

<https://powcoder.com>

Assignment Project Exam Help

Assignment Project Exam Help

Conditional Random Fields

<https://powcoder.com>

Add WeChat powcoder

## Conditional Random Fields

<https://powcoder.com>

### Assignment Project Exam Help

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

## The probability model

<https://powcoder.com>

- ▶ The basic probability model is

Assignment Project Exam Help

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

Assignment Project Exam Help

- ▶ This is almost identical to Logistic Regression, except that the label space is sequence of tags
- ▶ This requires efficient algorithms to *find the best tag sequence* in decoding, and for *summing over all tag sequences* 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 can be used to find the best tag sequence, just as it can be used for decoding HMM and Perpcetron models.
- ▶ The decoding algorithm is identical to that of perceptron, because the denominator is the same for all tag sequences and can thus be ignored:

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

## Learning in CRFs

<https://powcoder.com>

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

$$\begin{aligned}\mathcal{L} &= \frac{\lambda}{2} \|\boldsymbol{\theta}\|^2 - \sum_{i=1}^M \log P(\mathbf{y}^{(i)} | \mathbf{w}^{(i)}; \boldsymbol{\theta}) \\ &= \frac{\lambda}{2} \|\boldsymbol{\theta}\|^2 - \sum_{i=1}^M \boldsymbol{\theta} \cdot \mathbf{f}(\mathbf{w}^{(i)}, \mathbf{y}^{(i)}) + \log \sum_{\mathbf{y}' \in \mathcal{Y}^M} \exp(\boldsymbol{\theta} \cdot \mathbf{f}(\mathbf{w}^{(i)}, \mathbf{y}'))\end{aligned}$$

- ▶ 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://powcoder.com>

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

$$\frac{\partial \mathcal{L}}{\partial \theta_j} = \lambda \theta_j + \sum_{i=1}^N \mathbb{E}[f_j(\mathbf{w}^{(i)}, \mathbf{y})] - f_j(\mathbf{w}^{(i)}, \mathbf{y}^{(i)})$$

- ▶ Computing this gradient involves computing the posterior 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://powcoder.com>

## Assignment Project Exam Help

$$Pr(Y_{m-1} = y'_{m-1}, Y_m = y_m | w) = \frac{\sum_{y'_m: Y_m = y_m, Y_{m-1} = y'_{m-1}} \prod_{n=1}^{M+1} \exp s_n(y_n, y_{n-1})}{\sum_{y'} \prod_{n=1}^{M+1} \exp s_n(y'_n, y'_{n-1})}$$

How do we compute the numerator and denominator?

- ▶ The naive way to compute this probability is to enumerate all possible labelings for the sequence. Since there are  $\gamma^M$  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://powcoder.com>

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  $Y_m = k$ .

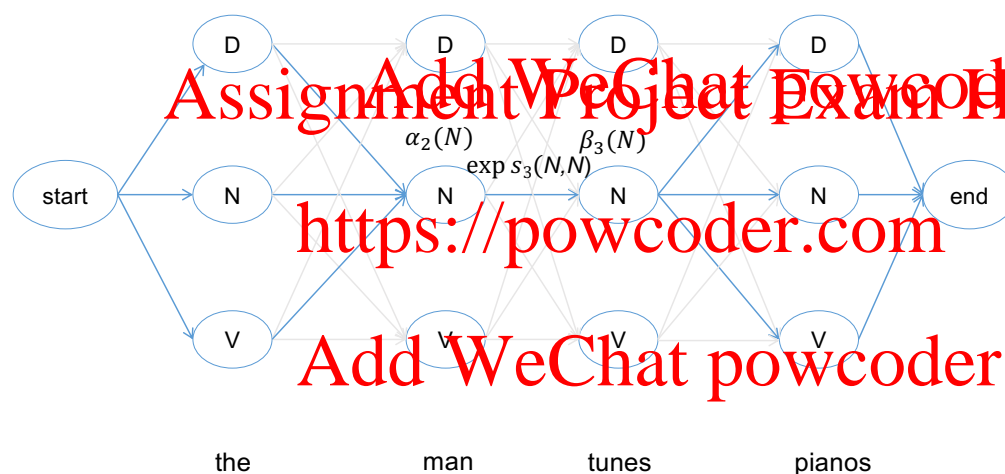
<https://powcoder.com>

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



## Trellis

We can illustrate the numerator graphically with a trellis:



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://powcoder.com>

- A **forward variable**  $\alpha_m(y_m)$  is equal to the sum of scores of all paths leading to the tag  $y_m$  at position  $m$ :

$$\begin{aligned}\alpha_m(y_m) &\triangleq \sum_{y_{1:m-1}} \exp \sum_{n=1}^{m-1} s_n(y_n, y_{n-1}) \\ &= \sum_{y_{1:m-1}} \prod_{n=1}^{m-1} \exp s_n(y_n, y_{n-1})\end{aligned}$$

