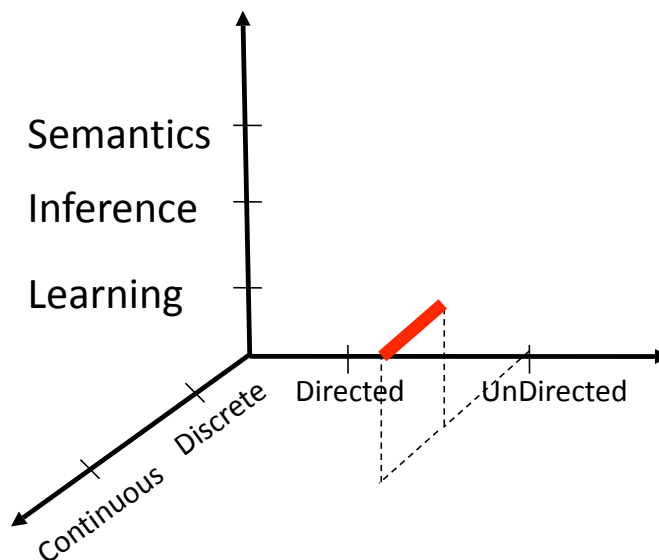


# Probabilistic Graphical Models (Cmput 651): Learning Undirected Models

Matthew Brown  
{14,17}/11/2008

Reading: Koller-Friedman Ch. 19

## Space of topics



## Learning Markov nets

	complete data	partial data
known structure	not as easy as for Bayes nets	hard
unknown structure	hard	very hard

↑  
This lecture

3

## Learning Markov nets is expensive

Markov factorization: 
$$P_{\mathcal{F}}(X) = \frac{1}{Z} \prod_{\phi_i \in \mathcal{F}} \phi_i(\text{scope}(\phi_i))$$

- partition function couples all the factors
  - cannot separate parameter estimation into local groups
- no closed form solutions
  - even for max. likelihood w/ complete data (recall this is easy for Bayes nets)
- learning based on iteration
  - inference required in each step -> expensive
  - but convex (for complete data)
- structure learning also expensive

4

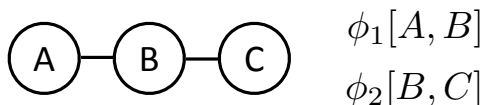
## Outline

### Parameter learning

### Structure learning

5

## Likelihood and the partition function (also see KF 19.2.1)



Log likelihood:  $\ln P(a, b, c) = \ln \left( \frac{1}{Z} \phi_1[a, b] \phi_2[b, c] \right)$   
 $= \ln \phi_1[a, b] + \ln \phi_2[b, c] - \ln Z$

Log likelihood for data  $D$  with  $M$  instances:

$$\begin{aligned} \ell(\theta : D) &= \sum_m (\ln \phi_1[a[m], b[m]] + \ln \phi_2[b[m], c[m]] - \ln Z) \\ &= \sum_{a,b} M[a, b] \ln \phi_1[a, b] + \sum_{b,c} M[b, c] \ln \phi_2[b, c] - M \ln Z(\theta) \end{aligned}$$

Counts of various assignments in  $D$

6

## Likelihood and the partition function (also see KF 19.2.1)

(Continued from last slide)

Log likelihood for data  $D$  with  $M$  instances:

$$\begin{aligned} \ell(\theta : \mathcal{D}) &= \sum_m (\ln \phi_1[a[m], b[m]] + \ln \phi_2[b[m], c[m]] - \ln Z) \\ &= \sum_{a,b} M[a, b] \ln \phi_1[a, b] + \sum_{b,c} M[b, c] \ln \phi_2[b, c] - M \ln Z(\theta) \end{aligned}$$

Counting terms involve only a single factor

Partition function term involves ALL factors:

$$Z(\theta) = \sum_{a,b,c} \phi_1[a, b] \phi_2[b, c]$$

7

## Maximum likelihood estimation (also see KF 19.2.1)

Bayes nets: estimate conditional distributions separately for each node

-> likelihood decomposable

vs.

