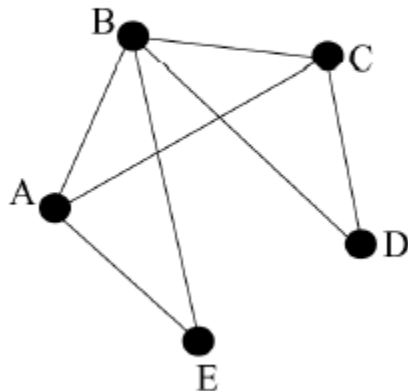


Reasoning with uncertainty

- (Very) basic review of probability and uncertainty
- Joint distribution and inference
- Exploiting independences
- Special case: Directed graphs
- General case: Undirected graphs, factor graphs
- Examples

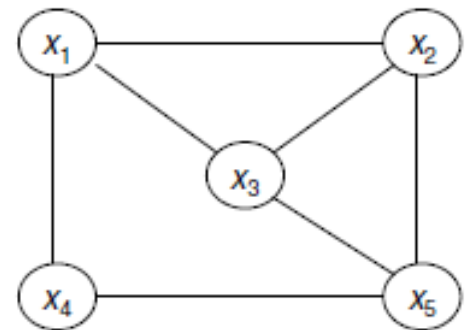
- Deterministic:
 - Represent facts and constraints
 - Find configuration that satisfies representation
- Examples:
 - Clauses $A \wedge B \Rightarrow C$
 - Satisfiability problems

$\varphi = \{(-C), (A \vee B \vee C), (-A \vee B \vee E), (-B \vee C \vee D)\}$.



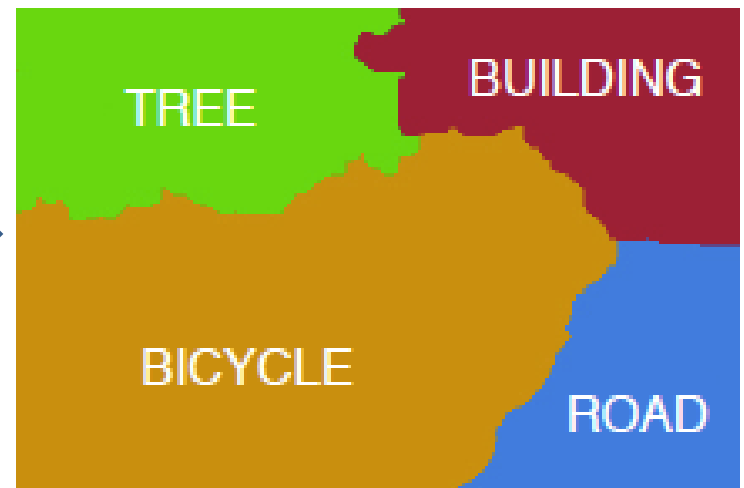
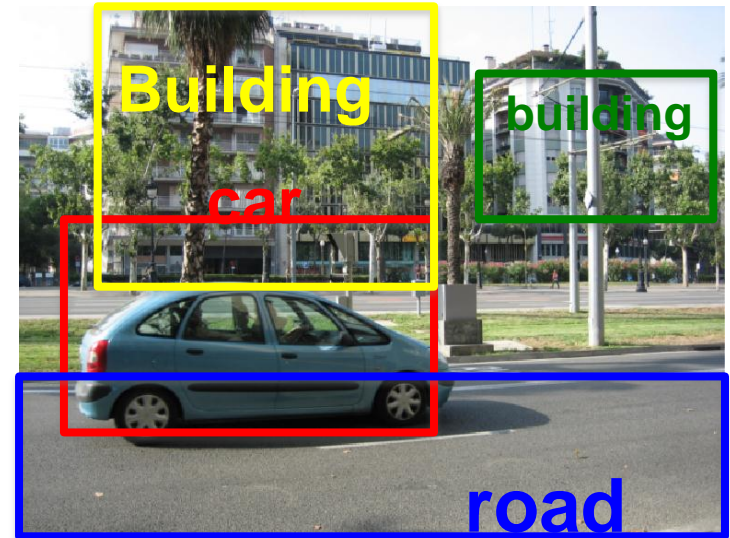
$$\begin{aligned} & \bar{f}_1(x_1, x_2, x_3) \\ & f_2(x_2, x_3, x_5) \\ & f_3(x_1, x_4) \\ & f_4(x_4, x_5) \end{aligned}$$

$$\min_{t \in \text{Sol}} \left\{ \sum_{i=1}^{m'} f_i(t) \right\}$$



- Generalization to include uncertainty due to imperfect knowledge
 - Variables: Deterministic \rightarrow Random variables
 - Constraints: Deterministic functions (e.g., CNF, CSP, SAT) \rightarrow continuous output
- Similar problems, generalization

- What we have: scores from noisy classifiers from local features ($P(\text{label} | \text{image features})$)
- What would like:



- What we get:



Reasoning

- Need to use knowledge about the world
- Need to integrate uncertainty in “sensing”
- Scenes:
 - “road scenes” contain “cars”, “building”....
 - “Office scenes” contain “desks”, “computers”...
- Co-occurrence:
 - “keyboard” implies “mouse”
- Location:
 - “Cars” are on top of “roads”
 - “Sky” is above “buildings”

Reasoning with uncertainty

- Need use *uncertain* knowledge about the world
- Scenes:
 - “road scenes” *is likely to* contain “cars”, “building”
 - “Office scenes” *is likely to* contain “desks”, “computers” ...
- Co-occurrence:
 - “keyboard” *usually* implies “mouse”
-
- Probably possible to represent uncertainty on each individual piece of knowledge
- Intractable to integrate them all to find the “optimal” interpretation

Probability Reminder

- Conditional probability for 2 events A and B:

$$P(A | B) = \frac{P(A, B)}{P(B)}$$

- Chain rule:

$$P(A, B) = P(A | B) P(B)$$

Probability Reminder

- Conditional probability for 2 variables X and Y:

$$P(X=x \mid Y=y) = \frac{P(X=x, Y=y)}{P(Y=y)}$$

- Chain rule:

$$P(X=x, Y=y) = P(X=x \mid Y=y) P(Y=y)$$

- For any values x,y

The Joint Distribution

- Joint distribution = collection of all the probabilities $P(X = x, Y = y, Z = z \dots)$ for all possible combinations of values.
- For m binary variables, size is 2^m
- Any query can be computed from the joint distribution

X	Y	Z	Prob
T	T	T	0.1
T	T	F	0.22
T	F	T	0.2
T	F	F	0.08
F	T	T	0.1
F	T	F	0.15
F	F	T	0.07
F	F	F	0.08

The Joint Distribution

- Any query can be computed from the joint distribution
- Marginal distribution

$$P(X = \text{True}), P(X = \text{False})$$

- Conditional distribution:

$$P(X = \text{True} \mid Y = \text{True}) = \\ P(X = \text{True}, Y = \text{True}) / P(Y = \text{True})$$

- In general:

$$P(E_1 \mid E_2) = P(E_1, E_2) / P(E_2)$$

$$P(E_2) = \sum_{\text{Entries that match } E_2} P(\text{Joint Entries})$$

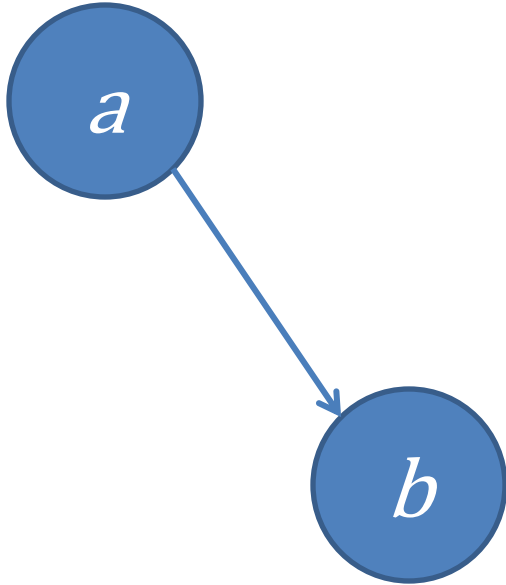
X	Y	Z	Prob
T	T	T	0.1
T	T	F	0.22
T	F	T	0.2
T	F	F	0.08
F	T	T	0.1
F	T	F	0.15
F	F	T	0.07
F	F	F	0.08

Summary

- Any query computable from
- Sum rule:
- Product rule

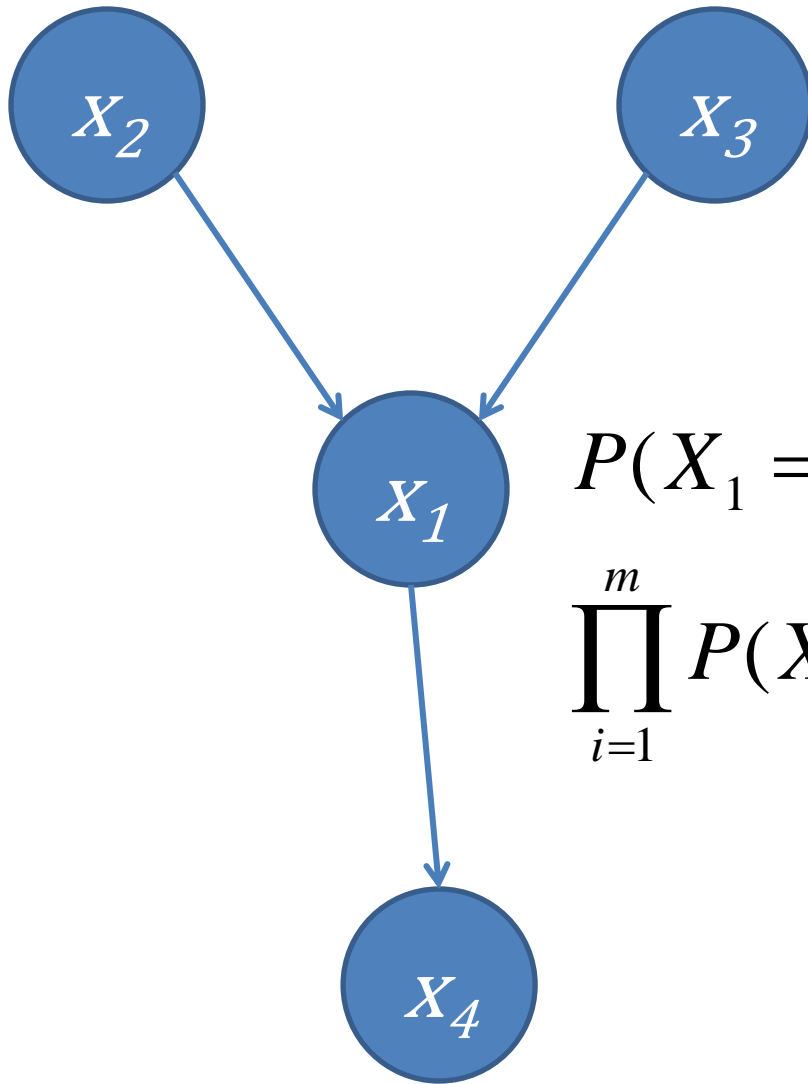
But requires entire joint distribution →
Represent dependencies between variables

