Lecture 9

Sequence Tagging

Part 1: Sequence Tagging

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

Outline:

Sequence Tagging Definition and examples

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

Evaluation

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

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

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

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

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

	the	old	man	the	boat
y _a	det	adj	noun	det	noun
y _b	det	noun	verb	det	noun

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

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

Example 2: Frame-level phone(me) classification (may be part of speech recognition)



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

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

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:

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

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

Output $y = (y_1, ..., y_n)$, e.g., a sequence of part-of-speech tags.

For each data point (sentence), |y| = |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:

 $\boldsymbol{y} = (y_1, \dots, y_n)$ 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), |y| = |x|; different data points have different lengths.

For fixed length *n*, some possible outputs:

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

How many in terms of *n*?

Part-Of-Speech Tags

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

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).

Lecture 9

Sequence Tagging

Part 2: Different Scoring Models

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

Outline:

Sequence Tagging Definition and example

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

Evaluation

Writing $\mathbf{y} = (y_1, \dots, y_n)$, take			det	noun	adi	verb
$\operatorname{score}(\boldsymbol{y}) = \sum_{j} a_{j,y_j}.$		the	5	0	0	0
A is a matrix of scores,		old	0	1	3	0
e.g., computed by a NN encoder.	A =	man	0	3	0	1
		the	5	0	0	0

boat 0 5 0 0

Writi score	ng y = (v) =	= (y ₁ , Σ.; a; y.,	., y _n), t	ake					det	noun	adj	verb
00010		$\Delta_j \simeq_{j,y_j}$						the	5	0	0	0
A is a	matr	ix of sc	ores,	naada			Λ_	old	0	1	3	0
e.g., (compu	ited by	ainne	encode	er.	,	H -	man	0	3	0	1
								the	5	0	0	0
y _a	the det	old adj	man noun	the det	boat noun			boat	0	5	0	0
y _b	det	noun	verb	det	noun							

Writi score	ng y = (v) =	= (y ₁ , Σ.; a; y:.	., y _n), t	ake				det	noun	adj	verb
		- ري- ري- ري- ري-					the	5	0	0	0
A is a	matr	ix of sc	ores,			Λ_	old	0	1	3	0
e.g., (compu	ited by	a inin e	encode	er.	A =	man	0	3	0	1
							the	5	0	0	0
V a	the det	old adj	man noun	the det	boat noun		boat	0	5	0	0
y _a y _b	det	noun	verb	det	noun						

 $score(\boldsymbol{y}_{a}) =$

Writi score	ng y = (y) =	$=(y_1,\ldots)_{j}a_{j,y_j}.$., y _n), t	ake					
A is a e.g., o	a matr compu	ix of sc ited by	ores, a NN e	ncode	er.				
y _a	the old man the boat y_a det adj noun det noun								
y _b	det	noun	verb	det	noun				

 $score(\boldsymbol{y}_a) = 21$



Writi score	ng y = (v) =	= (y ₁ , Σ.; a; y:.	., y _n), t	ake				det	noun	adj	verb
		- ري- ري- ري- ري-					the	5	0	0	0
A is a	matr	ix of sc	ores,			Λ_	old	0	1	3	0
e.g., (compu	ited by	a inin e	encode	er.	A =	man	0	3	0	1
							the	5	0	0	0
V a	the det	old adj	man noun	the det	boat noun		boat	0	5	0	0
y _a y _b	det	noun	verb	det	noun						

 $score(\boldsymbol{y}_a) = 21$

 $score(y_b) =$

Writi score	writing $\mathbf{y} = (y_1, \dots, y_n)$, take score $(\mathbf{y}) = \sum_j a_{j,y_j}$.								
A is a matrix of scores, e.g., computed by a NN encoder.									
	the	old	man	the	boat				
y _a	det	adj	noun	det	noun				
y _b	det	noun	verb	det	noun				

score(\boldsymbol{y}_a) = 21 score(\boldsymbol{y}_b) = 17



A first attempt: separate classifier for each position.

