

# Probabilistic Graphical Models (Cmput 651): Learning With Partial Data

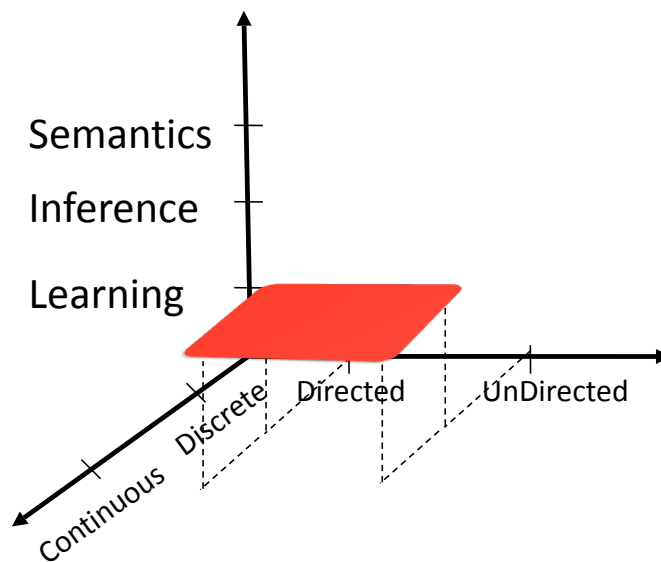
Matthew Brown

17/11/2008

Reading: Koller-Friedman Ch. 18, 19.3.3

1

## Space of topics



2

## Learning Markov nets

	complete data	partial data
known structure	easy (sort of)	hard
unknown structure	hard	very hard

↑  
This lecture

3

## Outline

### **Short review**

Markov net parameter learning

Structure learning

4

## Observation models (KF 18.1.1)

$$\mathbf{X} = \{X_1, \dots, X_n\}$$

$$O_{\mathbf{X}} = \{O_{X_1}, \dots, O_{X_n}\} \quad \text{observability variables}$$

$$P_{\text{missing}}(\mathbf{X}, O_{\mathbf{X}}) = P(\mathbf{X}) \cdot P_{\text{missing}}(O_{\mathbf{X}} | \mathbf{X})$$

### Can complicate likelihood function

no closed form for Bayes nets

(whereas exists closed form for complete data)

5

## Observation models (KF 18.1.1)

$$\mathbf{X} = \{X_1, \dots, X_n\}$$

$$O_{\mathbf{X}} = \{O_{X_1}, \dots, O_{X_n}\} \quad \text{observability variables}$$

$$P_{\text{missing}}(\mathbf{X}, O_{\mathbf{X}}) = P(\mathbf{X}) \cdot P_{\text{missing}}(O_{\mathbf{X}} | \mathbf{X})$$

Decoupling between  $\mathbf{X}$  and  $O_{\mathbf{X}}$  makes likelihood decomposable:

Missing completely at random (MCAR)  $P_{\text{missing}} \models (\mathbf{X} \perp O_{\mathbf{X}})$

-> allowed to ignore missing data

Missing at random (MAR)  $P_{\text{missing}} \models (O_{\mathbf{X}} \perp \mathbf{x}_{\text{hidden}}^{\mathbf{y}} | \mathbf{x}_{\text{obs}}^{\mathbf{y}})$

conditional decoupling (see next page)

eg: medical tests setting

6

## Observation models (KF 18.1.1)

### Theorem:

If  $P_{\text{missing}}$  satisfies MAR, the likelihood  $L(\theta, \psi : \mathcal{D})$  can be decomposed into (written as product of)

$$L(\theta : \mathcal{D}) \text{ and } L(\psi : \mathcal{D})$$

X parameters       $O_x$  parameters

7

## Identifiability (KF 18.1.4)

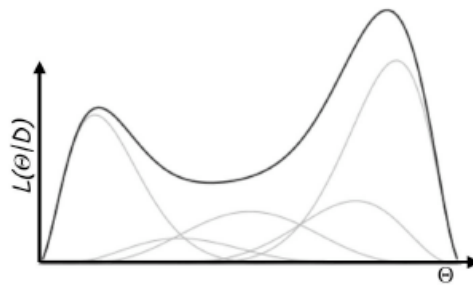
Can we uniquely “pin down” the model?

