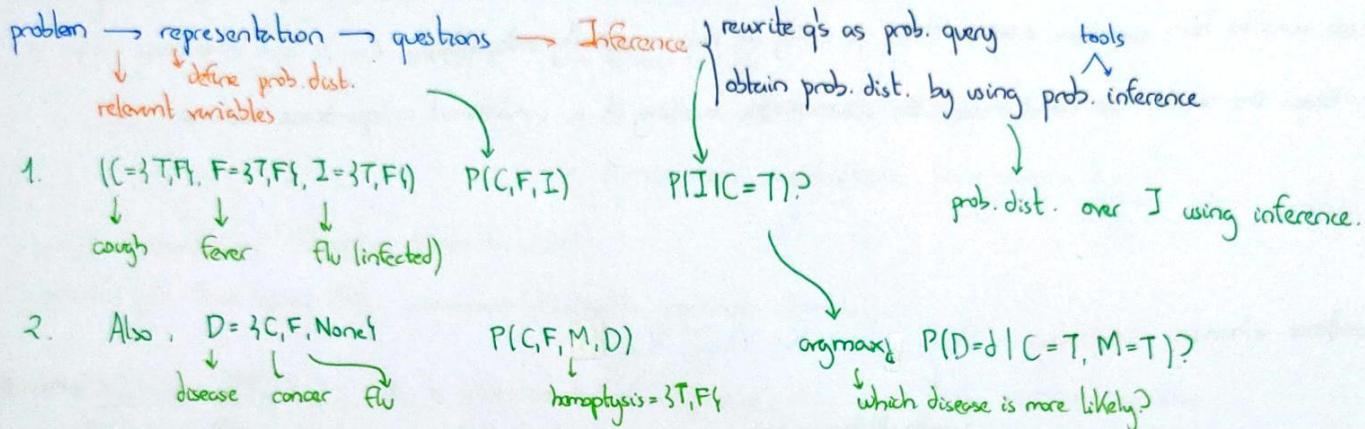
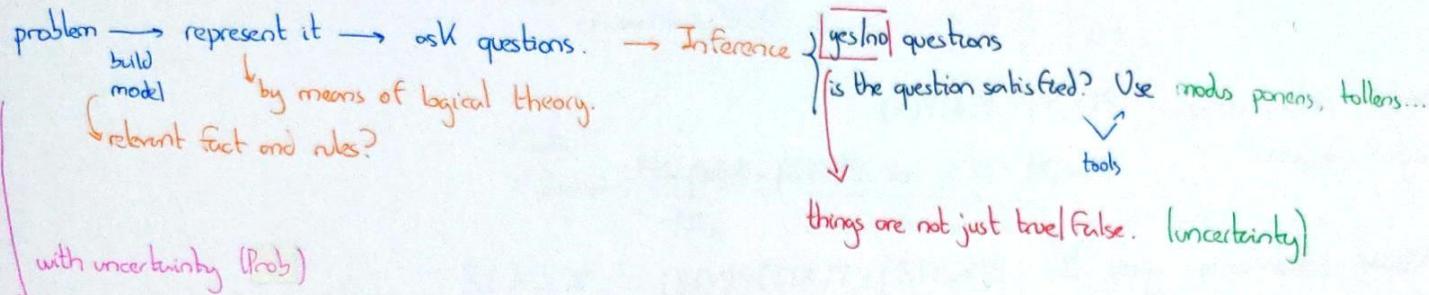


# Probabilistic Graphical Models

## Lecture 1. Introduction

- Better to give a probabilistic distribution than a static decision (this and just this)



- Types of queries:
  - Conditional probabilities queries: find marginal distribution  $(P(I | C=T))$
  - MAP queries: find value of that maximizes a prob distribution

PGMs get involved when: we don't have an expert to know the probabilistic distributions: just data

problem → representation → question → learning: choose most likely prob dist.

↓ set of possible prob. dist. of variables given data

- Bayesian approach: initial beliefs (prior) → refine.

Frequentist: coin HHTTTT:  $\theta = 5/6$

Bayesian: initial belief ( $1/2$ ) +  $\uparrow p(\theta)$

- Probabilistic approach to Data Science: Assumptions → Identify Patterns → Predict & Evaluate



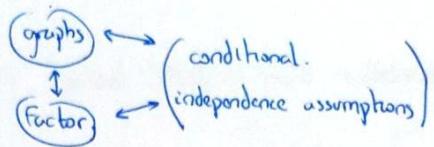
- Big/Huge Data: we can't reason on a joint prob. dist. over all variables at some time → PGM.

- encode prob. dist. over large number of variables in a compact way
- learn from available data
- efficiently answer q's

} PGM → use graphs!

easy to viz, math adj.  
nodes: var  
edges: possible relations.

- PGMs: graphs + set of possible prob. dists
- types: directed acyclic, undirected, bipartite
- it is a factorization of the prob. dists.



- Independent random variables:  $P(X, Y) = P(X)P(Y)$

$\downarrow$   
marginal independence       $P(X|Y) = P(X)$  or  $P(Y|X) = P(Y)$

$X \perp\!\!\!\perp Y$

- Conditional independence: given  $Z$ :  $P(X, Y|Z) = P(X|Z)P(Y|Z)$

- when two variables have common origin, they generally are not marginally independent.
- if we know the source, we could model the intermediates, making it a conditional independence scenario.



