

Readings:

K&F: 19.1, 19.2, 19.3, 19.4

Parameter learning in Markov nets

Graphical Models – 10708

Carlos Guestrin

Carnegie Mellon University

November 17th, 2008

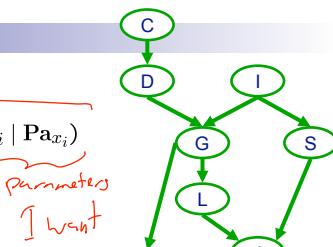
10-708 – ©Carlos Guestrin 2006-2008

1

Learning Parameters of a BN

- Log likelihood decomposes:

$$\ell(\mathcal{D} : \theta) = \log P(\mathcal{D} | \theta) = m \sum_i \sum_{x_i, \text{Pa}_{x_i}} \hat{P}(x_i, \text{Pa}_{x_i}) \log P(x_i | \text{Pa}_{x_i})$$



- Learn each CPT independently

- Use counts

$$P(x_i | \text{Pa}_{x_i} = u) \stackrel{\text{MLE}}{=} \frac{\text{Count}(x_i = x_i, \text{Pa}_{x_i} = u)}{\text{Count}(\text{Pa}_{x_i} = u)}$$

10-708 – ©Carlos Guestrin 2006-2008

2

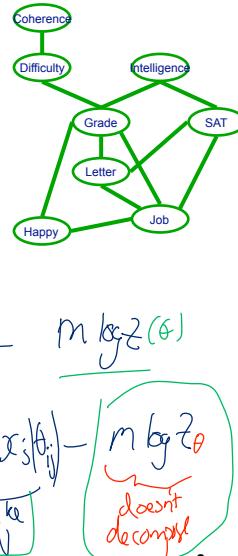
Log Likelihood for MN

$$\log Z_{\theta} = \log \sum_{ij} \prod_{k=1}^K \phi_{ijk}(x_i^{(k)}, x_j^{(k)} | \theta_{ijk})$$

- Log likelihood of the data:

$$\begin{aligned}
 l(D; \theta) &= \log P(D|\theta) \stackrel{iid}{=} \sum_k \log P(x^{(k)}|\theta) \\
 &= \sum_k \log \frac{1}{Z} \prod_{ij} \phi_{ijk}(x_i^{(k)}, x_j^{(k)} | \theta) \\
 &= \sum_k \sum_{ij} \log \phi_{ijk}(x_i^{(k)}, x_j^{(k)} | \theta) - \sum_{k=1}^m \log Z(\theta) \\
 &= \sum_{(i,j)} \text{Count}(x_i=x_i, x_j=x_j) \log \phi_{ijk}(x_i=x_i, x_j=x_j | \theta) - m \log Z(\theta) \\
 &= m \sum_{(i,j)} \hat{P}(x_i=x_i, x_j=x_j) \log \phi_{ijk}(x_i=x_i, x_j=x_j | \theta) - m \log Z(\theta)
 \end{aligned}$$

images
 decomposes nicely other factors just like BN
 doesn't decompose



10-708 - ©Carlos Guestrin 2006-2008

3

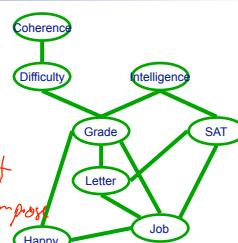
Log Likelihood doesn't decompose for MNs

$$\hat{P}(u) = \frac{\text{Count}(U=u)}{m}$$

- Log likelihood:

$$\ell(D : \theta) = \log P(D | \theta, \mathcal{G}) = m \sum_i \sum_{c_i} \hat{P}(c_i) \log \psi_i(c_i | \theta) - m \log Z(\theta)$$

decomposes nicely
 doesn't decompose



- A concave problem
- good news
- Can find global optimum!!



