

Machine Learning I

B. Andres, J. Irmay, J. Presberger, D. Stein, S. Zhao

Machine Learning for Computer Vision
TU Dresden



Winter Term 2023/2024

Conditional Graphical Models

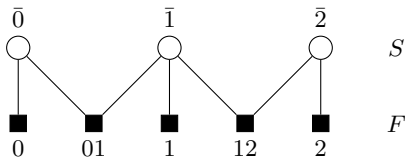
Contents. This part of the course is about supervised structured learning of conditional graphical models.

Conditional Graphical Models

Definition. For any factor graph $G = (S, F, E)$, a function $H : \{0, 1\}^S \rightarrow \mathbb{R}$ is said to **factorize** w.r.t. G iff, for every $f \in F$, there exists a function a function $h_f : \{0, 1\}^{S_f} \rightarrow \mathbb{R}$, called a **factor** of H , such that

$$\forall y \in \{0, 1\}^S : \quad H(y) = \sum_{f \in F} h_f(y_{S_f}) . \quad (1)$$

Example: A function $H : \{0, 1\}^S \rightarrow \mathbb{R}$ factorizes w.r.t. the factor graph



iff there exist suitable functions $h_0, h_{01}, h_1, h_{12}, h_2$ such that, for any $y \in \{0, 1\}^S$: $H(y) = h_0(y_{\bar{0}}) + h_1(y_{\bar{1}}) + h_2(y_{\bar{2}}) + h_{01}(y_{\bar{0}}, y_{\bar{1}}) + h_{12}(y_{\bar{1}}, y_{\bar{2}})$.

Definition. A tuple $(S, F, E, \{X_f\}_{f \in F}, \Theta, \{h_f\}_{f \in F})$ is called a **conditional graphical model** with attribute space $X := \prod_{f \in F} X_f$ and parameter space Θ iff the following conditions hold:

- ▶ (S, F, E) is a factor graph
- ▶ $\Theta \neq \emptyset$
- ▶ For every $f \in F$:
 - ▶ X_f is non-empty, called the **attribute space** of f
 - ▶ $h_f : \Theta \rightarrow \mathbb{R}^{X_f \times \{0,1\}^{S_f}}$, called a **factor**.

The family $H : \Theta \rightarrow \mathbb{R}^{X \times \{0,1\}^S}$ such that

$$\forall \theta \in \Theta \quad \forall x \in X \quad \forall y \in \{0,1\}^S : \quad H_\theta(x, y) = \sum_{f \in F} h_{f\theta}(x_f, y_{S_f}) \quad (2)$$

is called the family of **energy functions** of the conditional graphical model.

Family of Functions

- ▶ We consider a conditional graphical model $(S, F, E, \{X_f\}_{f \in F}, \Theta, \{h_f\}_{f \in F})$ and its family H of energy functions.
- ▶ We assume that Θ is a finite-dimensional, real vector space, i.e., there exists a finite, non-empty set J and $\Theta = \mathbb{R}^J$.
- ▶ We assume that every function h_f is linear in Θ , i.e., for every $f \in F$, there exists a $\varphi_f : X_f \times \{0, 1\}^{S_f} \rightarrow \mathbb{R}^J$ such that for any $x_f \in X_f$, any $y_{S_f} \in \{0, 1\}^{S_f}$ and any $\theta \in \Theta$:

$$h_{f\theta}(x_f, y_{S_f}) = \langle \theta, \varphi_f(x_f, y_{S_f}) \rangle \quad (3)$$

Conditional Graphical Models

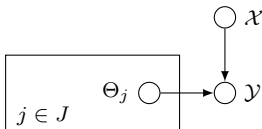
For convenience, we define $\xi : X \times \{0, 1\}^S \rightarrow \mathbb{R}^J$ such that for any $x \in X$ and any $y \in \{0, 1\}^S$:

$$\xi(x, y) = \sum_{f \in F} \varphi_f(x_f, y_{S_f}) \quad (4)$$

Thus, we obtain for any $\theta \in \Theta$, any $x \in X$ and any $y \in Y$:

$$\begin{aligned} H_\theta(x, y) &= \sum_{f \in F} h_{f\theta}(x_f, y_{S_f}) \\ &= \sum_{f \in F} \langle \theta, \varphi_f(x_f, y_{S_f}) \rangle \\ &= \left\langle \theta, \sum_{f \in F} \varphi_f(x_f, y_{S_f}) \right\rangle \\ &= \langle \theta, \xi(x, y) \rangle \end{aligned} \quad (5)$$

