Probabilistic Graphical Models

Lecture 11 – CRFs, Exponential Family

CS/CNS/EE 155 Andreas Krause

Announcements

- Homework 2 due today
- Project milestones due next Monday (Nov 9)
 - About half the work should be done
 - 4 pages of writeup, NIPS format
 - http://nips.cc/PaperInformation/StyleFiles

So far

- Markov Network Representation
 - Local/Global Markov assumptions; Separation
 - Soundness and completeness of separation
- Markov Network Inference
 - Variable elimination and Junction Tree inference work exactly as in Bayes Nets
- How about Learning Markov Nets?

MLE for Markov Nets

Log likelihood of the data

$$\log P(D(0) = \sum_{e} \log P(x^{(e)}|\theta)$$

$$= \sum_{e=i}^{n} \log \frac{1}{2} \prod_{i} \gamma_{i} (C_{i}^{(e)})$$

$$= \sum_{e=i}^{n} \log \gamma_{i} (C_{i}^{(e)}) - m \log 2$$

$$= m \sum_{i} \sum_{e=i}^{n} P(C_{i}) \log \gamma_{i} (C_{i}) - m \log 2$$

$$= \log \sum_{i} \frac{1}{2} \gamma_{i} (C_{i}) \log 2 = \log \sum_{i} \frac{1}{2} \gamma_{i} (C_{i})$$

$$= \sum_{i} \sum_{i} \frac{1}{2} \gamma_{i} (C_{i}) \log 2 = \log \sum_{i} \frac{1}{2} \gamma_{i} (C_{i})$$

$$= \sum_{i} \frac{1}{2} \gamma_{i} (C_{i}) \log 2 = \log \sum_{i} \frac{1}{2} \gamma_{i} (C_{i})$$

Log-likelihood doesn't decompose

Log likelihood

•
$$I(D \mid \theta) = m \sum_{i} \sum_{c_i} \hat{P}(c_i) \log \psi_i(c_i) - m \log Z(\theta)$$

 $decomposes nicely does not decompose
 $I(D \mid \theta) \text{ is concave function! by } P(0,0)$
No local optima!
Gradient ascent von't get stuck!$

• Log Partition function log $Z(\theta)$ doesn't decompose

Computing the derivative

С

D

G

Derivative



• Computing $P(c_i | \theta)$ requires inference! Can do this Using VE

Can optimize using conjugate gradient etc.

S

Alternative approach: Iterative Proportional Fitting (IPF)

At optimum, it must hold that

$$(c_{i}) = (c_{i}) = (c_{i})$$

Must recompute parameters every iteration

Parameter learning for log-linear models

- Feature functions $\phi_i(C_i)$ defined over cliques
- Log linear model over undirected graph G
 - Feature functions $\phi_1(C_1),...,\phi_k(C_k)$
 - Domains C_i can overlap
- Joint distribution

$$P(X_1, \dots, X_n) = \frac{1}{Z} \exp\left(\sum_i w_i^T \phi_i(C_i)\right)$$

How do we get weights w_i?

Optimizing parameters

Gradient of log-likelihood

$$\frac{\partial \log P(\mathcal{D} \mid w)}{\partial w_i} = m \sum_{\mathbf{c}_i} \hat{P}(\mathbf{c}_i) \phi_i(\mathbf{c}_i) - m \sum_{\mathbf{c}_i} P(\mathbf{c}_i \mid w) \phi_i(\mathbf{c}_i)$$

$$\underbrace{\hat{P}(\mathbf{c}_i)}_{\text{F}(\mathbf{c}_i)} \qquad \underbrace{\hat{P}(\mathbf{c}_i \mid w)}_{\text{F}(\mathbf{c}_i)} \quad \underbrace{\hat{P}(\mathbf{c}_i \mid$$

• Thus, w is MLE $\Leftrightarrow \ \hat{\mathbb{E}}[\phi_i] = \mathbb{E}_w[\phi_i]$

Regularization of parameters

PW

Put prior on parameters w

$$\frac{\partial \log P(\mathcal{D} \mid w) P(w)}{\partial w_i} = m \sum_{\mathbf{c}_i} \hat{P}(\mathbf{c}_i) \phi_i(\mathbf{c}_i) - m \sum_{\mathbf{c}_i} P(\mathbf{c}_i \mid w) \phi_i(\mathbf{c}_i) + \underbrace{\frac{\partial \log P(w)}{\partial w_i}}_{\mathcal{L}}$$

