

<https://powcoder.com>

Assignment Project Exam Help

Assignment Project Exam Help
Generative approaches: Hidden Markov Models

<https://powcoder.com>

Add WeChat powcoder

Hidden Markov Models (HMM): the generative story

<https://powcoder.com>

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

Assignment Project Exam Help

$y_0 \leftarrow \diamond, m \leftarrow 1$

repeat

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

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

until $y_m = \blacklozenge$

\triangleright sample the current tag

\triangleright sample the current word

\triangleright terminate when the stop symbol is generated

The independence assumptions of HMM

<https://powcoder.com>
In addition to the usual independence assumptions for Naïve Bayes, 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:

[Add WeChat powcoder](https://powcoder.com)

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

- Each tag y_m depends only on its predecessor

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

when $y_m = \diamond$ in all cases.

Parameter estimation of HMMs

<https://powcoder.com>

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 of parameters are typically computed from relative frequency estimation on a labeled corpus. The unsmoothed probabilities are,

<https://powcoder.com>
Add WeChat powcoder

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

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

Inference for HMMs

<https://powcoder.com>

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

$$\mathbf{y} = \operatorname{argmax}_y \log P(\mathbf{y} | \mathbf{w})$$

- As $P(\mathbf{y}, \mathbf{w}) = P(\mathbf{y} | \mathbf{w}) \times P(\mathbf{w}) \propto P(\mathbf{y} | \mathbf{w})$, the inference problem can be reformulated as finding the joint probability of \mathbf{w} and \mathbf{y} :

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

Inference for HMMs

- Applying the independent assumptions, we get

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

- 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, one disadvantage of HMMs is that 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:

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

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

<https://powcoder.com>

Assignment Project Exam Help

Assignment Project Exam Help
Structured Perceptron for sequence labeling

<https://powcoder.com>

Add WeChat powcoder

Structured Perceptron

- ▶ Each tagging sequence is assigned a score with a linear model model

$$\begin{aligned}\Psi(\mathbf{w}, \mathbf{y}) &= \sum_{m=1}^{M+1} \psi(\mathbf{w}, y_m, y_{m-1}, m) \\ &= \sum_{m=1}^{M+1} \theta \cdot \mathbf{f}(\mathbf{w}, y_m, y_{m-1}, m) \\ &= \theta \cdot \sum_{m=1}^{M+1} \mathbf{f}(\mathbf{w}, y_m, y_{m-1}, m) \\ &= \theta \cdot \mathbf{f}^{(global)}(\mathbf{w}, \mathbf{y}_{1:M})\end{aligned}$$

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://powcoder.com>

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 same as the correct tag sequence, the parameters are updated.

<https://powcoder.com>

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

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

The averaged perceptron algorithm

```
1: procedure AVE_PERCEPTRON( $\mathbf{x}^{1:N}, \mathbf{y}^{1:N}$ )
2:    $t \leftarrow 0, \boldsymbol{\theta}^{(0)} \leftarrow \mathbf{0}, m \leftarrow 0$ 
3:   repeat
4:      $t \leftarrow t + 1$ 
5:     Select a sequence  $i$  ▷ Online training
6:      $\hat{\mathbf{y}} \leftarrow \operatorname{argmax}_{\mathbf{y}} \boldsymbol{\theta}^{(t-1)} \cdot \mathbf{f}(\mathbf{w}^{(i)}, \mathbf{y})$  ▷ Decoding by Viterbi
7:     if  $\hat{\mathbf{y}} \neq \mathbf{y}^{(i)}$  then
8:        $\boldsymbol{\theta}^{(t)} \leftarrow \boldsymbol{\theta}^{(t-1)} + \mathbf{f}^{(global)}(\mathbf{w}^{(i)}, \mathbf{y}^{(i)}) - \mathbf{f}^{(global)}(\mathbf{w}^{(i)}, \hat{\mathbf{y}})$ 
9:     else ▷ Feature count for entire sequence
10:       $\boldsymbol{\theta}^{(t)} \leftarrow \boldsymbol{\theta}^{(t-1)}$ 
11:       $m \leftarrow m + \boldsymbol{\theta}^{(t)}$ 
12:   until tired
13:    $\bar{\boldsymbol{\theta}} = \frac{1}{t} m$ 
14:   return  $\bar{\boldsymbol{\theta}}$ 
```

Parameter Update example

<https://powcoder.com>

- ▶ Correct tag sequence:
 - ▶ The_DT man_NN who_WP whistles_VBZ tunes_VBZ pianos_NNS
- ▶ Highest scoring sequence under the current model:
 - ▶ The_DT man_NN who_WP whistles_VBZ tunes_NNS pianos_NNS
- ▶ 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://powcoder.com>

Assignment Project Exam Help

Assignment Project Exam Help
Structured Support Vector Machines

<https://powcoder.com>

Add WeChat powcoder

Structured Support Vector Machines

<https://powcoder.com>

Assignment Project Exam Help

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

Add WeChat powcoder

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

- This can be minimally extended to sequence labeling as

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

Extending SVMs to sequences

<https://powcoder.com>

- ▶ A “structured” Support Vector Machine (SVM) outputs a structured object such as a sequence.
- ▶ When extending SVMs to sequence labeling, two issues need to be addressed.
 - ▶ Some errors are more serious than others. Two labelings might have the same “margin” from the correct labeling, but have different number of tokens that are incorrectly labeled. We would prefer to train “against” the latter rather than the former.
 - ▶ Having a fixed margin of 1 would suggest we need to enumerate 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, we use a cost function that reflects how serious the errors are

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

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

$$c(\mathbf{y}, \mathbf{y}^{(i)}) = \sum_{m=1}^{M+1} \delta(y_m \neq y_m^{(i)})$$

- ▶ Instead of training against all labelings \mathbf{y} that have a margin 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 margin constraint, we get:

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

- ▶ This means that the constraint will be met (and training will be complete) when the score for $f(\mathbf{w}^{(i)}, \mathbf{y}^{(i)})$ is equal or greater than the cost-augmented scores for all alternatives.
- ▶ In the training process, we identify predictions that are strong (scores high according to the model) and wrong (incurs a large cost according to the truth), and reducing the scores of 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://powcoder.com>

Assignment Project Exam Help

- ▶ In the training process, the perceptron algorithm simply finds the prediction that has the highest score according to the model to train against, while Structured SVM needs to find the prediction that has a high score with a heavy cost. Both can be identified with the Viterbi algorithm
- ▶ No cost needs to (or can be) computed during decoding for both models
- ▶ In practice, with a large training set, the perceptron algorithm works pretty well, obviating the need for more complicated SVM models.