

<https://eduassistpro.github.io/>

Assignment Project Exam Help

Assignment Project Add WeChat edu_assist_pro

<https://eduassistpro.github.io/>

Add WeChat edu_assist_pro

Sequence labeling problems

- ▶ Many problems in NLP can be formulated as sequence labeling problems
 - ▶ POS tagging:
 - ▶ The_DT man_NN who_WP whi _ VBZ
plans_NN to_VB get_VB edu_assist_pro
 - ▶ Named
 - ▶ T _ _ _ _ Microsoft_B-ORG
co _ _ _ _ _O venture_O
capitalist_O Andre _B- _ R
 - ▶ Time expression detection
 - ▶ Bedford_O police_O said_O they_O received_O a_O call_O
about_O 3:45_B-TIMEX p.m._I-TIMEX Monday_B-TIMEX
 - ▶ Spoken language understanding
 - ▶ Which_O flights_FLIGHT arrive_ARRIVE in_O Burbank_CITY
from_O Denver_CITY on_ON Saturday_Day
 - ▶

Search and Learning

Recall most natural language processing problems can be modeled mathematically as optimization:

$$\hat{y} = \underset{y \in \mathcal{Y}(x)}{\operatorname{argmax}} \Psi(x, y)$$

There are two modules:

- ▶ Search, the module that is responsible for finding the optimal y for a given x and parameters θ
- ▶ Learning, the module that is responsible for finding optimal parameters θ

For simple text classification problems, the search module is fairly straightforward, and most of the work goes to learning. For sequence labeling and more complicated NLP problems, the search module is getting more complicated.

Sequence labeling: first idea

<https://eduassistpro.github.io/>

Assignment Project Exam Help

Assignment Project Add WeChat edu_assist_pro

► Classify the se

<https://eduassistpro.github.io/>

Add WeChat edu_assist_pro

Sequence labeling example: POS tagging

<https://eduassistpro.github.io/>

- ▶ Let's use POS tagging as an example
- ▶ The most common used data set for training P the Penn TreeBank

DT					NNS
The					pianos

<https://eduassistpro.github.io/>

- ▶ DT: Determiner
- ▶ NN: uncountable noun or noun in singular
- ▶ WP: Wh-pronoun
- ▶ VBZ: 3rd person singular verb
- ▶ NNS: plural noun

How do we extract features from sequences in a linear model?

<https://eduassistpro.github.io/>

Assignment Project Exam Help

- ▶ We take as input a sequence of word tokens x , their corresponding POS tags y , as well as a position m , and return a set of features
- ▶ Typically we consider the word at position m and its surrounding words. We define a window w centered at position m of size k , and only extract contextual information from this window.

Extracting features from a window size of 1

<https://eduassistpro.github.io/>

- Assignment Project Exam Help
- ▶ Assuming a window of 1, the features we will be extracting from the example sentence will be:

Add WeChat edu_assist_pro

$$f((\mathbf{w} = \text{the man who whistles tunes pi}, m = 1), DT)$$

<https://eduassistpro.github.io/>

$$= (w$$
$$f((\mathbf{w} = \text{the man who whistles tunes pi}, m = 2), NN)$$

Add WeChat edu_assist_pro

$$= (w_0 = \text{man}, NN)$$

.....

How many features will we extract if we use a window of size 1?

Weights

<https://eduassistpro.github.io/>

We can then train a classifier using these features and get a weight for each feature:

	DT	ANN	MLP	MLP	MLP
$w_0 = \text{the}$	-0.05	-3.9	-	-	4.6
$w_0 = \text{ma}$	-	-	-	4	-3.5
$w_0 = \text{wh}$	-	-	-	6	-4.6
$w_0 = \text{whistles}$	-4.6	-4.6	-	-	-0.63
$w_0 = \text{tunes}$	-4.6	-4.6	-	-	-0.6
$w_0 = \text{pianos}$	-4.6	-4.6	-4.6	-3.0	-0.08

For example, the weight $\theta_1 = -0.05$ for the feature $f_1(w_0 = \text{the}, DT)$

