

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

Assignment Project Exam Help

Assignment Project Add WeChat edu_assist_pro

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

Add WeChat edu_assist_pro

Conditional Random Fields

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

Assignment Project Exam Help

- ▶ The **Conditional Random Field** is a probabilistic model for sequence labeling based on logistic regression.
- ▶ The name is derived from **Conditional Random Fields**, a class of models in which the joint probability distribution is proportional to a product of scores across edges in a factor graph.

The probability model

- ▶ The basic probab

Assignment Project Exam Help

$$P(\mathbf{y}|\mathbf{w}) = \frac{\exp(\Psi(\mathbf{w}, \mathbf{y}))}{\sum_{\mathbf{y}' \in \mathcal{Y}(\mathbf{w})} \exp(\Psi(\mathbf{w}, \mathbf{y}'))}$$

- ▶ This is almost identical to Logistic Regression
label space is se

- ▶ This requires t tag sequence
in decoding, and for *summing over* in
training

- ▶ The usual locality assumption on the scoring function:

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

Decoding in CRFs

- ▶ The Viterbi algorithm finds the best tag sequence, just as it can be used for decoding HMM and Perceptron models.
- ▶ The decoding algorithm is identical to that of p because the denominator is the same for all tags and can thus be ignored:

$$\begin{aligned}\hat{\mathbf{y}} &= \underset{\mathbf{y}}{\operatorname{argmax}} \Psi(\mathbf{y}, \mathbf{w}) \\ &= \underset{\mathbf{y}}{\operatorname{argmax}} \Psi(\mathbf{y}, \mathbf{w}) - \log \left(\sum_{\mathbf{y}' \in \mathcal{Y}(\mathbf{w})} \Psi(\mathbf{y}', \mathbf{w}) \right) \\ &= \underset{\mathbf{y}}{\operatorname{argmax}} \Psi(\mathbf{y}, \mathbf{w}) \\ &= \underset{\mathbf{y}}{\operatorname{argmax}} \sum_{m=1}^{M+1} s_m(y_m, y_{m-1})\end{aligned}$$

Learning in CRFs

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

- ▶ As with logistic regression, the weights can be learned by minimizing the regularized negative log-p

$$\mathcal{L} = \frac{\lambda}{2} \|\boldsymbol{\theta}\|^2 - \sum_{i=1}^n \log p(\boldsymbol{\theta} \cdot \mathbf{f}(\mathbf{w}^{(i)}, \mathbf{y}^{(i)}))$$

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

- ▶ The second term requires computing *the sum of all possible labelings*. There are $|\mathcal{Y}|^M$ possible labeling for an input of length M , so an efficient algorithm is required to compute this sum.

Computing the gradients of the loss function

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

- ▶ As in logistic regression, the gradient of the negative log likelihood with respect to the parameters is the difference between observed and expected feature co

$$\frac{\partial}{\partial \mathbf{w}} \sum_{i=1}^N \ell(\mathbf{w}; \mathbf{x}^{(i)}, \mathbf{y}^{(i)})$$

- ▶ Computing this gradient involves computing the expected count of a sequence that includes the transition $Y_{m-1} \rightarrow Y_m$
- ▶ Recall the feature function for bigram tag sequences is of the form $f(Y_{m-1}, Y_m, \mathbf{w}, m)$. To compute the expected count of a feature, we need $P(Y_{m-1} = k', Y_m = k | \mathbf{w})$.

Marginal probabilities over tag bigrams

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

Assignment Project Exam Help

$$Pr(Y_{m-1} = y'_{m-1}, Y_m = y'_m | Y_{n-1} = y'_{n-1}, Y_n = y'_n) = \frac{\sum_{y'_m: Y_m = k} \exp s_n(y_n, y'_{n-1})}{\sum_{y'_n: Y_n = k} \exp s_n(y_n, y'_{n-1})}$$

How do we compute <https://eduassistpro.github.io/>

- ▶ The naive way to compute this probability is to sum over all possible labelings for the sequence. Since there are $|V|^n$ possible labelings, this is prohibitively expensive for a typical tag set and sentence length.
- ▶ So we need find a more efficient way of doing this.

Computing the numerator

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

The numerator sums over all tag sequences that includes the transition $(Y_{m-1} = k') \rightarrow (Y_m = k)$. This sum can be decomposed into three parts: the prefix $\mathbf{y}_{1:m-1}$, terminating in $Y_{m-1} = k'$, the transition $(Y_{m-1} = k') \rightarrow (Y_m = k)$, and the suffix $\mathbf{y}_{m:M}$, beginning with the tag $Y_m = k$.