First case: Directed



$$P(a, b) = P(b|a)P(a)$$

First case: Directed

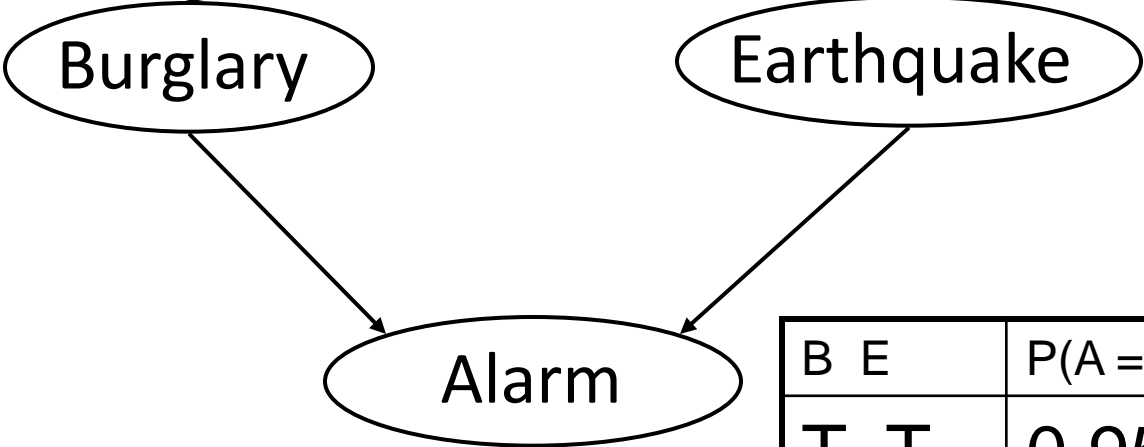


$$P(X_1 = x_1, X_2 = x_2, \dots, X_m = x_m) = \prod_{i=1}^m P(X_i = x_i \mid \text{Parents}(X_i))$$

Graphical Representation

$$P(B=True) = 0.001$$

$$P(E=True) = 0.002$$



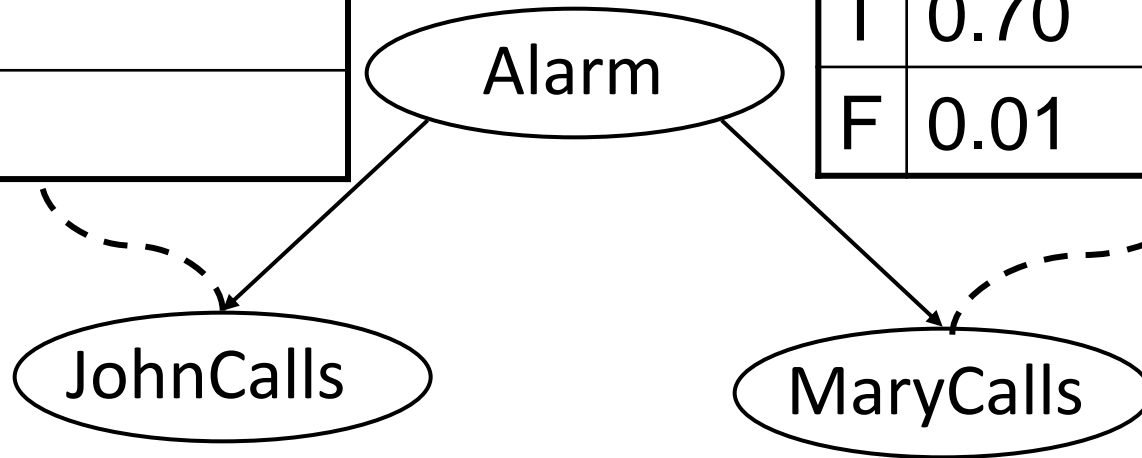
B	E	$P(A = True B=b,E=e)$
T	T	0.95
T	F	0.94
F	T	0.29
F	F	0.001

$$P(A,J,M) = P(A|B,E)P(B)P(E)$$

Graphical Representation

A	$P(J = \text{True} A=a)$
T	0.90
F	0.05

A	$P(M = \text{True} A= a)$
T	0.70
F	0.01



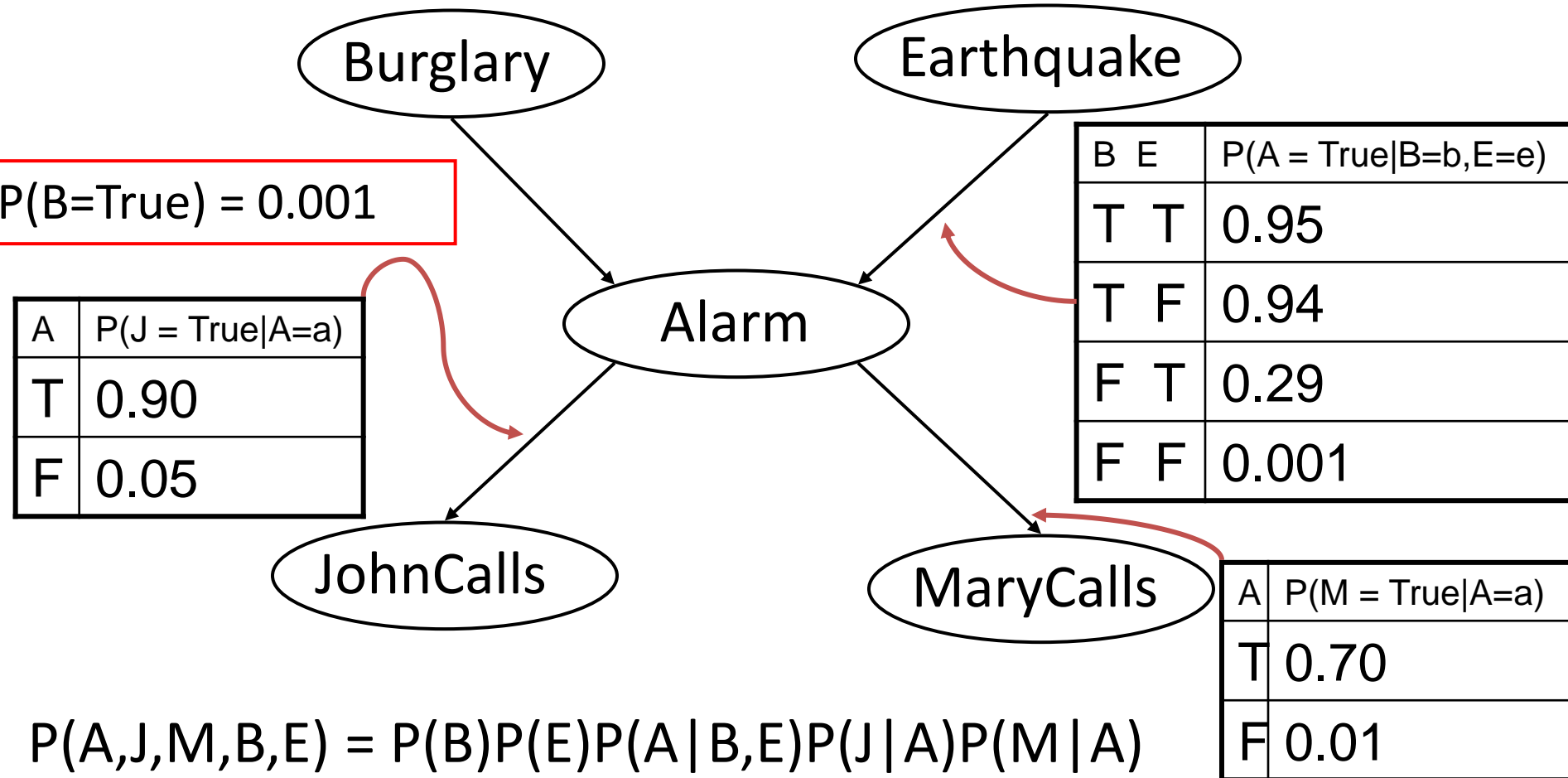
Given knowledge of A, knowing anything else in the diagram won't help with J and M

$$P(A, J, M) = P(A)P(J | A)P(M | A)$$

Inference

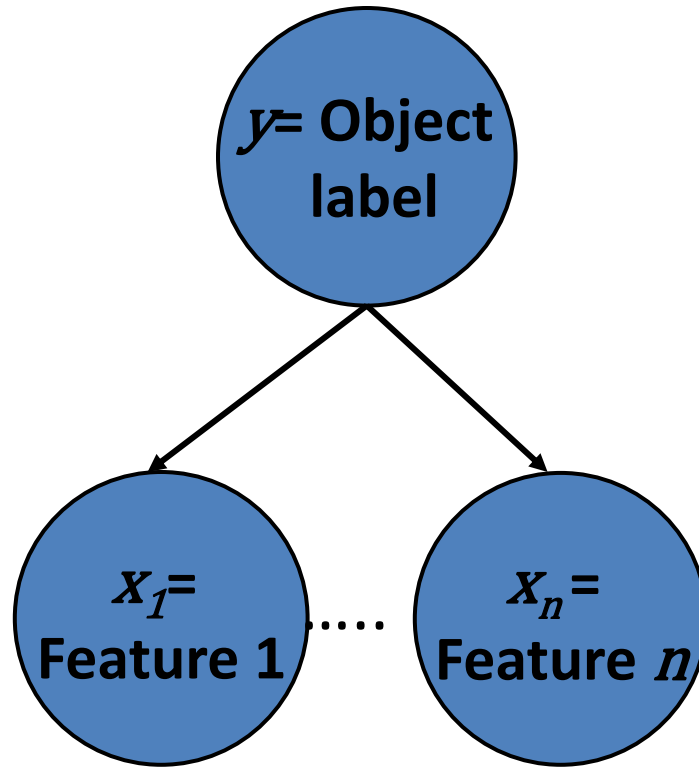
- Any inference operation of the form $P(\text{values of some variables} \mid \text{values of the other variables})$ can be computed

$$P(E=\text{true}) = 0.002$$



$$P(A, J, M, B, E) = P(B)P(E)P(A | B, E)P(J | A)P(M | A)$$

Example Naïve Bayes classification



$$P(x_1, \dots, x_n | y) = \prod_i P(x_i | y)$$

Example



$$f_1(0, 0) = \#5710$$

$$P_1(\#5710, 0, 0 \mid \mathbf{F}) = 0.53$$

$$P_1(\#5710, 0, 0 \mid \mathbf{not F}) = 0.56$$



$$f_1(0, 1) = \#3214$$

$$P_1(\#3214, 0, 1 \mid \mathbf{F}) = 0.57$$

$$P_1(\#3214, 0, 1 \mid \mathbf{not F}) = 0.48$$

$y = 1$ if face

- Lots of (discretized) features from local filters
- Estimate likelihood ratio

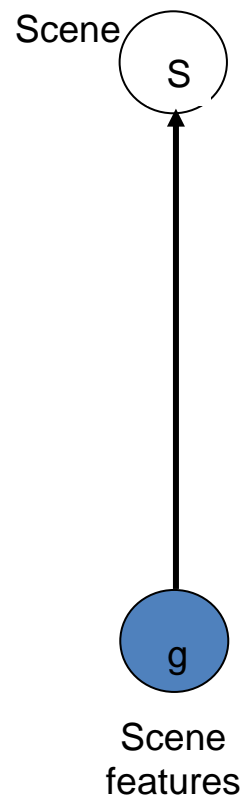
$$\frac{P(x_1, \dots, x_n \mid y = \text{face})}{P(x_1, \dots, x_n \mid y = \text{not face})}$$



