

Markov Random Fields: Inference and Estimation

SPiNCOM reading group April 24th , 2017

Dimitris Berberidis

Ack: Juan-Andres Bazerque

Probabilistic graphical models

- \Box Set of random variables $\mathbf{x} = [x_1 x_2 \dots x_N]$
 - \succ Graph ${\cal G}$ represents joint $p({f x})$
 - Nodes correspond to random variables
 - Edges imply relations between rv's
- Some applications
 - Speech recognition, computer vision
 - Decoding
 - Gene reg. networks, disease diagnosis
- Key idea: Graph models <u>conditional independencies</u>
 - Two main tasks: Inference and Estimation

Inference: Given observed \mathbf{x}_s , obtain (marginal) conditionals $p(x_i|\mathbf{x}_s)$, $orall i\in \mathcal{S}^c$

Estimation: Given samples $\{\mathbf{x}^{(t)}\}_{t=1}^T$ estimate \mathcal{G} (and thus $p(\mathbf{x})$)



Roadmap

- Bayesian networks basics
- Markov Random Fields
- Continuous valued MRFs
 - Inference using Harmonic solution
 - Structure estimation through I-1 penalized MLE
- □ Binary valued MRFs (Ising model)
 - Inference
 - Gaussian approximation Random walk interpretation
 - MCMC
 - Structure estimation
 - Pseudo MLE
 - Logistic regression

Conclusions

Directed Acyclical GMs (Bayesian networks)

□ Ordered Markov property : $x_s \perp \mathbf{x}_{\text{pred}(s) \mid pa(s)} | \mathbf{x}_{pa(s)} |$

- Complete independence: Markov "Blanket" (Parents+children+co-parents)
- Joint pdf modeled as product of conditionals: $p(\mathbf{x}_{1:V}|G) = \prod p(x_t|\mathbf{x}_{pa(t)})$



Basic building blocks of Bayesian nets



Undirected GMs (Markov random fields)

- More natural in some domains (e.g. special statistics, relational data)
 - Simple rule: Nodes not connected w. edge are conditionally independent
- Joint pdf parametrized and modeled as product of factors(<u>not conditionals</u>)
 - > Each factor $\psi_c(\mathbf{x}_c; \theta_c)$ or potential corresponds to a maximal clique c
- Hamersley-Clifford theorem
 - $\succ p(\mathbf{x})$ satisfies the CI properties of an undirected graph $\underline{\textit{iff}}$

$$p(\mathbf{x};\theta) = \frac{1}{Z(\theta)} \prod_{c \in \mathcal{C}} \psi_c(\mathbf{x}_c;\theta_c) \quad \text{where} \quad Z(\theta) := \sum_{\mathbf{x}} \prod_{c \in \mathcal{C}} \psi_c(\mathbf{x}_c;\theta_c)$$

Example $p(y|\theta) = \frac{1}{Z(\theta)}\psi_{123}(y_1, y_2, y_3)\psi_{234}(y_2, y_3, y_4)\psi_{35}(y_3, y_5)$ Partition Function: $Z = \sum_{y} \psi_{123}(y_1, y_2, y_3)\psi_{234}(y_2, y_3, y_4)\psi_{35}(y_3, y_5)$ Generally NP-hard to compute

Equivalence of DGMs and UGMs

- Moralization: Transition from directed to undirected GM
 - Drop directionality and connect "unmarried" parents
 - Information may be lost during transition (see example)





5|2|

 $4 \perp$

MRFs with energy functions

Clique potentials usually represented using an "energy" function $E_c(\mathbf{x}_c)$

$$\psi_c(\mathbf{x}_c; \theta_c) = \exp\left(-E(\mathbf{x}_c; \theta_c)\right)$$

- Joint (Gibbs distribution) $p(\mathbf{x}; \theta) = \frac{1}{Z(\theta)} \exp\left(-\sum_{c} E(\mathbf{x}_{c}; \theta_{c})\right)$
- High probability states correspond to low energy configurations
- Any MRF can be decomposed to pairwise potentials (and energy functions)

$$p(\mathbf{x};\theta) = \frac{1}{Z(\theta)} \exp\left(-\sum_{i,j} \phi(x_i, x_j)\theta_{ij}\right)$$

