

Conditional Random Fields

1 Introduction

Conditional random fields (CRFs) define a distribution over target structures \mathbf{y} given an input structure \mathbf{x} , using a global feature function Φ and a weight vector \mathbf{w} :

$$p(\mathbf{y}|\mathbf{x}) = \frac{\exp(\mathbf{w} \cdot \Phi(\mathbf{x}, \mathbf{y}))}{\sum_{\mathbf{y}'} \exp(\mathbf{w} \cdot \Phi(\mathbf{x}, \mathbf{y}'))} \quad (1)$$

CRFs are powerful because Φ can consider an *entire* target structure \mathbf{y} . We will focus on the case where the structure is a sequence: \mathbf{x} is a sequence of tokens $x_1 \dots x_n$ and \mathbf{y} is a sequence of labels $y_1 \dots y_n$ where label $y_i \in \mathcal{Y}$ corresponds to token $x_i \in \mathcal{X}$. A key assumption in CRFs to keep learning and decoding efficient is the following.

Assumption 1.1. *The global feature function Φ is a sum of local functions ϕ . Throughout, we will assume that $\Phi(\mathbf{x}, \mathbf{y}) \in \mathbb{R}^d$ has the form*

$$\Phi(\mathbf{x}, \mathbf{y}) = \sum_i \phi(\mathbf{x}, i, y_{i-1}, y_i) \quad (2)$$

where $\phi(\mathbf{x}, i, y_{i-1}, y_i) \in \mathbb{R}^d$ is the local representation of \mathbf{y} at position i .

2 Decoding

Suppose that we know the weight vector \mathbf{w} . Given an input sequence \mathbf{x} , how do we find a target sequence \mathbf{y}^* such that $p(\mathbf{y}^*|\mathbf{x}) \geq p(\mathbf{y}|\mathbf{x})$ for all possible \mathbf{y} ? Rather than calculating Eq. (1) for exponentially many candidates $\mathbf{y} \in \mathcal{Y}^n$, we make the following observation based on Assumption 1.1.

$$\begin{aligned} \mathbf{y}^* &= \arg \max_{\mathbf{y}} p(\mathbf{y}|\mathbf{x}) \\ &= \arg \max_{\mathbf{y}} \frac{\exp(\mathbf{w} \cdot \Phi(\mathbf{x}, \mathbf{y}))}{\sum_{\mathbf{y}'} \exp(\mathbf{w} \cdot \Phi(\mathbf{x}, \mathbf{y}'))} \\ &= \arg \max_{\mathbf{y}} \exp(\mathbf{w} \cdot \Phi(\mathbf{x}, \mathbf{y})) \\ &= \arg \max_{\mathbf{y}} \mathbf{w} \cdot \Phi(\mathbf{x}, \mathbf{y}) \\ &= \arg \max_{\mathbf{y}} \sum_i \mathbf{w} \cdot \phi(\mathbf{x}, i, y_{i-1}, y_i) \end{aligned}$$

Note that if Φ is not decomposed, we cannot avoid explicitly enumerating \mathbf{y} . Thanks to the decomposition, however, a Viterbi algorithm can be used to find \mathbf{y}^* such that

$$\mathbf{y}^* = \arg \max_{\mathbf{y}} \sum_i \mathbf{w} \cdot \phi(\mathbf{x}, i, y_{i-1}, y_i)$$

In the algorithm, table π is used to record for each position i and state y

$$\pi[i, y] = \max_{\substack{y_1 \dots y_i \\ y_i = y}} \sum_{j=1}^i \mathbf{w} \cdot \phi(\mathbf{x}, j, y_{j-1}, y_j)$$

Then \mathbf{y}^* can be retrieved from $\max_y \pi[n, y]$ using a backtracker κ .