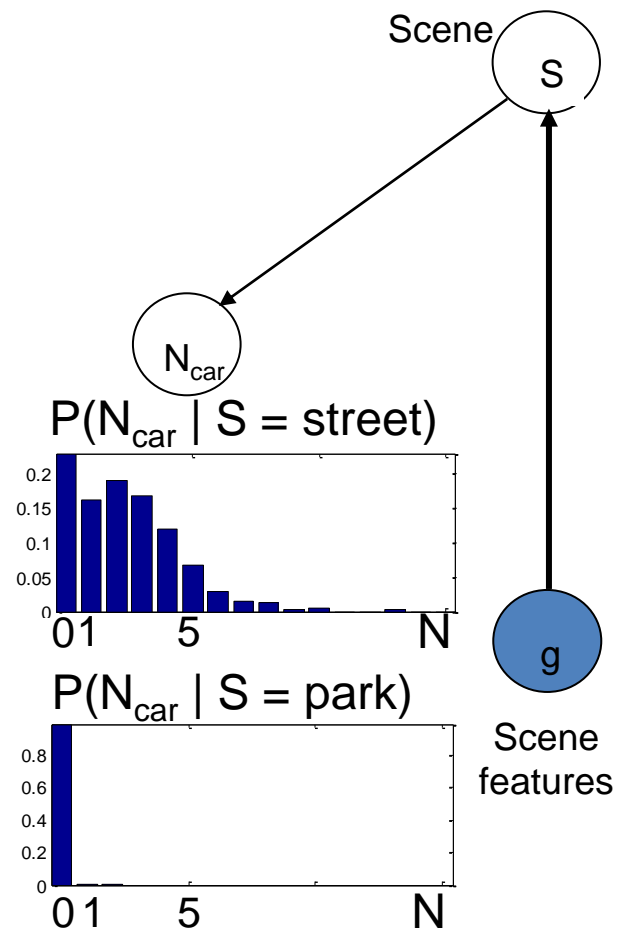
Example



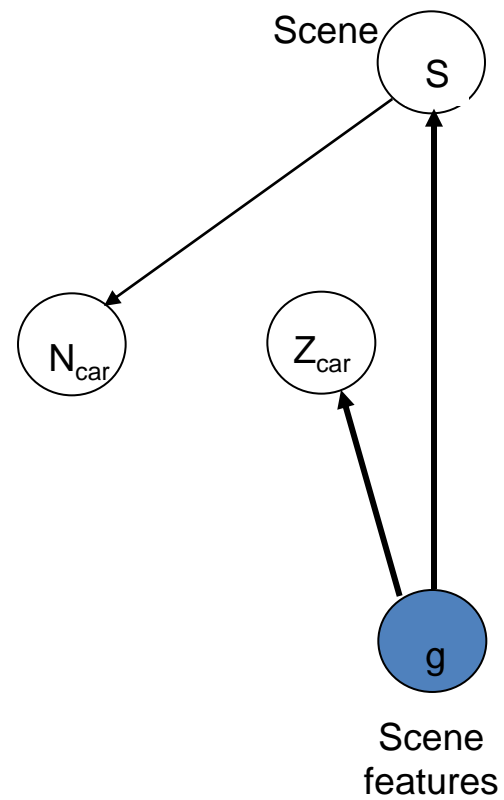
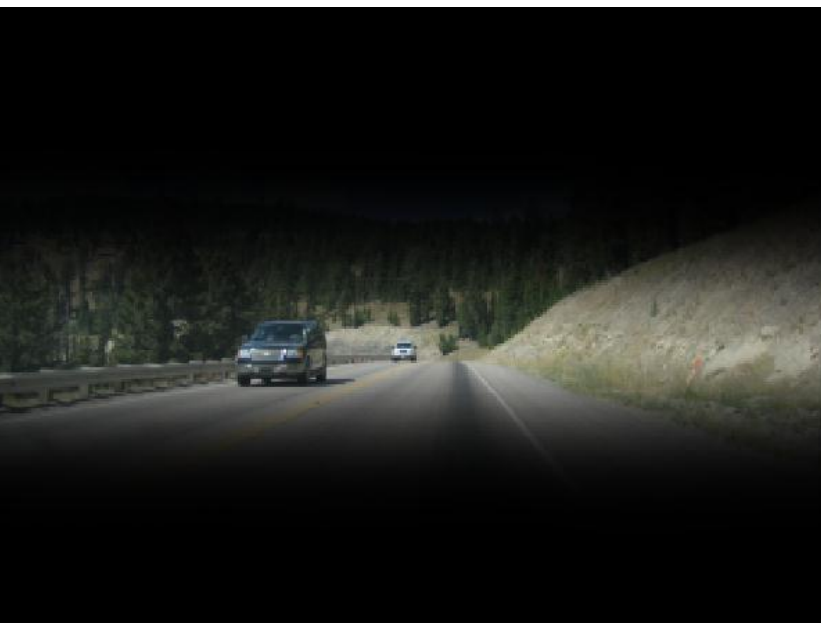
- Need use *uncertain* knowledge about the world
- Scenes:
 - “road scenes” *is likely to* contain “cars”, “building”
 - “Office scenes” *is likely to* contain “desks”, “computers” ...
- Co-occurrence:
 - “keyboard” *usually* implies “mouse”
-
- Probably possible to represent uncertainty on each individual piece of knowledge
- Intractable to integrate them all to find the “optimal” interpretation

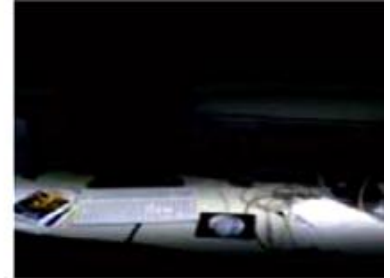


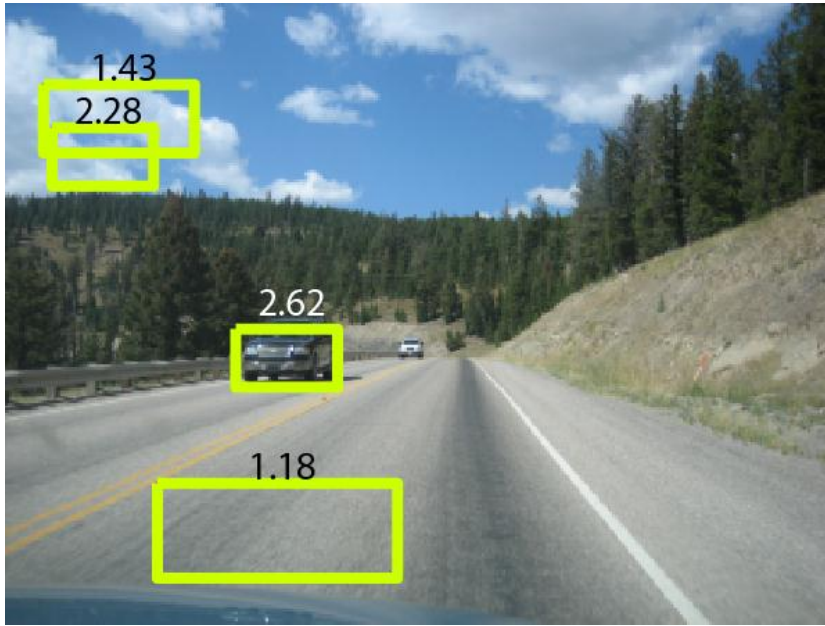
Murphy, Torralba, Freeman; NIPS 2003. Torralba, Murphy, Freeman, CACM 2010.



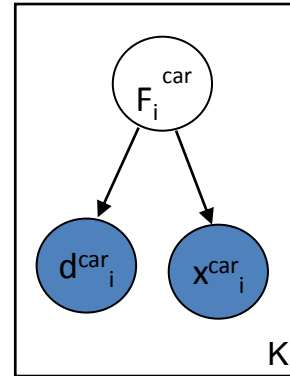
Murphy, Torralba, Freeman; NIPS 2003. Torralba, Murphy, Freeman, CACM 2010.





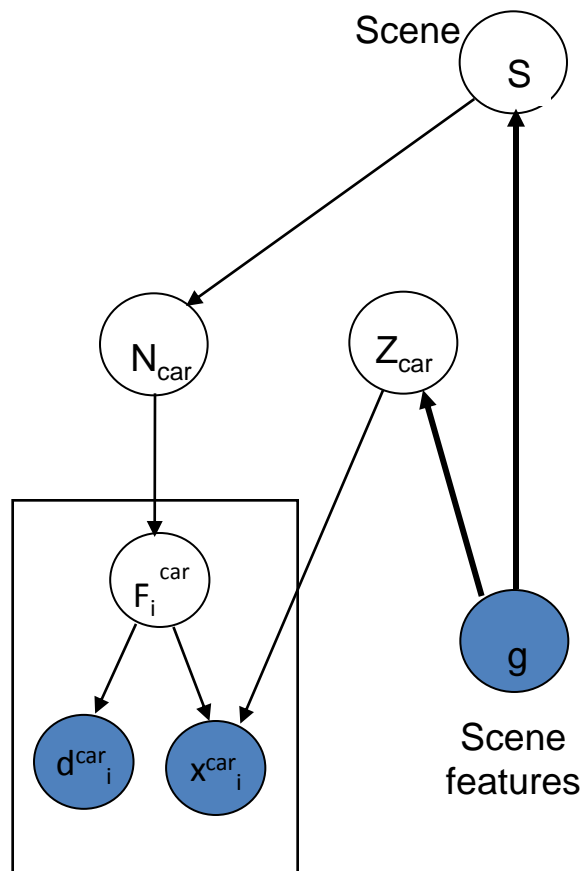


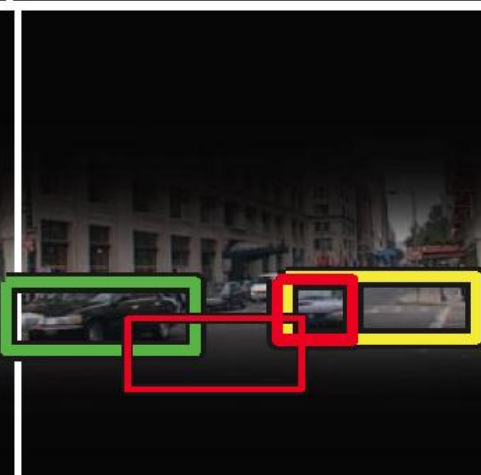
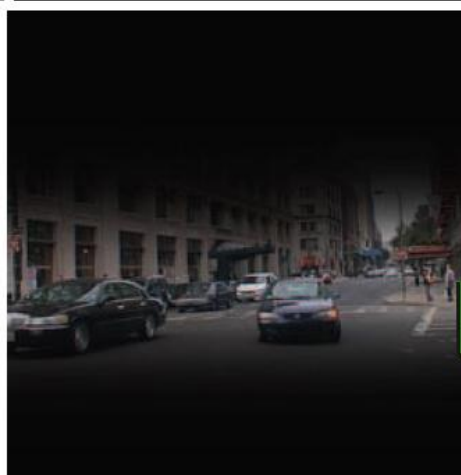
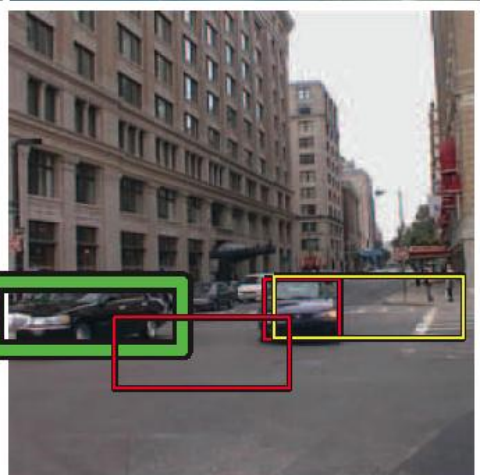
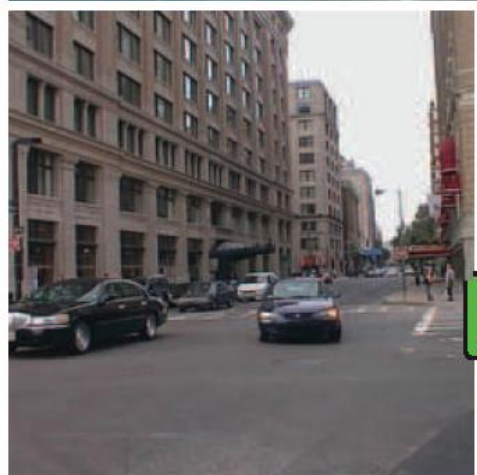
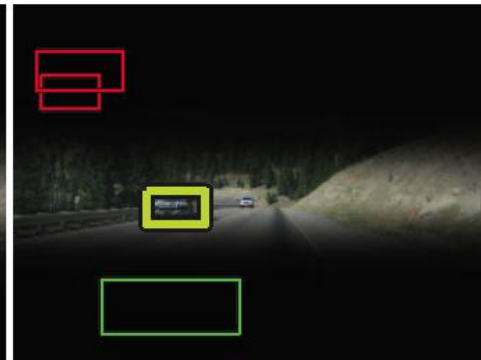
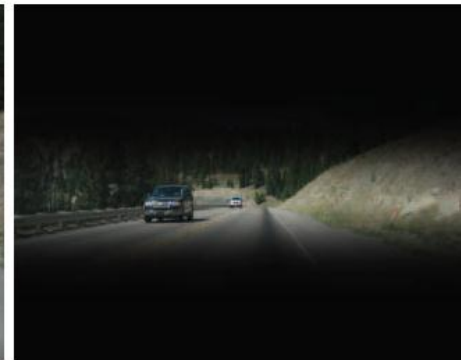
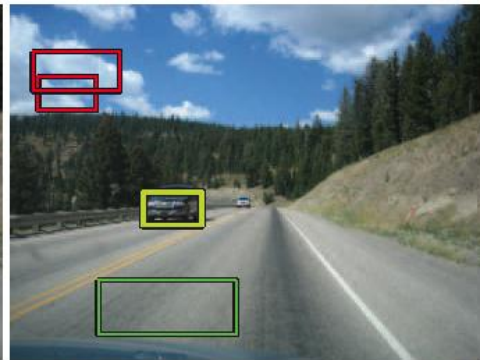
Multiview car detector.