Using these weights we can classify each word in the sequence

<https://eduassistpro.github.io/>

$$\begin{aligned} \psi((\mathbf{w} = \text{the man who whistles tunes pianos}, m = 1), DT) \\ &= \sum_i f_i \theta_i = -0.05 \\ \psi((\mathbf{w} = \text{the man who whistles tunes pianos}, m = 1), N) \\ &= \sum_i f_i \theta_i = -0.05 \\ \psi((\mathbf{w} = \text{tunes pianos}, m = 1), WP) \\ &= \sum_i f_i \theta_i = -4.6 \\ \psi((\mathbf{w} = \text{the man who whistles tunes pianos}, m = 1), VBZ) \\ &= \sum_i f_i \theta_i = -4.6 \\ \psi((\mathbf{w} = \text{the man who whistles tunes pianos}, m = 1), NNS) \\ &= \sum_i f_i \theta_i = -4.6 \end{aligned}$$

Predicting the tag for each word in the sequence

<https://eduassistpro.github.io/>

After finding the argmax_y at all positions of the sentence we get:

	DT	NN	WP	NNS		
	The	man	who	whistles		s
$w_0 = \text{the}$					6	-4.6
$w_0 = \text{man}$		-4.6	-0.35			-3.5
$w_0 = \text{who}$		-4.6	-4.6			-4.6
$w_0 = \text{whistles}$		-4.6	-4.6	-4.6	-0.8	-0.63
$w_0 = \text{tunes}$		-4.6	-4.6	-4.6	-0.8	-0.6
$w_0 = \text{pianos}$		-4.6	-4.6	-4.6	-3.0	-0.08

Extracting features from a larger window

<https://eduassistpro.github.io/>

- If we increase the window size to 2 and also include the previous word in the context

$$f((\mathbf{w} = \text{the man who whistles tunes piano}, m=1), VT) \\ = \{(w_{-1} = \text{man}, VT), (w_0 = \text{piano}, VT)\}$$

$$f((\mathbf{w} = \text{the man who whistles tunes pianos}, m=2), NN) \\ = \{(w_0 = \text{man}, NN), (w_{-1} = \text{piano}, NN), (w_1 = \text{pianos}, NN)\}$$

..... [Add WeChat edu_assist_pro](#)

$$f((\mathbf{w} = \text{the man who whistles tunes pianos}, m=4), VBZ) \\ = \{(w_0 = \text{whistles}, VBZ), (w_{-1} = \text{who}, VBZ), (w_1 = \text{pianos}, VBZ), (w_2 = \text{tunes}, VBZ)\}$$

Include weights for the new features

<https://eduassistpro.github.io/>

	DT	NN	WP	VBZ	NNS
$w_0=\text{the}$	-0.05	-3.9	-4.6	-4.6	-4.6
$w_0=\text{man}$	-4.6	-0.35	-	-	.5
$w_0=\text{who}$	-1.6	-4.6	-	-	.6
$w_0=\text{wh}$				8	-0.63
$w_0=\text{tu}$				8	-0.6
$w_0=\text{pia}$				0	-0.08
$w_{-1}=\text{START}$	-0.92	-3.9			-0.92
$w_{-1}=\text{the}$	-4.6	-0.7			-0.75
$w_{-1}=\text{man}$	-1.6	-2.3			-2.3
$w_{-1}=\text{who}$	-1.8	-4.6	-4.6	-0.2	-4.6
$w_{-1}=\text{whistles}$	-2.3	-4.6	-4.6	-1.6	-0.4
$w_{-1}=\text{tunes}$	-1.6	-4.6	-4.6	-4.6	-0.26

Classification with the new weights

$$\psi((\mathbf{w} = \text{th} \quad \text{nes pianos}), DT)$$

$$= \sum_i f_i \theta_i = -4.6 - 1.8 = -6.4$$

$$\psi((\mathbf{w} = \text{the man who whi} \quad \text{es pianos} \quad N)$$

$$= \sum_i f_i \theta_i = -4.6 - 4.6 = -9.2$$

$$\psi((\mathbf{w} = \text{tunes pianos} \quad = 4), WP)$$

$$= \sum_i f_i \theta_i = -4.6 - 4.6 = -9.2$$

$$\psi((\mathbf{w} = \text{the man who w} \quad \text{nes pianos}), VBZ)$$

$$= \sum_i f_i \theta_i = -0.8 - 0.2 = -1$$

$$\psi((\mathbf{w} = \text{the man who whistles tunes pianos}, m = 4), NNS)$$

$$= \sum_i f_i \theta_i = -0.63 + -4.6 = -5.23$$

So VBZ receives the highest score when classifying position 4.

