

<https://powcoder.com>

Assignment Project Exam Help

Assignment Project Exam Help

Sequence labeling problems

<https://powcoder.com>

Add WeChat powcoder

Sequence labeling problems

<https://powcoder.com>

- ▶ Many problems in NLP can be formulated as sequence labeling problems
 - ▶ POS tagging:
 - ▶ The_DT man_NN who_WP whistles_VBZ tunes_VBZ pianos_NN
 - ▶ Named Entity Recognition (NER)
 - ▶ The_O company_O is_O backed_O by_O Microsoft_B-ORG cofounder_O Bill_B-PER Gates_I-PER and_O venture_O capitalist_O Andressen_B-PER Horowitz_I-PER
 - ▶ 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 problems can be formulated mathematically as optimization:

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

There are two modules to this:

- ▶ Search, the module that is responsible for finding the argmax of the score function ψ
- ▶ Learning, the module that is responsible for finding the 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://powcoder.com>

Assignment Project Exam Help

Assignment Project Exam Help

- Classify the sequence one element at a time

<https://powcoder.com>

Add WeChat powcoder

Sequence labeling example: POS tagging

<https://powcoder.com>

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

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

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

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

<https://powcoder.com>

Assignment Project Exam Help

- ▶ We take as input a sequence of word tokens x and their corresponding POS tags y , as well as a position m , and return a set of features associated with that position.
- ▶ Typically we make the assumption that the context that matters for classifying the word at position m are its surrounding words. We define a window that is centered on position m of size k , and only extract contextual information from this window.

Extracting features from a window size of 1

<https://powcoder.com>

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

Assignment Project Exam Help

$$\begin{aligned} f((\mathbf{w} = \text{the man who whistles tunes pianos}, m = 1), DT) \\ &= (w_0 = \text{the}, DT) \\ f((\mathbf{w} = \text{the man who whistles tunes pianos}, m = 2), NN) \\ &= (w_0 = \text{man}, NN) \\ &\dots\dots \end{aligned}$$

<https://powcoder.com>

Add WeChat powcoder

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

Weights

<https://powcoder.com>

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

	DT	MN	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	-4.6	-1.4	-3.5
$w_0 = \text{who}$	-4.6	-4.6	-0.05	-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

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://powcoder.com>

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

$$\begin{aligned}\psi((\mathbf{w} = \text{the man who whistles tunes pianos}, m = 1), W/V) \\ = \sum_i f_i \theta_i = -3.9\end{aligned}$$

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

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

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

Predicting the tag for each word in the sequence

<https://powcoder.com>

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

Assignment Project Exam Help
Add WeChat powcoder

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

	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	-4.6	-1.4	-3.5
$w_0=\text{who}$	-4.6	-4.6	-0.05	-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://powcoder.com>

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

$$\begin{aligned} f((\mathbf{w} = \text{the man who whistles tunes pianos}, m = 1), DT) \\ = \{(w_0 = \text{the}, DT), (w_{-1} = \text{START}, DT)\} \end{aligned}$$

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

..... Add WeChat powcoder

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

Include weights for the new features

<https://powcoder.com>

	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	-4.6	-1.4	-3.5
$w_0=\text{who}$	-1.6	-4.6	-0.9	-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
$w_{-1}=\text{START}$	-0.92	-3.9	-1.9	-3.5	-0.92
$w_{-1}=\text{the}$	-4.6	-0.7	-4.6	-4.6	-0.75
$w_{-1}=\text{man}$	-1.6	-2.3	-0.9	-1.6	-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

$$\begin{aligned}\psi((\mathbf{w} = \text{the man who whistles tunes pianos}, m = 4), DT) \\&= \sum_i f_i \theta_i = -4.6 - 1.8 = -6.4 \\ \psi((\mathbf{w} = \text{the man who whistles tunes pianos}, m = 4), NN) \\&= \sum_i f_i \theta_i = -4.6 - 4.6 = -9.2 \\ \psi((\mathbf{w} = \text{the man who whistles tunes pianos}, m = 4), WP) \\&= \sum_i f_i \theta_i = -4.6 - 4.6 = -9.2 \\ \psi((\mathbf{w} = \text{the man who whistles tunes pianos}, m = 4), VBZ) \\&= \sum_i f_i \theta_i = -0.8 - 0.2 = -1 \\ \psi((\mathbf{w} = \text{the man who whistles tunes pianos}, m = 4), NNS) \\&= \sum_i f_i \theta_i = -0.63 + -4.6 = -5.23\end{aligned}$$

So VBZ receives the highest score when classifying position 4.

