

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

Assignment Project Exam Help

Generativ Assignment Project Add WeChat edu_assist_pro models

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

Add WeChat edu_assist_pro

Hidden Markov Models (HMM): the generative story

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

Assignment Project Exam Help

The generative story: first, the tags are drawn from a prior distribution; next, the tokens are drawn from a conditional likelihood

[Add WeChat edu_assist_pro](#)

$y_0 \leftarrow \diamond, m \leftarrow 1$ <https://eduassistpro.github.io/>

repeat

$y_m \sim \text{Categorical}(\lambda_{y_{m-1}})$

$w_m \sim \text{Categorical}(\phi_{y_m})$

until $y_m = \blacklozenge$ \triangleright terminate when the stop symbol is generated

The independence assumptions of HMM

In addition to the usual in yes, two additional independence assumptions are needed for HMM:

- The probability of each word token only depends on its tag, not on any other element in the sequence:

- Each tag y_m depends only on its pred

$$P(\mathbf{y}) = \prod_{m=1}^M P(y_m | y_{m-1})$$

when $y_m = \diamond$ in all cases.

Parameter estimation of HMMs

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

The hidden Markov model has two groups of parameters:

- ▶ **Emission probabilities** ϕ . The probability $Pr(w_m | y_m; \phi)$ is the emission probability
- ▶ **Transition probabilities** λ . The probability $Pr(y_m | y_{m-1}; \lambda)$

Both of these groups

relative frequency of the probabilities are, <https://eduassistpro.github.io/>

Add WeChat edu_assist_pro

$$\phi_{k,i} \triangleq Pr(W_m = i | Y_m = k) = \frac{\text{count}(W_m = i, Y_m = k)}{\text{count}(Y_m = k)}$$

$$\lambda_{k,k'} \triangleq Pr(Y_m = k' | Y_{m-1} = k) = \frac{\text{count}(Y_m = k', Y_{m-1} = k)}{\text{count}(Y_{m-1} = k)}$$

Inference for HMMs

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

- The goal of the inference in the Hidden Markov model is to find the highest probability sequence:

[Add WeChat edu_assist_pro](https://eduassistpro.github.io/)

$$\hat{\mathbf{y}} = \operatorname{argmax}_{\mathbf{y}} \log$$

- As $P(\mathbf{y}, \mathbf{w})$ the inference problem can be reformulated as finding $\hat{\mathbf{y}}$ for \mathbf{w} and \mathbf{y} :

$$\hat{\mathbf{y}} = \operatorname{argmax}_{\mathbf{y}} \log P(\mathbf{y}, \mathbf{w})$$

Inference for HMMs

- Applying the independence assumption

$$\begin{aligned}\log P(\mathbf{y}, \mathbf{w}) &= \log P(\mathbf{y}) + \log P(\mathbf{w} | \mathbf{y}) \\ &= \sum_{m=1}^{M+1} \log P(y_m | y_{m-1}) + \sum_{m=1}^M \log P(w_m | y_m)\end{aligned}$$

$$= \sum_{m=1}^{M+1} s_m(y_{m-1}, y_m, w_m)$$

- This allows us to apply a variant of the Viterbi algorithm, where the only parameters are the transition probabilities and emission probabilities, which correspond to two features at each position in the sequence. This limitation explains the performance disadvantage of HMMs.

Discriminative alternatives to HMMs

As with Naïve Bayes, additional features (e.g., suffixes) cannot be applied without violating the independent assumptions. Discriminative models do not have this problem and additional features can be included:

$$\begin{aligned} f(\mathbf{w} = \text{the man who whistles tunes pianos}, \mathbf{y} = \text{DT NN WP VBZ VBZ NNS}) \\ = f(w_0 = \text{the}, y_0 = \text{DT}) \\ + f(w_0 = \text{man}, y_0 = \text{NN}) \\ + f(w_0 = \text{who}, y_0 = \text{WP}) + f(y_0 = \text{WP}, y_{-1} = \text{NN}) \\ + f(w_0 = \text{whistles}, y_0 = \text{VBZ}) + f(y_0 = \text{VBZ}, y_{-1} = \text{WP}) + f(\text{suffix}_0 = \text{es}) \\ + f(w_0 = \text{tunes}, y_0 = \text{VBZ}) + f(y_0 = \text{VBZ}, y_{-1} = \text{VBZ}) + f(\text{suffix}_0 = \text{es}) \\ + f(w_0 = \text{pianos}, y_0 = \text{NNS}) + f(y_0 = \text{NNS}, y_{-1} = \text{VBZ}) \\ + f(y_0 = \blacklozenge, y_{-1} = \text{NNS}) \end{aligned}$$

Note that you do not need to add the same morphological feature for each word token.

