

* Probabilistic Graphical Models (Koller - Staufend)

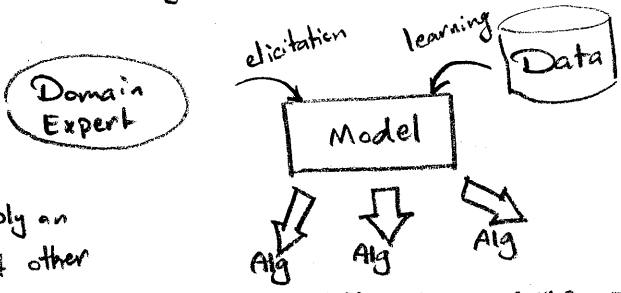
MODEL
PROBABILISTIC
GRAPHICAL

- First: Medical Diagnosis → predisposing factors, symptoms, test results, diseases, treatment outcome
- Image Segmentation → pixels ⇒ superpixels ⇒ label

- In common: Large # of Factors
- Huge amount of uncertainty about results

• Model: Declarative Representation of our understanding from the world

- what variables are
- How vars interact with each other
- Declarative → standalone
 - ↳ no matter what algorithm apply on
 - ↳ we can understand it without other knowledge
- Algs can ask different questions or the same question with different accuracy-complexity trade off



• Uncertainty

- Partial knowledge of state of the world
- Noisy observations
- Phenomena not covered by our own model → modeling limitation
- Inherent stochastic → micro level and sometimes even macro level

• Probability Theory

- Declarative representation with clear semantics
- Powerful reasoning patterns → Toolbox of conditioning, reasoning, decision making
- Establish learning methods

• Complex Systems

- Random variables → X_1, \dots, X_n → simplest case: all binary
- Distribution over R.V. → simplest case 2^n states → intractable
- Goal: to capture the uncertainty in form of Joint Prob → to exploit the structure in distribution and represent and manipulate in effective way

• Graphical Model

- Bayesian Networks → directed graph
 - ↳ nodes: random variables
 - ↳ edges: probabilistic connection between R.V.
- ✓ CPDs
- Markov Network → undirected graph
 - ✓ Image segmentation

• Graphical Representation

- Intuitive & Compact data structure for capturing high dimensional distribution
- Efficient reasoning using general-purpose algorithms → exploiting graphical structure
- Sparse parameterization
 - ↳ feasible elicitation ← by hand
 - ↳ learning from data ← automatic
 - ⇒ Reduction in the number of Params.

Image Segmentation

Pixels → Superpixels → Labeled : Machine Learning

Many Applications
Lecture 17

- Medical Care
- Textual Information Extraction
- Multi-Sensor Integration

Sensors → Learned Models → Current & Future Data

Biologic Network Reconstruction

• Roadplan

- Representation
 - ↳ Directed and Undirected
 - ↳ Temporal and Plate model → high level & complex scenarios
- Inference (= Reasoning)
 - ↳ Exact & Approximate
 - ↳ Decision Making under Uncertainty
- Learning
 - ↳ Parameters and Structure
 - ↳ With and without complete data

* Joint Distribution

Student Example

- I: Intelligence → i^0 : low, i^1 : high
- D: Difficulty → d^0 : easy, d^1 : hard
- G: Grade → g^1 : A, g^2 : B, g^3 : C

$2 \times 2 \times 3 = 12$ possible combinations
↳ free parameters

• Independent variables: vars whose values are not determined by the values of other parameters

↳ here: $\sum P_i = 1 \Rightarrow$ if we know 11 of the values, the 12th is obvious → 11

			$P(I, D, G)$
I	D	G	Prob
i^0	d^0	g^1	0.126
i^0	d^0	g^2	0.168
i^0	d^0	g^3	0.126
i^0	d^1	g^1	0.009
i^0	d^1	g^2	0.045
i^0	d^1	g^3	0.126
i^1	d^0	g^1	0.252
i^1	d^0	g^2	0.0224
i^1	d^0	g^3	0.0056
i^1	d^1	g^1	0.06
i^1	d^1	g^2	0.036
i^1	d^1	g^3	0.024
Sum			1

• Conditioning: assume that student got an A

$A \sim g^1 \rightarrow$ eliminates all rows of table that has $G \neq g^1 \rightarrow$ reduced prob. distribution
 \hookrightarrow not consistent with the observation

This operation is called reduction: reduce away stuff that is not consistent with observation

\hookrightarrow doesn't give a probability distribution

$$\sum P_i(I, D, g^1) \neq 1$$

$\hookrightarrow P(I, D, g^1)$ is unnormalized measure

\hookrightarrow Renormalization needed.

\hookrightarrow divide each of them by 0.447

$$\hookrightarrow P(I, D | g^1)$$

I	D	G	$P(I, D, G)$
i^0	d^0	g^1	0.126
i^0	d^1	g^1	0.009
i^1	d^0	g^1	0.252
i^1	d^1	g^1	0.06
			0.447

\Downarrow divide by 0.447

I	D	$P(I, D g^1)$
i^0	d^0	0.282
i^0	d^1	0.02
i^1	d^0	0.564
i^1	d^1	0.134
		1

• Marginalization: operation that takes the probability dist of a larger subset of variables and produces a prob. dist over a subset of these

to throw away I from $P(I, D)$ we should sum up over I.

\Downarrow sum up over I

D	$P(D)$
d^0	0.846
d^1	0.156
1	

* Factors

• A function or a table

input: random variable \rightarrow output: related to assignments of r.v. a real value

\hookrightarrow accepts all possible assignment in the cross product space of r.v.

\hookrightarrow all possible combinations

$$\phi : \text{val}(X_1, \dots, X_k) \rightarrow \mathbb{R}$$

$$\hookrightarrow \text{Scope}(\phi) = \{X_1, \dots, X_k\}$$

• A joint distribution is a factor. \rightarrow but factors doesn't require $\sum p_i = 1$

an unnormalized measure is a factor. $\rightarrow P(I, D, g^1) \sim \text{scope}(\phi) = \{I, D\}$

conditional probability distribution

• CPD: Conditional Probability Distribution

- a factor

□ $P(G | I, D)$

* for every combination of values of I & D we have a probability dist. over G.

* each row sums to 1 → distribution over G
 ↳ context

	g^1	g^2	g^3
i^0, d^0	0.3	0.4	0.3
i^0, d^1	0.05	0.25	0.7
i^1, d^0	0.9	0.08	0.02
i^1, d^1	0.5	0.3	0.2

• General Factor

□ Scope = {A, B}

A	B	ϕ
a^0	b^0	30
a^0	b^1	5
a^1	b^0	1
a^1	b^1	10

• Operations on Factors

- Factor Product

$\phi_1(A, B) \times \phi_2(B, C)$

scope = {A, B} \cap {B, C} = {A, B, C}

$\phi_3(A, B, C)$

A	B	$\phi_1(A, B)$	B	C	$\phi_2(B, C)$
a^1	b^1	0.5	b^1	c^1	0.5
a^1	b^2	0.8	b^1	c^2	0.7
a^2	b^1	0.1	b^2	c^1	0.1

$T_1(\text{row } 1) \times T_2(\text{row } 1) \leftarrow a^1, b^1 \& b^1, c^1$
 $T_1(\text{row } 1) \times T_2(\text{row } 2) \leftarrow a^1, b^1 \& b^1, c^2$

A	B	C	$\phi_3(A, B, C)$
a^1	b^1	c^1	$0.5 \times 0.5 = 0.25$
a^1	b^1	c^2	$0.5 \times 0.7 = 0.35$
a^1	b^2	c^1	$0.8 \times 0.1 = 0.08$

- Factor Marginalization

$\phi_3(A, B, C) \rightarrow \phi_4(A, C)$

find different values of B and sum them up

A	C	$\phi_4(A, C)$
a^1	c^1	$0.25 + 0.08 = 0.33$

- Factor Reduction

focus on context c^1
 look for rows with values $C=c^1$

$\phi_3(A, B, C) \rightarrow \phi_5(A, B)$

A	B	C	$\phi_5(A, B)$
a^1	b^1	c^1	0.75
a^1	b^2	c^1	0.08
a^2	b^1	c^1	0.05

• why factors?

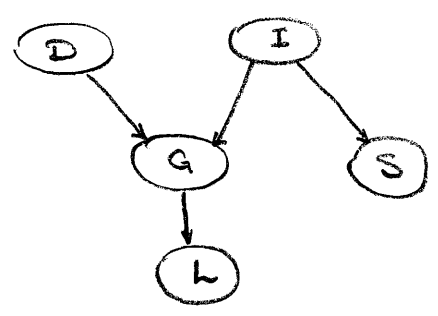
- Fundamental building blocks for defining distributions in high-dimensional spaces
- set of basic operations for manipulating these probability distributions

* Bayesian Networks

□ The student example

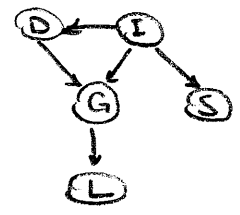
- Grade
- Difficulty of Course
- Intelligence of Student
- SAT
- Letter of Reference

$$P(G, D, I, S, L)$$



this is not the only model possible.

for example smarter student may go for harder classes.



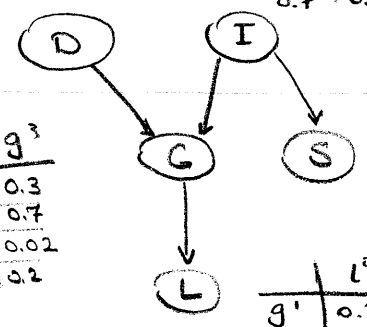
• Each node have a CPD

- $P(D)$
- $P(I)$
- $P(G|I, D)$
- $P(S|I)$
- $P(L|G)$

d^0	d^1
0.6	0.4

i^0	i^1
0.7	0.3

	g^1	g^2	g^3
i^0, d^0	0.3	0.4	0.3
i^0, d^1	0.05	0.25	0.7
i^1, d^0	0.9	0.08	0.02
i^1, d^1	0.5	0.3	0.2



	s^0	s^1
i^0	0.95	0.05
i^1	0.2	0.8

	l^0	l^1
g^1	0.1	0.9
g^2	0.4	0.6
g^3	0.99	0.01

• Chain Rule for BN

simply take all different CPDs in BN and multiply the together.

$$P(G, I, D, S, L) = P(D) \cdot P(I) \cdot P(G|I, D) \cdot P(S|I) \cdot P(L|G)$$

five factors with overlapping scopes \implies big factor with scope of all variables

$$\begin{aligned} \square P(d^0, i^1, g^3, s^1, l^1) &= P(d^0) \times P(i^1) \times P(g^3|d^0, i^1) \times P(s^1|i^1) \times P(l^1|g^3) \\ &= 0.6 \times 0.3 \times 0.02 \times 0.01 \times 0.8 \end{aligned}$$

• Definition:

- BN is a directed acyclic graph G whose nodes represent the random variables
 \hookrightarrow DAG $\hookrightarrow X_1, \dots, X_n$

- For each node X_i there is a CPD $P(X_i | \text{Pa}_G(X_i))$
 \hookrightarrow parents in graph G

- BN represents a joint dist. via the chain rule for BN

$$P(X_1, \dots, X_n) = \prod_i P(X_i | \text{Pa}_G(X_i))$$

• Proof of BN is a joint distribution

1- BN is a legal distribution: $P \geq 0$

P is multiply of CPDs $\Rightarrow P \geq 0$ ✓
 CPDs are non-negative

2-BN is legal distribution : $\sum P = 1$

$$\sum_{D,I,G,S,L} P(D,I,G,S,L) = \sum_{D,I,G,S,L} \underbrace{P(D) P(I) P(G|I,D) P(S|I) P(L|G)}_{\text{chain rule}}$$

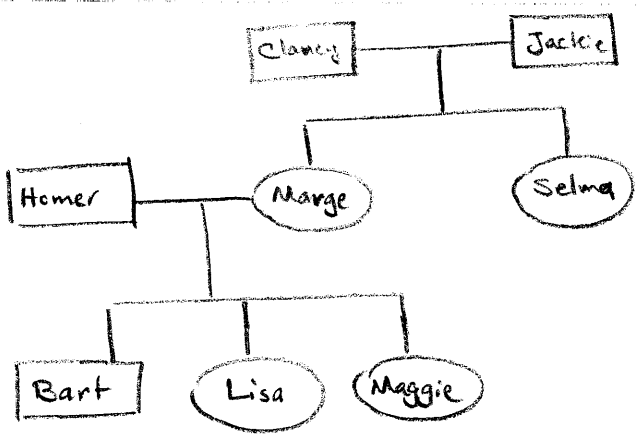
$$= \sum_{D,I,G,S} P(D) P(I) P(G|I,D) P(S|I) \underbrace{\sum_L P(L|G)}_{1} \quad \text{push summation inside}$$

$$= \sum_{D,I,G} P(D) P(I) P(G|I,D) \underbrace{\sum_S P(S|I)}_{1} = \dots = 1 \quad \checkmark$$

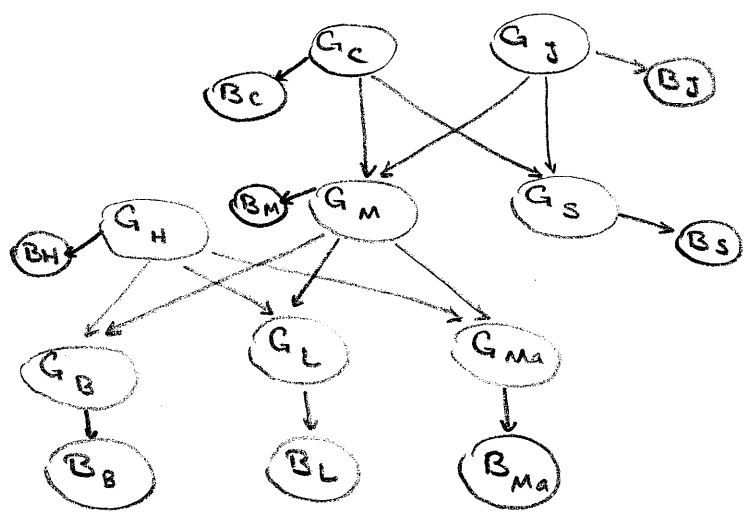
• P Factorizes over G

- We can represent it over graph G if we can encode it using the chain rule for BN
- G: Graph over $X_1, \dots, X_n \iff P$ factorizes over G if $P(X_1, \dots, X_n) = \prod_i P(X_i | \text{Par}_G(X_i))$

✓ The first BN ever! \Rightarrow Genetic Inheritance



Chromosomes : A, B, O
 \downarrow each person have two chromosomes
 Genotype : AA, AB, AO, BO, BB, OO
 \downarrow AA and AO show the same effect
 Phenotype : A, B, AB, O



* Reasoning Patterns

✓ Student Network

• Create graph → Joint dist. Using Chain Rule → Marginalization:

✓ $P(i') \approx 0.5$

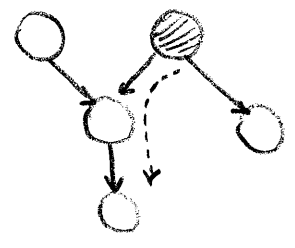
• Condition on the Parents of Query: Causal Reasoning

Joint Dist → Reduction on Evidence → Marginalization

✓ $P(i' | i^0) \approx 0.39$

↳ causal reasoning: reason from top to down

$P(i' | i^0, d^0) \approx 0.51$

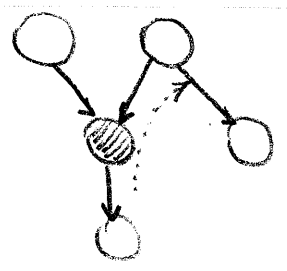


• Condition on the children of Query: Evidential Reasoning

✓ $P(d') = 0.4 \rightarrow P(d' | g^3) \approx 0.63$

$P(i') = 0.3 \rightarrow P(i' | g^3) \approx 0.08$

if the student get C there is more prob. that course is harder or student is less intelligent.



• Intercausal Reasoning: explain away effect

✓ class is hard, student get a C (or a B)
 $\downarrow_{g^3} \quad \downarrow_{g^2}$

$P(i') = 0.3$

$P(i' | g^3) \approx 0.08$

↳ increase again but not too much because

$P(i' | g^3, d') \approx 0.11$

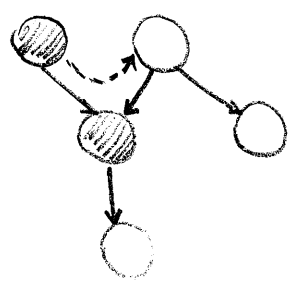
$P(g^3 | d')$ is still down

$P(i') = 0.3$

$P(i' | g^2) \approx 0.175$

↳ increase even more than $P(i')$

$P(i' | g^2, d') \approx 0.34$



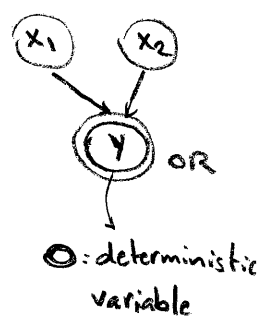
- explanation

✓ Value of Y is completely determined with X_1 and X_2 so other rows have 0 probability. we observe $Y=1$:

$P(X_1=1) = \frac{2}{3} \quad P(X_2=1) = \frac{2}{3}$

$P(X_1=1 | X_2=1) = 0.5$

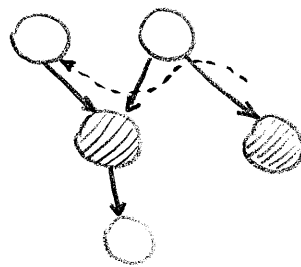
	X_1	X_2	Y	Prob
$Y=1 \leftarrow$	0	0	0	0.25
	0	1	1	0.25
	1	0	1	0.25
	1	1	1	0.25



If I know $Y=1$ the chance of X_1 and X_2 is 66%

If I further know $X_2=1$ so that explains why $Y=1$ and it doesn't matter if X_1 is equal to 1 or not. It is explained away.

Student example: student aced the SAT but get C grade.
 so we may say that the class is difficult, because SAT explained away that student is not intelligent.



$P(d') = 0.4$

$P(i') = 0.3$

$P(d' | g^3) \approx 0.63$

$P(i' | g^3) \approx 0.08$

$P(d' | g^3, s') \approx 0.73$

$P(i' | g^3, s') \approx 0.58$

increase dramatically

* Flow of Probabilistic Influence

The question is when can X influence Y? →

Influence: condition on X changes belief on Y

- $X \rightarrow Y$ ✓: connected, X is parent of Y

- $X \leftarrow Y$ ✓: connected, Y is child of X

- $X \rightarrow W \rightarrow Y$ ✓: causal chain

- $X \leftarrow W \leftarrow Y$ ✓: evidential chain

- $X \leftarrow W \rightarrow Y$ ✓

- $X \rightarrow W \leftarrow Y$ X: v-structure

In general probabilistic influence is symmetrical. If X can influence Y then Y can influence X.

□ If I say the student took a difficult class, does that tell me anything about the student intelligence?

• Active trail with no evidence

- Trail: a sequence of nodes that are connected to each other by single edges in graph
 ↳ undirected

- Influence can flow from one variable to another except in v-structure → this flow can continue

- A trail $X_1 - \dots - X_n$ is active if it has no block
 ↳ v-structure $X_{i-1} \rightarrow X_i \leftarrow X_{i+1}$

• The question is when can X influence Y given evidence about Z?

- In general if we have an edge $X \rightarrow Y$ in some BN, there is there ^{any} set of other variables that we can condition on to make X and Y independent? NO!

- $X \rightarrow Y$ ✓

- $X \leftarrow Y$ ✓

- $X \rightarrow W \rightarrow Y$

- $X \leftarrow W \leftarrow Y$

- $X \leftarrow W \rightarrow Y$

- $X \rightarrow W \leftarrow Y$

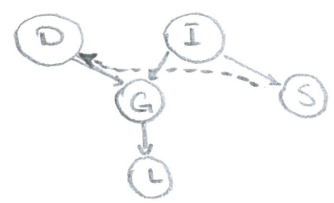
$W \in Z$	$W \notin Z$	
✓	X	I know W, so Y only depend on W
✓	X	symmetry
✓	X	the root is already known, so observation can't change that
X	✓	→ explain away, but this row holds only if W and all of its descendants not in Z

☑ S-I-D-G allow influence

I observed X

I not observed, nothing else X

I not observed, G observed ✓



• Active trail $X_1 \dots X_k$ given Z

- the v-structure is activated \rightarrow node or one of descendants is observed $\hookrightarrow e \in Z$
 $\hookrightarrow X_{i-1} \rightarrow X_i \leftarrow X_{i+1}$
 $\hookrightarrow X_i$

- no other X_i in Z
 \hookrightarrow not in the middle of v-structure $X_{i-1} \rightarrow X_i \leftarrow X_{i+1}$

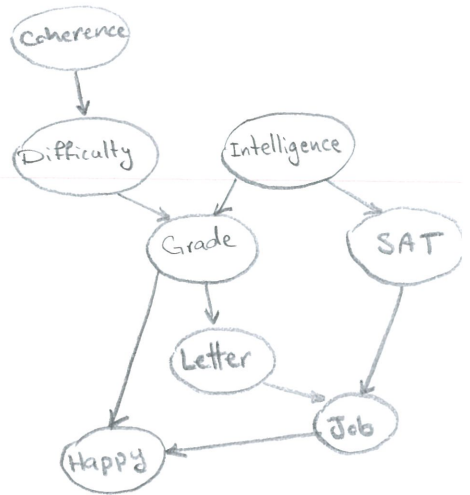
☑ Active if we observe G

$C \rightarrow D \rightarrow G \leftarrow I \leftarrow S$ ✓

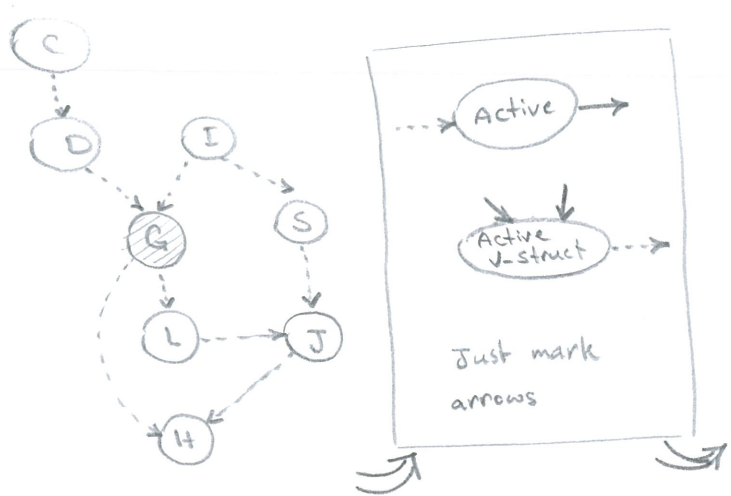
$I \rightarrow G \rightarrow L \rightarrow J \rightarrow H$ X

$I \rightarrow S \rightarrow J \rightarrow H$ ✓

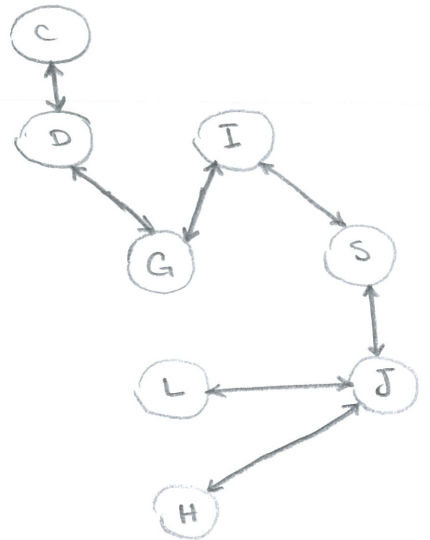
$C \rightarrow D \rightarrow G \leftarrow I \rightarrow S \rightarrow J \leftarrow L$ X



* my idea: re-drawing the graph



bad idea, have problems with last example



* Independencies

• so far: graphs to encode prob. dist.

complementary view: representation of a set of independencies that a dist. must satisfy

• Independence in probability distribution

- for events α and β , $P \models \alpha \perp \beta$ (read P satisfies α is independent of β) if:

$\hookrightarrow P(\alpha, \beta) = P(\alpha) \times P(\beta)$

\rightarrow conjunction

$\hookrightarrow P(\alpha | \beta) = P(\alpha)$

\rightarrow flow of information, β doesn't effect prob in α

$1 \ P(\alpha | \beta) = P(\alpha)$

\rightarrow information flow is symmetrical

- for random variables X, Y , $P \models X \perp Y$ if:

$\hookrightarrow P(X, Y) = P(X) P(Y)$

$\hookrightarrow P(X|Y) = P(X)$

$\hookrightarrow P(Y|X) = P(Y)$

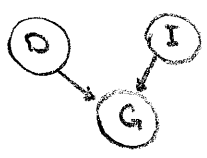
universal statement: $\forall x, y \ P(x, y) = P(x) \times P(y)$
factor interpretation: factor $\Phi_3(x, y)$ is the product of two lower dimension factors $\Phi_1(x)$ and $\Phi_2(y)$

□

$P(I, D) = P(I) \cdot P(D)$



$\sum_G P(G|I, D)$



• Conditional Independence

- for (set of) random variables X, Y, Z , $P \models (X \perp Y | Z)$ if:

$\hookrightarrow P(X, Y | Z) = P(X | Z) P(Y | Z)$

$\hookrightarrow P$ satisfies X is independent of Y given Z

$\hookrightarrow P(X | Y, Z) = P(X | Z)$ \rightarrow given Z , Y gives me no additional information about X

$\hookrightarrow P(Y | X, Z) = P(Y | Z)$

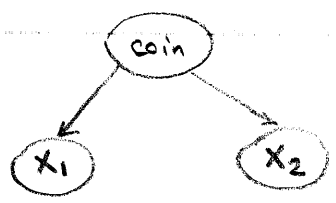
$\hookrightarrow P(X, Y, Z) \propto \Phi_1(X, Z) \Phi_2(Y, Z)$

□ two coins, one is fair, one had 90% of heads
we pick a coin and toss it twice

If the first toss is head \rightarrow the prob of 2nd toss to be H rises!