Updated classification results

<https://powcoder.com>

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

	DT	NN	WP	VBZ	NNS
w_0 =the	-0.05	-3.9	-4.6	-4.6	-4.6
w_0 =man	-4.6	-0.35	-4.6	-1.4	-3.5
w_0 =who	-4.6	-4.6	-0.05	-4.6	-4.6
w_0 =whistles	-4.6	-4.6	-4.6	-0.8	-0.63
w_0 =tunes	-4.6	-4.6	-4.6	-0.8	-0.6
w_0 =pianos	-4.6	-4.6	-4.6	-3.9	-0.08
w_{-1} =START	-0.92	-3.9	-1.9	-3.5	-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

Sequence labeling as structured prediction

<https://powcoder.com>

Assignment Project Exam Help

- ▶ Enlarging the window to include more context helps, to a degree
- ▶ To further improve the classifier requires also evaluating sequences of tags. For example, the tag sequence “NNS NNS” should receives a very low score as it rarely appears. Incorporating such information in the model would help 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://powcoder.com>

- Consider all possible label sequences for the input sequence, and choose the one that has the highest score

[Assignment Project Exam Help](#)

$$\Psi(\mathbf{w}, (DT, NN, WP, VBZ, NNS, NNS)) =$$

$$\Psi(\mathbf{w}, (DT, NN, WP, VBZ, VBZ, NNS)) =$$

<https://powcoder.com>

- For a sequence of M elements with a tagset of size N , 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

[Add WeChat powcoder](#)

Sequence labeling as structured prediction

- The goal of the model is to find the tag sequence that yields the highest score for the input sequence:

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

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

- To make the computation tractable we factor the score for the entire sequence into the sum of local scores at each position m

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

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

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

Feature representation for sequences

<https://powcoder.com>

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

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

$$= f(\mathbf{w}, \text{DT}, \diamond, 1) + f(\mathbf{w}, \text{NN}, \text{DT}, 2) + f(\mathbf{w}, \text{WP}, \text{NN}, 3) \\ + f(\mathbf{w}, \text{VBZ}, \text{WP}, 4) + f(\mathbf{w}, \text{VBZ}, \text{VBZ}, 5) + f(\mathbf{w}, \text{NNS}, \text{VBZ}, 6) \\ + f(\mathbf{w}, \diamond, \text{NNS}, 7)$$

$$= f(w_0 = \text{the}, y_0 = \text{DT}) + f(y_0 = \text{DT}, y_{-1} = \diamond) \\ + f(w_0 = \text{man}, y_0 = \text{NN}) + f(y_0 = \text{NN}, y_{-1} = \text{DT}) \\ + f(w_0 = \text{who}, y_0 = \text{WP}) + f(y_0 = \text{WP}, y_{-1} = \text{NN}) \\ + f(w_0 = \text{whistles}, y_0 = \text{VBZ}) + f(y_0 = \text{VBZ}, y_{-1} = \text{WP}) \\ + f(w_0 = \text{tunes}, y_m = \text{VBZ}) + f(y_0 = \text{VBZ}, y_{-1} = \text{VBZ}) \\ + f(w_0 = \text{pianos}, y_0 = \text{NNS}) + f(y_0 = \text{NNS}, y_{-1} = \text{VBZ}) \\ + f(y_0 = \diamond, y_{-1} = \text{NNS})$$

Decoding for sequences: The Viterbi algorithm

<https://powcoder.com>

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

Assignment Project Exam Help

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

$= \operatorname{argmax}_{y_{1:M}} \sum_{m=1}^{M+1} \psi(w, y_m, y_{m-1}, m)$

$= \operatorname{argmax}_{y_{1:M}} \sum_{m=1}^{M+1} s_m(y_m, y_{m-1})$

- 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://powcoder.com>

Assignment Project Exam Help

$$\begin{aligned} & \max_{\mathbf{y}_{1:M}} \Psi(\mathbf{x}, \mathbf{y}_{1:M}) \\ &= \max_{\mathbf{y}_{1:M}} \sum_{m=1}^{M+1} s_m(y_m, y_{m-1}) \\ &= \left(\max_{y_M} s_{M+1}(\diamond, y_M) \right) + \left(\max_{\mathbf{y}_{1:M-1}} \sum_{m=1}^M s_m(y_m, y_{m-1}) \right) \end{aligned}$$

Viterbi variable

<https://powcoder.com>

Caching Viterbi variables as intermediate results:

Assignment Project Exam Help

Assignment Project Exam Help

$$v_m(y_m) \triangleq \max_{\mathbf{y}_{1:m-1}} \sum_{n=1}^m s_n(y_n, y_{n-1})$$

<https://powcoder.com>

$$\begin{aligned} &= \max_{y_{m-1}} s_m(y_m, y_{m-1}) + \max_{\mathbf{y}_{1:m-2}} \sum_{n=1}^{m-1} s_n(y_n, y_{n-1}) \\ &= \max_{y_{m-1}} s_m(y_m, y_{m-1}) + v_{m-1}(y_{m-1}) \end{aligned}$$