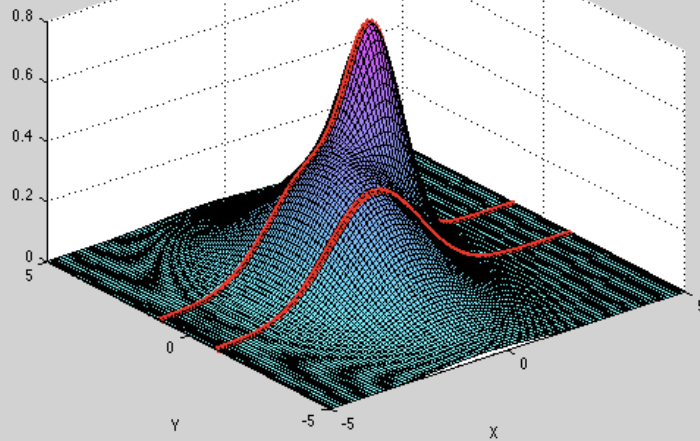
Markov nets: partition function involves all factors

changing  $\phi_1$  could change optimal value for  $\phi_2$

8

## Maximum likelihood estimation (also see KF 19.2.1)

### Maximum likelihood surface



Max w.r.t.  $X$   
depends on  
value of  $Y$

9

## Log-likelihood (also see KF 19.2.2)

Log linear model: 
$$P(X_1, \dots, X_n : \theta) = \frac{1}{Z} \exp \sum_{i=1}^k \theta_i \phi_i[D_i]$$

$\theta_i$  → weight  
 $\phi_i[D_i]$  → indicator function

Log-likelihood: 
$$\ell(\theta : \mathcal{D}) = \sum_i \theta_i \left( \sum_m \phi_i[\xi[m]] \right) - M \ln Z(\theta)$$

$\xi[m]$  →  $m^{\text{th}}$  data point

10

Log-likelihood (also see KF 19.2.2)

Log-likelihood:  $\ell(\theta : \mathcal{D}) = \sum_i \theta_i \left( \sum_m \phi_i[\xi[m]] \right) - M \ln Z(\theta)$

Divide by M (no. data points)

$$\frac{1}{M} \ell(\theta : \mathcal{D}) = \sum_i \theta_i E_{\mathcal{D}}[\phi_i[d_i]] - \ln Z(\theta)$$

Empirical expectation of  $\phi_i$

Partition function:  $\ln Z(\theta) = \ln \sum_{\xi} \exp \left\{ \sum_i \theta_i \phi_i[\xi] \right\}$

Sum over data points

Log-likelihood (also see KF 19.2.3)

$$\frac{1}{M} \ell(\theta : \mathcal{D}) = \sum_i \theta_i E_{\mathcal{D}}[\phi_i[d_i]] - \ln Z(\theta)$$

Linear in  $E_{\mathcal{D}}[\phi_i[d_i]]$

Partition function:  $\ln Z(\theta) = \ln \sum_{\xi} \exp \left\{ \sum_i \theta_i \phi_i[\xi] \right\}$

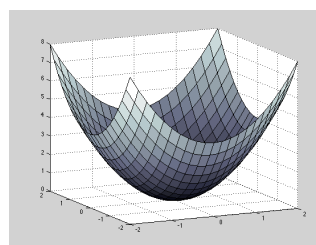
- is convex

⇒ log-likelihood is concave

Convex function:

$$f(\alpha \vec{x} + (1 - \alpha) \vec{y}) \leq \alpha f(\vec{x}) + (1 - \alpha) f(\vec{y})$$

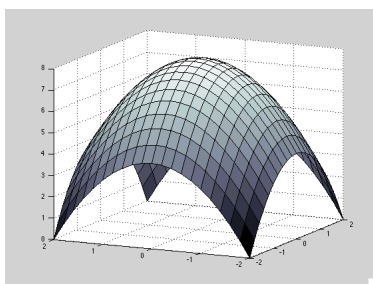
$$0 \leq \alpha \leq 1$$



## Log-likelihood (also see KF 19.2.3)

Log-likelihood is concave

- $\exists$  global maximum
- $\nexists$  local maxima
- global maximum may not be unique
  - Markov net parameterization may be redundant
  - i.e. multiple representations of same distribution



include simple examples based on diamond net if there's time

13

## Maximum likelihood parameter estimation (KF 19.3.1)

Want to find parameters that maximize log-likelihood

Gradient = 0 at maximum

$$\frac{\partial}{\partial \theta_i} \frac{1}{M} \ell(\theta : \mathcal{D}) = \sum_i \theta_i \mathbf{E}_{\mathcal{D}}[\phi_i[\mathbf{d}_i]] - \ln Z(\theta)$$

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

-> @ maximum:  $\mathbf{E}_{\mathcal{D}}[\phi_i[\mathcal{X}]] = \mathbf{E}_{\theta}[\phi_i]$

BUT  $\nexists$  closed form solution!