If we are told that we pick the fair coin \rightarrow no info about 2nd toss

If we are told that we pick the biased coin \rightarrow again no info about 2nd toss!
anymore!



so: X_1 and X_2 are not independent: $P \not\models X_1 \perp X_2$

but $P \models X_1 \perp X_2 | C$

□

$P(S, G | I^0) = P(S | I^0) \times P(G | I^0)$

do we still need to check $P(S, G | I^1) = P(S | I^1) \times P(G | I^1)$ to say $P \models S \perp G | I$?

YES! we need to check for all values of I .

• Conditioning can lose independence!

□

$P \models I \perp D$

$P \not\models I \perp D | G^1$

* Independence in BN

• why independence and factorization are related to each other?

- $P(X, Y) = P(X) P(Y)$ \rightarrow definition of independence
 \searrow factorization over two factors

$P(X, Y, Z) \propto \phi_1(X, Z) \phi_2(Y, Z) \iff (X \perp Y | Z)$

- factorization of a distribution P implies independencies that hold in P
- if P factorizes over G , can we read these independencies from the structure of G ?

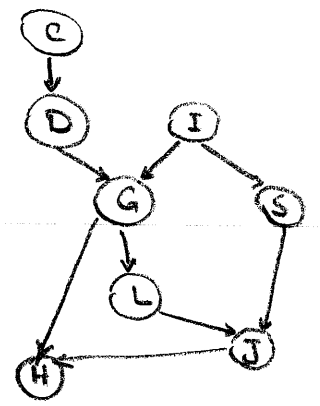
• Flow of Influence and d-Separation

- what happens when there is no active trails in graph?
 - ↳ influence can't flow → notion of d-separation

- X and Y are d-separated in G given Z if there is no active trail in G between X and Y given Z

↳ $d\text{-sep}_G(X, Y | Z)$

- ☑ $d\text{-sep}(D, I | L) \rightarrow v\text{-structure in } G \text{ active}$
- $d\text{-sep}(D, J | L) \rightarrow v\text{-structure in } G \text{ active}$
- $d\text{-sep}(D, J | L, I) \checkmark$
- $d\text{-sep}(D, J | L, H, I) \rightarrow v\text{-structure in } H \text{ active}$

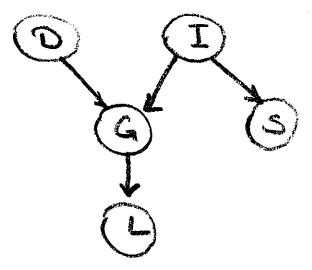


- Proof of (BN: Factorization \implies independence)

Theorem: If P factorizes over G , and $d\text{-sep}_G(X, Y | Z)$ then P satisfies $(X \perp Y | Z)$

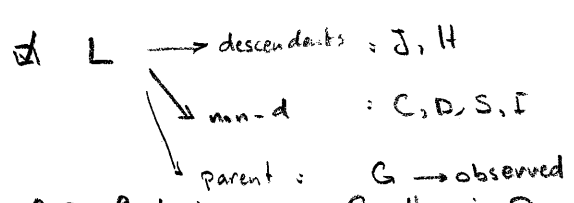
☑ $D \perp S \rightarrow$ no active trail from X to Y given Z

$P(D, I, G, S, L) = \underbrace{P(D) P(I) P(G|D, I) P(S|I) P(L|G)}_{\text{chain rule for BN}}$



$$\begin{aligned}
 P(D, S) &= \sum_{G, L, I} P(D) P(I) P(G|D, I) P(S|I) P(L|G) \\
 &= \sum_I P(D) P(I) P(S|I) \underbrace{\left(\sum_G (P(G|D, I) \sum_L P(L|G)) \right)}_1 \\
 &= P(D) \underbrace{\left(\sum_I \frac{P(I) P(S|I)}{P(S, I)} \right)}_{P(S)} = \underbrace{P(D)}_{\phi_1(D)} \underbrace{P(S)}_{\phi_2(S)} \implies P \models D \perp S \checkmark
 \end{aligned}$$

- Any node is d-separated from its non-descendants given its parents



from S to J → trail blocked
 from S to H →

If P factorizes over G , then in P , any variable is independent of its non-descendants given its parents

• I-maps \rightarrow Independency Map

d-separation \Rightarrow P satisfies corresponding independence statement

$$I(G) = \{ (X \perp Y | Z) : d\text{-sep}_G(X, Y | Z) \}$$

all of the independence statements

correspond to d-separation statements in graph

The set of all independencies that are derived from these separation statements, and those are the ones that graph implies hold for P.

- If P satisfies $I(G)$ we say that G is an I-map of P

□ P_1 :

I	D	Prob
i^0	d^0	0.42
i^0	d^1	0.18
i^1	d^0	0.28
i^1	d^1	0.12

P_2 :

I	D	Prob
i^0	d^0	0.282
i^0	d^1	0.02
i^1	d^0	0.564
i^1	d^1	0.134



$I(G_1) = \{ D \perp I \}$

$I(G_2) = \{ \}$

$P_1 \models D \perp I \Rightarrow G_1 = \text{I-map}(P_1) \checkmark$

$P_2 \Rightarrow G_2 = \text{I-map}(P_2) \checkmark$

$P_2 \not\models D \perp I \Rightarrow G_1 \neq \text{I-map}(P_2) \times$

G_2 has no independence assumption

\hookrightarrow both P_1 and P_2 satisfy the empty set

$\hookrightarrow P_1$ also satisfies additional indep. assumption but I-map don't require that the graph exactly represents the independencies

$P_1 \Rightarrow G_2 = \text{I-map}(P_1) \checkmark$

- Proof of (BN: Factorization \Rightarrow Independence)

Theorem: If P factorizes over G, then G is an I-map for P

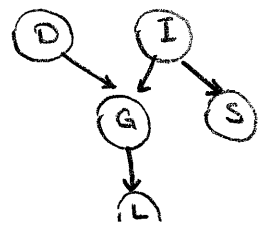
= I can read from G, independencies in P regardless of parameters

Conversely: (BN: Independence \Rightarrow Factorization)

Theorem: If G is an I-map for P, then P factorizes over G

= If the dist. satisfies the independence assumption that are implicit in graph, we can take that dist and represent it as a BN over the graph

□ $P(D, I, G, S, L) = P(D)P(I|D)P(G|D, I)P(S|D, I, G)P(L|D, I, G, S)$
 chain rule for probability
 L ... still don't know if it is BN



$P(I|D) \xrightarrow{I \perp D} P(I)$

$P(S|D, I, G) \xrightarrow{\text{last theorem}} P(S|I)$

$P(L|D, I, G, S) \xrightarrow{\text{last theorem}} P(L|I, G)$

$\Rightarrow P(D, I, G, S, L) = P(D) P(I) P(G|D, I) P(S|I) P(L|I, G)$

last theorem: remove non-descendants

- Two equivalent view of graph structure:

↳ Factorization: G allow P to be represented

↳ a data structure that show us how dist can be broken up to pieces (factors=CPDs)

↳ I-map: Independencies encoded by G hold in P

↳ structure in G encode independencies that doesn't necessarily hold in P.

- If P factorizes over a graph G we can read from the graph independencies that must hold in P (and I-map)

- Independency is valuable → tells sth about what influence what → structure

↳ if I observe sth what else might change → observation

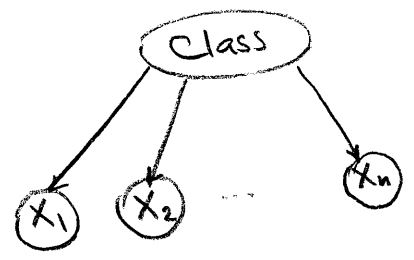
* Naïve Bayes → idiot bayes

• Reason of naming: because the model make independence assumptions that indeed are very naïve and overly simplistic.

• This model is typically used for classification

- Given: features

- Infer: class



- If we observe the class variable, influence cannot flow between any of the other variables, because the only path between two other variables is through the class variable. If we don't observe the class variable, however, all the variables are dependent.

$(X_i \perp X_j | C)$ for all X_i, X_j

$P(C, X_1, \dots, X_n) = P(C) \prod_{i=1}^n P(X_i | C)$

• Better understanding: looking at ratio of probabilities of two different classes given the particular observation

$$\frac{P(C = c_1 | x_1, \dots, x_n)}{P(C = c_2 | x_1, \dots, x_n)} = \frac{P(C = c_1)}{P(C = c_2)} \times \prod_{i=1}^n \frac{P(x_i | C = c_1)}{P(x_i | C = c_2)}$$

odds ratios

Text classification & Categorization

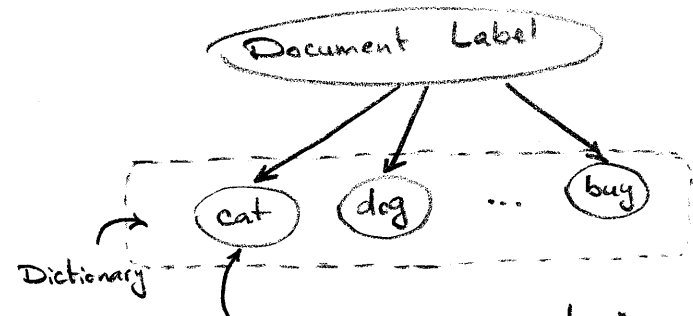
1. Bernoulli: Naive Bayes For Text

- It treats every word in the dictionary as r.v.
- For each word we have a binary variable
 - ↳ 1 : word appears in a document
 - ↳ 0 : otherwise

5000 words in dict \Rightarrow 5000 binary var.

- CPD associated to each r.v. \longrightarrow

	cat	dog	buy	sell
financial	0.001	0.001	0.2	0.3
pets	0.3	0.4	0.02	0.0001



$P(\text{"cat" appears in doc} \mid \text{label of doc})$

Bernoulli \rightarrow binary

Naive Bayes \rightarrow strong independence assumption

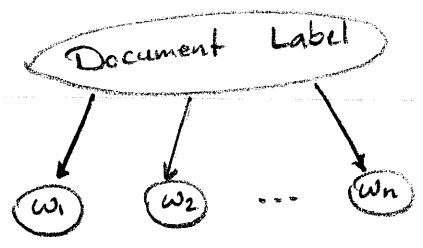
2. Multinomial Naive Bayes For Text

- Features are words in the document
- n: length of document \rightarrow number of r.v.
- values of r.v. are actual words \rightarrow if there are D words in dictionary cardinality of each r.v. is D
- assuming that prob. dist over the words is the same for every position

- CPD associated to all random variables

	cat	dog	buy	sell
financial	0.001	0.001	0.2	0.3
pet	0.3	0.4	0.02	0.0001

\rightarrow multinomial } all of the entries sum to one.



- Naive \rightarrow strong independence assumption \rightarrow word in position i is independent of word in pos j
 \rightarrow don't support two word phrases

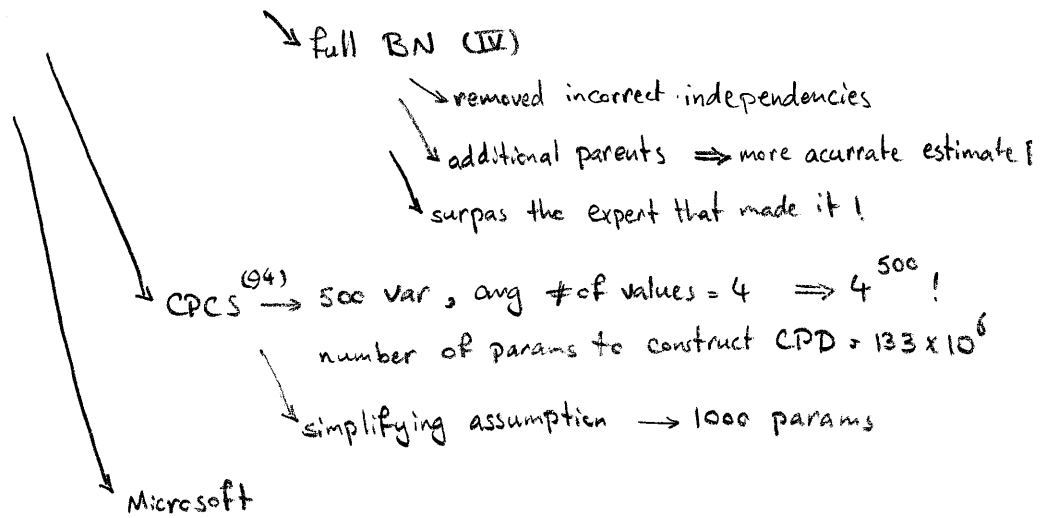
• Simple Approach \rightarrow Computationally efficient
 \rightarrow Easy to construct

- Surprisingly efficient in domains with many weakly relevant features
- Strong independence assumption reduce performance when many features are strongly correlated

Medical Diagnosis \rightarrow Path finder (92) \rightarrow rule based system (I)
 \rightarrow naive Bayes (II)
 \rightarrow naive Bayes + knowledge engineering (III)

* the product of any number and 0 is 0, once a zero prob. is (wrongly) assigned to an event that is actually possible in reality, nothing can change it again. It is difficult to be sure which event had real 0 prob. and which have really small prob. Wrongly assign 0 is dangerous, but small value for impossible events is safer.

* better calibration of conditional prob. $\rightarrow P(\text{finding}_1 | \text{disease}_1)$ to $P(\text{finding}_1 | \text{disease}_2)$ / $P(\text{finding}_2 | \text{disease}_1)$ to $P(\text{finding}_2 | \text{disease}_2)$ X P.15



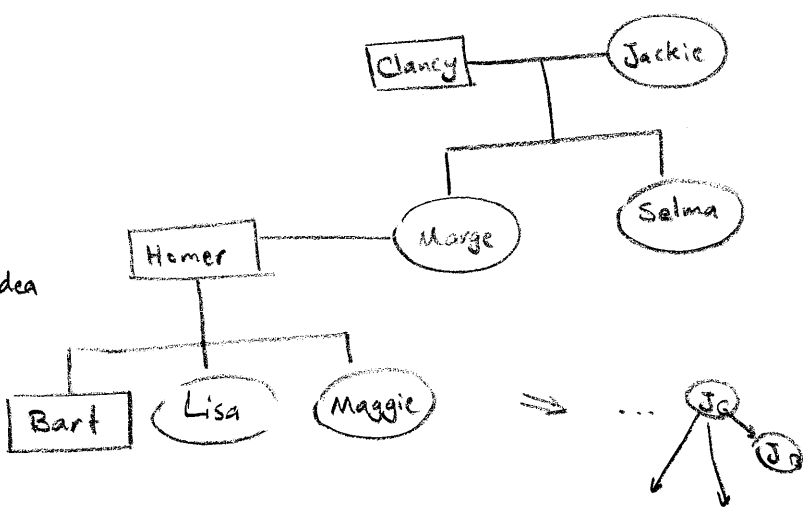
✓ Fault Diagnosis \rightarrow Microsoft troubleshooters
 \rightarrow Car Repair

Benefits: flexible GUI, easy to design and maintain
 \downarrow skip question... \downarrow change prob., add an edge...

* Template Models

Genetic Inheritance

- Input: Pedigree
- Output: reasoning about family trait
- If we add nodes to family tree, we want have new network the same idea about Genotypes and Phenotypes



↳ Sharing between Models

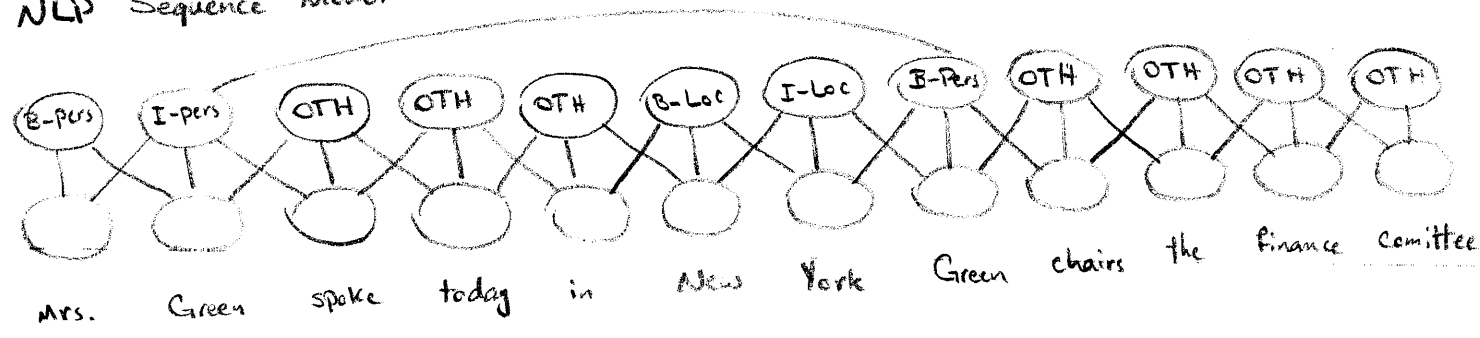
- This model is too redundant

↳ Sharing within model

* CPD that tell how Selma G effects Selma P is the same with other

* The structure of each person has two parents repeated in model

NLP Sequence Model



- task: Named entity recognition

- shared: relation between observation & latent var that are indep. from the place in the sequence.

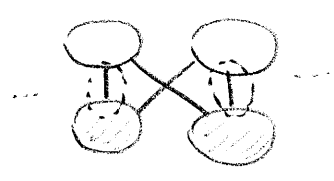
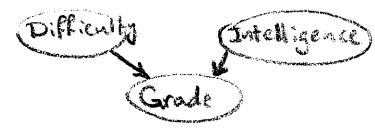
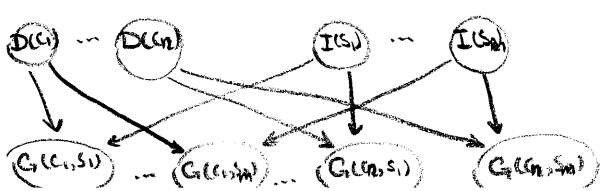


Image Segmentation

- no separate model for every superpixel in image
- sharing across superpixels
- the model that relates class label of superpixel to image features
- parameters that involves adjacent super pixels are shared.
- pairs of super pixels shared
- for other images it works → between model

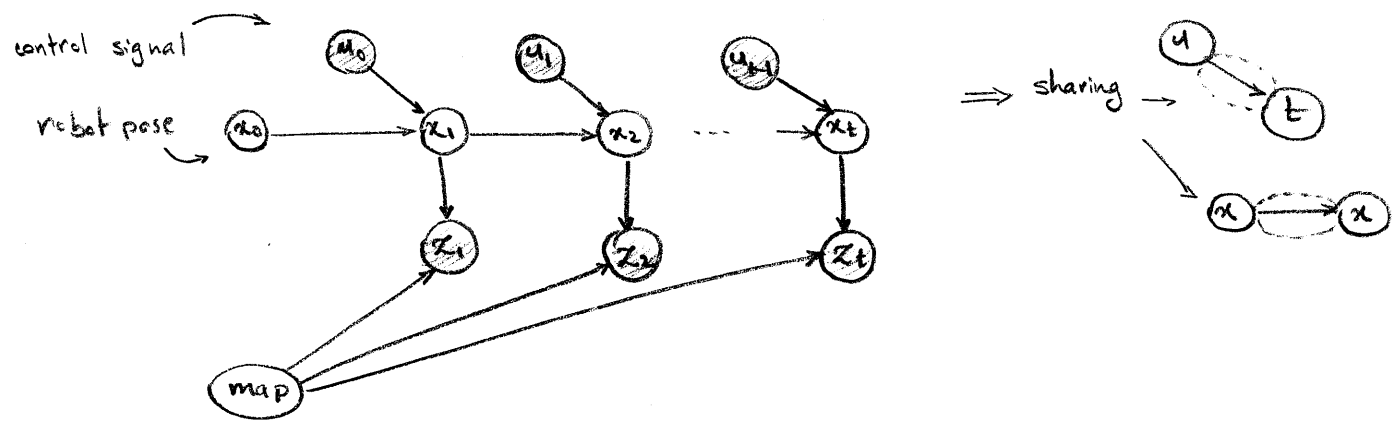
Student Example For whole university



- the grade of a student within a course, depends on the relevant difficulty and intelligence of relevant student. → sharing of both structure and parameters

Robot localization

- time series
- position at time t changes over time, we expect the dynamics of the robot to be fixed.



• Template variable: something that we are replicating again and again within single model as well as across model.

- template variable → like a function that take arguments $X(u_1, \dots, u_k)$
↳ instantiated / duplicated multiple times

- time points → Location (t), Sonar (t)
- person → Genotype (person), Phenotype (person)
- pixel → Label (pixel)

Difficulty (course), Intelligence (student), Grade (course, student)

• Template model is a language that tells us the dependency model of template variables and how concrete instantiation of variables (= ground variable) inherit dependency from template.

- Benefits → CPDs in TMs can often be copied many times.
- ↳ TMs can often capture events that occur in time series.
- ↳ TMs can capture parameter sharing within the model.

Different languages developed for TMs

- Dynamic BN → for temporal processes (replication over time)
- Object-relational Models → Directed → plate models, ... → for multiple objects
↳ Undirected

* Temporal Models

- For systems that evolve over time → distribution over trajectories
- When representing over continuous time → most cases forget the time is continuous → discretize
 - pick time granularity Δ → usually dictated by sensors
 - set of template r.v → $X^{(t)}$: variable X at time $t \times \Delta$

- $X^{(t:t')}$ = $\{X^{(t)}, \dots, X^{(t')}\}$ $t \leq t'$

- want to represent $P(X^{(t:t')})$ for any t, t'

• Markov Assumption

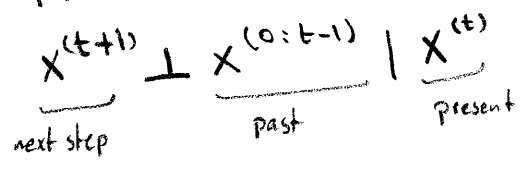
- A type of conditional independence assumption \rightarrow helps to compactify time course data

\rightarrow state @ $t+1$ given state at previous times

$$P(X^{(0:T)}) = P(X^{(0)}) \prod_{t=0}^{T-1} P(X^{(t+1)} | X^{(0:t)})$$

chain rule for probabilities

- Assumption: Future is independent of the past given the present



\rightarrow a forgetting assumption \rightarrow don't care anymore about past

$$P(X^{(0:T)}) = P(X^{(0)}) \prod_{t=0}^{T-1} P(X^{(t+1)} | X^{(t)})$$

\rightarrow Is this assumption true?

\checkmark X = location of robot $L^{t+1} \perp L^{t-1} | L^t$

in most cases does not hold \rightarrow ignores velocity (direction + speed)

to fix it \rightarrow enrich the state description by adding V^t, \dots

\rightarrow make approximation better

\rightarrow Markov assumption become much more warranted

\rightarrow move away Markov assumption

\rightarrow adding dependencies that go further back in time

\rightarrow Semi-Markov assumption

• Time Invariance Assumption

- we still need a probabilistic model for every P .

\rightarrow we can use template probability model

- $P(X' | X) = P(\text{next step} | \text{current time point}) \rightarrow$ we assume that model is replicated for every single time point

\rightarrow Time Invariance

$$P(X^{(t+1)} | X^{(t)}) = P(X' | X)$$

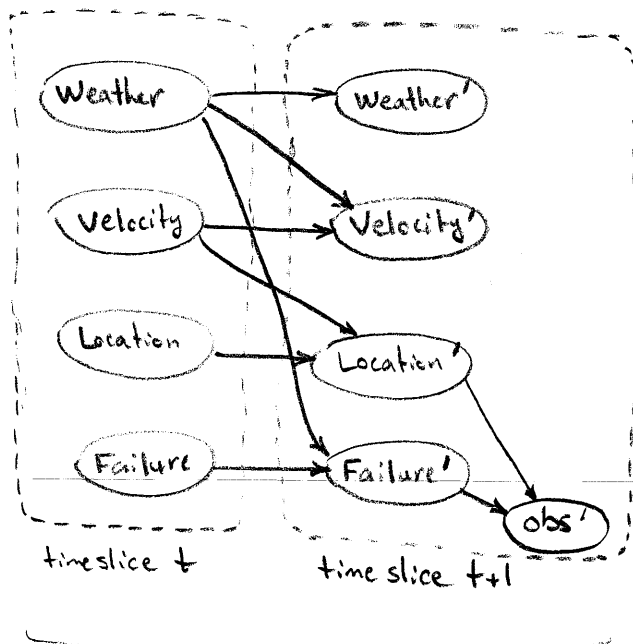
- we assume that the dynamics of the system don't depend on current time point t .