Updated classification results

<https://eduassistpro.github.io/>

DT	NN	WP	VBZ	NNS	NNS
The	man	who	whistles	tunes	pianos

Assignment Project Exam Help

	DT	NN			NS
w_0 =the	-4.6	-3.9	-4.6	-1.4	-3.5
w_0 =man	-4.6	-0.35	-4.6	6	-4.6
w_0 =wh				8	0.63
w_0 =wh					
w_0 =tunes	-4.6	-4.6			-0.6
w_0 =pianos	-4.6	-4.6			-0.08
w_{-1} =START	-0.92	-3.9			-0.92
w_{-1} =the	-4.6	-0.7	-4.6	-4.6	-0.75
w_{-1} =man	-1.6	-2.3	-0.9	-1.6	-2.3
w_{-1} =who	-1.8	-4.6	-4.6	-0.2	-4.6
w_{-1} =whistles	-2.3	-4.6	-4.6	-1.6	-0.4
w_{-1} =tunes	-1.6	-4.6	-4.6	-4.6	-0.26

<https://eduassistpro.github.io/>

Add WeChat edu_assist_pro

Sequence labeling as structured prediction

<https://eduassistpro.github.io/>

Assignment Project Exam Help

- ▶ Enlarging the window to include more context helps, to a degree

- ▶ To further improve the classifier requires also sequences of t

NNS" should <https://eduassistpro.github.io/>

Incorporating such information in the model improve tagging accuracy

- ▶ The tags are of course not observable in the data, and they need to be predicted together.

Sequence labeling: Computing a global score for the entire sequence

<https://eduassistpro.github.io/>

- Assignment Project Exam Help
- ▶ Consider all possible label sequences for the input sequence, and choose the one that has the highest score

Add WeChat edu_assist_pro

$\Psi(\mathbf{w}, (DT, NN, WP, VB$

$\Psi(\mathbf{w}, (DT, NN, WP, VB, NS)) =$

<https://eduassistpro.github.io/>

- Add WeChat edu_assist_pro
- ▶ For a sequence of M elements with N possible labels, there are N^M possible sequences, a very large number!
 - ▶ To find the sequence with the highest score, we need to do this efficiently
 - ▶ The common solution is the Viterbi Algorithm

Sequence labeling as structured prediction

- ▶ The goal of the model is to find the sequence \mathbf{y} that has the highest score for a given input \mathbf{w} .

$$\hat{\mathbf{y}} = \underset{\mathbf{y} \in \mathcal{Y}(\mathbf{w})}{\operatorname{argmax}} \Psi(\mathbf{w}, \mathbf{y})$$

The score $\Psi(\mathbf{w}, \mathbf{y})$ is the sum of the local scores $\psi(\mathbf{w}_m, y_m, y_{m-1}, m)$ over the entire sequence \mathbf{y} .

- ▶ To make the computation tractable, we assume that the local score $\psi(\mathbf{w}_m, y_m, y_{m-1}, m)$ depends only on the current input \mathbf{w}_m and the current and previous labels y_m and y_{m-1} .

$$\Psi(\mathbf{w}, \mathbf{y}) = \sum_{m=1}^{M+1} \psi(\mathbf{w}_m, y_m, y_{m-1}, m)$$

- ▶ The local score is a weighted sum of the local features at position m .

$$\psi(\mathbf{w}_{1:M}, y_m, y_{m-1}, m) = \boldsymbol{\theta} \cdot \mathbf{f}(\mathbf{w}_m, y_m, y_{m-1}, m)$$

