

Readings:

K&F: 18.1, 18.2, 18.3, 18.4

Dynamic Bayesian Networks

Beyond 10708

Graphical Models – 10708

Carlos Guestrin

Carnegie Mellon University

December 1st, 2006

1

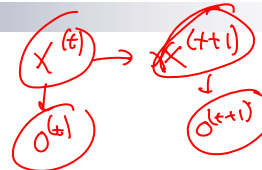
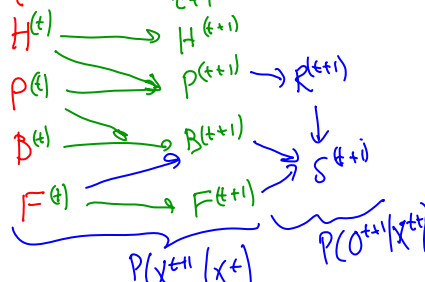
Dynamic Bayesian network (DBN)

- HMM defined by

- Transition model $P(X^{(t+1)}|X^{(t)})$
- Observation model $P(O^{(t)}|X^{(t)})$
- Starting state distribution $P(X^{(0)})$

- DBN – Use Bayes net to represent each of these compactly

- Starting state distribution $P(X^{(0)})$ is a BN
- (silly) e.g. performance in grad. school DBN
 - Vars: Happiness, Productivity, HiraBility, Fame
 - Observations: Paper, Schmooze



$P(X^{(t+1)} | X^{(t)})$
 how many params $(2^2-1)2^2$
 without DBN $2^8 - 2^2$
 with DBN
 $P(H^{t+1} | H^t)$ $(2-1) \cdot 2$
 $P(P^{t+1} | P^t, H^t)$ $(2-1) \cdot 2^2$
 $P(B^{t+1} | P^t, B^t, F^t)$ $(2-1) \cdot 2^3$
 $P(F^{t+1} | F^t)$ $(2-1) \cdot 2$

