

Probabilistic Graphical Models

Lectures 23,24

Learning MRFs

Contrastive Divergence

Learning CRFs

MRF Log-likelihood



K. N. Toosi
University of Technology

$$P_{\theta}(A, B, C) = 1/Z \psi_{\alpha}(A, B) \varphi_{\beta}(B, C), \quad \theta = (\alpha, \beta)$$

$$Z(\theta) = Z(\alpha, \beta) = \sum_A \sum_B \sum_C \psi_{\alpha}(A, B) \varphi_{\beta}(B, C)$$

$$\ell(\alpha, \beta) = \sum_i \log P_{\theta}(A^i, B^i, C^i)$$

$$= \sum_i -\log(Z) + \log(\psi_{\alpha}(A^i, B^i)) + \log(\varphi_{\beta}(B^i, C^i))$$

$$= -m \log(Z) + \sum_i \log(\psi_{\alpha}(A^i, B^i)) + \sum_i \log(\varphi_{\beta}(B^i, C^i))$$

$$= -m \log\left(\sum_A \sum_B \sum_C \psi_{\alpha}(A, B) \varphi_{\beta}(B, C)\right)$$

$$+ \sum_i \log(\psi_{\alpha}(A^i, B^i)) + \sum_i \log(\varphi_{\beta}(B^i, C^i))$$

Data:

$$X^1 = (A^1, B^1, C^1)$$

$$X^2 = (A^2, B^2, C^2)$$

:

$$X^m = (A^m, B^m, C^m)$$

MRF Log-likelihood



K. N. Toosi
University of Technology

$$P_{\theta}(A, B, C) = 1/Z \psi_{\alpha}(A,B) \varphi_{\beta}(B,C), \quad \theta = (\alpha, \beta)$$

$$\ell(\alpha, \beta) = -m \log(Z) + \sum_i \log(\psi_{\alpha}(A^i, B^i)) + \sum_i \log(\varphi_{\beta}(B^i, C^i))$$

$$= \underbrace{-m \log(\sum_A \sum_B \sum_C \psi_{\alpha}(A,B) \varphi_{\beta}(B,C))}_{f(\alpha, \beta)} + \underbrace{\sum_i \log(\psi_{\alpha}(A^i, B^i))}_{g(\alpha, \text{data})} + \underbrace{\sum_i \log(\varphi_{\beta}(B^i, C^i))}_{h(\beta, \text{data})}$$

entangles parameters

MRF log-likelihood - Exponential form



$$P_{\theta}(X) = \frac{1}{Z(\theta)} \tilde{P}_{\theta}(X) \quad X \in \mathbb{R}^n$$

$$P_{\theta}(X) = \frac{1}{Z(\theta)} \tilde{P}_{\theta}(X) = \frac{1}{Z(\theta)} e^{F_{\theta}(X)}$$

CV 23 (I)

Data
 $X^1, X^2, X^3, \dots, X^m$
 $X^i \in \mathbb{R}^n$

$$\ell(\theta) = \log \prod_{i=1}^m P_{\theta}(X^i) = \sum_{i=1}^m \log \frac{1}{Z(\theta)} \tilde{P}_{\theta}(X^i) = \sum_{i=1}^m (-\log Z(\theta) + \log \tilde{P}_{\theta}(X^i))$$

$\log \equiv \ln$

$$= -m \log Z(\theta) + \sum_{i=1}^m \log \tilde{P}_{\theta}(X^i)$$

$$= -m \log Z(\theta) + \sum_{i=1}^m \log \exp(F_{\theta}(X^i))$$

$$\Rightarrow \ell(\theta) = -m \log Z(\theta) + \sum_{i=1}^m F_{\theta}(X^i)$$

MRF log-likelihood - Exponential form



$$\Rightarrow \ell(\theta) = -m \log Z(\theta) + \sum_{i=1}^m F_{\theta}(x^i)$$

$\theta = (\theta_1, \theta_2, \dots, \theta_p)$ parameters

$$\theta^* = \operatorname{argmax} \ell(\theta)$$

$$\frac{\partial \ell(\theta)}{\partial \theta_k} = -m \frac{\partial}{\partial \theta_k} \log Z(\theta) + \sum_{i=1}^m \frac{\partial}{\partial \theta_k} F_{\theta}(x^i)$$

$$= -m \frac{\frac{\partial}{\partial \theta_k} Z(\theta)}{Z(\theta)} + \sum_{i=1}^m \frac{\partial}{\partial \theta_k} F_{\theta}(x^i)$$

$$P_{\theta}(x) = \frac{1}{Z(\theta)} \exp(F_{\theta}(x)) \quad Z(\theta) = \sum_X \exp(F_{\theta}(x))$$

$$\frac{\partial Z(\theta)}{\partial \theta_k} = \sum_X \frac{\partial}{\partial \theta_k} \exp(F_{\theta}(x)) = \sum_X \left[\frac{\partial}{\partial \theta_k} F_{\theta}(x) \right] \exp(F_{\theta}(x))$$

