

10-708 Probabilistic Graphical Models

Learning Fully Observed Undirected Graphical Models

Readings:

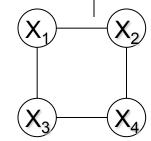
Jordan Chap. 9.3 Jordan Chap. 20 Sha & Pereira (2003) Matt Gormley Lecture 6 February 1, 2016

Housekeeping

Homework 1

- Was released on Saturday 'morning'
- Due date: Feb 15, 12:00 noon





ML Parameter Est. for completely observed MRFs of given structure

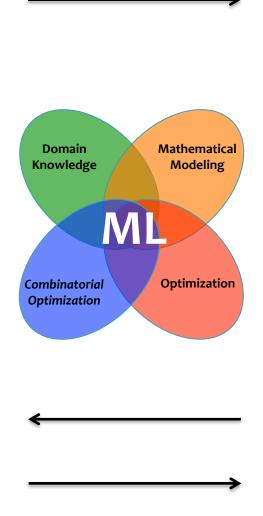
Machine Learning

The data inspires
the structures
we want to
predict

Inference finds

{best structure, marginals, partition function} for a new observation

(Inference is usually called as a subroutine in learning)



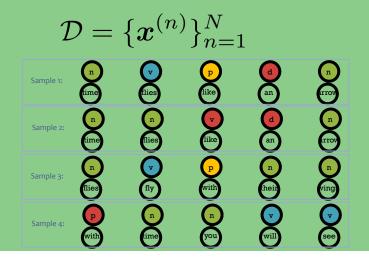
Our **model**defines a score
for each structure

It also tells us what to optimize

Learning tunes the parameters of the model

MLE for Undirected GMs

1. Data



2. Model

$$p(\boldsymbol{x}\mid\boldsymbol{\theta}) = \frac{1}{Z(\boldsymbol{\theta})} \prod_{C\in\mathcal{C}} \psi_C(\boldsymbol{x}_C)$$

3. Objective $_N$

$$\ell(\theta; \mathcal{D}) = \sum_{n=1} \log p(\boldsymbol{x}^{(n)} \mid \boldsymbol{\theta})$$

5. Inference

1. Marginal Inference

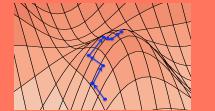
$$p(\boldsymbol{x}_C) = \sum_{\boldsymbol{x}': \boldsymbol{x}_C' = \boldsymbol{x}_C} p(\boldsymbol{x}' \mid \boldsymbol{\theta})$$

2. Partition Function

$$Z(oldsymbol{ heta}) = \sum_{oldsymbol{x}} \prod_{C \in \mathcal{C}} \psi_C(oldsymbol{x}_C)$$

4. Learning

$$\boldsymbol{\theta}^* = \operatorname*{argmax}_{\boldsymbol{\theta}} \ell(\boldsymbol{\theta}; \mathcal{D})$$



1. Data

Given training examples:

2. Model

$$\mathcal{D} = \{oldsymbol{x}^{(n)}\}_{n=1}^N$$



2. Model

Define the model to be an MRF:

$$p(\boldsymbol{x} \mid \boldsymbol{\theta}) = \frac{1}{Z(\boldsymbol{\theta})} \prod_{C \in \mathcal{C}} \psi_C(\boldsymbol{x}_C)$$

$$T_1 \qquad T_2 \qquad T_3 \qquad T_4 \qquad T_5$$

$$W_1 \qquad W_2 \qquad W_3 \qquad W_4 \qquad W_5$$

3. Objective

Choose the objective to be log-likelihood:

(Assign high probability to the things we observe and low probability to everything else)

$$\ell(\theta; \mathcal{D}) = \sum_{n=1}^{N} \log p(\boldsymbol{x}^{(n)} \mid \boldsymbol{\theta})$$

3. Objective

Choose the objective to be log-likelihood:

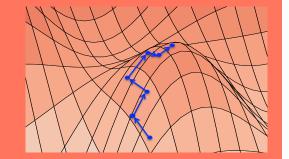
(Assign high probability to the things we observe and low probability to everything else)

$$\ell(\theta; \mathcal{D}) = \sum_{n=1}^{N} \log p(\boldsymbol{x}^{(n)} \mid \boldsymbol{\theta})$$

4. Learning

Tune the parameters to maximize the objective function

$$m{ heta}^* = \operatorname*{argmax}_{m{ heta}} \ell(m{ heta}; \mathcal{D})$$



3. Objective

Choose the objective to be log-likelihood:

(Assign high probability to the things we observe and low probability to everything else) \mathcal{N}

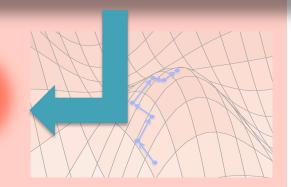
Goals for Today's Lecture

76-1

- Optimize this objective function
- 2. Characterize the applicability of different optimizers

Tune the parameter function

$$oldsymbol{ heta}^* = rgmax \, \ell(oldsymbol{ heta}; \mathcal{D})$$



5. Inference

Three Tasks:

1. Marginal Inference

Compute marginals of variables and cliques

$$p(x_i) = \sum_{\boldsymbol{x}': x_i' = x_i} p(\boldsymbol{x}' \mid \boldsymbol{\theta}) \qquad p(\boldsymbol{x}_C) = \sum_{\boldsymbol{x}': \boldsymbol{x}_C' = \boldsymbol{x}_C} p(\boldsymbol{x}' \mid \boldsymbol{\theta})$$

2. Partition Function

Compute the normalization constant

$$Z(\boldsymbol{\theta}) = \sum_{\boldsymbol{x}} \prod_{C \in \mathcal{C}} \psi_C(\boldsymbol{x}_C)$$

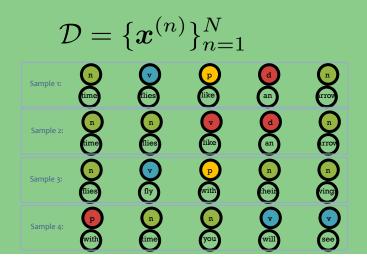
3. MAP Inference

Compute variable assignment with highest probability

$$\hat{\boldsymbol{x}} = \underset{\boldsymbol{x}}{\operatorname{argmax}} p(\boldsymbol{x} \mid \boldsymbol{\theta})$$

MLE for Undirected GMs

1. Data



2. Model

$$p(\boldsymbol{x}\mid\boldsymbol{\theta}) = \frac{1}{Z(\boldsymbol{\theta})} \prod_{C\in\mathcal{C}} \psi_C(\boldsymbol{x}_C)$$

3. Objective $_N$

$$\ell(\theta; \mathcal{D}) = \sum_{n=1} \log p(\boldsymbol{x}^{(n)} \mid \boldsymbol{\theta})$$

5. Inference

1. Marginal Inference

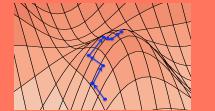
$$p(\boldsymbol{x}_C) = \sum_{\boldsymbol{x}': \boldsymbol{x}_C' = \boldsymbol{x}_C} p(\boldsymbol{x}' \mid \boldsymbol{\theta})$$

2. Partition Function

$$Z(\boldsymbol{ heta}) = \sum_{oldsymbol{x}} \prod_{C \in \mathcal{C}} \psi_C(oldsymbol{x}_C)$$

4. Learning

$$\boldsymbol{\theta}^* = \operatorname*{argmax}_{\boldsymbol{\theta}} \ell(\boldsymbol{\theta}; \mathcal{D})$$



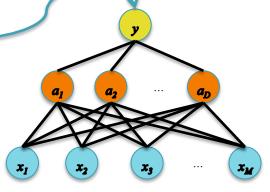
MLE for Undirected GMs

- Today's parameter estimation assumptions:
 - The graphical model structure is given
 - 2. Every variable appears in the training examples

Questions

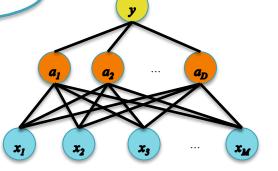
- 1. What does the likelihood objective accomplish?
- 2. Is likelihood the **right objective** function?
- 3. How do we optimize the objective function (i.e. learn)?
- 4. What guarantees does the optimizer provide?
- 5. (What is the mapping from data → model? In what ways can we incorporate our domain knowledge? How does this impact learning?)