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

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

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

$$\boldsymbol{a}_j = \boldsymbol{W}^\top \boldsymbol{z}_j + \boldsymbol{b}$$

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

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]
```

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

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

			det	noun	adj	verb
$\max_{\boldsymbol{y} \in \mathcal{Y}} \operatorname{score}(\boldsymbol{y})$		the	5	0	0	0
	Δ –	old	0	1	3	0
	A =	man	0	3	0	1
		the	5	0	0	0
		boat	0	5	0	0

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

			det	noun	adj	verb
max score(y) y∈y		the	5	0	0	0
$= \max_{y_n \in [K]} \operatorname{score}\left([y_1, \ldots, y_n]\right)$	4 -	old	0	1	3	0
$y_1 \in [\Lambda],, y_n \in [\Lambda]$	A =	man	0	3	0	1
		the	5	0	0	0
		boat	0	5	0	0

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

			det	noun	adj	verb
$\max_{\mathbf{y} \in \mathcal{Y}} \operatorname{score}(\mathbf{y})$		the	5	0	0	0
$= \max_{y_n \in [K]} \operatorname{score}\left([y_1, \dots, y_n]\right)$	Λ_	old	0	1	3	0
$y_1 \in [\Lambda], \dots, y_n \in [\Lambda]$	A –	man	0	3	0	1
$= \max_{y_1 \in [K], \dots, y_n \in [K]} \sum_{i} a_{j, y_i}$		the	5	0	0	0
5		boat	0	5	0	0

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

		det	noun	adj	verb
$\max_{\mathbf{y}\in\mathcal{Y}}\operatorname{score}(\mathbf{y})$	the	5	0	0	0
$= \max_{y_n \in [K]} \operatorname{score}\left([y_1, \dots, y_n]\right)$	old	0	1	3	0
$\sum_{n=1}^{n} \sum_{i=1}^{n} \sum_{j=1}^{n} \sum_{j=1}^{n} \sum_{j=1}^{n} \sum_{i=1}^{n} \sum_{j=1}^{n} \sum_{i=1}^{n} \sum_{j=1}^{n} \sum_{j$	man	0	3	0	1
$= \max_{y_1 \in [K], \dots, y_n \in [K]} \sum_{i} a_{j, y_j}$	the	5	0	0	0
$=\sum \max_{i=1}^{n} \max_{i=1}^{n} \sum_{j=1}^{n} \max_{i=1}^{n} \sum_{j=1}^{n} \sum_{i=1}^{n} \sum_{j=1}^{n} \sum_{j=1}^{n} \sum_{j=1}^{n} \sum_{i=1}^{n} \sum_{j=1}^{n} \sum_{j=1}^{n} \sum_{j=1}^{n} \sum_{i=1}^{n} \sum_{j=1}^{n} \sum_{i=1}^{n} \sum_{j=1}^{n} \sum_{j=1}^{n} \sum_{j=1}^{n} \sum_{i=1}^{n} \sum_{j=1}^{n} \sum_{i=1}^{n} \sum_{j=1}^{n} \sum$	boat	0	5	0	0
$= \sum_{j} y_{j} \in [K]$					

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

			det	noun	adj	verb
max score(𝔥) 𝖕∈𝒴		the	5	0	0	0
$= \max_{y_1 \in [K]} \operatorname{score}\left([y_1, \ldots, y_n]\right)$	Δ -	old	0	1	3	0
$\sum_{i=1}^{n}$	A –	man	0	3	0	1
$= \max_{y_1 \in [K], \dots, y_n \in [K]} \sum_j a_{j, y_j}$		the	5	0	0	0
$=\sum_{i}\max_{j}a_{j,y_i}$		boat	0	5	0	0
$\sum_{j} y_{j} \in [K]$						

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

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

$$\log \sum_{\boldsymbol{y} \in \boldsymbol{\mathcal{Y}}} \exp\left(\mathsf{score}(\boldsymbol{y})\right)$$

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



		det	noun	adj	verb
	the	5	0	0	0
Λ -	old	0	1	3	0
A –	man	0	3	0	1
	the	5	0	0	0
	boat	0	5	0	0

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



		det	noun	adj	verb
	the	5	0	0	0
Λ -	old	0	1	3	0
A –	man	0	3	0	1
	the	5	0	0	0
	boat	0	5	0	0

With our score(y) = $\sum_{j} a_{j,y_i}$, can we compute:



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

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

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

		det	noun	adj	verb
	the	5	0	0	0
Λ -	old	0	1	3	0
A –	man	0	3	0	1
	the	5	0	0	0
	boat	0	5	0	0