Positive and Negative forces



$$P_{\theta}(X) = \frac{1}{Z(\theta)} \exp(F_{\theta}(X)) \quad Z(\theta) = \sum_X \exp(F_{\theta}(X))$$

$$\frac{\partial Z(\theta)}{\partial \theta_k} = \sum_X \frac{\partial}{\partial \theta_k} \exp(F_{\theta}(X)) = \sum_X \left[\frac{\partial}{\partial \theta_k} F_{\theta}(X) \right] \exp(F_{\theta}(X))$$

$$\begin{aligned} \frac{\partial \ln Z(\theta)}{\partial \theta_k} &= \frac{-1}{Z(\theta)} \sum_X \left[\frac{\partial}{\partial \theta_k} F_{\theta}(X) \right] \exp(F_{\theta}(X)) + \sum_{i=1}^m \frac{\partial}{\partial \theta_k} F_{\theta}(X^i) \\ &= m \left(\frac{1}{m} \sum_{i=1}^m \left[\frac{\partial}{\partial \theta_k} F_{\theta}(X^i) \right] - \sum_X \left[\frac{\partial}{\partial \theta_k} F_{\theta}(X) \right] \frac{1}{Z(\theta)} \exp(F_{\theta}(X)) \right) \\ &= m \left(E_P \left\{ \frac{\partial}{\partial \theta_k} F_{\theta}(X^i) \right\} - E_{P_{\theta}(X)} \left\{ \frac{\partial}{\partial \theta_k} F_{\theta}(X) \right\} \right) \\ &\quad \text{positive phase} \qquad \qquad \qquad \text{negative phase} \end{aligned}$$

Log-linear models



$$P_{\theta}(X) = \frac{1}{z(\theta)} e^{F_{\theta}(X)} \quad z(\theta) = \sum_X e^{F_{\theta}(X)} \quad \text{pgm 23}$$

$$\text{Data} = D = \{X^1, X^2, \dots, X^m\} \quad \theta = (\theta_1, \theta_2, \dots, \theta_p)$$

$$\frac{\partial}{\partial \theta_k} \ell(\theta) = m \left[E_D \left\{ \frac{\partial}{\partial \theta_k} F_{\theta}(X) \right\} - E_{P_{\theta}(X)} \left\{ \frac{\partial}{\partial \theta_k} F_{\theta}(X) \right\} \right]$$

log-linear MRF / No shared parameters

$$\text{MRF: } F_{\theta}(X) = \sum_{j=1}^p \theta_j f_j(X)$$

$$\frac{\partial F_{\theta}(X)}{\partial \theta_k} = \frac{f_k(X)}{e_j}$$

$$\frac{\partial}{\partial \theta_k} \ell(\theta) = m \left[E_D \left\{ \frac{f_k(X)}{e_j} \right\} - E_{P_{\theta}(X)} \left\{ \frac{f_k(X)}{e_j} \right\} \right]$$

log-linear models

A

B

C



K. N. Toosi
University of Technology

log-linear $\phi_{\theta_1}(A, B) = \exp\left(\sum_{k=1}^P w_k f_k(A, B)\right)$ $\theta_1 = (w_1, \dots, w_p)$

$\psi_{\theta_2}(B, C) = \exp\left(\sum_{k=1}^Q u_k g_k(B, C)\right)$ $\theta_2 = (u_1, \dots, u_q)$

$\phi_w(A, B) = \exp(w f(A, B))$

$\psi_u(B, C) = \exp(u g(B, C))$

$P(A, B, C) = \frac{1}{Z(w, u)} e^{w f(A, B) + u g(B, C)}$

$\ell(\theta) = -m \log\left(\sum_A \sum_B \sum_C e^{w f(A, B) + u g(B, C)}\right) + \sum_{i=1}^m w f(A^i, B^i) + \sum_{i=1}^m u g(B^i, C^i)$

$= -m \log\left(\sum_A \sum_B \sum_C e^{w f(A, B) + u g(B, C)}\right) + w \sum_{i=1}^m f(A^i, B^i) + u \sum_{i=1}^m g(B^i, C^i)$

$\frac{\partial \ell(w, u)}{\partial w} = -m \frac{\sum_A \sum_B \sum_C f(A, B) e^{w f(A, B) + u g(B, C)}}{\sum_A \sum_B \sum_C e^{w f(A, B) + u g(B, C)}} + \sum_{i=1}^m f(A^i, B^i)$

