

Probabilistic Graphical Models

Lecture 6: Learning and Inference. EM Algorithm

Volkan Cevher, Matthias Seeger
Ecole Polytechnique Fédérale de Lausanne

17/10/2011

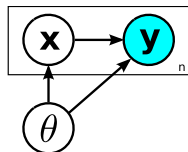


- 1 Learning and Inference
- 2 Bounding with Convexity
- 3 Expectation Maximization
- 4 Learning Markov Random Fields. Log-Linear Models

Learning is Estimation

$$P(\mathbf{y}, \mathbf{x}, \theta) = P(\mathbf{y}|\mathbf{x}, \theta)P(\mathbf{x}|\theta)P(\theta)$$

- \mathbf{y} Observed variable
- \mathbf{x} Latent variable (nuisance)
- θ Parameters (latent query, “higher up”)



- Learning: What is a (single) good value for θ ?
 For which θ does model fit data $D = \{\mathbf{y}_1, \dots, \mathbf{y}_n\}$? \Rightarrow [argmax...]
Estimation

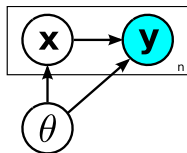
Learning is Estimation

$$P(\mathbf{y}, \mathbf{x}, \theta) = P(\mathbf{y}|\mathbf{x}, \theta)P(\mathbf{x}|\theta)P(\theta)$$

\mathbf{y} Observed variable

\mathbf{x} Latent variable (nuisance)

θ Parameters (latent query, “higher up”)



- **Learning:** What is a (single) good value for θ ? [argmax...]

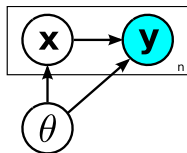
For which θ does model fit data $D = \{\mathbf{y}_1, \dots, \mathbf{y}_n\}$? \Rightarrow **Estimation**

- Maximum likelihood (ML) estimation: $\hat{\theta} = \operatorname{argmax} \log P(D|\theta)$
- Maximum a posteriori (MAP) estimation:
 $\hat{\theta} = \operatorname{argmax} \log P(\theta|D) = \operatorname{argmax}(\log P(D|\theta) + \log P(\theta))$

Learning is Estimation

$$P(\mathbf{y}, \mathbf{x}, \theta) = P(\mathbf{y}|\mathbf{x}, \theta)P(\mathbf{x}|\theta)P(\theta)$$

- \mathbf{y} Observed variable
- \mathbf{x} Latent variable (nuisance)
- θ Parameters (latent query, “higher up”)



- **Learning:** What is a (single) good value for θ ? [argmax ...]
 For which θ does model fit data $D = \{\mathbf{y}_1, \dots, \mathbf{y}_n\}$? \Rightarrow **Estimation**
 - Maximum likelihood (ML) estimation: $\hat{\theta} = \operatorname{argmax} \log P(D|\theta)$
 - Maximum a posteriori (MAP) estimation:
 $\hat{\theta} = \operatorname{argmax} \log P(\theta|D) = \operatorname{argmax}(\log P(D|\theta) + \log P(\theta))$
- **Inference:** What is posterior $P(\theta|D)$? [∫ ...]
 - Range / shape of “good values” mass
 - Uncertainty in estimates

Terms like “MAP inference”: Just wrong

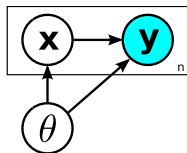
Learning is Estimation

$$P(\mathbf{y}, \mathbf{x}, \theta) = P(\mathbf{y}|\mathbf{x}, \theta)P(\mathbf{x}|\theta)P(\theta)$$

\mathbf{y} Observed variable

\mathbf{x} Latent variable (nuisance)

θ Parameters (latent query, “higher up”)



- **Learning:** What is a (single) good value for θ ? [argmax ...]