Feature representation for sequences

<https://eduassistpro.github.io/>

$f(\mathbf{w} = \text{the man who whistles tunes pianos}, \mathbf{y} = \text{DT NN WP VBZ VBZ NNS})$

$$= \sum_{m=1}^{M+1} f(\mathbf{w}, y_m, y_{m-1}, m)$$

$$= f(\mathbf{w}, \text{DT}, \diamond, 1) + f(\mathbf{w}, \text{NN}, \text{DT}, 2) + f(\mathbf{w}, \text{WP}, \text{NN}, 3) + f(\mathbf{w}, \text{VBZ}, \text{WP}, 4) + f(\mathbf{w}, \text{VBZ}, \text{VBZ}, 5) + f(\mathbf{w}, \text{NNS}, \text{VBZ}, 6)$$

$$= f(w_0 = \text{the}, y_0 = \text{DT}) + f(y_0 = \text{DT}, y_{-1} = \diamond)$$

$$+ f(w_0 = \text{man}, y_0 = \text{NN}) + f(y_0 = \text{NN}, y_{-1} = \text{DT})$$

$$+ f(w_0 = \text{who}, y_0 = \text{WP}) + f(y_0 = \text{WP}, y_{-1} = \text{NN})$$

$$+ f(w_0 = \text{whistles}, y_0 = \text{VBZ}) + f(y_0 = \text{VBZ}, y_{-1} = \text{WP})$$

$$+ f(w_0 = \text{tunes}, y_m = \text{VBZ}) + f(y_0 = \text{VBZ}, y_{-1} = \text{VBZ})$$

$$+ f(w_0 = \text{pianos}, y_0 = \text{NNS}) + f(y_0 = \text{NNS}, y_{-1} = \text{VBZ})$$

$$+ f(y_0 = \diamond, y_{-1} = \text{NNS})$$

Assignment Project Exam Help

Add WhatsApp <https://eduassistpro.github.io/>

<https://eduassistpro.github.io/>

Add WhatsApp <https://eduassistpro.github.io/>

Decoding for sequences: The Viterbi algorithm

<https://eduassistpro.github.io/>

- The goal is to find the sequence of tags with the highest score:

Assignment Project Exam Help

$\hat{y} = \underset{y \in \mathcal{Y}(w)}{\operatorname{argmax}} \Psi(w, y)$
Assignment Project edu_assist_pro

<https://eduassistpro.github.io/>

$= \underset{y_{1:M}}{\operatorname{argmax}} \sum_{m=1}^{M+1} s_{1:m}$
Add WeChat edu_assist_pro

- Instead of finding the argmax for the entire sequence directly, we start by finding the max up to position m and keep a sequence of back pointers

Finding the max score for the sequence

<https://eduassistpro.github.io/>

Assignment Project Exam Help

$\max_{y_{1:M}} \Psi(x, y_{1:M})$

$= \max_{y_{1:M}} \sum_{m=1}^M s_m(\diamond, y_m)$

$= \left(\max_{y_M} s_{M+1}(\diamond, y_M) \right) + \left(\sum_{m=1}^{M-1} s_m(\diamond, y_m) \right)$

Viterbi variable

<https://eduassistpro.github.io/>

Caching Viterbi variables as intermediate results:

Assignment Project Exam Help

Add WeChat edu_assist_pro

$$v_m(y_m) \triangleq \max_{\mathbf{y}_{1:m-1}} s_m(y_m, \mathbf{y}_{1:m-1})$$

<https://eduassistpro.github.io/>

Add WeChat edu_assist_pro

$$= \max_{y_{m-1}} s_m(y_m, y_{m-1})$$

Note that $v_1(y_1) \triangleq s_1(y_1, \diamond)$ and the maximum overall score for the sequence is the final Viterbi variable

$$\max_{\mathbf{y}_{1:M}} \Psi(\mathbf{w}_{1:M}, \mathbf{y}_{1:M}) = v_{M+1}(\diamond)$$