\rightarrow Is this assumption true?

traffic on road → does the dynamics of that traffic depends on time of system?

in most cases does not hold → depends on time of day, day of week, big football match
to correct the inaccuracies → enriching the model by including these variable

template transition model



$P(W', V', L', F', O' | W, V, L, F) =$
 $P(W'|W) \times P(V'|V, W) \times P(L'|L, V) \times P(F'|F, W) P(O'|L')$

dependency across time $t \rightarrow t+1$ within $t: t+1 \rightarrow t+1$

rapidly-acting before the next time step

- * intra time-slice edge
- * inter time-slice edge

net work fragment → gives conditional dist given t

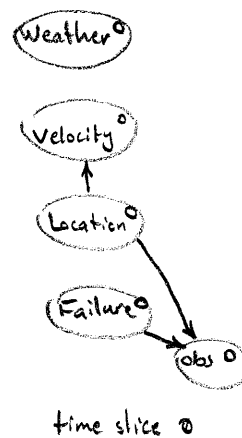
* inter time-slice edges from $X^{(t)} \rightarrow X^{(t+1)}$ are often called persistence edges

* model don't have CPD for time t, only tries to represent the probability of the next time slice given the previous one.

* we need initial state distribution

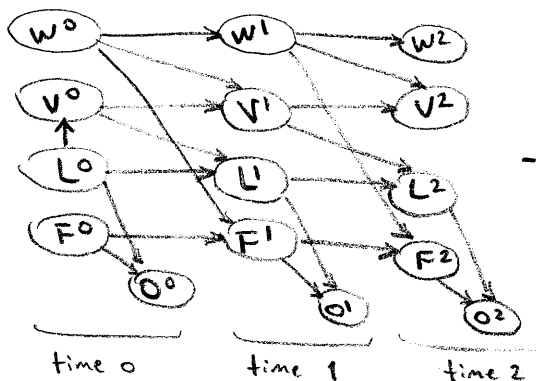
$P(W^{(0)}, V^{(0)}, L^{(0)}, F^{(0)}, O^{(0)}) = P(W^{(0)}) P(V^{(0)} | L^{(0)}) P(L^{(0)}) P(F^{(0)}) P(O^{(0)} | F^{(0)}, L^{(0)})$

chain rule for BN



* ground BN

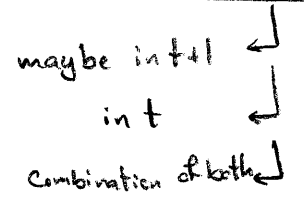
- ↳ copy time slice 0
- ↳ copy prob. dist. @ t1 given t0
- ↳ copy prob. dist @ t2 given t1



• 2-Time-Slice BN

A transition model (2TBN) over X_1, \dots, X_n is specified as a BN fragment such that:
 - the nodes have 2 copies: X_1', \dots, X_n' and a subset of X_1, \dots, X_n

- this subset of X_1, \dots, X_n at time t are variables that directly affect the state at $t+1$
- only the nodes X_1', \dots, X_n' have parents and CPD
- the 2TBN defines a conditional distribution using chain rule $P(X' | X) = \prod_{i=1}^n P(X_i' | Pa(X_i'))$



• Dynamic BN

A Dynamic Bayesian Network (DBN) over X_1, \dots, X_n is defined by a:

- 2TBN BN over X_1, \dots, X_n → dynamics
- a Bayesian network $BN^{(0)}$ over $X_1^{(0)}, \dots, X_n^{(0)}$ → time 0 → initial state

• Ground Network

For trajectory over $0, \dots, T$ we define a ground (unrolled network) such that:

- the dependency model for $X_1^{(0)}, \dots, X_n^{(0)}$ is copied from $BN^{(0)}$ → time 0
- the dependency model for $X_1^{(t)}, \dots, X_n^{(t)}$ for all $t > 0$ is copied from BN → transition

• Summary

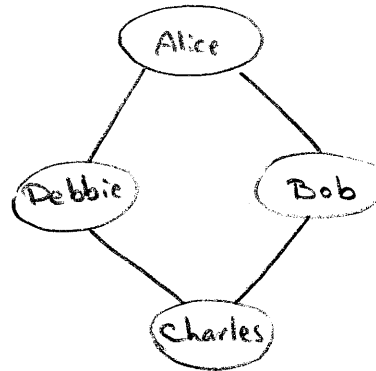
- DBN → compact representation → encoding structured dist.
 - ↳ arbitrarily long temporal trajectories

- Assumptions
 - ↳ Markov assumption
 - ↳ Time invariance

↳ some times are not true → need model redesign for better approximation

* Markov Network Fundamentals:

- Undirected graphical models are typically called MN. → Markov Random Field
- Simplest case: Pairwise Markov Network



* Pairwise Markov Network

• Misconception Example

- Alice & Charles don't get along, Bob & Debbie had breakup
- Random variables indicates whether the students have a particular misconception
- Connection: if two students study together they influence each other
↳ influence goes both ways → undirected
- Parameterization → undirected: no conditional prob. dist.

using factors → general factor: numbers not in [0,1]

↳ names: affinity functions, compatibility function
soft constraints

↳ means: local happiness of vars A and B in $\phi_i[A,B]$ in a joint assignment

↳ the happiest assignment for A and B in isolation of everything else is $a^0 b^0$ ⇒ neither have misconception

↳ the second happiest is $a^1 b^1$ where students both have misconception

↳ B & C really like to agree with each other. A & D like to agree the same of B & C.

↳ C & D likes to argue with each other all the time.

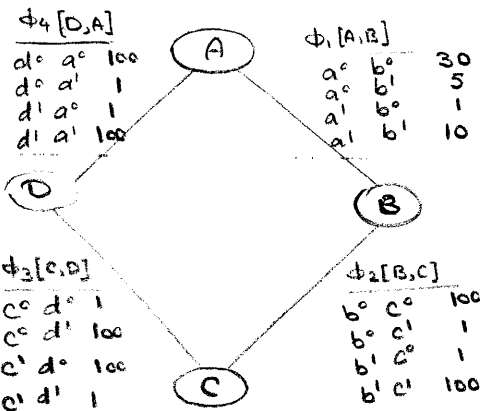
describe the whole situation by a bunch of little pieces together → using product of factors

$$\tilde{P}(A,B,C,D) = \phi_1(A,B) \times \phi_2(B,C) \times \phi_3(C,D) \times \phi_4(D,A)$$

not in range of [0,1] and $\sum \tilde{P}(A,B,C,D) \neq 1$

so \tilde{P} is unnormalized measure.

assignment	unnormalized
$a^0 b^0 c^0 d^0$	300,000
$a^1 b^1 c^0 d^1$	200,000



$\phi_4[D,A]$		
$d^0 a^0$	100	
$d^0 a^1$	1	
$d^1 a^0$	1	
$d^1 a^1$	100	

$\phi_1[A,B]$		
$a^0 b^0$	30	
$a^0 b^1$	5	
$a^1 b^0$	1	
$a^1 b^1$	10	

$\phi_3[C,D]$		
$c^0 d^0$	1	
$c^0 d^1$	100	
$c^1 d^0$	100	
$c^1 d^1$	1	

$\phi_2[B,C]$		
$b^0 c^0$	100	
$b^0 c^1$	1	
$b^1 c^0$	1	
$b^1 c^1$	100	

$\phi_1[A,B]$		
$a^0 b^0$	30	
$a^0 b^1$	5	
$a^1 b^0$	1	
$a^1 b^1$	10	

$\phi_2[B,C]$		
$b^0 c^0$	100	
$b^0 c^1$	1	
$b^1 c^0$	1	
$b^1 c^1$	100	

$\phi_3[C,D]$		
$c^0 d^0$	1	
$c^0 d^1$	100	
$c^1 d^0$	100	
$c^1 d^1$	1	

↳ to normalize it

$$P(A, B, C, D) = \frac{1}{Z} \tilde{P}(A, B, C, D)$$

Z: partition function \hookrightarrow come from statistical physics

\hookrightarrow constant factor that help P to sum to 1

$$\hookrightarrow Z = \sum_{A, B, C, D} \tilde{P}(A, B, C, D)$$

↳ $\Phi_i(A, B)$ = local happiness between A and B, how it related to prob. dist?

\hookrightarrow the marginal prob. $P(A, B)$ X

\hookrightarrow the conditional prob. $P(A|B)$ X

\hookrightarrow the conditional prob. $P(A, B | C, D)$ X

A	B	$P_{\Phi}(A, B)$
a^0	b^0	0.13
a^0	b^1	0.69
a^1	b^0	0.14
a^1	b^1	0.04

$\rightarrow P_{\Phi}$ derives from the set of factors $\Phi = \{\phi_1, \phi_2, \phi_3, \phi_4\}$

\rightarrow this does not agree by the fact A & B likes to agree

\rightarrow highest assignment: $a^0 b^1$

\hookrightarrow this P_{Φ} comes from multiplying all four

B really likes to agree with C

A really likes to agree with D

C and D strongly like to disagree

these 3 assignments are bigger than values in $\Phi_i(A, B)$

so $P(A, B)$ should break the loop of $D \rightarrow A \rightarrow B \rightarrow C \rightarrow D$

because it's weaker, so it is a complicated aggregate of

these different factors that are used to compose the MN

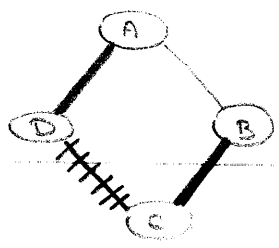
\hookrightarrow there isn't a natural mapping between $P(A, B)$ and

the factors that are used to compose it.

\hookrightarrow In contrast with BN

\hookrightarrow Im possible to look at the prob. dist and say eg. this piece of it is what ϕ_1 ought to be.

\hookrightarrow Makes it harder to learn factors from data because we can't extract them directly from the prob. dist



• Pairwise MN is an undirected graph whose nodes are X_1, \dots, X_n and each edge $X_i - X_j$ associated with a factor $\phi_{ij}(X_i, X_j)$

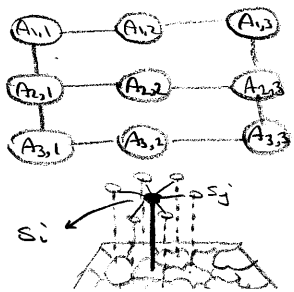
✓ 2D lattice

Images \rightarrow variables = pixels, edges = neighbor

✓ Image segmentation

no regular grid

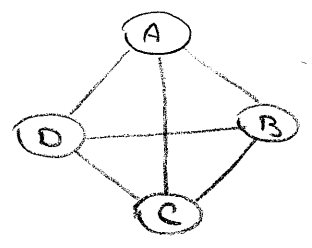
... = superpixels



* General Gibbs Distribution

- More general than the pairwise case
- Most expressive pairwise MN: fully connected

$P(A, B, C, D)$



↳ is this fully expressive? (= can I represent any prob. dist. over four random variable?)

↳ If each X_i has d values in fully connected pairwise MN over $X_1, \dots, X_n \rightarrow \#$ of parameters?

$$\begin{array}{l} \# \text{ of edges in network: } \binom{n}{2} \\ \# \text{ of parameters per edge: } d^2 \end{array} \quad \Bigg| \quad \Rightarrow O(n^2 d^2)$$

can we represent any prob. dist. using $O(n^2 d^2)$ parameters? No!

↳ How many prob. dist. over N random variables where each has D values? $O(d^n)$

$O(d^n) \gg O(n^2 d^2) \rightarrow$ pairwise MN is not sufficiently expressive to have all prob. dist. $\underbrace{dx dx \dots dx}_n$

to increase the expressiveness \rightarrow move away from pairwise edges.

• General Gibbs Distribution

- parameterize using general factors, each with scope that may contain more than two variables

$\phi_i(D_i) \rightarrow$ general factors

$\Phi = \{\phi_i(D_i)\} \rightarrow$ set of factors

- can we represent any prob. dist.? Yes!

↳ prob. dist. is a factor whose scope is the exponential of X_n .

- Formal: A Gibbs dist is parameterized by a set of factors Φ

$\Phi = \{\phi_1(D_1), \dots, \phi_k(D_k)\} \rightarrow$ set of factors

$\tilde{P}_\Phi(X_1, \dots, X_n) = \prod_{i=1}^k \phi_i(D_i) \rightarrow$ factor product, not a prob. dist., unnormalized measure

$Z_\Phi = \sum_{X_1, \dots, X_n} \tilde{P}_\Phi(X_1, \dots, X_n) \rightarrow$ partition function, normalizing constant

$$P_\Phi(X_1, \dots, X_n) = \frac{1}{Z_\Phi} \tilde{P}_\Phi(X_1, \dots, X_n)$$

- what is the MN that we would like to have for a Gibbs dist. of the certain set of fact Φ ?

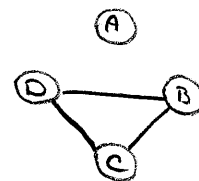
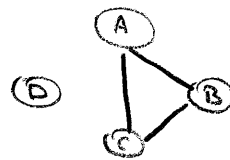
Two factors

$\Phi_1(A, B, C) \rightarrow A, B$ and C all get to interact with each other

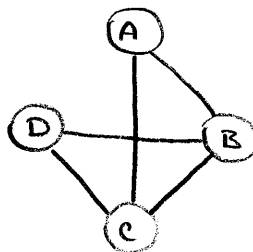
\hookrightarrow links between AB, BC and AC

\hookrightarrow direct prob. relationships between all of them

$\Phi_2(B, C, D)$



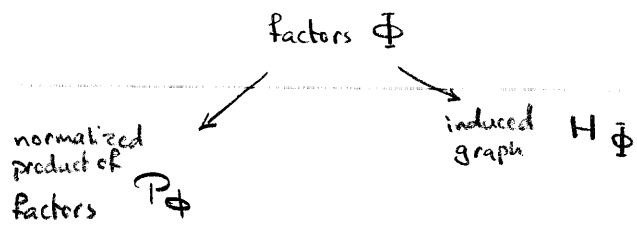
Induced MN by factors Φ_1 and Φ_2



\hookrightarrow Generally if we have a set of factors Φ , where each Φ_i has a scope \mathcal{D}_i the induced Markov Network H_Φ has an edge $X_i - X_j$ whenever there exist a $\Phi_m \in \Phi$ such that $X_i, X_j \in \mathcal{D}_m$.

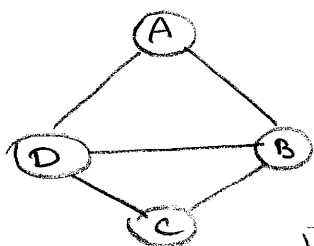
Factorization \rightarrow Answer the question: is there set of factors Φ that are going to let me represent P ?

P factorizes over H if there exist $\Phi = \{\Phi_1(\mathcal{D}_1), \dots, \Phi_k(\mathcal{D}_k)\}$ such that $P = P_\Phi$ and H is induced graph for Φ .



- If we have graph H can we tell what factorization of dist. P over that graph would be?

Which Gibbs distribution would induce the graph H ?



$\hookrightarrow \Phi_1(A, B, D), \Phi_2(B, C, D) \checkmark$

$\hookrightarrow \Phi_1(A, B), \Phi_2(B, C), \Phi_3(C, D), \Phi_4(D, A), \Phi_5(B, D) \checkmark$

$\hookrightarrow \Phi_1(A, B, D), \Phi_2(B, C), \Phi_3(C, D) \checkmark$

We cannot read the factorization from the graph

\hookrightarrow different factorization \rightarrow different expressive power \rightarrow all may induce the same graph.

- \hookrightarrow Pairwise $O(n^2 d^2)$
- \hookrightarrow one giant factor $O(d^n)$
- $\hookrightarrow \dots$

\hookrightarrow why the graph is the same? what does it tell us? it's not telling us the structure of the factorization so what?

- Flow of Influence

↳ when can one var. influence another?

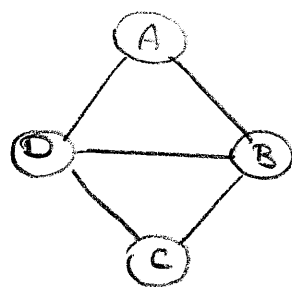
* when B influence D?

Case 1: $\phi_1(A, B, D), \phi_2(B, C, D)$

↳ tie B & D directly

Case 2: $\phi_1(A, B), \phi_2(B, C), \phi_3(C, D), \phi_4(D, A), \phi_5(B, D)$

tie B & D directly



* when A influence C?

Case 1: $A \rightarrow D \rightarrow C : \phi_1(A, B, D) \rightarrow \phi_2(B, C, D)$

Case 2: $A \rightarrow D \rightarrow C : \phi_4(A, D) \rightarrow \phi_3(C, D)$

↳ the parameterization of the two dist. are different BUT

the trail in the graph through which influence can flow is the same REGARDLESS finer grain structure of the factorization → this why graphs are the same

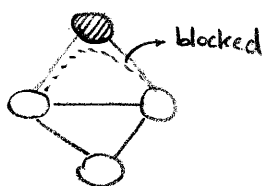
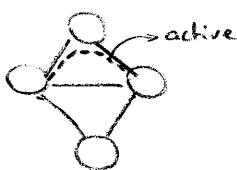
∴ Influence can flow along any trail, regardless of the form of the factors.

- Active trail in MN

A trail $X_1 - \dots - X_n$ is active given \mathcal{Z} if no X_i is in \mathcal{Z}

no block

✓



↳ Difference of active trails in BN and MN

↳ V-structure → BN: \mathcal{Z} observed → create active trail

↳ MN: \mathcal{Z} observed → block the trail

- Summary

1. Gibbs dist represent dist as a product of factors
2. Induced MN connects every pair of nodes that are in the same factor
3. MN structure doesn't fully specify the factorization of \mathcal{P}
4. But active trails depend only on graph structure

* Conditional Random Fields (CRF)

• This class of models are intended to deal with "task specific prediction"

- input / observed variables X

target variables Y

task: predict $Y \rightarrow$ classic models: same type of variables in input and output

☑ Image segmentation X : pixel values & processed features (histograms, colors, textures, ...)

Y : class for every pixel (e.g. grass, sky, cow, water, ...)

Task: going from $X \rightarrow Y$

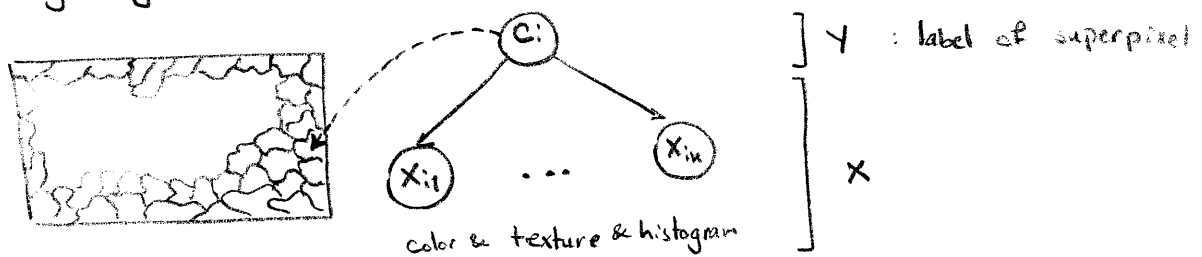
☑ Text processing

X : words in sentence


Y : labels of the words (e.g. person, location, organization, ...)

- Why don't use the models we know so far?

↳ Image Seg. task \rightarrow predict C_i :



↳ high level features \rightarrow more informative
 \rightarrow more correlated with each other

↳ we represented it in naive bayes model () in which features should be independent given label \rightarrow we ignore the correlation structure \rightarrow we count the same feature again and again \rightarrow push us toward very skewed dist \rightarrow far from true belief because of indep. assumption

↳ to correct indep. assumption \rightarrow add some edges to capture the correlations

↳ hard to figure out the correlations

↳ gives rise to densely connected models

↳ Alternative solution: I don't care about the image features!

I don't want to predict the prob. dist. over pixels!

I don't want to do image synthesis!

I don't care about prob. of having green pixel next to another green pixel

I only use X to model distribution over Y !

Reformulate: Instead of modeling $P(X, Y)$

Model $P(Y|X) \rightarrow$ not trying to capture dist over $X \rightarrow$ I should not care about

لا أريد أن أصنع صور

• CRF Representation

- Just like a Gibbs dist

- $\Phi = \{ \phi_1(D_1), \dots, \phi_k(D_k) \}$ \rightarrow set of factors ϕ_i with scope D_i

$$\tilde{P}_{\Phi}(X, Y) = \prod_{i=1}^k \phi_i(D_i) \quad \rightarrow \text{unnormalized measure} = \text{product of factors}$$

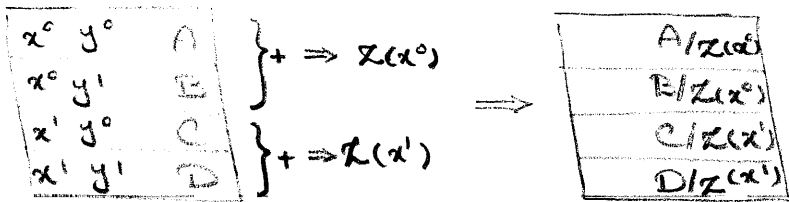
difference with Gibbs

$$Z_{\Phi}(X) = \sum_Y \tilde{P}_{\Phi}(X, Y)$$

\rightarrow partition function, a function of X
 difference with Gibbs
 \rightarrow means: for any given X , we sum over all Y 's corresponding to that X

$$P_{\Phi}(Y|X) = \frac{1}{Z_{\Phi}(X)} \tilde{P}_{\Phi}(X, Y)$$

□



□ consider a new factor $\phi(D)$ such that $D \subseteq X$. what is the effect of adding this factor ϕ on distribution $P(Y|X)$?

$$P_{\Phi}(Y|X) = \frac{\tilde{P}_{\Phi}(X, Y)}{Z_{\Phi}(X)} = \frac{\prod_i \phi_i(D_i)}{\sum_Y \tilde{P}_{\Phi}(X, Y)} = \frac{\prod_i \phi_i(D_i)}{\sum_Y \prod_i \phi_i(D_i)}$$

adding a new factor \rightarrow multiply to numerator & denominator \rightarrow cancels out

• CRFs and Logistic Model

- X_i and Y are binary r.v.

$\mathbb{1}$ indicator function $\rightarrow 1$ if all conditions are met, 0 otherwise

using log-linear model

$$\phi_i(x_i, y) = \exp \{ \omega_i \mathbb{1} \{ x_i = 1, y = 1 \} \}$$

$$\phi_i(x_i, y=1) = \exp \{ \omega_i \mathbb{1} \{ x_i = 1, y = 1 \} \} = \exp \{ \omega_i \mathbb{1} \{ x_i = 1 \} \} = \exp \{ \omega_i x_i \}$$

\rightarrow always 1 binary $x_i \in \{0, 1\}$

$$\phi_i(x_i, y=0) = \exp \{ \omega_i \mathbb{1} \{ x_i = 1, \text{false} \} \} = \exp \{ \omega_i x_i 0 \} = e^0 = 1$$

$$\tilde{P}_{\Phi}(X, Y=1) = \exp \left\{ \sum_i (\omega_i X_i) \right\}$$

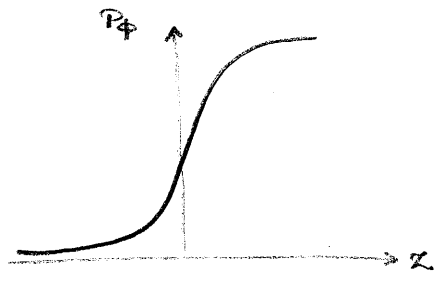
↙ $\prod_i e^{\omega_i X_i}$

$$\tilde{P}_{\Phi}(X, Y=0) = 1$$

↙ product of ones

$$Z_{\Phi}(X) = \sum_Y \tilde{P}_{\Phi}(X, Y) = \tilde{P}_{\Phi}(X, Y=1) + \tilde{P}_{\Phi}(X, Y=0) = 1 + \exp \left\{ \sum_i \omega_i X_i \right\}$$

$$P_{\Phi}(Y=1 | X) = \frac{\exp \left\{ \sum_i \omega_i X_i \right\}}{1 + \exp \left\{ \sum_i \omega_i X_i \right\}} = \frac{e^z}{1 + e^z}$$



Logistic is very simple CRF.