Prior:
$$P(w) = \mathcal{N}(w; 0, \mathbb{I}) \propto \exp(-\sum_{i} w_{i}^{2})$$

$$(x) = \log P(w) = -\sum_{i}^{2} w_{i}^{2}$$

$$\frac{\partial}{\partial w_{i}} \log P(w) = -2w_{i}^{2}$$

Summary: Parameter learning in MN

- MLE in BN is easy (score decomposes)
- MLE in MN requires inference (score doesn't decompose)
- Can optimize using gradient ascent or IPF

Generative vs. discriminative models

- Often, want to predict Y for inputs X
 - Bayes optimal classifier: Predict according to P(Y | X)
- Generative model
 - Model P(Y), P(X|Y)
 - Use Bayes' rule to compute P(Y | X)
- Discriminative model
 - Model P(Y | X) directly!
 - Don't model distribution P(X) over inputs X
 - Cannot "generate" sample inputs
 - Example: Logistic regression

Example: Logistic Regression



Log-linear conditional random field

- Define log-linear model over outputs Y
 - No assumptions about inputs X
- Feature functions $\phi_i(C_i,x)$ defined over cliques and inputs $C_i \subseteq \mathcal{V}$
- Joint distribution

$$P(Y_1, \dots, Y_n \mid x) = \frac{1}{Z(x_i)} \exp\left(\sum_i w_i^T \phi_i(C_i, x)\right)$$

Example: CRFs in NLP



Mrs. Greene spoke today in New York. Green chairs the finance committee

Classify into Person, Location or Other

Example: CRFs in vision





Parameter learning for log-linear CRF

$$P(Y_1, \dots, Y_n \mid x) = \frac{1}{Z(x)} \exp\left(\sum_i w_i^T \phi_i(C_i, x)\right)$$

• Conditional log-likelihood of data $\log P(D_y | w_i, D_x) = \sum_{i} \sum_{i} w_i^T \phi_i (C_{i}, e^{(\ell)}) - \sum_{i} \log \frac{1}{2} (e^{(\ell)})$

Can maximize using conjugate gradient

Gradient of conditional log-likelihood

Partial derivative

$$\frac{\partial \log P(\mathcal{D}_Y \mid w, \mathcal{D}_X)}{\partial w_i} = \sum_j \left[\phi_i(\mathbf{c}_i^{(j)}, x^{(j)}) + \sum_{\mathbf{c}_i} P(\mathbf{c}_i \mid w, x^{(j)}) \phi_i(\mathbf{c}_i, x^{(j)}) \right]$$

Requires one inference per feature and per data point
Can be very expensive
Con do "pseudo" likelihood ost. (approximate)

Can optimize using conjugate gradient

Exponential Family Distributions

Distributions for log-linear models

$$P(X_1, \dots, X_n) = \frac{1}{Z} \exp\left(\sum_i w_i^T \phi_i(C_i)\right)$$

- More generally: Exponential family distributions $P(x) = h(x) \exp\left(w^T \phi(x) - A(w)\right)$
 - h(x): Base measure < Offer constant</p>
 - w: natural parameters
 - $\phi(\mathbf{x})$: Sufficient statistics

A(w): log-partition function
 A(w) = log & (w)
 Hereby x can be continuous (defined over any set)

Examples

• Exp. Family:

$$P(x) = h(x) \exp\left(w^T \phi(x) - A(w)\right)$$

Gaussian distribution

h(x): Base measure w: natural parameters ϕ (x): Sufficient statistics A(w): log-partition function

$$P(x) = \frac{1}{\sqrt{2\pi\sigma^{2}}} \exp\left(-\frac{(x-\mu)^{2}}{\sqrt{2\sigma^{2}}}\right)$$

$$h(x) = \frac{1}{\sqrt{2\pi}} \qquad -\frac{x^{2}}{2\sigma^{2}} + \frac{x\mu}{\sigma^{2}} - \frac{\mu^{2}}{2\sigma^{2}}$$

$$\phi(x) = \left[-\frac{x^{2}}{2}, x\right], \qquad A(w) = \frac{\mu^{2}}{2\sigma^{2}} - \log \sigma^{2}$$

$$w = \left[-\frac{1}{\sigma^{2}}, \frac{M}{\sigma^{2}}\right]$$

 Other examples: Multinomial, Poisson, Exponential, Gamma, Weibull, chi-square, Dirichlet, Geometric, ...