- Setting I: $\psi_C(oldsymbol{x}_C) = heta_{C,oldsymbol{x}_C}$
 - A. MLE by inspection (Decomposable Models)
 - B. Iterative Proportional Fitting (IPF)
- Setting II: $\psi_C(m{x}_C) = \exp(m{ heta} \cdot m{f}(m{x}_C))$
 - C. Generalized Iterative Scaling
 - D. Gradient-based Methods
- Setting III: $\psi_C(oldsymbol{x}_C) = 0$
 - D. Gradient-based Methods

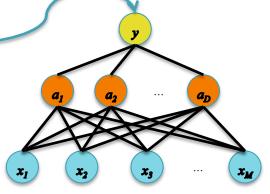


Today's Lecture

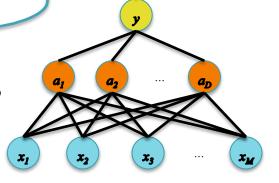
- Setting I: $\psi_C(oldsymbol{x}_C) = heta_{C,oldsymbol{x}_C}$
 - A. MLE by inspection (Decomposable Models)
 - B. Iterative Proportional Fitting (IPF)
- Setting II: $\psi_C(m{x}_C) = \exp(m{ heta} \cdot m{f}(m{x}_C))$
 - C. Generalized Iterative Scaling
 - D. Gradient-based Methods
- Setting III: $\psi_C(oldsymbol{x}_C) = 0$
 - D. Gradient-based Methods



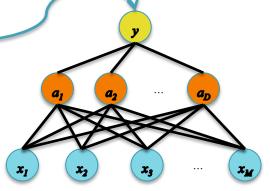
- Setting I: $\psi_C({m x}_C) = heta_{C,{m x}_C}$
 - A. MLE by inspection (Decomposable Models)
 - B. Iterative Proportional Fitting (IPF)
- Setting II: $\psi_C(m{x}_C) = \exp(m{ heta} \cdot m{f}(m{x}_C))$
 - C. Generalized Iterative Scaling
 - D. Gradient-based Methods
- Setting III: $\psi_C(m{x}_C) = 0$
 - D. Gradient-based Methods



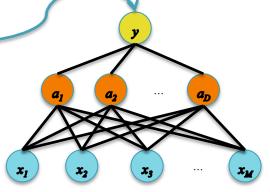
- Setting I: $\psi_C({m x}_C) = heta_{C,{m x}_C}$
 - A. MLE by inspection (Decomposable Models)
 - B. Iterative Proportional Fitting (IPF)
- Setting II: $\psi_C(m{x}_C) = \exp(m{ heta} \cdot m{f}(m{x}_C))$
 - C. Generalized Iterative Scaling
 - D. Gradient-based Methods
- Setting III: $\psi_C(m{x}_C) = 0$
 - D. Gradient-based Methods



- Setting I: $\psi_C(oldsymbol{x}_C) = heta_{C,oldsymbol{x}_C}$
 - A. MLE by inspection (Decomposable Models)
 - B. Iterative Proportional Fitting (IPF)
- Setting II: $\psi_C(m{x}_C) = \exp(m{ heta} \cdot m{f}(m{x}_C))$
 - C. Generalized Iterative Scaling
 - D. Gradient-based Methods
- Setting III: $\psi_C(m{x}_C) = -$
 - D. Gradient-based Methods



- Setting I: $\psi_C({m x}_C) = heta_{C,{m x}_C}$
 - A. MLE by inspection (Decomposable Models)
 - B. Iterative Proportional Fitting (IPF)
- Setting II: $\psi_C(m{x}_C) = \exp(m{ heta} \cdot m{f}(m{x}_C))$
 - C. Generalized Iterative Scaling
 - D. Gradient-based Methods
- Setting III: $\psi_C(m{x}_C) = 0$
 - D. Gradient-based Methods



Whiteboard

- Log-likelihood of MRFs with discrete variables (i.e. tabular clique potentials)
- Derivative of log-likelihood with respect to potentials

Log Likelihood for UGMs with slides from 2015

• Sufficient statistics: for a UGM (V,E), the number of times that a configuration x (i.e., $X_V=x$) is observed in a dataset $D=\{x_1,...,x_N\}$ can be represented as follows:

$$m(\mathbf{x}) = \sum_{n} \delta(\mathbf{x}, \mathbf{x}_{n})$$
 (total count), and $m(\mathbf{x}_{c}) = \sum_{\mathbf{x}_{V \setminus c}} m(\mathbf{x})$ (clique count)

In terms of the counts, the log likelihood is given by:

$$p(D|\theta) = \prod_{n} \prod_{\mathbf{x}} p(\mathbf{x} \mid \theta)^{\delta(\mathbf{x}, \mathbf{x}_n)}$$

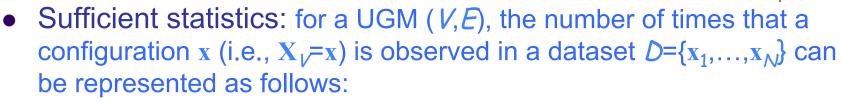
$$\log p(D|\theta) = \sum_{n} \sum_{\mathbf{x}} \delta(\mathbf{x}, \mathbf{x}_n) \log p(\mathbf{x} \mid \theta) = \sum_{\mathbf{x}} \sum_{n} \delta(\mathbf{x}, \mathbf{x}_n) \log p(\mathbf{x} \mid \theta)$$

$$\ell = \sum_{\mathbf{x}} m(\mathbf{x}) \log \left(\frac{1}{Z} \prod_{c} \psi_c(\mathbf{x}_c)\right)$$

$$= \sum_{c} \sum_{\mathbf{x}_c} m(\mathbf{x}_c) \log \psi_c(\mathbf{x}_c) - N \log Z$$

There is a nasty log Z in the likelihood

Log Likelihood for UGMs white stickes from 2015



$$m(\mathbf{x}) = \sum_{n} \delta(\mathbf{x}, \mathbf{x}_n)$$
 (total count), and $m(\mathbf{x}_c) = \sum_{\mathbf{x}_{V \setminus c}} m(\mathbf{x})$ (clique count)

In terms of the counts, the log likelihood is given by:

$$\log p(D|\theta) = \sum_{c} \sum_{\mathbf{x}_{c}} m(\mathbf{x}_{c}) \log \psi_{c}(\mathbf{x}_{c}) - N \log Z$$

There is a nasty log Z in the likelihood



Derivative of log Likelihood

- Log-likelihood: $\ell = \sum \sum m(\mathbf{x}_c) \log \psi_c(\mathbf{x}_c) - \mathcal{N} \log Z$
- $\frac{\partial \ell_1}{\partial \psi_c(\mathbf{x}_c)} = \frac{m(\mathbf{x}_c)}{\psi_c(\mathbf{x}_c)}$ First term:
- Second term: $\frac{\partial \log Z}{\partial \psi_c(\mathbf{x}_c)} = \frac{1}{Z} \frac{\partial}{\partial \psi_c(\mathbf{x}_c)} \left(\sum_{\mathbf{x}} \prod_{d} \psi_d(\widetilde{\mathbf{x}}_d) \right)$