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

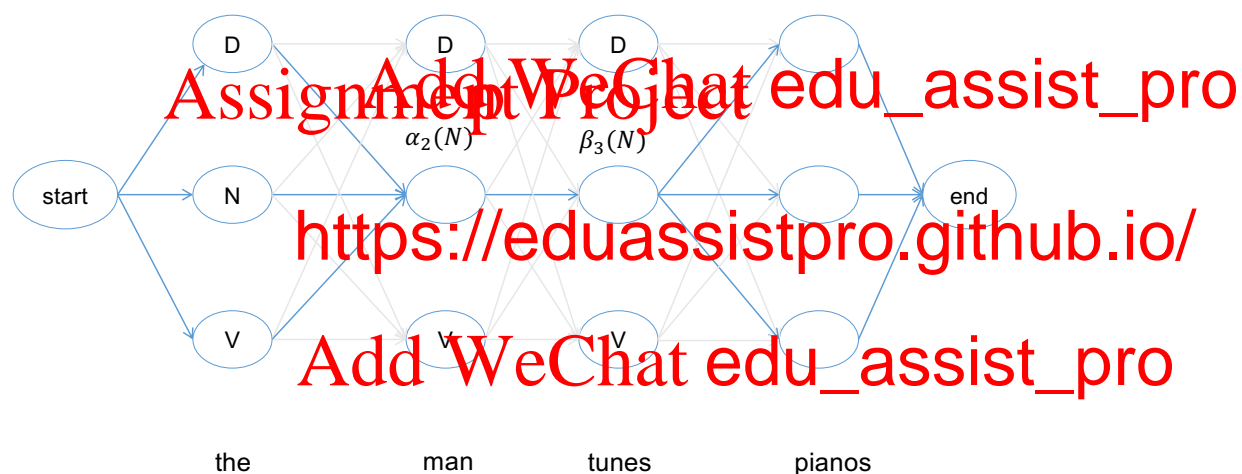
$$\sum_{\mathbf{y}: Y_m=k, Y_{m-1}=k'} \prod_{n=1}^{m-1} \exp s_n(y_n, y_{n-1}) \times \exp s_m(k, k') \times \sum_{\mathbf{y}_{m:M}: Y_m=k} \prod_{n=m+1}^{M+1} \exp s_n(y_n, y_{n-1})$$

Trellis

We can illustrate the nu

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

Assignment Project Exam Help



We compute the numerator by positing the first term as a **forward variable** $\alpha_{m-1}(k')$, and the third term as a **backward variable** $\beta_m(k)$.

Defining a forward variable that can be cached

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

- A **forward variable** $\alpha_n(y_n)$ is equal to the sum of scores of all paths leading to the tag y_n at position n :

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

$$\alpha_n(y_n) \triangleq \sum_{y_{1:n-1}} \exp \left(\sum_{t=1}^n \phi(y_t, y_{1:t-1}) \right)$$

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

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

- The forward recurrence is also known as the **sum-product algorithm** and can be computed through a recurrence

The forward recurrence

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

- Computing the forward variable at position m

$$\alpha_m(y_m) = \sum_{y_{m-1}} \alpha_{m-1}(y_{m-1}) \exp s_m(y_m, y_{m-1})$$

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

$$= \sum_{y_{m-1}} (\exp s_m(y_m, y_{m-1})) \times \alpha_{m-1}(y_{m-1})$$

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

$$= \sum_{y_{m-1}} (\exp s_m(y_m, y_{m-1})) \times \alpha_{m-1}(y_{m-1})$$

The backward recurrence

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

Assignment Project Exam Help

$$\begin{aligned}
 \beta_m(k) &\triangleq \sum_{\mathbf{y}_{m:M}} \exp s_m(k, \mathbf{y}_{m:M}) \beta_{m+1}(\mathbf{y}_m) \\
 &= \sum_{k' \in \mathcal{Y}} \exp s_m(k', k) \beta_{m+1}(k') \\
 &= \sum_{k' \in \mathcal{Y}} \exp s_m(k', k) \times \beta_{m+1}(k')
 \end{aligned}$$

where k' is the label at position $m + 1$.