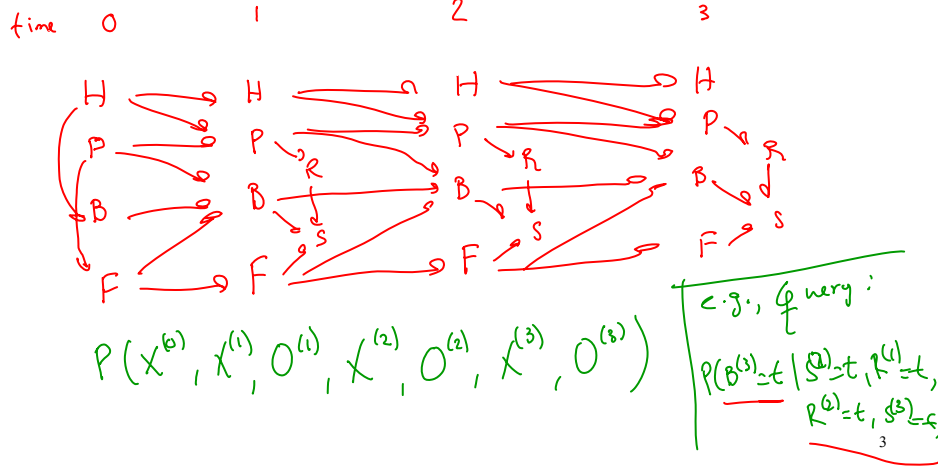
2

Unrolled DBN

X
H, P, B, F

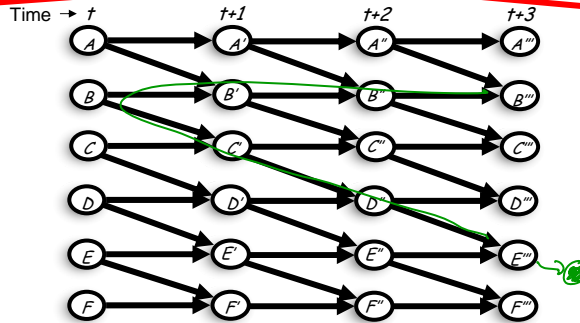
O
R, S

- Start with $P(X^{(0)})$
- For each time step, add vars as defined by 2-TBN



"Sparse" DBN and fast inference

~~"Sparse" DBN ☹️ Fast inference~~

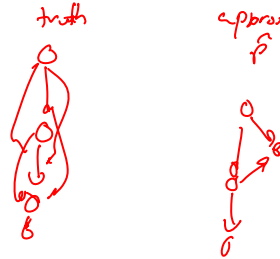
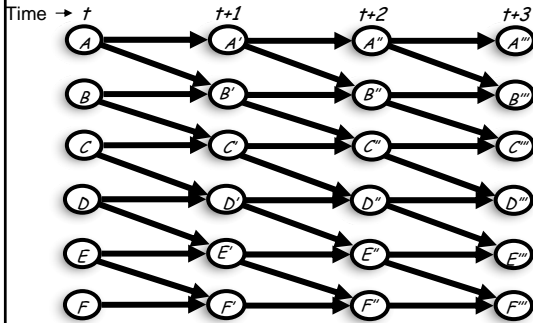


$A''' \perp E''$ true
 $B'' \perp E''$ no!
 $A''' \perp E''$ no!!
 \vdots

BK Algorithm for approximate DBN inference [Boyen, Koller '98]

- Assumed density filtering:

- Choose a factored representation \hat{P} for the belief state
- Every time step, belief not representable with \hat{P} , project into representation



7

A simple example of BK: Fully-Factorized Distribution

- Assumed density:

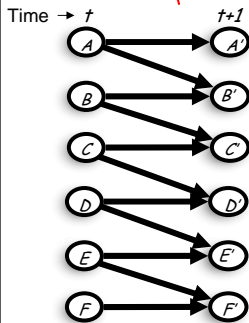
- Fully factorized

$$\hat{P}^{(t+1)}(X^{(t+1)}) = \prod_i \hat{P}_i(X_i^{(t+1)})$$

True $P(X^{(t+1)})$:

fully connected graph!!
i

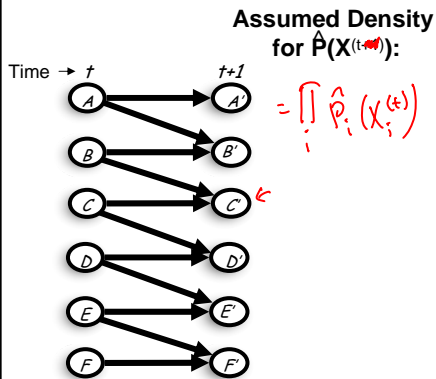
Assumed Density for $\hat{P}(X^{(t+1)})$:



8

Computing Fully-Factorized Distribution at time t+1

- Assumed density:
 - Fully factorized



Computing for $\hat{P}(X^{(t+1)})$:

$= \prod_i \hat{P}_i(X_i^{(t+1)})$

how would you compute $\hat{P}(C')$

$\hat{P}(B') \rightarrow P(C^{(t+1)} | B^{(t)}, C^{(t)})$

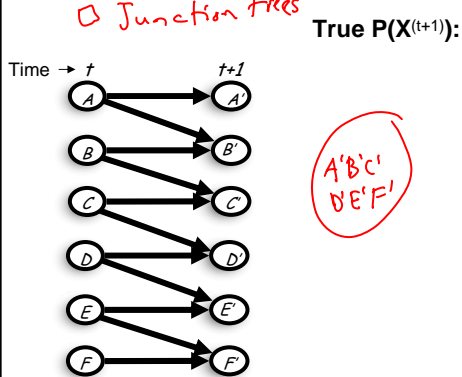
$\hat{P}(C) \rightarrow P(C) \xrightarrow{\text{transition model}} P(C^{(t+1)} | B^{(t)}, C^{(t)})$

$\hat{P}(C) = \sum_{bc} \hat{P}(b) \cdot \hat{P}(c) \cdot P(C^{(t+1)} | b, c)$

Variable elimination inference!