- independence enhances simplicity:  $P(X)$  with  $X = Y \cup Z$ ,  $Y \perp\!\!\!\perp Z$

$$\begin{array}{c} \downarrow \\ X \text{ with 10 binary vars} \end{array} \qquad \begin{array}{c} \downarrow \\ \text{each has 2 bin vars} \end{array} \Rightarrow \begin{array}{c} \frac{1024}{2^{10}} \ggg \frac{64}{2 \cdot 2^5} \\ \text{parameters.} \\ \text{efficient.} \end{array}$$

- Structural Statistical Model: two variables  $X, Y \xrightarrow{\text{two poss.}}$

1.  $M_1 = \emptyset \rightarrow$  they are not independent.

2.  $M_2 = \{X \perp\!\!\!\perp Y\}$

} generalize to  $n$ .

- A distribution respects model  $M$  if it satisfies all independencies included in  $M$ .

- Factorization: reexpression into product of simpler factors.

prob. dist.  $\downarrow$

$$P(X_1, \dots, X_{10}) = P(X_1, \dots, X_5)P(X_6, \dots, X_{10}) \text{ using independence.}$$

- also: we can move to smaller spaces:  $P(X, Y, Z) = P(X|Z)P(Y|Z)P(Z)$

- Inference: queries, prob. based operations.
  - exact marginalization, conditioning, belief propagation...
  - approx: random sampling

- Learning
  - parametric
  - structural (NP)

For robust results: tradeoff  $\rightarrow$  model complexity  $\curvearrowleft$  available data.

## Lecture 2. Probability Overview.

- Learn  $P(X|X)$  instead of  $F:X \rightarrow X$
- Joint distribution of M variables:  $2^M$  combinations  $\rightarrow$  learn a value per each (probabilities)

$$\downarrow P(E) = \sum_{r \text{ rows that match } E} P(r)$$

x	y	z	prob
0	0	0	0.1
0	1	0	
0	0	1	
0	1	1	

$$P(E_1 | E_2) = \frac{P(E_1 \wedge E_2)}{P(E_2)} = \frac{\sum \text{rows that match both}}{\sum \text{rows that match } E_2}$$

- learn joint distribution from train  $\rightarrow$  use it to calculate using test data.

problem: sometimes data is not enough.  $\rightarrow$  # rows in table

solutions  
↓

# samples needed to learn faithfully

# uncommon combinations (not enough observed)

smart representations (Bayesian networks, PGM)

estimate prob. from sparse data (maximum likelihood/a posteriori estimates)

- If  $p(x_A|x_B, x_C) = p(x_A|x_B)$ ,  $x_A$  is independent from  $x_B$  given  $x_C$ .  
not  $\equiv x_A \perp\!\!\!\perp x_B | x_C$

not: uppercase = variables  
lowercase = value

- prod:  $x_A \perp\!\!\!\perp x_B | x_C \Rightarrow p(x_A, x_B | x_C) = p(x_A | x_C) \cdot p(x_B | x_C)$

- Chain rule:  $p(x, y, z) = p(x|y, z) \cdot p(y|z) \cdot p(z)$

↳ allows to simplify the factorization

$$x = x_1, \dots, x_5 \text{ with } 3 \perp\!\!\!\perp 4 \perp\!\!\!\perp 1, 5 \Rightarrow \text{order } \overbrace{1, 4, 5, 3, 2}^Z \times x$$

$$p(x) = p(x_{1,4,5}) p(x_3 | x_{1,4,5}) p(x_1 | x_{1,3,4,5})$$

(is redundant. (due to independency))

- Let A, B be a partition of V. Then,  $\forall x_A, x_B, p(x_A) = \sum_{x_B} p(x_A, x_B)$

We can have a marginal dist for each value  $X_B = x_B : \sum_{x_A} p(x_A | x_B) = 1$

Also,  $p(x_A | x_B) = \frac{p(x_A, x_B)}{p(x_B)} = \frac{p(x_A, x_B)}{\sum_{x_A} p(x_A, x_B)}$   
↓  
post definition.

} collaborates by reducing number free params.  
 $p(x) = r_x - 1$   
 $p(x,y) = r_x \cdot r_y - 1$   
 $p(x|y) = (r_x - 1) \cdot r_y$

} compute number of free params.

- Count parameters without ind.

$$\begin{array}{ccccc} A \setminus B & + & - & & \\ & + & a & b & \\ & - & c & 1-a-b-c & \end{array} \Rightarrow 3$$

with ind:

$$\begin{array}{ccccc} A \setminus B & + & - & & \\ & + & b & & \\ & - & (1-b) & & \\ & & a & (1-a) & \end{array}$$

$$\begin{array}{ccccc} A \setminus B & + & - & & \\ & + & a \cdot b & b(1-a) & \\ & - & (1-a \cdot b) & (1-a) \cdot (1-b) & \end{array}$$

Less complexity, avoiding modelling irrelevant params. linked to soft conditional dependencies, helps manage trade-off complexity  $\hookrightarrow$  data  
↓  
overfitting

## &gt; Exercise:

We want  $P(H=+ | V=-)$

$$\text{We have } P(V=+ | H=+) = 0.43 \\ P(V=- | H=+) = 0.57$$

$$P(V) \left| \begin{array}{l} + = 0.84 \\ - = 0.16 \end{array} \right.$$

• Bayes' rule:

$$P(X|V) = \frac{P(V|X) \cdot P(X)}{P(V)}$$

$$= P(V|X) \cdot P(X) = P(X|V) \cdot P(V) =$$

chain rule by one side

chain rule by the other

$$P(H=+ | V=-) = \frac{P(V=- | H=+) P(H=+)}{P(V=-)}$$

missing

STILL, we can tell the increased risk by not getting married:

$$\frac{P(H=+ | V=-)}{P(H=+ | V=+)} \rightarrow P(H=+) \text{ disappears}$$

||  
7 (7 times more)

## Lecture 3. Bayesian Networks.

- Directed acyclic graph  $G = (V, E)$  + parameters  $\theta$

variables

chain rule

table of cond. prob. dist.

