

Machine Learning I

Bjoern Andres

Machine Learning for Computer Vision
TU Dresden

Conditional Graphical Models II

Contents. This part of the course introduces algorithms for supervised structured learning of conditional graphical models.

Conditional Graphical Models II

On the one hand, supervised structured learning of conditional graphical models whose factors are linear functions is a **convex** optimization problem.

Thus, it can be solved exactly by means of the **steepest descent algorithm** with a tolerance parameter $\epsilon \in \mathbb{R}_0^+$:

```
 $\theta := 0$   
repeat  
   $d := \nabla_{\theta} L(H_{\theta}(x, \cdot), y)$   
   $\eta := \operatorname{argmin}_{\eta' \in \mathbb{R}} L(H_{\theta - \eta' d}(x, \cdot), y)$  (line search)  
   $\theta := \theta - \eta d$   
  if  $\|d\| < \epsilon$   
    return  $\theta$ 
```

Conditional Graphical Models II

On the other hand, computing the gradient naïvely takes time $O(2^{|S|})$:

$$\begin{aligned}
 -\frac{\partial}{\partial \theta_j} \ln Z &= \mathbb{E}_{y' \sim p_{Y|X, \theta}} (\xi_j(x, y')) \\
 &= \frac{1}{Z(x, \theta)} \sum_{y' \in \{0,1\}^S} \xi_j(x, y') e^{-\langle \theta, \xi(x, y') \rangle} \\
 &= \frac{1}{Z(x, \theta)} \sum_{y' \in \{0,1\}^S} \sum_{f \in F} \varphi_{fj}(x_f, y'_{S_f}) e^{-\langle \theta, \xi(x, y') \rangle} \\
 &= \frac{1}{Z(x, \theta)} \sum_{f \in F} \sum_{y'_{S(f)} \in \{0,1\}^{S(f)}} \sum_{y'_{S \setminus S(f)} \in \{0,1\}^{S \setminus S(f)}} \varphi_{fj}(x_f, y'_{S(f)}) e^{-\langle \theta, \xi(x, y') \rangle} \\
 &= \sum_{f \in F} \sum_{y'_{S(f)} \in \{0,1\}^{S(f)}} \varphi_{fj}(x_f, y'_{S(f)}) \frac{1}{Z(x, \theta)} \sum_{y'_{S \setminus S(f)} \in \{0,1\}^{S \setminus S(f)}} e^{-\langle \theta, \xi(x, y') \rangle} \\
 &= \sum_{f \in F} \sum_{y'_{S(f)} \in \{0,1\}^{S(f)}} \varphi_{fj}(x_f, y'_{S(f)}) p_{Y_{S(f)}|X, \theta}(y'_{S(f)} | x, \theta) \\
 &= \sum_{f \in F} \mathbb{E}_{y'_{S(f)} \sim p_{Y_{S(f)}|X, \theta}} (\varphi_{fj}(x_f, y'_{S(f)}))
 \end{aligned}$$

Conditional Graphical Models II

Computing the gradient requires that we compute

- ▶ the partition function

$$Z(x, \theta) = \sum_{y' \in \{0,1\}^S} e^{-\langle \theta, \xi(x, y') \rangle} \quad (1)$$

- ▶ for every factor $f \in F$, the so-called **factor marginal**

$$p_{\mathcal{Y}_{S(f)} | \mathcal{X}, \Theta}(y'_{S(f)} | x, \theta) = \frac{1}{Z(x, \theta)} \sum_{y'_{S \setminus S(f)} \in \{0,1\}^{S \setminus S(f)}} e^{-\langle \theta, \xi(x, y') \rangle} \quad (2)$$

- ▶ for every factor $f \in F$, the expectation value

$$\sum_{y'_{S(f)} \in \{0,1\}^{S(f)}} \varphi_{fj}(x_f, y'_{S(f)}) p_{\mathcal{Y}_{S(f)} | \mathcal{X}, \Theta}(y'_{S(f)} | x, \theta) . \quad (3)$$

Conditional Graphical Models II

The challenge is to sum the function

$$\psi_{\theta}(x, y') := e^{-\langle \theta, \xi(x, y') \rangle} \quad (4)$$

over assignments of 0 or 1 to linearly many (2) or all (1) variables y' .