- Term $\log Z$ doesn't decompose!!
- bad news

10-708 - ©Carlos Guestrin 2006-2008

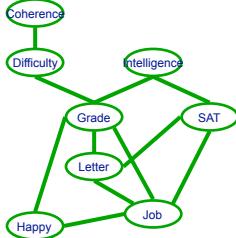
4

Derivative of Log Likelihood for MNs

$\hat{P}(u) = \frac{\text{Count}(U = u)}{m}$

$$\ell(\mathcal{D} : \theta) = \log P(\mathcal{D} | \theta, \mathcal{G}) = m \sum_i \sum_{c_i} \hat{P}(c_i) \log \psi_i(c_i) - m \log Z$$

$$\frac{\partial \ell(\mathcal{D} : \theta)}{\partial \psi_i(c_i)} = m \hat{P}(c_i) \frac{\partial \log \psi_i(c_i)}{\partial \psi_i(c_i)} - m \frac{\partial \log Z}{\partial \psi_i(c_i)}$$

$$= m \frac{\hat{P}(c_i)}{\psi_i(c_i)} - m \frac{\partial \log Z}{\partial \psi_i(c_i)}$$


10-708 - ©Carlos Guestrin 2006-2008

5

Derivative of Log Likelihood for MNs 2

$\hat{P}(u) = \frac{\text{Count}(U = u)}{m}$

$$\ell(\mathcal{D} : \theta) = \log P(\mathcal{D} | \theta, \mathcal{G}) = m \sum_i \sum_{c_i} \hat{P}(c_i) \log \psi_i(c_i) - m \log Z$$

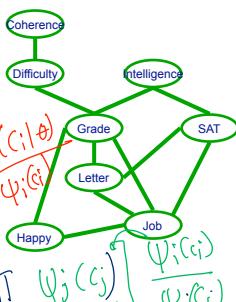
$$\frac{\partial \ell}{\partial \psi_i(c_i)} = m \frac{\hat{P}(c_i)}{\psi_i(c_i)} - m \frac{\partial \log Z}{\partial \psi_i(c_i)}$$

$$\frac{\partial \log Z}{\partial \psi_i(c_i)} = \frac{\partial Z}{\partial \psi_i(c_i)} = \frac{1}{Z} \sum_{x \sim c_i} \prod_j \psi_j(c_j) = \frac{P(c_i | \theta)}{\psi_i(c_i)}$$

$\frac{1}{Z} \sum_{x \sim c_i} \prod_j \psi_j(c_j)$
 if x not consistent with c_i
 $= \left(\sum_{x \sim c_i} \prod_{j \neq i} \psi_j(c_j) \right) \frac{1}{\psi_i(c_i)}$

$$\frac{\partial Z}{\partial \psi_i(c_i)} = \frac{\partial}{\partial \psi_i(c_i)} \left(\sum_{x \sim c_i} \prod_j \psi_j(c_j) \right) = \frac{\partial}{\partial \psi_i(c_i)} \sum_{x \sim c_i} \prod_{j \neq i} \psi_j(c_j)$$

consistent with c_i



10-708 - ©Carlos Guestrin 2006-2008

6

Derivative of Log Likelihood for MNs

$$\hat{P}(u) = \frac{\text{Count}(U = u)}{m}$$

$$\ell(\mathcal{D} : \theta) = \log P(\mathcal{D} | \theta, \mathcal{G}) = m \sum_i \sum_{c_i} \hat{P}(c_i) \log \psi_i(c_i) - m \log Z$$

- Derivative:

$$\frac{\partial \ell}{\partial \psi_i(c_i)} = \frac{m \hat{P}(c_i)}{\psi_i(c_i)} - \frac{m P_F^\psi(c_i)}{\psi_i(c_i)}$$

- Computing derivative requires inference:

$$\frac{\partial \ell}{\partial \psi(G=A, J=t)} = m \underbrace{\frac{\text{Count}(G=A, J=t)}{m}}_{P(G=A, J=t)} - m \underbrace{\frac{P_F^\psi(G=A, J=t)}{\psi(G=A, J=t)}}_{\psi(G=A, J=t)}$$

- Can optimize using gradient ascent

- Common approach
- Conjugate gradient, Newton's method, ...