Viterbi
Input: input sequence \mathbf{x} of length n , weight vector \mathbf{w}
Output: optimal target sequence \mathbf{y}^*

- For all $y \in \mathcal{Y}$,

$$\pi[1, y] \leftarrow \mathbf{w} \cdot \phi(\mathbf{x}, 1, y_0, y)$$
 where y_0 is a start symbol
- For $i = 2 \dots n$, for all $y \in \mathcal{Y}$,

$$\pi[i, y] \leftarrow \max_{y'} \pi[i-1, y'] + \mathbf{w} \cdot \phi(\mathbf{x}, i, y', y)$$

$$\kappa[i, y] \leftarrow \arg \max_{y'} \pi[i-1, y'] + \mathbf{w} \cdot \phi(\mathbf{x}, i, y', y)$$
- Return $\mathbf{y}^* = y_1^* \dots y_n^*$ where

$$y_n^* \leftarrow \arg \max_y \pi[n, y]$$

$$y_i^* \leftarrow \kappa[i, y_{i+1}^*] \text{ for } i = n-1 \dots 1$$

3 Learning

We calculate \mathbf{w} that maximizes the regularized log likelihood of the training data $\{(\mathbf{x}^q, \mathbf{y}^q)\}_{q=1}^Q$ with some parameter $\lambda > 0$:

$$L(\mathbf{w}) = \sum_{q=1}^Q \log p(\mathbf{y}^q | \mathbf{x}^q) - \frac{\lambda}{2} \|\mathbf{w}\|^2$$

$L(\mathbf{w})$ is a convex function in \mathbf{w} , so we can resort to some hill climbing method to find an optimal weight vector \mathbf{w}^* such that

$$\mathbf{w}^* = \arg \max_{\mathbf{w}} L(\mathbf{w})$$

Below is a gradient-based algorithm in its roughest form, but in practice one will have to be more clever.

GradientAscent**Input:** training sequences $\{(\mathbf{x}^q, \mathbf{y}^q)\}_{q=1}^Q$, feature function $\Phi : \mathcal{X} \times \mathcal{Y} \rightarrow \mathbb{R}^d$ **Output:** weight vector $\mathbf{w}^* \in \mathbb{R}^d$ close to optimal

- $\mathbf{w} \leftarrow \mathbf{0}$
- Until convergence,
 - For $k = 1 \dots d$, calculate the gradient $\Delta_k \leftarrow \frac{\partial L(\mathbf{w})}{\partial w_k}$ and set

$$w_k = w_k + \alpha \times \Delta_k$$

where $\alpha > 0$ is some learning rate.

- Return \mathbf{w} .

We need to be able to compute the gradient of the function $L(\mathbf{w})$ with respect to each component w_k . Since CRFs are just log-linear models, we have the well-known form

$$\frac{\partial L(\mathbf{w})}{\partial w_k} = \sum_{q=1}^Q \Phi_k(\mathbf{x}^q, \mathbf{y}^q) - \sum_{q=1}^Q \sum_{\mathbf{y} \in \mathcal{Y}^n} p(\mathbf{y}|\mathbf{x}^q) \Phi_k(\mathbf{x}^q, \mathbf{y}) - \lambda w_k$$

The first and last terms are easy to compute. The middle term, which involves a sum over exponentially many sequences $\mathbf{y} \in \mathcal{Y}^n$, can be again computed efficiently by the following observation based on Assumption 1.1. For any $\mathbf{x} \in \mathcal{X}^n$,