Conditional Graphical Models



Probabilistic Model

- ▶ Let \mathcal{X} be a random variable whose value is an element $x \in X$ of the attribute space.
- ▶ Let \mathcal{Y} be a random variable whose value is a combination of decisions $y \in \{0, 1\}^S$
- ▶ For any $j \in J$, let Θ_j a random variable whose value is a parameter $\theta_j \in \mathbb{R}$

Factorization

► We assume:

$$P(\mathcal{X}, \mathcal{Y}, \Theta) = P(\mathcal{Y} | \mathcal{X}, \Theta) P(\mathcal{X}) \prod_{j \in J} P(\Theta_j) \quad (6)$$

► Thus:

$$\begin{aligned} P(\Theta | \mathcal{X}, \mathcal{Y}) &= \frac{P(\mathcal{X}, \mathcal{Y}, \Theta)}{P(\mathcal{X}, \mathcal{Y})} \\ &= \frac{P(\mathcal{Y} | \mathcal{X}, \Theta) P(\mathcal{X}) \prod_{j \in J} P(\Theta_j)}{P(\mathcal{X}, \mathcal{Y})} \\ &\propto P(\mathcal{Y} | \mathcal{X}, \Theta) \prod_{j \in J} P(\Theta_j) \end{aligned} \quad (7)$$

Distributions

Definition. For any conditional graphical model, the **partition function** $Z: X \times \Theta \rightarrow \mathbb{R}$ and **Gibbs distribution** $p: X \times \{0, 1\}^S \times \Theta \rightarrow [0, 1]$ are defined by the forms

$$Z(x, \theta) = \sum_{y \in \{0, 1\}^S} e^{-H_\theta(x, y)} \quad (8)$$

$$p(y, x, \theta) = \frac{1}{Z(x, \theta)} e^{-H_\theta(x, y)} \quad (9)$$

We consider a $\sigma \in \mathbb{R}^+$ and

$$p_{Y|X, \Theta}(y, x, \theta) = \frac{1}{Z(x, \theta)} e^{-H_\theta(x, y)} \quad (10)$$

$$\forall j \in J: \quad p_{\Theta_j}(\theta_j) = \frac{1}{\sigma\sqrt{2\pi}} e^{-\theta_j^2/2\sigma^2} . \quad (11)$$

Conditional Graphical Models

Lemma. Estimating maximally probable parameters θ , given attributes x and decisions y , i.e.,

$$\operatorname{argmax}_{\theta \in \mathbb{R}^J} p_{\Theta|X,Y}(\theta, x, y)$$

is identical to the supervised structured learning problem w.r.t. L , R and λ such that

$$L(H_\theta(x, \cdot), y) = H_\theta(x, y) + \ln Z(x, \theta) \quad (12)$$

$$= H_\theta(x, y) + \ln \sum_{y' \in \{0,1\}^S} e^{-H_\theta(x, y')} \quad (13)$$

$$= \langle \theta, \xi(x, y) \rangle + \ln \sum_{y' \in \{0,1\}^S} e^{-\langle \theta, \xi(x, y') \rangle} \quad (14)$$

$$R(\theta) = \|\theta\|_2^2 \quad (15)$$

$$\lambda = \frac{1}{2\sigma^2} \quad (16)$$

Conditional Graphical Models

Lemma: The first and second partial derivatives of the logarithm of the partition function have the forms

$$\frac{\partial}{\partial \theta_j} \ln Z = \frac{1}{Z(x, \theta)} \sum_{y' \in \{0,1\}^S} (-\xi_j(x, y')) e^{-\langle \theta, \xi(x, y') \rangle} \quad (17)$$

$$= \mathbb{E}_{y' \sim p_{\mathcal{Y}|\mathcal{X}, \theta}} (-\xi_j(x, y')) \quad (18)$$

$$\begin{aligned} \frac{\partial^2}{\partial \theta_j \partial \theta_k} \ln Z &= \mathbb{E}_{y' \sim p_{\mathcal{Y}|\mathcal{X}, \theta}} (\xi_j(x, y') \xi_k(x, y')) \\ &\quad - \mathbb{E}_{y' \sim p_{\mathcal{Y}|\mathcal{X}, \theta}} (\xi_j(x, y')) \mathbb{E}_{y' \sim p_{\mathcal{Y}|\mathcal{X}, \theta}} (\xi_k(x, y')) \\ &= \text{COV}_{y' \sim p_{\mathcal{Y}|\mathcal{X}, \theta}} (\xi_j(x, y'), \xi_k(x, y')) \end{aligned} \quad (19)$$

Lemma: Supervised structured learning of a conditional graphical model is a convex optimization problem.