 \Box MRF is associative if $\phi(x_i, x_j)$ measures difference btw x_i and x_j , and $\theta \ge 0$

- > Gaussian MRF: $\phi(x_i, x_j) = (x_i x_j)^2$
- > Ising (binary +1,-1) model: $\phi(x_i, x_j) = x_i x_j$

Gaussian MRFs

□ Joint Gaussian fully parametrized by covariance and mean

$$p(\boldsymbol{x};\boldsymbol{\Sigma}) = \frac{1}{2^{N/2} \operatorname{det}(\boldsymbol{\Sigma})^{1/2}} e^{-\frac{1}{2}(\boldsymbol{x}-\boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1}(\boldsymbol{x}-\boldsymbol{\mu})}$$

GMRF structure given by precision matrix (inv. Cov.)

Also viewed as the Laplacian of the graph

$$\Theta = \Sigma^{-1} = ((\theta_{ij}))$$

$$p(x; \Theta) = \frac{\det(\Theta)^{1/2}}{2^{N/2}} e^{-\frac{1}{2}(x-\mu)^T \Theta(x-\mu)}$$

$$x_i \perp x_j | \mathbf{x}_{\mathcal{V}-i,j} \longleftrightarrow \theta_{i,j} = 0$$



lacksquare Assume for simplicity (and wlog) that $\mu = oldsymbol{0}$

lacksim Inference: Given known $m\Theta$ and observed \mathbf{x}_s , find $p(\mathbf{x}_{s^c}|\mathbf{x}_s)$

Inference via Harmonic solution

Negative log-likelihood of joint

$$-\log p(\mathbf{x}_{s^{c}}, \mathbf{x}_{s}) \propto [\mathbf{x}_{s^{c}}^{T} \ \mathbf{x}_{s}^{T}] \begin{bmatrix} \Theta_{s^{c}, s^{c}} & \Theta_{s^{c}, s} \\ \Theta_{s, s^{c}} & \Theta_{s, s} \end{bmatrix} \begin{bmatrix} \mathbf{x}_{s^{c}} \\ \mathbf{x}_{s} \end{bmatrix}$$

$$= \mathbf{x}_{s^{c}}^{T} \Theta_{s^{c}, s^{c}} \mathbf{x}_{s^{c}} + 2\mathbf{x}_{s}^{T} \Theta_{s, s^{c}} \mathbf{x}_{s^{c}} + \mathbf{x}_{s}^{T} \Theta_{s, s} \mathbf{x}_{s} + \mathbf{x}_{s^{c}}^{T} \Theta_{s^{c}, s} \Theta_{s, s}^{-1} \Theta_{s, s^{c}} \mathbf{x}_{s^{c}} - \mathbf{x}_{s^{c}}^{T} \Theta_{s^{c}, s} \Theta_{s, s}^{-1} \Theta_{s, s^{c}} \mathbf{x}_{s^{c}} - \mathbf{x}_{s^{c}}^{T} \Theta_{s^{c}, s} \Theta_{s, s}^{-1} \Theta_{s, s^{c}} \mathbf{x}_{s^{c}} = \begin{bmatrix} \mathbf{x}_{s^{c}}^{T} \Theta_{s^{c}, s^{c}} \mathbf{x}_{s^{c}} + 2\mathbf{x}_{s}^{T} \Theta_{s, s^{c}} \mathbf{x}_{s^{c}} + \mathbf{x}_{s}^{T} \Theta_{s, s^{c}} \Theta_{s^{c}, s} \Theta_{s^{c}, s} \Theta_{s^{c}, s} \mathbf{x}_{s} + \mathbf{x}_{s}^{T} \Theta_{s^{c}, s^{c}} \Theta_{s^{c}, s} \Theta_{s^{c}, s} \mathbf{x}_{s} + \mathbf{x}_{s}^{T} (\Theta_{s, s} - \Theta_{s, s^{c}} \Theta_{s^{c}, s^{c}} \Theta_{s^{c}, s}) \mathbf{x}_{s} + \log p(\mathbf{x}_{s^{c}} | \mathbf{x}_{s}) = (\mathbf{x}_{s^{c}} + \Theta_{s^{c}, s^{c}}^{-1} \Theta_{s^{c}, s^{c}} \Theta_{s^{c}, s}) \mathbf{x}_{s} + \log p(\mathbf{x}_{s}) = (\mathbf{x}_{s^{c}} | \mathbf{x}_{s}) - \log p(\mathbf{x}_{s})^{T} \Theta_{s^{c}, s^{c}} (\mathbf{x}_{s^{c}} + \Theta_{s^{c}, s^{c}}^{-1} \Theta_{s^{c}, s} \mathbf{x}_{s})$$