Moments and gradients

 $P(x) = h(x) \exp\left(w^T \phi(x) - A(w)\right)$

 Correspondence between moments and log-partition function (just like in log-linear models)

$$\frac{\partial A(w)}{\partial w_i} = \int p(x \mid w)\phi_i(x)dx = \mathbb{E}[\phi_i \mid w]$$

$$\frac{\partial^2 A(w)}{\partial w_i \partial w_j} = Cov(\phi_i, \phi_j \mid w)$$

- Can compute moments from derivatives, and derivatives from moments!
- MLE 🗇 moment matching

Recall: Conjugate priors

- Consider parametric families of prior distributions:
 - $P(\theta) = f(\theta; \alpha)$
 - α is called "hyperparameters" of prior
- A prior $P(\theta) = f(\theta; \alpha)$ is called **conjugate** for a likelihood function $P(D \mid \theta)$ if $P(\theta \mid D) = f(\theta; \alpha')$
 - Posterior has same parametric form
 - Hyperparameters are updated based on data D
- Obvious questions (answered later):
 - How to choose hyperparameters??
 - Why limit ourselves to conjugate priors??

Conjugate priors in Exponential Family

$$P(x \mid w) = h(x) \exp(w^T \phi(x) - A(w))$$

Any exponential family likelihood has a conjugate prior $P(w|\alpha,\beta) = \exp(\alpha^T w - \beta A(w) - B(\alpha,\beta))$

Maximum Entropy interpretation

Theorem: Exponential family distributions maximize the entropy over all distributions satisfying

$$\max_{p} - \int p(x) \log p(x) dx$$

such that
$$t_i = \mathbb{E}[\phi_i(x)] = \int p(x)\phi_i(x)dx$$

Summary exponential family

Distributions of the form

$$P(x) = h(x) \exp\left(w^T \phi(x) - A(w)\right)$$

- Most common distributions are exponential family
 - Multinomial, Gaussian Poisson, Exponential, Gamma, Weibull, chisquare, Dirichlet, Geometric, ...
 - Log-linear Markov Networks
- All exponential family distributions have conjugate prior in EF
- Moments of sufficient stats = derivatives of log-partition function
- Maximum Entropy distributions ("most uncertain" distributions with specified expected sufficient statistics)

Exponential family graphical models

- So far, only defined graphical models over discrete variables.
- Can define GMs over continuous distributions!
- For exponential family distributions:

$$p(X_1, \dots, X_n) = \prod_i h_i(C_i) \exp\left(\sum_i w_i^T \phi_i(C_i) - A(w)\right)$$
$$\exp A(w) = \int \int \dots \int \prod_i h_i(C_i) \exp\left(\sum_i w_i^T \phi_i(C_i) - A(w)\right) dx_1 \dots dx_n$$

- Can do much of what we discussed (VE, JT, parameter learning, etc.) for such exponential family models
- Important example: Gaussian Networks

Gaussian distribution



- σ = Standard deviation
- = mean μ

Bivariate Gaussian distribution

$$\frac{1}{2\pi\sqrt{|\Sigma|}}\exp\left(-\frac{1}{2}(x-\mu)^T \underline{\Sigma}^{-1}(x-\mu)\right) \qquad \Sigma = \begin{pmatrix}\sigma_1^2 & \sigma_{12}\\\sigma_{21} & \sigma_2^2\end{pmatrix} \quad \mu = \begin{pmatrix}\mu_1\\\mu_2\end{pmatrix}$$





28

Multivariate Gaussian distribution

$$\mathcal{N}(x;\Sigma,\mu) = \frac{1}{(2\pi)^{n/2}\sqrt{|\Sigma|}} \exp\left(-\frac{1}{2}(x-\mu)^T \Sigma^{-1}(x-\mu)\right)$$

$$\Sigma = \begin{pmatrix} \sigma_1^2 & \sigma_{12} & \dots & \sigma_{1n} \\ \vdots & & \vdots \\ \sigma_{n1} & \sigma_{n2} & \dots & \sigma_n^2 \end{pmatrix} \qquad \mu = \begin{pmatrix} \mu_1 \\ \mu_2 \\ \vdots \\ \mu_n \end{pmatrix}$$
Solution The second second