With our score(y) = $\sum_{j} a_{j,y_i}$, can we compute:



 $\mathbf{A} = \begin{array}{ccccccc} \text{det noun adj verb} \\ \text{the 5 0 0 0} \\ \text{old 0 1 3 0} \\ \text{man 0 3 0 1} \\ \text{the 5 0 0 0} \\ \text{boat 0 5 0 0} \end{array}$

Probabilistic interpretation: independence

$$\log \Pr(\mathbf{y}) = \operatorname{score}(\mathbf{y}) - \log \sum_{\mathbf{y}' \in \mathcal{Y}} \exp \operatorname{score}(\mathbf{y}')$$
$$= \sum_{j} \underbrace{\left(a_{j,y_{j}} - \log \sum_{k \in [K]} \exp a_{j,k}\right)}_{\log \Pr(y_{j})}$$

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

$$score(\boldsymbol{y}) = \sum_{j} a_{j,y_j}$$

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(y) taking into account the sequential structure of y?

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

у	$score(\mathbf{y})$
---	---------------------

-1000	det det det det det
-940	det det det det noun
-800	det det det det verb
400	det noun verb det noun
-1100	verb verb verb verb verb

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

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

у	score	(y)
---	-------	--------------

-1000	det det det det det
-940	det det det det noun
-800	det det det det verb
400	det noun verb det noun
-1100	verb verb verb verb verb

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

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

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

у	score	(y)
---	-------	--------------

-1000	det det det det det
-940	det det det det noun
-800	det det det det verb
400	det noun verb det noun
-1100	verb verb verb verb verb

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.

Scoring Transitions Between Tags

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

scoring transitions between adjacent tags

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

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

Sequence Modeling With Transition Scores

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

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

- could be a learned parameter. (size does not depend on *n*)
- could be predicted by the neural net as a function of *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:

Sequence Tagging

Definition and examples

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

Evaluation

Sequence Tagging As A DAG



G = (V, E, w) where: $V = \{(i, c): i \in [n], c \in [K]\}$ \cup {*s*, *t*} $E = \{(i - 1, c') \rightarrow (i, c) : i \in [2, n], c, c' \in [K]\}$ $\cup \{ s \rightarrow (1, c) \colon c \in [K] \}$ $\cup \{(n, c) \rightarrow t : c \in [K]\}$ $w\left((j-1,c')\to(j,c)\right)=a_{j,c}+t_{c',c}$ $w(s \rightarrow (1, c)) = a_{1,c}$ $w((n, c) \rightarrow t) = 0$

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

Topological ordering?

Viterbi For Sequence Tagging



General Viterbi (reminder sketch)

initialize $m_1 \leftarrow 0$ for i = 2, ..., 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 $A(n \times K \text{ array})$ Transition scores $T(K \times K \text{ array})$

Forward: compute scores recursively $m_{1c} = a_{1c}$ 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 $y^* = [y_1, ..., y_n]$

19/∞

 $m_{j,c}$ 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
Λ -	old	0	1	3	0
A –	man	0	3	0	1
	the	5	0	0	0
	boat	0	5	0	0
		det	noun	adj	verb
	det	-4	3	2	-1
T =	noun	-3	-2	$^{-1}$	2
	adj	-2	2	1	1