Conditional mean of \mathbf{X}_{s^c} contains all information from observed \mathbf{X}_s

GMRF structure estimation via maximum likelihood

Given
$$\{x^{(t)}\}_{t=1}^{T}$$
, goal is to estimate Θ and μ

$$p(x; \Theta) = \frac{\det(\Theta)^{1/2}}{2^{N/2}} e^{-\frac{1}{2}(x-\mu)^{T}\Theta(x-\mu)}$$

$$\log p(x; \Theta) = \frac{1}{2}\log \det(\Theta) - \frac{1}{2}(x-\mu)^{T}\Theta(x-\mu) + c$$

Log-likelihood

$$\mathcal{L}(\Theta) = \max_{\mu} \frac{1}{T} \sum_{n=1}^{T} \left(\log \det(\Theta) - \left(x^{(t)} - \mu \right)^{T} \Theta \left(x^{(t)} - \mu \right) \right)$$
$$= \log \det(\Theta) - \frac{1}{T} \sum_{t=1}^{N} \left(x^{(t)} - \hat{\mu} \right)^{T} \Theta \left(x^{(t)} - \hat{\mu} \right)$$
$$= \log \det(\Theta) - \operatorname{trace} \left(\underbrace{\frac{1}{T} \sum_{t=1}^{T} \left(x^{(t)} - \hat{\mu} \right) \left(x^{(t)} - \hat{\mu} \right)^{T} \Theta}_{S} \Theta \right)$$

$$\hat{\Theta}_{ML} = rg\max_{oldsymbol{\Theta} \succ oldsymbol{0}} \log \det(oldsymbol{\Theta}) - \mathrm{trace}(oldsymbol{S}oldsymbol{\Theta})$$

ℓ_1 -penalized MLE of Θ

$$\hat{\Theta}_{ML} = \arg \max_{\Theta \succ 0} \log \det(\Theta) - \operatorname{trace}(S\Theta)$$

- lacksim Closed-form solution: $\Theta^{-1} S = 0 \implies \hat{\Theta}_{ML} = S^{-1}$
- $\square \hat{\Theta}_{ML}$ generally is full matrix

 \Box Idea: Add ℓ_1 constrain on to enforce (sparse) graph structure

$$\hat{\Theta}_{\ell_1} = \arg\min_{\Theta \succ 0} - \log\det(\Theta) + \operatorname{trace}(S\Theta) + \lambda ||\Theta||_1$$
$$||\Theta||_1 = \sum_{i=1}^{P} \sum_{j=1}^{P} |\theta_{ij}|$$

D Problem is convex and for $\mathbf{W} = \mathbf{\Theta}^{-1}$ is equivalent to

$$\widehat{W} = rg\max_{oldsymbol{W} \succ oldsymbol{0}} |\log \det(oldsymbol{W})|$$
 s.to $||oldsymbol{W} - oldsymbol{S}||_{\infty} \leq \lambda$

Solvable via Graphical Lasso

O. Banerjee, L. El Ghaoui, and A. d'Aspremont, "Model selection through sparse maximum likelihood estimation for multivariate Gaussian or binary data," *J. Machine Learning Research*, vol. 9, pp. 485-516, June 2008.