☐ CRFs for Image Segmentation

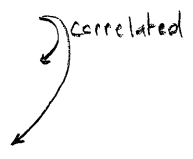
- Node Factors can use any features of the image
 - Color histograms
 - Texture features
 - Discriminative Patches → like looking for the eye on a cow
- Features can be in and out of the superpixel
 - ↳ green under this superpixel ☐ on grass
 - ↳ counting redundant superpixels? I don't care
- Correlation doesn't matter
- Can train a discriminative classifier (SVM, Boosting) to improve performance
 - ↳ e.g. train very strong classifiers for node potentials then adding pairwise or higher order features between the Y 's.

$$P(Y_i | X_1, \dots, n)$$

features of log-linear model not image features

☐ CRFs for Languages

- High level features → high correlation
 - ↳ word capitalized
 - ↳ word in Atlas or NameList
 - ↳ previous word is "Mrs."
 - ↳ next word is "Times"



- $P(\text{Labels | words})$

• Summary

- CRF is parameterize the same as a Gibbs dist. but normalized differently → $P(Y|X)$
 - ↳ Generalizes logistic regression models
- Don't need to model dist. over vars. we don't care about
- Allows models with highly expressive features without worrying about wrong independencies.
 - ↳ Enable design really expressive predictors of pieces of the model

* Markov Networks and Independencies

what kind of independencies are encoded by the structural graph of MN

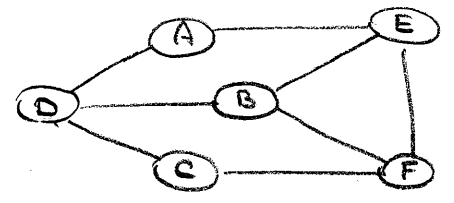
• Separation in MNs

- X and Y are separated in H given Z if there is no active trail in H between X and Y.

- simpler than d-separation → no different cases
↳ directed

What does it take to separate A from E?

A and E separated given B, D
given D
given C



• Factorization ⇒ Independence

- Theorem: If P factorizes over H, and $sep_H(X, Y | Z)$ then P satisfies $(X \perp Y | Z)$

$$I(H) = \{ (X \perp Y | Z) : sep_H(X, Y | Z) \}$$

If P satisfies $I(H)$ we say that H is an I-map of P

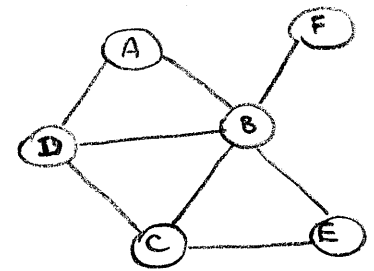
↳ all the independencies that we can read from graph H

- Theorem: If P factorizes over H, then H is an I-map of P

P factorizes over H → satisfy all of the independencies that one can read from H.

For the graph H shown here, which independence assertions are in $I(H)$

- A ⊥ C | B, D ✓
- A ⊥ E | B ✗
- C ⊥ F | B ✓
- A ⊥ F ⊥ B, D ✓



- In BN the converse holds, in MN the converse sort of holds.

• Independence ⇒ Factorization

- Theorem (Hammersley-Clifford): For a positive dist. P, if H is an I-map for P then

P factorizes over H.

↳ difference with BN
 $P(x_i) > 0$

↳ if distribution involves deterministic relationship this property no longer holds.

$$P(x_i) \in \{0, 1\}$$

• Summary:

- Two almost equivalent views of graph structure: Factorization, I-map (for positive dist.)

- If P factorizes over a graph H, we can read from the graph independencies that must hold in P

I-map(P)

* I-maps and Perfect Maps :

• A graph structure $\xrightarrow{\text{encode}}$ set of independencies $\xrightarrow{\text{necessarily hold}}$ every dist. that can be encoded as BN over that graph

↳ How to take a distribution that has a certain set of independencies that satisfies and encode it within a graph structure?

• what is the independencies of a dist? Capturing indep. in P

$$I(P) = \{ (X \perp Y | Z) : P \models (X \perp Y | Z) \}$$

↳ set of all independence that hold for dist P

↳ write all of the exponential numbers of possible indep. statements $\rightarrow I(P)$ is the ones that hold

if P factorizes over G \Rightarrow G is I-map for P

↳ every independence that hold in G that is implied by d-separation properties of G also holds in P

$$\rightarrow I(G) \subseteq I(P)$$

↳ there can be indep. in $I(P)$ that are not in $I(G)$

• We have a graph that doesn't capture some of the independencies in P then it's unnecessarily complicated

↳ we want a sparser graph \rightarrow encodes more indep.
fewer parameters \rightarrow more efficient inference
more informative

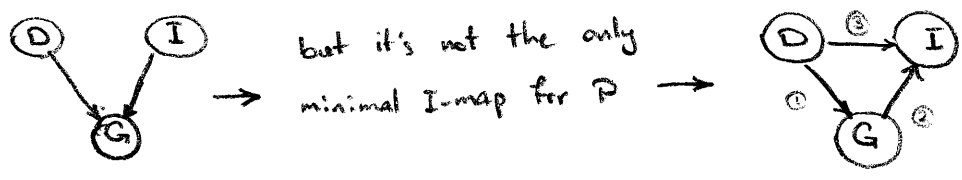
- want a graph that captures as much of the structure in P as possible.

- Minimal I-map : I-map without redundant edges

□ If we have $(X \rightarrow Y)$ but Y doesn't really depend on X in the sense $P(Y|x^0) = P(Y|x^1)$ we can remove the edge and it remains an I-map, so the first I-map was not minimal

↳ Minimal I-map may still not capture $I(P)$!

□ we have dist P that corresponds to left minimal I-map



it is a minimal I-map because no edge could be removed.

remove:

- ① $\Rightarrow D \perp G \mid X$
- ② $\Rightarrow G \perp I \mid D, X$
- ③ $\Rightarrow D \perp I \mid G, X$

- so Imap is not the best tool for capturing structure in data distribution

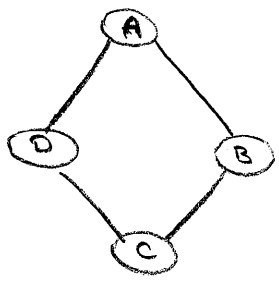
• Perfect Map

- Perfect map: $I(G) = I(P) \rightsquigarrow$ BN as a perfect map

↳ G perfectly captures independencies in P

↳ usually hard to come by perfect maps \rightarrow sometime dist. doesn't have perfect maps.

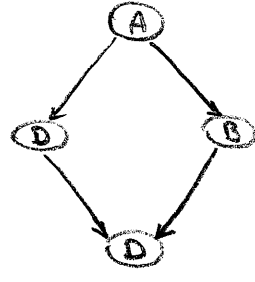
✓



$P \models A \perp C \mid B, D$

$P \models B \perp D \mid A, C$

attempt 1

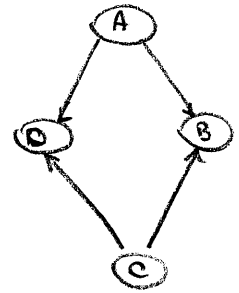


$G: B \perp D \mid A$

$P \not\models B \perp D \mid A$

not an I-map

attempt 2



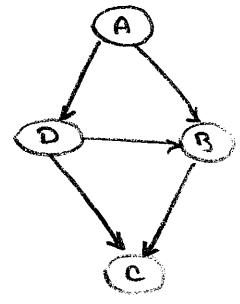
$G: B \perp D \mid A, C$

$G: A \perp C$

$P \not\models A \perp C$

not an I-map

attempt 3



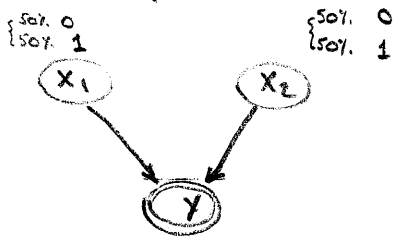
$G: A \perp C \mid B, D$

$I(G) \subset I(P)$

I-map ✓

perfect map ✗

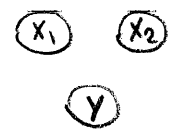
✓ XOR example



X_1	X_2	Y	Prob
0	0	0	0.25
0	1	1	0.25
1	0	1	0.25
1	1	0	0.25

$X_1 \perp X_2$
 $X_1 \perp Y$
 $X_2 \perp Y$

no I-map for this graph because the only way is:



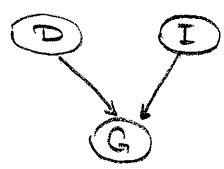
that is not an I-map

- perfect map: $I(H) = I(P) \rightsquigarrow$ MN as a perfect map

↳ H perfectly captures independencies in P

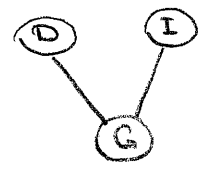
↳ can we capture all possible dist. in terms of a MN as a perfect map? NO!

✓



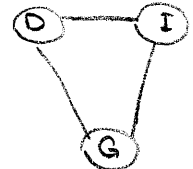
$P \models D \perp I$

has perfect map as a BN
 doesn't have perfect map as a MN
 v-structure example



$H: G \perp I \mid G$

$P \not\models G \perp I \mid G$



$P \not\models D \perp I$

✗

- Uniqueness of Perfect Map

✓ simplest case



$I(G_1) = \emptyset$

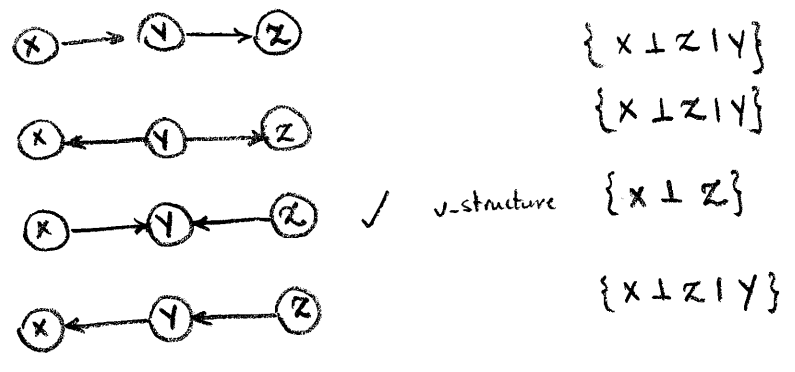


$I(G_2) = \emptyset$

\Rightarrow

can represent exactly the same set of all of distributions

Which of the following graphs does not encode the same indep. as the others?



- I-equivalence

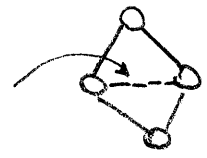
Two graphs G_1 and G_2 over X_1, \dots, X_n are I-equivalent if $I(G_1) = I(G_2)$

Two graphs G_1 and G_2 over X_1, \dots, X_n are I-equivalent if $X \rightarrow Y \rightarrow Z$ and $X \leftarrow Y \rightarrow Z$ and $X \leftarrow Y \leftarrow Z$ are I-equivalent

- ↳ tells us that there's certain aspects of the graphical model that are unidentifiable.
- ↳ without prior knowledge or preference of who is the parent of who we have no single choice
- ↳ most graphs have large numbers of I-equivalents.
- ↳ complicated e.g. for learning

• Summary

- Graphs that capture more of $I(P)$ are more compact and provide more insight
- A minimal I-map may fail to capture alot of structure even if present and representable as a PG
- A perfect map is great, but may not exist
- Converting $BN \rightleftharpoons MN$ loses independencies
 - ↳ $BN \rightarrow MN$: loses indep. in v-structures
 - ↳ $MN \rightarrow BN$: must add triangulating edges to loops



* Log-Linear Models

- Local structure that doesn't need full table representations is important in both directed and undirected models
 - ↳ how do we incorporate input structure into undirected models? → Log-linear models!

• Log-Linear Representation

$\tilde{P} = \prod_i \phi_i(D_i)$ → potentially a full table



$\tilde{P} = \exp(-\sum_j w_j f_j(D_j))$ → linear model → because logarithm is linear fun

coefficients features scope of feature j

But different features can have the same scope. Each feature have specific coeff.

↓

$$\tilde{P} = \prod_j \underbrace{\exp(-w_j f_j(D_j))}_{\text{factor}} \rightarrow \text{product of factors that each has a single variable } w_j$$

• Representing Table Factors

$$\Phi(X_1, X_2) = \begin{bmatrix} a_{00} & a_{01} \\ a_{10} & a_{11} \end{bmatrix}$$

factor Φ over two binary r.v. X_1 and X_2
 ↳ Full table has 4 params: $a_{00}, a_{01}, a_{10}, a_{11}$

we can capture this model using log-linear model using a set of such features
 ↳ indicator functions

$$f_{12}^{00} = \mathbb{1} \{X_1=0, X_2=0\}$$

$$f_{12}^{01} = \mathbb{1} \{X_1=0, X_2=1\}$$

$$f_{12}^{10} = \mathbb{1} \{X_1=1, X_2=0\}$$

$$f_{12}^{11} = \mathbb{1} \{X_1=1, X_2=1\}$$

$$\Phi(X_1, X_2) = \exp\left(-\sum_{k=0,1} \sum_{l=0,1} w_{kl} f_{ij}^{kl}(X_1, X_2)\right)$$

↳ when $X_1=0$ and $X_2=0$: $\Phi = \exp(-w_{00})$
 ↳ when $X_1=0$ and $X_2=1$: $\Phi = \exp(-w_{01})$

if we assume $w_{kl} = -\log a_{kl}$ we get the original factor back.
 but generally we want a much finer grain set of features.

□ Language model

Y : represent the annotations for each word in the sequence corresponding to what category ...
 ↳ Begin-Person, Intermediate-Person, OTH, B-location, I-loc
 nothing → no person or location, other
 ↳ value Y tells us what category it belongs to

X : actual words in the sentence

Full table rep. : tries to relate each and every Y that has a feature, to every word in language

features: $f(Y_i, X_i) = \mathbb{1} \{Y_i = \text{B-person}, X_i = \text{capitalized}\}$
 ↳ how important is capitalization for recognizing that sth is a person

$= \mathbb{1} \{Y_i = \text{B-location}, X_i \text{ appears in Atlas}\}$

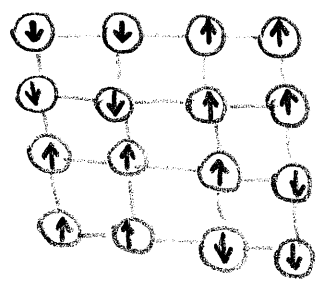
□ Ising Model : from statistical physics

- looks at the pairs of variables → Pairwise MW
 looks at pairs of adjacent variables and gives a coefficient for their product

$$E(x_1, \dots, x_n) = -\sum_{i < j} w_{ij} \cdot x_i \cdot x_j - \sum_i u_i x_i$$

$$f(x_i, x_j) = x_i \cdot x_j \quad x_i \in \{-1, +1\}$$

- modeling spins of electrons in a grid



↳ P = the joint spins in the model

↳ depends on whether adjacent atoms have the same spins or opposite spin.

$$P(X) \propto e^{-\frac{1}{T} E(X)}$$

↳ depends on ω

↳ systems that atoms spin at same direction VS. system that atoms have opposite spin.

ferro magnetic anti ferro magnetic

↳ temperature T : how strong is this connection

↳ T ↑ : $\frac{\omega_{ij}}{T} \rightarrow 0 \rightarrow$ atoms become decoupled from each other

T ↓ : $\frac{\omega_{ij}}{T} \uparrow \rightarrow$ atoms impose a much stronger effect on each other

• Metric MRF

- another kind of features that are used in real applications is metric MRF

↳ sth that comes up mostly when random var. X_i all take values in label space V

↳ all binary, all $\in \{1, 2, 3, 4\}$

↳ we want X_i and X_j that are connected to each other take similar values.

↳ we need similarity relation

- Distance function : $\mu : V \times V \rightarrow \mathbb{R}^+$
for X_i for X_j

↳ need to satisfy the standard condition on a distance function or metric

* Reflexivity $\mu(v, v) = 0, \forall v$

* Symmetry $\mu(v_1, v_2) = \mu(v_2, v_1), \forall v_1, v_2$

* Triangle Inequality $\mu(v_1, v_2) \leq \mu(v_1, v_3) + \mu(v_3, v_2), \forall v_1, v_2, v_3$

↳ if a function just satisfies the first two \rightarrow semi-metric

↳ put in MRF $f_{i,j}(x_i, x_j) = \mu(x_i, x_j)$

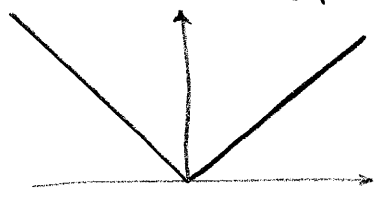
$$\phi(x_i, x_j) = \exp(-\omega_{ij} \cdot f_{i,j}(x_i, x_j)), \omega_{ij} > 0$$

↳ lower distance $\mu \downarrow \Rightarrow \phi \uparrow$

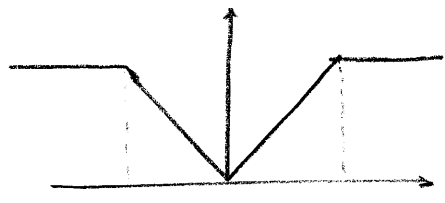
values of x_i and x_j are far \Rightarrow lower probability

$$\mu(v_k, v_l) = \begin{cases} 0 & v_k = v_l \\ 1 & \text{otherwise} \end{cases} \Rightarrow \text{step function} \Rightarrow \Phi = \begin{bmatrix} 0 & 1 & 1 & \dots & 1 \\ 1 & 0 & 1 & & \\ \vdots & 1 & 0 & & \\ \vdots & & \vdots & \ddots & \\ 1 & & & & 0 \end{bmatrix}_{V \times V}$$

$$\mu(v_k, v_l) = |v_k - v_l|$$



$$\mu(v_k, v_l) = \min(|v_k - v_l|, d)$$



$\mu(v_k, v_l) = \begin{cases} 0 & v_k = v_l \\ 1 & \text{otherwise} \end{cases}$
 Image Segmentation: we want the adjacent super pixels to take the same class

Image Denoising \rightarrow Gaussian Noise
 X : noisy pixels
 Y : clean pixels

we want y_i to be close to $x_i \rightarrow$ recover original image
 adjacent pixels tend to have the same value \rightarrow constraint

\rightarrow linear penalty \rightarrow fragile model \rightarrow tend to set all the pixels equal
 \rightarrow truncated linear \rightarrow only close pixels are important

Stereo Reconstruction

y_i : depth disparity for a given pixel in the image \rightarrow we have spatial continuity.
 \rightarrow depth of a pixel is close to the depth of adjacent pixel

we don't want to push this constraint too hard \rightarrow truncated linear model + some tricks
 tricks \rightarrow pixel appearance: color, texture, ... \rightarrow adj' pixel are nearly the same \Rightarrow same object
 \rightarrow adj' pixel are very different \Rightarrow another object

Shared Features in Log-Linear Models

\rightarrow In most MRFs, same feature and weight are used over many scopes

Ising Model

$$E(x_1, \dots, x_n) = - \sum_{(i,j) \in \text{edges}} \underbrace{w_{ij}}_{\text{weights}} \underbrace{x_i x_j}_{f(x_i, x_j)} - \sum_i d_i x_i$$

\rightarrow same weight for every adjacent pair $\rightarrow w = w_{ij}$
 \rightarrow feature $f(x_i, x_j) = x_i x_j$ is reused throughout the model

Natural Language Processing

\rightarrow same energy terms $w_k f_k(x_i, y_i)$ repeat for all positions i in the sequence
 \rightarrow same energy terms $w_m f_m(y_i, y_{i+1})$ also repeat for all positions i

Image Segmentation

↳ same features and weights for all superpixels in the image

- Need to specify for each feature f_k a set of scopes $Scopes[f_k]$

↳ for each $D_k \in Scopes[f_k]$ we have a term $w_k f_k(D_k)$ in the energy function

↳ e.g. if we want f_k to be applied to adjacent superpixels in the image

$$Scope[f_k] = \{Y_i, Y_j : i, j \text{ adjacent}\}$$

↳ Weight depends on the feature, but not in the scope
 w_k f_k D_k

$$w_k \sum_{D_k \in Scopes(f_k)} f_k(D_k)$$

Summary

↳ Same feature & weight can be used for multiple subsets of variables

↳ pairs of adjacent pixels / atoms / words

↳ occurrence of same word in the document

↳ Can provide a single template for multiple MNs

↳ different images

↳ different sentences

↳ Parameters and structure are reused within an MN and across different MNs.

↳ Need to specify set of scopes for each feature

* Knowledge Engineering

↳ How to build a graphical model for some application having all of this info → Art, Black Magic

• Important Distinctions

1- Template based vs. Specific

↳ specific model for fixed set of r.v

2- Directed vs. Undirected

3- Generative vs. Discriminative

⇒ Hybrids are also common

Hybrids

template-based pieces + stuff that is not in template level
directed components + undirected ones

• Template - Based vs. Specific

Specific

Medical Diagnosis

↳ particular set of symptoms, diseases, ...

Template - Based

Image Segmentation

↳ different images with same model

Fault diagnosis

unique elements → generic printer

particular printer → shared elements

- the place decided to sit on this spectrum changes knowledge engineering problem

- template-based → usually have a small number of variable types

↳ Image segmentation → class label only

↳ most effort goes to features

↳ the main design challenge goes to figuring out that which features are most predictive

↳ feature engineering

- specific → usually a large number of unique vars. → each requires its own model

• Generative vs. Discriminative

Generative

- Don't have a predetermined task

↳ task shifts

↳ e.g. medical diagnosis → every patient present differently

- Flexibility

↳ to measure different vars. and predict others

- Easier to train in certain regimes

↳ the case that data are not fully labeled

? what are the key decisions about designing a graphical model

• Variable Types

- Target

↳ regardless of whether we have a fix or varying task in hand we have a set of target vars.

- Observed

↳ including complex, constructed features

- Latent / Hidden

Discriminative

- Particular Prediction task

↳ Richly expressive features

↳ Avoid dealing with correlations

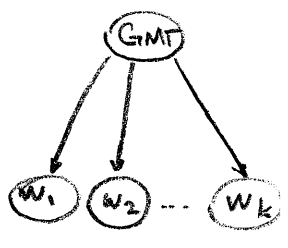
↳ Higher performance

↳ we care about 'e

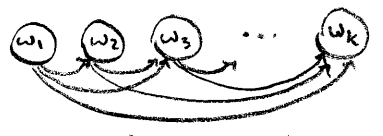
↳ we don't care about predicting necessarily if it is not observed
↳ WTF it is used for then!?

What time does your watches show?

w_i : the time watch show → all correlated with each other
↳ not because every body set it with each other
↳ using GMT



naive bayes ⇒ $w_i \perp w_j \mid GMT$ → GMT is latent
↳ if we don't have GMT ⇒



Fully connected model required to capture correlation

↳ latent variables can simplify our structure

Structure

Causal vs. non-causal ordering

↳ in BN the question is: do arrows show causality?

$(X) \rightarrow (Y) \stackrel{?}{\Rightarrow}$ causal connection between X and Y → YES and NO!

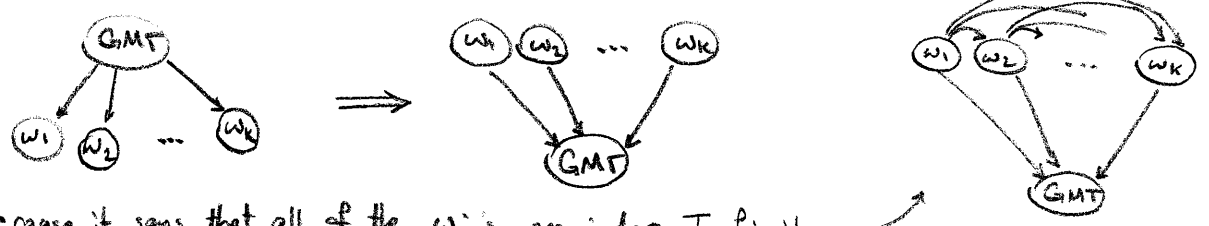
↳ Any dist that I model this way ($X \rightarrow Y$) I can equally model this way ($Y \rightarrow X$)

↳ we can reverse the edges and have a model that is equally expressive → for any distribution

↳ nasty in case of v-structure $X_1 \rightarrow Y \leftarrow X_2 \Rightarrow X_1 \rightarrow X_2$
if we reverse edge $X_2 \rightarrow Y$ we should add extra link $X_1 \rightarrow X_2$

↳ causal directionality is often simpler

Can we invert the edges? Is this the correct model?



No, because it says that all of the w_i 's are indep. To fix it ...

↳ So, causal ordering although is not more correct than a non causal ordering, it's sparser

↳ Causal ordering → sparser, more intuitive, easier to parameterize

Extending the Conversation

↳ usually we don't know how to construct distribution using variables $X_1 \dots X_n$.

↳ we need to figure out how to encode them using a graph

↳ first step is to have set of variables that we wish to reason about

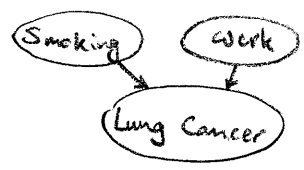
Lung Cancer

1.

Lung Cancer

2. what influences whether somebody's going to get lung cancer?

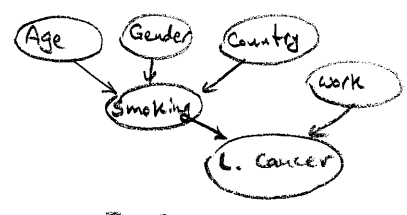
↳ ask doctor about reasons



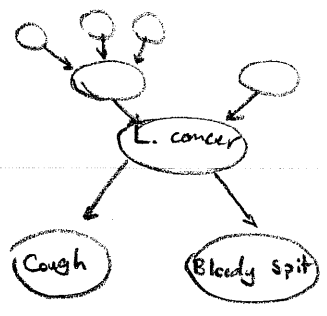
3. continue questioning: what is the prob. of someone smokes?