For which θ does model fit data $D = \{\mathbf{y}_1, \dots, \mathbf{y}_n\}$? \Rightarrow **Estimation**

- Maximum likelihood (ML) estimation: $\hat{\theta} = \operatorname{argmax} \log P(D|\theta)$
- Maximum a posteriori (MAP) estimation:
 $\hat{\theta} = \operatorname{argmax} \log P(\theta|D) = \operatorname{argmax}(\log P(D|\theta) + \log P(\theta))$

- **Inference:** What is posterior $P(\theta|D)$? [$\int \dots$]

- Range / shape of “good values” mass
- Uncertainty in estimates

Terms like “MAP inference”: Just wrong

- Still: We’ll need inference [$P(\mathbf{x}|\mathbf{y}, \theta)$] for learning

Why Learning can be Hard

- Data from $N(\mathbf{y}|\boldsymbol{\mu}, \sigma^2\mathbf{I})$. Learn mean $\boldsymbol{\mu} \Rightarrow$ That's not hard. Why?

Why Learning can be Hard

- Data from $N(\mathbf{y}|\boldsymbol{\mu}, \sigma^2 \mathbf{I})$. Learn mean $\boldsymbol{\mu} \Rightarrow$ That's not hard. Why?
 - Model directed graph. CPTs: Nice form (Gaussian)
 - No latent variables except parameters

\Rightarrow No inference required
- Learning gets hard if you need inference: Nice-form distributions become nasty through marginalization
 - Latent nuisance variables
 - Undirected models (MRFs)

Why Learning can be Hard

- Data from $N(\mathbf{y}|\boldsymbol{\mu}, \sigma^2 \mathbf{I})$. Learn mean $\boldsymbol{\mu} \Rightarrow$ That's not hard. Why?
 - Model directed graph. CPTs: Nice form (Gaussian)
 - No latent variables except parameters

\Rightarrow No inference required
- Learning gets hard if you need inference: Nice-form distributions become nasty through marginalization
 - Latent nuisance variables
 - Undirected models (MRFs)
- Marginalization creates **log partition functions**

Bayes net, latent variables

$$\underbrace{\log P(\mathbf{y}|\boldsymbol{\theta})}_{\text{coupled}} = \log \int \underbrace{P(\mathbf{y}, \mathbf{x}|\boldsymbol{\theta})}_{\text{decoupled}} d\mathbf{x}$$

Markov random field

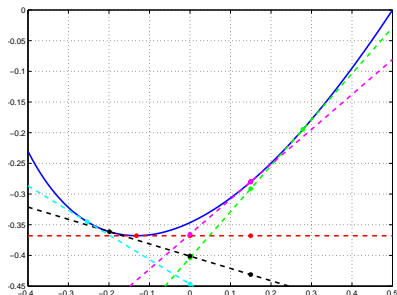
$$\log Z = \log \sum_{\mathbf{x}} \prod_j \Phi_j(\mathbf{x}_{C_j})$$

\Rightarrow Optimization of log partition functions needs inference

Convex Functions. Jensen's Inequality

- Convex set \mathcal{C} :
 $\mathbf{x}_1, \mathbf{x}_2 \in \mathcal{C}, \lambda \in [0, 1]$
 $\Rightarrow \lambda \mathbf{x}_1 + (1 - \lambda) \mathbf{x}_2 \in \mathcal{C}$
- Convex function $f : \mathcal{C} \rightarrow \mathbb{R}$:
 $f(\lambda \mathbf{x}_1 + (1 - \lambda) \mathbf{x}_2)$
 $\leq \lambda f(\mathbf{x}_1) + (1 - \lambda) f(\mathbf{x}_2)$
- Same concept: f convex \Leftrightarrow
 $\text{epi}(f) := \{(\mathbf{x}, y) \mid f(\mathbf{x}) \leq y\}$ convex

F3



Dig for yourself about convexity:

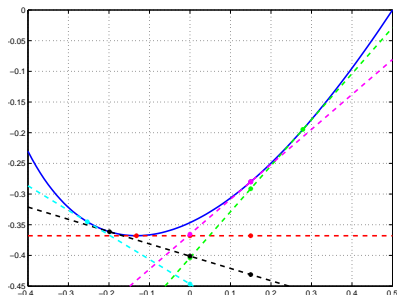
- Boyd, Vandenberghe: Convex Optimization (2004)