General case for BK: Junction Tree Represents Distribution

- Assumed density: *too simple*
 - Fully factorized
 - Junction trees



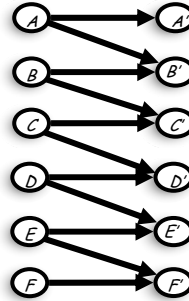
Assumed Density for $\hat{P}(X^{(t+1)})$:

choose JT for example

$A'B'$
|
 $B'C'$
|
 $C'D'$
|
 $D'E'$
|
 $E'F'$

Computing factored belief state in the next time step

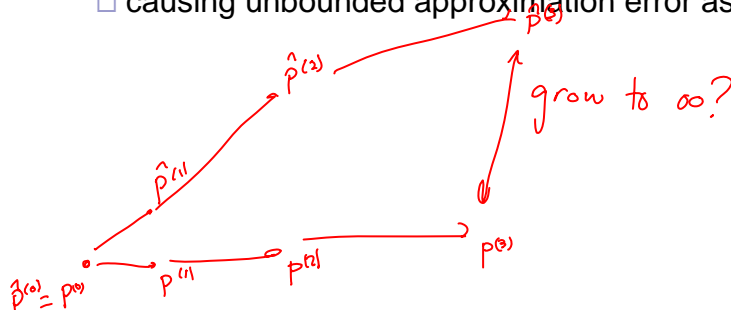
- Introduce observations in current time step
 - Use J-tree to calibrate time t beliefs
- Compute $t+1$ belief, project into approximate belief state
 - marginalize into desired factors
 - corresponds to KL projection
- Equivalent to computing marginals over factors directly
 - For each factor in $t+1$ step belief
 - Use variable elimination



KL-Projection
 start $\hat{p}^{(t)}$
 compute true $P^{(t+1)}$
 again $KL(p^{(t+1)}, \hat{z})$
 $\hat{z} \in \text{family of JT for } P^{(t+1)}$
 $= \hat{p}^{(t+1)}$
 if $\hat{p}^{(t+1)}$ is a JT &
 $\hat{p}_i(C_i^{(t+1)})$ computed exactly with V.E.

Error accumulation

- Each time step, projection introduces error
- Will error add up?
 - causing unbounded approximation error as $t \rightarrow \infty$



Contraction in Markov process

$d(t) > d(t+1)$

transition model:

- $s \xrightarrow{(1-\gamma)} P(s|s)$
- $s \xrightarrow{\gamma} \text{Uniform dist.}$

transition model:

- $p(t) \xrightarrow{(1-\gamma)} P_{\text{n.u.}}(t+1)$
- $p(t) \xrightarrow{\gamma} \text{uniform}$

transition model:

- $\hat{p}(t) \xrightarrow{(1-\gamma)} \hat{P}_{\text{n.u.}}(t+1)$
- $\hat{p}(t) \xrightarrow{\gamma} \text{Uniform}$

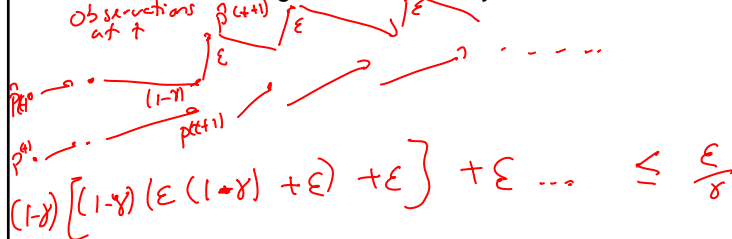
$\Rightarrow (1-\gamma) \text{KL}(p(t), \hat{p}(t)) \geq \text{KL}(p(t+1), \hat{p}(t+1))$

side note: if make obs. o $\text{KL}(p^{(t)}, \hat{p}^{(t)}) \geq E[\text{KL}(p^{(t+1)}, \hat{p}^{(t+1)})]$
 no expansion in expectation