Set the value of variables
$$\mathbf{X}$$
 \mathbf{X} $=\frac{1}{Z}\sum_{\mathbf{X}}\delta(\mathbf{\widetilde{x}}_c,\mathbf{x}_c)\frac{\partial}{\partial\psi_c(\mathbf{x}_c)}\left(\prod_d\psi_d(\mathbf{\widetilde{x}}_d)\right)$ $=\sum_{\mathbf{X}}\delta(\mathbf{\widetilde{x}}_c,\mathbf{x}_c)\frac{1}{\psi_c(\mathbf{\widetilde{x}}_c)}\frac{1}{Z}\prod_d\psi_d(\mathbf{\widetilde{x}}_d)$ $=\frac{1}{\psi_c(\mathbf{x}_c)}\sum_{\mathbf{X}}\delta(\mathbf{\widetilde{x}}_c,\mathbf{x}_c)p(\mathbf{\widetilde{x}})=\frac{p(\mathbf{x}_c)}{\psi_c(\mathbf{x}_c)}$





Derivative of log-likelihood

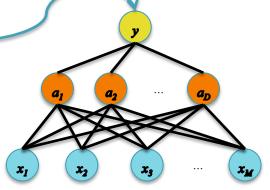
$$\frac{\partial \ell}{\partial \psi_c(\mathbf{x}_c)} = \frac{m(\mathbf{x}_c)}{\psi_c(\mathbf{x}_c)} - \mathcal{N} \frac{p(\mathbf{x}_c)}{\psi_c(\mathbf{x}_c)}$$

Hence, for the maximum likelihood parameters, we know that:

$$p_{MLE}^*(\mathbf{x}_c) = \frac{m(\mathbf{x}_c)}{N}^{\text{def}} = \widetilde{p}(\mathbf{x}_c)$$

- In other words, at the maximum likelihood setting of the parameters, for each clique, the model marginals must be equal to the observed marginals (empirical counts).
- This doesn't tell us how to get the ML parameters, it just gives us a condition that must be satisfied when we have them.

- Setting I: $\psi_C(oldsymbol{x}_C) = heta_{C,oldsymbol{x}_C}$
 - A. MLE by inspection (Decomposable Models)
 - B. Iterative Proportional Fitting (IPF)
- Setting II: $\psi_C(m{x}_C) = \exp(m{ heta} \cdot m{f}(m{x}_C))$
 - C. Generalized Iterative Scaling
 - D. Gradient-based Methods
- Setting III: $\psi_C(m{x}_C) = 0$
 - D. Gradient-based Methods



Whiteboard

- MLE by Guessing
 - Example 1: linear-chain on three variables
 - Example 2: "decomposable" with four variables

MLE by Guessing

 Definition: Graph is decomposable if it can be recursively subdivided into sets A, B, and S such that S separates A and B.

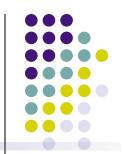
MLE by Guessing

 Definition: Graph is decomposable if it can be recursively subdivided into sets A, B, and S such that S separates A and B.

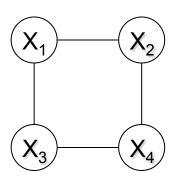
Recipe for MLE by Guessing:

- Three conditions:
 - 1. Graphical model is decomposable
 - 2. Potentials defined on maximal cliques
 - 3. Potentials are are parameterized as: $\psi_C(\boldsymbol{x}_C) = \theta_{C,\boldsymbol{x}_C}$
- Step 1: set each clique potential to its empirical marginal
- Step 2: divide out every non-empty intersection between cliques exactly once

Non-decomposable and/or with non-maximal clique potentials



• If the graph is non-decomposable, and or the potentials are defined on non-maximal cliques (e.g., ψ_{12} , ψ_{34}), we could not equate empirical marginals (or conditionals) to MLE of cliques potentials.



$$p(x_1, x_2, x_3, x_4) = \prod_{\{i, j\}} \psi_{ij}(x_i, x_j)$$

$$\exists (i, j) \quad \text{s.t.} \quad \psi_{ij}^{\text{MLE}}(x_i, x_j) \neq \begin{cases} \widetilde{p}(x_i, x_j) \\ \widetilde{p}(x_i, x_j) / \widetilde{p}(x_i) \\ \widetilde{p}(x_i, x_j) / \widetilde{p}(x_j) \end{cases}$$

- Setting I: $\psi_C(oldsymbol{x}_C) = heta_{C,oldsymbol{x}_C}$
 - A. MLE by inspection (Decomposable Models)
 - B. Iterative Proportional Fitting (IPF)
- Setting II: $\psi_C(m{x}_C) = \exp(m{ heta} \cdot m{f}(m{x}_C))$
 - C. Generalized Iterative Scaling
 - D. Gradient-based Methods
- Setting III: $\psi_C(m{x}_C) = 0$
 - D. Gradient-based Methods

- Fixed point iteration is a general tool for solving systems of equations
- It can also be applied to optimization.

$$J(\boldsymbol{\theta})$$

$$\frac{dJ(\boldsymbol{\theta})}{d\theta_i} = 0 = f(\boldsymbol{\theta})$$

$$0 = f(\boldsymbol{\theta}) \Rightarrow \theta_i = g(\boldsymbol{\theta})$$

$$\theta_i^{(t+1)} = g(\boldsymbol{\theta}^{(t)})$$

Given objective function:

2. Compute derivative, set to zero (call this function f).

Rearrange the equation s.t. one of parameters appears on the LHS.

1. Initialize the parameters.

5. For i in $\{1,...,K\}$, update each parameter and increment t:

6. Repeat #5 until convergence

- Fixed point iteration is a general tool for solving systems of equations
- It can also be applied to optimization.

$$J(x) = \frac{x^3}{3} + \frac{3}{2}x^2 + 2x$$

$$\frac{dJ(x)}{dx} = f(x) = x^2 - 3x + 2 = 0$$

$$\Rightarrow x = \frac{x^2 + 2}{3} = g(x)$$

$$x \leftarrow \frac{x^2 + 2}{3}$$

Given objective function:

Compute derivative, set to zero (call this function f).

Rearrange the equation s.t. one of parameters appears on the LHS.

