

15-780: Grad AI

Lecture 19: Graphical models, Monte Carlo methods

Geoff Gordon (this lecture)

Tuomas Sandholm

TAs Erik Zawadzki, Abe Othman

Admin



- Reminder: midterm March 29
- Reminder: project milestone reports due March 31

Review: scenarios

- Converting QBF+ to PBI/MILP by scenarios
 - ▶ Replicate decision variables for each scenario
 - ▶ Replicate clauses: share first stage vars; set scenario vars by scenario index; replace decision vars by replicates
 - ▶ Sample random scenarios
- Example: PSTRIPS

Review: dynamic programming

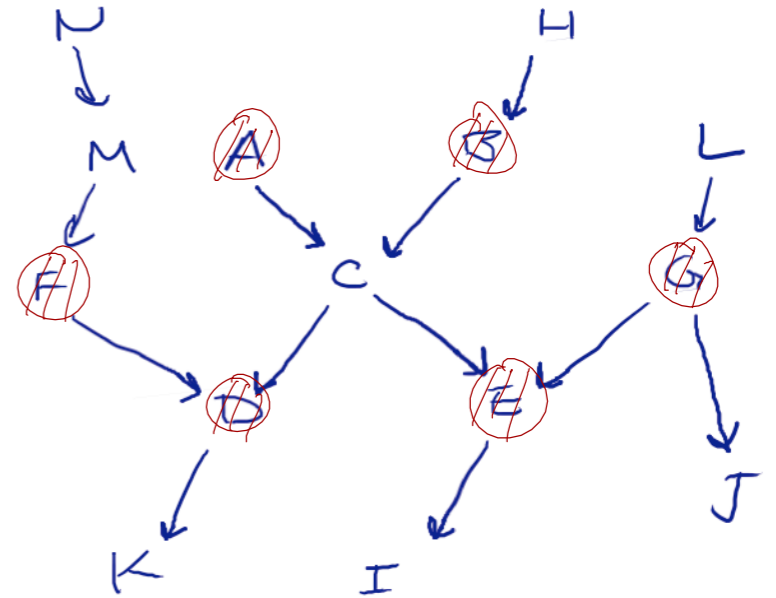
- Solving #SAT by dynamic programming (variable elimination)
 - ▶ repeatedly move sums inward, combine tables, sum out
 - ▶ treewidth and runtime/space

Review: graphical models

- Bayes net = DAG + CPTs
 - ▶ For each RV (say X), there is one CPT specifying $P(X \mid \text{pa}(X))$
 - ▶ Can simulate with propositional logic + random causes
- Inference: similar to #SAT DP—move sums inward
 - ▶ Can do partly analytically
 - ▶ Allows us to prove independences and conditional ind's from DAG alone

Review: graphical models

- Blocking, explaining away
- Markov blanket
- Learning: counting, Laplace smoothing
 - ▶ if hidden variables: take 10-708 or use a toolbox