Defining

$$\psi_{f\theta}(x_f, y'_{S(f)}) = e^{-\langle \theta, \varphi_f(x_f, y'_{S(f)}) \rangle} \quad (5)$$

we obtain

$$\begin{aligned} \psi_{\theta}(x, y') &= e^{-\langle \theta, \xi(x, y') \rangle} \\ &= e^{-\sum_{f \in F} \langle \theta, \varphi_f(x_f, y_{S(f)}) \rangle} \end{aligned} \quad (6)$$

$$= \prod_{f \in F} e^{-\langle \theta, \varphi_f(x_f, y_{S(f)}) \rangle} \quad (7)$$

$$= \prod_{f \in F} \psi_{f\theta}(x_f, y_{S(f)}) \cdot \quad (8)$$

Conditional Graphical Models II

Thus, the challenge in (2) and (1) is to compute a sum of a product of functions. Specifically:

$$Z(x, \theta) = \sum_{y' \in \{0,1\}^S} \prod_{f \in F} \psi_{f\theta}(x_f, y_{S(f)}) \quad (9)$$

$$p_{\mathcal{Y}_{S(f)} | \mathcal{X}, \Theta}(y'_{S(f)} | x, \theta) = \frac{1}{Z(x, \theta)} \sum_{y'_{S \setminus S(f)} \in \{0,1\}^{S \setminus S(f)}} \prod_{f \in F} \psi_{f\theta}(x_f, y_{S(f)}) \quad (10)$$

- ▶ One approach to tackle this problem is to sum over variables recursively.
- ▶ In order to avoid redundant computation, Kschischang et al. (2001) define partial sums.

Definition (Kschischang et al. (2001)) For any variable node $s \in \mathcal{S}$ and any factor node $f \in F$, the functions

$$m_{s \rightarrow f}, m_{f \rightarrow s} : \{0, 1\} \rightarrow \mathbb{R} , \quad (11)$$

called **messages**, are defined such that for all $y_s \in \{0, 1\}$:

$$m_{s \rightarrow f}(y_s) = \prod_{f' \in F(s) \setminus \{f\}} m_{f' \rightarrow s}(y_s) \quad (12)$$

$$m_{f \rightarrow s}(y_s) = \sum_{y_{S(f) \setminus \{s\}}} \psi_{f\theta}(x_f, y_{S(f)}) \prod_{s' \in S(f) \setminus \{s\}} m_{s' \rightarrow f}(y_{s'}) \quad (13)$$

Conditional Graphical Models II

Lemma. If the factor graph is acyclic, messages are defined recursively by (12) and (13), beginning with the messages from leaves. Moreover, for any $s \in S$ and any $f \in F$:

$$Z(x, \theta) = \sum_{y_s \in \{0,1\}} \prod_{f' \in F(s)} m_{f' \rightarrow s}(y_s) \quad (14)$$

$$p_{\mathcal{Y}_{S(f)} | \mathcal{X}, \Theta}(y'_{S(f)} | x, \theta) = \frac{1}{Z(x, \theta)} \psi_{f\theta}(x_f, y_{S(f)}) \prod_{s' \in S(f)} m_{s' \rightarrow f}(y_{s'}) \quad (15)$$

The recursive computation of messages is known as **message passing**.

Summary

- ▶ For conditional graphical models whose factor graph is **acyclic**, the supervised structured learning problem can be solved efficiently by means of the steepest descent algorithm and message passing.
- ▶ For conditional graphical models whose factor graph is **cyclic**, the definition of messages is cyclic as well. The partition function and marginals cannot be computed by message passing in general.
- ▶ A heuristic without guarantee of correctness or even convergence is to initialize all messages as normalized constant functions and to update messages according to some schedule, e.g., synchronously. This heuristic is commonly known as **loopy belief propagation**.