- use gradient ascent instead

KF Proposition 19.2.3

$$\frac{\partial}{\partial \theta_i} \ln Z(\theta) = \mathbf{E}_{\theta}[\phi_i]$$

14

## Maximum likelihood parameter estimation

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

Show numerical example using A-B-C network

15

## Maximum likelihood parameter estimation (KF 19.3.1)

Want to find parameters that maximize log-likelihood

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

Easy to compute

Expected counts over parameter space

- requires inference
- once at each step of gradient ascent
- typically expensive

16



## Conditionally-trained models (also see KF 19.3.2)

Discriminative (instead of generative)

Conditional random field (CRF)

encodes  $P(Y | X)$

Maximize log conditional likelihood

$$\ell_{Y|X}(\theta : \mathcal{D}) = \ln P(\mathbf{y}[1, \dots, M] | \mathbf{x}[1, \dots, M], \theta) = \sum_{m=1}^M \ln P(\mathbf{y}[m] | \mathbf{x}[m], \theta)$$

- is concave

- use gradient ascent

17

## Conditionally-trained models (also see KF 19.3.2)

Gradient ascent on log conditional likelihood

$$\ell_{Y|X}(\theta : \mathcal{D}) = \ln P(\mathbf{y}[1, \dots, M] | \mathbf{x}[1, \dots, M], \theta) = \sum_{m=1}^M \ln P(\mathbf{y}[m] | \mathbf{x}[m], \theta)$$

Gradient:

$$\frac{\partial}{\partial \theta_i} \ell_{Y|X}(\theta : \mathcal{D}) = \sum_{m=1}^M \left[ \phi_i[\mathbf{y}[m], \mathbf{x}[m]] - \mathbf{E}_{\theta}[\phi_i | \mathbf{x}[m]] \right]$$

Counts on dataset

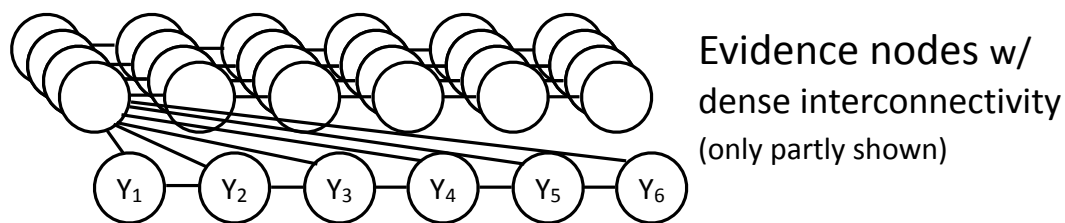
Expectation w.r.t. model conditioned on  $m^{\text{th}}$  data point

-> must run inference M times for each gradient step!

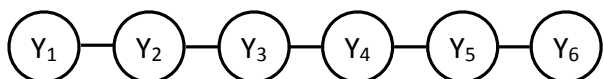
V. expensive

But, can be offset by simpler model (comp. to generative) <sup>18</sup>

### Conditionally-trained models (also see KF Example 19.3.3)



### Discriminative model conditioned on evidence nodes



Inference on chain is much easier

### Maximum entropy & maximum likelihood (KF 19.3.4)

Given some data  $D$ , find distribution  $Q$  that matches it without a lot of extra structure / assumptions

Maximum-Entropy

Find  
that maximize  
subject to

$$Q(\mathcal{X})$$

$$H_Q(\mathcal{X})$$

Entropy  
high value  $\rightarrow$  little structure

$$E_Q[\phi_i] = E_D[\phi_i] \quad i = 1, \dots, k$$

expectation constraints

## Maximum entropy & maximum likelihood (KF 19.3.4)

**Theorem:** Max entropy solution  $Q^*$  (to problem from previous slide) satisfies

$$Q^* = P_{\hat{\theta}}(\mathcal{X}) = \frac{1}{Z(\hat{\theta})} \exp \left\{ \sum_i \hat{\theta}_i \phi_i[\mathcal{X}] \right\}$$