	verb	1	$^{-1}$	0	0

 $m_{j,c}$ 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
4	old	0	1	3	0
A –	man	0	3	0	1
	the	5	0	0	0
	boat	0	5	0	0
		det	noun	adj	verb
	det	-4	3	2	-1
T =	noun	-3	-2	-1	2
	adj	-2	2	1	1
	verb	1	$^{-1}$	0	0

 $m_{j,c}$ 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
4	old	0	1	3	0
A –	man	0	3	0	1
	the	5	0	0	0
	boat	0	5	0	0
		det	noun	adj	verb
	det	-4	3	2	-1
T =	noun	-3	-2	$^{-1}$	2
	adj	-2	2	1	1
	verb	1	$^{-1}$	0	0

 $m_{j,c}$ 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.

 $M = \begin{array}{cccc} \det & \operatorname{noun} & \operatorname{adj} & \operatorname{verb} \\ \operatorname{the} & 5 & 0 & 0 & 0 \\ \operatorname{old} & 1 & & \\ \operatorname{man} & & \\ \operatorname{the} & & \\ \operatorname{boat} & & \end{array}$

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

 $m_{j,c}$ 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	0	1	3	0
A –	man	0	3	0	1
	the	5	0	0	0
	boat	0	5	0	0
		det	noun	adj	verb
	det	-4	3	2	-1
T =	noun	-3	-2	$^{-1}$	2
	adj	-2	2	1	1
	verb	1	$^{-1}$	0	0

 $m_{j,c}$ 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.

 $M = \begin{array}{cccc} \det & \operatorname{noun} & \operatorname{adj} & \operatorname{verb} \\ the & 5 & 0 & 0 & 0 \\ old & 1 & 9 & \\ man & & \\ the & \\ boat & \end{array}$

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

 $m_{j,c}$ 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.

 $M = \begin{array}{cccc} \det & \operatorname{noun} & \operatorname{adj} & \operatorname{verb} \\ \operatorname{the} & 5 & 0 & 0 & 0 \\ \operatorname{old} & 1 & 9 & 10 & 4 \\ \operatorname{man} & & & \\ \operatorname{the} & & \\ \operatorname{boat} & & & \end{array}$

		det	noun	adj	verb
	the	5	0	0	0
Δ	old	0	1	3	0
A –	man	0	3	0	1
	the	5	0	0	0
	boat	0	5	0	0
		det	noun	adj	verb
	det	-4	3	2	$^{-1}$
T =	noun	-3	-2	-1	2
	adj	-2	2	1	1
	verb	1	$^{-1}$	0	0

 $m_{j,c}$ 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
M —	old	1	9	10	4
111 -	man	8	15	11	12
	the	18	13	14	17
	boat	18	26	20	17

		det	noun	adj	verb
	the	5	0	0	0
Λ_	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
Γ =	noun	-3	-2	-1	2
	adj	-2	2	1	1
	verb	1	$^{-1}$	0	0

$m_{j,c}$	is store	d as a	a matr	ix M	, same shape as A .	una	ary and	trans	ition sc	ores:
App	ly <i>m</i> _{1,c} =	= a _{1,c}	to ge	t the	first row: (copied from A))		det	noun	adj
Ther	Then iteratively: $m_{i,c} = \max_{a' \in [K]} m_{i,1,a'} + a_{i,c} + t_{a',c}$					the	5	0	0	
			J,c			Λ -	old	0	1	3
At th	he end,	take 1	the ma	axim	um over the last row.	A –	man	0	3	0
		det	noun	adi	verb		the	5	0	0
		E	noun	auj	0		boat	0	5	0
	the	5	0	0	0					
NA _	old	1	9	10	4			det	noun	adi
101 -	man	8	15	11	12			4		
	the	10	12	1/	17		det	-4	3	2
	une	10	13	14	17	T =	noun	-3	-2	-1
	boat	18	26	20	17		adj	-2	2	1
							-			

To find the best tag sequence y^* , keep track of the path.

-1

-1

verb

verb

verb

-1

 $m_{i,c}$ 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: $m_{i,c} = \max_{c' \in [K]} m_{i-1,c'} + a_{i,c} + t_{c',c}$ At the end, take the maximum over the last row.

det adi verb noun 5 the 0 0 9 10 4 M =15 8 11 12 man 13 14 17

26

the

boat

18

18

To find the best tag sequence y^* , keep track of the path.

17

20

unary and transition scores:

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

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 A ($n \times K$ array) Transition scores 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'})
```