BK Theorem

$$\sum_{i=0}^{\infty} (1-\gamma)^i = \frac{1}{\gamma}$$

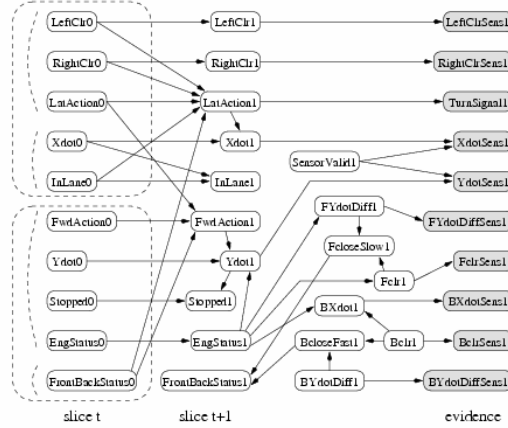
- Error does not grow unboundedly!



Error ϵ is error of approx. one step by JT. (or assumed density)

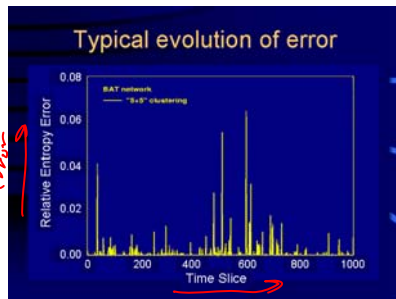
- **Theorem:** If Markov chain **contracts** at a rate of γ (usually very small), and **assumed density projection** at each time step has **error bounded by ϵ** (usually large) then the **expected error at every iteration is bounded by ϵ/γ** .

Example – BAT network [Forbes et al.]



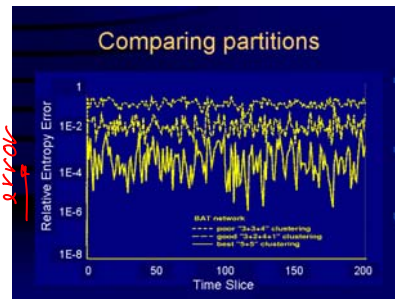
15

BK results [Boyen, Koller '98]



time
bounded error

Spikes because some observations introduced lots of error (that gets contracted later)



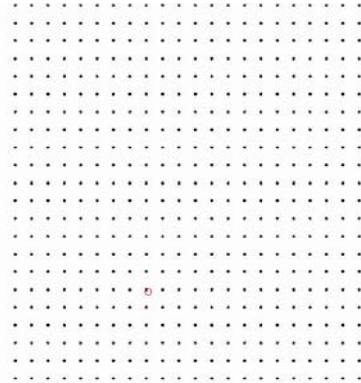
time

different clusters of S.T.

16

Thin Junction Tree Filters [Paskin '03]

- BK assumes fixed approximation clusters
- TJTF adapts clusters over time
 - attempt to minimize projection error



17

Hybrid DBN (many continuous and discrete variables)

- DBN with large number of discrete and continuous variables
- # of mixture of Gaussian components blows up in one time step!
- Need many smart tricks...
 - e.g., see Lerner Thesis



Figure 10.1: The prototype RWGS system

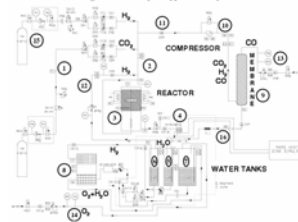


Figure 10.2: The RWGS schematic

Reverse Water Gas Shift System (RWGS) [Lerner et al. '02]

18

DBN summary

- **DBNs**

- factored representation of HMMs/Kalman filters
- sparse representation does not lead to efficient inference

- **Assumed density filtering**

- BK – factored belief state representation is assumed density
- Contraction guarantees that error does not blow up (but could still be large)
- Thin junction tree filter adapts assumed density over time
- Extensions for hybrid DBNs

19

This semester...