$\hat{\theta}$  = maximum likelihood solution relative to data set D

i.e. max likelihood and max entropy are related

21

## Priors and regularization (also see KF 19.4)

Max. likelihood estimation prone to overfitting

Use priors to constrain  $\theta$  parameters

from log linear model  $P(X_1, \dots, X_n : \theta) = \frac{1}{Z} \exp \sum_{i=1}^k \theta_i \phi_i[D_i]$

closed form like Bayes formula not available

use max. a posteriori (MAP) instead - maximize  $P(\theta)P(\mathcal{D} | \theta)$

22

## Gaussian prior (also see KF 19.4.1)

$$P(\theta | \sigma^2) = \prod_{i=1}^k \frac{1}{\sqrt{2\pi}} \exp \left\{ -\frac{\theta_i^2}{2\sigma^2} \right\}$$

$\sigma^2$  hyperparameter for variance

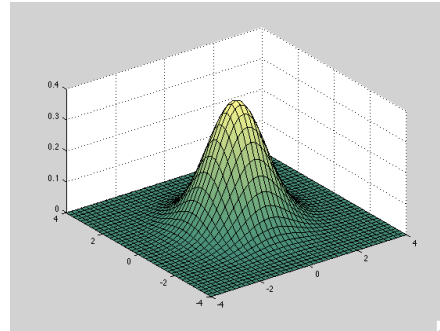
Amount of regularization

Can use different  $\sigma_i$

Typically use  $\text{Cov}(\theta_i, \theta_j) = 0$  for  $i \neq j$

i.e. assume  $\theta_i$  independent

Predisposes  $\theta$ 's to be small



23

## Laplacian prior (also see KF 19.4.1)

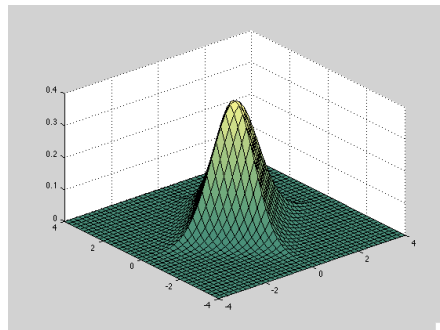
$$P_{Laplacian}(\theta | \beta) = \frac{1}{2\beta} \exp \left\{ -\frac{|\theta|}{\beta} \right\}$$

$\beta$  hyperparameter

Can use different  $\beta_i$

Assume  $\theta_i$  independent

Also predisposes  $\theta$ 's to be small



24

## Priors and regularization (also see KF 19.4.1)

$$P(\theta, \mathcal{D}) = P(\theta)P(\mathcal{D} | \theta)$$

↑ Prior    ↑ Likelihood

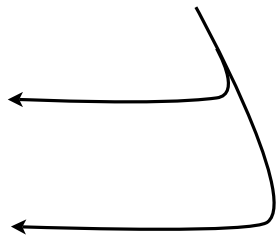
Penalty term regularizes MLE

Taking logarithms:

$$\ell(\theta : \mathcal{D}) = \sum_i \theta_i \left( \sum_m \phi_i[\xi[m]] \right) - M \ln Z(\theta) + \text{penalty term}$$

$$-\frac{1}{2\sigma^2} \sum_{i=1}^k \theta_i^2 \quad \begin{array}{l} \text{Gaussian penalty} \\ \text{L}_2 \text{ regularization} \end{array}$$

$$-\frac{1}{\beta} \sum_{i=1}^k |\theta_i| \quad \begin{array}{l} \text{Laplacian penalty} \\ \text{L}_1 \text{ regularization} \end{array}$$



## L1 vs L2 regularization (also see KF 19.4.1)

$$-\frac{1}{2\sigma^2} \sum_{i=1}^k \theta_i^2 \quad \begin{array}{l} \text{L}_2 \text{ regularization} \\ \text{larger } \theta \text{ values penalized more heavily} \end{array}$$

$$-\frac{1}{\beta} \sum_{i=1}^k |\theta_i| \quad \begin{array}{l} \text{L}_1 \text{ regularization} \\ \text{uniform penalty} \\ \text{-> better at driving } \theta \text{ all the way to zero} \\ \text{-> sparser models learned (more } \theta_i=0 \end{array}$$

## Learning with approximation (also see KF 19.5)

How to learn when inference is hard?

eg: grid networks

Approach 1: approximate inference inside learning loop

generalized belief propagation

particle-based methods

Approach 2: approximate cost function

inference easier

In many cases, approaches 1 and 2 are formally equivalent.