$$\text{Identifiability: } P(X|\theta) = P(X|\theta') \iff \theta = \theta'$$

$$\text{Local identifiability: } P(X|\theta) = P(X|\theta') \iff \theta = \theta' \\ \text{as long as } \|\theta - \theta'\| < \epsilon$$

8

## Likelihood function with partial data (KF 18.1.3)



partial data -> likelihood becomes marginal probability over unobserved variables

$$P(x_{obs}|\theta) = \int_{X_{missing}} P(x_{obs}, x_{missing}|\theta) dx_{missing}$$

i.e. sum of complete data likelihood functions

-> multimodal (not concave)

9

## Learning with partial data in Bayes nets

### Likelihood

no closed form

non-concave

### Iterative approaches:

Gradient methods

EM

alternately compute expected values of unobserved variables and then maximize parameters in "complete" data fashion

Gibbs sampling

10

## Outline

Short review

**Markov net parameter learning**

Structure learning

11

## Learning Markov nets w/ partial data (KF 19.3.3)

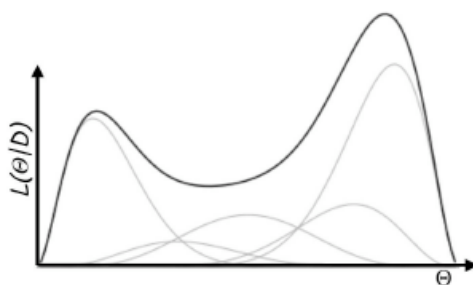
MAR, MCAR, ...

Issues with identifiability

Likelihood

no closed form to begin with (even w/ complete data)

no longer concave w/ partial data



12

### MN gradient ascent w/ partial data (KF 19.3.3.1)

#### Gradient ascent:

Assume MAR

$o[m]$  = observed entries in  $m^{\text{th}}$  data point

$\mathcal{H}[m]$  = missing entries in  $m^{\text{th}}$  data point

$(o[m], h[m])$  = complete assignment to  $\mathcal{X}$

#### Average log-likelihood

$$\begin{aligned} \frac{1}{M} \ln P(\mathcal{D} | \theta) &= \frac{1}{M} \ln \left( \sum_{m=1}^M \sum_{\mathbf{h}[m]} P(o[m], \mathbf{h}[m] | \theta) \right) \\ (\sim \text{means unnormalized}) &= \frac{1}{M} \ln \left( \sum_{m=1}^M \sum_{\mathbf{h}[m]} \tilde{P}(o[m], \mathbf{h}[m] | \theta) \right) - \ln Z \end{aligned}$$

### MN gradient ascent w/ partial data (KF 19.3.3.1)

#### Average log-likelihood

$$\begin{aligned} \frac{1}{M} \ln P(\mathcal{D} | \theta) &= \frac{1}{M} \ln \left( \sum_{m=1}^M \sum_{\mathbf{h}[m]} P(o[m], \mathbf{h}[m] | \theta) \right) \\ (\sim \text{means unnormalized}) &= \frac{1}{M} \ln \left( \sum_{m=1}^M \sum_{\mathbf{h}[m]} \tilde{P}(o[m], \mathbf{h}[m] | \theta) \right) - \ln Z \end{aligned}$$

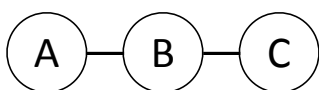
$$\frac{\partial}{\partial \theta_i} \frac{1}{M} \ell(\theta : \mathcal{D}) = \frac{1}{M} \left[ \sum_{m=1}^M \mathbf{E}_{\mathbf{h}[m] \sim P(\mathcal{H}[m] | o[m], \theta)} [\phi_i] \right] - \mathbf{E}_{\theta} [\phi_i]$$

Expectation conditional on observation  $m$       M inferences      1 inference

### MN gradient ascent w/ partial data (KF 19.3.3.1)

$$\frac{\partial}{\partial \theta_i} \frac{1}{M} \ell(\theta : \mathcal{D}) = \frac{1}{M} \left[ \sum_{m=1}^M \mathbf{E}_{h[m] \sim P(\mathcal{H}[m] | \mathbf{o}[m], \theta)}[\phi_i] \right] - \mathbf{E}_{\theta}[\phi_i]$$

For each data point m,  
estimate  $P(H[m] | \mathbf{o}[m], \theta)$   
requires inference