$Z(\theta) = Z(w, u)$

sufficient statistics

log-linear models

A

B

C



K. N. Toosi
University of Technology

29 (II)

$$\Rightarrow \frac{\partial \ell(\theta)}{\partial w} = m \left(\frac{1}{m} \sum_{i=1}^m f(A^i, B^i) - \sum_A \sum_B \sum_C f(A, B) \frac{1}{z(\theta)} e^{w f(A, B) + u g(B, C)} \right)$$

$$\frac{1}{z(\theta)} \phi_w(A, B) \psi_u(B, C)$$

$$= m \left(\frac{1}{m} \sum_{i=1}^m f(A^i, B^i) - \sum_A \sum_B \sum_C f(A, B) P_{\theta}(A, B, C) \right)$$

$$\frac{\partial \ell(\theta)}{\partial w} = m \left(\mathbb{E}_D \{ f(A^i, B^i) \} - \mathbb{E}_{\theta} \{ f(A, B) \} \right)$$

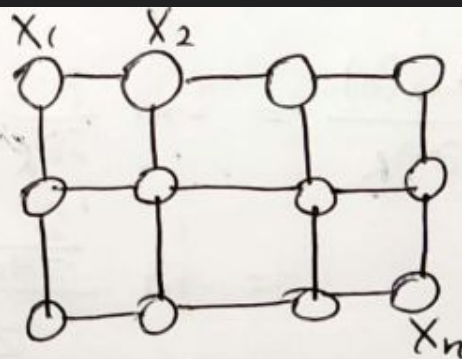
$$\frac{\partial \ell(\theta)}{\partial u} = m \left(\mathbb{E}_D \{ g(B^i, C^i) \} - \mathbb{E}_{\theta} \{ g(B, C) \} \right)$$

Example: Pairwise MRFs



Example: Pairwise MRF
No shared parameters

$$P_{\theta}(X) = \frac{1}{Z} \prod_{i=1}^m \phi_j(X_j) \prod_{(i,j) \in E} \phi_{ij}(X_i, X_j)$$



$$P_{\theta}(X) = \frac{1}{Z(\theta)} \exp\left(\sum w_i f_i(X_i) + \sum w_{ij} f_{ij}(X_i, X_j)\right)$$

$$\theta = (\{w_i\}, \{w_{ij}\})$$

Data:

$$X^1 = (X_1^1, X_2^1, \dots, X_n^1)$$

$$X^2 = (X_1^2, X_2^2, \dots, X_n^2)$$

$$\vdots$$

$$X^m = (X_1^m, X_2^m, \dots, X_n^m)$$

$$\frac{\partial}{\partial w_i} \log \ell(\theta) = m E_D \{ f_i(X_i) \} - m E_{P_{\theta}} \{ f_i(X_i) \}$$

Example: Pairwise MRFs



$$P_{\theta}(X) = \frac{1}{Z(\theta)} \exp\left(\sum w_i f_i(X_i) + \sum w_{ij} f_{ij}(X_i, X_j)\right)$$

$$\theta = (\{w_i\}, \{w_{ij}\})$$

Data:

$$X^1 = (X_1^1, X_2^1, \dots, X_n^1)$$

$$X^2 = (X_1^2, X_2^2, \dots, X_n^2)$$

$$\vdots$$

$$X^m = (X_1^m, X_2^m, \dots, X_n^m)$$

$$\frac{\partial}{\partial w_i} \ell(\theta) = m E_D \{ f_i(X_i) \} - m E_{P_{\theta}} \{ f_i(X_i) \}$$

$$= \sum_{k=1}^m f_i(X_i^k) - m \sum_X P_{\theta}(X_1, \dots, X_m) f_i(X_i)$$

$$= \text{''} - m \sum_{X_1, X_2} \dots \sum_{X_m} P_{\theta}(X_1, X_2, \dots, X_m) f_i(X_i)$$

$$= \text{''} - m \sum_{X_i} \left[\sum_{X_1, X_2} \dots \sum_{X_{i-1}, X_{i+1}, \dots, X_n} P_{\theta}(X_1, \dots, X_m) \right] f_i(X_i)$$

$$= \text{''} - m \sum_{X_i} P_{\theta}(X_i) f_i(X_i) \rightarrow \text{marginal distr}$$