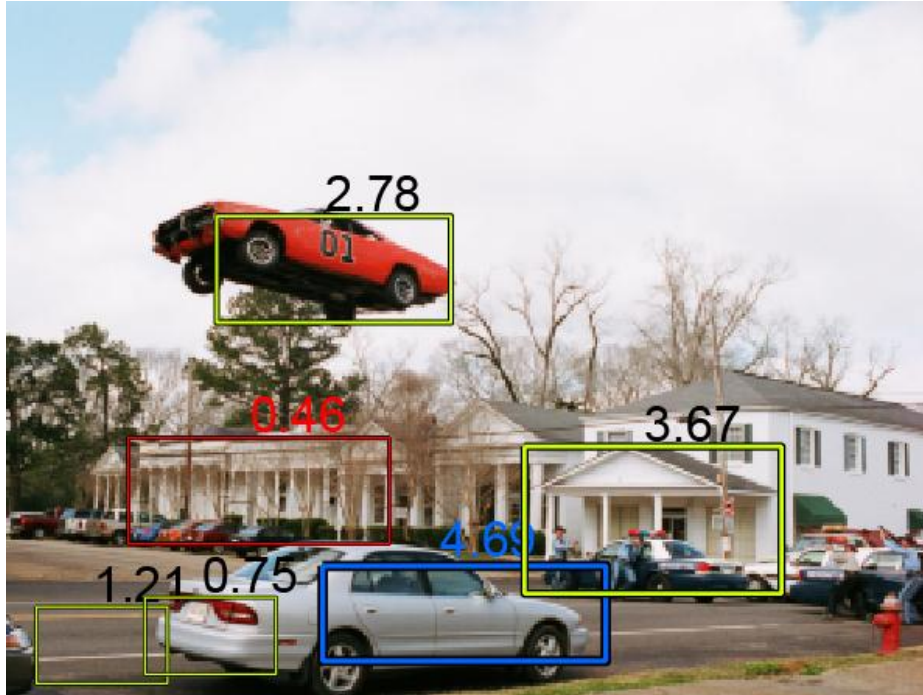


$F = 1$ if car present in box
 $p(d | F=1)$

An integrated model of Scenes, Objects, and Parts



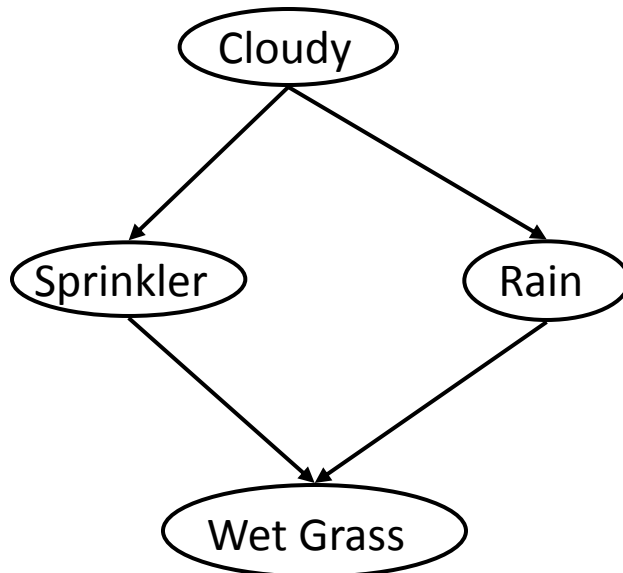




- No miracle: Fancy representation can only model the knowledge that we encoded.

Inference

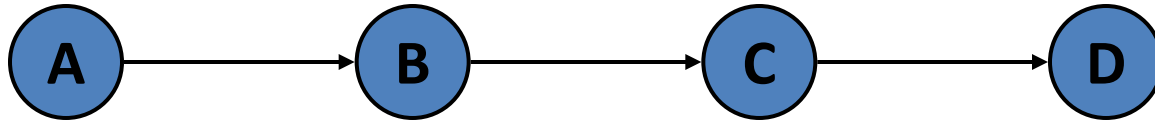
- Can answer any query but : Need to sum over the possible assignments of the hidden variables.
 - Variable elimination
 - Separation
- Query variables: E_1
- Evidence variables: E_2
- The rest, E_3



$P(W \mid \text{Cloudy} = \text{True})$

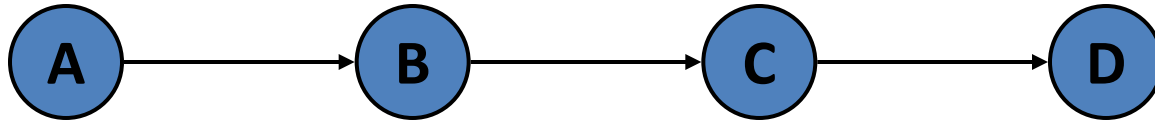
- $E_1 = \{W\}$
- $E_2 = \{\text{Cloudy} = \text{True}\}$
- $E_3 = \{\text{Sprinkler}, \text{Rain}\}$

Inference: A Simple Case



- Suppose that we want to compute $P(D = d)$ from this network.

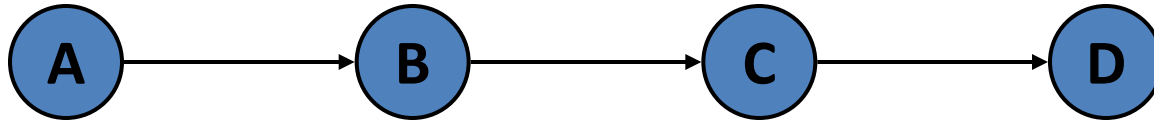
A Simple Case



- Compute $P(D = d)$ by summing the joint probability over all possible values of the remaining variables A , B , and C :

$$P(D = d) = \sum_{a,b,c} P(A = a, B = b, C = c, D = d)$$

A Simple Case

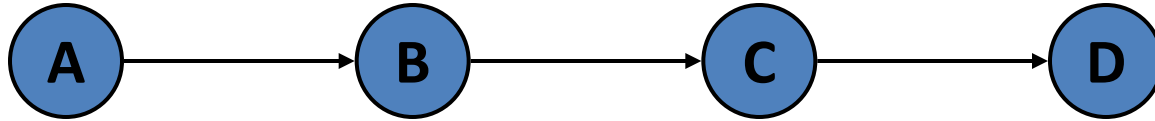


- Decompose the joint by using the fact that it is the product of terms of the form:

$$P(X \mid \text{Parents}(X))$$

$$P(D = d) = \sum_{a,b,c} P(D = d \mid C = c)P(C = c \mid B = b)P(B = b \mid A = a)P(A = a)$$

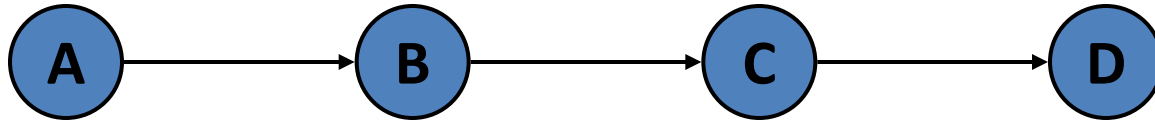
A Simple Case



- We can avoid computing the sum for all possible triplets (A, B, C) by distributing the sums inside the product

$$P(D = d) = \sum_c P(D = d | C = c) \sum_b P(C = c | B = b) \sum_a P(B = b | A = a) P(A = a)$$

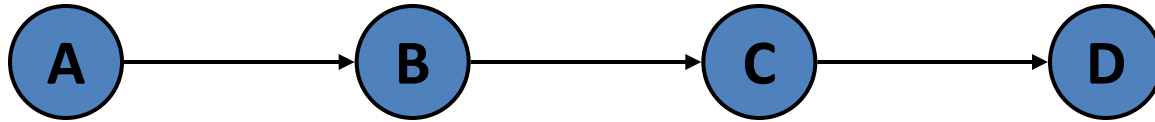
A Simple Case



This term depends only on **B** and can be written as a function $f_A(b)$

$$P(D = d) = \sum_c P(D = d | C = c) \sum_b P(C = c | B = b) \sum_a P(B = b | A = a) P(A = a)$$

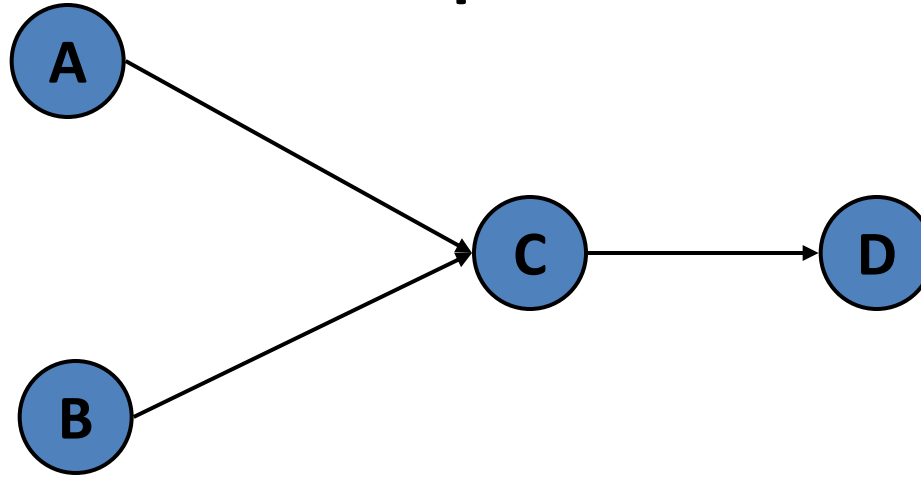
A Simple Case



This term depends only on c and can be written as a function $f_B(c)$

$$P(D = d) = \sum_c P(D = d | C = c) \sum_b P(C = c | B = b) f_A(b)$$

Example



$$P(D=d) = \sum_{a,b,c} P(D=d | C=c) P(C=c | B=b, A=a) P(B=b) P(A=a)$$

$$= \sum_c P(D=d | C=c) \sum_b P(B=b) \sum_a P(C=c | B=b, A=a) P(A=a)$$

$$= \sum_c P(D=d | C=c) \sum_b P(B=b) \sum_a f_1(a, b, c)$$

$$= \sum_c P(D=d | C=c) \sum_b f_2(b, c)$$

General Case: Variable Elimination

- Write the desired probability as a sum over all the unassigned variables

$$P(D = d) = \sum_{a,b,c} P(A = a, B = b, C = c, D = d)$$

- Distribute the sums inside the expression

- Pick a variable
- Group together all the terms that contain this variable

$$P(D = d) = \sum_c P(D = d | C = c) \sum_b P(C = c | B = b) \sum_a P(B = b | A = a) P(A = a)$$