	boat	the	man	old	the	
$score(y_a) = 25$	noun	det	noun	adj	det	y a
$score(y_b) = 26$	noun	det	verb	noun	det	y _b
$score(y_c) = 1$	noun	noun	noun	noun	noun	\boldsymbol{y}_c

Applying the Forward algorithm yields

		det	noun	adj	verb
	the	5.00	0.00	0.00	0.00
0 -	old	1.73	9.00	10.00	4.19
Q –	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
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

V a	the det	old adj	man noun	the det	boat noun	$score(\mathbf{v}_2) = 25$	ur	hary an	d tran	sition s	cores:	
у _b	det	noun	verb	det	noun	$score(y_b) = 26$			det	noun	adj	verb
y _c	noun	noun	noun	noun	noun	$score(y_c) = 1$		the	5	0	0	0
Apply	ing the	Forward	d algoritl	hm yield	S		Λ —	old	0	1	3	0
	•		-				A =	man	0	3	0	1
		det	noun	adj	verb			the	5	0	0	0
	the	5.00	0.00	0.00	0.00			boat	0	5	0	0
0 -	old	1.73	9.00	10.00	4.19							
Q –	man	8.18	15.01	11.05	12.70				det	noun	adj	verb
	the	18.88	13.92	14.37	17.03			det	-4	3	2	$^{-1}$
	boat	18.08	26.88	20.90	18.38		T =	noun	-3	-2	$^{-1}$	2
								adj	-2	2	1	1
					log	$Z \approx 26.885$		verb	1	$^{-1}$	0	0

V a	the det	old adi	man noun	the det	boat noun	$score(v_{a}) = 25$	ur	hary an	d tran	sition s	cores:	
у _b	det	noun	verb	det	noun	$score(y_b) = 26$			det	noun	adj	verb
y _c	noun	noun	noun	noun	noun	$score(y_c) = 1$		the	5	0	0	0
Apply	ing the	Forward	d algoritl	nm yield	s		Δ_	old	0	1	3	0
	-						A -	man	0	3	0	1
		det	noun	adj	verb			the	5	0	0	0
	the	5.00	0.00	0.00	0.00	I		boat	0	5	0	0
0 -	old	1.73	9.00	10.00	4.19							
Q –	man	8.18	15.01	11.05	12.70	1			det	noun	adj	verb
	the	18.88	13.92	14.37	17.03			det	-4	3	2	$^{-1}$
	boat	18.08	26.88	20.90	18.38		T =	noun	-3	-2	-1	2
								adj	-2	2	1	1
					log	$Z \approx 26.885$		verb	1	-1	0	0
lo	$\log P(y_a) = \text{score}(y_a) - \log Z = 25 - 26.885 = -1.885$											

y _a	the det	old adj	man noun	the det	boat noun	$score(y_a) = 25$	ur	nary an	d tran	sition s	cores:	
y _b	det	noun	verb	det	noun	$score(y_b) = 26$			det	noun	adj	verb
y _c	noun	noun	noun	noun	noun	$score(y_c) = 1$		the	5	0	0	0
Apply	ing the	Forward	d algoritl	nm yield	S		Λ -	old	0	1	3	0
							A –	man	0	3	0	1
		det	noun	adj	verb			the	5	0	0	0
	the	5.00	0.00	0.00	0.00			boat	0	5	0	0
0 -	old	1.73	9.00	10.00	4.19							
Q –	man	8.18	15.01	11.05	12.70				det	noun	adj	verb
	the	18.88	13.92	14.37	17.03			det	-4	3	2	$^{-1}$
	boat	18.08	26.88	20.90	18.38		T =	noun	-3	-2	$^{-1}$	2
								adj	-2	2	1	1
					log	$Z \approx 26.885$		verb	1	-1	0	0
lo	$\log P(y_a) = \text{score}(y_a) - \log Z = 25 - 26.885 = -1.885$											
lo	$P(\mathbf{y}_b)$	= score	$(\boldsymbol{y}_b) - \boldsymbol{lc}$	$\log Z = 26$	5 – 26.8	85 = -0.885						