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

Assignment Project Exam Help

Assignment Project
Structur

Add WeChat edu_assist_pro
ling

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

Add WeChat edu_assist_pro

Structured Perceptron

- Each tagging sequence model

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

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

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

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

$$= \theta \cdot \mathbf{f}^{(global)}(\mathbf{w}, \mathbf{y}_{1:M})$$

where $y_{M+1} = \blacklozenge$ and $y_0 = \blacklozenge$ by construction.

- The best tagging sequence can be found efficiently with the Viterbi algorithm.
- As a discriminative model, Perceptron can handle an arbitrary number of features at each position.

Parameter estimation for structured perceptron

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

Assignment Project Exam Help

- ▶ In the training phase, the sentences in the training set are decoded one a time (with the Viterbi Algorithm)
- ▶ If the highest scoring tagging sequence is not the correct tag sequence

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

$\hat{y} = \underset{y \in \mathcal{Y}(w)}{\operatorname{argmax}} \theta \cdot f(w, y)$

$$\theta^{(t+1)} \leftarrow \theta^{(t)} + f(w, y) - f(w, \hat{y})$$

The averaged perceptron

```
1: procedure AVE_PERCEPTRON( $\mathbf{w}^{(0)}$ ,  $\mathbf{y}$ )
2:    $t \leftarrow 0$ ,  $\boldsymbol{\theta}^{(0)} \leftarrow \mathbf{0}$ ,  $m \leftarrow 0$ 
3:   repeat
4:      $t \leftarrow t + 1$ 
5:     Select a sequence  $i$  line training
6:      $\hat{\mathbf{y}} \leftarrow$  encoding by Viterbi
7:     if  $\hat{\mathbf{y}} \neq \mathbf{y}(i)$ 
8:        $\boldsymbol{\theta}^{(t)} \leftarrow \boldsymbol{\theta}^{(t-1)} + \mathbf{f}^{(global)}(\mathbf{w}^{(i)}, \hat{\mathbf{y}})$ 
9:     else  $\boldsymbol{\theta}^{(t)} \leftarrow \boldsymbol{\theta}^{(t-1)}$  a sequence
10:     $m \leftarrow m + \boldsymbol{\theta}^{(t)}$ 
11:  until tired
12:   $\bar{\boldsymbol{\theta}} = \frac{1}{t} m$ 
13:  return  $\bar{\boldsymbol{\theta}}$ 
```

Parameter Update example

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

- ▶ Correct tag sequence:
 - ▶ The_DT man_NN who_WP whistles_VBZ tunes_VBZ pianos_NNS
- ▶ Highest scoring sequence under the current
 - ▶ The_D - - - es_NNS pianos_
- ▶ Which features need to be updated?

$$\begin{aligned}\theta_{(tunes, VBZ)} &\leftarrow \theta_{(tunes, VBZ)} + 1, & \theta_{(tunes, NNS)} &\leftarrow \theta_{(tunes, NNS)} - 1 \\ \theta_{(VBZ, VBZ)} &\leftarrow \theta_{(VBZ, VBZ)} + 1 & \theta_{(VBZ, NNS)} &\leftarrow \theta_{(VBZ, NNS)} - 1 \\ \theta_{(VBZ, NNS)} &\leftarrow \theta_{(VBZ, NNS)} + 1 & \theta_{(NNS, NNS)} &\leftarrow \theta_{(NNS, NNS)} - 1\end{aligned}$$

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

Assignment Project Exam Help

Assignment Project Add WeChat edu_assist_pro
Stru es

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

Add WeChat edu_assist_pro

Structured Support Vector Machines

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

Assignment Project Exam Help

- Classification with SVMs enforces a large-margin constraint that requires that there is a margin of at least 1 bet score of the correct label and all incorrect label this means

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

- This can be minimally extended to seque

$$\forall \mathbf{y} \neq \mathbf{y}^{(i)}, \theta \cdot \mathbf{f}(\mathbf{w}, \mathbf{y}^{(i)}) - \theta \cdot \mathbf{f}(\mathbf{w}, \mathbf{y}) \geq 1$$

Extending SVMs to sequences

- ▶ A “structured” Support Vector Machine (SVM) outputs a structured object such as a sequence.
- ▶ When extending SVMs to sequence labeling, there are several issues to be addressed.
 - ▶ Some errors are more serious than others. They have the same effect on the overall score, but different errors would lead to different sequences. For example, having a fixed margin of 1 would suggest enumerating all possible sequences, which is infeasible as the number of sequences is exponential to the length of the sequence.
- ▶ The solution requires an adjustment of how the margin is computed.