- Represent as a single function of the variables appearing in the group

$$P(D = d) = \sum_c P(D = d | C = c) \sum_b P(C = c | B = b) f_A(b)$$

- Repeat until no more variables are left

General Case: Variable Elimination

- Write the desired probability as a sum over all the unassigned variables

Computation exponential in the size of the largest group \rightarrow The order in which the variables are selected is important.

on

– Pick a variable

– Group together all the terms that contain this variable

$$P(D=d) = \sum_c P(D=d | C=c) \sum_b P(C=c | B=b) \sum_a P(B=b | A=a) P(A=a)$$

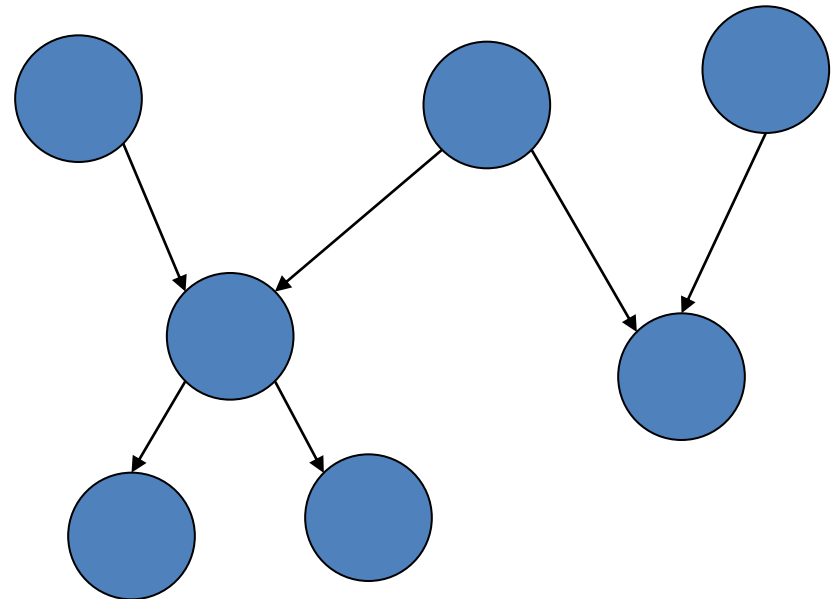
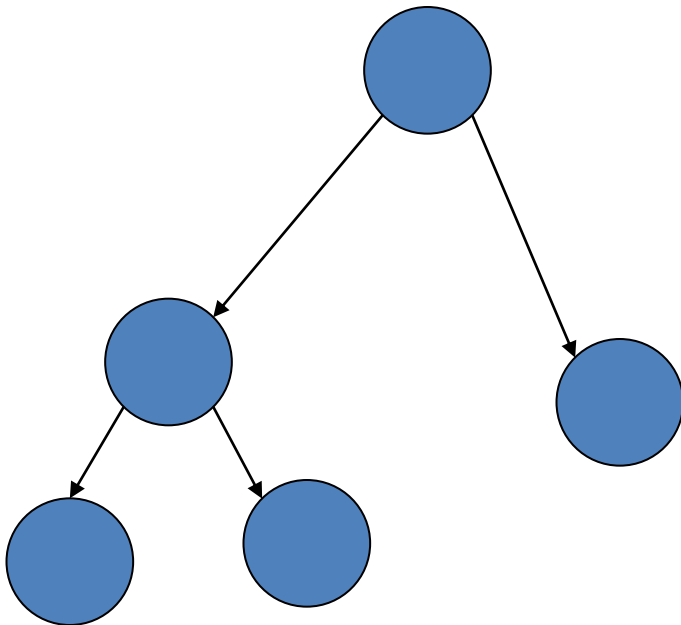
– Represent as a single function of the variables appearing in the group

$$P(D=d) = \sum_c P(D=d | C=c) \sum_b P(C=c | B=b) f_A(b)$$

– Repeat until no more variables are left

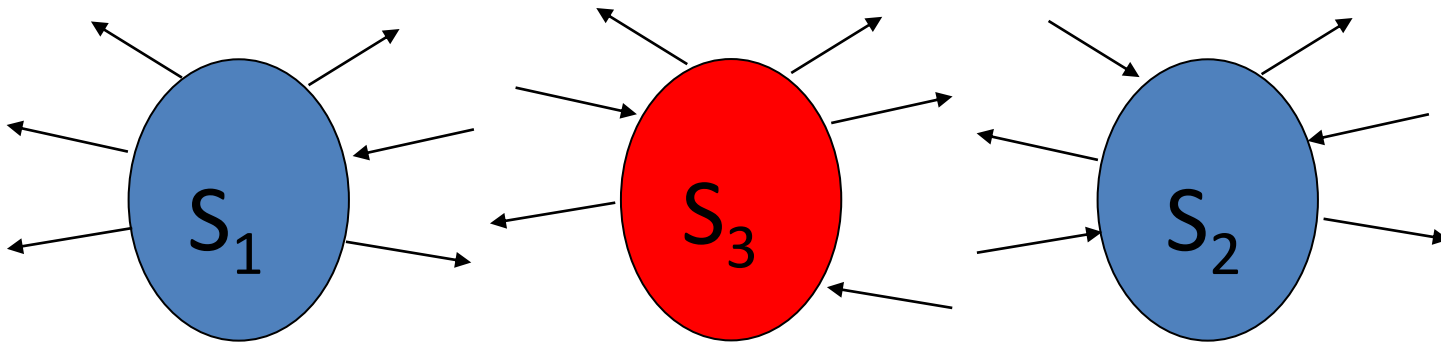
Special Case

- Polytrees: Undirected version of the graph is a tree = there is a single undirected path between two nodes
- In this case: Inference linear in the number of nodes ($d^{k+1}n$)
- General case: See later approximate inference (e.g., sampling)



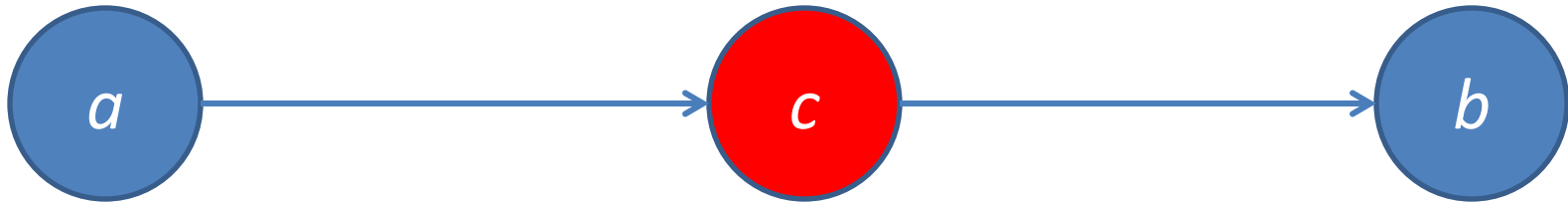
Conditional independence

- $P(\text{any assignments to } S_1 \mid \text{any assignments to } S_2, \text{any assignments to } S_3) = P(\text{assignment to } S_1 \mid \text{assignments to } S_3)$
- $P(\text{any assignments to } S_1, \text{any assignments to } S_2 \mid \text{any assignments to } S_3) = P(\text{assignment to } S_1)P(\text{assignments to } S_2)$

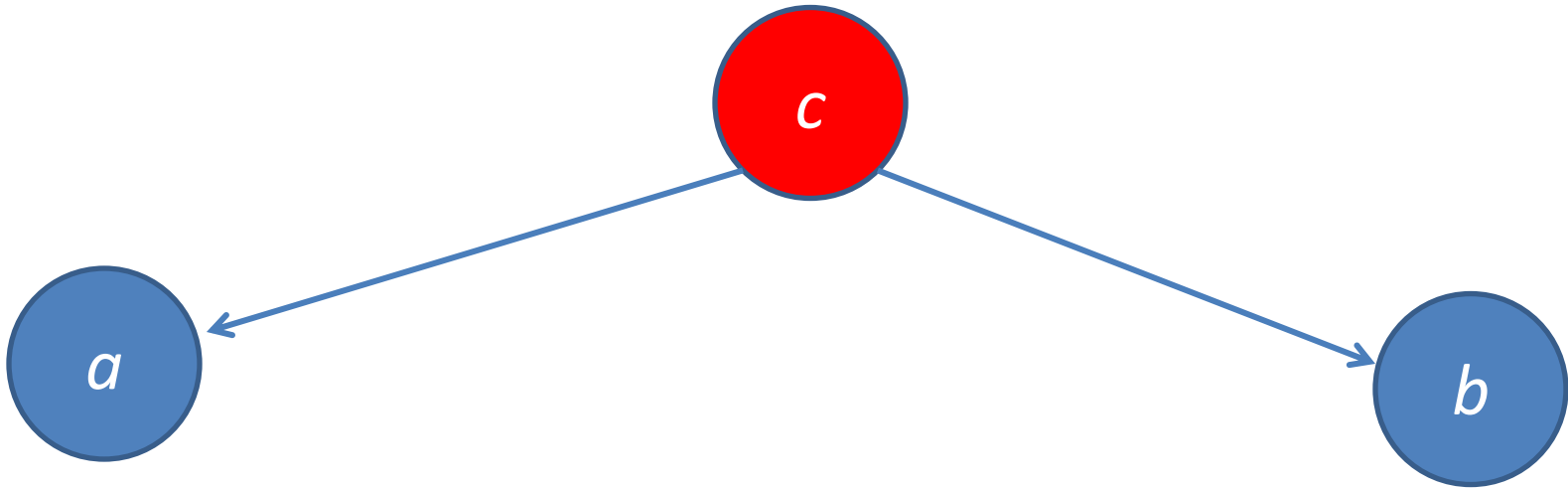


Finding independences

- The more *independence* relations we can find, the faster the inference → Test to find independences?



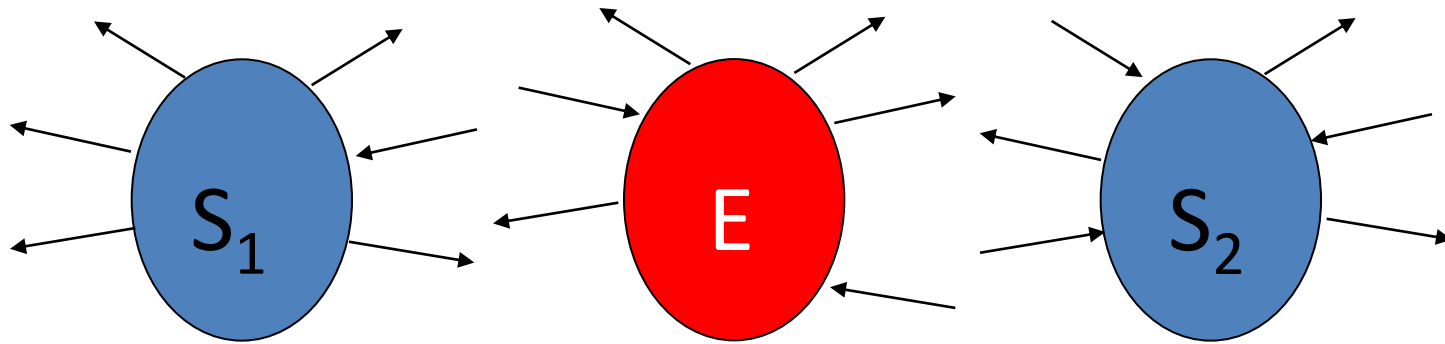
$$p(b|c, a) = P(b|c) \quad a \perp b | c$$



$$p(a, b|c) = P(a|c) P(b|c) \quad a \perp b | c$$

More General

- How can we find if S_1 and S_2 are conditionally independent given E ?