Joint distribution over n random variables P(X₁,...X_n)

- $\sigma_{jk} = E[(X_j \mu_j)(X_k \mu_k)] < Cov(X_j, X_k)$
- X_j and X_k independent $\Leftrightarrow \sigma_{jk}=0$

171

Marginalization

- Suppose $(X_1,...,X_n) \sim N(\mu, \Sigma)$

• What is $P(X_1)$? • $P(X_1) = \iiint \left(\begin{array}{cc} \sigma_1^2 & \sigma_{12} \\ \vdots \\ \sigma_{n1} & \sigma_{n2} \end{array} \right) \left(\begin{array}{cc} \sigma_1^2 & \sigma_{1n} \\ \vdots \\ \sigma_{n1} & \sigma_{n2} \end{array} \right) \left(\begin{array}{cc} \mu_1 \\ \mu_2 \\ \vdots \\ \mu_n \end{array} \right) \left(\begin{array}{cc} \mu_1 \\ \mu_2 \\ \vdots \\ \mu_n \end{array} \right) \left(\begin{array}{cc} \mu_1 \\ \mu_2 \\ \vdots \\ \mu_n \end{array} \right) \left(\begin{array}{cc} \mu_1 \\ \mu_2 \\ \vdots \\ \mu_n \end{array} \right) \left(\begin{array}{cc} \mu_1 \\ \mu_2 \\ \vdots \\ \mu_n \end{array} \right) \left(\begin{array}{cc} \mu_1 \\ \mu_2 \\ \vdots \\ \mu_n \end{array} \right) \left(\begin{array}{cc} \mu_1 \\ \mu_2 \\ \vdots \\ \mu_n \end{array} \right) \left(\begin{array}{cc} \mu_1 \\ \mu_2 \\ \vdots \\ \mu_n \end{array} \right) \left(\begin{array}{cc} \mu_1 \\ \mu_2 \\ \vdots \\ \mu_n \end{array} \right) \left(\begin{array}{cc} \mu_1 \\ \mu_2 \\ \vdots \\ \mu_n \end{array} \right) \left(\begin{array}{cc} \mu_1 \\ \mu_2 \\ \vdots \\ \mu_n \end{array} \right) \left(\begin{array}{cc} \mu_1 \\ \mu_2 \\ \vdots \\ \mu_n \end{array} \right) \left(\begin{array}{cc} \mu_1 \\ \mu_2 \\ \vdots \\ \mu_n \end{array} \right) \left(\begin{array}{cc} \mu_1 \\ \mu_2 \\ \vdots \\ \mu_n \end{array} \right) \left(\begin{array}{cc} \mu_1 \\ \mu_2 \\ \vdots \\ \mu_n \end{array} \right) \left(\begin{array}{cc} \mu_1 \\ \mu_2 \\ \vdots \\ \mu_n \end{array} \right) \left(\begin{array}{cc} \mu_1 \\ \mu_2 \\ \vdots \\ \mu_n \end{array} \right) \left(\begin{array}{cc} \mu_1 \\ \mu_2 \\ \vdots \\ \mu_n \end{array} \right) \left(\begin{array}{cc} \mu_1 \\ \mu_2 \\ \vdots \\ \mu_n \end{array} \right) \left(\begin{array}{cc} \mu_1 \\ \mu_2 \\ \vdots \\ \mu_n \end{array} \right) \left(\begin{array}{cc} \mu_1 \\ \mu_2 \\ \vdots \\ \mu_n \end{array} \right) \left(\begin{array}{cc} \mu_1 \\ \mu_2 \\ \vdots \\ \mu_n \end{array} \right) \left(\begin{array}{cc} \mu_1 \\ \mu_2 \\ \vdots \\ \mu_n \end{array} \right) \left(\begin{array}{cc} \mu_1 \\ \mu_2 \\ \vdots \\ \mu_n \end{array} \right) \left(\begin{array}{cc} \mu_1 \\ \mu_2 \\ \vdots \\ \mu_n \end{array} \right) \left(\begin{array}{cc} \mu_1 \\ \mu_2 \\ \vdots \\ \mu_n \end{array} \right) \left(\begin{array}{cc} \mu_1 \\ \mu_2 \\ \vdots \\ \mu_n \end{array} \right) \left(\begin{array}{cc} \mu_1 \\ \mu_2 \\ \vdots \\ \mu_n \end{array} \right) \left(\begin{array}{cc} \mu_1 \\ \mu_2 \\ \vdots \\ \mu_n \end{array} \right) \left(\begin{array}{cc} \mu_1 \\ \mu_2 \\ \mu_2 \\ \vdots \\ \mu_n \end{array} \right) \left(\begin{array}{cc} \mu_1 \\ \mu_2 \\ \mu_2 \\ \vdots \\ \mu_n \end{array} \right) \left(\begin{array}{cc} \mu_1 \\ \mu_2 \\ \mu_2 \\ \vdots \\ \mu_n \end{array} \right) \left(\begin{array}{cc} \mu_1 \\ \mu_2 \\ \mu_2 \\ \vdots \\ \mu_n \end{array} \right) \left(\begin{array}{cc} \mu_1 \\ \mu_2 \\ \mu_2 \\ \vdots \\ \mu_n \end{array} \right) \left(\begin{array}{cc} \mu_1 \\ \mu_2 \\ \mu_2 \\ \mu_2 \\ \mu_n \end{array} \right) \left(\begin{array}{cc} \mu_1 \\ \mu_2 \\ \mu_2 \\ \mu_2 \\ \mu_n \end{array} \right) \left(\begin{array}{cc} \mu_1 \\ \mu_2 \\ \mu_2 \\ \mu_2 \\ \mu_2 \\ \mu_2 \\ \mu_2 \end{array} \right) \left(\begin{array}{cc} \mu_1 \\ \mu_2 \\$