4. Initialize the parameters.

For i in $\{1,...,K\}$, update each parameter and increment t:

6. Repeat #5 until convergence

We can implement our example in a few lines of python.

$$J(x) = \frac{x^3}{3} + \frac{3}{2}x^2 + 2x$$

$$\frac{dJ(x)}{dx} = f(x) = x^2 - 3x + 2 = 0$$

$$\Rightarrow x = \frac{x^2 + 2}{3} = g(x)$$

$$x \leftarrow \frac{x^2 + 2}{3}$$

```
def f1(x):
    ''''f(x) = x^2 - 3x + 2'''
    return x^{**2} - 3.*x + 2.
def g1(x):
    '''g(x) = \frac{x^2 + 2}{3}'''
    return (x**2 + 2.) / 3.
def fpi(g, x0, n, f):
    '''Optimizes the 1D function g by fixed point iteration
    starting at x0 and stopping after n iterations. Also
    includes an auxiliary function f to test at each value.'''
    x = x0
    for i in range(n):
        print("i=\%2d x=\%.4f f(x)=\%.4f" % (i, x, f(x)))
        x = g(x)
    i += 1
    print("i=%2d x=%.4f f(x)=%.4f" % (i, x, f(x)))
    return x
if __name__ == "__main__":
    x = fpi(g1, 0, 20, f1)
```

$$J(x) = \frac{x^3}{3} + \frac{3}{2}x^2 + 2x$$

$$\frac{dJ(x)}{dx} = f(x) = x^2 - 3x + 2 = 0$$

$$\Rightarrow x = \frac{x^2 + 2}{3} = g(x)$$

$$x \leftarrow \frac{x^2 + 2}{3}$$

```
$ python fixed-point-iteration.py
i = 0 \times 0.0000 f(x) = 2.0000
i = 1 \times -0.6667 f(x) = 0.4444
i = 2 \times 0.8148 f(x) = 0.2195
i = 3 \times 0.8880 f(x) = 0.1246
i = 4 \times -0.9295 f(x) = 0.0755
i = 5 \times 0.9547 f(x) = 0.0474
i = 6 \times 0.9705 f(x) = 0.0304
i = 7 \times 0.9806 f(x) = 0.0198
i = 8 \times 0.9872 f(x) = 0.0130
i = 9 \times -0.9915 f(x) = 0.0086
i=10 x=0.9944 f(x)=0.0057
i=11 x=0.9963 f(x)=0.0038
i=12 x=0.9975 f(x)=0.0025
i=13 x=0.9983 f(x)=0.0017
i=14 \times =0.9989 f(x)=0.0011
i=15 x=0.9993 f(x)=0.0007
i=16 \times 0.9995 f(x)=0.0005
i=17 x=0.9997 f(x)=0.0003
i=18 \times -0.9998 f(x)=0.0002
i=19 \times -0.9999 f(x)=0.0001
i=20 x=0.9999 f(x)=0.0001
```

Iterative Proportional Fitting (IPF)

IPF applies fixed point iteration to the derivative of the likelihood objective

$$\ell(\theta; \mathcal{D}) = \sum_{n=1}^{N} \log p(\boldsymbol{x}^{(n)} \mid \boldsymbol{\theta})$$

$$\frac{\partial \ell}{\partial \psi_c(\mathbf{x}_c)} = \frac{\mathbf{m}(\mathbf{x}_c)}{\psi_c(\mathbf{x}_c)} - \mathcal{N} \frac{\mathbf{p}(\mathbf{x}_c)}{\psi_c(\mathbf{x}_c)}$$

$$\psi_c(\boldsymbol{x}_c) = \psi_c(\boldsymbol{x}_c) \frac{\tilde{p}(\boldsymbol{x}_c)}{p(\boldsymbol{x}_c)}$$

$$\psi_c^{(t+1)}(\mathbf{x}_c) = \psi_c^{(t)}(\mathbf{x}_c) \frac{\widetilde{p}(\mathbf{x}_c)}{p_{\mathbf{x}}^{(t)}(\mathbf{x}_c)}$$

Compute derivative, set to zero

Given likelihood objective

- Rearrange the equation s.t. one of potentials appears on the LHS.
- Initialize the potential tables.
- For each clique c in C, update each potential table and increment t:
- Repeat #5 until convergence

Need to do inference here

$$p(\boldsymbol{x}_C) = \sum_{\boldsymbol{x}': \boldsymbol{x}_C' = \boldsymbol{x}_C} p(\boldsymbol{x}' \mid \boldsymbol{\theta})$$





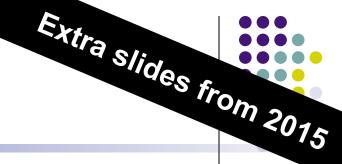
Applies only when potentials are parameterized as:

$$\psi_C(\boldsymbol{x}_C) = \theta_{C,\boldsymbol{x}_C}$$

IPF iterates a set of fixed-point equations:

$$\psi_c^{(t+1)}(\mathbf{X}_c) = \psi_c^{(t)}(\mathbf{X}_c) \frac{\tilde{p}(\mathbf{X}_c)}{p^{(t)}(\mathbf{X}_c)}$$

- However, we can prove it is also a coordinate ascent algorithm (coordinates = parameters of clique potentials).
- Hence at each step, it will increase the log-likelihood, and it will converge to a global maximum.



KL Divergence View

- IPF can be seen as coordinate ascent in the likelihood using the way of expressing likelihoods using KL divergences.
- We can show that maximizing the log likelihood is equivalent to minimizing the KL divergence (cross entropy) from the observed distribution to the model distribution:

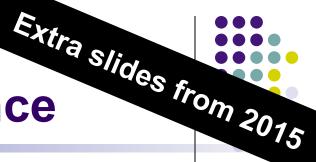
$$\max \ell \Leftrightarrow \min KL(\widetilde{p}(x) || p(x | \theta)) = \sum_{x} \widetilde{p}(x) \log \frac{\widetilde{p}(x)}{p(x | \theta)}$$

• Using a property of KL divergence based on the conditional chain rule: $p(x) = p(x_a)p(x_b|x_a)$:

$$KL(q(x_{a}, x_{b}) \parallel p(x_{a}, x_{b})) = \sum_{x_{a}, x_{b}} q(x_{a})q(x_{b} \mid x_{a}) \log \frac{q(x_{a})q(x_{b} \mid x_{a})}{p(x_{a})p(x_{b} \mid x_{a})}$$

$$= \sum_{x_{a}, x_{b}} q(x_{a})q(x_{b} \mid x_{a}) \log \frac{q(x_{a})}{p(x_{a})} + \sum_{x_{a}, x_{b}} q(x_{a})q(x_{b} \mid x_{a}) \log \frac{q(x_{b} \mid x_{a})}{p(x_{b} \mid x_{a})}$$

$$= KL(q(x_{a}) \parallel p(x_{a})) + \sum_{x_{a}, x_{b}} q(x_{a})KL(q(x_{b} \mid x_{a}) \parallel p(x_{b} \mid x_{a}))$$
© Eric Xing @ CMU, 2005-2015



IPF minimizes KL divergence

Putting things together, we have

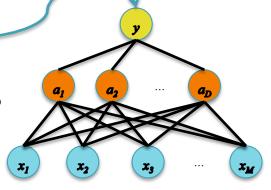
$$KL(\widetilde{p}(\mathbf{x}) \parallel p(\mathbf{x} \mid \theta)) = KL(\widetilde{p}(\mathbf{x}_{c}) \parallel p(\mathbf{x}_{c} \mid \theta)) + \sum_{\mathbf{x}_{c}} \widetilde{p}(\mathbf{x}_{c}) KL(\widetilde{p}(\mathbf{x}_{-c} \mid \mathbf{x}_{c}) \parallel p(\mathbf{x}_{-c} \mid \mathbf{x}_{c}))$$

It can be shown that changing the clique potential ψ_c has no effect on the conditional distribution, so the second term in unaffected.

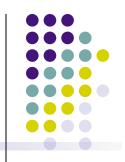
- To minimize the first term, we set the marginal to the observed marginal, just as in IPF.
 - Note that this is only good when the model is decomposable!
- We can interpret IPF updates as retaining the "old" conditional probabilities $p^{(t)}(\mathbf{x}_{-c}|\mathbf{x}_{c})$ while replacing the "old" marginal probability $p^{(t)}(\mathbf{x}_c)$ with the observed marginal $\widetilde{p}(\mathbf{x}_c)$.