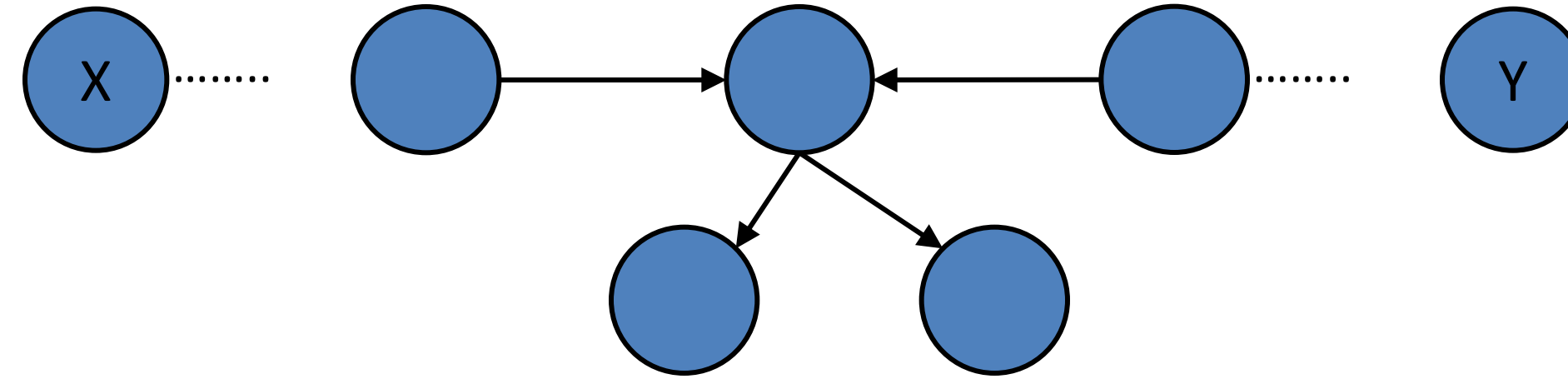
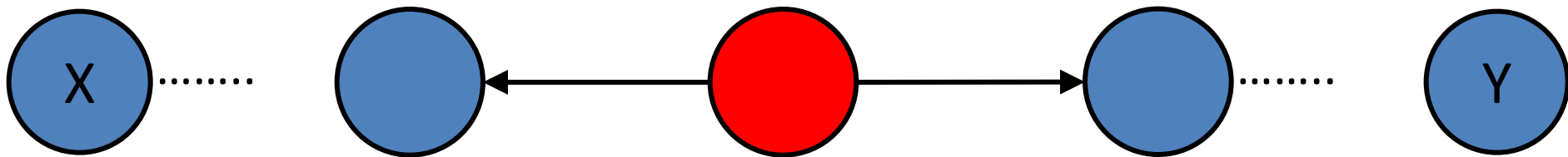
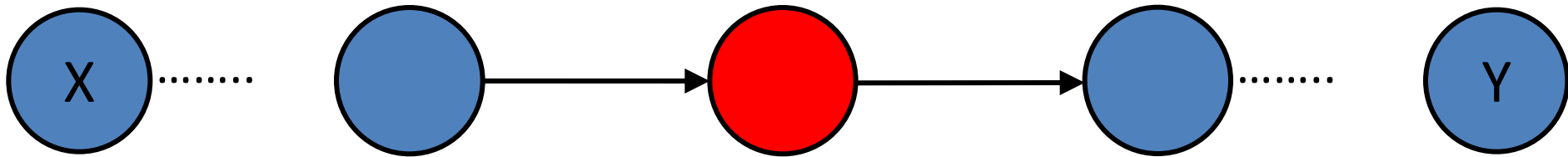


$$P(\text{assignments to } S_1 \mid E \text{ and assignments to } S_2) = P(\text{assignments to } S_1 \mid E)$$

- Why is it important and useful?

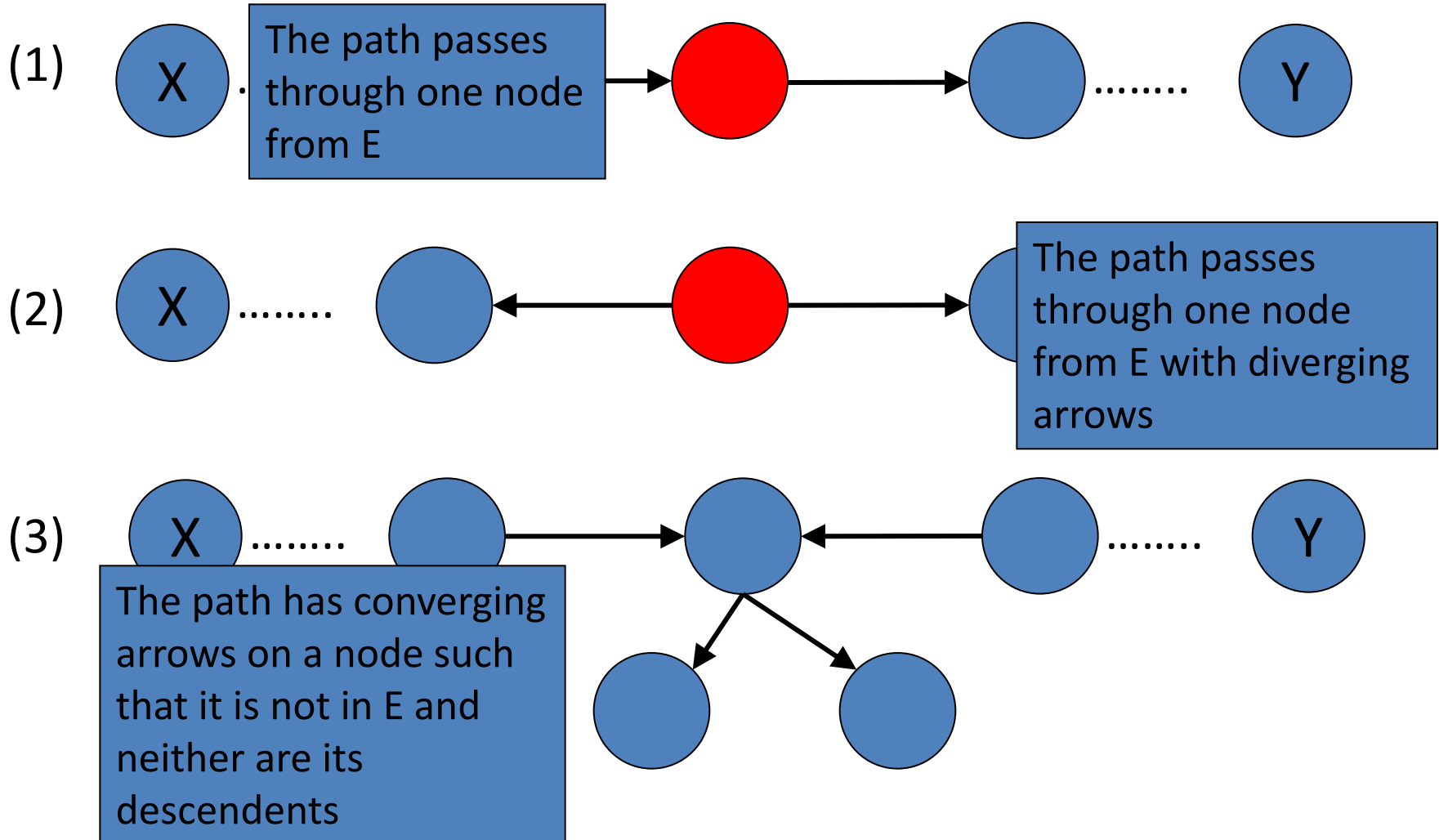
We can simplify any computation that contains something like $P(S_1 \mid E, S_2)$ by $P(S_1 \mid E)$

Intuitively E stands in between or “blocks” S_1 from S_2



Blockage: Formal Definition

- A path from a node X to a node Y is *blocked* by a set E if either:

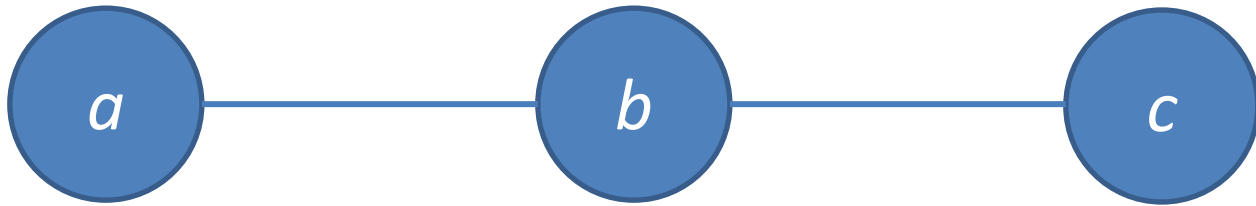


General case: Undirected

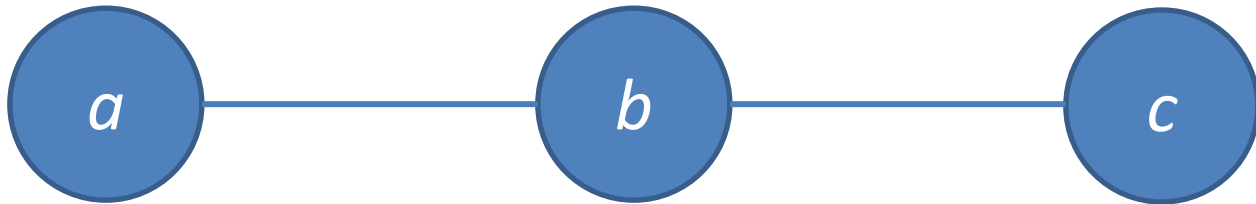


$$P(a, b) = \varphi(a, b)$$

$$P(a, b) = \varphi_1(a)\varphi_2(b)$$

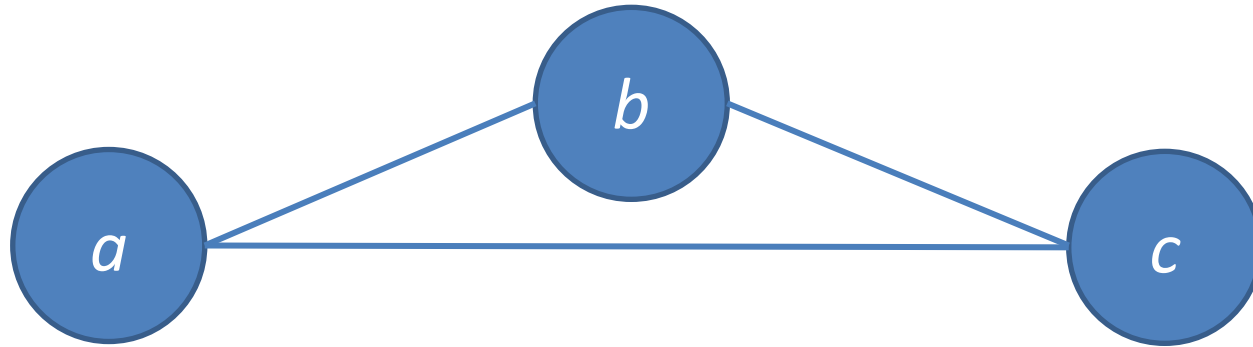


$$P(a, b, c) = \varphi_1(a, b) \varphi_2(b, c)$$



$$P(a, b, c) = \varphi_1(a, b) \varphi_2(b, c)$$

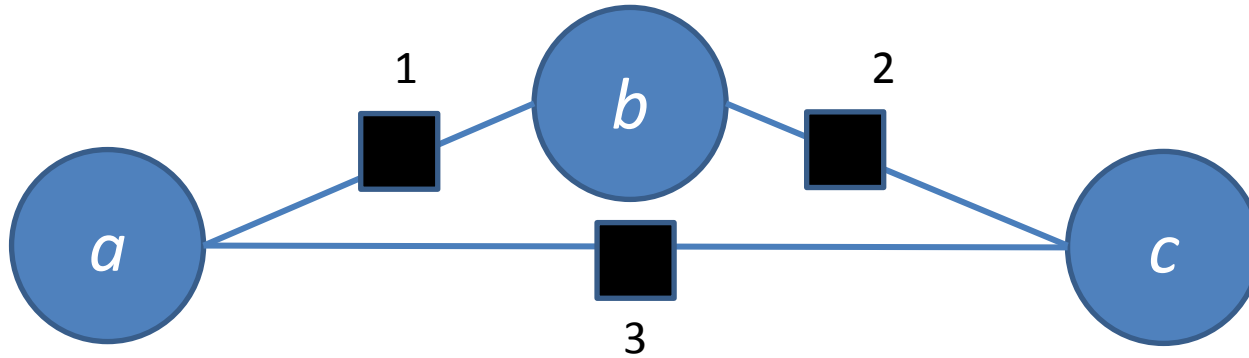
$a \perp c | b$ because all paths between a and c go through b