- Let's also look at a simpler solution

10-708 - ©Carlos Guestrin 2006-2008

7

Iterative Proportional Fitting (IPF)

$$\hat{P}(u) = \frac{\text{Count}(U = u)}{m}$$

$$\frac{\partial \ell}{\partial \psi_i(c_i)} = \frac{m \hat{P}(c_i)}{\psi_i(c_i)} - \frac{m P_F^\psi(c_i)}{\psi_i(c_i)} = 0$$

- Setting derivative to zero:

$$\frac{m \hat{P}(c_i)}{\psi_i(c_i)} - \frac{m P_F^\psi(c_i)}{\psi_i(c_i)} = 0$$

- Fixed point equation:

$$\psi_i(c_i) = \psi_i^{(t)} \frac{\hat{P}(c_i)}{P_F^\psi(c_i)}$$

must initialize $\psi_i(c_i) > 0$

- Iterate and converge to optimal parameters

- Each iteration, must compute:

$$\psi_i^{(t+1)}(c_i) \leftarrow \frac{\hat{P}(c_i)}{P_F^\psi(c_i)} \psi_i^{(t)}$$

Compute
 $P_F^\psi(c_i)$
 at iteration t

10-708 - ©Carlos Guestrin 2006-2008

8

Log-linear Markov network (most common representation)

- **Feature** is some function $\phi[\mathbf{D}]$ for some subset of variables \mathbf{D}
 - e.g., indicator function $\mathbf{D} \subseteq \{X_1, \dots, X_n\}$
- **Log-linear model** over a Markov network H :
 - a set of features $\phi_1[\mathbf{D}_1], \dots, \phi_k[\mathbf{D}_k]$
 - each \mathbf{D}_i is a subset of a clique in H
 - two ϕ 's can be over the same variables
 - a set of weights w_1, \dots, w_k
 - usually learned from data
 - $P(X_1, \dots, X_n) = \frac{1}{Z} \exp \left[\sum_{i=1}^k w_i \phi_i(\mathbf{D}_i) \right]$

10-708 - ©Carlos Guestrin 2006-2008

9

Learning params for log linear models – Gradient Ascent

$$P(X_1, \dots, X_n) = \frac{1}{Z} \exp \left[\sum_{i=1}^k w_i \phi_i(\mathbf{D}_i) \right]$$

- Log-likelihood of data:

$$\begin{aligned} Q(D; w) &= \log \prod_{j=1}^m \frac{1}{Z} \exp \left[\sum_{i=1}^k w_i \phi_i(d_i^{(j)}) \right] \\ &= \sum_{j=1}^m \log \frac{1}{Z} \exp \left[\sum_{i=1}^k w_i \phi_i(d_i^{(j)}) \right] \end{aligned}$$

- Compute derivative & optimize
 - usually with conjugate gradient ascent

10-708 - ©Carlos Guestrin 2006-2008

10

Derivative of log-likelihood 1 – log-linear models

$$P(X_1, \dots, X_n) = \frac{1}{Z} \exp \left[\sum_{i=1}^k w_i \phi_i(\mathbf{D}_i) \right]$$

$$\ell(\mathcal{D} : \mathbf{w}) = \log P(\mathcal{D} | \mathbf{w}, \mathcal{G}) = \sum_{j=1}^m \log \frac{1}{Z} \exp \left[\sum_{i=1}^k w_i \phi_i(\mathbf{d}_i^{(j)}) \right]$$

$$\begin{aligned} \frac{\partial \ell}{\partial w_i} &= \frac{\partial}{\partial w_i} \left[\sum_{j=1}^m \sum_{i=1}^k w_i \phi_i(\mathbf{d}_i^{(j)}) - m \log Z \right] \\ &= \sum_{j=1}^m \phi_i(\mathbf{d}_i^{(j)}) - m \frac{\partial \log Z}{\partial w_i} \\ &\quad \text{--- } \underbrace{\sum_{d_i} \hat{P}(\mathbf{d}_i) \phi_i(\mathbf{d}_i)}_{\hat{E}_{\mathcal{D}}[\phi_i]} \text{ --- } \underbrace{m \phi_i(\mathbf{d}_i)}_{\text{will be expected feature value in model}} \\ &\quad \text{--- } \underbrace{\hat{E}_{\mathcal{D}}[\phi_i]}_{\text{expected feature value in data}} \end{aligned}$$