[<http://www.stanford.edu/~boyd/cvxbook/>]

Convex Functions. Jensen's Inequality

- Convex set \mathcal{C} :
 $\mathbf{x}_1, \mathbf{x}_2 \in \mathcal{C}, \lambda \in [0, 1]$
 $\Rightarrow \lambda \mathbf{x}_1 + (1 - \lambda) \mathbf{x}_2 \in \mathcal{C}$
- Convex function $f : \mathcal{C} \rightarrow \mathbb{R}$:
 $f(\lambda \mathbf{x}_1 + (1 - \lambda) \mathbf{x}_2)$
 $\leq \lambda f(\mathbf{x}_1) + (1 - \lambda) f(\mathbf{x}_2)$
- Same concept: f convex \Leftrightarrow
 $\text{epi}(f) := \{(\mathbf{x}, y) \mid f(\mathbf{x}) \leq y\}$ convex
- Equivalent: For each $\mathbf{x}_0 \in \mathcal{C}$, there exists \mathbf{u} s.t.

$$f(\mathbf{x}) \geq \mathbf{u}^T (\mathbf{x} - \mathbf{x}_0) + f(\mathbf{x}_0) \quad \text{for all } \mathbf{x} \in \mathcal{C}$$



F3b

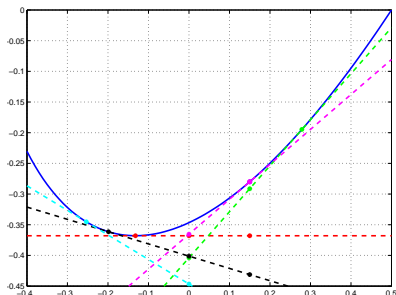
Convex Functions. Jensen's Inequality

- Convex set \mathcal{C} :
 $\mathbf{x}_1, \mathbf{x}_2 \in \mathcal{C}, \lambda \in [0, 1]$
 $\Rightarrow \lambda \mathbf{x}_1 + (1 - \lambda) \mathbf{x}_2 \in \mathcal{C}$
- Convex function $f : \mathcal{C} \rightarrow \mathbb{R}$:
 $f(\lambda \mathbf{x}_1 + (1 - \lambda) \mathbf{x}_2)$
 $\leq \lambda f(\mathbf{x}_1) + (1 - \lambda) f(\mathbf{x}_2)$
- Same concept: f convex \Leftrightarrow
 $\text{epi}(f) := \{(\mathbf{x}, y) \mid f(\mathbf{x}) \leq y\}$ convex
- Equivalent: For each $\mathbf{x}_0 \in \mathcal{C}$, there exists \mathbf{u} s.t.

$$f(\mathbf{x}) \geq \mathbf{u}^T (\mathbf{x} - \mathbf{x}_0) + f(\mathbf{x}_0) \quad \text{for all } \mathbf{x} \in \mathcal{C}$$

- Jensen's inequality: $f : \mathcal{C} \rightarrow \mathbb{R}$ convex, P distribution over \mathcal{C}

$$\mathbb{E}_P[f(\mathbf{x})] \geq f(\mathbb{E}_P[\mathbf{x}])$$



Bounding Log Partition Functions

Recall the problem:

$$\max_{\theta} \log \int \Phi(\mathbf{x}|\theta) d\mathbf{x}, \quad \Phi(\mathbf{x}|\theta) = \prod_j \Phi_j(\mathbf{x}_{C_j}|\theta)$$

- $t \mapsto -\log(t)$ convex function: For positive f :
 $\log \mathbb{E}_Q[f(\mathbf{x})] \geq \mathbb{E}_Q[\log f(\mathbf{x})]$ (by Jensen)

Bounding Log Partition Functions

Recall the problem:

$$\max_{\theta} \log \int \Phi(\mathbf{x}|\theta) d\mathbf{x}, \quad \Phi(\mathbf{x}|\theta) = \prod_j \Phi_j(\mathbf{x}_{C_j}|\theta)$$