Example: Pairwise MRFs



$$P_{\theta}(X) = \frac{1}{Z(\theta)} \exp\left(\sum_{i=1}^n w_i f_i(X_i) + \sum_{(i,j) \in \mathcal{E}} w_{ij} f_{ij}(X_i, X_j)\right)$$

$$\frac{\partial \ell(\theta)}{\partial w_i} = \sum_{k=1}^m f_i(X_i^k) - \frac{1}{m} \sum_{X_i} P_{\theta}(X_i) f_i(X_i)$$

$$\frac{\partial \ell(\theta)}{\partial w_{ij}} = \sum_{k=1}^m f_{ij}(X_i^k, X_j^k) - \frac{1}{m} \sum_{X_i, X_j} P_{\theta}(X_i, X_j) f_{ij}(X_i, X_j)$$

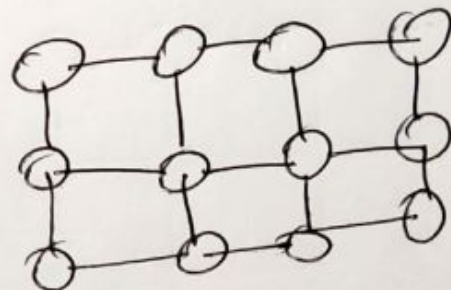
To compute the gradient we need
the marginal distribution over nodes and
edges, i.e. $P_{\theta}(X_1), P_{\theta}(X_2), \dots, P_{\theta}(X_n)$
 $P_{\theta}(X_i, X_j) \quad \forall (i,j) \in \mathcal{E}$

Example: Pairwise MRFs



To compute the gradient we need
the marginal distribution over nodes and
edges, i.e. $P_{\theta}(X_i), P_{\theta}(X_2), \dots, P_{\theta}(X_n)$
 $P_{\theta}(X_i, X_j) \quad \forall (i, j) \in \mathcal{E}$

Need inference
exact, approximate



Optimize Log-linear models



$$P_{\theta}(X) = \frac{1}{Z(\theta)} \exp\left(\sum_{c \in C} \theta_c f(X_c)\right)$$

$$\frac{\partial \log \ell(\theta)}{\partial \theta_c} = \sum_{k=1}^m f(X_c^k) - m \sum_{X_c} P_{\theta}(X_c) f(X_c)$$

need inference

init $\theta = \theta_0$

~~while~~ while not converged

inference \Rightarrow find $P_{\theta}(X_c)$ for all "c"

$$\theta_t = \theta_{t-1} + \frac{\partial \log \ell(\theta)}{\partial \theta} \rightarrow \nabla_{\theta} \log \ell(\theta)$$

Learning MRFs, log-linear models



$$H_{ij} = \frac{\partial^2 \ln(Z)}{\partial \theta_i \partial \theta_j}$$

positive definite $u^T H u > 0$ for all $u \neq 0$



convex function



unique
 \Rightarrow single global
minimum

$-\ln(Z)$

concave



How to compute?

$$E_{\theta} \{f(A, B)\} = \sum_A \sum_B \sum_C f(A, B) P_{\theta}(A, B, C)$$

$$= \sum_A \sum_B f(A, B) \sum_C P_{\theta}(A, B, C)$$

$$= \sum_A \sum_B f(A, B) P(A, B)$$

large no. of variables

\rightarrow marginal distribution


Learning MRFs, log-linear models



How to compute? $-\ln(Z)$ concave $###$

$$\begin{aligned} \mathbb{E}_{\theta} \{f(A, B)\} &= \sum_A \sum_B \sum_C f(A, B) P_{\theta}(A, B, C) \\ &= \sum_A \sum_B f(A, B) \sum_C P_{\theta}(A, B, C) \\ &= \sum_A \sum_B f(A, B) P(A, B) \end{aligned}$$

large no. of variables \rightarrow marginal distribution

$$\begin{aligned} \mathbb{E}_{\theta} \{f(A, B)\} &= \sum_A \sum_B \sum_{C, \dots, Z} f(A, B) P_{\theta}(A, B, \dots, Z) \\ &= \sum_A \sum_B f(A, B) P(A, B) \end{aligned}$$


\rightarrow compute using inference (exact (Approximate))

- \downarrow VE, Junction tree
- \rightarrow Loopy BP / Variational inference
- Sample-Based (MCMC)

Learning MRFs - General Case



MRF

$$P_{\theta}(X) = P_{\theta}(X_1, \dots, X_n) = \frac{1}{Z(\theta)} \prod_{c \in \mathcal{C}} \phi_c^{\theta}(X_c)$$

$$Z(\theta) = \sum_{X_1} \sum_{X_2} \dots \sum_{X_n} \prod_{c \in \mathcal{C}} \phi_c^{\theta}(X_c)$$

$$\ell(\theta) = -m \log Z(\theta) + \sum_{c \in \mathcal{C}} \sum_{i=1}^m \log \phi_c^{\theta}(X_c^i)$$