Options for MLE of MRFs

- Setting I: $\psi_C(oldsymbol{x}_C) = heta_{C,oldsymbol{x}_C}$
 - A. MLE by inspection (Decomposable Models)
 - B. Iterative Proportional Fitting (IPF)
- Setting II: $\psi_C(m{x}_C) = \exp(m{ heta} \cdot m{f}(m{x}_C))$
 - C. Generalized Iterative Scaling
 - D. Gradient-based Methods
- Setting III: $\psi_C(m{x}_C) = 0$
 - D. Gradient-based Methods



Feature-based Clique Potentials



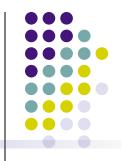
- So far we have discussed the most general form of an undirected graphical model in which cliques are parameterized by general "tabular" potential functions $\psi_c(\mathbf{x}_c)$.
- But for large cliques these general potentials are exponentially costly for inference and have exponential numbers of parameters that we must learn from limited data.
- One solution: change the graphical model to make cliques smaller. But this changes the dependencies, and may force us to make more independence assumptions than we would like.
- Another solution: keep the same graphical model, but use a less general parameterization of the clique potentials.
- This is the idea behind feature-based models.

Features



- Consider a clique x_c of random variables in a UGM, e.g. three consecutive characters $c_1c_2c_3$ in a string of English text.
- How would we build a model of $p(c_1c_2c_3)$?
 - If we use a single clique function over $c_1c_2c_3$, the full joint clique potential would be huge: 26^3-1 parameters.
 - However, we often know that some particular joint settings of the variables in a clique are quite likely or quite unlikely. e.g. ing, ate, ion, ?ed, qu?, jkx, zzz,...
- A "feature" is a function which is vacuous over all joint settings except a few particular ones on which it is high or low.
 - For example, we might have $f_{ing}(c_1c_2c_3)$ which is 1 if the string is 'ing' and 0 otherwise, and similar features for '?ed', etc.
- We can also define features when the inputs are continuous.
 Then the idea of a cell on which it is active disappears, but we might still have a compact parameterization of the feature.





- By exponentiating them, each feature function can be made into a "micropotential". We can multiply these micropotentials together to get a clique potential.
- Example: a clique potential $\psi(c_1c_2c_3)$ could be expressed as:

$$\psi_{c}(c_{1}, c_{2}, c_{3}) = e^{\theta_{ing}f_{ing}} \times e^{\theta_{ing}f_{ing}} \times \dots$$

$$= \exp\left\{\sum_{k=1}^{K} \theta_{k}f_{k}(c_{1}, c_{2}, c_{3})\right\}$$

- This is still a potential over 26³ possible settings, but only uses *K* parameters if there are *K* features.
 - By having one indicator function per combination of x_c , we recover the standard tabular potential.





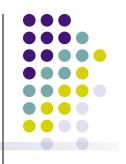
- Each feature has a weight θ_k which represents the numerical strength of the feature and whether it increases or decreases the probability of the clique.
- The marginal over the clique is a generalized exponential family distribution, actually, a GLIM:

$$p(c_1, c_2, c_3) \propto \exp \begin{cases} \theta_{\text{ing}} f_{\text{ing}}(c_1, c_2, c_3) + \theta_{\text{?ed}} f_{\text{?ed}}(c_1, c_2, c_3) + \theta_{\text{qu?}} f_{\text{qu?}}(c_1, c_2, c_3) + \theta_{\text{zzz}} f_{\text{zzz}}(c_1, c_2, c_3) + \cdots \end{cases}$$

 In general, the features may be overlapping, unconstrained indicators or any function of any subset of the clique variables:

$$\psi_c(\mathbf{x}_c) \stackrel{\text{def}}{=} \exp \left\{ \sum_{i \in I_c} \theta_k f_k(\mathbf{x}_{c_i}) \right\}$$





We can multiply these clique potentials as usual:

$$p(\mathbf{x}) = \frac{1}{Z(\theta)} \prod_{c} \psi_{c}(\mathbf{x}_{c}) = \frac{1}{Z(\theta)} \exp \left\{ \sum_{c} \sum_{i \in I_{c}} \theta_{k} f_{k}(\mathbf{x}_{c_{i}}) \right\}$$

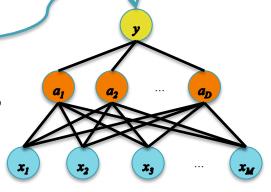
 However, in general we can forget about associating features with cliques and just use a simplified form:

$$p(\mathbf{x}) = \frac{1}{Z(\theta)} \exp \left\{ \sum_{i} \theta_{i} f_{i}(\mathbf{x}_{c_{i}}) \right\}$$

- This is just our friend the exponential family model, with the features as sufficient statistics!
- Learning: recall that in IPF, we have $\psi_c^{(t+1)}(\mathbf{x}_c) = \psi_c^{(t)}(\mathbf{x}_c) \frac{p(\mathbf{x}_c)}{p^{(t)}(\mathbf{x}_c)}$
 - Not obvious how to use this rule to update the weights and features individually !!!

Options for MLE of MRFs

- Setting I: $\psi_C(oldsymbol{x}_C) = heta_{C,oldsymbol{x}_C}$
 - A. MLE by inspection (Decomposable Models)
 - B. Iterative Proportional Fitting (IPF)
- Setting II: $\psi_C(m{x}_C) = \exp(m{ heta} \cdot m{f}(m{x}_C))$
 - C. Generalized Iterative Scaling
 - D. Gradient-based Methods
- Setting III: $\psi_C(oldsymbol{x}_C)=$
 - D. Gradient-based Methods



Generalized Iterative Scaling (GIS)

Key idea:

- Define a function which lower-bounds the loglikelihood
- Observe that the bound is tight at current parameters
- Increase lower-bound by fixed-point iteration in order to increase log-likelihood

Side note: This idea is akin to a standard derivation of the Expectation-Maximization (EM) algorithm

Generalized Iterative Scaling (GIS)

GIS applies fixed point iteration to the derivative of a lower-bound of the likelihood objective

$$\tilde{\ell}(\theta; \mathcal{D}) = \frac{1}{N} \sum_{n=1}^{N} \log p(\mathbf{x}^{(n)} \mid \boldsymbol{\theta})$$

$$\tilde{\ell}(\theta; \mathcal{D}) \ge \Lambda(\theta)$$

$$\frac{\partial \Lambda}{\partial \theta_{i}} = \sum_{\mathbf{x}} \tilde{\boldsymbol{p}}(\mathbf{x}) f_{i}(\mathbf{x}) - \exp(\Delta \theta_{i}^{(t)}) \sum_{\mathbf{x}} \boldsymbol{p}(\mathbf{x} \mid \boldsymbol{\theta}^{(t)}) f_{i}(\mathbf{x})$$

$$\theta_{i}^{(t+1)} = \theta_{i}^{(t)} + \log \left(\frac{\sum_{x} \tilde{\boldsymbol{p}}(x) f_{i}(x)}{\sum_{x} \boldsymbol{p}^{(t)}(x) f_{i}(x)}\right)$$

- Given avg. likelihood objective
- Derive lower bound
- 3. Compute derivative of bound,set to zero
 - Rearrange the equation s.t. one parameter appears on the LHS.
 - Initialize the parameters.
 - For each i in $\{1,...K\}$, update each parameter and increment t:
- 7. Repeat #6 until convergence

The lower bound is obtained by linearizing a log and applying Jensen-Shannon.