27

## Learning with belief propagation (also see KF 19.5.1)

Theorem: When using an approximate inference algorithm on the trained model, it is best to do the training with the same inference algorithm.

**BUT**: generalized belief propagation inside a gradient-based learning loop can cause problems:

- marginals are only approximate -> noise in gradient
- possible non-convergence -> unstable gradient
- can use heuristics & manual testing to address these
- one solution on next slide

28

## Learning with belief propagation (also see KF 19.5.1.2)

Recall: Generalized belief propagation equivalent to optimization on approximate factored energy functional:

$$\tilde{F}[P_{\mathcal{F}}, Q] = \sum_i E_{C_i \sim \beta_i} [\ln \psi_i] + \sum_{C_i \in \mathcal{T}} H_{\beta_i}(C_i) - \sum_{(C_i - C_j) \in \mathcal{T}} H_{\mu_{i,j}}(S_{i,j})$$

which comes from KL-Divergence.

Can use similar approach for learning  
(go to next slide)

29

## Learning with belief propagation (also see KF 19.5.1.2)

Want tractable approximation to:

Maximum-Entropy

**Find**  $Q(\mathcal{X})$   
**that maximize**  $H_Q(\mathcal{X})$

**subject to**

$$E_Q[\phi_i] = E_{\mathcal{D}}[\phi_i] \quad i = 1, \dots, k$$

(equivalent to max. likelihood)

Solution = constrained optimization of factored form of entropy:

$$H_Q(\mathcal{X}) \approx \sum_{C_i \in \mathcal{K}} H_{\beta_i}(C_i) - \sum_{(C_i - C_j) \in \mathcal{K}} H_{\mu_{i,j}}(S_{i,j})$$

30

## Learning with belief propagation (also see KF 19.5.1.2)

Constrained optimization of factored form of entropy:

$$H_Q(\mathcal{X}) \approx \sum_{C_i \in \mathcal{K}} H_{\beta_i}(C_i) - \sum_{(C_i - C_j) \in \mathcal{K}} H_{\mu_{i,j}}(S_{i,j})$$

valid objective function

- avoid (non)convergence issues
- use whatever optimization method you want

reformulation exact when cluster graph = tree

- approximate otherwise
- e.g. generalized belief propagation

31

## Alternate objective functions (also see KF 19.6)

Goals:

easier objective than likelihood

**VALID** objective

no convergence issues

32



### Alternate objective functions (also see KF 19.6)

Likelihood of one data point:

$$\ell(\theta : \xi) = \ln \tilde{P}(\xi | \theta) - \ln Z(\theta) = \ln \tilde{P}(\xi | \theta) - \ln \left( \sum_{\xi'} \tilde{P}(\xi' | \theta) \right)$$

Want to make this large

Want to make this small

In 2<sup>nd</sup> term, summation over all assignments to Val(X) requires inference -> expensive

Approach:

in 2<sup>nd</sup> term, use more tractable set than all of Val(X)

### Pseudo-likelihood (also see KF 19.6.1)

Likelihood  $P(\xi) = \prod_{j=1}^n P(x_j | x_1, \dots, x_{j-1})$  (from chain rule)

$$P(\xi) \approx \prod_j P(x_j | x_1, \dots, x_{j-1}, x_{j+1}, \dots, x_n)$$



Pseudo-likelihood

$$\ell_{\text{pseudo}}(\theta : \mathcal{D}) = \frac{1}{M} \sum_m \sum_j \ln P(x_j[m] | x_{-j}[m], \theta)$$

$x_{-j}$  means  $x_1, \dots, x_{j-1}, x_{j+1}, \dots, x_n$

Note:  $P(X_j | X_{-j}) = P(X_j | \text{Neighbours}_{X_j})$

### Pseudo-likelihood (also see KF 19.6.1)

$$\ell_{\text{pseudo}}(\theta : \mathcal{D}) = \frac{1}{M} \sum_m \sum_j \ln P(x_j[m] | \mathbf{x}_{-j}[m], \theta)$$

$\mathbf{x}_{-j}$  means  $x_1, \dots, x_{j-1}, x_{j+1}, \dots, x_n$

Why do we care?