The Viterbi Algorithm

<https://eduassistpro.github.io/>

Viterbi Algorithm: Each $s_m(k, k')$ is a local score for tag $y_m = k$ and $y_{m-1} = k'$

```
1: for  $k \in \{0, \dots, K\}$  do
2:    $v_1(k) \leftarrow s_1(k, \diamond)$ 
3: for  $m \in \{2, \dots, M\}$  do
4:   for  $k \in \{0, \dots, K\}$  do
5:      $v_m(k) \leftarrow \max_{k'} s_m(k, k') + v_{m-1}(k')$ 
6:      $b_m(k) \leftarrow \operatorname{argmax}_{k'} s_m(k, k')$ 
7:  $y_M \leftarrow \operatorname{argmax}_k s_{M+1}(\diamond, k) + v_M(k)$ 
8: for  $m \in \{M-1, \dots, 1\}$  do
9:    $y_m \leftarrow b_m(y_{m+1})$ 
10: return  $y_{1:M}$ 
```

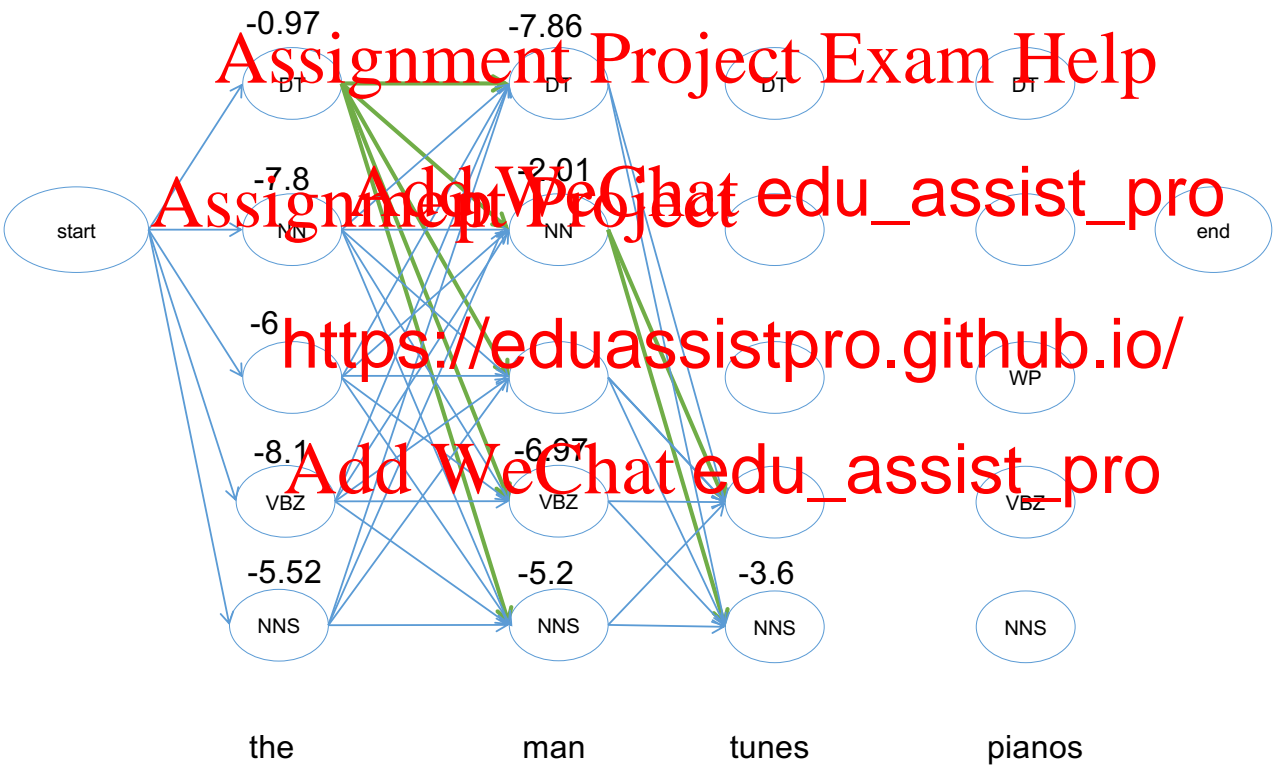
Assuming these parameters

<https://eduassistpro.github.io/>

	DT	NN	WP	VBZ	NNS	◆
$w_0 = \text{the}$	-0.05	-3.9	-4.6	-4.6	-4.6	$-\infty$
$w_0 = \text{man}$	-4.6	-0.35	-4.6			$-\infty$
$w_0 = \text{who}$	-4.6	-1.0	-0.05			$-\infty$
$w_0 = \text{whistles}$	-4.6	-4.6	-4.6			$-\infty$
$w_0 = \text{tunes}$					6	$-\infty$
$w_0 = \text{pianos}$					0.08	$-\infty$
$t_{-1} = \diamond$	-0.92	-3.9	-1.9			$-\infty$
$t_{-1} = \text{DT}$	-2.3	-0.69	-4.6			-4.6
$t_{-1} = \text{NN}$	-4.6	-1.6	-0.3	-0.36	-1.0	-0.7
$t_{-1} = \text{WP}$	-3.8	-4.6	-4.6	-0.2	-4.6	-4.6
$t_{-1} = \text{VBZ}$	-0.2	-1.3	-1.6	-4.6	-0.92	-2.3
$w_{-1} = \text{NNS}$	-4.6	-4.6	-0.1	-4.6	-3.9	-1.2

Example Viterbi computation

<https://eduassistpro.github.io/>



Additional features (and their weights) can be added

<https://eduassistpro.github.io/>

	DT	NN	WP	VBZ	NNS	◆
$w_0 = \text{the}$	-0.05	-3.9	-4.6	-4.6	-4.6	$-\infty$