Binary random variables

lacksquare Ising model for $m{x} \in \{-1,1\}^N$ or $m{x} \in \{0,1\}^N$

$$P(\boldsymbol{x}; \boldsymbol{\Theta}) = \exp\left(\boldsymbol{x}^T \boldsymbol{\Theta} \boldsymbol{x} - Z(\boldsymbol{\Theta})\right) \quad \theta_{ii} = 0$$

Log partition function: $Z(\boldsymbol{\Theta}) = \log\left(\sum_{\boldsymbol{x} \in \{-1,1\}^N} \exp\left(\boldsymbol{x}^T \boldsymbol{\Theta} \boldsymbol{x}\right)\right)$

Estimation: I-1 penalized maximum likelihood for Θ

$$\hat{\Theta}_{ML} = \arg\min_{\Theta} Z(\Theta) - \operatorname{trace}(\Theta S) + \lambda ||\Theta||_1$$

Problem: $Z(\Theta)$ combinatorialy complex to compute

Two alternatives: $Z(\Theta)$ is upper-bounded or avoided

C Similar problem for inference: $p(\mathbf{x}_{s^c}|\mathbf{x}_s)$ can only be approximated





The role of Θ in Ising model

Claim:
$$\theta_{ij} = 0 \iff x_i \perp x_j | \{X_k, k \neq i, j\}$$

D Proof: consider
$$x = (x_i, x_{-i})$$
 and $\Theta = \begin{bmatrix} 0 & \theta_i^T \\ \theta_i & \Theta_{-i} \end{bmatrix}$

$$P(x_{i} = a | \boldsymbol{x}_{-i}) = \frac{P(x_{i} = a, \boldsymbol{x}_{-i})}{P(x_{i} = a, \boldsymbol{x}_{-i}) + P(x_{i} = -a, \boldsymbol{x}_{-i})} \quad \forall a \in \{-1, 1\}$$
$$= \frac{1}{1 + P(x_{i} = -a, \boldsymbol{x}_{-i}) / P(x_{i} = a, \boldsymbol{x}_{-i})} = \gamma$$

Use the Ising model

$$\gamma = \frac{\exp(-a\theta_i^T x_{-i} + \frac{1}{2}x_{-i}^T \Theta_{-i}x_{-i} - Z(\Theta))}{\exp(a\theta_i^T x_{-i} + \frac{1}{2}x_{-i}^T \Theta_{-i}x_{-i} - Z(\Theta))} = \exp\left(-2a\theta_i^T x_{-i}\right)$$

Plug
$$\gamma$$
 in the expression above

$$P(x_i = a | \boldsymbol{x}_{-i}) = \frac{1}{1 + \exp(-2a\boldsymbol{\theta}_i^T \boldsymbol{x}_{-i})} = \frac{1}{1 + \exp(-2a\sum_{j \neq i} \theta_{ij} x_j)}$$

Example: Image segmentation

Use 2-D HMM (Ising as hidden layer) to infer "meaning" of image pixels



Hidden layer: Pixel Class (water, sky, etc)

Inference via Gaussian field approximation

- Exact inference NP-hard
- □ Use surrogate continuous-values Gaussian random field: $\tilde{\mathbf{x}}_{s^c} \sim \mathcal{N}(\mu_{s^c}, \Theta_{s^c, s^c})$ > Compute exact Harmonic solution: $\mu_{s^c} = -\Theta_{s^c, s^c}^{-1}\Theta_{s^c, s^c}\mathbf{X}_s$
- **D** Predictor of unknown labels via GMRF mean: $\hat{x}_i = \begin{cases} 1 & \mu_i > 1/2 \\ 0 & \text{else} \end{cases} \quad \forall i \in S^c$
- **D** Approximation of marginal posteriors: $p(x_i = 1 | \mathbf{x}_s) \approx \mathbb{E}[\tilde{x}_i] = \mu_i \quad \forall i \in \mathcal{S}^c$
- Random walk interpretation
 - Imagine particle performing a random walk on (unobserved) graph
 - > Let normalized Laplacian $\mathbf{P} = \mathbf{D}^{-1} \mathbf{\Theta}_{s^c,s^c}$ be transition probability matrix
 - Observed variables act as sink nodes where the walk ends
 - Starting from node i, probability that walk ends in +1 node is μ_i

Inference via MCMC

Collect samples $\{\mathbf{x}_{s^c}^{(t)}\}_{t=1}^T$ from MC with $p(\mathbf{x}_{s^c}|\mathbf{x}_s)$ as stationary distribution