$$P(x_j | \mathbf{x}_{-j}) = \frac{P(x_j, \mathbf{x}_{-j})}{P(\mathbf{x}_{-j})} = \frac{\tilde{P}(x_j, \mathbf{x}_{-j})}{\tilde{P}(\mathbf{x}_{-j})} = \frac{\tilde{P}(x_j, \mathbf{x}_{-j})}{\sum_{x'_j} \tilde{P}(x'_j, \mathbf{x}_{-j})}$$

-> no global partition function!

only local partition function

much cheaper to evaluate

~ means unnormalized

### Gradient of pseudo-likelihood (KF Definition 19.6.1)

$$\frac{\partial}{\partial \theta_i} \ell_{\text{pseudo}}(\theta : \mathcal{D}) = \sum_{j: X_j \in \text{Scope}[\phi_i]} \left( \frac{1}{M} \sum_m \phi_i[\xi[m]] - \mathbf{E}_{x'_j \sim P_\theta(X_j | \mathbf{x}_{-j}[m])} [\phi_i[x'_j, \mathbf{x}_{-j}[m]]] \right)$$

Much cheaper than likelihood's gradient:

$$\mathbf{E}_{x'_j \sim P_\theta(X_j | \mathbf{x}_{-j}[m])} [\phi_i[x'_j, \mathbf{x}_{-j}[m]]]$$

-> summation only over  $X_j$  conditioned on its neighbours

(not over all of  $X$  as in partition function  $Z$ )

i.e. inference over only small part of graph ( $X_j$  and neighbours)

## Pseudo-likelihood (also see KF 19.6.1)

Pseudo-likelihood is concave  
-> unique, global maximum

37

## Likelihood & pseudo-likelihood (KF Theorem 19.6.2)

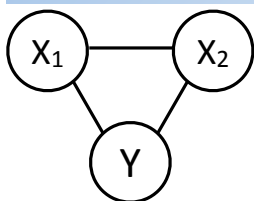
### **Theorem:**

Assuming data generated by log-linear model  $P_{\theta^*}$ ,  
as no. data points  $M \rightarrow \infty$ ,  $P(\theta_{pl} = \theta^*) \rightarrow 1$ ,  
where  $\theta_{pl}$  = global optimum of pseudo-likelihood  
objective

i.e. pseudo-likelihood converges to likelihood with large  $M$   
BUT, this assumes sufficiently expressive model and large  $M$   
- these assumptions typically do not hold

38

## Likelihood & pseudo-likelihood (KF Example 19.6.3)



$$X_1 \simeq X_2$$

$X_1, X_2$  somewhat correlated with  $Y$

Pseudo-likelihood will overestimate  $X_1$ - $X_2$  parameters and underestimate  $X_i$ - $Y$  parameters.  
Okay for  $P(X_2|X_1)$  but not  $P(X_2|Y)$

In general, pseudo-likelihood assumes  $X$ 's neighbourhood is observed  
"ignores" weaker or longer-ranger range dependencies

39

## Contrastive optimization (also see KF 19.6.2)

log-likelihood:

$$\ell(\theta; \xi) = \ln \tilde{P}(\xi | \theta) - \ln Z(\theta) = \ln \tilde{P}(\xi | \theta) - \ln \left( \sum_{\xi'} \tilde{P}(\xi' | \theta) \right)$$

Want to maximize contrast between these

Similarly motivated methods:

contrastive divergence

max. margin training

40

## Outline

Parameter learning

**Structure learning**

41

## Structure learning (also see KF 19.7)

constraint-based

constrain structure to reflect independencies in  $P(X)$

score-based

score each structure, optimize score

42

## Constraint-based structure learning (KF 19.7.1)

Similar to case for Bayes net.

Want  $H^*$  to factorize  $P^*$

assume:  $H^*$  perfect map for  $P^*$ , degree  $H^* \leq d^*$

Test for independencies:

Markov independence  $(X \perp \mathcal{X} - \{X\} - \mathcal{N}_{H^*}(X) \mid \mathcal{N}_{H^*}(X)) \quad \forall X$

Pairwise independence  $(X \perp Y \mid \mathcal{X} - \{X, Y\}) \quad \forall (X, Y) \notin \mathcal{H}$

BUT, testing Markov or pairwise involves all variables

-> exponential in num. nodes

see next slide

43

## Testing for independencies

Independence testing using pairwise independence:

$$(X \perp Y \mid \mathcal{X} - \{X, Y\}) \quad \forall (X, Y) \notin \mathcal{H}$$

Suppose  $N$  binary nodes:

$\exists \binom{N}{2}$  sets  $\{X, Y\}$  to test

for each of 4 assignments  $\{x, y\}$ , must check equality  
under  $2^{(N-2)}$  assignments to other nodes in  $\mathcal{X} - \{X, Y\}$

-> exponential!