$$\begin{aligned} \sum_{\mathbf{y} \in \mathcal{Y}^n} p(\mathbf{y}|\mathbf{x}) \Phi_k(\mathbf{x}, \mathbf{y}) &= \sum_{\mathbf{y} \in \mathcal{Y}^n} p(\mathbf{y}|\mathbf{x}) \sum_{i=1}^n \phi_k(\mathbf{x}, i, y_{i-1}, y_i) \\ &= \sum_{i=1}^n \sum_{a, b \in \mathcal{Y}} \sum_{\substack{\mathbf{y} \in \mathcal{Y}^n: \\ y_{i-1}=a, y_i=b}} p(\mathbf{y}|\mathbf{x}) \phi_k(\mathbf{x}, i, y_{i-1}, y_i) \\ &= \sum_{i=1}^n \sum_{a, b \in \mathcal{Y}} \phi_k(\mathbf{x}, i, a, b) \sum_{\substack{\mathbf{y} \in \mathcal{Y}^n: \\ y_{i-1}=a, y_i=b}} p(\mathbf{y}|\mathbf{x}) \\ &= \sum_{i=1}^n \sum_{a, b \in \mathcal{Y}} \phi_k(\mathbf{x}, i, a, b) g(i, a, b|\mathbf{x}) \end{aligned}$$

where we define

$$g(i, a, b|\mathbf{x}) = \sum_{\substack{\mathbf{y} \in \mathcal{Y}^n: \\ y_{i-1}=a, y_i=b}} p(\mathbf{y}|\mathbf{x})$$

We can compute $g(i, a, b|\mathbf{x})$ for all $i \in [n]$ and $a, b \in \mathcal{Y}$ efficiently using a forward-backward algorithm. Thus we have shown that the gradient can be computed efficiently for training.

3.1 Forward-Backward

Given $\mathbf{x} \in \mathcal{X}^n$, we compute for every position $i \in [n]$ and label $y \in \mathcal{Y}$

$$\alpha[i, y] = \sum_{\substack{y_1 \dots y_i: \\ y_i = y}} \prod_{j=1}^i \exp(\mathbf{w} \cdot \phi(\mathbf{x}, j, y_{j-1}, y_j))$$

$$\beta[i, y] = \sum_{\substack{y_i \dots y_n: \\ y_i = y}} \prod_{j=i+1}^n \exp(\mathbf{w} \cdot \phi(\mathbf{x}, j, y_{j-1}, y_j))$$

Then it can be verified that

$$g(i, a, b|\mathbf{x}) = \sum_{\substack{\mathbf{y} \in \mathcal{Y}^n: \\ y_{i-1} = a, y_i = b}} p(\mathbf{y}|\mathbf{x})$$

$$= \frac{\alpha[i-1, a] \times \exp(\mathbf{w} \cdot \phi(\mathbf{x}, i, a, b)) \times \beta[i, b]}{\sum_{y \in \mathcal{Y}} \alpha[n, y]}$$

ForwardBackward

Input: input sequence \mathbf{x} of length n

Output: $g(i, a, b|\mathbf{x})$ for all $i \in [n]$ and $a, b \in \mathcal{Y}$

- For all $y \in \mathcal{Y}$,

$$\alpha[1, y] \leftarrow \exp(\mathbf{w} \cdot \phi(\mathbf{x}, 1, y_0, y)) \text{ where } y_0 \text{ is a start symbol}$$

- For $i = 2 \dots n$, for all $y \in \mathcal{Y}$,

$$\alpha[i, y] \leftarrow \sum_{y' \in \mathcal{Y}} \alpha[i-1, y'] \times \exp(\mathbf{w} \cdot \phi(\mathbf{x}, i, y', y))$$

- For all $y \in \mathcal{Y}$,

$$\beta[n, y] \leftarrow 1$$

- For $i = n-1 \dots 1$, for all $y \in \mathcal{Y}$,

$$\beta[i, y] \leftarrow \sum_{y' \in \mathcal{Y}} \exp(\mathbf{w} \cdot \phi(\mathbf{x}, i+1, y, y')) \times \beta[i+1, y']$$

- Return $g(i, a, b|\mathbf{x})$ where

$$g(i, a, b|\mathbf{x}) = \frac{\alpha[i-1, a] \times \exp(\mathbf{w} \cdot \phi(\mathbf{x}, i, a, b)) \times \beta[i, b]}{\sum_{y \in \mathcal{Y}} \alpha[n, y]}$$