \log Z = 1 - 26.885 = -25.885$$

Putting It All Together

At this point, we have all the ingredients needed to train a probabilistic sequence tagger with transition scores!

1. Receiving an input sequence \mathbf{x} , the model returns unary and transition scores \mathbf{A} and \mathbf{T} .
2. If we're at test time:
run Viterbi to get predicted sequence; compute accuracies etc.
3. If training time:
run Forward algorithm to compute the training objective
 $-\log P(\mathbf{y} | \mathbf{x}) = -\text{score}(\mathbf{y}) + \log \sum_{\mathbf{y}' \in \mathcal{Y}} \exp \text{score}(\mathbf{y}')$.

This probabilistic model is often known as a Linear-Chain Conditional Random Field.

(Historically, Linear-Chain CRFs didn't use neural net scorers, but the math doesn't change. Today I prefer to teach it this way.)