10-708 - ©Carlos Guestrin 2006-2008

11

Derivative of log-likelihood 2 – log-linear models

$$P(X_1, \dots, X_n) = \frac{1}{Z} \exp \left[\sum_{i=1}^k w_i \phi_i(\mathbf{D}_i) \right]$$

$$\frac{\partial \ell(\mathcal{D} : \mathbf{w})}{\partial w_i} = m \sum_{\mathbf{d}_i} \hat{P}(\mathbf{d}_i) \phi_i(\mathbf{d}_i) - m \frac{\partial \log Z(w)}{\partial w_i}$$

$$\begin{aligned} \frac{\partial \log Z(w)}{\partial w_i} &= \frac{\partial}{\partial w_i} \left[\sum_x e^{\sum_{j=1}^k w_j \phi_j(\mathbf{d}_j)} \right] \\ &= \frac{1}{Z} \sum_x \sum_{d_i} \phi_i(\mathbf{d}_i) e^{\sum_{j=1}^k w_j \phi_j(\mathbf{d}_j)} \\ &= \sum_{\mathbf{d}_i} \phi_i(\mathbf{d}_i) P(\mathbf{d}_i | w) = \hat{E}_{f,w}[\phi_i] \end{aligned}$$

10-708 - ©Carlos Guestrin 2006-2008

12

Learning log-linear models with gradient ascent

- Gradient: $\frac{\partial \ell(\mathcal{D} : \mathbf{w})}{\partial w_i} = m \sum_{\mathbf{d}_i} \hat{P}(\mathbf{d}_i) \phi_i(\mathbf{d}_i) - m \sum_{\mathbf{d}_i} P(\mathbf{d}_i | \mathbf{w}) \phi_i(\mathbf{d}_i)$
- Requires one inference computation per fixture
- Theorem: \mathbf{w} is maximum likelihood solution iff $E[\phi_i] = E[\hat{\phi}_i]$ ← MN learning often called moment matching
- Usually, must regularize
- $\ell(\mathcal{D} : \mathbf{w}) \sim \lambda \mathbf{w}^2 = \ell^{\text{reg}}(\mathcal{D} : \mathbf{w})$ same as Gaussian prior e.g., $\phi_i = x$
- $\frac{\partial \ell(\mathcal{D} : \mathbf{w})}{\partial w_i} = \frac{\partial \ell(\mathcal{D} : \mathbf{w})}{\partial w_i} - 2\lambda w_i$!

10-708 - ©Carlos Guestrin 2006-2008

13

What you need to know about learning MN parameters?

- BN parameter learning easy
- MN parameter learning doesn't decompose!
- Learning requires inference!
- Apply gradient ascent or IPF iterations to obtain optimal parameters
 - Learning MNs if inference is intractable (read book)
 - two approaches : 1. use approx. inf. to compute derivative & hope for the best
 - 2. approximate log-likelihood (e.g., pseudo likelihood) and hope for the best

10-708 - ©Carlos Guestrin 2006-2008

14

Readings:
K&F: 19.3.2

Conditional Random Fields

Graphical Models – 10708

Carlos Guestrin

Carnegie Mellon University

November 17th, 2008

10-708 – ©Carlos Guestrin 2006-2008

15

Generative v. Discriminative classifiers – A review

- **Want to Learn:** $h: \mathbf{X} \mapsto \mathbf{Y}$
 - \mathbf{X} – features
 - \mathbf{Y} – target classes
- **Bayes optimal classifier** – $P(\mathbf{Y}|\mathbf{X})$
- **Generative classifier**, e.g., Naïve Bayes:
 - Assume some **functional form** for $P(\mathbf{X}|\mathbf{Y})$, $P(\mathbf{Y})$
 - Estimate parameters of $P(\mathbf{X}|\mathbf{Y})$, $P(\mathbf{Y})$ directly from training data
 - Use Bayes rule to calculate $P(\mathbf{Y}|\mathbf{X} = \mathbf{x})$
 - This is a '**generative**' model
 - Indirect computation of $P(\mathbf{Y}|\mathbf{X})$ through Bayes rule
 - But, can generate a sample of the data, $P(\mathbf{X}) = \sum_y P(y) P(\mathbf{X}|y)$
- **Discriminative classifiers**, e.g., Logistic Regression:
 - Assume some **functional form** for $P(\mathbf{Y}|\mathbf{X})$
 - Estimate parameters of $P(\mathbf{Y}|\mathbf{X})$ directly from training data
 - This is the '**discriminative**' model
 - Directly learn $P(\mathbf{Y}|\mathbf{X})$
 - But cannot obtain a sample of the data, because $P(\mathbf{X})$ is not available