Challenges

* $Z(\theta)$ couples the parameters

* $Z(\theta)$ is a sum of an exponentially

many terms,
large no. of

Recap



K. N. Toosi
University of Technology

$$p_{\theta}(X) = \frac{1}{Z(\theta)} \exp(F(X, \theta)) \quad \mathcal{D} = \{X^1, X^2, \dots, X^m\}$$

$$\ell(\theta) = \sum_{i=1}^m F(X^i, \theta) - m \ln(Z(\theta))$$

Recap



K. N. Toosi

$$p_{\theta}(\mathbf{X}) = \frac{1}{Z(\theta)} \exp(F(\mathbf{X}, \theta)) \quad \mathcal{D} = \{\mathbf{X}^1, \mathbf{X}^2, \dots, \mathbf{X}^m\}$$

$$\ell(\theta) = \sum_{i=1}^m F(\mathbf{X}^i, \theta) - m \ln(Z(\theta))$$

$$\begin{aligned} \frac{\partial}{\partial \theta_j} \ell(\theta) &= \sum_{i=1}^m \frac{\partial}{\partial \theta_j} F(\mathbf{X}^i, \theta) - m \sum_{\mathbf{X}} p_{\theta}(\mathbf{X}) \frac{\partial}{\partial \theta_j} F(\mathbf{X}, \theta) \\ &= m \mathbb{E}_{\mathcal{D}} \left\{ \frac{\partial}{\partial \theta_j} F(\mathbf{X}, \theta) \right\} - m \mathbb{E}_{p_{\theta}(\mathbf{X})} \left\{ \frac{\partial}{\partial \theta_j} F(\mathbf{X}, \theta) \right\} \end{aligned}$$

Recap



$$F(\mathbf{X}, \boldsymbol{\theta}) = \sum_k \theta_k f_k(\mathbf{X}_{c_k}) \quad \mathbf{X}_{c_k} \subset \mathbf{X}$$

$$\begin{aligned} \frac{\partial}{\partial \theta_j} \ell(\boldsymbol{\theta}) &= \sum_{i=1}^m \frac{\partial}{\partial \theta_j} F(\mathbf{X}^i, \boldsymbol{\theta}) - m \sum_{\mathbf{X}} p_{\boldsymbol{\theta}}(\mathbf{X}) \frac{\partial}{\partial \theta_j} F(\mathbf{X}, \boldsymbol{\theta}) \\ &= m \mathbb{E}_{\mathcal{D}} \left\{ \frac{\partial}{\partial \theta_j} F(\mathbf{X}, \boldsymbol{\theta}) \right\} - m \mathbb{E}_{p_{\boldsymbol{\theta}}(\mathbf{X})} \left\{ \frac{\partial}{\partial \theta_j} F(\mathbf{X}, \boldsymbol{\theta}) \right\} \\ &= m \mathbb{E}_{\mathcal{D}} \{ f_j(\mathbf{X}_{c_j}) \} - m \mathbb{E}_{p_{\boldsymbol{\theta}}(\mathbf{X})} \{ f_j(\mathbf{X}_{c_j}) \} \\ &= \sum_{i=1}^m f_j(\mathbf{X}_{c_j}^i) - m \sum_{\mathbf{X}} p_{\boldsymbol{\theta}}(\mathbf{X}) f_j(\mathbf{X}_{c_j}) \\ &= \sum_{i=1}^m f_j(\mathbf{X}_{c_j}^i) - m \sum_{\mathbf{X}_{c_j}} p_{\boldsymbol{\theta}}(\mathbf{X}_{c_j}) f_j(\mathbf{X}_{c_j}) \end{aligned}$$

Gradient Ascent



OOSI
Technology

$$P(X) = \frac{1}{Z(\theta)} e^{\frac{F(X, \theta)}{Z(\theta)}}$$

$$\frac{\partial}{\partial \theta_j} \log \ell(\theta)$$

$$\begin{bmatrix} \frac{\partial}{\partial \theta_1} \log \ell(\theta) \\ \frac{\partial}{\partial \theta_2} \log \ell(\theta) \\ \vdots \\ \frac{\partial}{\partial \theta_p} \log \ell(\theta) \end{bmatrix} \Rightarrow \nabla \log \ell(\theta)$$

$$= \frac{\partial}{\partial \theta} \log \ell(\theta)$$

gradient ascent

$$\frac{\partial}{\partial \theta_j} \log \ell(\theta) = m E_P \left\{ \frac{\partial}{\partial \theta_j} F(X, \theta) \right\} - m E_{P_\theta(X)} \left\{ \frac{\partial}{\partial \theta_j} F(X, \theta) \right\}$$