- Gibbs sampler: One variable (node) sampled at every round t (the rest are fixed)
 - Exploits (sparse) conditional dependence structure of MRF
 - Observed nodes used as (fixed) boundary conditions

$$\begin{aligned} x_i^{(t)} &\sim \operatorname{Ber}(p_i^{(t)}) \\ p_i^{(t)} &= \begin{cases} \left(1 + \exp(-2\sum_{j \in N(i)} \theta_{ij} x_j^{(t-1)})\right)^{-1}, & i \in \mathcal{S}^c \\ \mathbb{I}_{\{x_i=1\}}, & i \in \mathcal{S} \end{cases} \\ p(x_i &= 1 | \mathbf{x}_{s^c}) &\approx \frac{1}{T} \sum_{t=1}^T \mathbb{I}_{\{x_i^{(t)}=1\}}, \ \forall i \in \mathcal{S} \end{cases} \end{aligned}$$

Experiments indicate Gibbs smpl offers better inference in rect. Ising models

More sophisticated MCMC methods achieve faster mixing (e.g. Wolfs algorithm)

Towards estimation: Bounding the partition function

Goal: Find $\overline{Z}(\Theta) \ge Z(\Theta)$ computable with polynomial complexity

$$Z(\Theta) = \log \left(\sum_{\boldsymbol{x} \in \{0,1\}^N} \exp(\boldsymbol{x}^T \Theta \boldsymbol{x}) \right)$$

□ Consider partition $\mathcal{B}_k \subset \{0,1\}^N$ such that $x \in \mathcal{B}_k \iff ||x||_0 = k$

$$Z(\Theta) = \log\left(\sum_{k=0}^{N} \sum_{x \in \mathcal{B}_{k}} \exp(x^{T} \Theta x)\right) \leq \log\left(\sum_{k=0}^{N} \sum_{x \in \mathcal{B}_{k}} \exp\left(\max_{x \in \mathcal{B}_{k}} x^{T} \Theta x\right)\right)$$
$$= \log\left(\sum_{k=0}^{N} {N \choose k} \exp\left(\phi(\Theta)\right)\right)$$

Computing
$$\phi(\Theta) = \max_{x \in \mathcal{B}_k} x^T \Theta x$$
 is still hard

L. El Ghaoui, A. Gueye. "A Convex Upper Bound on the Log-Partition Function for Binary Graphical Models," *Journal of Machine Learning Research*, vol. 9, pp. 485–516, Mar. 2008.

Relaxation of the bound

D Relax
$$\phi(\Theta) = \max_{x \in \mathcal{B}_k} x^T \Theta x = \max_{x \in \mathcal{B}_k} \operatorname{trace}(\Theta x x^T)$$

Add redundant constrains

$$\phi(\Theta) = \max_{X \in S_+, x \in B_k} \operatorname{trace}(\Theta X)$$

s. to $X \succeq xx^T$, $x^T x = k$, $1^T m X 1 = k^2$, diag $(X) = 1$, rank $(X) = 1$

Relax

$$\begin{split} \psi(\Theta) &= \max_{X \in \mathcal{S}_+, x} \text{ trace}(\Theta X) \\ \text{s. to } X \succeq x x^T, \ x^T x = k, \ \mathbf{1}^T X \mathbf{1} = k^2, \ \text{diag}(X) = x \end{split}$$

lacksquare Upper-bound $\psi(\Theta) \geq \phi(\Theta)$

$$\Longrightarrow Z(\Theta) \le \log\left(\sum_{k=0}^{N} {N \choose k} \exp\left(\phi(\Theta\right)\right) \le \log\left(\sum_{k=0}^{N} {N \choose k} \exp\left(\psi(\Theta\right)\right) \ge \bar{Z}(\Theta)$$

Claim: bound quality $0 \leq \overline{Z}(\Theta) - Z(\Theta) \leq \min_{\mu,\lambda} ||\Theta - \mu I - \lambda \mathbf{1} \mathbf{1}^T||_1$

Pseudo Maximum Likelihood

□ Want to solve: $\hat{\Theta}_{PML} = \arg\min_{\Theta} \overline{Z}(\Theta) - \operatorname{trace}(S\Theta) + \lambda ||\Theta||_1$

$$\hat{\Theta}_{PML} = \arg\min_{\Theta} \log\left(\sum_{k=0}^{N} {N \choose k} \exp(\psi(\Theta))\right) - \operatorname{trace}(S\Theta) + \lambda ||\Theta||_{1}$$