$$\ell^{\sim}(\theta; \mathcal{D}) \ge \sum_{i} \theta_{i} \sum_{x} \widetilde{p}(x) f_{i}(x) - \sum_{x} p(x \mid \theta^{(t)}) \sum_{i} f_{i}(x) \exp(\Delta \theta_{i}^{(t)}) - \log Z(\theta^{(t)}) + 1 = \Lambda(\theta)$$



MLE of Feature Based UGMs

Scaled likelihood function

$$\ell^{\sim}(\theta; D) = \ell(\theta; D) / N = \frac{1}{N} \sum_{n} \log p(x_n | \theta)$$

$$= \sum_{x} \widetilde{p}(x) \log p(x | \theta)$$

$$= \sum_{x} \widetilde{p}(x) \sum_{i} \theta_{i} f_{i}(x) - \log Z(\theta)$$

- Instead of optimizing this objective directly, we attack its lower bound
 - The logarithm has a linear upper bound ...

$$\log Z(\theta) \le \mu Z(\theta) - \log \mu - 1$$

This bound holds for all μ , in particular, for

$$\mu = Z^{-1}(\theta^{(t)})$$

Thus we have

This bound holds for all
$$\mu$$
, in particular, for $\mu = Z^{-1}(\theta^{(t)})$ hus we have
$$\ell^{\sim}(\theta; \mathcal{D}) \geq \sum_{x} \widetilde{p}(x) \sum_{i} \theta_{i} f_{i}(x) - \frac{Z(\theta)}{Z(\theta^{(t)})} - \log Z(\theta^{(t)}) + 1$$

Extra slides from 2015 Generalized Iterative Scaling (GIS)

Lower bound of scaled loglikelihood

$$\ell^{\sim}(\theta; \mathcal{D}) \ge \sum_{x} \widetilde{p}(x) \sum_{i} \theta_{i} f_{i}(x) - \frac{Z(\theta)}{Z(\theta^{(t)})} - \log Z(\theta^{(t)}) + 1$$

• Define $\Delta \theta_i^{(t)} \stackrel{\text{def}}{=} \theta_i - \theta_i^{(t)}$

$$\ell^{\sim}(\theta; \mathcal{D}) \geq \sum_{x} \widetilde{p}(x) \sum_{i} \theta_{i} f_{i}(x) - \frac{1}{Z(\theta^{(t)})} \sum_{x} \exp\left\{\sum_{i} \theta_{i} f_{i}(x)\right\} - \log Z(\theta^{(t)}) + 1$$

$$= \sum_{i} \theta_{i} \sum_{x} \widetilde{p}(x) f_{i}(x) - \frac{1}{Z(\theta^{(t)})} \sum_{x} \exp\left\{\sum_{i} \theta_{i}^{(t)} f_{i}(x)\right\} \exp\left\{\sum_{i} \Delta \theta_{i}^{(t)} f_{i}(x)\right\} - \log Z(\theta^{(t)}) + 1$$

$$= \sum_{i} \theta_{i} \sum_{x} \widetilde{p}(x) f_{i}(x) - \sum_{x} p(x \mid \theta^{(t)}) \exp\left\{\sum_{i} \Delta \theta_{i}^{(t)} f_{i}(x)\right\} - \log Z(\theta^{(t)}) + 1$$

- Relax again

 - Assume $f_i(x) \ge 0$, $\sum_i f_i(x) = 1$ Convexity of exponential: $\exp\left(\sum_i \pi_i x_i\right) \le \sum_i \pi_i \exp(x_i)$
- We have:

$$\ell^{\sim}(\theta; \mathcal{D}) \ge \sum_{i} \theta_{i} \sum_{x} \widetilde{p}(x) f_{i}(x) - \sum_{x} p(x \mid \theta^{(t)}) \sum_{i} f_{i}(x) \exp(\Delta \theta_{i}^{(t)}) - \log Z(\theta^{(t)}) + 1 = \Lambda(\theta)$$
© Eric Xing @ CMU, 2005-2015





GIS

Lower bound of scaled loglikelihood

$$\ell^{\sim}(\theta; \mathcal{D}) \ge \sum_{i} \theta_{i} \sum_{x} \widetilde{p}(x) f_{i}(x) - \sum_{x} p(x \mid \theta^{(t)}) \sum_{i} f_{i}(x) \exp(\Delta \theta_{i}^{(t)}) - \log Z(\theta^{(t)}) + 1 \stackrel{\text{def}}{=} \Lambda(\theta)$$

- Take derivative: $\frac{\partial \Lambda}{\partial \theta} = \sum_{i} \widetilde{p}(x) f_i(x) \exp(\Delta \theta_i^{(t)}) \sum_{i} p(x \mid \theta^{(t)}) f_i(x)$
- Set to zero

$$e^{\Delta \theta_i^{(t)}} = \frac{\sum_{x} \widetilde{p}(x) f_i(x)}{\sum_{x} p(x \mid \theta^{(t)}) f_i(x)} = \frac{\sum_{x} \widetilde{p}(x) f_i(x)}{\sum_{x} p^{(t)}(x) f_i(x)} Z(\theta^{(t)})$$

- where $p^{(t)}(x)$ is the unnormalized version of $p(x|\theta^{(t)})$
- Update

$$\theta_i^{(t+1)} = \theta_i^{(t)} + \Delta \theta_i^{(t)} \Rightarrow p^{(t+1)}(x) = p^{(t)}(x) \prod_i e^{\Delta \theta_i^{(t)} f_i(x)}$$

$$p^{(t+1)}(x) = \frac{p^{(t)}(x)}{Z(\theta^{(t)})} \prod_{i} \left(\frac{\sum_{x} \tilde{p}(x) f_{i}(x)}{\sum_{x} p^{(t)}(x) f_{i}(x)} Z(\theta^{(t)}) \right)^{T_{i}(x)}$$

$$\Rightarrow \qquad = \frac{p^{(t)}(x)}{Z(\theta^{(t)})} \prod_{i} \left(\frac{\sum_{x} \tilde{p}(x) f_{i}(x)}{\sum_{x} p^{(t)}(x) f_{i}(x)} \right)^{f_{i}(x)} \left(Z(\theta^{(t)}) \right)^{\sum_{i} f_{i}(x)}$$

$$= p^{(t)}(x) \prod_{i} \left(\frac{\sum_{x} \tilde{p}(x) f_{i}(x)}{\sum_{x} p^{(t)}(x) f_{i}(x)} \right)^{f_{i}(x)}$$

$$= p^{(t)}(x) \prod_{i} \left(\frac{\sum_{x} \tilde{p}(x) f_{i}(x)}{\sum_{x} p^{(t)}(x) f_{i}(x)} \right)^{f_{i}(x)}$$

$$= p^{(t)}(x) \prod_{i} \left(\frac{\sum_{x} \tilde{p}(x) f_{i}(x)}{\sum_{x} p^{(t)}(x) f_{i}(x)} \right)^{f_{i}(x)}$$

$$= p^{(t)}(x) \prod_{i} \left(\frac{\sum_{x} \tilde{p}(x) f_{i}(x)}{\sum_{x} p^{(t)}(x) f_{i}(x)} \right)^{f_{i}(x)}$$

$$= p^{(t)}(x) \prod_{i} \left(\frac{\sum_{x} \tilde{p}(x) f_{i}(x)}{\sum_{x} p^{(t)}(x) f_{i}(x)} \right)^{f_{i}(x)}$$

$$= p^{(t)}(x) \prod_{i} \left(\frac{\sum_{x} \tilde{p}(x) f_{i}(x)}{\sum_{x} p^{(t)}(x) f_{i}(x)} \right)^{f_{i}(x)}$$

$$\psi_c^{(t+1)}(\mathbf{x}_c) = \psi_c^{(t)}(\mathbf{x}_c) \frac{\widetilde{\boldsymbol{p}}(\mathbf{x}_c)}{\boldsymbol{p}^{(t)}(\mathbf{x}_c)}$$

Contrast of IPF and GIS



- IPF is a general algorithm for finding MLE of UGMs.
 - a **fixed-point equation** for ψ_c over single cliques, coordinate ascent
 - Requires the potential to be fully parameterized
 - The clique described by the potentials do not have to be max-clique
 - For fully decomposable model, reduces to a single step iteration

GIS

- Iterative scaling on general UGM with feature-based potentials
- IPF is a special case of GIS which the clique potential is built on features defined as an indicator function of clique configurations.