Binary variables A, B, C  
data point  $(a^0, b^?, c^1)$   
-> need  $P(B | a^0, c^1, \theta)$

### MN gradient: Complete vs. partial data (KF 19.3.3.1)

Complete data  $\frac{\partial}{\partial \theta_i} \frac{1}{M} \ell(\theta : \mathcal{D}) = \mathbf{E}_{\mathcal{D}}[\phi_i[\mathcal{X}]] - \mathbf{E}_{\theta}[\phi_i]$

1 inference / gradient step

Partial data  $\frac{\partial}{\partial \theta_i} \frac{1}{M} \ell(\theta : \mathcal{D}) = \frac{1}{M} \left[ \sum_{m=1}^M \mathbf{E}_{h[m] \sim P(\mathcal{H}[m] | \mathbf{o}[m], \theta)}[\phi_i] \right] - \mathbf{E}_{\theta}[\phi_i]$

Expectation conditional on observation m

M+1 inferences / gradient step

(M = # data points)

-> much more expensive



## Expectation Maximization (KF 19.3.3.2)

$$\text{E-step: } \bar{M}_{\theta^{(t)}}[\phi_i] = \frac{1}{M} \left[ \sum_{m=1}^M \mathbf{E}_{h^{[m]} \sim P(\mathcal{H}^{[m]} | \mathbf{o}^{[m]}, \theta)}[\phi_i] \right]$$

↑  
 expected complete assignment      expectation over missed variables in data point m -> requires inference for each m

E-step requires M (no. data points) inferences

M-step: treat expected variables from E-step as complete data, then max. likelihood estimation

- for Bayes nets, closed form
  - for Markov nets, need inference inside gradient ascent
- gradient form  $\bar{M}_{\theta^{(t)}}[\phi_i] - \mathbf{E}_{\theta^{(t,k)}}[\phi_i]$

17

## Comparison (KF 19.3.3.2)

Gradient ascent on log-likelihood

M+1 inferences / gradient step

$$\frac{\partial}{\partial \theta_i} \frac{1}{M} \ell(\theta : \mathcal{D}) = \frac{1}{M} \left[ \sum_{m=1}^M \mathbf{E}_{h^{[m]} \sim P(\mathcal{H}^{[m]} | \mathbf{o}^{[m]}, \theta)}[\phi_i] \right] - \mathbf{E}_{\theta}[\phi_i]$$

Expectation Maximization

E step: computes 1st term (observation counts) and “caches” them

M step: gradient ascent in “complete” setting

only 1 inference / gradient step

“cached” counts become less relevant as ascent proceeds

often better not to run ascent to convergence

18

## Outline

Short review

Markov net parameter learning

**Structure learning**

19

## Structure learning with partial data

### Case 1:

- know all the nodes
- want to find the edges
- partially-observed data

### Case 2:

- do not know all the nodes
  - i.e. hidden variables
- want to find edges and hidden nodes
- partially-observed data

20

## Greedy MN structure learning (also see KF Fig. 19.3)

Total feature set  $\Omega$

Initial feature set  $\Phi_0$

at all times:  $\theta_i = 0, \forall \phi_i \notin \Phi$

Iterate {

Optimize  $\theta_\Phi$  (parameter optimization)

Iterate over modification operators  $\mathcal{O}$  to structure {

$\mathcal{O}$  creates  $\Phi_{mod}$

$\hat{\Delta}_{\mathcal{O}}$  = improvement in score

}

choose set of modifications  $\mathcal{O}$  based on  $\hat{\Delta}_{\mathcal{O}}$

-> new structure  $\Phi$

}

21

## Case 1 - no hidden variables

Somewhat like structure learning with complete data  
(see previous slide)

But scores much more expensive with partial data  
marginal log likelihood expensive  
inside structure modification loop

Various methods to approximate / speed things up

Log-likelihood & overfitting

Parameters & structure

Regularization methods

22

## Structure learning with hidden variables

Hidden variables are MCAR (trivially)

Not identifiable

Local maxima in likelihood

V. hard, generally

If wrong choice of hidden variables

-> underconstrained ( $\# \text{ parameters} > \# \text{ data}$ )

Best when model overconstrained to start

hidden variables then “useful” in modeling dependencies  
among observed variables

Proper initialization crucial