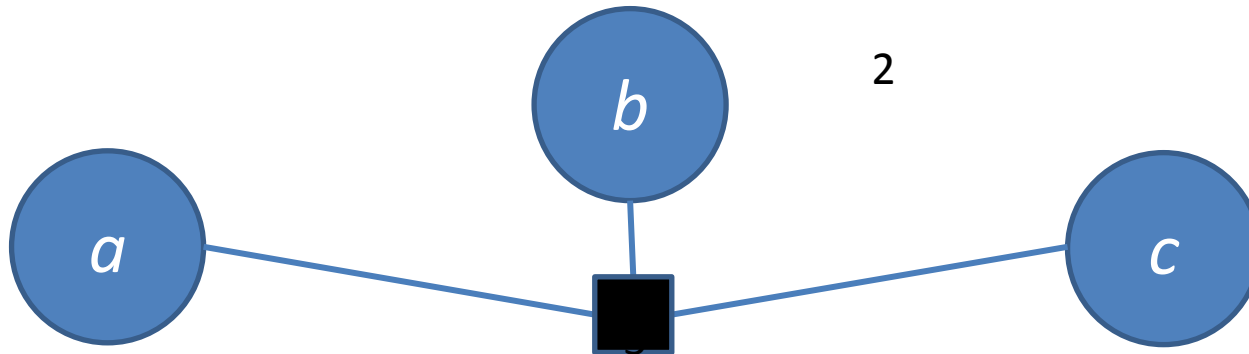
$$P(a, b, c) = \varphi_1(a, b) \varphi_2(b, c) \varphi_3(a, c)$$

$$P(a, b, c) = \varphi(a, b, c)$$

Factor graphs

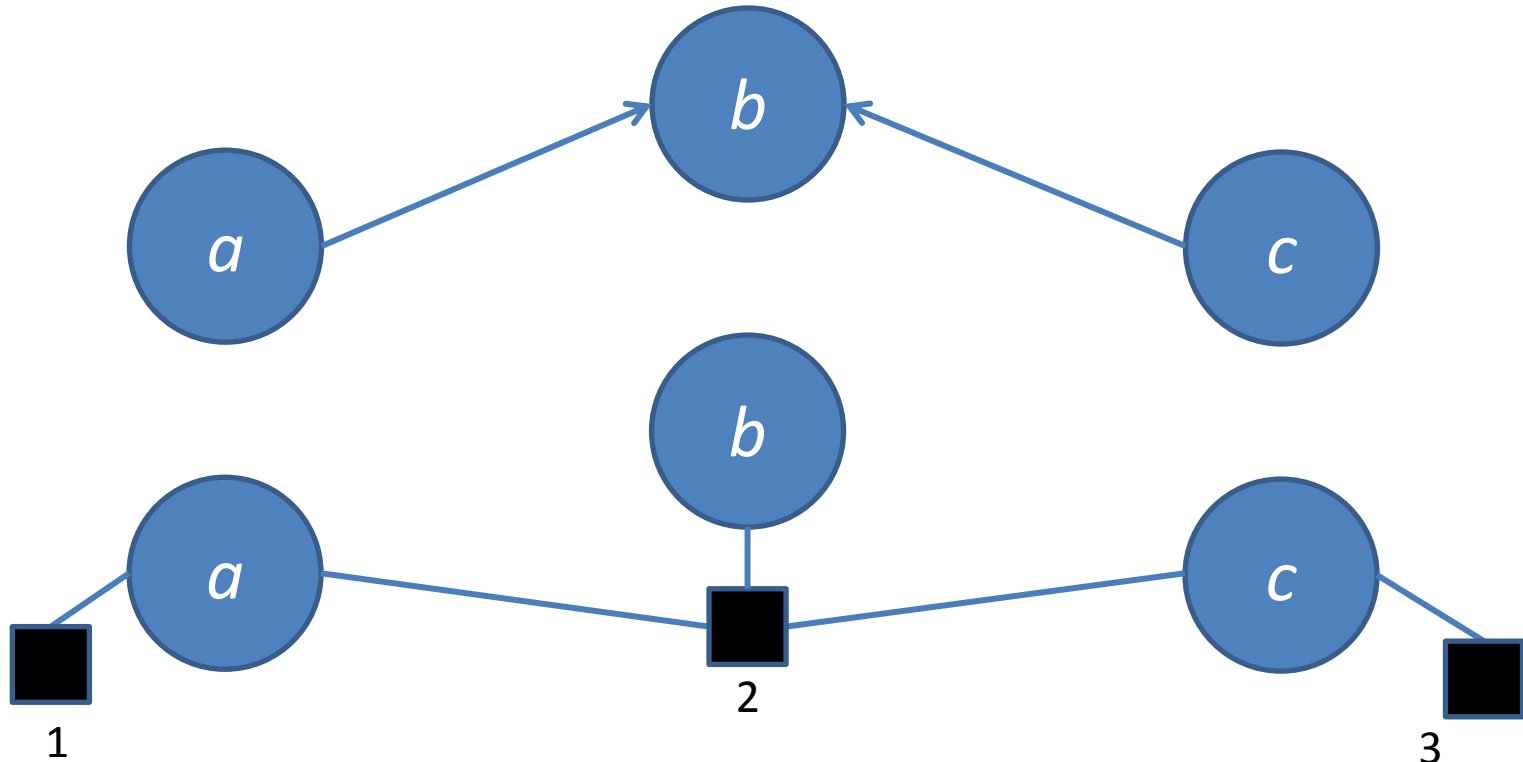


$$P(a, b, c) = \varphi_1(a, b) \varphi_2(b, c) \varphi_3(a, c)$$



$$P(a, b, c) = \varphi(a, b, c)$$

Directed vs. undirected



$$P(a, b, c) = \varphi_1(a, b) \varphi_2(b, c) \varphi_3(a, c)$$

$$\varphi_1(a) = P(a) \quad \varphi_3(c) = P(c)$$

$$\varphi_2(a, b, c) = P(b|a, c)$$



- N regions
- M possible labels
- Somehow, there is a way to estimate how likely a label is given image features $P(l_i|f)$
- We want to find the assignment of labels that optimizes $P(l_1, \dots, l_N|f)$



- Everything is independent:

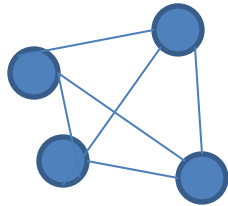


$$P(l_1, \dots, l_N | f) = \prod P(l_i | f)$$

Gives really stupid results because it does not take into account the distribution of likely relative occurrence of the labels



- Everything is dependent:

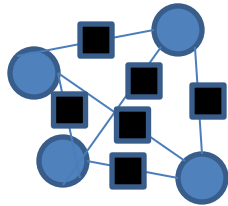


$$P(l_1, \dots, l_N | f) \propto \prod P(f | l_i) P(l_1, \dots, l_N)$$

Hard to learn or represent $P(l_1, \dots, l_N)$



- Factor pairwise dependencies:

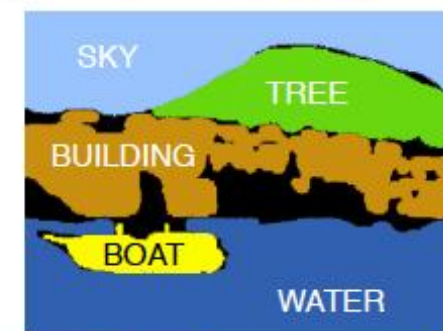
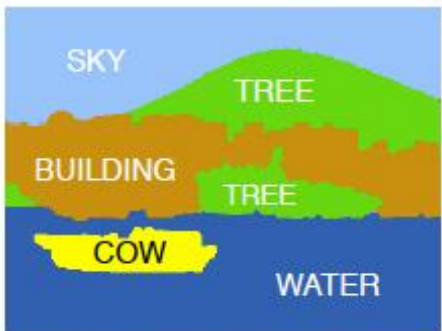
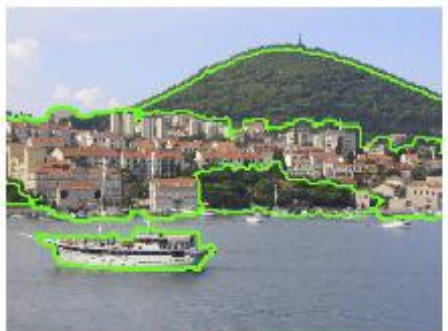
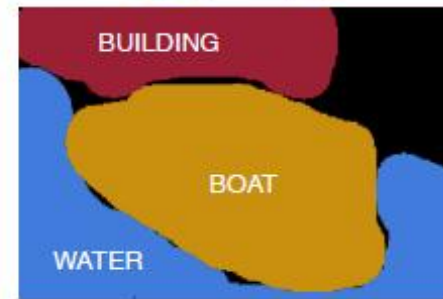