Lemma: Estimating maximally probable decisions y , given attributes x and parameters θ , i.e.

$$\operatorname{argmax}_{y \in \{0,1\}^S} p_{\mathcal{Y}|\mathcal{X},\Theta}(x, y, \theta) \quad (20)$$

is identical to the structured inference problem with $\hat{H}(x, y) = H_{\theta}(x, y)$.

Conditional Graphical Models

Summary. Supervised structured learning of conditional graphical models whose factors are linear functions is a convex optimization problem.

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 $O(2^{|S|})$:

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

Conditional Graphical Models II

Computing the gradient requires that we compute

- ▶ the partition function

$$Z(x, \theta) = \sum_{y' \in \{0,1\}^S} e^{-\langle \theta, \xi(x, y') \rangle} \quad (21)$$

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

$$p_{\mathcal{Y}_{S(f)} | \mathcal{X}, \Theta}(y'_{S(f)} | x, \theta) = \frac{1}{Z(x, \theta)} \sum_{y'_{S \setminus S(f)} \in \{0,1\}^{S \setminus S(f)}} e^{-\langle \theta, \xi(x, y') \rangle} \quad (22)$$

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

$$\sum_{y'_{S(f)} \in \{0,1\}^{S(f)}} \varphi_{fj}(x_f, y'_{S(f)}) p_{\mathcal{Y}_{S(f)} | \mathcal{X}, \Theta}(y'_{S(f)} | x, \theta) . \quad (23)$$

Conditional Graphical Models II

The challenge is to sum the function

$$\psi_{\theta}(x, y') := e^{-\langle \theta, \xi(x, y') \rangle} \quad (24)$$

over assignments of 0 or 1 to linearly many (22) or all (21) variables y' .

Defining

$$\psi_{f\theta}(x_f, y'_{S(f)}) = e^{-\langle \theta, \varphi_f(x_f, y'_{S(f)}) \rangle} \quad (25)$$

we obtain

$$\begin{aligned} \psi_{\theta}(x, y') &= e^{-\langle \theta, \xi(x, y') \rangle} \\ &= e^{-\sum_{f \in F} \langle \theta, \varphi_f(x_f, y_{S(f)}) \rangle} \end{aligned} \quad (26)$$

$$= \prod_{f \in F} e^{-\langle \theta, \varphi_f(x_f, y_{S(f)}) \rangle} \quad (27)$$

$$= \prod_{f \in F} \psi_{f\theta}(x_f, y_{S(f)}) . \quad (28)$$

Thus, the challenge in (22) and (21) is to compute a sum of a product of functions. Specifically:

$$Z(x, \theta) = \sum_{y' \in \{0,1\}^S} \prod_{f \in F} \psi_{f\theta}(x_f, y_{S(f)}) \quad (29)$$

$$p_{\mathcal{Y}_{S(f)} | \mathcal{X}, \Theta}(y'_{S(f)} | x, \theta) = \frac{1}{Z(x, \theta)} \sum_{y'_{S \setminus S(f)} \in \{0,1\}^{S \setminus S(f)}} \prod_{f \in F} \psi_{f\theta}(x_f, y_{S(f)}) \quad (30)$$

- ▶ 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 $s \in S$ and any factor node $f \in F$, the functions

$$m_{s \rightarrow f}, m_{f \rightarrow s} : \{0, 1\} \rightarrow \mathbb{R} , \quad (31)$$

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

$$m_{s \rightarrow f}(y_s) = \prod_{f' \in F(s) \setminus \{f\}} m_{f' \rightarrow s}(y_s) \quad (32)$$

$$m_{f \rightarrow s}(y_s) = \sum_{y_{S(f) \setminus \{s\}}} \psi_{f\theta}(x_f, y_{S(f)}) \prod_{s' \in S(f) \setminus \{s\}} m_{s' \rightarrow f}(y_{s'}) \quad (33)$$

Lemma. If the factor graph is acyclic, messages are defined recursively by (32) and (33), beginning with the messages from leaves. Moreover, for any $s \in S$ and any $f \in F$:

$$Z(x, \theta) = \sum_{y_s \in \{0,1\}} \prod_{f' \in F(s)} m_{f' \rightarrow s}(y_s) \quad (34)$$

$$p_{\mathcal{Y}_{S(f)} | \mathcal{X}, \Theta}(y'_{S(f)} | x, \theta) = \frac{1}{Z(x, \theta)} \psi_{f\theta}(x_f, y_{S(f)}) \prod_{s' \in S(f)} m_{s' \rightarrow f}(y_{s'}) \quad (35)$$

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**.

