

ML for UGM:

data $D = \{ \vec{X}^1, \dots, \vec{X}^N \}$

I am using $\vec{}$ on the top to indicate it is a multivariate variable (i.e. $\dim(\vec{X}) > 1$). I use superindex to represent observations. $\vec{X}^1, \dots, \vec{X}^N \stackrel{i.i.d.}{\sim} P(\vec{X}; \theta)$

$$\max_{\theta} \log P(D; \theta) = \max_{\theta} \sum_{n=1}^N \log P(\vec{X}^n; \theta) \quad (*)$$

↑
all Params

$$P(\vec{X}; \theta) = \frac{1}{Z(\theta)} \prod_c \phi_c(X_c; \theta)$$

normalizer factors

$$(*) \Rightarrow \sum_{n=1}^N \sum_c \log \phi_c(X_c^n; \theta) - N \log Z(\theta)$$

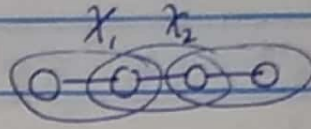
the trouble maker!!!
why?

$$L(\theta) = \sum_{n=1}^N \sum_c \log \phi_c(X_c^n; \theta) - N \log \left(\sum_{\vec{X}} \prod_c \phi_c(X_c; \theta) \right)$$

↙ not dependent on data

Derivative of ML (w.r.t. θ_c):

$$\nabla_{\theta_c} L(\theta) =$$



$$\sum_{n=1}^N \frac{\partial}{\partial \theta_c} \log \phi(x_n; \theta_c) - N \frac{\partial}{\partial \theta_c} \log Z(\theta)$$

$$\frac{\partial}{\partial \theta_c} Z(\theta) = \frac{\partial}{\partial \theta_c} \sum_y \prod_{c' \neq c} \phi(y; \theta_{c'})$$

$$= \sum_y \frac{\partial \phi_c(y; \theta_c)}{\partial \theta_c} \prod_{c' \neq c} \phi_{c'}(y; \theta_{c'})$$

Mult/divide $\phi_c(y; \theta_c)$

$$= \sum_y \frac{\partial \phi_c(y; \theta_c)}{\partial \theta_c} \cdot \frac{1}{\phi_c(y; \theta_c)} \prod_{c' \neq c} \phi_{c'}(y; \theta_{c'})$$

Mult/divide $Z(\theta)$

$$= Z(\theta) \sum_y \frac{\partial \phi_c(y; \theta_c)}{\partial \theta_c} \cdot \frac{1}{\phi_c(y; \theta_c)} \cdot \frac{\prod_{c' \neq c} \phi_{c'}(y; \theta_{c'})}{Z(\theta)}$$

$$= Z(\theta) \underbrace{E_{y \sim P(y)} \left[\frac{\partial}{\partial \theta_c} \log \phi_c(y; \theta_c) \right]}_{\text{marginal}}$$

Substitute it back:

$$\frac{\partial}{\partial \theta} L(\theta) = \underbrace{\sum_{n=1}^N \frac{\partial}{\partial \theta} \log \phi_c(x_c^n; \theta)}_{\text{empirical mean}} - \underbrace{N \mathbb{E}_{y_c \sim P(y)} \left[\frac{\partial \log \phi_c(y; \theta)}{\partial \theta} \right]}_{\text{actual mean}}$$

at the optimum ($\frac{\partial}{\partial \theta} L(\theta) = 0$), the empirical and actual moments match!



For tabular clique Potential:

$$L(\phi) = \sum_n \sum_c \sum_{y_c} \mathbb{I}(y_c = x_c^n) \log \phi_c(y_c) - N \log Z(\phi)$$

$$\nabla_{\phi_c} L(\phi) = \sum_n \mathbb{I}(y_c = x_c^n) \frac{1}{\phi_c(y_c)} - N \nabla_{\phi_c} \log Z(\theta)$$

$$\nabla_{\phi_c} \log Z(\theta) = \nabla_{\phi_c} \log \left(\sum_{y_c} \pi_{c'} \phi_{c'}(y_c) \right)$$

$$= \frac{\sum_{y_c} \pi_{c'} \phi_{c'}(y_c)}{\sum_{y_c} \pi_{c'} \phi_{c'}(y_c)} = \frac{\sum_{y_c} \pi_{c'} \phi_{c'}(y_c)}{Z(\theta) \phi_c(y_c)} = \frac{P(y_c)}{\phi_c(y_c)}$$

$$\nabla_{\phi_c} L(\phi) = \sum_n \mathbb{I}(y_c = x_c^{\wedge}) \frac{1}{\phi_c(y_c)} - N \frac{P(y_c)}{\phi_c(y_c)}$$

Let's introduce this notation for empirical mean:

$$\varepsilon(x_c^{\wedge}) = \frac{1}{N} \sum \mathbb{I}(y_c = x_c^{\wedge})$$

at the optimal param $\nabla_{\phi_c} L(\phi)$:

$$\varepsilon(x_c^{\wedge}) = P(y_c)$$

✘

$$\log P(x; \theta) = \sum_c \theta_c^T f(x_c) - \log Z(\theta)$$

$$\nabla_{\theta_c} (\log P(x; \theta)) = f(x_c) - \frac{\nabla_{\theta_c} Z(\theta)}{Z(\theta)}$$

$$Z(\theta) = e^{A(\theta)}$$

$$\nabla_{\theta_c} Z(\theta) = \frac{\partial A(\theta)}{\partial \theta_c} e^{A(\theta)}$$

$$= \mathbb{E}\left[\frac{\partial f_c(x_c)}{\partial \theta_c}\right] Z(\theta)$$

$$\text{log-likelihood } L(\theta) = \log P(D; \theta) = \sum_{N_c} \sum_{n=1}^N \theta_c^T f(x_c^n) - \log Z(\theta)$$

$$\nabla_{\theta_c} L(\theta) = \sum_{N_c} \sum_{n=1}^N f_c(x_c^n) - \mathbb{E}\left[\frac{\partial f_c(x_c)}{\partial \theta_c}\right]$$

Page 5

$$= \frac{1}{n} \mathbb{E}_{X_e \sim \hat{P}(\cdot; \theta)} [f(X_e)] - \mathbb{E}_{X_e \sim P(\cdot; \theta)} [f(X_e)]$$

~~~~~  
Empirical