المراتب

Gradient Ascent



K. N. Toosi
University of Technology

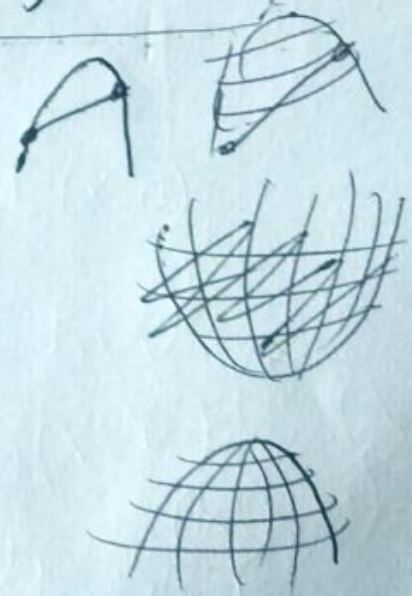
$$\frac{\partial}{\partial \theta_j} \ell(\theta) = m E_P \left\{ \frac{\partial}{\partial \theta_j} F(X, \theta) \right\} - m E_{P_\theta(X)} \left\{ \frac{\partial}{\partial \theta_j} F(X, \theta) \right\}$$

$t \leftarrow 0$
 θ_0

while not converged

$$\theta_{t+1} = \theta_t + \lambda \frac{\partial}{\partial \theta} \ell(\theta) \Big|_{\theta = \theta_t}$$

$t \leftarrow t+1$



Learning By Sampling



problem: How to compute $E_{P_{\theta}(X)} \left\{ \frac{\partial}{\partial \theta_j} F(X, \theta) \right\}$

take samples X^1, X^2, \dots, X^M from $P_{\theta}(X) = \frac{1}{Z(\theta)} e^{F(X, \theta)}$

$$E_{P_{\theta}(X)} \left\{ \frac{\partial}{\partial \theta_j} F(X, \theta) \right\} \approx \frac{1}{M} \sum_{i=1}^M \frac{\partial}{\partial \theta_j} F(X^i, \theta)$$

$$\frac{\partial}{\partial \theta_j} \ell(\theta) = \sum_{i=1}^m \frac{\partial}{\partial \theta_j} F(X^i, \theta) - \frac{m}{M} \sum_{i=1}^m \frac{\partial}{\partial \theta_j} F(X^i, \theta)$$

training data X^1, X^2, \dots, X^m

samples from $P_{\theta}(X)$ X^1, X^2, \dots, X^m

usually
 $M=m$

$$\frac{\partial}{\partial \theta_j} \ell(\theta) = \underbrace{\sum_{i=1}^m \frac{\partial F}{\partial \theta_j} F(X^i, \theta)}_{\text{positive force}} - \underbrace{\sum_{i=1}^m \frac{\partial F}{\partial \theta_j} F(X^i, \theta)}_{\text{negative force}}$$

positive force

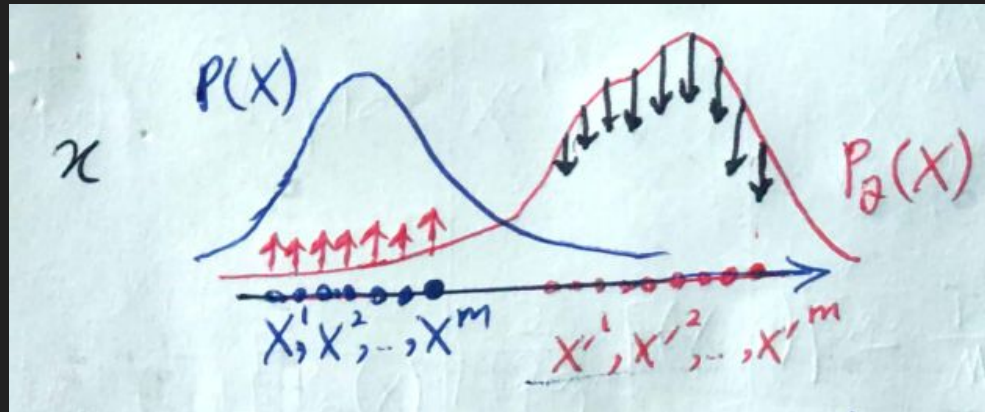
negative force

Positive and Negative Forces



M. N. Toosi
University of Technology

$$\frac{\partial}{\partial \theta_j} \ell(\theta) = \underbrace{\sum_{i=1}^m \frac{\partial F}{\partial \theta_j} F(X^i, \theta)}_{\text{positive force}} - \underbrace{\sum_{i=1}^m \frac{\partial F}{\partial \theta_j} F(X^i, \theta)}_{\text{negative force}}$$