• More generally: Let $A = \{i_1, \dots, i_k\} \subseteq \{1, \dots, N\}$

•
$$X_A \sim N(\mu_A, \Sigma_{AA})$$

 $\Sigma_{AA} = \begin{pmatrix} \sigma_{i_1}^2 & \sigma_{i_1 i_2} & \dots & \sigma_{i_1 i_k} \\ \vdots & & & \vdots \\ \sigma_{i_k i_1} & \sigma_{i_k i_2} & \dots & \sigma_{i_k}^2 \end{pmatrix} \qquad \mu_A = \begin{pmatrix} \mu_{i_1} \\ \mu_{i_2} \\ \vdots \\ \mu_{i_k} \end{pmatrix}$

Conditioning

- Suppose $(X_1,...,X_n) \sim N(\mu, \Sigma)$
- Decompose as (X_A, X_B)
- What is $P(X_A | X_B)$??

$$\sum = \begin{pmatrix} \sum_{AA} & \sum_{AB} \\ \sum_{BA} & \sum_{BB} \end{pmatrix}$$

•
$$P(X_A = x_A | X_B = x_B) = N(x_A; \mu_{A|B}, \Sigma_{A|B})$$
 where

$$\mu_{A|B=g} = \mu_A + \sum_{AB} \sum_{BB}^{-1} (x_B - \mu_B)$$

$$\Sigma_{A|B} = \Sigma_{AA} - \Sigma_{AB} \sum_{BB}^{-1} \Sigma_{BA}$$
• Computable using linear algebra!

Conditioning



Conditional linear Gaussians

$$\mu_{A|B} = \mu_{A} + \underbrace{\sum_{AB} \sum_{BB}^{-1} (x_{B} - \mu_{B})}_{\bigvee}$$

$$\sum_{A|B} = \sum_{AA} - \sum_{AB} \sum_{BB}^{-1} \sum_{BA}$$

$$P(X_{A}|X_{B}) = \sqrt{(X_{A} \ i \ M_{AB} \ i \ \mathcal{D}_{AB})}$$

$$x_{A} = M_{A} + \Im x_{B} - \Im M_{B} + \varepsilon \quad \varepsilon \sim \mathcal{M}_{0}, \sum_{AB})$$