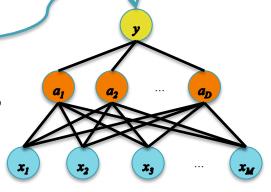
GIS:

$$p^{(t+1)}(x) = p^{(t)}(x) \prod_{i} \left(\frac{\sum_{x} \widetilde{p}(x) f_{i}(x)}{\sum_{x} p^{(t)}(x) f_{i}(x)} \right)^{f_{i}(x)} \qquad \qquad \psi_{c}^{(t+1)}(\mathbf{x}_{c}) = \psi_{c}^{(t)}(\mathbf{x}_{c}) \frac{\widetilde{p}(\mathbf{x}_{c})}{p^{(t)}(\mathbf{x}_{c})}$$

$$\theta_{i}^{(t+1)} = \theta_{i}^{(t)} + \log \left(\frac{\sum_{x} \widetilde{p}(x) f_{i}(x)}{\sum_{x} p^{(t)}(x) f_{i}(x)} \right)$$

Options for MLE of MRFs

- Setting I: $\psi_C({m x}_C) = heta_{C,{m x}_C}$
 - A. MLE by inspection (Decomposable Models)
 - B. Iterative Proportional Fitting (IPF)
- Setting II: $\psi_C(m{x}_C) = \exp(m{ heta} \cdot m{f}(m{x}_C))$
 - C. Generalized Iterative Scaling
 - D. Gradient-based Methods
- Setting III: $\psi_C(m{x}_C) = 0$
 - D. Gradient-based Methods



Recipe for Gradient-based Learning

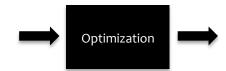
- 1. Write down the objective function
- Compute the partial derivatives of the objective (i.e. gradient, and maybe Hessian)
- Feed objective function and derivatives into black box



4. Retrieve optimal parameters from black box

Optimization Algorithms

What is the black box?



- Newton's method
- Hessian-free / Quasi-Newton methods
 - Conjugate gradient
 - L-BFGS
- Stochastic gradient methods
 - Stochastic gradient descent (SGD)
 - Stochastic meta-descent
 - AdaGrad

Whiteboard

- Gradient of MRF log-likelihood for featurebased potentials
- Gradient of CRF log-likelihood for featurebased potentials [next time]
- L1 and L2 regularization

Practical Considerations for Gradient-based Methods

- Overfitting
 - L2 regularization
 - L1 regularization
 - Regularization by early stopping
- For SGD: Sparse updates

"Empirical" Comparison of Parameter Estimation Methods

- Example NLP task: CRF dependency parsing
- Suppose: Training time is dominated by inference
- Dataset: One million tokens
- Inference speed: 1,000 tokens / sec
- → 0.27 hours per pass through dataset

	# passes through data to converge	# hours to converge
GIS	1000+	270
L-BFGS	100+	27
SGD	10	~3

Summary

Setting I: $\psi_C(oldsymbol{x}_C) = heta_C$

A. MLE by inspection (Decomposable Models)

- Very limited applicability
- Exemplifies the need for general algorithms

B. Iterative Proportional Fitting (IPF)

- Guaranteed to converge
- Only applies to "tabular" potential functions

C. Generalized Iterative Scaling (GIS)

- Maximizes a lower-bound of log-likelihood
- Iterative algorithm (like IPF), but more broadly applies to exponential family potentials

D. Gradient-based Methods

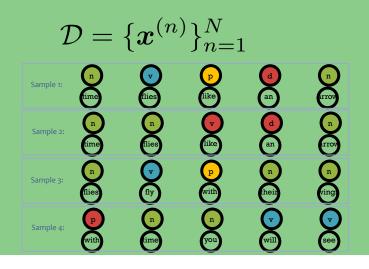
- Doesn't require fancy optimization algorithms (i.e. SGD works great)
- Faster convergence than GIS
- Applies to arbitrary potentials [later in the course]

Setting II: $\psi_C(oldsymbol{x}_C) = \exp(oldsymbol{ heta} \cdot oldsymbol{f}(oldsymbol{x}_C))$

58

MLE for Undirected GMs

1. Data



2. Model

$$p(\boldsymbol{x}\mid\boldsymbol{\theta}) = \frac{1}{Z(\boldsymbol{\theta})} \prod_{C\in\mathcal{C}} \psi_C(\boldsymbol{x}_C)$$

3. Objective $_N$

$$\ell(\theta; \mathcal{D}) = \sum_{n=1} \log p(\boldsymbol{x}^{(n)} \mid \boldsymbol{\theta})$$

5. Inference

1. Marginal Inference

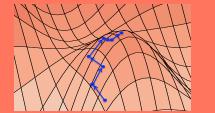
$$p(\boldsymbol{x}_C) = \sum_{\boldsymbol{x}': \boldsymbol{x}_C' = \boldsymbol{x}_C} p(\boldsymbol{x}' \mid \boldsymbol{\theta})$$

2. Partition Function

$$Z(\boldsymbol{ heta}) = \sum_{oldsymbol{x}} \prod_{C \in \mathcal{C}} \psi_C(oldsymbol{x}_C)$$

4. Learning

$$\boldsymbol{\theta}^* = \operatorname*{argmax}_{\boldsymbol{\theta}} \ell(\boldsymbol{\theta}; \mathcal{D})$$



Contrast of MLE for directed / undirected GMs



- For <u>directed graphical models</u>, the log-likelihood decomposes into a sum of terms, one per family (node plus parents).
- For <u>undirected graphical models</u>, the log-likelihood does not decompose, because the normalization constant Z is a function of all the parameters

$$P(x_1, \dots, x_n) = \frac{1}{Z} \prod_{c \in C} \psi_c(\mathbf{x}_c) \qquad Z = \sum_{x_1, \dots, x_n} \prod_{c \in C} \psi_c(\mathbf{x}_c)$$

 In general, we will need to do inference (i.e., marginalization) to learn parameters for undirected models, even in the fully observed case.