(similar argument for Markov independence)

44

## Independence testing (KF 19.7.1)

(Assuming degree (max # edges / node)  $H^* \leq d^* \ll N$ )

Consider X,Y

no edge  $\rightarrow \mathcal{N}_{H^*}(X)$  and  $\mathcal{N}_{H^*}(Y)$  separate X and Y

i.e.  $\exists$  set Z with  $|Z| \leq \min(|\mathcal{N}_{H^*}(X)|, |\mathcal{N}_{H^*}(Y)|)$

such that  $sep_{H^*}(X; Y \mid Z)$

so:  $X-Y \notin H^*$  if and only if  $\exists Z, |z| \leq d^* \& P^* \models (X \perp Y \mid Z)$

For each pair X,Y test for edge using  $\sum_{k=0}^{d^*} \binom{n-2}{k}$  tests

Polynomial number of tests

Each test involves  $\leq d^* + 2$  variables

$\leq 2^{d^*+2}$  assignments to check for binary nodes

Tractable for small  $d^*$

45

## Constraint-based structure learning (KF 19.7.1)

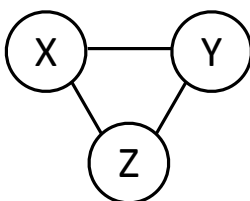
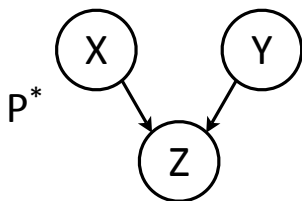
Limitations:

assume  $H^*$  perfect map for  $P^*$

assume bounded order for  $H^*$

assume enough data for reliable independence tests

Example



Correct  $H^*$  unreachable because  $X \perp Y \mid \{\}$

$\rightarrow$  discard X-Y edge  
i.e. no Markov net is perfect map for  $P^*$

46

## Constraint-based structure learning (KF 19.7.1)

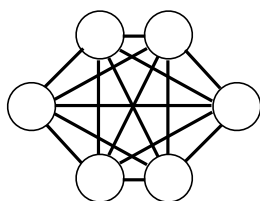
Limitations (cont'd):

For Markov nets, global independence structure not necessarily useful

eg: fully connected network with pairwise potentials

-> complex connectivity but compact factorization

constraint-based learning does not help find factorization



47

## Score-based learning (KF 19.7.2)

### Hypothesis space

- log-linear model

$$P(\mathcal{X} \mid \mathcal{M}, \theta) = \frac{1}{Z} \exp \left\{ \sum_{i \in \Phi[\mathcal{M}]} \theta_i \phi_i[\xi] \right\} = \frac{1}{Z} \exp \left\{ \phi^T \theta \right\}$$

Also:

- basic graph structure
- factor graph
- > different levels of model “granularity”

48

stru



## Score-based learning (KF 19.7.2)

### Hypothesis space

- log-linear model

$$P(\mathcal{X} \mid \mathcal{M}, \theta) = \frac{1}{Z} \exp \left\{ \sum_{i \in \Phi[\mathcal{M}]} \theta_i \phi_i[\xi] \right\} = \frac{1}{Z} \exp \left\{ \phi^T \theta \right\}$$

- Given set of features  $\Omega$ ,  
derive model  $\mathcal{M}$  from features  $\Phi[\mathcal{M}] \subseteq \Omega$   
by setting  $\theta_i = 0$  if  $\phi_i \notin \Phi[\mathcal{M}]$   
(Also, optimize other  $\theta_i$  )
- Structure implicit in  $\mathcal{M}$   
connect all  $X \in \text{scope}(\phi_i), \forall \phi_i \in \Phi[\mathcal{M}]$