$$= \Im x_{B} + b + \varepsilon$$

$$\Im = \sum_{AB} \sum_{BB}^{-1} b$$

$$b = M_{A} - \Im M_{B}$$

Canonical representation of Gaussians

Canonical Representation

$$p(X_1, \dots, X_n) = \frac{1}{(2\pi)^{n/2} \sqrt{|\Sigma|}} \exp\left(-\frac{1}{2}(x-\mu)^T \Sigma^{-1}(x-\mu)\right)$$
$$\propto \exp(\eta^T x - \frac{1}{2}x^T \Lambda x) = \operatorname{sp}\left(\sum_{i} \gamma_i X_i - \frac{1}{2}\sum_{i} \lambda_i Y_i \lambda_i\right)$$

- Multivariate Gaussians in exponential family!
- Standard vs canonical form:

$$\mu = \Lambda^{-1} \eta$$

 $\Sigma = \Lambda^{-1}$

- In standard form: Marginalization is easy
- Will see: In canonical form, multiplication/conditioning is easy!

Gaussian Networks

$$p(X_{1},...,X_{n}) \propto \exp\left(-\frac{1}{2}\sum_{i,j}\lambda_{i,j}x_{i}x_{j} + \sum_{i}\eta_{i}x_{i}\right)$$

$$= \mathcal{L}_{i,j}\left(\sum_{\lambda_{i,j}}\lambda_{i,j}\phi_{i,j}(x_{i},y_{j}) + \sum_{\lambda_{i}}\eta_{i}\phi_{i}(x_{i},y_{j})\right)$$

$$\varphi_{i,j}(x_{i,j}y_{j})$$

$$\varphi_{i,j}(x_{i,j}y_{j}) = -\frac{1}{2}x_{i}x_{j}$$

$$\varphi_{i}(x_{i}) = x_{i}$$
No elge between $X_{i}X_{j}$

$$\chi_{i,j} = 0$$

Zeros in precision matrix Λ indicate missing edges in log-linear model!

Inference in Gaussian Networks

- Can compute marginal distributions in O(n³)!
- For large numbers n of variables, still intractable
- If Gaussian Network has low treewidth, can use variable elimination / JT inference!
- Need to be able to multiply and marginalize factors!

$$g = \int_{x_i} \prod_j f_j$$

$$\begin{aligned} & \text{Multiplying factors in Gaussians} \\ & P(X_{A}) \leftarrow \mathcal{N} \{ k_{R} ; \Lambda_{i}, \gamma_{i} \} \xrightarrow{(A_{i} \models k, B) = m} \Lambda_{i} \in \mathbb{R}^{k \times k} \\ & P(X_{B} \mid X_{A}) \neq \mathcal{N} (x_{B}, x_{A}; \Lambda_{2}, \gamma_{2}) \leftarrow \frac{CLG \text{ representation}}{\Lambda_{2} \in \mathbb{R}^{m \times m}} \\ & P(X_{A}, X_{B}) = P(X_{A}) P(X_{B} \mid X_{A}) \\ & \downarrow \exp((x_{A}^{\top} \Lambda_{i} \times_{A} + \gamma_{i}^{\top} \times_{A}) \exp(((x_{A}, x_{B})^{\top} \Lambda_{2} \cdot k_{B} \times_{B}) + \gamma_{i}^{\top}) \\ & = \mathcal{N} (x_{A}, x_{B}; \Lambda_{i}, \gamma) \\ & \Lambda = \Lambda_{i} + \Lambda_{2} = (\Lambda_{B}) + (\Lambda_{2}) \\ & \gamma = \gamma_{i} + \gamma_{2} \end{aligned}$$

Marginalizing in canonical form

Recall conversion formulas

•
$$\mu$$
 = Λ -1 η

- $\Sigma = \Lambda^{-1}$
- Marginal distribution

$$\eta_A^m = \eta_A - \Lambda_{AB} \Lambda_{BB}^{-1} \eta_B$$
$$\Lambda_{AA}^m = \Lambda_{AA} - \Lambda_{AB} \Lambda_{BB}^{-1} \Lambda_{BA}$$

Variable elimination



 In Gaussian Markov Networks, Variable elimination = Gaussian elimination (fast for low bandwidth = low treewidth matrices)

Tasks

Read Koller & Friedman Chapters 4.6.1, 8.1-8.3