- Bayesian networks, Markov networks, factor graphs, decomposable models, junction trees, parameter learning, structure learning, semantics, exact inference, variable elimination, context-specific independence, approximate inference, sampling, importance sampling, MCMC, Gibbs, variational inference, loopy belief propagation, generalized belief propagation, Kikuchi, Bayesian learning, missing data, EM, Chow-Liu, IPF, GIS, Gaussian and hybrid models, discrete and continuous variables, temporal and template models, Kalman filter, linearization, switching Kalman filter, assumed density filtering, DBNs, BK, Causality,...

■ **Just the beginning...** 😊

20

Quick overview of some hot topics...

- Conditional Random Fields
- **Maximum Margin Markov Networks**
- **Relational Probabilistic Models**
 - e.g., the parameter sharing model that you learned for a recommender system in HW1
- **Hierarchical Bayesian Models**
 - e.g., Khalid's presentation on Dirichlet Processes
- **Influence Diagrams**

21

Generative v. Discriminative models – Intuition

- **Want to Learn:** $h: X \mapsto Y$
 - X – features
 - Y – set of variables
- **Generative classifier**, e.g., Naïve Bayes, Markov networks:
 - Assume some **functional form for $P(X|Y)$, $P(Y)$**
 - Estimate parameters of $P(X|Y)$, $P(Y)$ directly from training data
 - Use Bayes rule to calculate $P(Y|X=x)$ *$P(Y|X=\text{text})$*
 - This is a '**generative**' model *web page*
 - **Indirect** computation of $P(Y|X)$ through Bayes rule
 - But, can **generate a sample of the data**, $P(X) = \sum_y P(y) P(X|y)$
- **Discriminative classifiers**, e.g., Logistic Regression, Conditional Random Fields:
 - Assume some **functional form for $P(Y|X)$** *y^* : argmax $P(Y|X=z)$*
 - Estimate parameters of $P(Y|X)$ directly from training data
 - This is the '**discriminative**' model
 - Directly learn $P(Y|X)$, can have **lower sample complexity**
 - But **cannot obtain a sample of the data**, because $P(X)$ is not available

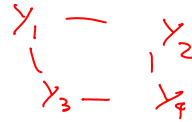
22

Conditional Random Fields

[Lafferty et al. '01]

- Define a Markov network using a log-linear model for $P(Y|X)$:

$$P(Y|X) = \frac{1}{Z(X)} e^{\sum_i w_i f_i(x,y)}$$



- Features, e.g., for pairwise CRF:

$$f_{17}(y_1, y_3, X)$$

- Learning: maximize conditional log-likelihood

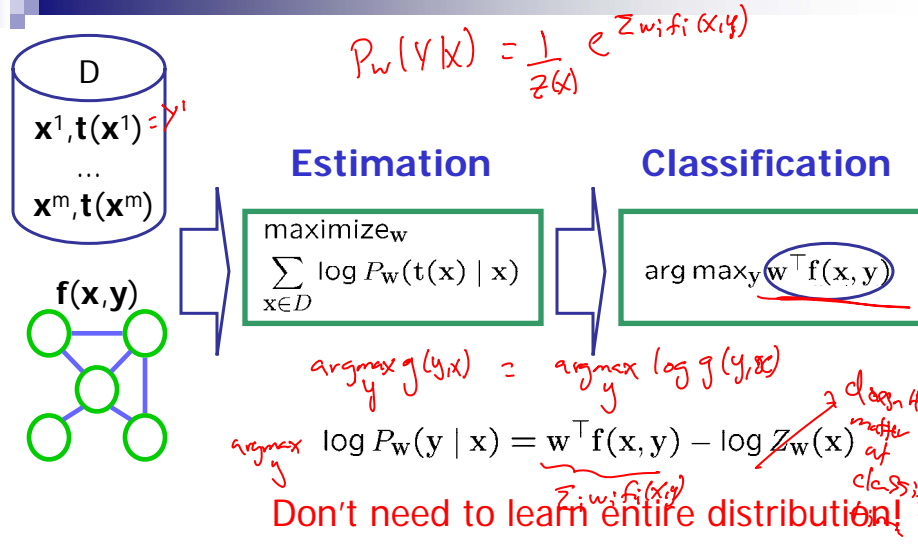
$$\hat{w} = \underset{w}{\operatorname{argmax}} \log P_w(Y|X) \quad \text{? in data}$$