49

stru  
c

## Score function (also see KF 19.7.3.1)

### log-likelihood

$$\text{score}_L(\mathcal{M} : \mathcal{D}) = \max_{\theta \in \Theta[\mathcal{M}]} \ln P(\mathcal{D} \mid \mathcal{M}, \theta) = \ell(\langle \mathcal{M}, \hat{\theta}_{\mathcal{M}} \rangle : \mathcal{D})$$

↑
↑

model            data

### Overfitting problem

$$\Phi[\mathcal{M}_1] \subset \Phi[\mathcal{M}_2] \rightarrow \text{score}_L(\mathcal{M}_1 : \mathcal{D}) \leq \text{score}_L(\mathcal{M}_2 : \mathcal{D})$$

more expressive model fits noise in  $\mathcal{D}$

must restrict factors' expressiveness or regularize

50

## 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  $\Phi_{mod}$  (see next slide)

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

}

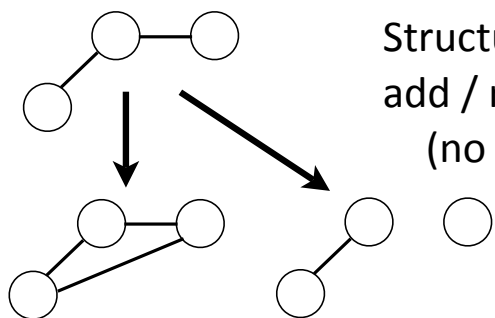
choose set of modifications  $\mathcal{O}$  based on  $\hat{\Delta}_{\mathcal{O}}$

-> new structure  $\Phi$

}

51

## Structure modification



Structure modification steps  $\mathcal{O}$  :  
add / remove single edge  
(no directionality!)

Then need score for each modification  $\hat{\Delta}_{\mathcal{O}}$   
(see below)

52

## Structure scoring (also see KF 19.7.4.2)

### Bayes nets:

structure score evaluation easy

closed form available

score decomposes based on structure

change to structure changes only one term in the score

changes to different parts of structure do not interact in the score

-> efficiency: dynamic programming, caching, etc.

53

## Structure learning performance (also see KF 19.7.4.2)

### Markov nets:

structure score evaluation **harder**

Score for each modification  $\hat{\Delta}_{\mathcal{O}}$

must optimize  $\theta_{\Phi_{mod}}$

requires inferences inside gradient ascent loop

can start from  $\theta_{\Phi_{current}}$  to speed things up

still expensive (& inside overall structure iteration loop)

Cheaper: rank order  $\hat{\Delta}_{\mathcal{O}}$  (instead of full evaluation )

54

## Structure learning performance (also see KF 19.7.4.2)

### Markov nets: structure score evaluation **harder**

partition function couples everything

computing likelihood score requires inference

structure score requires parameter estimation (no closed form)

structure score does not decompose

-> expensive!

### Good news:

structure learning is convex (with fully observed data)

structure learning more expensive for Markov vs. Bayes

55

## Regularization & structure learning (also see KF 19.7.3.2)

Bayesian score  $\text{score}_B(\mathcal{G} : \mathcal{D}) = \log P(\mathcal{D} | \mathcal{G}) + \log P(\mathcal{G})$

$$P(\mathcal{D} | \mathcal{G}) = \int P(\mathcal{D} | \mathcal{M}, \theta) P(\theta | \mathcal{M}) d\theta$$

likelihood      prior

marginal likelihood

average based on parameter prior

regularizes parameters

avoids overfitting

efficient for Bayes nets

but too hard to evaluate for Markov nets

(because of partition function)

use BIC score instead (next slide)

56

## Regularization & structure learning (also see KF 19.7.3.2)

$$\text{score}_{BIC}(\mathcal{M} : \mathcal{D}) = \ell(\langle \mathcal{M}, \hat{\theta}_{\mathcal{M}} \rangle : \mathcal{D}) - \frac{\dim(\mathcal{M})}{2} \ln M.$$

Asymptotic approximation to marginal likelihood

$\dim(\mathcal{M})$  = dimension of model

degrees of freedom

-> penalizes more complex models (i.e. more D.O.F.)

57

## Other regularizations for structure learning (KF 19.7.4)

MAP score

$L_1$  regularization

$$\text{score}_{L_1}(\theta : \mathcal{D}) == \ell(\langle \mathcal{M}, \theta \rangle : \mathcal{D}) - \|\theta\|_1$$

58