↳ continue backward up to the point there's no way to extend the conversation

↳ e.g. what is the prob. of Gender being male vs. female

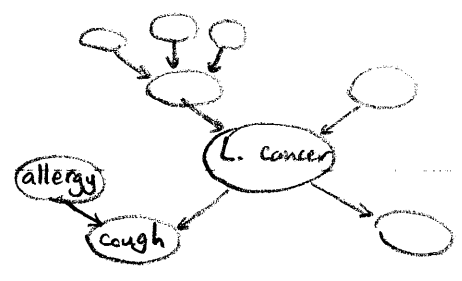


4. factors that help us indicate whether someone has lung cancer.



5. what is the probability P(cough | lung cancer)?

↳ extend the conversation backwards



Parameters

— values

↳ what matters?

↳ Zeros : make big difference

cause many of the mistakes → giving zero to unlikely but not impossible event

↳ Orders of Magnitude : big difference between 0.1 and 0.01

↳ Relative Values : relative values between conditional prob. make much bigger difference to the answer than the absolute prob.

↳ compare values in each CPD is useful to evaluate the PSM

↳ values used for relative ratios should make sense

— structured CPDs

↳ CPD tables are rare → small application

↳ context-specific: vars make a difference in some circumstance but not in others

↳ aggregating of multiple weak influences

	Context Specific	Aggregating
Discrete	Tree CPD	Sigmoid CPD Noisy OR
Continuous	Regression Tree CLG	Linear Gaussian

↳ regression tree : continuous version of tree CPD

↳ breaks up context based on some thresholds on continuous variables

↳ CLG : conditional linear Gaussian

• Iterative Refinement

↳ the model is rarely done the first time it is written → start → test → improve over time

- Model Testing

↳ ask queries, check results

- Sensitivity Analysis for parameters

↳ look at query and ask which parameters have the biggest difference on the value of the query

↳ those are probably the one we should fine tune

- Error Analysis

↳ add features → maybe help to eliminate certain types of error

↳ add dependency → give the model the ability to capture the structure

☑ For which of the following applications might a generative model be a better choice than discriminative model?

x making fast predictions of Y given x where the dist. of x is not important

✓ medical diagnosis where obtaining ground-truth labels for data is expensive

x image segmentation where the extracted features are high dimensional

✓ modeling different music composers: for classification of music and creating new music pieces

↳ discriminative models directly model $P(Y|X)$ and avoid dealing with joint distribution

↳ generative models handles missing data or partially labeled data better so they are preferable if obtaining labels is hard

↳ High-D variables are in general more difficult to model, and disc. models avoid this by conditioning on the

↳ generative models can be used to generate new observable data

* Inference :: Conditional Probability Queries

• set of observations $E = e \rightarrow$ evidence

a subset of variables $Y \rightarrow$ query

compute $P(Y | E=e) \rightarrow$ task

□ applications

- medical / fault diagnosis \rightarrow observations: symptoms, test results

\rightarrow task: predicting the probability of different failure modes / diseases

- Pedigree analysis \rightarrow observations: phenotypes, genotypes of some individuals

\rightarrow task: predicting about other individuals

• NP-Hard

- Inference in PGM is NP-Hard

\rightarrow if the problem is shown to be NP-Hard \Rightarrow it is extremely unlikely to have an efficient solution.

\rightarrow means that all algorithms that people have devised for this problem (and a bunch of problems that are equally hard) require at the very least time exponential in the size of representation of the problem

\rightarrow unlikely to be solvable for any problem that is bigger than small number of vars.

- Which particular problems in PGM is NP-Hard?

1. Exact Inference: Given a PGM P_{Φ}
a variable X and a value $x \in \text{Val}(X)$

compute $P_{\Phi}(X=x)$

2. $P_{\Phi}(X=x) > 0$

ok! exact inference is NP-hard, but can we approximate it? NO!

3. Approximate Inference: Given a PGM P_{Φ}
a variable X and a value $x \in \text{Val}(X)$

and observation $e \in \text{Val}(E)$

problem of finding an approximation of answer P that we guarantee

$$|P_{\Phi}(X=x | E=e) - P| < \epsilon, \quad \epsilon < 0.5$$

\rightarrow the case $\epsilon = 0.5$ can be obtained by simply random guessing

\rightarrow should we get depressed? This is a worst case result!

\rightarrow we can construct this bizarre PGMs with exponential time at worst case but in general case one can't do better

\rightarrow so we are looking at where to get better results...

• joint dist. defined by chain rule for BN.

$$\underbrace{\phi_c(C)}_{P(C)} \phi_D(D) \phi_I(I) \underbrace{\phi_G(G, I, D)}_{P(G|I, D)} \phi_S(S, I) \phi_L(L, G) \phi_J(J, L, S) \phi_H(H, G, J)$$

- Transform a BN to a set of factors.

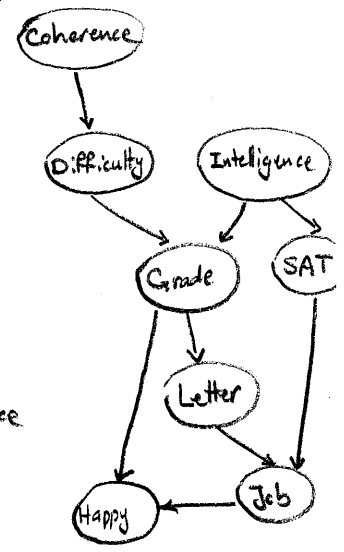
↳ each CPD to a factor of variable and its parents

- goal: compute $P(J)$

↳ $\sum_{C, D, I, G, S, L, H}$: marginalization

↳ this is why the problem of inference in conditional prob. inference

is called Sum-Product : sum over a product of factors
 ↳ root of becoming intractable



• The same for MN

$$\phi_1(A, B) \cdot \phi_2(B, C) \cdot \phi_3(C, D) \cdot \phi_4(A, D)$$

goal: compute $P(D)$

↳ $\sum_{A, B, C} \phi_1(A, B) \phi_2(B, C) \phi_3(C, D) \phi_4(A, D)$

↳ but it is not quite right, because it's not equal to $P(A, B, C, D)$

↳ it is unnormalized measure $\tilde{P}(A, B, C, D)$

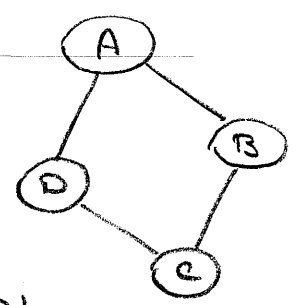
↳ should be divided by partition function

$$\frac{1}{Z} \sum_{A, B, C} \phi_1(A, B) \phi_2(B, C) \phi_3(C, D) \phi_4(A, D)$$

↳ how do we compute Z ? to goal is $P(D)$ not Z !

↳ $P(D) = \frac{1}{Z} \tilde{P}(D)$

so we calculate $\tilde{P}(D)$ and do renormalization



• Evidence: Reduced Factors

- only requires simple preprocessing step of factor reduction

$$P(Y | E=e) = \frac{P(Y, E=e)}{P(E=e)}$$

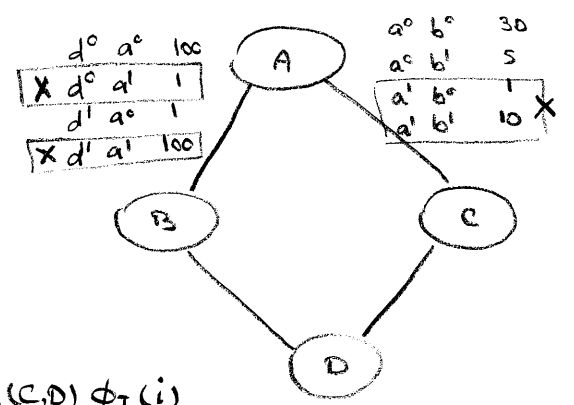
$W = \{X_1, \dots, X_n\} - Y - E$ → variables that are neither the query nor evidence

$P(Y, E=e) = \sum_W P(Y, W, E=e)$ → marginalizing out the W , the form of Sum-Product

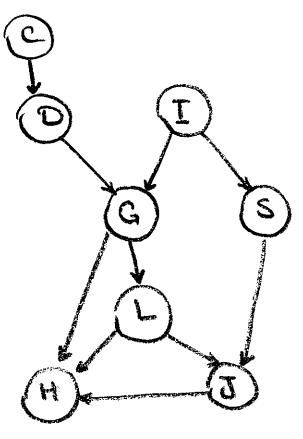
$= \sum_W \frac{1}{Z} \prod_{\kappa} \phi_{\kappa}(D_{\kappa}, E=e)$ → only the factors that are consistent with $E=e$
 ↳ reduce the factors by evidence
 $= \prod_{\kappa} \phi_{\kappa}(D_{\kappa}, E=e)$

✓ $A = a^0$

after reducing the factors we can treat it exactly the same way we treat it before.



□ $P(J, I=i, H=h)$



$P(J, I=i, H=h) \propto$

$$\sum_{C, D, I, G, S, L, H} \phi_C(C) \phi_D(C, D) \phi_I(i) \phi_G(G, i, D) \phi_S(S, i) \phi_L(L, G) \cdot \phi_J(J, L, S) \phi_H(h, G, J)$$

then renormalize it to have

$P(J | I=i, H=h)$

• Sum-Product Algorithm Summary

1. conditional prob \triangleright ratio

$$P(Y | E=e) = \frac{P(Y, E=e)}{P(E=e)}$$

2. numerator = product of reduced factors

$$P(Y, E=e) = \sum_w \frac{1}{Z} \prod_k \phi'_k(D'_k)$$

3. denominator = numerator summed up over the vars and query vars Y .

$$P(E=e) = \sum_Y \sum_w \frac{1}{Z} \prod_k \phi'_k(D'_k)$$

↳ OR compute $\sum_w \prod_k \phi'_k(D'_k)$ and normalize at the end

- many algorithms for computing conditional prob. queries

↳ Push summation into factor products \Rightarrow variable elimination

↳ a special case of Dynamic Programming

↳ form of exact inference

↳ Message passing over a graph

↳ exact variations : Belief propagation

↳ variational approximations

↳ Random sampling instantiations

↳ Markov Chain Monte Carlo (MCMC), Importance Sampling

↳ approximation inference

- # of entries in a table describing the full joint dist, where each of the n variables has k possible values?

○ n^k ○ $n^2 k$ ○ $n k^2$ ✓ k^n $\frac{k}{\text{var } 1} \frac{k}{\text{var } 2} \frac{k}{\text{var } 3} \dots \frac{k}{\text{var } n} = k \times k \times \dots \times k = k^n$

* Inference :: MAP Inference

• Maximum a Posterior: (MAP)

Evidence : $E = e$

Query : all other variables Y ($Y = \{X_1, \dots, X_n\} - E$)

Task : map assignment $MAP(Y|E=e) = \arg \max_y P(Y=y | E=e)$

↳ finding assignment y to variables Y that maximizes the conditional prob. of Y given evid

↳ the solution might not be unique → several different assignment that give exact prob.

✓ application

- message decoding : most likely transmitted message → assignment to the transmitted bits

- image segmentation : most likely segmentation

• MAP ≠ Max over Marginals

✎



	$P(A)$	$P(B A)$	$P(A, B)$
a^0	0.4	$a^0 b^0$ 0.1	$a^0 b^0$ 0.04
a^1	0.6	$a^0 b^1$ 0.9	$a^0 b^1$ 0.36
		$a^1 b^0$ 0.5	$a^1 b^0$ 0.3
		$a^1 b^1$ 0.5	$a^1 b^1$ 0.3

→ MAP assignment = highest probab:

look separately at A and B ⇒ most likely assignment is a^1 X

we can't look separately at the marginals over A and over B

↳ we are looking for a single assignment over all of the variables that together has the highest prob.

• NP-Hard

- which problems are NP-Hard?

1. Given a PGM P_Φ

Find a joint assignment α with highest probability $P_\Phi(\alpha)$

2. Given a PGM P_Φ

and a probability P

decide if there is an assignment α such that $P_\Phi(\alpha) > P$

↳ Should we give up? NO!

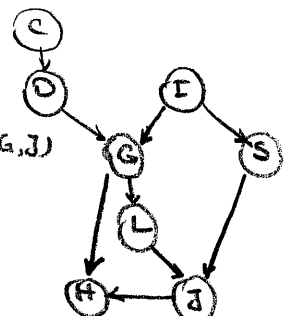
↳ we have even broader set of solutions compared to conditional prob. queries case

• Max-Product in BN

↳ now instead of marginalizing and summing out, we have max

$$\arg \max_{c, D, I, G, S, L, J, H} \Phi_c(c) \Phi_D(D, c) \Phi_I(I) \Phi_G(G, I, D) \Phi_S(S, I) \Phi_L(L, G) \Phi_J(J, L, S) \Phi_H(H, G, J)$$

↳ max of product ~ Max-Product



• Max-Product Algorithm Summary

1. conditional probability \triangleright ratio

$$P(Y|E=e) = \frac{P(Y, E=e)}{P(E=e)} \quad , Y = \{x_1, \dots, x_n\} - E$$

$\underbrace{\hspace{10em}}_{\text{constant w.r.t } Y}$

2. numerator = unnormalized measure

$$P(Y; E=e) = \frac{1}{Z} \prod_k \phi'_k(D'_k)$$

$\underbrace{\hspace{10em}}_{\text{normalizing by partition function constant w.r.t } Y}$

3. find max assignment

$$\underset{Y}{\operatorname{argmax}} P(Y|E=e)$$

OR compute $\underset{Y}{\operatorname{argmax}} \prod_k \phi'_k(D'_k)$

- many algorithms to compute MAP

\hookrightarrow Push maximization into factor product \Rightarrow Variable Elimination

$\underbrace{\hspace{2em}}_{\text{max Product}}$

\hookrightarrow Message passing over a graph

\hookrightarrow Max Product Belief Propagation

\hookrightarrow It is an optimization problem : we can use specific algorithms for optimization

\hookrightarrow Integer Programming : general class of optimization over discrete spaces

\hookrightarrow state-of-the-art algs are based on this especially in approximate case

\hookrightarrow for certain types of networks there are specific very efficient algorithms

\hookrightarrow Graph Cuts

\hookrightarrow Combinatorial Search

• Summary

- MAP : single coherent assignment of highest probability

\hookrightarrow not the same as maximizing individual marginal probs.

- Maxing over factor product

- Combinatorial optimization problem

- Many exact and approximate solutions

* Inference :: Variable Elimination Algorithm

- simplest and most fundamental algorithm

✓ Elimination in Chains

$$P(E) \propto \sum_D \sum_C \sum_B \sum_A \tilde{P}(A, B, C, D, E)$$



$$= \sum_D \sum_C \sum_B \sum_A \phi_1(A, B) \phi_2(B, C) \phi_3(C, D) \phi_4(D, E)$$

→ pairwise factors

→ move factors that doesn't mention a particular var. in their scope out of summation over it.

$$= \sum_D \sum_C \sum_B \phi_2(B, C) \phi_3(C, D) \phi_4(D, E) \sum_A \phi_1(A, B)$$

↳ summation over pairwise factor

↳ result is a factor over single variable B → $\tau_1(B)$

↳ eliminated A

$$= \sum_D \sum_C \sum_B \phi_2(B, C) \phi_3(C, D) \phi_4(D, E) \tau_1(B)$$



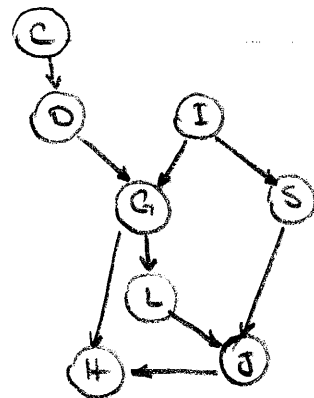
$$= \sum_D \sum_C \phi_3(C, D) \phi_4(D, E) \left(\sum_B \phi_2(B, C) \tau_1(B) \right)$$

↳ $\tau_2(C)$ → eliminated B

$$= \dots = \tau_5(E)$$

✓ Elimination in BN

$P(J) = ?$, eliminate C, D, I, G, L, S, H



$$\sum_{L, S, G, H, I, D, C} \phi_J(J, L, S) \phi_L(L, G) \phi_S(S, I) \phi_G(G, I, D) \phi_H(H, G, J) \phi_I(I) \phi_D(C, D) \phi_C(C)$$

$$\sum_{L, S, G, H, D, I} \phi_J(J, L, S) \phi_L(L, G) \phi_S(S, I) \phi_G(G, I, D) \phi_H(H, G, J) \phi_I(I) \underbrace{\sum_C \phi_D(C, D) \phi_C(C)}_{\tau_1(D)}$$

$$\sum_{L, S, G, H, D, I} \phi_J(J, L, S) \phi_L(L, G) \phi_S(S, I) \phi_G(G, I, D) \phi_H(H, G, J) \phi_I(I) \tau_1(D)$$

$$\sum_{L, S, G, H, I} \phi_J(J, L, S) \phi_L(L, G) \phi_S(S, I) \phi_H(H, G, J) \phi_I(I) \underbrace{\sum_D \phi_G(G, I, D) \tau_1(D)}_{\tau_2(G, I)}$$

$$\sum_{L, S, G, H} \phi_J(J, L, S) \phi_L(L, G) \phi_H(H, G, J) \sum_I \phi_S(S, I) \phi_I(I) \tau_2(G, I)$$

$$\sum_{L, S, G} \phi_J(J, L, S) \phi_L(L, G) \tau_3(G, S) \underbrace{\sum_H \phi_H(H, G, J)}_{\tau_4(G, J)}$$

$\phi_H(H, G, J) = P(H|G, J) \rightarrow \sum_H \dots = 1$
not every algorithm is clever enough to do this

$$\sum_{L, S} \phi_J(J, L, S) \sum_G \phi_L(L, G) \tau_3(G, S) \tau_4(G, J)$$

$$\sum_{L, S} \phi_J(J, L, S) \cdot \tau_5(L, S, J) \rightarrow \text{two factors} \rightarrow \text{multiply them}$$

Variable Elimination with Evidence

$P(J, I=i, H=h)$? \rightarrow eliminate C, D, G, S, L (no H , no I) \rightarrow have single value
 reduce the factors to correspond to scope

$$\sum_{L, S, G, D, C} \phi_J(J, L, S) \phi_L(L, G) \underbrace{\phi'_S(S) \phi'_G(G, D)}_{P(G|D, I=i)} \phi'_H(H, J) \underbrace{\phi'_I(I)}_{\phi_I(i) = cte} \phi_D(C, D) \phi_C(C)$$

elimination as before

renormalize it to have $P(J, I=i, H=h)$ by dividing it to $P(I=i, H=h)$

Variable Elimination in MN

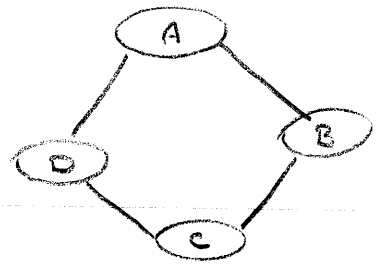
$P(D)$? \rightarrow eliminate A, B, C

$$\sum_{A, B, C} \phi_1(A, B) \phi_2(B, C) \phi_3(C, D) \phi_4(A, D)$$

$$\sum_{BC} \phi_2(B, C) \phi_3(C, D) \underbrace{\sum_A \phi_1(A, B) \phi_4(A, D)}_{T_1(B, D)}$$

$$\sum_B T_1(B, D) \underbrace{\sum_C \phi_2(B, C) \phi_3(C, D)}_{T_2(B, D)}$$

$$\sum_B T_1(B, D) T_2(B, D) = T_3(D) = \tilde{P}(D) \rightarrow \text{renormalize}$$



Routine: Eliminate var Z from Φ

- look in Φ , and find Φ' which is all factors that involve Z
- multiply all of those factors
- sum out the variable we want to eliminate, Z .
- the factors in Φ' is used and we don't want to reuse them, but we create T so introduce it.

$$\Phi' = \{ \phi_i \in \Phi : Z \in \text{Scope}[\phi_i] \}$$

$$\psi = \prod_{\phi_i \in \Phi'} \phi_i$$

$$T = \sum_Z \psi$$

$$\Phi = \Phi - \Phi' \cup \{T\}$$

Variable Elimination Algorithm Summary

- Reduce all factors by evidence
 - \rightarrow remove inconsistent rows of CPD
 - \rightarrow get us a set of factors Φ
- For each non-query variable Z
 - \rightarrow Eliminate var Z from $\Phi \rightarrow$ removes $\Phi' \Rightarrow$ adds T
- Multiply all remaining factors
- Renormalize to get a distribution

• VE Summary

- Simple algorithm
- Works for both BN and MN
- Factor product and summation steps can be done in any order, subject to:
 - ↳ when Z is eliminated, all factors involving Z have been multiplied in → ensure correctness

• Complexity Analysis of VE

- Basic operations in eliminating Z

↳ factor product

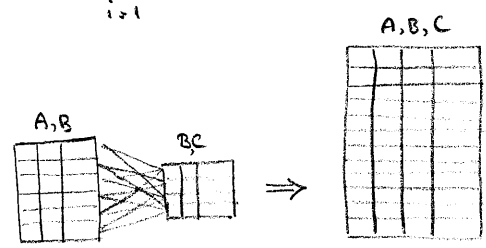
* every row of result = one row of each table

$$\Psi_k(X_k) = \prod_{i=1}^{m_k} \phi_i$$

* $N_k = |\text{Val}(X_k)|$ → row of results

↳ need to multiply m_k different factors

↳ $m_k - 1$ product



* total cost : $(m_k - 1) N_k$ multiplications

↳ factor marginalizations

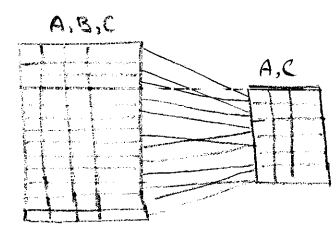
* every row of result = sum of all values of marginalized variable

$$\tau_k(X_k - \{Z\}) = \sum_Z \Psi_k(X_k)$$

* each number of original table is used exactly once

* $N_k = |\text{Val}(X_k)|$

* total cost : $\sim N_k$ additions → upper bound



- Complexity of VE

↳ start with m factors

↳ $m \leq n$ for BN

↳ $m = n$ reduced by evidence
↳ one factor for every variable = CPD

↳ can be larger for MN → grid MN, fully connected MN

↳ at each elimination step

↳ get in some factors

↳ generate exactly 1 factor

↳ we have at most n elimination steps

↳ worst case is when each elimination step use exactly 1 factor and generate 1

↳ total number of factors that we ever produce : m^*

↳ $m^* \leq$ original factors + generated factors = $m + n$

↳ size of the largest factor ever created N

1. $n_1 - \max(n_1)$

↳ full calculation

↳ product operations : $\sum_k (m_k - 1) N_k$

each factor is multiplied in at most once

↳ as soon as we multiply it in, it goes away

↳ $\sum_k (m_k - 1) \ll \# \text{ of factors} = m^*$

$\sum_k (m_k - 1) N_k \leq N \sum_k (m_k - 1) \leq \underline{\underline{N \cdot m^*}}$

↳ sum operations : $\leq \sum_k N_k \leq N (\# \text{ of elimination steps}) \leq \underline{\underline{N \cdot n}}$

↳ total work is linear in N and m*

↳ $N_k = |\text{Val}(X_k)| = O(d^{r_k}) \rightarrow$ total number of values in a factor

↳ $d = \max(|\text{Val}(X_i)|) \rightarrow$ all values have maximum d values in their scope
binary $\Rightarrow d=2$

↳ $r_k = |X_k| \rightarrow$ cardinality of the scope of the kth factor

↳ the source of exponential blow up

Complexity Example

$\tau_1(D) = \sum_c \phi_c(c) \phi_D(c, D)$

$\tau_2(G, I) = \sum_D \phi_G(G, I, D) \tau_1(D)$

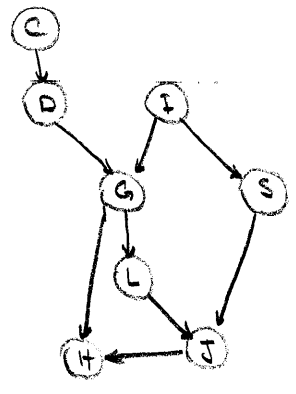
$\tau_3(S, G) = \sum_I \phi_S(S, I) \phi_I(I) \tau_2(G, I)$

$\tau_4(G, J) = \sum_H \phi_H(H, G, J)$

$\tau_5(J, L, S) = \sum_G \phi_L(L, G) \tau_3(S, G) \tau_4(G, J)$

$\tau_6(J) = \sum_{L, S} \phi(J, L, S) \tau_5(J, L, S)$

	# of vars in scope
c, D	2
G, I, D	3
S, I, G	3
H, G, J	3
L, G, S, J	4
J, L, S	3



Elimination Order

↳ variables can be eliminated in any order \rightarrow if we multiply things correctly and correct time
↳ order can affect the complexity

✓ if we start elimination by G!

$$\sum \phi_J \phi_L \phi_S \phi_G \phi_H \phi_I \phi_D \phi_C$$

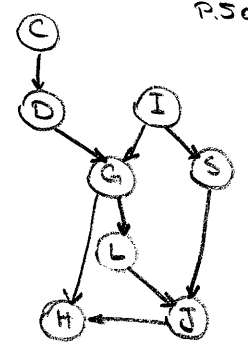
eliminate G

↳ which factors do we need to multiply to eliminate G?

$$\sum_G \phi_L(L,G) \times \phi_G(G,I,D) \times \phi_H(H,G,I)$$

$$\psi_1(L,G,I,D,H,J)$$

6 variables → previously d=4
 ↳ this order d ≥ 6



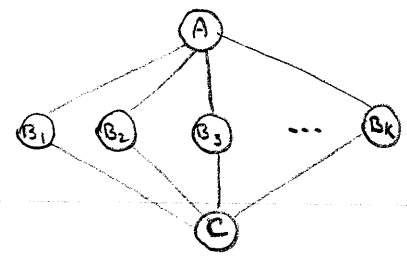
✓

Eliminate A first

↳ factors involved $\phi_1(A,B_1) \dots \phi_k(A,B_k)$

↳ scope of factor generated $\{A, B_1, \dots, B_k\} \rightarrow k+1$

↳ size of the factor is exponential in k



Eliminate B_i first

↳ factors involved $\phi_{A_i}(A, B_i) \times \phi_{C_i}(C, B_i) \rightarrow \psi_i(A, B_i, C) \rightarrow \tau_i(A, C)$

↳ scope of factor generated $\{A, B_i, C\} \rightarrow 3$

↳ size of the factor is exponential in 3

Summary

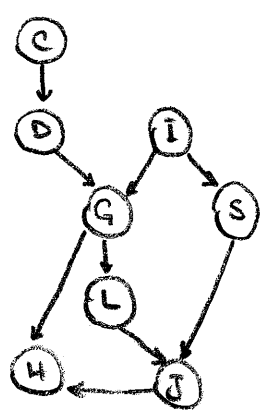
- ↳ complexity of variable elimination linear in
 - * size of model (# factors, # variables)
 - * size of largest factor generated

↳ size of factor is exponential in its scope

↳ complexity of algorithm depends heavily on elimination ordering

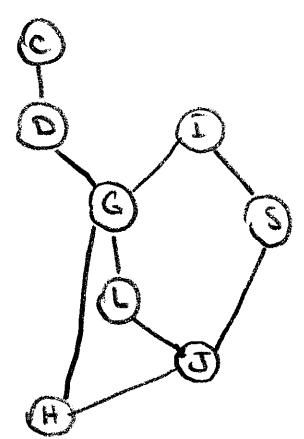
Graph Based Perspective

✓



$$\phi_J(J,L,S) \times \phi_L(L,G) \times \phi_S(S,I) \times \phi_G(G,I,D) \times \phi_H(H,G,J) \times \phi_I(I) \times \phi_D(C,D) \times \phi_C(C)$$

structure of graph that corresponds to this factors → make all edges undirected



that's not enough!