$$P(l_1, \dots, l_N) = \prod \varphi(l_i, l_j)$$

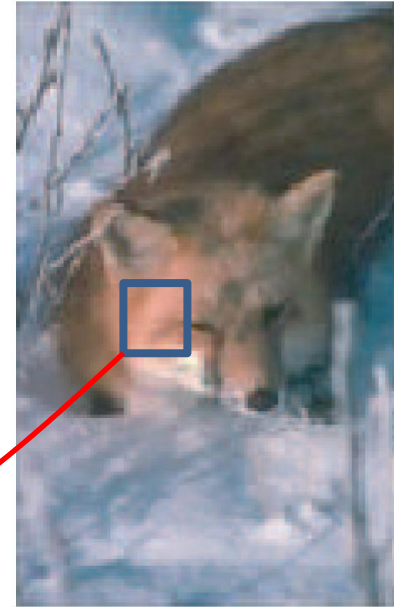
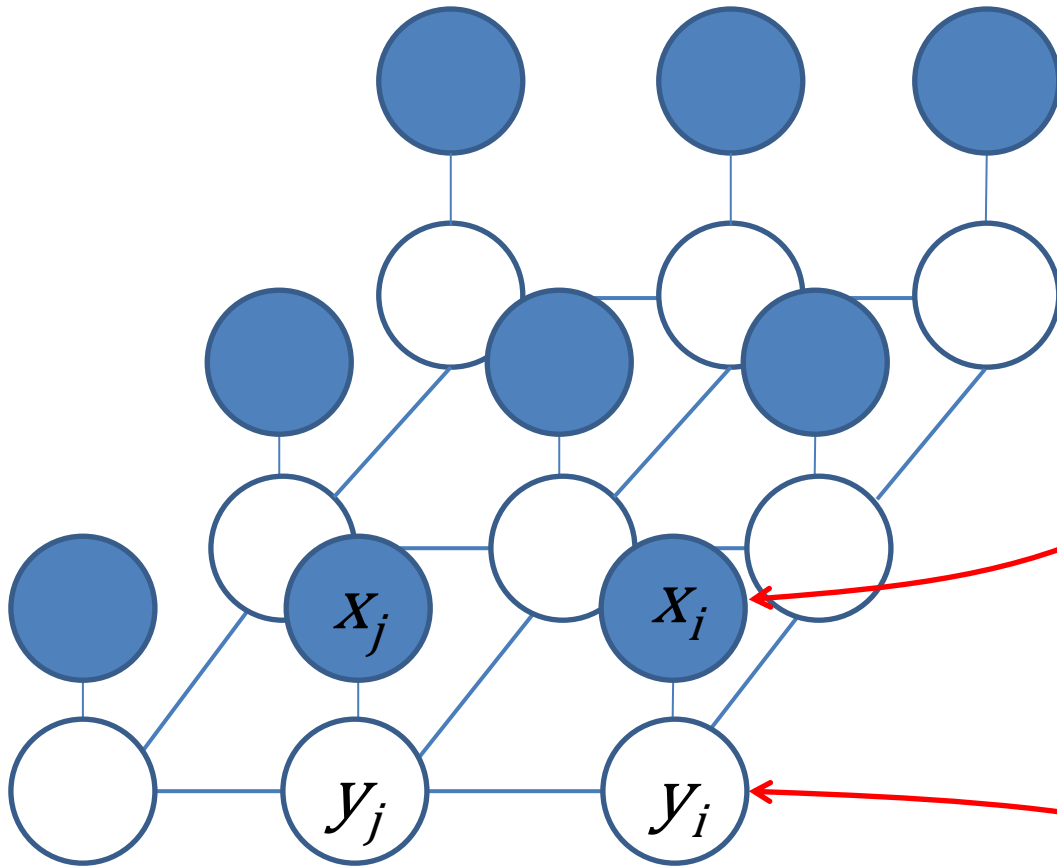
$\varphi(l_i, l_j)$ can be estimated from co-occurrence statistics from training data

MSRC training data

building	75	18	29			33	6	9	7	18	10		2	1		43		1	9	6	
grass	18	93	38	23	15	39	14	7	7		3	1		1		4	15		2	8	
tree	29	38	68	6		43	6	12	9	4		1	2			1	19			11	8
cow		23	6	23		4		4													
sheep			15		15				1										2		
sky	33	39	43	4		86	15	18	4	3			5	4		25				8	11
aeroplane	6	14	6			15	15									5					
water	9	7	12	4	1	18		43	4	1				7		6				6	16
face	7	7	9			4		4	28	1	1	1			3	7				28	1
car	18		4			3		1	1	20						19				1	
bike	10	3							1		15					12				1	
flower		1	1						1				1							1	
sign	2		2			5							8							1	1
bird	1	1				4		7						14		3					
book									3						3						3
chair		4	1													7	3				
road	43	15	19			2	25	5	6	7	19	12		1	3	3	86	7	10	8	1
cat																7	7				
dog	1	2														10			13		1
body	9	8	11			8		6	28	1	1	1	1		3	8				32	2
boat	6		8			11		18	1							1			1	2	19
	building	grass	tree	cow	sheep	sky	aeroplane	water	face	car	bike	flower	sign	bird	book	chair	road	cat	dog	body	boat



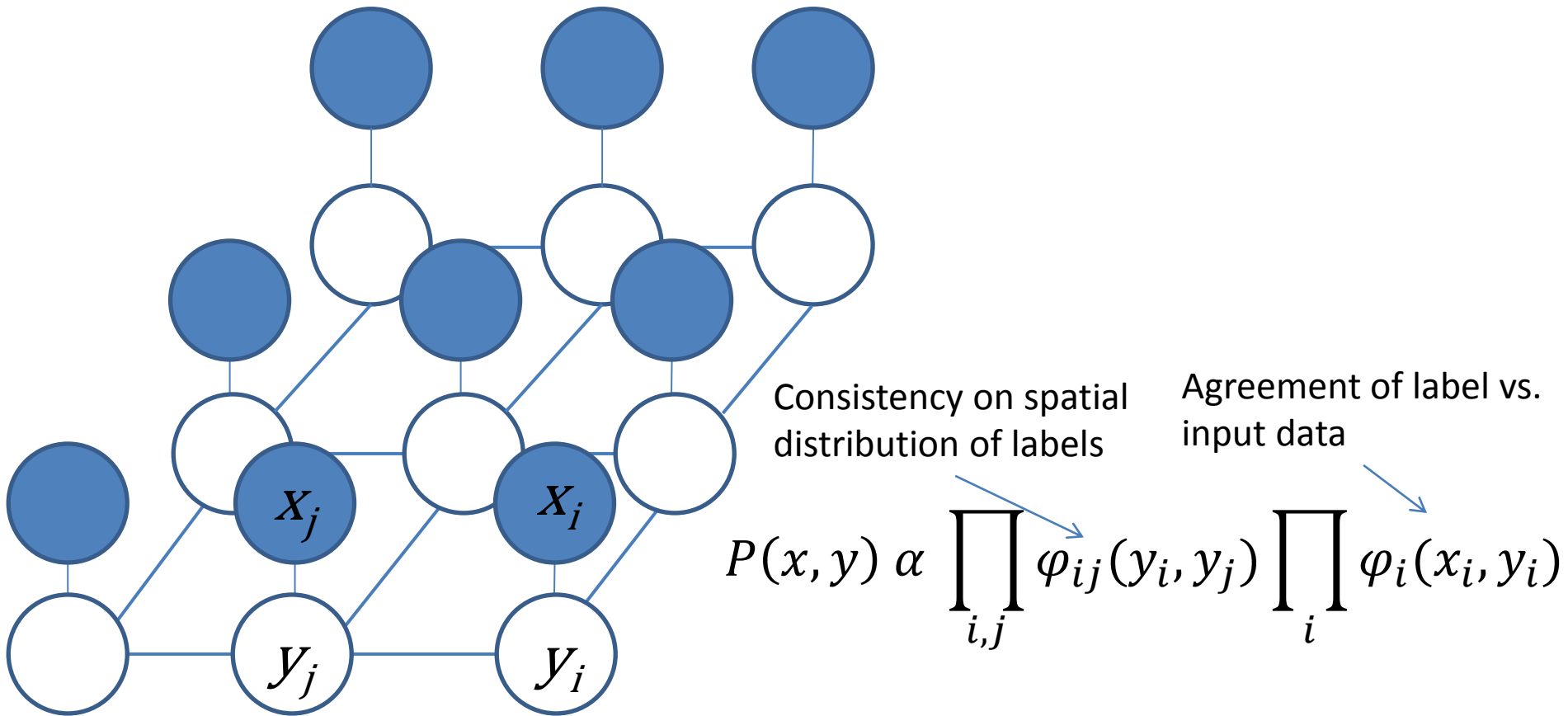
Example: MRF for image labeling



snow	snow	snow	fox
snow	fox	fox	fox
fox	fox	fox	fox
fox	fox	fox	fox
snow	snow	snow	snow
snow	snow	snow	snow

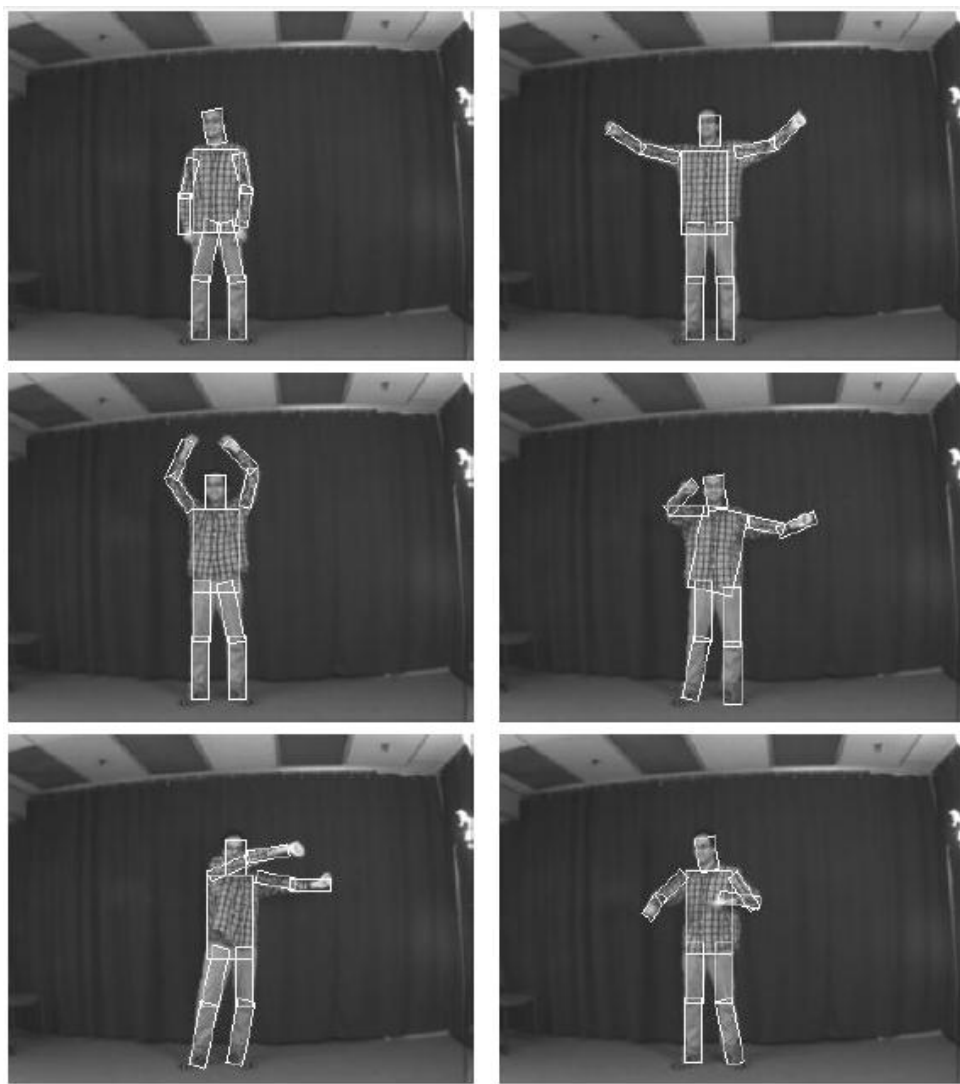
- y = class label (e.g., road, car, etc...)
- x = image data at local patch

Example: MRF for image labeling



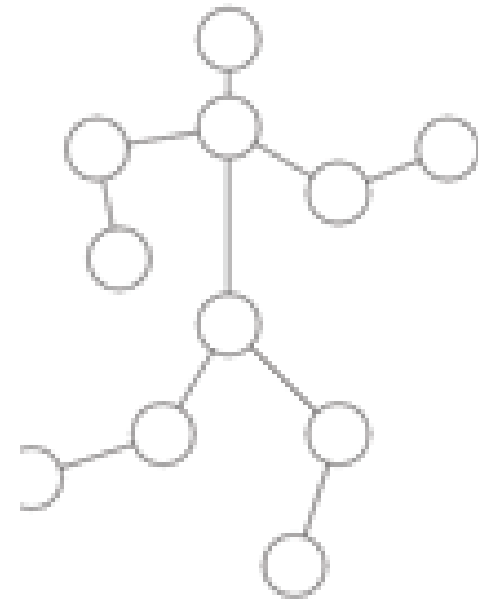
- y = class label (e.g., road, car, etc...)
- x = image data at local patch

Example: Inferring human poses



x_i = Input image data at limb i

y_i = Pose (location and orientation) of limb i



Example from Felzenszwalb'04

- Note: Efficient because tree-structured