- sum of log-likelihoods you know and love...
- learning algorithm based on gradient descent, very similar to learning MNs

$$\log P_w(Y_D | X_D) = \sum_i \log P(y^{(i)} | X = x^{(i)})$$

23

Max (Conditional) Likelihood



24

OCR Example

- We want:

$$\operatorname{argmax}_{\text{word}} \mathbf{w}^T \mathbf{f}(\text{brace}, \text{word}) = \text{"brace"}$$

don't need $\gamma(x)$

- Equivalently:

$$\mathbf{w}^T \mathbf{f}(\text{brace}, \text{"brace"}) > \mathbf{w}^T \mathbf{f}(\text{brace}, \text{"aaaaa"})$$

$$\mathbf{w}^T \mathbf{f}(\text{brace}, \text{"brace"}) > \mathbf{w}^T \mathbf{f}(\text{brace}, \text{"aaaab"})$$

...

$$\mathbf{w}^T \mathbf{f}(\text{brace}, \text{"brace"}) > \mathbf{w}^T \mathbf{f}(\text{brace}, \text{"zzzzz"})$$

a lot!

25

Max Margin Estimation

- Goal: find \mathbf{w} such that

$$\mathbf{w}^T \mathbf{f}(\mathbf{x}, \mathbf{t}(\mathbf{x})) > \mathbf{w}^T \mathbf{f}(\mathbf{x}, \mathbf{y}) \quad \forall \mathbf{x} \in D \quad \forall \mathbf{y} \neq \mathbf{t}(\mathbf{x})$$

$$\mathbf{w}^T [\mathbf{f}(\mathbf{x}, \mathbf{t}(\mathbf{x})) - \mathbf{f}(\mathbf{x}, \mathbf{y})] > 0$$

$$\mathbf{w}^T \Delta \mathbf{f}_{\mathbf{x}}(\mathbf{y}) \geq \gamma \Delta t_{\mathbf{x}}(\mathbf{y})$$

- Maximize margin γ
- Gain over \mathbf{y} grows with # of mistakes in \mathbf{y} : $\Delta t_{\mathbf{x}}(\mathbf{y})$

$$\Delta t_{\text{brace}}(\text{"craze"}) = 2$$

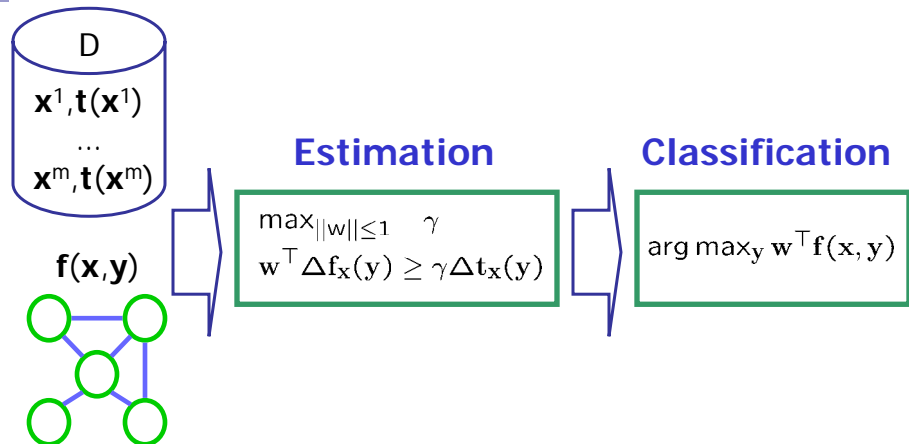
$$\Delta t_{\text{brace}}(\text{"zzzzz"}) = 5$$

$$\mathbf{w}^T \Delta \mathbf{f}_{\text{brace}}(\text{"craze"}) \geq 2\gamma$$

$$\mathbf{w}^T \Delta \mathbf{f}_{\text{brace}}(\text{"zzzzz"}) \geq 5\gamma$$

26

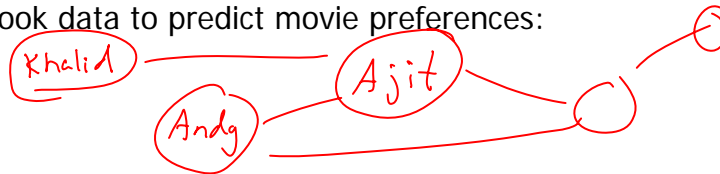
M³Ns: Maximum Margin Markov Networks [Taskar et al. '03]