↳ $\phi_G(G, I, D)$ so we need to add an edge $D-I$ to represent they are together in a factor

↳ $\phi_J(J, L, S) \rightarrow L-S \rightarrow J$'s parents

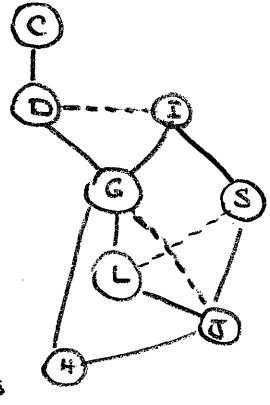
↳ $\phi_H(H, G, J) \rightarrow G-J \rightarrow H$'s parents

↳ for any v -structure we need to add an edge between parents

↳ called moralization: marrying the parents of child

↳ marry all of them no matter how many of them

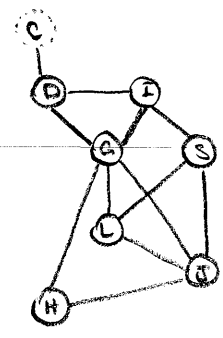
↳ the result is induced MN for the current set of factors



*eliminate: C

$$\tau_1(D) = \sum_C \phi_C(C) \phi_D(C, D)$$

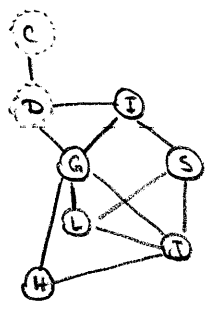
$$\phi_J(J, L, S) \phi_L(L, G) \phi_S(S, I) \phi_G(G, I, D) \phi_H(H, G, J) \tau_1(D)$$



*eliminate: D

$$\tau_2(G, I) = \sum_D \phi_G(G, I, D) \tau_1(D)$$

$$\phi_J(J, L, S) \phi_L(L, G) \phi_S(S, I) \phi_I(I) \phi_H(H, G, J) \tau_2(G, I)$$



*eliminate: I

$$\tau_3(S, G) = \sum_I \phi_S(S, I) \phi_I(I) \tau_2(G, I)$$

↳ the factor that we introduced $\tau_3(S, G)$ doesn't have an edge in the graph

↳ we should introduce $G-S$ edge in graph

↳ called fill edge: an edge we have to introduce because of elimination order

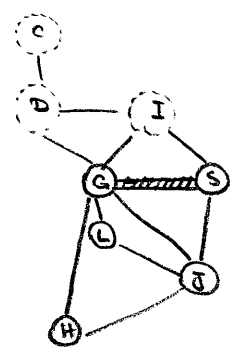
↳ all variables connected to I become connected directly

$$\phi_J(J, L, S) \phi_L(L, G) \phi_H(H, G, J) \tau_3(S, G)$$

*eliminate: H

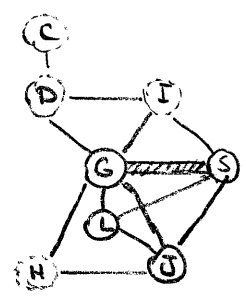
$$\tau_4(G, J) = \sum_H \phi_H(H, G, J)$$

$$\phi_J(J, L, S) \phi_L(L, G) \tau_3(S, G) \tau_4(G, J)$$



*eliminate: G

*eliminate: L & S



this VE order only introduce only 1 fill edge

↳ conservative → pretty good elimination ordering

we can put all of those edges in a graph and call it Induced Graph.

- Induced Graph

induced graph $I_{\Phi, \alpha}$ over factors Φ and ordering α

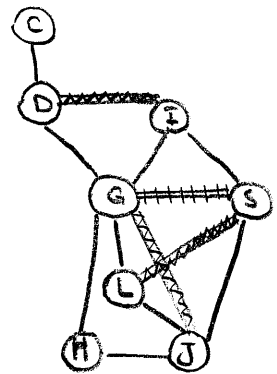
↳ undirected graph

↳ X_i and X_j are connected if they appear in the same factor in a run of VE algorithm using α as the ordering

↳ — begin with

▨ moralization

▨ fill edge



this ordering has 1 fill edge

- Cliques in the Induced tree

* Theorem: Every factor produced during VE is a clique in the induced graph

↳ clique: maximal fully connected subgraph e.g. DIG ✓

cliques

generated factors

GLJ \times → fully connected

↳ not maximal
↳ we can add S

DIG → ②

$$\textcircled{1} \tau_1(D) = \sum_c \phi_c \phi_D$$

GIS → ③

$$\textcircled{2} \tau_2(G, I) = \sum_D \phi_G \tau_1$$

GLSJ → ⑤

$$\textcircled{3} \tau_3(S, G) = \sum_I \phi_S \phi_I \tau_2$$

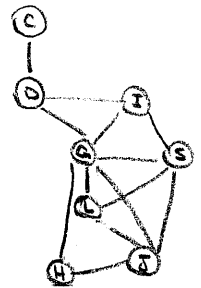
CD → ①

$$\textcircled{4} \tau_4(G, J) = \sum_H \phi_H$$

GHJ → ④

$$\textcircled{5} \tau_5(L, J, S) = \sum_G \phi_L \tau_3 \tau_4$$

$$\textcircled{6} \tau_6(J) = \sum_{L, S} \phi_J \tau_5 \longrightarrow \text{no clique}$$



* Theorem: Every (maximal) clique in the induced graph is a factor produced during VE

e.g. GSLJ as a max clique

↳ one of the variables has to be the one who is eliminated first

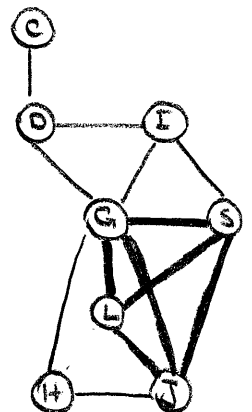
↳ once it eliminated it remove all corresponding edges

↳ so no new neighbors adds to it later

⇒ when eliminated it had all the neighbors that it had in clique.

⇒ it participated in factors with all other variables

⇒ when multiplied together, we have a factor over all of them



- Induced width

↳ the number of nodes in the largest clique in the graph - 1

↳ Minimal Induced width of a graph k is $\min(\text{width}(I_{k, \alpha}))$

↳ Provides a lower bound on best performance of VE to a model factorizing over K

- Summary

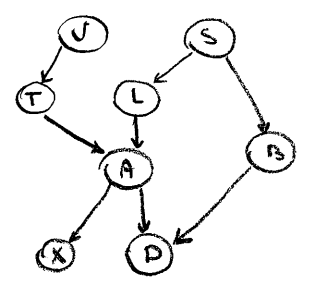
↳ VE can be viewed as transformations on undirected graph

↳ elimination connects all node's current neighbors \Rightarrow gives induced graph

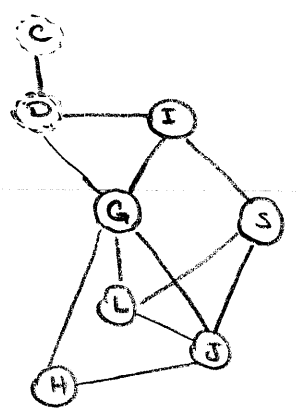
↳ cliques in resulting induced graph directly correspond to alg. complexity

⊠ How many edges should be added to the graph to make it moralized?

= # of v-structures = 2



⊠ what edge if any should be added to the graph when I is eliminated?



$$I \rightarrow \Phi(I, G, S) \xrightarrow{\text{elimination}} T(G, S) \rightarrow G-S$$

• Finding Elimination Orderings

- How good of an elimination ordering can we construct?

Theorem: For a graph H, determining whether there exists an elimination ordering for H with induced width $\leq K$ is NP-complete

↳ This NP-hardness result is distinct from the NP-hardness result of inference.

↳ even given the optimal ordering, inference may still be exponential

- Greedy search using heuristic cost function

↳ at each point: eliminate node with smallest cost

↳ possible cost functions

↳ min-neighbors: # neighbors in current graph \rightarrow produces smallest factor

↳ min-weight: weight (# values) of factors formed

⊠ 2 var with 100 values each

5 var with binary values ✓

↳ these two look at the size of factor \rightarrow ignores the fact that some of them are inevitable

↳ min-fill: number of new fill edges \rightarrow counting add complexity introduced by eliminating

↳ weighted min-fill: total weight of new fill edge

↳ edge weight = product of weights of the two nodes

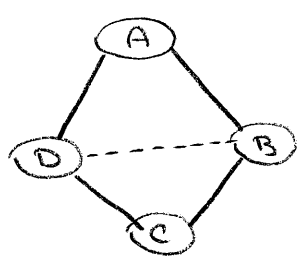
Practice

- Theorem: The induced graph is triangulated

↳ regardless of elimination ordering

↳ triangulated: no loop of length > 3 with a "bridge"

□ proof via example



assume we have a loop of size $4 > 3$

assume that C is the first to be eliminated

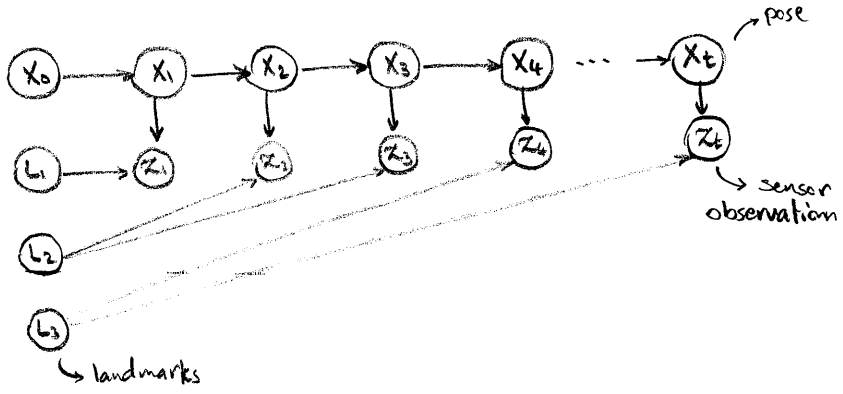
↳ after eliminating C we end up introducing B-D edge
↳ fill edge

↳ we can find elimination ordering by finding a low-width triangulation of original graph H_{Φ}

↳ using graph triangulation literature

□ Robot Localization & Mapping

robot sees landmark and distinguish them uniquely → can measure the distance to each landmark

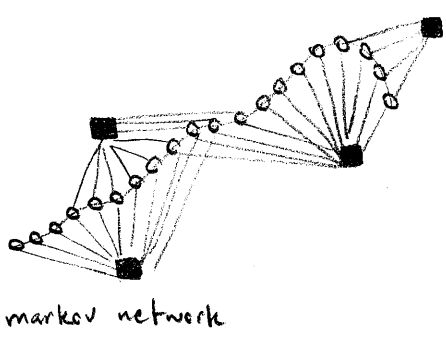
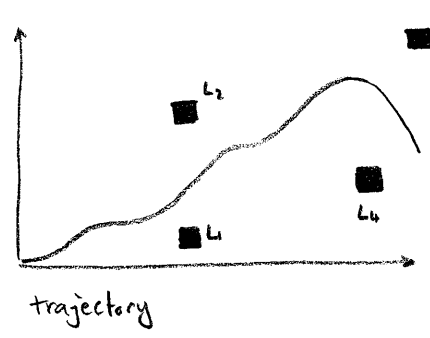


@ t_1 : robot sees L_1

@ t_2 and t_3 : robot sees L_2

@ t_4 : robot sees L_3

observed distance is a function of robot pose and landmark distance



- robot pose at every time
- edge between $p_t - p_{t+1}$
- landmark positions

eliminate robot poses then landmarks? → induced graph big spaghetti:

↳ every landmark connected to all other landmarks

eliminate landmarks then poses? → still densely connected

↳ but not as bad as the spaghetti before

min-fill elimination → very sparse

- Summary

↳ finding the optimal elimination ordering is NP-hard

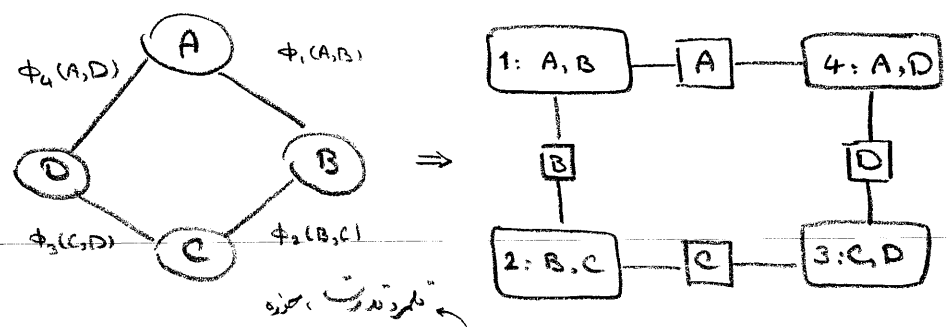
↳ simple heuristics that try to keep induced graph small often provide reasonable performance

* Message Passing :: Belief Propagation

- Message Passing class of algorithms
 - ↳ closely related to VE
 - ↳ more flexible → how to do summation and factor product steps
 - ↳ lower complexity compared even to minimal VE ordering

• Cluster Graph

- a data structure to get bits of knowledge from a PGM

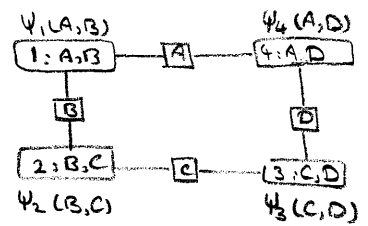


نكرتة سيرة

- cluster 1: a cluster whose jurisdiction of influence is pair of variables A, B
- clusters are going to talk to each other → convince each other that what they think a var. that they both have under jurisdiction is correct
 - ↳ cluster 1 talk to cluster 2 about B → might tell 2 sth about B that it become more informed about dist. over B

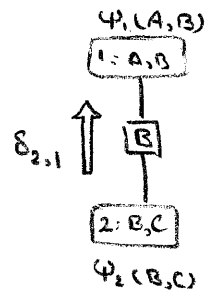
- $\phi_1(A, B) \rightarrow c1: A, B$
- $\phi_2(B, C) \rightarrow c2: B, C$
- ...

- variables are communicating via messages.
- name the initial set of beliefs (= evidence that a cluster has over vars. in jurisdiction) → ψ
 - ↳ can be more complicated than single ϕ 's.



- cluster 2 wants to send a msg to cluster 1

- ↳ initially we are going to initialize with a totally uninformed msg.
 - ↳ $\delta_{2,1} = 1$ → from 2 to 1
 - ↳ they haven't even start to talk with each other
 - ↳ cluster 1 get that msg and multiply it with current thoughts about var. A and B



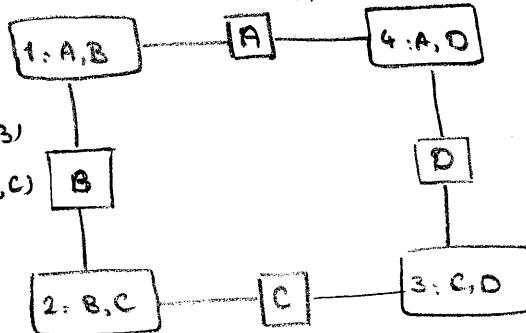
↳ now that cluster 1 get some info, it's gonna send it to c4.

↳ c4 doesn't care about B

$$\delta_{1,4} = \sum_B \delta_{2,1}(B) \psi_1(A, B) \rightarrow \sum \text{incoming msg} \times \text{initial beliefs}$$

$$\delta_{1,4} = \sum_B \delta_{2,1}(B) \psi_1(A, B)$$

$$\delta_{4,1} = \sum_D \delta_{3,4}(D) \psi_4(A, D)$$



$$\delta_{1,2} = \sum_A \delta_{4,1}(A) \psi_1(A, B)$$

$$\delta_{2,1} = \sum_C \delta_{3,2}(C) \psi_2(B, C)$$

$$\delta_{3,4} = \sum_C \delta_{2,3}(C) \psi_3(C, D)$$

$$\delta_{4,3} = \sum_A \delta_{1,4}(A) \psi_4(A, D)$$

$$\delta_{2,2} = \sum_B \delta_{1,2}(B) \psi_2(B, C)$$

$$\delta_{3,2} = \sum_D \delta_{4,3}(D) \psi_3(C, D)$$

- each cluster send messages to its adjacent cluster that reflect this exact same process
 - ↳ the msg that c3 send to c4 doesn't take into consideration, the information that it got from c4 → *میرا پیغام، میرا پیغام! عجب برہنہ سنی!*
 - ↳ avoid repeating back a number that you just told me or else we all become more convinced about the truth of this number

↳ restrict our attention that come in from other sources. → communication protocol

• Cluste Graphs (ctd)

- Undirected graph such that

↳ nodes are clusters $C_i \subseteq \{X_1, \dots, X_n\}$ → subsets of variables

↳ edge between C_i and C_j associated with the sepset $S_{i,j} \subseteq C_i \cap C_j$ → variable they both know about
↳ the variables they choose to talk about

- Given set of factors Φ

Initialize the model by giving each cluster a certain amount of information

↳ each of initial factors ϕ_k is going to be assigned to one and only one cluster $C_{\alpha(k)}$

↳ it needs to be at least one → to be taken to account

↳ shouldn't be more than one → avoid double counting the evidence

Assign ϕ_i to $C_{\alpha(k)} \rightarrow \text{Scope}[\phi_k] \subseteq C_{\alpha(k)}$

↳ where do we put the information corresponding to factor k?

↳ we can only put it in a cluster that understands every single var. in factor

Define initial beliefs

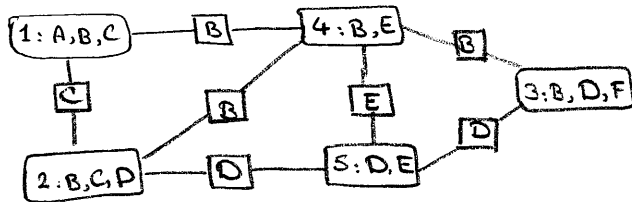
↳ product of all factors that are assigned to it

$$\psi_i(C_i) = \prod_{k: \alpha(k)=i} \phi_k$$

↳ some cluster have just one factor → $\psi = \phi$

some have several factors → multiply those factors → total informed belief of cluster

- $\Phi_1(A, B, C), \Phi_2(B, C), \Phi_3(B, D), \Phi_4(D, E),$
 $\Phi_5(B, E), \Phi_6(B, D), \Phi_7(B, D, F)$



find for each factor a cluster in which to put it:

- $\Phi_1(A, B, C) \rightarrow$ only one choice $\rightarrow C_1 \rightarrow \alpha(1) = 1$
- $\Phi_2(B, C) \rightarrow$ two choices: $C_1, C_2 \rightarrow$ let's choose $C_2 \rightarrow \alpha(2) = 2$
- $\Phi_3(B, D) \rightarrow$ two choices: $C_2, C_3 \rightarrow$ let's choose $C_2 \rightarrow \alpha(3) = 2$
- $\Phi_4(D, E) \rightarrow$ one choice $\rightarrow C_5 \rightarrow \alpha(4) = 5$
- $\Phi_5(B, E) \rightarrow$ one choice $\rightarrow C_4 \rightarrow \alpha(5) = 4$
- $\Phi_6(B, D) \rightarrow$ two choices: $C_2, C_3 \rightarrow$ let's choose $C_3 \rightarrow \alpha(6) = 3$
- $\Phi_7(B, D, F) \rightarrow$ only one choice $\rightarrow C_3 \rightarrow \alpha(7) = 3$

there are other alternative assignments

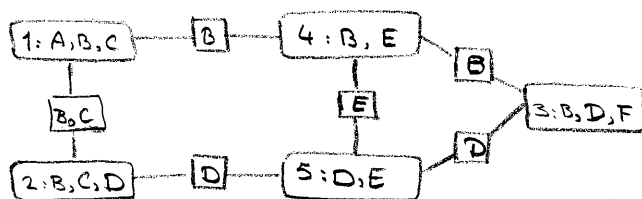
$\Psi_2 = \Phi_2(B, C) \times \Phi_3(B, D)$
 $\Psi_1 = \Phi_1(A, B, C)$
 $\Psi_3 = \Phi_6(B, D) \times \Phi_7(B, D, F)$

for alternative assignments:

$\Psi_1 = \Phi_1$
 $\Psi_2 = \Phi_2 \Phi_3 \Phi_6$
 $\Psi_3 = \Phi_7$
 $\Psi_4 = \Phi_5$
 $\Psi_5 = \Phi_4$

Same factors, different Cluster Graph

clusters haven't changed
 only sepsets between them changed



- e.g. $C_1: A, B, C$ & $C_2: B, C, D \rightarrow$ sepset above: C
 \rightarrow sepset now: B, C
- e.g. $C_2: B, C, D$ & $C_4: B, E \rightarrow$ sepset above: B
 \rightarrow sepset now: \emptyset

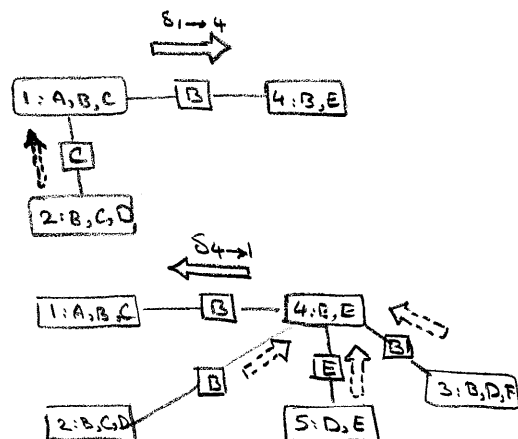
Message Passing

$$\delta_{1 \rightarrow 4}(B) = \sum_{A, C} \Psi_1(A, B, C) \delta_{2 \rightarrow 1}(C)$$

$$\delta_{4 \rightarrow 1}(B) = \sum_E \Psi_4(B, E) \times \delta_{2 \rightarrow 4}(B) \times \delta_{5 \rightarrow 4}(E) \times \delta_{3 \rightarrow 4}(B)$$

current most informed beliefs about B & E

\rightarrow message from C_1 itself was not used.



$$\delta_{i \rightarrow j}(S_{i,j}) = \sum_{C_i - S_{i,j}} \psi_i \times \prod_{k \in (\mathcal{N}_i - \{j\})} \delta_{k \rightarrow j}$$

from: i to: j sepset: $S_{i,j}$
 the vars cluster j doesn't know about
 factors initially assigned to cluster i
 all of the incoming messages other than the one from j

• Belief Propagation Algorithm

↳ clusters propagates informed beliefs together

– Algorithm

I. Assign each factor $\phi_k \in \Phi$ to a cluster $C_{\kappa(k)}$

II. Construct initial potentials $\psi_i(C_i) = \prod_{k: \kappa(k)=i} \phi_k$

III. Initialize all messages to be 1

VI. Repeat

(i) select edge (i,j) and pass message $\delta_{i \rightarrow j}(S_{i,j}) = \sum_{C_i - S_{i,j}} \psi_i \times \prod_{k \in (\mathcal{N}_i - \{j\})} \delta_{k \rightarrow i}$

V. Compute final belief $\beta_i(C_i) = \psi_i \times \sum_{k \in \mathcal{N}_i} \delta_{k \rightarrow i}$

↳ at the end of process a cluster needs to know what to believe about the variable that it understands

↳ takes its own initial beliefs and all of the information it got from all of neighbors

↳ produce beliefs (β_i)

– Repeat until when?