- Ancestral set: todos los antepasados de un nodo.

- Ancestral ordering: from older to younger (once you visit one node, all ancestors were visited)

- Moral: convert to undirected graph (remove directions + unite nodes with common children)

## &gt; Exercise.

1. Yes.

2. parents of  $X_3$ :  $X_5$

children:  $X_2 / X_6$

parents of  $X_5$ :  $\emptyset$

children:  $X_2, X_3$

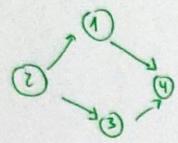
3.  $X_5 \ X_3 \ X_2 \ X_1 \ X_6 \ X_4$

No

- Factorization in Bay. Network:

$$P_M(x) = \prod_{i=1}^n p(x_i | \text{pa}_i; \theta_i)$$

↓  
set of variables pointing towards  $x_i$



$$p(x) = p(X_4 | X_1, X_3) p(X_1 | X_2) p(X_3 | X_2) p(X_2)$$

(chain rule)

ancestral ordering

$$p(x) = p(X_4 | X_1, X_2, X_3) p(X_1 | X_2, X_3) p(X_2)$$

$$X_1 \perp\!\!\!\perp X_3 | X_2 \\ X_4 \perp\!\!\!\perp X_2 | X_1, X_3$$

} independence of the rest  
given the parents.

> Exercise 1.  $X_2, X_1, X_5, X_6, X_4, X_3$

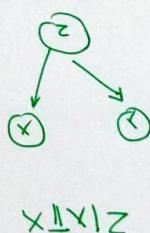
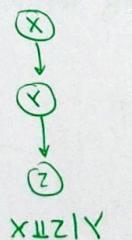
$$2. P(x) = P(X_3 | X_1, X_5) P(X_4 | X_6, X_1) P(X_6 | X_2) P(X_5 | X_1) P(X_1 | X_2) P(X_2)$$

(prob[each node | parents])

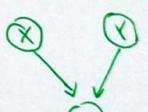
3. c)

d) also, but with way too many repeating parameters.

- Influence flow is stopped



Influence starts if Z is observed.

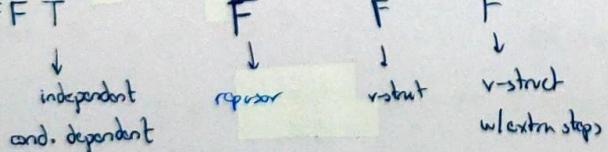


$$X \not\perp\!\!\!\perp Z$$

$$\text{while } (X \not\perp\!\!\!\perp Y) \\ \Rightarrow X \perp\!\!\!\perp Z$$

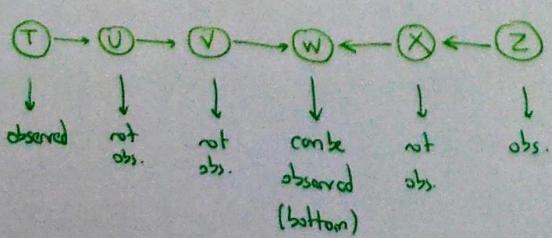
net gross example.

> Exercise: FFT



## Lecture 4.

- Active trail: trail  $X_1 \rightleftharpoons \dots \rightleftharpoons X_m$  of probabilistic influence flow  
it is active when



only time you can have an observation in middle of an active track (except v-structs) \*

↳ this observation activates the flow.

- $\delta$ -separation: generalization of conditional independence to set of variables

↳ independent if there is no activations. (

how to check: 1) identify ancestors of x,y,z

2) remove the rest, moralize graph

3) does Z block all the paths from X to Y?

1

\* observation = no activation  $\Rightarrow$  independence holds.

no obs. = active trail

↳ generates new paths from  $r$ -structures.

(by moralizing you can see all the trials)

- **Markov's blanket:** <sup>minimal</sup> set of variables you need to observe in order to make a variable independent from the rest.  
 ↳ parents, children and children's parents. useful to avoid useless observations → ML efficiency.  
 Block all the paths + the influence of v-structs.

- **Redundancy:** two DAGs produce some dependence model if they have the same edges (mod. direction) and some v-structures.

- $I(G)$  is the set of independence statements that hold on  $G$ .    }  
 $I(P)$     "    "    "    "    "    "    "    }  $G$  is an  $I$ -map of  $P$  if  $I(G) \subseteq I(P)$   
 $\downarrow$   
 $I(G)$

The other way around may not be true (wasted params.)

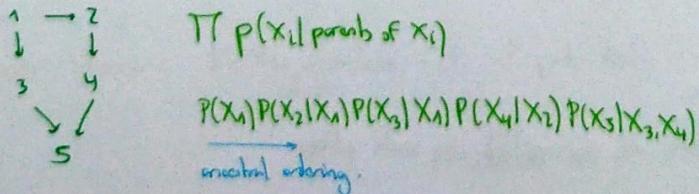
assume dependencies that are not real

> J-map quiz: ②)

- Given a Bay. Net.  $G$ , a prob. dist  $P$ .

$P$  factorizes according to  $G \Leftrightarrow P$  factorizes.