27

Propositional Models and Generalization

- Suppose you learn a model for social networks for CMU from FaceBook data to predict movie preferences:



- How would you apply when new people join CMU?
- Can you apply it to make predictions a some "little technical college" in Cambridge, Mass?

28

Generalization requires Relational Models (e.g., see tutorial by Getoor)

- Bayes nets defined specifically for an instance, e.g., CMU FaceBook today
 - fixed number of people
 - fixed relationships between people
 - ...
- Relational and first-order probabilistic models
 - talk about objects and relations between objects
 - allow us to represent different (and unknown) numbers
 - generalize knowledge learned from one domain to other, related, but different domains

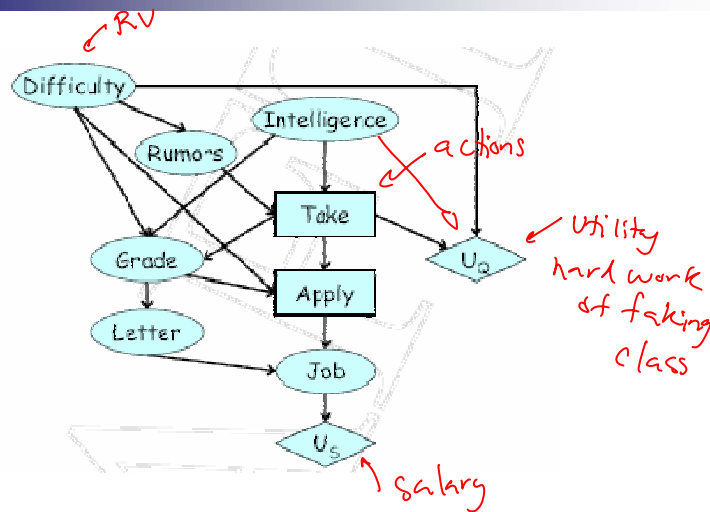
29

Reasoning about decisions K&F Chapters 20 & 21

- So far, graphical models only have random variables
- What if we could make decisions that influence the probability of these variables?
 - e.g., steering angle for a car, buying stocks, choice of medical treatment
- How do we choose the best decision?
 - the one that maximizes the expected long-term utility
- How do we coordinate multiple decisions?

30

Example of an Influence Diagram



31

Many, many, many more topics we didn't even touch, e.g.,...

- Active learning
- Non-parametric models
- Continuous time models
- Learning theory for graphical models
- Distributed algorithms for graphical models
- Graphical models for reinforcement learning
- Applications
- ...

32

What next?

■ Seminars at CMU:

- Machine Learning Lunch talks: <http://www.cs.cmu.edu/~learning/>
- Intelligence Seminar: <http://www.cs.cmu.edu/~iseminar/>
- Machine Learning Department Seminar: <http://calendar.cs.cmu.edu/cald/seminar>
- Statistics Department seminars: <http://www.stat.cmu.edu/seminar>
- ...

■ Journal:

- JMLR – Journal of Machine Learning Research (free, on the web)
- JAIR – Journal of AI Research (free, on the web)
- ...

■ Conferences:

- UAI: Uncertainty in AI
- NIPS: Neural Information Processing Systems
- Also ICML, AAAI, IJCAI and others

■ Some MLD courses:

- 10-705 Intermediate Statistics (Fall)
- 10-702 Statistical Foundations of Machine Learning (Spring)
- 10-801 Advanced Topics in Graphical Models: statistical foundations, approximate inference, and Bayesian methods (Spring)
- ...

33