Computing the denominator

The denominator, the score of the entire sequence, can be computed via forward recurrence, or at any given position m :

- ▶ The score of all possible labelings for the entire sequence is the value of the forward variable at position M in state \diamond :

$$\sum_{\mathbf{y} \in \mathcal{Y}(\mathbf{w})} \Psi(\mathbf{w}, \mathbf{y}) = \alpha_M(\diamond)$$

$$= \sum_{\mathbf{y}_{1:M}} (\exp s_{M+1}(\diamond, y_M) \prod_{m=1}^M p(s_m(y_m, y_{m-1}))$$

- ▶ It can also be computed via backward recurrence as the value of the backward variable at position 0:

$$\sum_{\mathbf{y} \in \mathcal{Y}(\mathbf{w})} \Psi(\mathbf{w}, \mathbf{y}) = \beta_0(\diamond)$$

Features and their weights

f_k	y_{m-1}	y_m						w_k
f_1	D	D	-	-	-	-	-	-0.5
f_2	N	D	-	-	-	-	-	-0.4
f_3	V	D	-	-	-	-	-	-0.4
f_4	D	N	-	-	-	-	-	-0.4
f_5	N	N	-	-	-	-	-	-0.4
f_6	V	N	-	-	-	-	-	3
f_7	D	V	-	-	-	-	-	0.4
f_8	N	V	-	-	-	-	-	0.4
f_9	V	V	-	-	-	-	-	6
f_{10}	-	D	-	-	man	-	-	-0.5
f_{11}	-	N	-	-	man	-	-	2
f_{12}	-	V	-	-	man	-	-	1
f_{13}	-	D	-	the	-	-	-	-4
f_{14}	-	N	-	the	-	-	-	5
f_{15}	-	V	-	the	-	-	-	-2

Computing local transition matrices

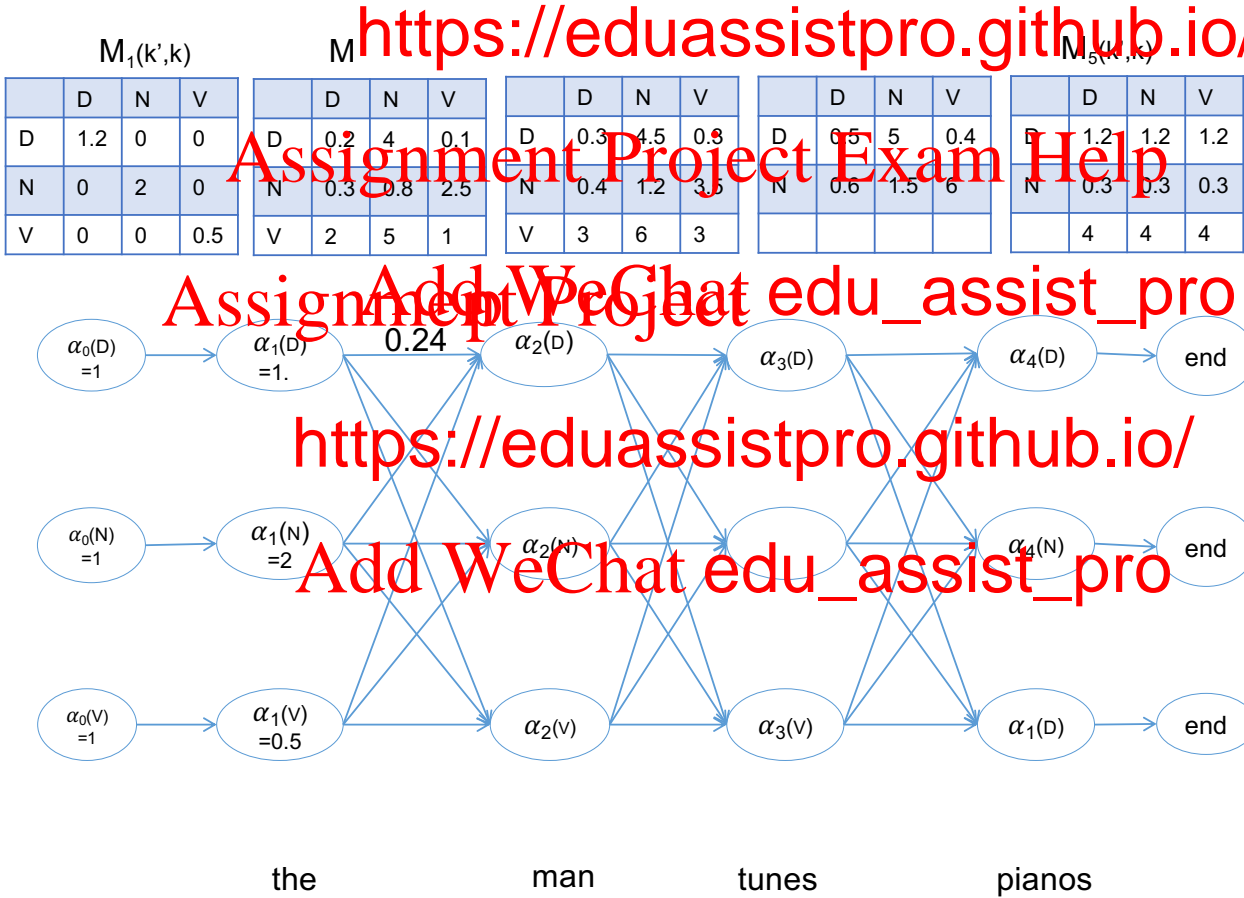
Recall the local score between two tags y_{m-1} and y_m is $\exp(s_m(y_m, y_{m-1}))$. The transition scores from each tag y_{m-1} to the tag y_m can be arranged in a $K \times K$ matrix, where K is the number of tags. Suppose we are computing the score at the position “man” based on the feature weights in the previous slide.