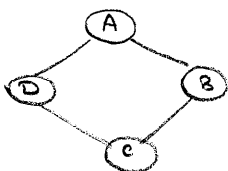
– How do we select edges to pass messages along?

↳ prespecified order? → round robin

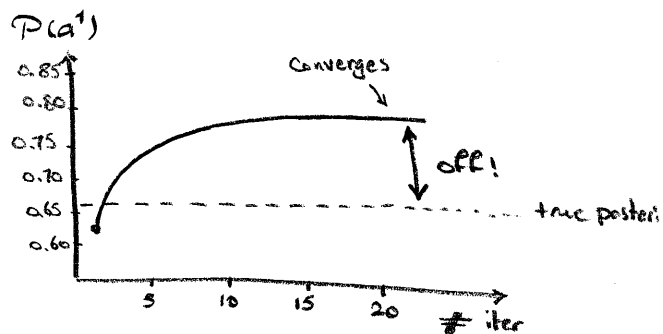
↳ there are better strategies

– Is this algorithm good? Yes & No → ! slips, C_i

□



convergence achieved
 but the result is a little bit off!
 ↳ the answer is not exact

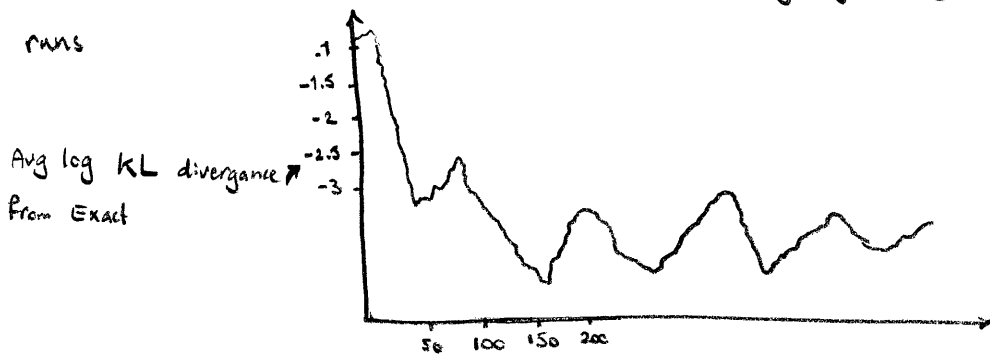


↳ In general this is an approximate algorithm → no free lunch (NP-hard)

↳ have exceptions

↳ no easy way to solve it

↳ different BP runs



- So why bother using?

↳ it misbehaves in loops

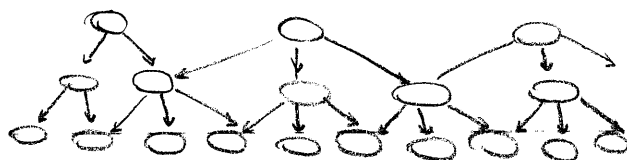
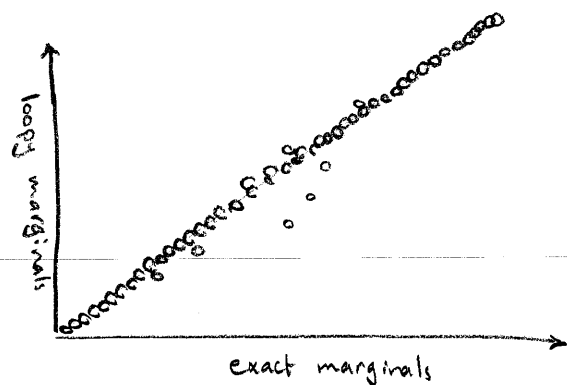
↳ pyramid network

↳ happens in image analysis

↳ compute exact marginals

compute loopy BP

with a few outliers the results are the same



- Summary

↳ Graph of clusters connected by sepsets

↳ Adjacent clusters pass information to each other about variables in sepsets

↳ The msg from i to j summarizes everything i knows, except information obtained from j .

↳ Algorithm may not converge

↳ may have oscillatory behavior

↳ The resulting beliefs are pseudo marginals → in theory

↳ but actually performs well in practical applications

• Cluster graph properties

- In order to → support reasonable message passing

Reminder: Undirected graph → nodes: clusters $C_i \subseteq \{X_1, \dots, X_n\}$
 ↘ edges: sepsets $S_{ij} \subseteq C_i \cap C_j$

- Family Preservation Property

↳ given a set of factors Φ

we need to be able to assign each factor ϕ_k to some cluster $C_{\alpha(k)}$ such that the cluster can accommodate ϕ_k ($\text{Scope}[\phi_k] \subseteq C_{\alpha(k)}$)

↳ we need an appropriate cluster in cluster graph

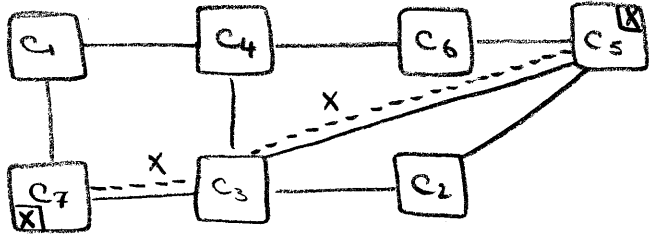
↳ for each factor $\phi_k \in \Phi$, there exists a cluster C_i such that $\text{Scope}[\phi_k] \subseteq C_i$

- Running Intersection Property

↳ let's assume that we have a pair of clusters C_i and C_j , and a variable X that belongs to both

e.g. X in C_5 & C_7

the property says that there exist a unique path between C_5 & C_7 for which all clusters and sepsets that is along the path contains X



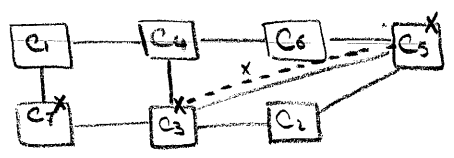
e.g. if the path is $C_7 - C_3 - C_5$ there should be X all along the way

↳ For each pair of clusters C_i, C_j and variable $X \in C_i \cap C_j$ there exist a unique path between C_i and C_j for which all clusters and sepsets contain X

↳ existence? uniqueness?

↳ existence proof: assume that there is no path \Rightarrow no way

for C_7 to communicate C_5 about the variable X . \Rightarrow separated communities each know some information about X and they can never talk to each other about $X \Rightarrow$ never going to get agree about X



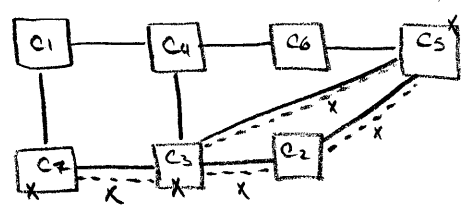
↳ uniqueness proof: assume that there are two paths involved X

↳ trace msg passing algorithm

$C_3 \xrightarrow{X} C_5$ eg. I suggest $X = \alpha'$

$C_5 \xrightarrow{\delta_{55}} C_2$

$C_2 \xrightarrow{\delta_{53}} C_3$ the C_3 suggestion about $X = \alpha'$ reinforces \rightarrow prob $\uparrow \rightarrow$ again to loop $C_3 - C_5 - C_2$



↳ self-reinforcing loop

↳ skewed prob. in many examples

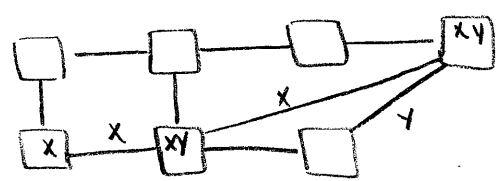
↳ we should avoid that \rightarrow only reduce the problem not eliminate it

✓ X and Y are strongly correlated

$C_3 \xrightarrow{X} C_5$

C_5 translates that to information about Y

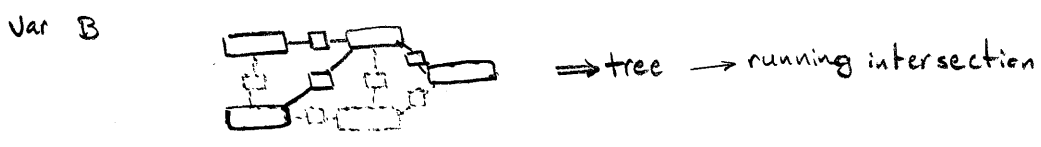
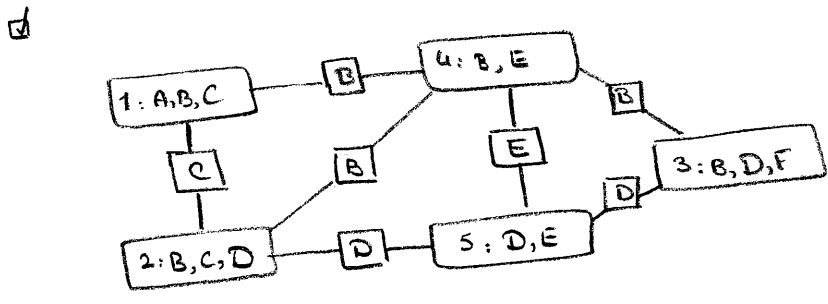
$C_5 \xrightarrow{Y} C_2 \xrightarrow{Y} C_3$ reinforced!



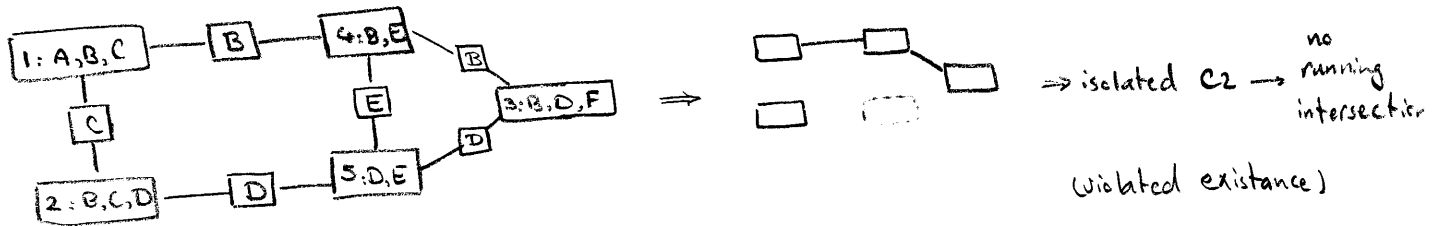
↳ BP does very poorly when we have strong correlations

↳ The more skew probabilities are, the harder for BP to get the results

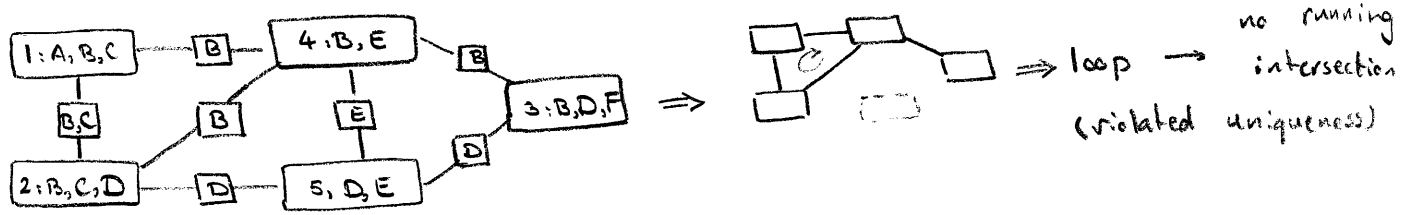
- ↳ equivalently: for any X , the set of clusters and sepsets containing X form a tree
- ↳ has to be connected = existence of the path
- ↳ tree, no loop = uniqueness
- ↳ each variable induce its own little tree across which info about that var. flows in the graph



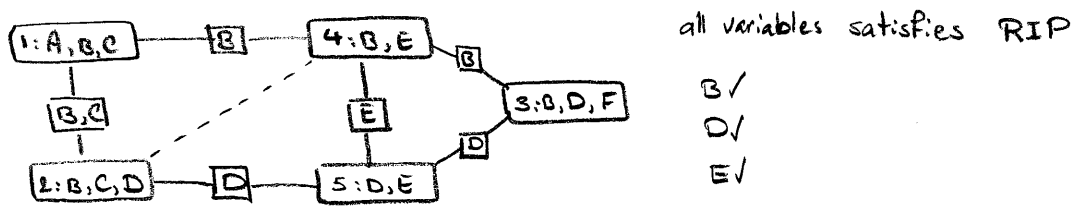
✓ Illegal cluster graph



✓ Illegal cluster graph



if we want to take the cluster graph and still allow communication between B & C
 ↳ e.g. want C_1 & C_2 to talk about B and C
 ↳ just eliminate one edge



- ↳ how do we construct a cluster graph that has desired properties?
- ↳ one very simple and in some ways degenerate version → Bethe Cluster graph

- Bethe Cluster Graph

- ↳ simple → often used
- ↳ from statistical physics: people use this kind of approximation to energy functions
- ↳ degenerate

↳ two types of clusters

↳ big clusters : correspond to factors in Φ

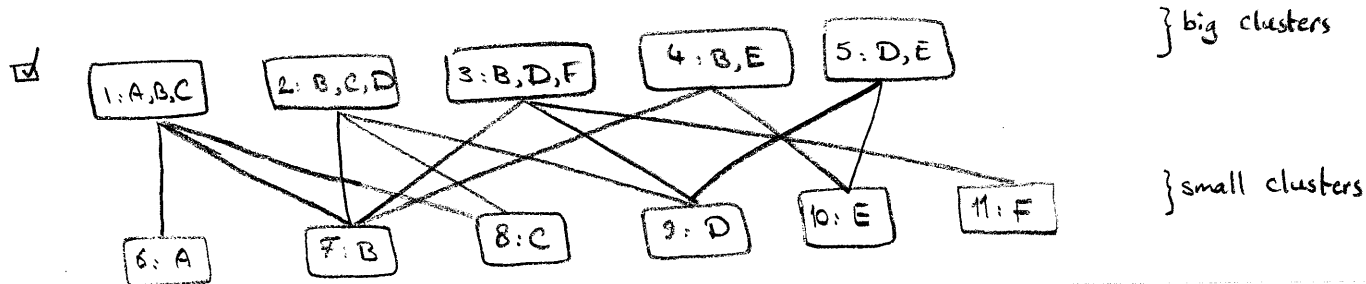
for each $\phi_k \in \Phi \Rightarrow$ a factor cluster $C_k = \text{Scope}[\phi_k]$

↳ little clusters : correspond to individual variables

for each $X_i \Rightarrow$ a singleton cluster

↳ connect C_k to X_i exactly when X_i is a member of C_k

edge $C_k - X_i$ if $X_i \in C_k$



∴ ABC is connected to A, B and C → how the msg's are passed

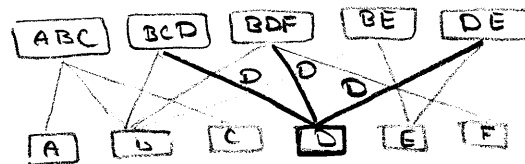
↳ degenerate : it only allows info about singletons to be passed

↳ in every msg passing step : loses any info about the correlation between variables

∴ simple to construct , guaranteed RIP

↳ why RIP is guaranteed in Bethe Cluster Graph

∴ in the above graph, any variable involve var D, sepsets have only D. It's a tree because D doesn't appear on any other sepsets except these ones



— Summary

↳ Cluster graph must satisfy two properties

↳ family preservation : allow Φ to be encoded

↳ running intersection : connects all info about any var, but without feedback loops

↳ Bethe Cluster Graph is often first default

↳ Richer cluster graph structure can offer different tradeoffs wrt computational cost and preservation of dependencies

↳ the amount of the sizes of the messages that are passed grow more expensive BUT allow us to improve the preservation of the dependencies as msg's are passed

↳ keep more info

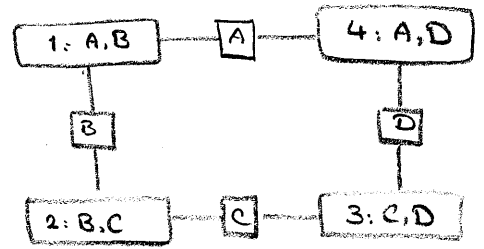
• Properties of BP Algorithm

- Calibration

Cluster beliefs

$$\beta_i(C_i) = \psi_i \times \prod_{k \in \mathcal{N}_i; \delta_{k \rightarrow i}}$$

$$\beta_1(A,B) = \psi_1(A,B) \times \delta_{4 \rightarrow 1}(A) \times \delta_{2 \rightarrow 1}(B)$$



↳ a cluster graph is calibrated if every pair of adjacent clusters C_i, C_j agree on their sepset $S_{i,j}$

↳ if we ask C_1 what it thinks about var B : $\beta_1(A, B)$

if we ask C_2 what it thinks about var B : $\beta_2(B, C)$

they should agree $\rightarrow \sum_A \beta_1(A, B) = \sum_C \beta_2(B, C)$

$$\sum_{C_i - S_{i,j}} \beta_i(C_i) = \sum_{C_j - S_{i,j}} \beta_j(C_j)$$

↳ convergence of BP algorithm implies calibration.

Convergence \blacktriangleright Calibration

Convergence of BP: msg from next time step equals to previous one.

$$\delta_{i \rightarrow j}(S_{i,j}) = \delta'_{i \rightarrow j}(S_{i,j}) \quad (1)$$

$$\delta'_{i \rightarrow j}(S_{i,j}) = \sum_{C_i - S_{i,j}} \left(\psi_i \times \prod_{k \in (\mathcal{N}_i - \{j\})} \delta_{k \rightarrow i} \right) \quad (2)$$

$$\beta_i(C_i) = \psi_i \times \prod_{k \in \mathcal{N}_i} \delta_{k \rightarrow i} \quad (3)$$

↳ all msgs

$$(2,3) \rightarrow \delta'_{i \rightarrow j}(S_{i,j}) = \sum_{C_i - S_{i,j}} \frac{\beta_i(C_i)}{\delta_{j \rightarrow i}(S_{i,j})} = \delta_{i \rightarrow j}(S_{i,j}) \quad (4)$$

$$\delta_{i \rightarrow j}(S_{i,j}) \times \delta_{j \rightarrow i}(S_{i,j}) = \sum_{C_i - S_{i,j}} \beta_i(C_i) \quad (5)$$

equally $\rightarrow \delta_{i \rightarrow j}(S_{i,j}) \times \delta_{j \rightarrow i}(S_{i,j}) = \sum_{C_j - S_{i,j}} \beta_j(C_j) \quad (6)$

$$(4,5) \rightarrow \sum_{C_i - S_{i,j}} \beta_i(C_i) = \sum_{C_j - S_{i,j}} \beta_j(C_j) \quad \checkmark$$

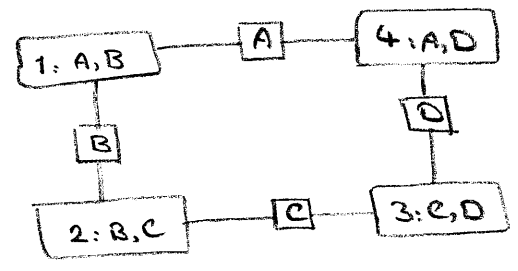
↳ sepset beliefs

$$\mu_{i,j}(S_{i,j}) = \delta_{j \rightarrow i} \delta_{i \rightarrow j} = \sum_{C_j \ni S_{i,j}} \beta_j(C_j)$$

- Reparameterization

$$\beta_i(C_i) = \psi_i \times \sum_{k \in \mathcal{N}_i} \delta_{k \rightarrow i}$$

$$\mu_{i,j}(S_{i,j}) = \delta_{j \rightarrow i} \delta_{i \rightarrow j}$$



$$\begin{aligned} \beta_1(A,B) &= \psi_1 \times \delta_{4 \rightarrow 1} \times \delta_{2 \rightarrow 1} & \mu_{1,2} &= \delta_{1 \rightarrow 2} \times \delta_{2 \rightarrow 1} \\ \beta_2(B,C) &= \psi_2 \times \delta_{1 \rightarrow 2} \times \delta_{3 \rightarrow 2} & \mu_{2,3} &= \delta_{2 \rightarrow 3} \times \delta_{3 \rightarrow 2} \\ \beta_3(C,D) &= \psi_3 \times \delta_{2 \rightarrow 3} \times \delta_{4 \rightarrow 3} & \mu_{3,4} &= \delta_{3 \rightarrow 4} \times \delta_{4 \rightarrow 3} \\ \beta_4(A,D) &= \psi_4 \times \delta_{1 \rightarrow 4} \times \delta_{3 \rightarrow 4} & \mu_{1,4} &= \delta_{1 \rightarrow 4} \times \delta_{4 \rightarrow 1} \end{aligned}$$

↳ lots of repetition $\rightarrow \delta_{4 \rightarrow 1}$ in β_1 and $\mu_{1,4}$
 $\searrow \delta_{1 \rightarrow 2}$ in β_2 and $\mu_{1,2}$
 \vdots

every δ appears exactly twice
 one in cluster belief
 one in sepset belief

$$\frac{\prod_i \beta_i}{\prod_{i,j} \mu_{i,j}} = \frac{\prod_i \psi_i \prod_{j \in \mathcal{N}_i} \delta_{i \rightarrow j}}{\prod_{i,j} \delta_{i \rightarrow j}} = \prod_i \psi_i = \tilde{P}_{\Phi}(X_1, \dots, X_n)$$

↳ $\tilde{P}_{\Phi}(X_1, \dots, X_n) = \frac{\prod \beta_i}{\prod_{i,j} \mu_{i,j}}$ \rightarrow shown with different sets of parameters

↳ no information loss as the result of BP

- Summary

- ↳ At convergence of BP, cluster graph beliefs are calibrated:
 - ↳ beliefs at adjacent clusters agree on sepsets
- ↳ Cluster graph beliefs are an alternative, calibrated parameterization of the original unnormalized density
 - ↳ we can read a variable from any clique in which it appears
- ↳ no info is lost by msg passing

☑ T or F? at convergence:

- the beliefs of neighboring clusters agree on all variables \rightarrow F
- the beliefs of neighboring clusters agree on all variables in the sepset \rightarrow T
- " " " " " " " " they have in common \rightarrow F
- " " " " " " " " in either cluster \rightarrow F

* Message Passing :: Clique Tree Algorithm

• BP → over a fairly general DS called cluster graph

Clique Tree Alg → over a special DS called Clique Tree

↳ special case of BP

↳ clique tree is a special case of cluster graph

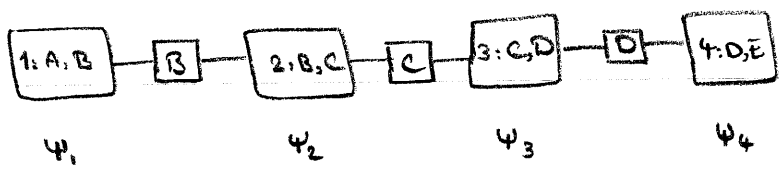
↳ brings stronger performance guarantee to BP

↳ faster, convergence to the exact answer 100%

✓ Passing Msg in a tree



pairwise MRF



pass msg from the left:

$$\delta_{1 \rightarrow 2}(B) = \sum_A \psi_1$$

$$\delta_{2 \rightarrow 3}(C) = \sum_B \psi_2 \times \delta_{1 \rightarrow 2}$$

$$\delta_{3 \rightarrow 4}(D) = \sum_C \psi_3 \times \delta_{2 \rightarrow 3}$$

pass on the other direction:

$$\delta_{4 \rightarrow 3}(D) = \sum_E \psi_4$$

$$\delta_{3 \rightarrow 2}(C) = \sum_D \psi_3 \times \delta_{4 \rightarrow 3}$$

$$\delta_{2 \rightarrow 1}(B) = \sum_C \psi_2 \times \delta_{3 \rightarrow 2}$$

regardless of when the msg get passed there is no other msg that C4 would multiply in when sending msg to C3.