- ## • Factorization of

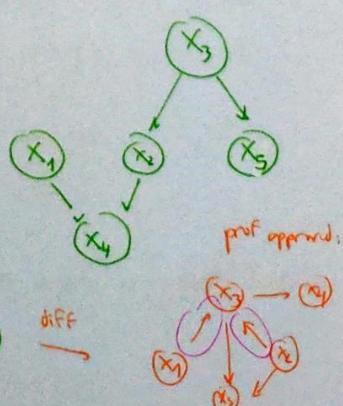


$$J = \{ x_3 \amalg x_1 | x_3; \quad x_2 \amalg x_1, \quad x_5 | x_3; \quad x_3 \amalg x_1; \quad x_4 \amalg x_3, x_3 | x_1, x_2 \}$$

$$\text{chain rule} \downarrow P(X_1|X_1, z_{1,3})P(X_2|X_1, X_3, x_8) \cdot P(X_3|X_{1,3})P(X_1) \cancel{P(X_3)}$$

↓  
start from longest

→ Deleted due to independencies.



For  $P(E)$  we need 2 parameters.  
<sub>0,1,2</sub>

↳ numbers of params in a Bay Network.

## Lecture 5. Markov Networks

Some scenarios are not representable by a Bayesian Network.

- Markov Networks:
  - undirected graph
  - Factors over cliques (instead of dist. per variable + parents)
  - requires normalization: partition function
  - independence determined by separation  $\xrightarrow{\text{Bay.}}$  instead of d-sep.
  - $F \xleftarrow[\text{only for pos. dist.}]{\text{always}} I \xrightarrow[\text{always}]{\text{always}}$
  - Markov blanket = neighbours

Lack of v-structures.

- $\phi$  is not a probability distribution; partition function

$\phi_1(A, B)$	graph	undirected	directed.
"strength"			
A, B			
0 0 30			
0 1 5			
1 0 1			
1 1 100			
	factorization	$P(x) = \frac{1}{Z} \prod_{i=1}^F \phi_i(x_{\phi_i})$	hard to compute
		Factorization constant: $Z = \sum_x$	
	Factors	$\phi \leftarrow \text{partition func.}$	

(instead of  $P(x) = \prod_{i=1}^V p(x_i | p_{a_i})$ )

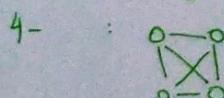
prob. dist for each value of the parents.

$X_1$	$X_2$	$P(X_1   X_2)$
a	b	0.3
b	b	0.7
a	a	0.6
b	a	0.4

- Clique: set of variables, where all <sup>vertices</sup> of them are connected to each other.

1-vertex: all single variables.

2- : connected variables



- maximal clique: a clique not contained by any other clique.

## &gt; Factorization in MNs Exercise

1. Valid (is over cliques)
2. Non-valid
3. Valid (over maximal cliques)
4. Valid
5. Valid (all pair-wise cliques)

&gt; 1. X

2. N: A,B,D not a clique

3. N:

4. N

- **Independence:** if there is no path between the variables.

- **Conditional independence:** if the condition variable blocks any connection between the two variables.

- Active trail  $\textcircled{X} \rightarrow \textcircled{Y} \rightarrow \textcircled{Z} \rightarrow \textcircled{W} \rightarrow \textcircled{T}$

$\underbrace{\quad\quad\quad}_{\text{if none is observed.}}$

&gt; Independence Ex

None

&gt; Conditional Ind. Ex.

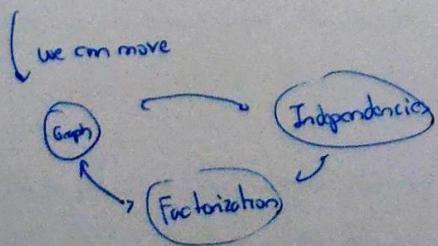
- a) F
- b) T
- c) T
- d) F

- **Markov blankets:** neighbours of a node (because we don't have v-structs)  
 ↳ if these are observed, observing any other doesn't influence at all)

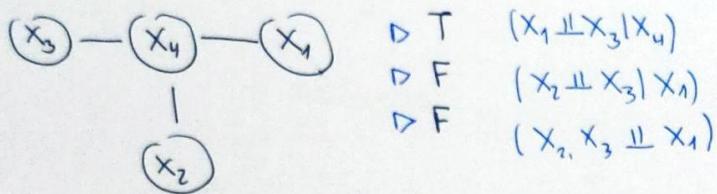
- **Hammersley Clifford Theorem:** Markov structure  $Z_l$ , P prob. dist' l.

P factorizes over  $Z_l \Rightarrow P$  satisfies independencies in  $Z_l$

$P \text{ satisfies ind. in } Z_l \\ P \text{ is positive} \quad \left. \right\} \Rightarrow P \text{ factorizes.}$



> Exercise



> Exercise (provide a Factorization)

$$P \propto \phi_{\{2, 6, 4, 7\}} \phi_{\{1, 3, 4\}} \phi_{\{1, 2, 3\}} \phi_{\{3, 5\}} \rightarrow \text{maximal cliques}$$

- prop: A BN without V-structs  $\equiv$  MN

↑  
normalization

- Number of parameters for a Markov Networks

$$\text{(Bayesian): } \sum_{\substack{\forall \text{ vars} \\ \text{all vars}}} \Omega_{PA} (\Omega_{Xv} - 1)$$

↓ all parent values

$X_4$	$X_2$	$P(X_4, X_2)$
a	a	0.4
a	b	0.6

need 2 params.

$$\sum_{\substack{F \\ \text{factors with Scope } S_\phi}} \Omega_{S_\phi} \quad \rightarrow \text{nothing else! It is not a prob. dist.}$$

↓  
all poss. values in the scope

↓  
MN;  
would need 4.

- Applications:
  - denoising: one observed variable  
one hidden variable } for each pixel
  - segmentation (similar)

→ input image  
→ models real values.