$w_0 = \text{man}$	-4.6	-0.35	-4.6	-1.4	-3.5	$-\infty$
$w_0 = \text{who}$	-4.6	-4.6	-0.05	-4.6		
$w_0 = \text{whistles}$	-4.6	-4.6	4.6	0.8		
$w_0 = \text{tunes}$	-4.6	-4.6	-1.6	4.8		
$w_0 = \text{pianos}$	-4.6	-4.6	-4.6	-3.0		∞
$t_{-1} = \diamond$	-					∞
$t_{-1} = \text{DT}$	-					6
$t_{-1} = \text{NN}$	-					.7
$t_{-1} = \text{WP}$	-3.8	-4.6	-4.6	-0.		
$t_{-1} = \text{VBZ}$	-0.2	1.3	-1.6	4.		
$w_{-1} = \text{NNS}$	-4.6	-4.6	-0.1	-4.		
$w_{-1} = \text{START}$	-0.92	-3.9	-1.9	-3.5	-0.92	$-\infty$
$w_{-1} = \text{the}$	-4.6	-0.7	-4.6	-4.6	-0.75	-10
$w_{-1} = \text{man}$	-1.6	-2.3	-0.9	-1.6	-2.3	-1
$w_{-1} = \text{who}$	-1.8	-4.6	-4.6	-0.2	-4.6	-9
$w_{-1} = \text{whistles}$	-2.3	-4.6	-4.6	-1.6	-0.4	-0.5
$w_{-1} = \text{tunes}$	-1.6	-4.6	-4.6	-4.6	-0.26	-0.3

Feature templates used in SoA models

<https://eduassistpro.github.io/>

State-of-the-art models tend to use richer set of features and high-order transitions

- ▶ current word, w_t
- ▶ previous words, w_{t-1}, w_{t-2}
- ▶ next words, w_{t+1}, w_{t+2}
- ▶ previous two tags, y_{t-1}, y_{t-2}
- ▶ for rare words:
 - ▶ first k characters, up to $K = 4$
 - ▶ last k characters, up to $k = 4$
 - ▶ whether w_m contains a number, uppercase character, or hyphen

Parameter estimation for sequence labeling

<https://eduassistpro.github.io/>

Assignment Project Exam Help

We can extend the text classification models to sequ

Assignment Project Add WeChat edu_assist_pro

Text classificat	beling
Naïve Bayes	Models (HMM)
Logistic Regression	dom Fields (CRF)
Perceptron	Add WeChat edu_assist_pro
Support Vector Machines (SVM)	achines (SVM)