- $t \mapsto -\log(t)$ convex function: For positive f :
 $\log \mathbb{E}_Q[f(\mathbf{x})] \geq \mathbb{E}_Q[\log f(\mathbf{x})]$ (by Jensen)
- Variational mean field inequality

F4

$$\log \int \Phi(\mathbf{x}) d\mathbf{x} \geq \sup_Q \mathbb{E}_Q \left[\log \frac{\Phi(\mathbf{x})}{Q(\mathbf{x})} \right]$$

Bounding Log Partition Functions

Recall the problem:

$$\max_{\theta} \log \int \Phi(\mathbf{x}|\theta) d\mathbf{x}, \quad \Phi(\mathbf{x}|\theta) = \prod_j \Phi_j(\mathbf{x}_{C_j}|\theta)$$

- $t \mapsto -\log(t)$ convex function: For positive f :
 $\log \mathbb{E}_Q[f(\mathbf{x})] \geq \mathbb{E}_Q[\log f(\mathbf{x})]$ (by Jensen)
- Variational mean field inequality

$$\log \int \Phi(\mathbf{x}) d\mathbf{x} \geq \sup_Q \mathbb{E}_Q \left[\log \frac{\Phi(\mathbf{x})}{Q(\mathbf{x})} \right]$$

- Pushing the log inside

$$\begin{array}{ll} \log \int \prod_j \dots & \text{Hard, no decoupling} \\ \int \log \prod_j \dots = \int \sum_j \log \dots & \text{Decoupling : Can be much easier} \end{array}$$

Towards Expectation Maximization

$$\log \int \Phi(\mathbf{x}) d\mathbf{x} \geq \sup_Q E_Q \left[\log \frac{\Phi(\mathbf{x})}{Q(\mathbf{x})} \right]$$

Here is a very simple question: What is the best $Q(\mathbf{x})$ I could choose?

F5

Towards Expectation Maximization

$$\log \int \Phi(\mathbf{x}) d\mathbf{x} = \sup_Q \mathbb{E}_Q \left[\log \frac{\Phi(\mathbf{x})}{Q(\mathbf{x})} \right]$$

Here is a very simple question: What is the best $Q(\mathbf{x})$ I could choose?
 $Q(\mathbf{x}) = \Phi(\mathbf{x})/Z$: The posterior in this situation

Expectation Maximization (Full Generality)

F5a

Goal: Maximize $\log \int \Phi(\mathbf{x}|\theta) d\mathbf{x}$. Iterate:

- Expectation (E step): Tightest lower bound
 $Q(\mathbf{x}) \leftarrow \Phi(\mathbf{x}|\theta)/Z(\theta)$
- Maximization (M step): Maximize lower bound
 $\theta \leftarrow \operatorname{argmax} \mathbb{E}_Q[\log \Phi(\mathbf{x}|\theta)]$ for **fixed** Q

EM Algorithm for Gaussian Mixtures

Gaussian mixture model: $P(\mathbf{y}|x) = N(\boldsymbol{\mu}_x, \mathbf{I})$, $P(x = k) = 1/K$

Observed data: $\mathbf{y}_1, \dots, \mathbf{y}_n \in \mathbb{R}^d$

Latent **indicators**: $x_1, \dots, x_n \in \{1, \dots, K\}$

How to find cluster centers $\boldsymbol{\mu}_k$?

EM Algorithm for Gaussian Mixtures

Gaussian mixture model: $P(\mathbf{y}|x) = N(\boldsymbol{\mu}_x, \mathbf{I})$, $P(x = k) = 1/K$

Observed data: $\mathbf{y}_1, \dots, \mathbf{y}_n \in \mathbb{R}^d$

Latent **indicators**: $x_1, \dots, x_n \in \{1, \dots, K\}$

How to find cluster centers $\boldsymbol{\mu}_k$?