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

The forward recurrence

<https://powcoder.com>

- Computing the forward variable at position  $m$

$$\begin{aligned}\alpha_m(y_m) &= \sum_{y_{1:m-1}} \prod_{n=1}^m \exp s_n(y_n, y_{n-1}) \\ &= \sum_{y_{m-1}} (\exp s_m(y_m, y_{m-1})) \sum_{y_{1:m-1}} \prod_{n=1}^{m-1} \exp s_n(y_n, y_{n-1}) \\ &= \sum_{y_{m-1}} (\exp s_m(y_m, y_{m-1})) \times \alpha_{m-1}(y_{m-1})\end{aligned}$$

The backward recurrence

<https://powcoder.com>

Assignment Project Exam Help

$$\begin{aligned}
 \beta_m(k) &\triangleq \sum_{\mathbf{y}_{m:M}: Y_m=k} \prod_{n=m}^{M+1} \exp s_n(y_n, y_{n-1}) \\
 &= \sum_{k' \in \mathcal{Y}} \exp s_m(k', k) \sum_{\mathbf{y}_{m+1:M}: Y_{m+1}=k'} \prod_{n=m+1}^{M+1} \exp s_n(y_n, y_{n-1}) \\
 &= \sum_{k' \in \mathcal{Y}} \exp s_m(k', k) \times \beta_{m+1}(k')
 \end{aligned}$$

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

## Computing the denominator

The denominator, the score of all possible labelings for the entire sequence, can be computed either via the forward recurrence or backward 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  $M + 1$  in state  $\diamond$ :

$$\begin{aligned}\sum_{\mathbf{y} \in \mathcal{Y}(\mathbf{w})} \Psi(\mathbf{w}, \mathbf{y}) &= \alpha_{M+1}(\diamond) \\ &= \sum_{\mathbf{y}_{1:M}} (\exp s_{M+1}(\diamond, y_M)) \prod_{m=1}^M \exp s_m(y_m, y_{m-1})\end{aligned}$$

- 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

$f_k$	$y_{m-1}$	$y_m$	$w_{-2}$	$w_{-1}$	$w_0$	$w_1$	$w_2$	$\theta_k$
$f_1$	D	D	-	-	-	-	-	-0.5
$f_2$	N	D	-	-	-	-	-	-0.8
$f_3$	V	D	-	-	-	-	-	1.2
$f_4$	D	N	-	-	-	-	-	2.5
$f_5$	N	N	-	-	-	-	-	0.5
$f_6$	V	N	-	-	-	-	-	3
$f_7$	D	V	-	-	-	-	-	0.4
$f_8$	N	V	-	-	-	-	-	4
$f_9$	V	V	-	-	-	-	-	0.6
$f_{10}$	-	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

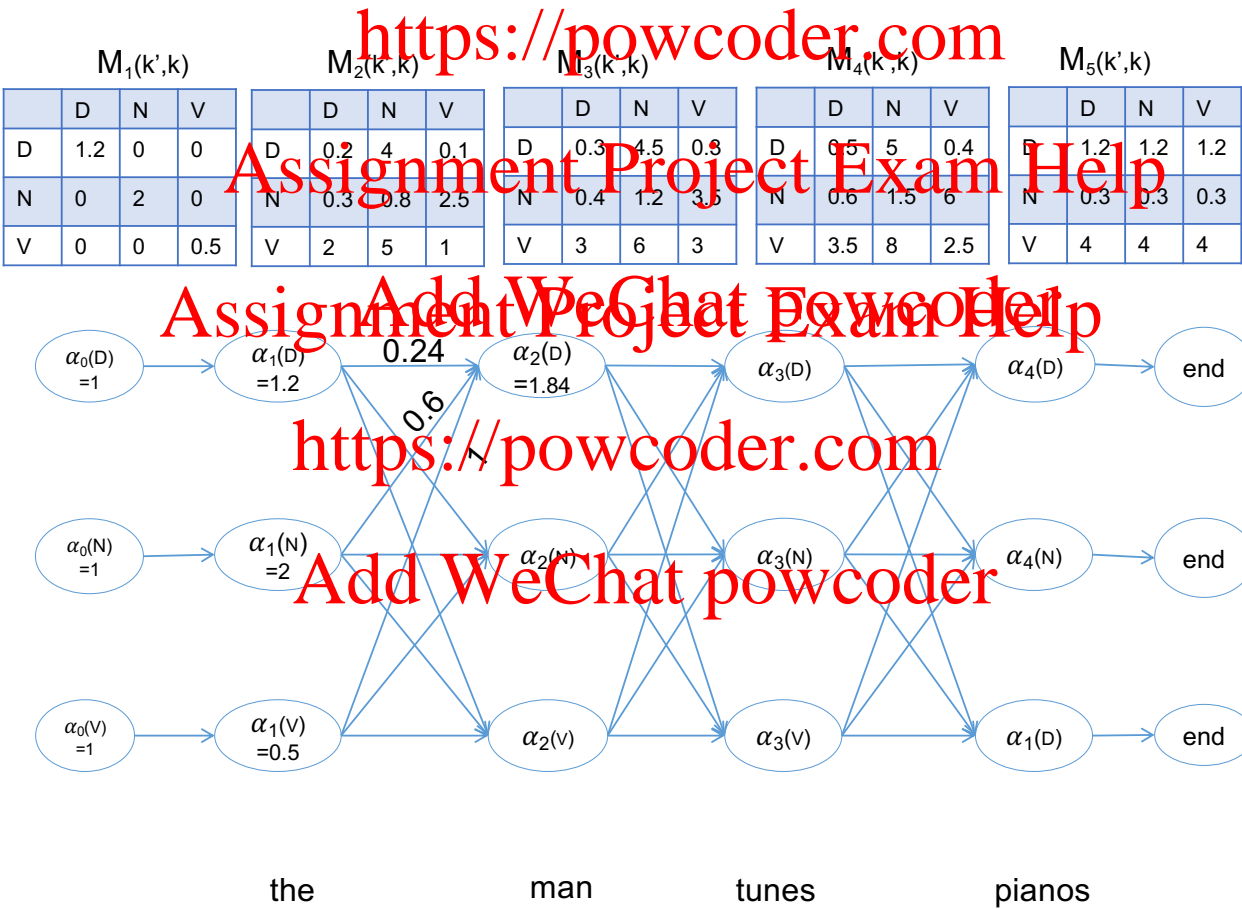
Recall the local score between two positions in a tag sequence is  $\exp s_m(y_m, y_{m-1})$ . The transition scores from each tag  $y_{m-1}$  to the tag  $y_m$  can be arranged in a  $K \times K$  matrix, where  $K$  is the number of tags. Suppose we are computing the score at the position “man” based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