Add WeChat powcoder

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://powcoder.com>

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:      $b_m(k) \leftarrow \operatorname{argmax}_{k'} s_m(k, k') + v_{m-1}(k')$ 
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  $y_{1:M}$ 
```

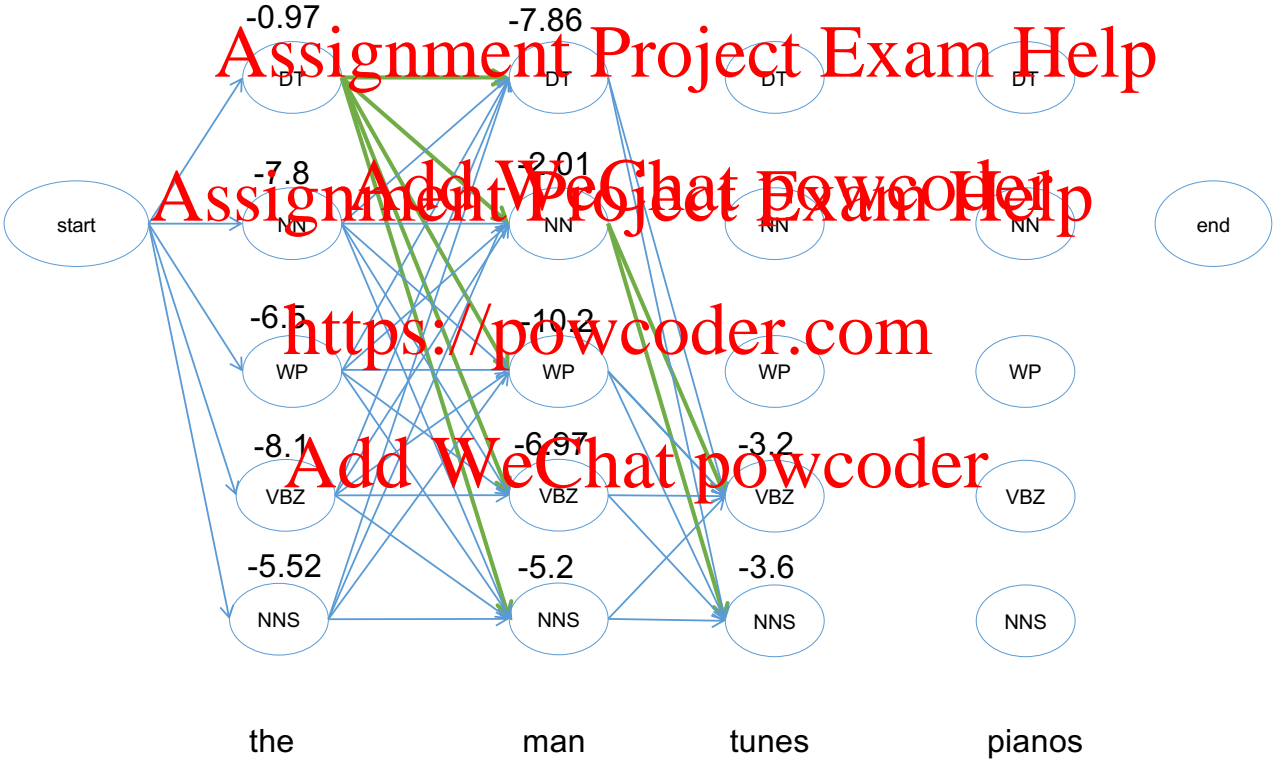
Assuming these parameters

<https://powcoder.com>

	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	-4.6	$-\infty$
$w_0 = \text{whistles}$	-4.6	-4.6	-4.6	-0.8	-0.63	$-\infty$
$w_0 = \text{tunes}$	-4.6	-4.6	-4.6	-0.8	-0.6	$-\infty$
$w_0 = \text{pianos}$	-4.6	-4.6	-4.6	-3.0	-0.08	$-\infty$
$t_{-1} = \diamond$	-0.92	-3.9	-1.9	-3.5	-0.92	$-\infty$
$t_{-1} = \text{DT}$	-2.3	-0.69	-4.6	-4.6	-0.76	-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://powcoder.com>



Additional features (and their weights) can be added

<https://powcoder.com>

	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	-4.6	$-\infty$
$w_0 = \text{whistles}$	-4.6	-4.6	4.6	0.8	-0.63	$-\infty$
$w_0 = \text{tunes}$	-4.6	-4.6	-1.6	4.8	0.6	$-\infty$
$w_0 = \text{pianos}$	-4.6	-4.6	-4.6	-3.0	-0.08	$-\infty$
$t_{-1} = \diamond$	-0.92	-3.9	-1.9	-3.5	-0.92	$-\infty$
$t_{-1} = \text{DT}$	-2.3	-0.69	-4.6	-4.6	-0.75	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
$w_{-1} = \text{START}$	-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://powcoder.com>

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

- ▶ current word, w_t
- ▶ previous words, w_{-1}, w_{-2}
- ▶ next words, w_1, w_2
- ▶ previous two tags, y_{-1}, y_{-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://powcoder.com>

Assignment Project Exam Help

We can extend the text classification models to sequence labeling:

Assignment Project Exam Help

Text classification	Sequence Labeling
Naïve Bayes	Hidden Markov Models (HMM)
Logistic Regression	Conditional Random Fields (CRF)
Perceptron	Perceptron
Support Vector Machines (SVM)	Support Vector Machines (SVM)