MCMC-based learning



while not converged **do**

Sample a minibatch of m examples $\{\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(m)}\}$ from the training set

$$\mathbf{g} \leftarrow \frac{1}{m} \sum_{i=1}^m \nabla_{\boldsymbol{\theta}} \log \tilde{p}(\mathbf{x}^{(i)}; \boldsymbol{\theta}).$$

Initialize a set of m samples $\{\tilde{\mathbf{x}}^{(1)}, \dots, \tilde{\mathbf{x}}^{(m)}\}$ to random values (e.g., from a uniform or normal distribution, or possibly a distribution with marginals matched to the model's marginals).

for $i = 1$ to k **do**

for $j = 1$ to m **do**

$$\tilde{\mathbf{x}}^{(j)} \leftarrow \text{gibbs_update}(\tilde{\mathbf{x}}^{(j)}).$$

end for

end for

$$\mathbf{g} \leftarrow \mathbf{g} - \frac{1}{m} \sum_{i=1}^m \nabla_{\boldsymbol{\theta}} \log \tilde{p}(\tilde{\mathbf{x}}^{(i)}; \boldsymbol{\theta}).$$

$$\boldsymbol{\theta} \leftarrow \boldsymbol{\theta} + \epsilon \mathbf{g}.$$

end while

Goodfellow et al. "Deep Learning," MIT Press, 2016, Chapter

Contrastive Divergence



```
while not converged do
```

```
  Sample a minibatch of  $m$  examples  $\{\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(m)}\}$  from the training set
```

```
   $\mathbf{g} \leftarrow \frac{1}{m} \sum_{i=1}^m \nabla_{\boldsymbol{\theta}} \log \tilde{p}(\mathbf{x}^{(i)}; \boldsymbol{\theta})$ .
```

```
  for  $i = 1$  to  $m$  do
```

```
     $\tilde{\mathbf{x}}^{(i)} \leftarrow \mathbf{x}^{(i)}$ .
```

```
  end for
```

```
  for  $i = 1$  to  $k$  do
```

```
    for  $j = 1$  to  $m$  do
```

```
       $\tilde{\mathbf{x}}^{(j)} \leftarrow \text{gibbs\_update}(\tilde{\mathbf{x}}^{(j)})$ .
```

```
    end for
```

```
  end for
```

```
   $\mathbf{g} \leftarrow \mathbf{g} - \frac{1}{m} \sum_{i=1}^m \nabla_{\boldsymbol{\theta}} \log \tilde{p}(\tilde{\mathbf{x}}^{(i)}; \boldsymbol{\theta})$ .
```

```
   $\boldsymbol{\theta} \leftarrow \boldsymbol{\theta} + \epsilon \mathbf{g}$ .
```

```
end while
```

Goodfellow et al. "Deep Learning," MIT Press, 2016, Chapter

Persistent Contrastive Divergence



Initialize a set of m samples $\{\tilde{\mathbf{x}}^{(1)}, \dots, \tilde{\mathbf{x}}^{(m)}\}$ to random values (e.g., from a uniform or normal distribution, or possibly a distribution with marginals matched to the model's marginals).

while not converged **do**

Sample a minibatch of m examples $\{\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(m)}\}$ from the training set

$\mathbf{g} \leftarrow \frac{1}{m} \sum_{i=1}^m \nabla_{\boldsymbol{\theta}} \log \tilde{p}(\mathbf{x}^{(i)}; \boldsymbol{\theta})$.

for $i = 1$ to k **do**

for $j = 1$ to m **do**

$\tilde{\mathbf{x}}^{(j)} \leftarrow \text{gibbs_update}(\tilde{\mathbf{x}}^{(j)})$.

end for

end for

$\mathbf{g} \leftarrow \mathbf{g} - \frac{1}{m} \sum_{i=1}^m \nabla_{\boldsymbol{\theta}} \log \tilde{p}(\tilde{\mathbf{x}}^{(i)}; \boldsymbol{\theta})$.

$\boldsymbol{\theta} \leftarrow \boldsymbol{\theta} + \epsilon \mathbf{g}$.

end while

Learning CRFs - Conditional Likelihood



K. N. Toosi
log

$$\text{CRF} = P_{\theta}(Y|X) = \frac{1}{Z(\theta, X)} e^{F_{\theta}(X, Y)}$$

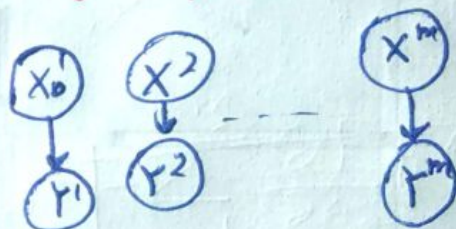
Conditional likelihood

Data
 $(X^1, Y^1), (X^2, Y^2), \dots, (X^m, Y^m)$

$$\text{cLL}(\theta) = \Pr(Y^1, Y^2, \dots, Y^m | X^1, X^2, \dots, X^m)$$

$$\prod_{i=1}^m \Pr(Y^i | X^1, \dots, X^m)$$

$$\prod_{i=1}^m \Pr(Y^i | X^i) \Rightarrow \text{cLL}(\theta) = \prod_{i=1}^m P_{\theta}(Y^i | X^i)$$