- Factor Algebra:

$$X \times \phi$$

0	0	3
---	---	---

$$X \times \psi$$

0	0	5
---	---	---

$$X \times X \times Z \times \phi \times \psi$$

0	0	0	15
0	0	1	12

→  
always larger  
or equal table

commutative, associative,  $\exists$  neutral element

- reduction: observe a set of variables (reduce table by fixing variables, and selecting rows where values apply)

multiple values = AND.

$$-\text{prop: } (\phi_1 \times \phi_2)[U=u] = \underbrace{\phi_1[U=u] \times \phi_2[U=u]}_{\text{uses smaller tables.}}$$

- marginalization: sum the rows that have the specified values

$$\begin{array}{cc} X & X \\ 0 & 0 \end{array} \quad \sum_{z=0}^1 \phi$$

sum of 000 and 001

- prop:  $\sum_x (\phi_1 \times \phi_2) = \phi_1 \times \sum_x \phi_2$

↓  
if  $x \notin \text{Scope}(\phi_1)$   
generalization:  $\sum_x (\phi_1 \times \phi_2 \times \phi_3) = \phi_1 \times \sum_x (\phi_2 \times \phi_3)$

> Reduction and marginalization exercise

$$p(X|Y=1)$$

reduce:	X	Z	$\phi[X=1]$
	0	0	11
	0	1	2
	0	2	1
	1	0	4
	1	1	1
	1	2	9

marginalize

X	$\sum_z$
0	14
1	14

normalize  
(for prob dist)

X	$p(X Y=1)$
0	0.5
1	0.5

> 1. Reduce  $\phi_1[B=b]$  ...

Reduce  $\phi_2[B=b]$ .

$$P(A,C|B) = \frac{1}{Z} \phi_1[B=b] \phi_2[B=b] \sum_D (\phi_3 \cdot \phi_4)$$

$$3. P(A|B=b) = \frac{1}{Z} \sum_B \left( \sum_A (\phi_1 \cdot \phi_4) \phi_2[C=c] \right) \phi_3[C=c]$$

|| → more efficient

$$\frac{1}{Z} \sum_A (\phi_4 \cdot \sum_B (\phi_1 \cdot \phi_2[C=c])) \cdot \phi_3[C=c]$$

## Lecture 6.

- normalization: it will always result in one variable and a value for all its possible values.

> Exercise "Factors in Markov Network": b), c), e)

product of everything and getting rid of all the variables (summing): normalizing factor to get a prob. dist.

## Plate Models.

notation: boxes mean that the variables inside are repeated n times. (also called plate)

stochastic variables those that are usually observed.

> Plate semantics exercise: a), b), d)

> Plate interpretation exercise: b)

## > Grounded Plates: b)

- limitations: no connections between copies of the same plates.  
no connections between specific variables depending on another variable.

• Temp<sub>plate</sub>(local) late models: future and past are independent given the present.  
time invariance  $P(X(t+1)|t)$  always the same  $\forall t$

- Dynamic Bayesian Network
  - 1)  $X(0)$  uses Bayesian Network
  - 2)  $X(t+1)$  given  $X(t)$   $\rightarrow$  Conditional BN

> Applications of DBNs:
 

- a) F  $\rightarrow$  because it's connecting to several connections
- b) T
- c) F
- d) T

> Independencies in DBNs:
 

- a)
- b)
- c)
- d)

## Exact Inference

• Conditional Probability Query: Find the distribution  $P(X|E=e)$  given a partition of variables  $X = (X, H, E)$

Evidence ↑  
 ↓ exponential complexity (even if P factorizes to graph  $T()$ )  
 Hidden ↓  
 not polynomial time except for some graphs (specific)

- Variable Elimination (for 1 query)
- Message Passing (for several queries)

- ↓
1. Eliminate all variables in  $X$   $\rightarrow$ 
    - a) for all variables that have said variable  $V$  in the scope
    - b) keep a copy without them
    - c) product between all them and  $V$
    - d) Marginalize  $V$
    - e) sum b) and d) and return
  2. Normalize

> Visualizing VE: start by  $X_1$  or  $X_8$  because they're in only one factor

$$\sum_{X_6} \phi(X_6, X_8) = z_7(X_6)$$

↳ Then continue with  $X_6$  ( $\text{scope}(X_6) \subseteq \text{scope}(X_4, X_6, X_7)$ )  $\rightarrow$  no large intermediate tables.  
but also continuing with nodes with least factors! ( $X_2$  step = bad)

- Induced graph: add those edges that appear during elimination.

↳ the more extra edges = the more complex our elimination order was.

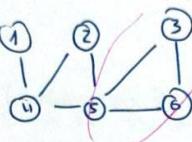
a good heuristic is the size of the biggest clique in the induced graph.

> VE ordering exercise: 1, 3, 6, 2, 8, 7, 5

- with evidence: observe and marginalize.

## > Intermediate Factors (Bayesian!)

moralize! →



$$\rightarrow P(X_6 | X_3, X_5) = f(X_5, X_3, X_6)$$

$X_5 \leftarrow \text{answer. } (F_1(X_4, X_2, X_5) \text{ and } f(X_3, X_5, X_6))$

## Lecture 7.

[20 mins late]

## • Variable Elimination For MAP (argmax)

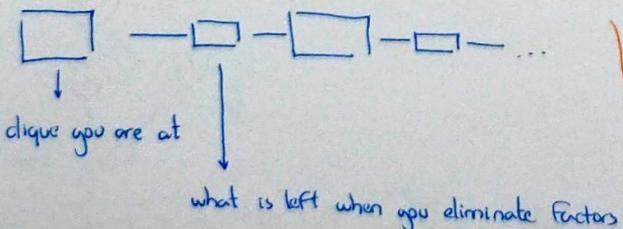
- sum of factors: like product but summing

- max-marginal: like marginalization but taking the maximum instead of summing.