10-708 – ©Carlos Guestrin 2006-2008

16

Log-linear CRFs

generalization of logistic regression

(most common representation)

- **Graph H :** only over hidden vars Y_1, \dots, Y_n
 - No assumptions about dependency on observed vars X
 - You must always observe all of X $P(Y|X)$
- **Feature** is some function $\phi[\mathbf{D}]$ for some subset of variables \mathbf{D}
 - e.g., indicator function, $D_i \subseteq \{Y_1, \dots, Y_n\}$ also depend on X
- **Log-linear model** over a CRF H :
 - a set of features $\phi_1[\mathbf{D}_1], \dots, \phi_k[\mathbf{D}_k]$
 - each \mathbf{D}_i is a subset of a clique in H
 - two ϕ 's can be over the same variables
 - a set of weights w_1, \dots, w_k
 - usually learned from data
 - $P(Y_1, \dots, Y_n | x) = \frac{1}{Z(x)} \exp \left[\sum_{i=1}^k w_i \phi_i(\mathbf{D}_i, x) \right]$

10-708 - ©Carlos Guestrin 2006-2008

17

Learning params for log linear CRFs – Gradient Ascent

$$P(Y_1, \dots, Y_n | x) = \frac{1}{Z(x)} \exp \left[\sum_{i=1}^k w_i \phi_i(\mathbf{D}_i, x) \right]$$

- Log-likelihood of data:

$$\begin{aligned} \ell(D; w) &= \log \prod_{j=1}^m P(y_j^{(j)}, \dots, y_n^{(j)} | x^{(j)}) \\ &= \sum_{j=1}^m \log \frac{1}{Z(x^{(j)})} \exp \left[\sum_{i=1}^k w_i \phi_i(d_i^{(j)}, x^{(j)}) \right] \end{aligned}$$

- Compute derivative & optimize
 - usually with conjugate gradient ascent

10-708 - ©Carlos Guestrin 2006-2008

18

Learning log-linear CRFs with gradient ascent

- Gradient: $\frac{\partial \ell(\mathcal{D} : \mathbf{w})}{\partial w_i} = \sum_{j=1}^m \left[\phi_i(\mathbf{d}_i^{(j)}, \mathbf{x}^{(j)}) - \sum_{\mathbf{d}_i} P(\mathbf{d}_i | \mathbf{x}^{(j)}, \mathbf{w}) \phi_i(\mathbf{d}_i, \mathbf{x}^{(j)}) \right]$
- $= \sum_{j=1}^m \phi_i(\mathbf{d}_i^{(j)}, \mathbf{x}^{(j)}) - \sum_{j=1}^m \sum_{\mathbf{d}_i} P(\mathbf{d}_i | \mathbf{x}^{(j)}, \mathbf{w}) \phi_i(\mathbf{d}_i, \mathbf{x}^{(j)})$
- "from data part" "from model part"
- "expected value of feature"
- "probability depends on $x^{(j)}$ "
- Requires one inference computation per data point j , per feature i
- Usually, must regularize
 - E.g., L₂ regularization on parameters
especially important

10-708 - ©Carlos Guestrin 2006-2008

19

What you need to know about CRFs

- Discriminative learning of graphical models
 - Fewer assumptions about distribution → often performs better than "similar" MN
 - Gradient computation requires inference over x per datapoint
→ Can be really slow!!

10-708 - ©Carlos Guestrin 2006-2008

20