↳ no msg recieved

↳ got a msg from C3 itself

• Correctness

- Surprisingly relevant to what we've seen before the resulting beliefs are correct

$$\beta_3(C, D) = \psi_3 \times \delta_{2 \rightarrow 3} \times \delta_{4 \rightarrow 3}$$

$$= \psi_3 \times \left(\sum_B \psi_2 \times \delta_{1 \rightarrow 2} \right) \times \left(\sum_E \psi_4 \right)$$

$$= \psi_3 \times \left(\sum_B \psi_2 \times \sum_A \psi_1 \right) \times \left(\sum_E \psi_4 \right)$$

we have a product of all the network $\psi_1 \dots \psi_4$ + we have a summation over all vars not in c

⇒ product of factors marginalized out unnecessary variables ~ same as VE

↳ legal order of operations

↳ correct result

• Clique tree

— Undirected tree

↳ nodes : clusters $C_i \subseteq \{X_1, \dots, X_n\}$

↳ edges : sepsets $C_i - C_j \Rightarrow S_{i,j} = C_i \cap C_j$

⦿ in cluster graphs sepsets could be a subset of $C_i \cap C_j$
 but in clique tree sepset should be equal to $C_i \cap C_j$

— Family Preservation

↳ Given set of factors Φ

assign each ϕ_k to a cluster $C_{\alpha(k)}$ such that $\text{Scope}[\phi_k] \subseteq C_{\alpha(k)}$

↳ each factor should go to a cluster → each piece of information should be represented ^{somewhere}

↳ For each factor $\phi_k \in \Phi$

there exist a cluster C_i such that $\text{Scope}[\phi_k] \subseteq C_i$

↳ there must exist one cluster whose scope is big enough to accommodate that factor

— Running Intersection Property

↳ For each pair of clusters C_i, C_j ↗ as variants of cluster graph

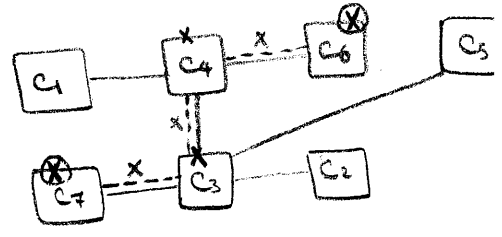
and variable $X \in C_i \cap C_j$

there exist a unique path between C_i and C_j
 for which all clusters and sepsets contain X

✓ C_6 and C_7

it's a tree → only one path between C_6 and C_7

so X has to be all over the path

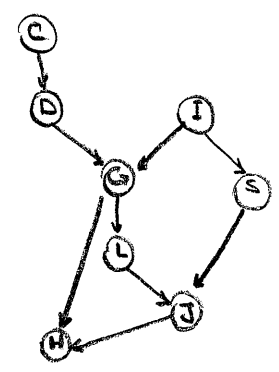
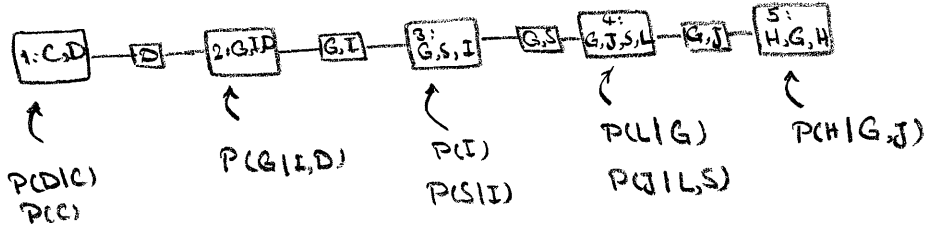


↳ (Clique tree version)

For each pair of clusters C_i, C_j
 and variable $X \in C_i \cap C_j$

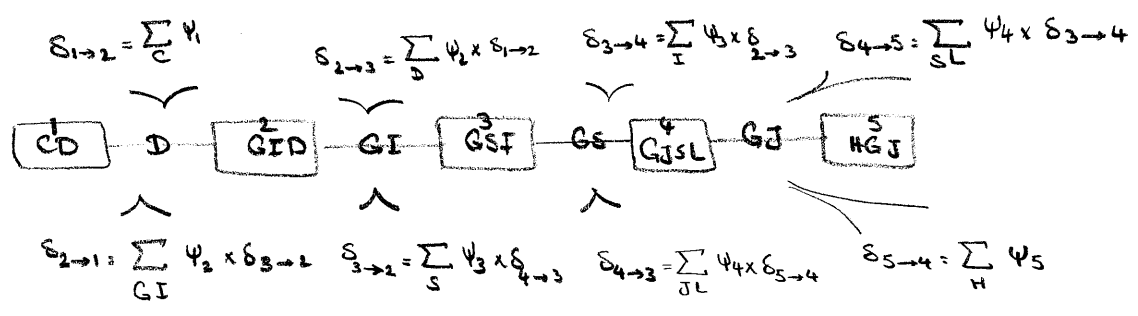
in the unique path between C_i and C_j
 all clusters and sepsets contain X

✓ more complicated clique tree



because of Family preservation property \rightarrow each factor can be assigned to some clique

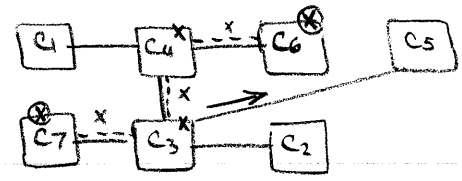
tree satisfies RIP $\rightarrow G : C_2 \sim C_5$



- RIP \Rightarrow Clique Tree Correctness

\hookrightarrow we want to prove that if we eliminated X while passing msg from C_i to C_j , then we have multiplied in all factors that involve X .

\square let's imagine that X is eliminated while passing msg from C_3 to C_5 .



X cannot be in C_5 , because it is not in the sepset $S_{3,5}$. \rightarrow if $X \in C_5 \xrightarrow{\text{RIP}} X \in S_{3,5} \Rightarrow X$ wouldn't be eliminated

\hookrightarrow so the RIP relates the structure of the graph to the time that those variable are eliminated.

$$\begin{aligned} \beta_3(G, S, I) &= \psi_3 \times \delta_{2 \rightarrow 3} \times \delta_{4 \rightarrow 3} \\ &= \psi_3 \times \sum_C \psi_1 \times \sum_D \psi_2 \times \sum_{J,L} \psi_4 \times \sum_H \psi_5 \\ &= \sum_{C,D,J,L,H} \psi_1 \times \psi_2 \times \psi_3 \times \psi_4 \times \psi_5 \end{aligned}$$

\hookrightarrow when we eliminate one variable, we multiply all of factors that involve the variable

- Summary

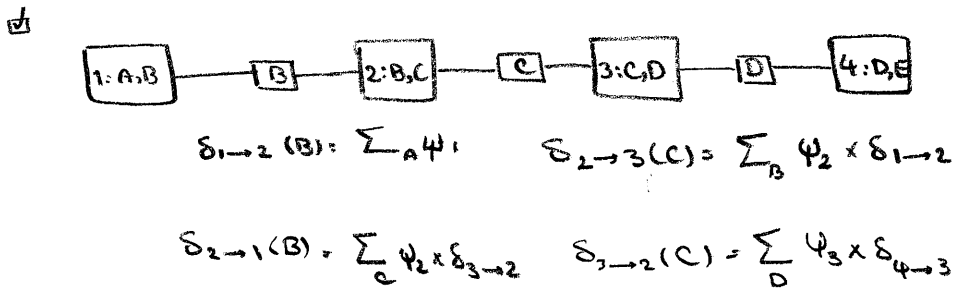
\hookrightarrow BP can be run over a tree-structured cluster graph

\hookrightarrow In this case, computation is a variant of variable elimination

\hookrightarrow The resulting beliefs are guaranteed to be correct marginals of the unnormalized density

• Computation of Clique Tree Alg

\hookrightarrow exploit tree structure to improve behavior of alg



$$\begin{aligned} \delta_{3 \rightarrow 4}(D) &= \sum_C \psi_3 \times \delta_{2 \rightarrow 3} \\ \delta_{4 \rightarrow 3}(D) &= \sum_E \psi_4 \end{aligned}$$

↳ observation

$\delta_{1 \rightarrow 2}$: once computed never changes \rightarrow instant convergence

$\delta_{2 \rightarrow 3}$: depend one when it is passed

↳ if passed first before C_2 get msg from C_1 then after getting that msg changes

↳ if C_2 is clever and it wait until $\delta_{1 \rightarrow 2}$ is passed it converges and never changes

$\delta_{3 \rightarrow 4}$: needs to wait longer

↳ if it wait for $\delta_{1 \rightarrow 2}$ then $\delta_{2 \rightarrow 3}$ now can fix $\delta_{3 \rightarrow 4}$ with one computation

From right-to-left it happens the same : $\delta_{4 \rightarrow 3}$ converge instantly then $\delta_{3 \rightarrow 2}$ then $\delta_{2 \rightarrow 1}$

\therefore so we can compute all the msg's in one pass in either direction ltr or rtl.

↳ for chains it has its own name : Forward-Backward Algorithm

↳ commonly used for HMMs

↳ to total effort that we had to compute all of these msg's is just one pass in each direction

↳ once all msg's have been passed we have beliefs that are marginals of $\tilde{P}_{\Phi}(C_i)$

↳ Convergence : once C_i received a final msg from all of its neighbors except C_j (if it waits enough to receive msg's from everyone else and then send it to j) then that msg $\delta_{i \rightarrow j}$ is final

↳ Can you always find a msg passing order to achieve that goal?

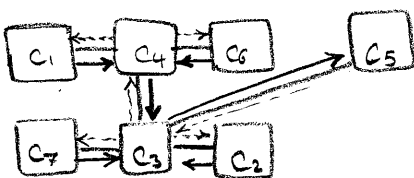
↳ is there any deadlock situations? = all ppl waiting for somebody else to send'em msg.

↳ msg's from leaves are immediately final \rightarrow then the parent of that $\rightarrow \dots$

↳ because it is a tree it is guaranteed to find a legal msg-passing order

↳ if we pass msg in correct order (= from leaves inward) only need to pass $2(k-1)$ msg's ($k = \#$ of clusters $\rightarrow k-1$ edges $\times 2$ msg for each edge)

message passing order I



\rightarrow forward
 \leftarrow backward

we can start with any leaf. say C_2

↳ C_3 can't pass : waiting

↳ we had to go for another leaf. say C_1

↳ C_3 & C_4 can't pass : waiting

↳ another leaf : C_6

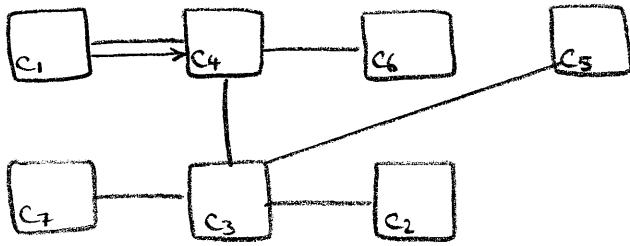
↳ now we have the option of activating C_4

↳ another leaf : C_7

↳ C_3 can be activated

at this point every one receive all the msg's but one

↳ another direction: $C_5 \rightarrow C_3 \rightarrow C_2 \rightarrow C_7 \rightarrow C_4 \rightarrow C_6 \rightarrow C_1$



$C_1 \rightarrow C_4$

$C_4 \rightarrow C_3$ x illegal

↳ does not have all the info it needs to pass a final msg.

↳ still waiting for C6

- Answering Query

↳ clique tree have good properties → make it a useful DS

↳ which kind of queries does it answer?

↳ Posterior distribution queries on variables that appear together in a clique

↳ find any clique that contains all of those variables

sum out irrelevant variables from that clique → results in \tilde{P}_Φ

renormalize them

↳ Introducing new evidence $\mathcal{Z} = z$ and querying X

↳ a form of incremental inference: you already calibrated clique tree

and say wait a minute! I made another observation

↳ clique trees are good for that

↳ Case 1: If X appears in clique with \mathcal{Z}

$X, \mathcal{Z} \in \Phi_i$ → what the posterior X would be if I further observe \mathcal{Z}

↳ $\tilde{P}_\Phi(\mathcal{Z}, X, \dots)$

↳ we can condition or reduce the clique

↳ just restricting the attention to the entries that are consistent with evidence \mathcal{Z}

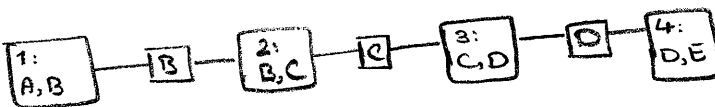
$= \Phi_i \times \mathbb{1}(\mathcal{Z} = z) = \tilde{P}_\Phi(\mathcal{Z} = z, X, \dots)$

sum out irrelevant variables

renormalize

↳ Case 2: If X does not share a clique with \mathcal{Z}

⊠



observe $A = a$
query D

↳ reduction of factors including \mathcal{Z}

↳ multiply by $\mathbb{1}(\mathcal{Z} = z)$

observe $A = a \rightarrow \psi_i = \psi_i \times \mathbb{1}(A = a) \rightarrow$ removing all entries inconsistent with $A = a$

what happens when we change Ψ_i ?

↳ strategy 1: recompute everything from new \rightarrow all msgs and beliefs

↳ strategy 2: look which msgs would change

Ψ_i change $\rightarrow \delta_{1 \rightarrow 2}$ change $\rightarrow \delta_{2 \rightarrow 3}$ change $\rightarrow \delta_{3 \rightarrow 4}$ change but since we are interested in C_3 we don't care

$\delta_{4 \rightarrow 3}$ doesn't change $\rightarrow \delta_{3 \rightarrow 2}$ and $\delta_{2 \rightarrow 1}$ doesn't change

↳ we can reuse at least half of messages since they don't change

↳ the only msgs that we need to recompute are the msgs that are along the path to clique containing X

- Summary

↳ In clique tree with k cliques, if messages are passed starting at leaves, $2(k-1)$ msgs suffice to compute all beliefs

↳ we get posterior over every single var in the model

↳ contrast that with VE complexity and it only gives posterior over 1 var in model

↳ Can compute marginals over all variables at only twice the cost of VE

↳ all marginals \rightarrow VE: $\# \text{ var} \times O(VE)$

↳ clique tree: $2 \times O(VE)$

↳ By storing msgs, inference can be reused in incremental queries.

• Clique tree & Independence

- clique tree alg properties \rightarrow guarantee we achieve the correct marginals at every single clique

\rightarrow marginals agree with each other

\rightarrow marginals can be computed using a single upward and downward pass

↳ Inference in PGM is NP-hard \rightarrow there should be a catch somewhere.

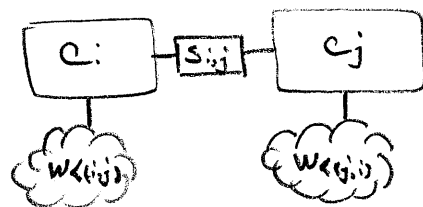
↳ a computational cost that occurs at least in certain cases \rightarrow where is it?

- For an edge (i,j) in T let break tree in 3 parts:

↳ $W_{\langle i,j \rangle}$ = all vars that appear only on C_i side of T

↳ $W_{\langle j,i \rangle}$ = all vars that appear only on C_j side of T

↳ variables on both sides are in the sepset S_{ij}



- Theorem: T satisfies RIP if and only if for every (i,j)

$$P_{\Phi} \models (W_{\langle(i,j)} \perp W_{\langle(j,i)} \mid S_{i,j})$$

variables on left side of edge are independent of the variables on the right side of the edge given the variables on the sepset \rightarrow sepset separate the left side from the right side.

\hookrightarrow RIP is the critical property to prove correctness of clique tree alg

\hookrightarrow so RIP is important to have all those nice properties

\checkmark focus on sepset $S_{3,4} (G,S)$

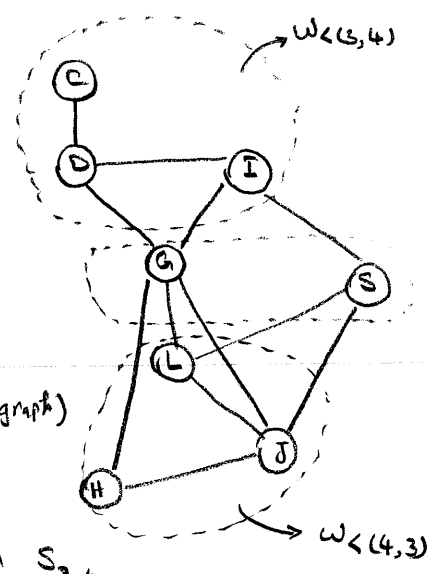


$$W_{\langle(3,4)} = \{I, D, C\}$$

$$W_{\rangle(4,3)} = \{J, L, H\}$$

$$S_{3,4} (G, S)$$

looking at the induced graph
(BN CPDs \rightarrow factors in MN \rightarrow induced graph)

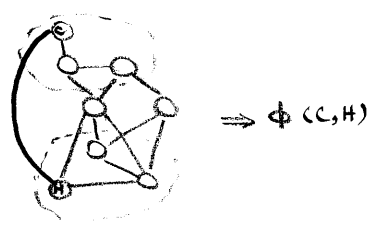
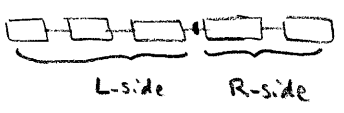


there are no path between $W_{\langle(3,4)}$ and $W_{\rangle(4,3)}$ that are go through $S_{3,4}$

\hookrightarrow C, I, D separated from H, L, J given G, S

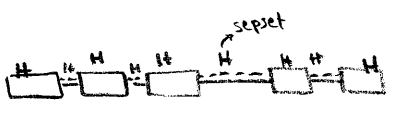
$$\hookrightarrow P_{\Phi} \models (\{C, I, D\} \perp \{J, L, H\} \mid \{G, S\}) \checkmark$$

\hookrightarrow Proof by contradiction: let's assume that vars on left side are not independent of vars on right side, given sepst.



$$\Phi(C, H) \rightarrow C_1(C, D, H)$$

\downarrow
left side



- \downarrow there needs to be path in the induced MN between L-side and R-side. that doesn't go through $S_{i,j}$
- \downarrow there has to be an edge where one node in one side, and one node in the other
- \downarrow there needs to be a factor Φ , that involves those two variables
- \downarrow because of family preservation, that factors should sit in one of the cliques.
- \downarrow that clique is either on the L-side or R-side.
- \downarrow now we have variable H in one clique in L-side and in one-clique in R-side
- \downarrow R.I.P tells that H should be everywhere in between and specifically in sepset
- \downarrow violation of the assumption that H is not in sepset

- Implications

- ↳ where are computational implications? what we conclude from the fact that sepsets divide the graph into conditionally indep. pieces?
- ↳ Each sepset needs to separate graph into two conditionally independent pieces.
 - ↳ In many graphs that implies a certain minimal complexity → sometimes quite large

↳ Case 1: complete bipartite graph

- two sets of vars
- edges only cross the sets

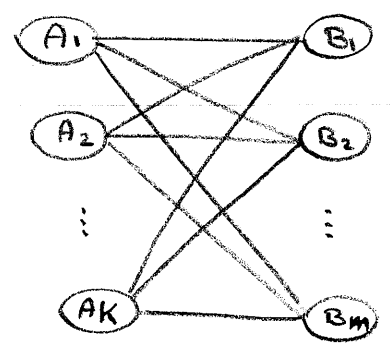
↳ what is the smallest sepset that we could construct?

↳ we need to separate graph into 2 conditionally indep pieces

↳ smallest sepset that we could construct that actually break up the graph into meaningful pieces is all of var on A side or all the variables on B side

↳ smallest sepset size $\geq \min(k, m)$

course difficulty & student intelligence

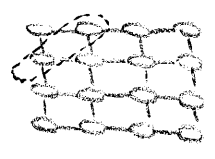


↳ Case 2: grid

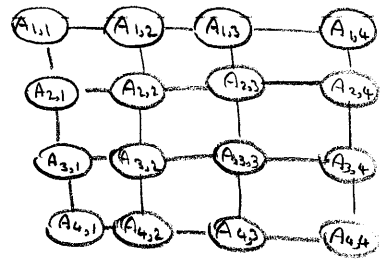
↳ how to break up graph into separate conditionally indep pieces

↳ we can construct clique trees to have smaller sepsets

e.g.

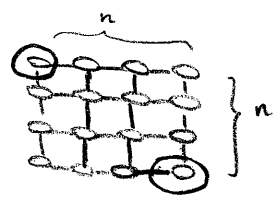


→ leaves a very large everything else



↳ if we want to break the grid such that A1,1 appears in one side and A4,4 in another

↳ any clique tree satisfying this separation must have a sepset of size $\geq n$



↳ breaking up in other ways doesn't make it any better

↳ In these cases we could put a lower bound on the size of the sepset that is required for running clique tree inference

↳ this is where we pay exponential blow-up

- Summary

↳ Correctness of clique tree inference relies on running intersection property

↳ RIP implies separation in original dist.

↳ Implies minimal complexity incurred by any clique tree

↳ (as seen before) Related to width of minimal induced graph

induced width talks about cliques but it is related to sepsets too

Consider the optimal clique tree for a pairwise MRF structured as $n \times n$ grid. what is the size (# var)

of the largest sepset in the clique tree?
 → smallest sepset possible
 → sepset that separates two parts yet to be smallest $\Rightarrow O(n)$

• Clique tree and VE

- VE → a way of proving clique tree alg. correctness
- help answering an important question: where clique trees come from?
- ↳ or we might construct a little clique tree?

- Variable Elimination

- I. each step create a large factor λ_i through factor product
- II. a variable is eliminated in λ_i to generate new factor τ_i
- III. τ_i is used in computing other factors λ_j

↳ each generated τ_i will be used in context of computing another large factor

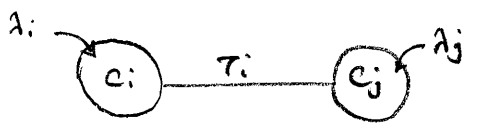
↳ Let's think about this process as passing msgs

- * Intermediate factors λ_i are cliques
- * τ_i 's are messages generated by clique λ_i and transmitted to another clique λ_j

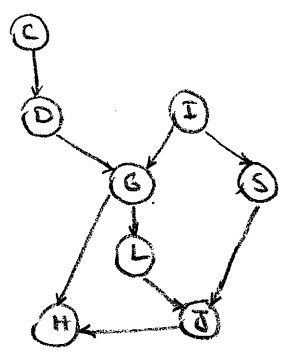
↳ VE defines a graph

↳ cluster C_i : for each factor λ_i used in computation

↳ draw edge $C_i - C_j$ if the factor generated from λ_i is used in the computation of λ_j

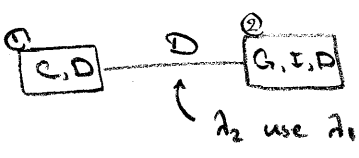


- ✓ $\phi_C(C) \phi_D(C,D) \phi_I(I) \phi_G(D,I,G) \phi_S(S,I) \phi_L(S,L,G)$
 $\phi_H(G,J,H) \phi_J(J,L,S)$



① eliminate C: $\tau_1(D) = \sum_C \frac{\phi_C(C) \phi_D(C,D)}{\lambda_1(G,D)}$

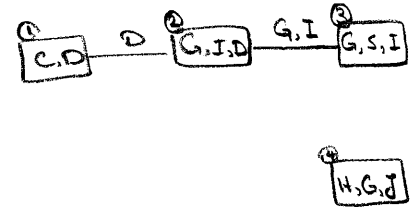
② eliminate D: $\tau_2(G,I) = \sum_D \frac{\phi_G(G,I,D) \tau_1(D)}{\lambda_2(G,I,D)}$



③ eliminate I: $\tau_3(G,S) = \sum_I \frac{\phi_I(I) \phi_S(S,I) \tau_2(G,I)}{\lambda_3(G,I,S)}$

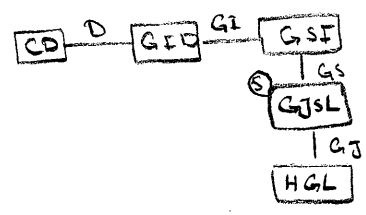
④ eliminate H: $T_4(G, J) = \sum_H \frac{\phi_H(H, G, J)}{\lambda_4(H, G, J)}$

↳ didn't use any factor produced by other factors so far



⑤ eliminate G: $T_5(J, L, S) = \sum_G \frac{\phi_L(L, G) T_3(G, S) T_4(G, J)}{\lambda_5(G, J, S, L)}$

↳ use factors that produced in two separate computations

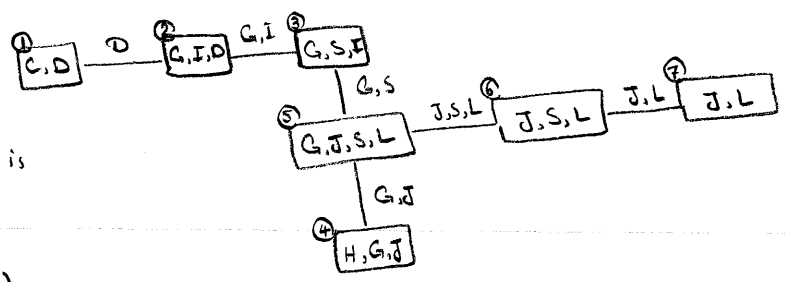


⑥ eliminate S: $T_6(J, L) = \sum_S \frac{\phi_J(J, L, S) T_5(J, L, S)}{\lambda_6(J, L, S)}$

⑦ eliminate L: $T_7(J) = \sum_L \frac{T_6(J, L)}{\lambda_7}$

↳ this clique tree is bigger than previous ones

↳ reason: we have 3 adjacent cliques where one is subset of another.



$C_7(J, L) \subset C_6(J, S, L) \subset C_5(J, G, S, L)$

↳ require post processing step: remove redundant cliques

↳ those whose scope is a subset of adjacent clique's scope