## • Message Passing:

we are going to use cluster graphs and clique tree

- clique tree:



} very related to Variable Elimination

- cluster graph each node is associated with a subset (cluster)

clusters are related to each other with the intersection or c<sub>n</sub> (overlap)

$$\text{set of beliefs } \left\{ \begin{array}{l} \text{one factor for each subset } B_i \\ \text{... " " " } \\ \text{subset } B_i \end{array} \right. \rightarrow P_B \frac{\prod B_i}{\prod \mu}$$

now  $p(X_j) = \frac{1}{Z} \sum_{C_i \ni x_j} \beta_i(C_i)$  such that  $x_j \in C_i$  (only cluster that contains  $x_j$  and marginalize one only)

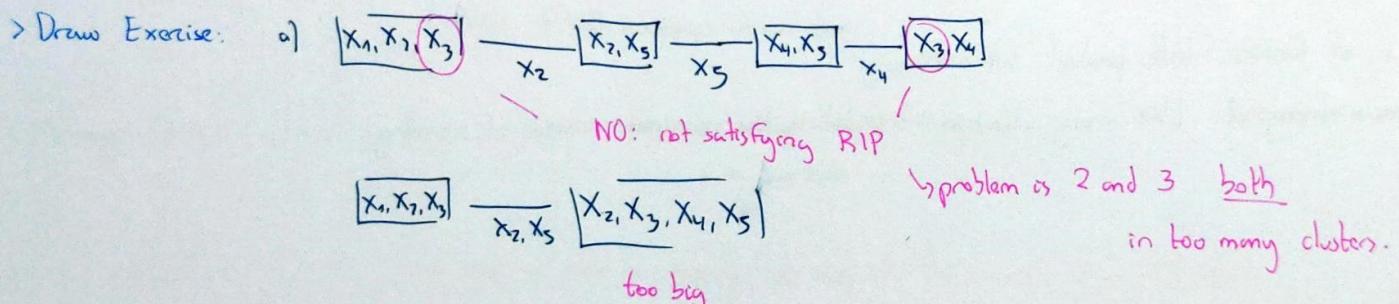
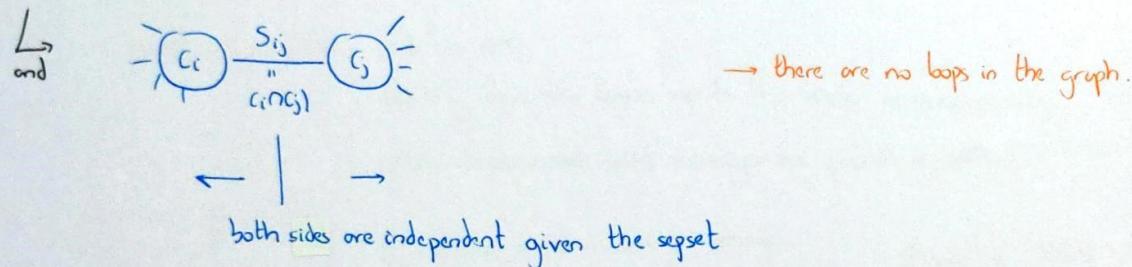
the transformation is expensive  
but the queries are really cheap } Great for many queries.

- building clusters:
  - \* **preservation:** all factors  $\phi_j$  need a cluster  $C_i$  such that  $\text{Scope}(\phi_j) \subseteq C_i$ .
  - \* **running intersection property:** if there is a variable in several clusters.
    - 1) There is a unique path among them.
    - 2) The variable is present in all intermediate clusters and sepsets.

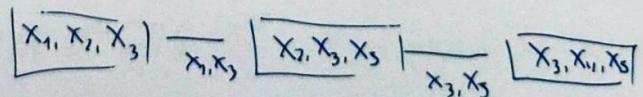
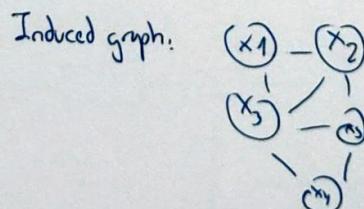
> Cluster Graph construction: a), b)

> Family perservation: c) (by definition of perservation)

in clique trees, the sepset is exactly the intersection

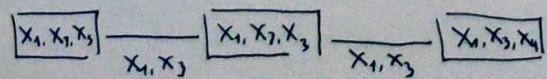


Think of Variable Elimination! Good order: 1, 4, 2, 3, 5



b) Similar. Order: 5, 2, 1, 3, 4

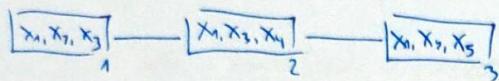
Induced Graph:



- cluster factor:  $\psi_i(C_i) = \prod \phi$

↓  
all the factors that were assigned to that cluster.

- message: From one cluster to another



$$\delta_{2 \rightarrow 3} = \sum \psi_{(2)} \cdot \delta_{1 \rightarrow 2}$$

↓                    ↓  
cluster Factor of    all incoming messages of the sender.

marginalize over  
all variables not  
in 2,3 scopes

belief propagation: select a root.

starting from the leaves up to the root: message passing.

select another root until messages are passed in both directions in all edges.

$B_i$  = Factor \* all incoming messages.

$m_{ij}$  = mult. of both messages in an edge.

> Message Ordering Exercise: d)

b)  $\rightarrow$  not enough messages ✗

a)  $\rightarrow$  two roots ✗

↓  
every time we send a message, we need all the incoming msgs. beforehand.

$C_1 \rightarrow C_2, C_2 \rightarrow C_3, (C_3 \rightarrow C_4, C_5 \rightarrow C_3)$  would be wrong.

> Message passing in cluster tree: a) we marginalize getting rid of the variables in 3 not in sepset.  
x all incoming msgs.