$$\exp(f_1\theta_1 + f_{10}\theta_{10} + f_{13}\theta_{13}) = \exp(-0.5 - 0.5 - 4) = 0.007$$

Add WeChat powcoder

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

Forward computation in CRF





## Representing features as matrices

Some features span multiple state (tag) transitions:

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

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

If represented as matrices, these features have multiple cells that have a value of 1:

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

$$f_2(k', k, \mathbf{w}, m) = \begin{matrix} & D & N & V \\ \begin{matrix} D \\ N \\ V \end{matrix} & \begin{bmatrix} 0 & 1 & 0 \\ 0 & 1 & 0 \\ 0 & 1 & 0 \end{bmatrix} \end{matrix}$$

## Feature expectations

- ▶ Let  $f$  be any local feature,  $f_m[k', k] = f(k', k, \mathbf{w}, m)$ . The count of the this feature in a particular sequence is

Assignment Project Exam Help

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

Assignment Project Exam Help

- ▶ Expectation for a single feature over a sequence is:

<https://powcoder.com>

$$\mathbb{E}[F(\mathbf{w}, \mathbf{y})] = \sum_m \Pr(Y_{m-1} = k', Y_m = k | \mathbf{w}) f(k', k, \mathbf{w}, m)$$

Add WeChat powcoder

$$= \sum_m \frac{\alpha_{m-1} (f \odot M_m) \beta_m^\top}{Z(\mathbf{w})}$$

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://powcoder.com>

- Assume the transition matrix at position  $m$ :

<https://powcoder.com>

$$\begin{matrix} D \\ N \\ V \end{matrix} \begin{matrix} M[k, P] \\ M[k, B] \\ M[k, W] \end{matrix} = \begin{bmatrix} 2.2 & 0.1 & 0.4 \\ 1.2 & 3.4 & 0.3 \\ 6.1 & 7.2 & 0.01 \end{bmatrix}$$

<https://powcoder.com>

- 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://powcoder.com>

Assignment Project Exam Help

$$\alpha_{m-1}(f_1 \odot M_m)\beta_m^\top = [40 \quad 30 \quad 65] \begin{bmatrix} 0 & 0.1 & 0 \\ 0 & 3.4 & 0 \\ 0 & 7.2 & 0 \end{bmatrix} \begin{bmatrix} 45 \\ 65 \\ 30 \end{bmatrix}$$

Assignment Project Exam Help

$$= [0 \quad 574 \quad 0.0] \begin{bmatrix} 45 \\ 65 \\ 30 \end{bmatrix} = 37310.0$$

Add WeChat powcoder

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

=?

## CRFs

- ▶ Avoids the strong independence assumptions of HMMs
- ▶ Allows arbitrary features over observations (but not hidden states)
- ▶ Sound probabilistic model – gives us a distribution over sequence labelings
- ▶ Uses the same efficient algorithm (the Viterbi Algorithm) for decoding as structured perceptron and HMMs
- ▶ Estimation/learning is harder (than say Perceptron) because we have to compute the posterior for a sequence
- ▶ 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.