Translation

general \rightarrow particular

$\Phi(\mathbf{x}) \rightarrow \prod_i P(\mathbf{y}_i|x_i)P(x_i)$ [joint likelihood]

$\theta \rightarrow \boldsymbol{\mu}_1, \dots, \boldsymbol{\mu}_K$ [cluster centers]

$Q(\mathbf{x}) \rightarrow Q(\mathbf{x}) = \prod_i Q(x_i)$

$Z(\theta) \rightarrow Z = \prod_i Z_i$, $Z_i = \sum_{x_i} P(\mathbf{y}_i|x_i)P(x_i)$

Note: Decoupling

$$\log \Phi(\mathbf{x}) = \sum_i \log[P(\mathbf{y}_i|x_i)P(x_i)], \quad \log Z = \sum_i \log Z_i$$

EM Algorithm for Gaussian Mixtures

Gaussian mixture model: $P(\mathbf{y}|x) = N(\boldsymbol{\mu}_x, \mathbf{I})$, $P(x = k) = 1/K$

Observed data: $\mathbf{y}_1, \dots, \mathbf{y}_n \in \mathbb{R}^d$

Latent **indicators**: $x_1, \dots, x_n \in \{1, \dots, K\}$

How to find cluster centers $\boldsymbol{\mu}_k$?

Iterate:

- 1 **Expectation**: Posterior distribution for each datapoint

$$Q(x_i = k) \leftarrow P(x_i = k | \mathbf{y}_i)$$

- 2 **Maximization**: Posterior average of all datapoints

$$\boldsymbol{\mu}_k \leftarrow n_k^{-1} \sum_i Q(x_i = k) \mathbf{y}_i = \operatorname{argmax} \sum_i Q(x_i = k) \log P(\mathbf{y}_i | x_i = k),$$

$$n_k = \sum_i Q(x_i = k). \text{ Posterior weighted maximum likelihood}$$

EM Algorithm for Bayesian Networks

$$P(\mathbf{x}) = \prod_j P(x_j | \mathbf{x}_{\pi_j}, \theta_j), \quad \pi_j : \text{parents of node } j, \quad j = 1, \dots, J$$

Parameters θ_j : CPT for $x_j | \mathbf{x}_{\pi_j}$. Data: $D = \{\mathbf{x}^{(i)} = (x_j^{(i)})\}$, $i = 1, \dots, n$.
 In each $\mathbf{x}^{(i)}$: Coefficients can be **missing**

- All $\mathbf{x}^{(i)}$ complete: Match CPTs to empirical averages (counts)
 \Rightarrow No EM needed

EM Algorithm for Bayesian Networks

$$P(\mathbf{x}) = \prod_j P(x_j | \mathbf{x}_{\pi_j}, \theta_j), \quad \pi_j : \text{parents of node } j, \quad j = 1, \dots, J$$

Parameters θ_j : CPT for $x_j | \mathbf{x}_{\pi_j}$. Data: $D = \{\mathbf{x}^{(i)} = (x_j^{(i)})\}$, $i = 1, \dots, n$.
In each $\mathbf{x}^{(i)}$: Coefficients can be **missing**

- All $\mathbf{x}^{(i)}$ complete: Match CPTs to empirical averages (counts)
⇒ No EM needed
- Partially observed $\mathbf{x}^{(i)}$: Your exercise sheet!

Log Partition Function: A Closer Look

Log partition function of $P(\mathbf{x}) = \Phi(\mathbf{x})/Z$:

F7

$$\log Z = \log \sum_{\mathbf{x}} e^{\sum_j \psi_j(\mathbf{x}_{C_j})}, \quad \Psi(\mathbf{x}) = \sum_j \psi_j(\mathbf{x}_{C_j}) = \log \Phi(\mathbf{x})$$

Note: Can have $\sum_{\mathbf{x}} \rightarrow \int \dots d\mathbf{x}$

Log Partition Function: A Closer Look

Log partition function of $P(\mathbf{x}) = \Phi(\mathbf{x})/Z$:

$$\log Z = \log \sum_{\mathbf{x}} e^{\sum_j \psi_j(\mathbf{x}_{C_j})}, \quad \Psi(\mathbf{x}) = \sum_j \psi_j(\mathbf{x}_{C_j}) = \log \Phi(\mathbf{x})$$

Note: Can have $\sum_{\mathbf{x}} \rightarrow \int \dots d\mathbf{x}$

① Moment-generating:

F7a

$$\nabla_{\theta} \log Z = E_P[\nabla_{\theta} \Psi(\mathbf{x})] = \sum_j E_P[\nabla_{\theta} \psi_j(\mathbf{x}_{C_j})]$$

Log Partition Function: A Closer Look

Log partition function of $P(\mathbf{x}) = \Phi(\mathbf{x})/Z$:

$$\log Z = \log \sum_{\mathbf{x}} e^{\sum_j \psi_j(\mathbf{x}_{C_j})}, \quad \Psi(\mathbf{x}) = \sum_j \psi_j(\mathbf{x}_{C_j}) = \log \Phi(\mathbf{x})$$

Note: Can have $\sum_{\mathbf{x}} \rightarrow \int \dots d\mathbf{x}$

① Moment-generating:

$$\nabla_{\theta} \log Z = E_P[\nabla_{\theta} \Psi(\mathbf{x})] = \sum_j E_P[\nabla_{\theta} \psi_j(\mathbf{x}_{C_j})]$$

② Convex: $(v_{\mathbf{x}}) \mapsto \log \sum_{\mathbf{x}} e^{v_{\mathbf{x}}}$

F7b

Log Partition Function: A Closer Look

Log partition function of $P(\mathbf{x}) = \Phi(\mathbf{x})/Z$:

$$\log Z = \log \sum_{\mathbf{x}} e^{\sum_j \psi_j(\mathbf{x}_{C_j})}, \quad \Psi(\mathbf{x}) = \sum_j \psi_j(\mathbf{x}_{C_j}) = \log \Phi(\mathbf{x})$$

Note: Can have $\sum_{\mathbf{x}} \rightarrow \int \dots d\mathbf{x}$

① Moment-generating:

$$\nabla_{\theta} \log Z = E_P[\nabla_{\theta} \Psi(\mathbf{x})] = \sum_j E_P[\nabla_{\theta} \psi_j(\mathbf{x}_{C_j})]$$

② Convex: $(v_{\mathbf{x}}) \mapsto \log \sum_{\mathbf{x}} e^{v_{\mathbf{x}}}$

Consequence of (1):

- Computing $\nabla_{\theta} \log Z$: Exactly same as E step (posterior moments over clique marginals)
- Can use any gradient-based optimizer instead of EM

F7c

Learning Markov Random Fields

$$P(\mathbf{x}) = Z^{-1} e^{\Psi(\mathbf{x})}, \quad \Psi(\mathbf{x}) = \sum_j \psi_j(\mathbf{x}_{C_j})$$

Note: All \mathbf{x} observed here $\rightarrow \tilde{\mathbf{x}}$

- Maximum likelihood:

$$\max_{\theta} \log P(\tilde{\mathbf{x}}) = \max_{\theta} (\Psi(\tilde{\mathbf{x}}) - \log Z)$$

Minus $\log Z$: EM won't do

Learning Markov Random Fields

$$P(\mathbf{x}) = Z^{-1} e^{\Psi(\mathbf{x})}, \quad \Psi(\mathbf{x}) = \sum_j \psi_j(\mathbf{x}_{C_j})$$

Note: All \mathbf{x} observed here $\rightarrow \tilde{\mathbf{x}}$

- Maximum likelihood:

$$\max_{\theta} \log P(\tilde{\mathbf{x}}) = \max_{\theta} (\Psi(\tilde{\mathbf{x}}) - \log Z)$$

Minus $\log Z$: EM won't do

- Gradient-based optimization:

$$\nabla_{\theta} \log P(\tilde{\mathbf{x}}) = \sum_j \left(\nabla_{\theta} \psi_j(\tilde{\mathbf{x}}_{C_j}) - E_P[\nabla_{\theta} \psi_j(\mathbf{x}_{C_j})] \right)$$

F8

Learning Markov Random Fields

$$P(\mathbf{x}) = Z^{-1} e^{\Psi(\mathbf{x})}, \quad \Psi(\mathbf{x}) = \sum_j \Psi_j(\mathbf{x}_{C_j})$$