>

b)

- cluster graphs are calibrated when the  $B_i$  marginalized for all variables in a sepset is the same for both clusters.

$$\sum_{G \setminus S_j} B_i = \sum_{G \setminus S_j} B_j$$

not met with loops, etc. (loop = lack of leaves = message passing is not constant) \*\*  
met automatically in clique trees.

- Exact inference:

1. Represent as clique tree

2. Calibrate

3. Answer

} not polynomial.

- alternatively, initialize all messages to 1  
pass messages until it is calibrated enough.  
it won't be 100% calibrated  $\Rightarrow$  not exact answers. ~~exact~~ inference.

## Lecture 8. Approximate Inference

Also sampling

$\hookrightarrow$  larger  $M$ , better estimations.

For conditional distribution: restrict to those samples that satisfy the condition.

> (At. Dist. c), d), a), d), b, b, b ...

- Forward sampling: follow ancestral ordering
  - $\hookrightarrow$  for each variable you sample it
- Conditional probability query: filter the sample to those that fulfill  $E=e$ .
  - $\downarrow$  you get one sample; for answering queries you need to repeat a bit
  - might have very low prob.  $\Rightarrow$  no samples.

> Rejecting samples: e)

$$\text{Imagine } P(e) = \frac{1}{10^6} \text{ very small } \Rightarrow M = M' \cdot 10^6 = 10^9 \rightarrow \text{samples you need to have.}$$

$\downarrow$

$M' = 1000$   
(samples you want)

instead of generating, we may weigh samples.

these samples will always fulfill evidence  $E=e$ .

Instead of summing samples, weighted-sum them to answer queries.

- CPQ for Markov Networks: how to sample? Factors are not prob. dist.
    - no ancestor ordering (some can have inner loops)
- $\} \longrightarrow$  other tools.

- Gibbs sampling:
    - 1) pick a sample initial state
    - 2) Keep iterating to update the distributions  $\leftarrow$  variables are not independent (contrary to assumptions)
- $\downarrow$
- most used.
- only those in the Markov blanket could be needing an update.

then remove samples  $E \neq e$ . Keep those  $X=y$

> Gibbs Sampling on BN: c)

not d: active trail (Markov Blanket definition on BNs)

> Gibbs BN (ii): F) (children, parents and children's parents.  
marginalize on  $x_{ij}$ , only unknown value)  $\leftarrow ?$

- Variational inference: against problematic samplings.

- 1) Take a simple family of distributions. ( $Q$ )
- 2) Find the closest one to  $P(Y|E=e)$   
 $\forall Y \in Q$
- 3) Use  $Q$  to answer queries

~ = will be defined further.

- Family of distribution: "Normal distribution"  $\rightarrow$  depends on  $\mu, \sigma$ .

- Closure: through 'cross-entropy' distance find the distribution in the family that is the closest to  $P$

Kullback-Liebler distance: positive, not symmetric, no triangular inequality.

$$Q^* = \underset{Q \in \mathcal{Q}}{\operatorname{arg\min}} \text{KL}(Q||P) \quad \begin{matrix} \leftarrow \text{maximize Evidence Lower Bound (ELBO)} \\ \downarrow \text{instead} \end{matrix}$$

- Mon-field algorithm: use  $Q$  that factorizes all variables independently:  $Q(X) = \prod_{i=1}^n Q_i(X_i)$

restrictive  
efficient after all

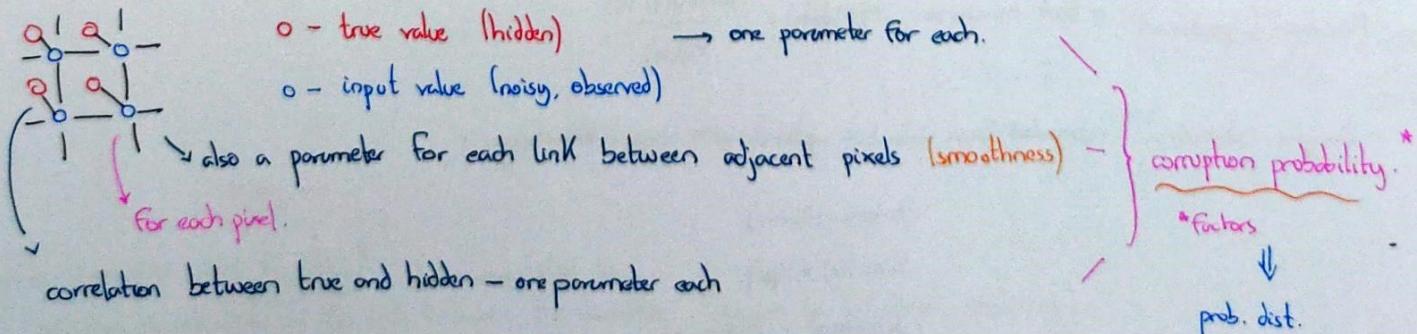
you can update one marginal at a time (given that the rest are fixed)  
without updating everything

} 'coordinate-descent'

- M-projection: takes care of probability mass.

J-projection: covers the mass of the distribution. (don't put prob. mass where dist. is not defined  $\Rightarrow$  high penalization)

## Lecture 9. Image denoising



- energy-diff-if-swap: how will the energy function change if we swap one value.

- swap values until energy function stops improving.

## Lecture 10. Probabilistic Learning

Can we learn from data? Without knowing the structure, the parameters...

- Model learning: we have a dataset and maybe domain experts.

↳ objective: combine them to infer a model  $M$  that approximates  $P^*$

restrict the set of possible models (needed assumption): PGMs.

extract info from dataset  $D$  to encode the structure and the parameters.

structural + parametric learning.