23/∞

	the	old	man	the	boat	unary and transition scores:						
Уa	det	adj	noun	det	noun	$score(y_a) = 25$						
y _b	det	noun	verb	det	noun	$score(y_b) = 26$			det	noun	adj	verb
у _с	noun	noun	noun	noun	noun	$score(y_c) = 1$		the	5	0	0	0
Apply	ing the	Forward	d algorith	nm yield	S		Λ –	old	0	1	3	0
							A -	man	0	3	0	1
		det	noun	adj	verb			the	5	0	0	0
	the	5.00	0.00	0.00	0.00			boat	0	5	0	0
0 -	old	1.73	9.00	10.00	4.19							
Q –	man	8.18	15.01	11.05	12.70				det	noun	adj	verb
	the	18.88	13.92	14.37	17.03			det	-4	3	2	$^{-1}$
	boat	18.08	26.88	20.90	18.38		T =	noun	-3	-2	$^{-1}$	2
								adj	-2	2	1	1
					log	$Z \approx 26.885$		verb	1	-1	0	C
lo	$g P(y_a)$) = score	$(\boldsymbol{y}_a) - \boldsymbol{lc}$	$\log Z = 2!$	5 – 26.8	85 = -1.885						
lo	$g P(y_b)$) = score	$(\boldsymbol{y}_b) - lc$	$\log Z = 26$	6 – 26.8	85 = -0.885						
le	$\log P(y_c) = \text{score}(y_c) - \log Z = 1 - 26.885 = -25.885$											

Putting It All Together

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

- Receiving an input sequence x, the model returns unary and transition scores A and 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

 $-\log P(\boldsymbol{y} \mid \boldsymbol{x}) = -\operatorname{score}(\boldsymbol{y}) + \log \sum_{\boldsymbol{y}' \in \mathcal{Y}} \exp \operatorname{score}(\boldsymbol{y}').$

Evaluation

Well, what would we do in the unstructured case?

Notation: Iverson Bracket

$$\llbracket p \rrbracket = \begin{cases} 1, & p \text{ is true,} \\ 0, & \text{otherwise.} \end{cases}$$

• Accuracy:

What fraction of test cases are correctly classified?

$$\mathsf{Acc} = \frac{1}{N} \sum_{i=1}^{N} \llbracket y^{(i)} = \widehat{y}^{(i)} \rrbracket$$

Structured evaluation: POS tagging

For sequential data, accuracy already becomes more complicated: sequence-level?

$$Acc_{seq} = \frac{\sum_{i=1}^{N} \left[\mathbf{y}^{(i)} = \hat{\mathbf{y}}^{(i)} \right]}{N}$$

or (micro-averaged) tag accuracy? (writing $n^{(i)} = |\mathbf{y}^{(i)}|$):
$$Acc_{tag} = \frac{\sum_{i=1}^{N} \sum_{j=1}^{n^{(i)}} \left[\left[\mathbf{y}_{j}^{(i)} = \hat{\mathbf{y}}_{j}^{(i)} \right] \right]}{\sum_{i=1}^{N} n^{(i)}}$$

(could also imagine a macro-averaged version, but it's not meaningful here)

Example:

true:	PRO	VERB	NUM	NOUN	ADV	
pred:	PRO	VERB	NUM	NOUN	PRO	
words:	there	are	70	children	there	
true: pred: words:	INTJ X eeeeek					

$$Acc_{seq} = \frac{0}{2} = 0$$
$$Acc_{tag} = \frac{4}{6} = 0.667$$

Historical Notes And References

- This probabilistic model is often known as a Linear-Chain Conditional Random Field and due to Lafferty et al. (2001). (Historically, Linear-Chain CRFs didn't use neural net scorers, but the math doesn't change. I prefer to teach it in a more general way.)
- On POS tagging: (Jurafsky and Martin, 2024) Speech and Language Processing [link]
- Structured prediction: (Daumé III, 2012), A Course In Machine Learning, ch. 17 [link]

References I

- Daumé III, Hal (2012). A Course in Machine Learning.
- Jurafsky, Daniel and James H. Martin (2024). Speech and Language Processing: An Introduction to Natural Language Processing, Computational Linguistics, and Speech Recognition with Language Models. 3rd. Online manuscript released August 20, 2024.
- Lafferty, John, Andrew McCallum, Fernando Pereira, et al. (2001). "Conditional random fields: Probabilistic models for segmenting and labeling sequence data". In: *Icml*. Vol. 1. 2. Williamstown, MA, p. 3.