Conditional Graphical Models III

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

The **inference problem** w.r.t. a **conditional graphical model** has the form of an unconstrained binary optimization problem:

$$\operatorname{argmin}_{y \in \{0,1\}^S} H_\theta(x, y) \quad (36)$$

It is NP-hard. (This can be shown, e.g., by reduction of binary integer programming, which is one of Karp's 21 problems).

Conditional Graphical Models III

We consider transformations that change one decision at a time:

Definition. For any $s \in S$, let $\text{flip}_s: \{0, 1\}^S \rightarrow \{0, 1\}^S$ such that for any $y \in \{0, 1\}^S$ and any $t \in S$:

$$\text{flip}_s[y](t) = \begin{cases} 1 - y_t & \text{if } t = s \\ y_t & \text{otherwise} \end{cases} . \quad (37)$$

The greedy local search algorithm w.r.t these transformations is known as **Iterated Conditional Modes**, or ICM (Besag 1986).

$$y' = \text{icm}(y)$$

$$\text{choose } s \in \underset{s' \in S}{\text{argmin}} H_\theta(x, \text{flip}_{s'}[y]) - H_\theta(x, y)$$

$$\text{if } H_\theta(x, \text{flip}_s[y]) < H_\theta(x, y)$$

$$y' := \text{icm}(\text{flip}_s[y])$$

else

$$y' := y$$

- ▶ The **inference problem** consists in computing the minimum of a sum of functions:

$$\begin{aligned} & \operatorname{argmin}_{y \in \{0,1\}^S} H_\theta(x, y) \\ &= \operatorname{argmin}_{y \in \{0,1\}^S} \sum_{f \in F} h_{f\theta}(x_f, y_{S(f)}) \end{aligned} \quad (38)$$

- ▶ This problem is analogous to that of computing the sum of a product of functions (from the previous lecture) in that both $(\mathbb{R}, \min, +)$ and $(\mathbb{R}, +, \cdot)$ are commutative semi-rings.
- ▶ This analogy is sufficient to transfer the idea of **message passing**, albeit with messages adapted to the $(\mathbb{R}, \min, +)$ semi-ring:

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

$$\mu_{s \rightarrow f}, \mu_{f \rightarrow s} : \{0, 1\} \rightarrow \mathbb{R} , \quad (39)$$

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

$$\mu_{s \rightarrow f}(y_s) = \sum_{f' \in F(s) \setminus \{f\}} \mu_{f' \rightarrow s}(y_s) \quad (40)$$

$$\mu_{f \rightarrow s}(y_s) = \min_{y_{S(f) \setminus \{s\}}} \psi_{f\theta}(x_f, y_{S(f)}) + \sum_{s' \in S(f) \setminus \{s\}} \mu_{s' \rightarrow f}(y_{s'}) \quad (41)$$

Lemma. If the factor graph is acyclic, messages are defined recursively by (40) and (41), beginning with the messages from leaves. Moreover, for any $s \in S$:

$$\begin{aligned}
 & \operatorname{argmin}_{y \in \{0,1\}^S} H_\theta(x, y) \\
 &= \min_{y \in \{0,1\}^S} \sum_{f \in F} h_{f\theta}(x_f, y_{S(f)}) \\
 &= \min_{y_s \in \{0,1\}} \sum_{f' \in F(s)} \mu_{f' \rightarrow s}(y_s)
 \end{aligned} \tag{42}$$

Proof. Analogous to that of Lemma 18 in the lecture notes.

Summary

- ▶ For conditional graphical models whose factor graph is **acyclic**, the inference problem can be solved efficiently by means of **min-sum message passing**.
- ▶ For conditional graphical models whose factor graph is **cyclic**, one local search algorithm for the inference problem is known as **Iterated Conditional Modes (ICM)**.