Learning CRFs



$$\begin{aligned} \text{c}ll(\theta) &= \log cl(\theta) = \sum_{i=1}^m \log P_{\theta}(Y^i | X^i) \\ &= \sum_{i=1}^m \log \frac{1}{Z_{\theta}(X^i)} e^{F_{\theta}(X^i, Y^i)} \\ &= -\sum_{i=1}^m \log Z_{\theta}(X^i) + \sum_{i=1}^m F_{\theta}(X^i, Y^i) \end{aligned}$$

$$\begin{aligned} \frac{\partial}{\partial \theta_j} \text{c}ll(\theta) &= -\sum_{i=1}^m \sum_Y P_{\theta}(Y | X^i) F_{\theta}(X^i, Y) + \sum_{i=1}^m F_{\theta}(X^i, Y^i) \\ &= -\sum_{i=1}^m \underbrace{E_{P_{\theta}(Y | X^i)} \{ F_{\theta}(X^i, Y) \}}_{\text{Needs inference per } X^i} + \sum_{i=1}^m E_D \{ F_{\theta}(X, Y) \} \end{aligned}$$

Learning CRFs



Conditional log-likelihood

$$\begin{aligned} \ell(\theta) &= \sum_{i=1}^m \log P_{\theta}(Y^i | X^i) = \sum_{i=1}^m \left(-\log Z_{\theta}(X^i) + \sum_{c \in C} \log \Phi_c^{\theta}(Y_c^i, X^i) \right) \\ &= -\sum_{i=1}^m \log Z_{\theta}(X^i) + \sum_{c \in C} \sum_{i=1}^m \log \Phi_c^{\theta}(Y_c^i, X^i) \end{aligned}$$

Log-linear models $\theta = (w_1, w_2, \dots)$

$$P_{\theta}(Y | X) = \frac{1}{Z(\theta, X)} \exp \left(\sum_{j \in J} w_j f_j(X, Y) \right)$$

$Y = \{Y_i\}_{i \in J}$

$$\ell(\theta) = -\sum_{i=1}^m \log Z(\theta, X^i) + \sum_{j \in J} w_j \sum_{i=1}^m f_j(X^i, Y^i)$$

Learning CRFs



K. N. Toosi
University of Technology

$$\ell(\theta) = -\sum_{i=1}^m \log z(\theta, X^i) + \sum_{j \in J} w_j \sum_{i=1}^m f_j(X^i, Y^i)$$

$$\frac{\partial}{\partial w_k} \ell(\theta) = -\sum_{i=1}^m \sum_Y \left(f_k(X^i, Y) P(Y|X^i) + \sum_{i=1}^m f_k(X^i, Y^i) \right)$$

$$= -\sum_{i=1}^m \mathbb{E}_{P_{\theta}(Y|X^i)} \left\{ f_k(X^i, Y) \right\} + m \mathbb{E}_D \left\{ f_k(X^i, Y^i) \right\}$$

↑ interaction

Example

$$\sum_Y f_k(X^i, Y_3, Y_5) P(Y|X^i) = \sum_{Y_3} \sum_{Y_5} f_k(X^i, Y_3, Y_5) P(Y_3, Y_5 | X^i)$$

Learning with Shared Parameters



K. N. Toosi
University of Technology

$$F(\theta) = \theta \left\{ f_1(x_1, x_2) + f_2(x_2, x_3) + f_3(x_3, x_4) \right\}$$
$$E \left\{ \frac{\partial F}{\partial \theta} \right\} = E \left\{ f_1(x_1, x_2) + f_2(x_2, x_3) + f_3(x_3, x_4) \right\}$$
$$E \left\{ f_1(x_1, x_2) \right\} + E \left\{ f_2(x_2, x_3) \right\} + E \left\{ f_3(x_3, x_4) \right\}$$

References



K. N. Toosi
University of Technology

- Goodfellow et al. "Deep Learning," MIT Press, 2016, Chapter 18
- Nowozin, et al. "Structured learning and prediction in computer vision," 2011, Chapter 5
- Koller, Daphne, and Nir Friedman. Probabilistic graphical models: principles and techniques. MIT press, 2009