Dual

$$\psi(\Theta) = \min_{t,\mu,\nu,\lambda} t + \mu k + \lambda k^2 \quad \text{s. to} \begin{pmatrix} \mathsf{Diag}(\nu) + \mu I + \lambda \mathbf{1} \mathbf{1}^T - \Theta & \frac{1}{2}\nu \\ \frac{1}{2}\nu^T & t \end{pmatrix} \succeq \mathbf{0}$$

lacksquare Substituting dual $\psi(\Theta)$ above

$$\begin{split} \hat{\Theta}_{PML} &= \arg\min_{\Theta, t, \mu, \nu, \lambda} \log \left(\sum_{k=0}^{N} {N \choose k} \exp\left(t + \mu k + \lambda k^{2}\right) \right) - \operatorname{trace}(S\Theta) + \lambda ||\Theta||_{1} \\ & \text{s. to} \left(\begin{array}{c} \operatorname{Diag}(\nu) + \mu I + \lambda 11^{T} - \Theta & \frac{1}{2}\nu \\ & \frac{1}{2}\nu^{T} & t \end{array} \right) \succeq 0 \end{split}$$

Logistic regression for Θ

Goal: Estimate Θ while avoiding computation of $Z(\Theta)$

- Idea: consider node i and its connections
 - > Separate $x = (x_i, x_{-i})$
 - \succ Use x_{-i} as input and x_i as output
 - > Logistic regression \longrightarrow parametric estimation of $\log P(x_i | x_{-i})$
 - > Estimate θ_i as a byproduct
- Problem statement: re-write problem bellow for the Ising model

$$\hat{\theta}_i = \arg\min_{\theta_i} - \sum_{t=1}^T \log P_{\theta_i} \left(x_i = x_i^{(t)} | \boldsymbol{x}_{-i} = \boldsymbol{x}_{-i}^{(t)} \right) + \lambda ||\theta_i||_1$$

P. Ravikumar, M. J. Wainwright and J. Lafferty. *High-dimensional Ising model selection using* -regularized lc_{y}^{ℓ} stic regression. To appear in the Annals of Statistics. Available at http://www.eecs.berkeley.edu



Estimation of Θ

We have:

$$P(x_i|\boldsymbol{x}_{-i}) = \frac{1}{1 + \exp(-2x_i\boldsymbol{\theta}_i^T\boldsymbol{x}_{-i})}$$
$$= \frac{\exp(x_i\boldsymbol{\theta}_i^T\boldsymbol{x}_{-i})}{\exp(x_i\boldsymbol{\theta}_i^T\boldsymbol{x}_{-i}) + \exp(-x_i\boldsymbol{\theta}_i^T\boldsymbol{x}_{-i})}$$

Taking the logarithm

$$| og P(x_i | \boldsymbol{x}_{-i}) = x_i \boldsymbol{\theta}_i^T \boldsymbol{x}_{-i} - log \left(exp(x_i \boldsymbol{\theta}_i^T \boldsymbol{x}_{-i}) + exp(-x_i \boldsymbol{\theta}_i^T \boldsymbol{x}_{-i}) \right)$$

Substituting the log-likelihood

$$\hat{\theta}_{i} = \arg\min_{\theta_{i}} - \sum_{t=1}^{T} \log P\left(x_{i} = x_{i}^{(t)} | x_{-i} = x_{-i}^{(t)}\right) + \lambda ||\theta_{i}||_{1}$$

$$= \arg\min_{\theta_{i}} \sum_{t=1}^{T} \left[\log\left(\exp(x_{i}^{(t)}\theta_{i}^{T}x_{-i}^{(t)}) + \exp(-x_{i}^{(t)}\theta_{i}^{T}x_{-i}^{(t)})\right) - x_{i}^{(t)}\theta_{i}^{T}x_{-i}^{(t)}\right] + \lambda ||\theta_{i}||_{1}$$

Convex problem

Conclusions

- Graphical models
 - Modeling pdfs using conditional dependencies
 - Undirected models (MRFs) naturally modeled by graphs
 - Inference in closed form for Gaussian MRFs
 - Estimation of GMRFs as Laplacian fitting problem
 - Inference and estimation approximations for binary MRFs (Ising model)
- Possible research directions
 - Active sampling on binary MRFs using MCMC
 - Active sampling for MRF structure estimation