Extending SVMs to sequences

- Instead of using a fixed margin γ , we use a margin that reflects how the sequence \mathbf{y} differs from the target sequence $\mathbf{y}^{(i)}$.

$$\forall \mathbf{y} \neq \mathbf{y}^{(i)} \quad \theta \cdot \mathbf{f}(\mathbf{w}^{(i)}, \mathbf{y}) - \theta \cdot \mathbf{f}(\mathbf{w}^{(i)}, \mathbf{y}^{(i)}) \geq c(\mathbf{y}, \mathbf{y}^{(i)})$$

- Next, instead of using a delta function $\delta(\mathbf{y}, \mathbf{y}^{(i)})$ to compute the error for the entire sequence, we use the **Hamming cost**, which counts the number of errors in \mathbf{y} :

$$c(\mathbf{y}, \mathbf{y}^{(i)}) = \sum_{t=1}^m \mathbb{1}_{y_t \neq y_t^{(i)}}$$

- Instead of training against all labelings \mathbf{y} , we have a margin γ that satisfies the above constraint, we focus on the prediction that **maximally** violates the margin constraint. We can identify this prediction by solving:

$$\begin{aligned} \hat{\mathbf{y}} &= \operatorname{argmax}_{\mathbf{y} \neq \mathbf{y}^{(i)}} \theta \cdot \mathbf{f}(\mathbf{w}^{(i)}, \mathbf{y}) - \theta \cdot \mathbf{f}(\mathbf{w}^{(i)}, \mathbf{y}^{(i)}) + c(\mathbf{y}, \mathbf{y}^{(i)}) \\ &= \operatorname{argmax}_{\mathbf{y} \neq \mathbf{y}^{(i)}} \theta \cdot \mathbf{f}(\mathbf{w}^{(i)}, \mathbf{y}) + c(\mathbf{y}, \mathbf{y}^{(i)}) \end{aligned}$$

Extending SVMs to sequence labeling

- ▶ Reformulating the constraint as $\theta \cdot f(\mathbf{w}^{(i)}, \mathbf{y}^{(i)}) - \max_{\mathbf{y} \in \mathcal{Y}(\mathbf{w})} (\theta \cdot f(\mathbf{w}^{(i)}, \mathbf{y}) + c(\mathbf{y}, \mathbf{y}^{(i)})) \geq 0$

- ▶ This means that the constraint will be met (and be complete) if $\theta \cdot f(\mathbf{w}^{(i)}, \mathbf{y}^{(i)})$ is equal or greater than $\max_{\mathbf{y} \in \mathcal{Y}(\mathbf{w})} (\theta \cdot f(\mathbf{w}^{(i)}, \mathbf{y}) + c(\mathbf{y}, \mathbf{y}^{(i)}))$ alternatives.
- ▶ In the training process, we identify predictions (scores highly according to the model) and large cost according to the truth), and reduce these predictions by adjusting weights.
- ▶ Note that Hamming cost can be reduced to local parts by adding a feature $f_m(y_m) = \delta(y_m \neq y_m^{(i)})$, and can be incorporated into the Viterbi algorithm for purposes of the identifying the prediction to train against.

A comparison between Structured Perceptron and
Structured SVM <https://eduassistpro.github.io/>

Assignment Project Exam Help

- ▶ In the training process, the perceptron algorithm finds the prediction that has the highest score across all models to train against, while Structured SVM finds the prediction that has the lowest score across all models. <https://eduassistpro.github.io/>
- ▶ No cost needs to (or can be) computed during training for both models. [Add WeChat edu_assist_pro](https://eduassistpro.github.io/)
- ▶ In practice, with a large training set, the perceptron algorithm works pretty well, obviating the need for more complicated SVM models.