Note: All \mathbf{x} observed here $\rightarrow \tilde{\mathbf{x}}$

- Maximum likelihood:

$$\max_{\theta} \log P(\tilde{\mathbf{x}}) = \max_{\theta} (\Psi(\tilde{\mathbf{x}}) - \log Z)$$

Minus $\log Z$: EM won't do

- Gradient-based optimization:

$$\nabla_{\theta} \log P(\tilde{\mathbf{x}}) = \sum_j \left(\nabla_{\theta} \Psi_j(\tilde{\mathbf{x}}_{C_j}) - E_P[\nabla_{\theta} \Psi_j(\mathbf{x}_{C_j})] \right)$$

- **Log-linear models:** Surprisingly often, $\Psi_j(\mathbf{x}_{C_j}) = \theta^T \mathbf{f}_j(\mathbf{x}_{C_j})$
 - $\nabla_{\theta} \log P(\tilde{\mathbf{x}}) = \sum_j (\mathbf{f}_j(\tilde{\mathbf{x}}_{C_j}) - E_P[\mathbf{f}_j(\mathbf{x}_{C_j})])$
 - **Convex** optimization problem

Iterative Proportional Fitting

Log-linear Markov random field, separable parameters:

$$P(\mathbf{x}) = Z^{-1} e^{\sum_j \theta_j^T \mathbf{f}_j(\mathbf{x}_{C_j})}, \quad \theta = (\theta_1, \theta_2, \dots)$$

- Assume: For each j , $Q(\mathbf{x}_{C_j})$: Can easily find $\Delta\theta_j$ s.t.
 $E_{Q_\Delta}[\mathbf{f}_j(\mathbf{x}_{C_j})] = \mathbf{f}_j(\tilde{\mathbf{x}}_{C_j})$, where $Q_\Delta(\mathbf{x}_{C_j}) \propto Q(\mathbf{x}_{C_j}) e^{(\Delta\theta_j)^T \mathbf{f}_j(\mathbf{x}_{C_j})}$

Iterative Proportional Fitting

Log-linear Markov random field, separable parameters:

$$P(\mathbf{x}) = Z^{-1} e^{\sum_j \theta_j^T \mathbf{f}_j(\mathbf{x}_{C_j})}, \quad \theta = (\theta_1, \theta_2, \dots)$$

- Assume: For each j , $Q(\mathbf{x}_{C_j})$: Can easily find $\Delta\theta_j$ s.t.
 $E_{Q_\Delta}[\mathbf{f}_j(\mathbf{x}_{C_j})] = \mathbf{f}_j(\tilde{\mathbf{x}}_{C_j})$, where $Q_\Delta(\mathbf{x}_{C_j}) \propto Q(\mathbf{x}_{C_j}) e^{(\Delta\theta_j)^T \mathbf{f}_j(\mathbf{x}_{C_j})}$
- Iterative proportional fitting (IPF)**: Iterate
 - Pick some potential j . Determine marginal $P(\mathbf{x}_{C_j})$ (inference)
 - Find $\Delta\theta_j$: $E_{P_\Delta}[\mathbf{f}_j(\mathbf{x}_{C_j})] = \mathbf{f}_j(\tilde{\mathbf{x}}_{C_j})$
 - Update $\theta_j \leftarrow \theta_j + (\Delta\theta_j)$ [Afterwards: $E_P[\mathbf{f}_j(\mathbf{x}_{C_j})] = \mathbf{f}_j(\tilde{\mathbf{x}}_{C_j})$]

Iterative Proportional Fitting

Log-linear Markov random field, separable parameters:

$$P(\mathbf{x}) = Z^{-1} e^{\sum_j \theta_j^T \mathbf{f}_j(\mathbf{x}_{C_j})}, \quad \theta = (\theta_1, \theta_2, \dots)$$