} remember overfitting.

trading performance, complexity, dataset.

- problem of freq. approach: you need infinite samples. (empirical distribution, tends to true distribution)

- Laplace smoothing: to avoid never seen events

$$\frac{N_1 + 1}{N_0 + N_1 + 2}$$

number of classes.

- sufficient statistics to get an estimate  
Coin:  $N_0$  and  $N_1$   
 $N_0$  and  $N$

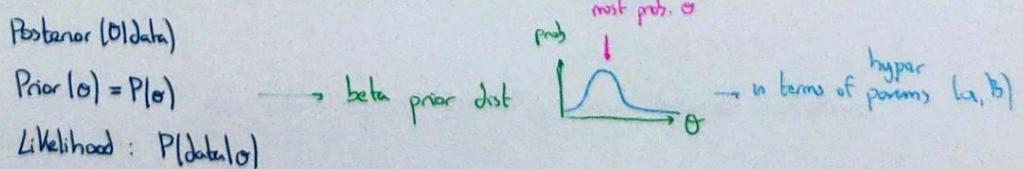
> Sufficient Statistics Exercise: c)

## Lecture 11.

> Parametric learning in MNs and BNs: a) → you can have info of a variable in several factors.  
not gonna be seen.

- Maximum likelihood:  $\theta$  that maximizes  $P(\text{data}|\theta)$
- Maximum a posteriori:  $\theta$  that maximizes  $P(\theta|\text{data}) = \frac{P(\text{data}|\theta)P(\theta)}{P(\text{data})}$

- Bayesian estimation: estimated from data but updated with expert knowledge.



Beta prior is appropriate when our data follows one of these distributions: Bernoulli, Binomial

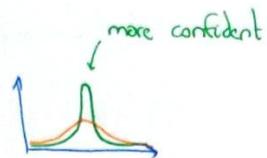
The higher the parameters, the more dense the prob. mass ↗

↳ The posterior is also a beta. (conjugate posterior) → with higher hyperparameters.

- map estimate: formula and def. in slides.

For multinomial dist. it is similar, but with a Dirichlet.

- The hyperparameters represent most prob. value and confidence  
in prior



- non-informative prior: bayesian approach without the frequentist problems.



- Equivalent sample size: both 'weights' sum up to 1 → they are weights (weighting data and prior means)

Posterior mean of Beta distribution  $\xrightarrow{n \rightarrow \infty}$  mean\_data

(Remember Laplace smoothing was doing the weighted sum with the 'ideal' distribution)

Dirichlet MAP estimate with  $B_{0,1} = 2$

## Lecture R. Structural Learning

↳ We want to approximate a prob. dist. from a dataset.

{ Removing edges: less params, better generalization, might not learn real distribution.  
Adding: the opposite.

> Counting params. >

▷ In the Bayesian Network Factorization, only  $P(X_4|X_1, X_6)$  changes to  $P(X_4|X_1, X_2, X_6)$

$$\prod_i P(X_i|PA_i)$$

↓  
parents

remember in this case, binary variables

formula counting params:  $(\Gamma_{X_i} - 1) \cdot \Gamma_{PA_i}$   
 ↓  
 $(2-1) \cdot 2^2$  vs.  $(2-1) \cdot 2^3$   
 parents  
 ↗ parents  
 possible values of  $X_i$   
 possible comb. of parents value

- (Why: we need the model to use it)

we may be interested in the link between the variables (this can be the goal: medicine) -

- Score based structural learning:
  - 1) Define scoring function (penalized) likelihood → refer to slides.
  - 2) Find structure that maximizes it Hill-Climbing, genetic alg.

- scoring function reduces to Mutual Information (if we use likelihood): it represents the strength of the dependence relationship  $X_i \perp\!\!\!\perp PA_i$

Mutual Information is "the part of the variable" that can be explained "with its parents".

- likelihood scoring function:
    - ⊕ pros:
      - decomposes in one term per variable.
    - ⊖ con:
      - MI is always  $\geq 0$ , and mostly  $> 0$
      - adding a new edge always adds score  $\Rightarrow$  best score is totally complete graph
- overfitting!

- > Likelihood score exercise:
- a) T: we only change one direction (no new edges) without involving v-structs.
  - b) F: adding one edge  $\Rightarrow$  increases score.
  - c) T: some reason
  - d) F: we can not know for all datasets (add one, remove one)
    - not directly related to score
    - but it is relevant to know if cond. ind. statements change or not (the num. of)

- avoid overfitting: explicitly control it (limit)
- penalize complexity in the score  $\rightarrow$  penalized log-like likelihood (BIC-Score)
- (Bayesian approach: compute prob. of all possible structures)  $\downarrow$  penalize #edges params.

\* BIC Score: in terms of params.

If  $N$  grows (number of samples)  $\rightarrow$  more emphasize on fitting to data.

- > BIC Score:
- a) T: some parameters and some  $N$  if some  $D$
  - b) F:
  - c) F: there can be a dataset that balances both parts of formula.
  - d) F

- Explore options and check which one maximizes:

[...]  $\rightarrow$  30 mins (find tree or forest algorithm  $\rightarrow$  all variables have  $\leq 1$  parent)

> Recovering directionality. d)

• Local search: moving across neighborhoods, choose best subgraph  $T_G$  (it may not be the global best)

- possible modifications: add, remove or reverse an edge  $\rightarrow$  evaluate all one-action possible modification. (Greedy Hill climbing)  
costly but use score decomposability

> Calculating Likelihood Differences:

remember  $\sum_v \underbrace{MI(X_v | PA_v)}$

c)  $\rightarrow$  identify the terms that are different in all graphs.

> Hidden variable