$$\exp(f_1 \theta_1 + f_2 \theta_2 + f_3 \theta_3 + f_4 \theta_4) = 0.007$$

Add WeChat edu_assist_pro

$$\begin{array}{c} D \quad N \quad V \\ \begin{array}{c} D \\ N \\ V \end{array} \begin{bmatrix} 0.007 & ? & ? \\ ? & ? & ? \\ ? & ? & ? \end{bmatrix} \end{array}$$

Forward computation in CRF



Representing features as matrices

Some features span mul

$$f_1(k', k, \mathbf{w}, m) = 1 \text{ iff } w_m = \text{run} \ \& \ k = V$$

$$f_2(k', k, \mathbf{w}, m) = 1 \text{ iff } w_m = \text{run} \ \& \ k = V$$

If represented as matrices, these features have mu have a value of 1:

$$f_1(k', k, \mathbf{w}, m) = \frac{D}{N} \begin{matrix} & D & N & V \\ \begin{matrix} D \\ N \\ V \end{matrix} & \begin{bmatrix} 0 & 1 & 0 \\ 0 & 1 & 0 \\ 0 & 1 & 0 \end{bmatrix} \end{matrix}$$

Feature expectations

- ▶ Let f be any local feature. The count of the this feature in a particular sequence is

$$F(\mathbf{y}, \mathbf{w}) = \sum_{m=1}^{M+1} f(\mathbf{y}_m, \mathbf{w})$$

- ▶ Expectation of the feature f in a sequence is:

$$\begin{aligned} \mathbb{E}[F(\mathbf{w}, \mathbf{y})] &= \sum_{m=1}^M \frac{\alpha_{m-1}(f) \odot \mathbf{w}_m}{Z(\mathbf{w})} \\ &= \sum_{m=1}^M \frac{\alpha_{m-1}(f) \odot \mathbf{w}_m}{Z(\mathbf{w})} \end{aligned}$$

where $Z(\mathbf{w})$ is the total score of all labelings, also known as the **partition function**.

- ▶ Given the feature expectations, we can now compute the gradient of each feature.

Example computation of feature expectations

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

- Assume the transition matrix at position m :

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

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

V 6.1

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

- Assume the following vectors:

- $\alpha_{m-1}(k') = [40 \quad 30 \quad 65]$

- $\beta_m(k) = [45 \quad 65 \quad 30]$

Computing Feature Expectations

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

Assignment Project Exam Help

$$\alpha_{m-1}(f_1 \odot M_m)\beta_m^\top = \begin{bmatrix} 40 & 30 & 65 & 0 & 0.1 & 0 \end{bmatrix} \begin{bmatrix} 45 \\ 65 \\ 30 \end{bmatrix}$$

Assignment Project Exam Help

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

Add WeChat edu_assist_pro

$$\alpha_{m-1}(f_1 \odot M_m)\beta_m^\top = \begin{bmatrix} 40 & 30 & 65 \end{bmatrix} \begin{bmatrix} 0 & 0 & 0.4 \\ 0 & 0 & 0.3 \\ 0 & 0 & 0.01 \end{bmatrix} \begin{bmatrix} 45 \\ 65 \\ 30 \end{bmatrix}$$

=?

CRFs

- ▶ Avoids the strong independence assumption
- ▶ Allows arbitrary features over observations (but not hidden states)
- ▶ Sound probabilistic model – gives us a distribution over sequence labelings
- ▶ Uses the same efficient forward-backward algorithm for decoding as state transition models
- ▶ Estimation/learning is harder (than state transition models) – we have to compute the posterior for a sequence of labels
- ▶ Empirical results generally show CRFs outperforms HMMs (and other classifiers)
- ▶ Feature estimation can be replaced with LSTM RNNs, resulting in what's called RNN-CRFs, of which LSTM-CRFs are the most widely used.