- Assume: For each j , $Q(\mathbf{x}_{C_j})$: Can easily find $\Delta\theta_j$ s.t.
 $E_{Q_\Delta}[\mathbf{f}_j(\mathbf{x}_{C_j})] = \mathbf{f}_j(\tilde{\mathbf{x}}_{C_j})$, where $Q_\Delta(\mathbf{x}_{C_j}) \propto Q(\mathbf{x}_{C_j}) e^{(\Delta\theta_j)^T \mathbf{f}_j(\mathbf{x}_{C_j})}$
- **Iterative proportional fitting** (IPF): Iterate
 - Pick some potential j . Determine marginal $P(\mathbf{x}_{C_j})$ (inference)
 - Find $\Delta\theta_j$: $E_{P_\Delta}[\mathbf{f}_j(\mathbf{x}_{C_j})] = \mathbf{f}_j(\tilde{\mathbf{x}}_{C_j})$
 - Update $\theta_j \leftarrow \theta_j + (\Delta\theta_j)$ [Afterwards: $E_P[\mathbf{f}_j(\mathbf{x}_{C_j})] = \mathbf{f}_j(\tilde{\mathbf{x}}_{C_j})$]
- Coordinate ascent: Simple, other algorithms can be faster
- Requires inference with changing potentials (e.g., belief propagation)
- Problem with general MRFs: Inference hard [not always: CRFs, next lecture]

Examples for Log-Linear Models

$\Psi_j(\mathbf{x}_{C_j}) = \theta^T \mathbf{f}_j(\mathbf{x}_{C_j})$: Seems special ... **No, it's not**: Very common!

- Discrete model, multinomial CPTs.

F10

Examples for Log-Linear Models

$\Psi_j(\mathbf{x}_{C_j}) = \theta^T \mathbf{f}_j(\mathbf{x}_{C_j})$: Seems special ... **No, it's not**: Very common!

- Discrete model, multinomial CPTs. $\pi_k = P(x = k)$

$$\theta_k = \log(\pi_k / \pi_K), \quad \mathbf{f}(x) = (\mathbf{I}_{\{x=k\}}), \quad k = 1, \dots, K-1,$$

$$\Rightarrow P(x) = Z^{-1} e^{\theta^T \mathbf{f}(x)}, \quad Z = 1 / \pi_K$$

Directed Bayesian networks are not log-linear models, but

- Feature based models:
 \mathbf{f}_j indicators for presence / strength of certain features

Examples for Log-Linear Models

$\psi_j(\mathbf{x}_{C_j}) = \theta^T \mathbf{f}_j(\mathbf{x}_{C_j})$: Seems special ... **No, it's not**: Very common!

- Discrete model, multinomial CPTs. $\pi_k = P(x = k)$

$$\theta_k = \log(\pi_k / \pi_K), \quad \mathbf{f}(x) = (\mathbf{I}_{\{x=k\}}), \quad k = 1, \dots, K-1,$$

$$\Rightarrow P(x) = Z^{-1} e^{\theta^T \mathbf{f}(x)}, \quad Z = 1 / \pi_K$$

Directed Bayesian networks are not log-linear models, but

- Feature based models:
 \mathbf{f}_j indicators for presence / strength of certain features
- Gaussian Markov random field
Note: Positive definiteness comes for free ($\log Z < \infty$), does not destroy convexity

F10a

Further Points

- Learning MRF with latent variables:
 - Use EM for latent variables (marginal \rightarrow joint likelihood)
 - Use ∇ -based optimization / IPF for M step (convex optimization)
[No need to maximize, just descent]

Further Points

- Learning MRF with latent variables:
 - Use EM for latent variables (marginal \rightarrow joint likelihood)
 - Use ∇ -based optimization / IPF for M step (convex optimization)
[No need to maximize, just descent]
- Learning with inner maximization (“Viterbi learning”)
 - Sometimes: MAP (argmax) easier than inference (f)
 - Learning with maximization can work well (K-Means, . . .)
 - In most cases: No equivalent guarantees to learning with inference

[Exceptions: Some work by Taskar, Altun, . . .]

Wrap-Up

- Learning requires (marginal) inference in most cases
⇒ Even frequentists need Bayesian inference
- Inequalities from convexity: Underlying very many ideas / methods
- Expectation Maximization: General-purpose algorithm for marginal likelihood maximization
- Log-linear Markov random fields: Learning is convex optimization