5. Inference

Three Tasks:

1. Marginal Inference

Compute marginals of variables and cliques

$$p(x_i) = \sum_{\boldsymbol{x}': x_i' = x_i} p(\boldsymbol{x}' \mid \boldsymbol{\theta}) \qquad p(\boldsymbol{x}_C) = \sum_{\boldsymbol{x}': \boldsymbol{x}_C' = \boldsymbol{x}_C} p(\boldsymbol{x}' \mid \boldsymbol{\theta})$$

2. Partition Function

Compute the normalization constant

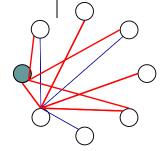
$$Z(\boldsymbol{\theta}) = \sum_{\boldsymbol{x}} \prod_{C \in \mathcal{C}} \psi_C(\boldsymbol{x}_C)$$

3. MAP Inference

Compute variable assignment with highest probability

$$\hat{\boldsymbol{x}} = \underset{\boldsymbol{x}}{\operatorname{argmax}} p(\boldsymbol{x} \mid \boldsymbol{\theta})$$





ML Structural Learning via Neighborhood Selection for completely observed MRF

Data

$$(x_1^{(1)},...,x_n^{(1)})$$

 $(x_1^{(2)},...,x_n^{(2)})$
...
 $(x_1^{(M)},...,x_n^{(M)})$





Multivariate Gaussian density:

$$p(\mathbf{x} \mid \mu, \Sigma) = \frac{1}{(2\pi)^{n/2} |\Sigma|^{1/2}} \exp\left\{-\frac{1}{2} (\mathbf{x} - \mu)^T \Sigma^{-1} (\mathbf{x} - \mu)\right\}$$

• WOLG: let $\mu=0$ $Q=\Sigma^{-1}$

$$p(x_1, x_2, \dots, x_p \mid \mu = 0, Q) = \frac{|Q|^{1/2}}{(2\pi)^{n/2}} \exp\left\{-\frac{1}{2} \sum_{i} q_{ii}(x_i)^2 - \sum_{i < j} q_{ij} x_i x_j\right\}$$

 We can view this as a continuous Markov Random Field with potentials defined on every node and edge:





 Assuming the nodes are discrete, and edges are weighted, then for a sample x_d, we have

$$P(\mathbf{x}_d|\Theta) = \exp\left(\sum_{i \in V} \theta_{ii}^t x_{d,i} + \sum_{(i,j) \in E} \theta_{ij} x_{d,i} x_{d,j} - A(\Theta)\right)$$

The covariance and the precision matrices



Covariance matrix ∑

$$\Sigma_{i,j} = 0 \implies X_i \perp X_j \quad \text{or} \quad p(X_i, X_j) = p(X_i)p(X_j)$$

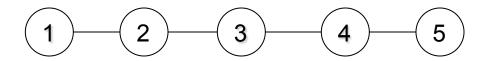
- Graphical model interpretation?
- Precision matrix $Q = \Sigma^{-1}$

$$Q_{i,j} = 0 \quad \Rightarrow \quad X_i \perp X_j | \mathbf{X}_{-ij} \quad \text{or} \quad p(X_i, X_j | \mathbf{X}_{-ij}) = p(X_i | \mathbf{X}_{-ij}) p(X_j | \mathbf{X}_{-ij})$$

Graphical model interpretation?

Sparse precision vs. sparse covariance in GGM





$$\Sigma^{-1} = \begin{pmatrix} 1 & 6 & 0 & 0 & 0 \\ 6 & 2 & 7 & 0 & 0 \\ 0 & 7 & 3 & 8 & 0 \\ 0 & 0 & 8 & 4 & 9 \\ 0 & 0 & 0 & 9 & 5 \end{pmatrix}$$

$$\Sigma^{-1} = \begin{pmatrix} 1 & 6 & 0 & 0 & 0 \\ 6 & 2 & 7 & 0 & 0 \\ 0 & 7 & 3 & 8 & 0 \\ 0 & 0 & 8 & 4 & 9 \\ 0 & 0 & 0 & 9 & 5 \end{pmatrix} \qquad \Sigma = \begin{pmatrix} 0.10 & 0.15 & -0.13 & -0.08 & 0.15 \\ 0.15 & -0.03 & 0.02 & 0.01 & -0.03 \\ -0.13 & 0.02 & 0.10 & 0.07 & -0.12 \\ -0.08 & 0.01 & 0.07 & -0.04 & 0.07 \\ 0.15 & -0.03 & -0.12 & 0.07 & 0.08 \end{pmatrix}$$

$$\Sigma_{15}^{-1} = 0 \Leftrightarrow X_1 \perp X_5 | X_{nbrs(1) \text{ or } nbrs(5)}$$

$$\Leftrightarrow$$

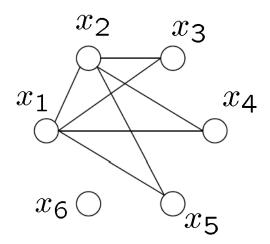
$$X_1 \perp X_5 \Leftrightarrow \Sigma_{15} = 0$$





$$Q = \begin{pmatrix} * & * & * & * & * & 0 \\ * & * & * & * & * & 0 \\ * & * & * & 0 & 0 & 0 \\ * & * & 0 & * & 0 & 0 \\ * & * & 0 & 0 & 0 & * & 0 \\ 0 & 0 & 0 & 0 & 0 & * \end{pmatrix}$$





- How to estimate this MRF?
- What if *p* >> *n*
 - MLE does not exist in general!
 - What about only learning a "sparse" graphical model?
 - This is possible when s=o(n)
 - Very often it is the structure of the GM that is more interesting ...



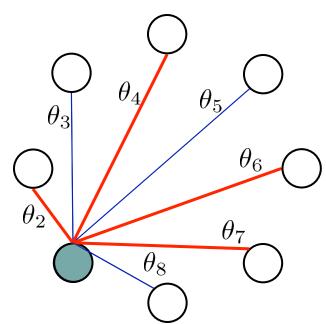


$$\hat{\theta}_i = \arg\min_{\theta_i} l(\theta_i) + \lambda_1 || \theta_i ||_1$$

where
$$l(\theta_i) = \log P(y_i|\mathbf{x}_i, \theta_i)$$
.

Graph Regression





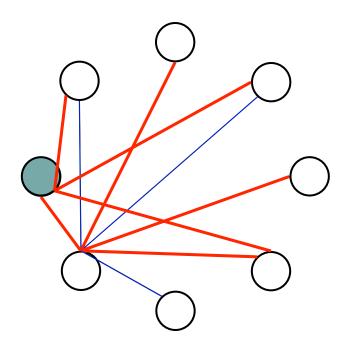
Neighborhood selection

Lasso:

$$\hat{\theta} = \arg\min_{\theta} \sum_{t=1}^{T} l(\theta) + \lambda_1 ||\theta||_1$$

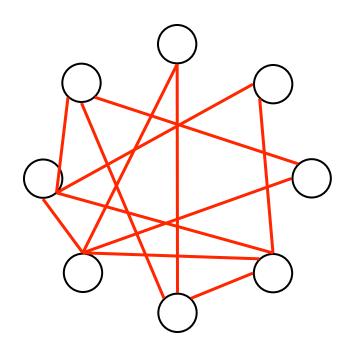












It can be shown that: given *iid* samples, and under several technical conditions (e.g., "irrepresentable"), the recovered structured is "sparsistent" even when p >> n

Learning Ising Model (i.e. pairwise MRF)



 Assuming the nodes are discrete, and edges are weighted, then for a sample x_d, we have

$$P(\mathbf{x}_d|\Theta) = \exp\left(\sum_{i \in V} \theta_{ii}^t x_{d,i} + \sum_{(i,j) \in E} \theta_{ij} x_{d,i} x_{d,j} - A(\Theta)\right)$$

 It can be shown following the same logic that we can use L_1 regularized logistic regression to obtain a sparse estimate of the neighborhood of each variable in the discrete case.

Consistency



• **Theorem**: for the graphical regression algorithm, under certain verifiable conditions (omitted here for simplicity):

$$\mathbb{P}\left[\hat{G}(\lambda_n) \neq G\right] = \mathcal{O}\left(\exp\left(-Cn^{\epsilon}\right)\right) \to 0$$

Note the from this theorem one should see that the regularizer is not actually used to introduce an "artificial" sparsity bias, but a devise to ensure consistency under finite data and high dimension condition.