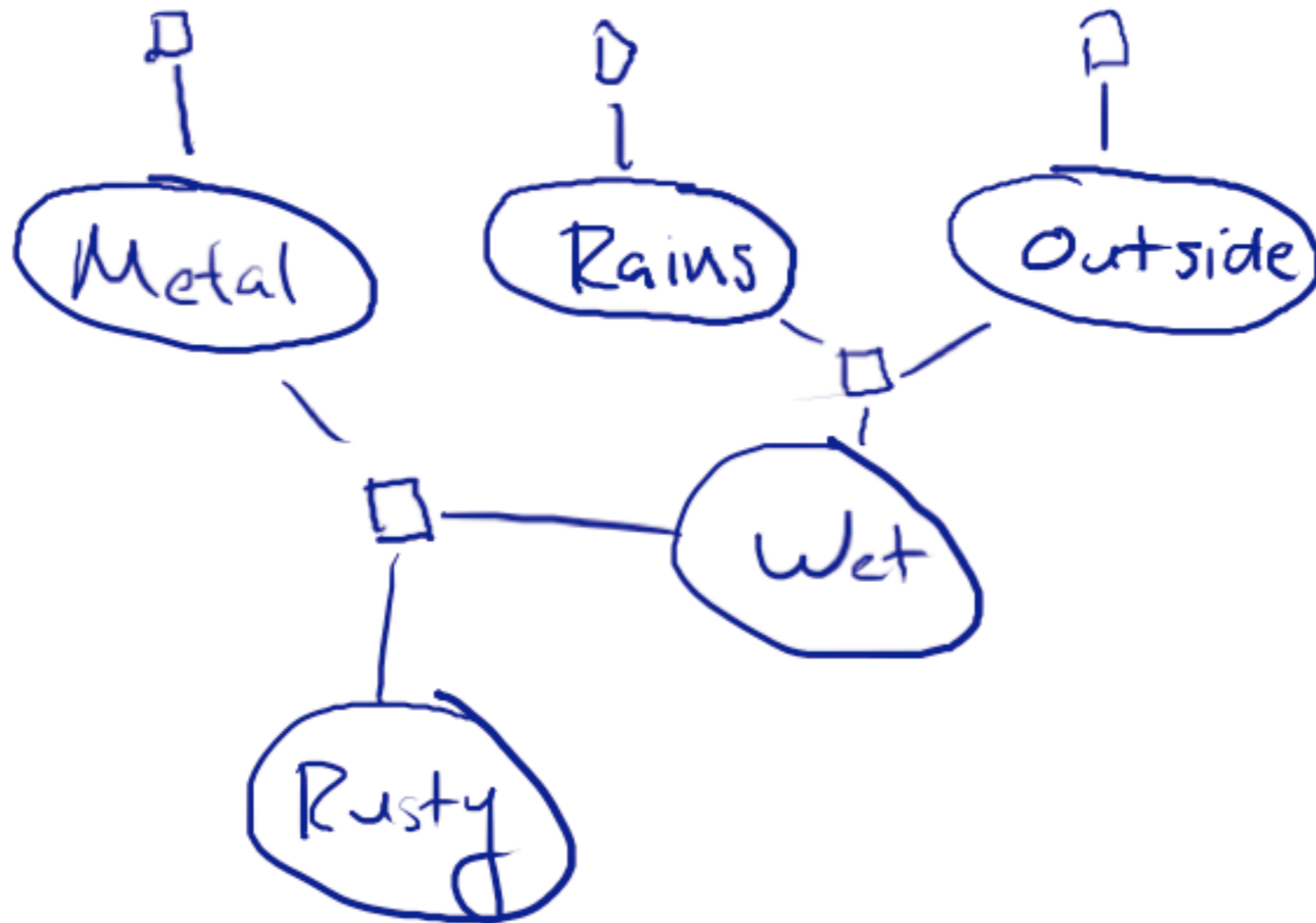


Factor graphs



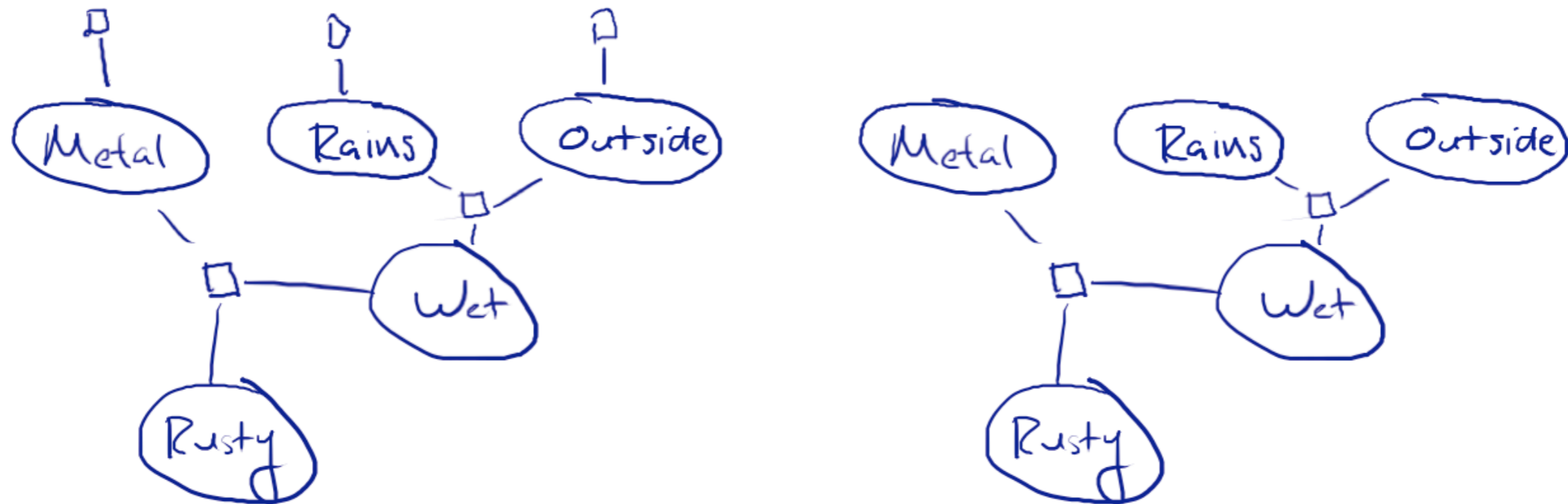
- Another common type of graphical model
- Uses ***undirected, bipartite*** graph instead of DAG

Rusty robot: factor graph



$$P(M) P(Ra) P(O) P(W|Ra, O) P(Ru|M, W)$$

Convention



- Don't need to show unary factors
- Why? They don't affect algorithms below.

Non-CPT factors

- Just saw: easy to convert Bayes net \rightarrow factor graph
- In general, factors need not be CPTs: any nonnegative #s allowed
- In general, $P(A, B, \dots) =$

- $Z =$

Hard v. soft factors

Hard

X

	0	1	2
0	0	0	0
1	0	0	1
2	0	1	1

Soft

X

	0	1	2
0	1	1	1
1	1	1	3
2	1	3	3

Factor graph \rightarrow Bayes net

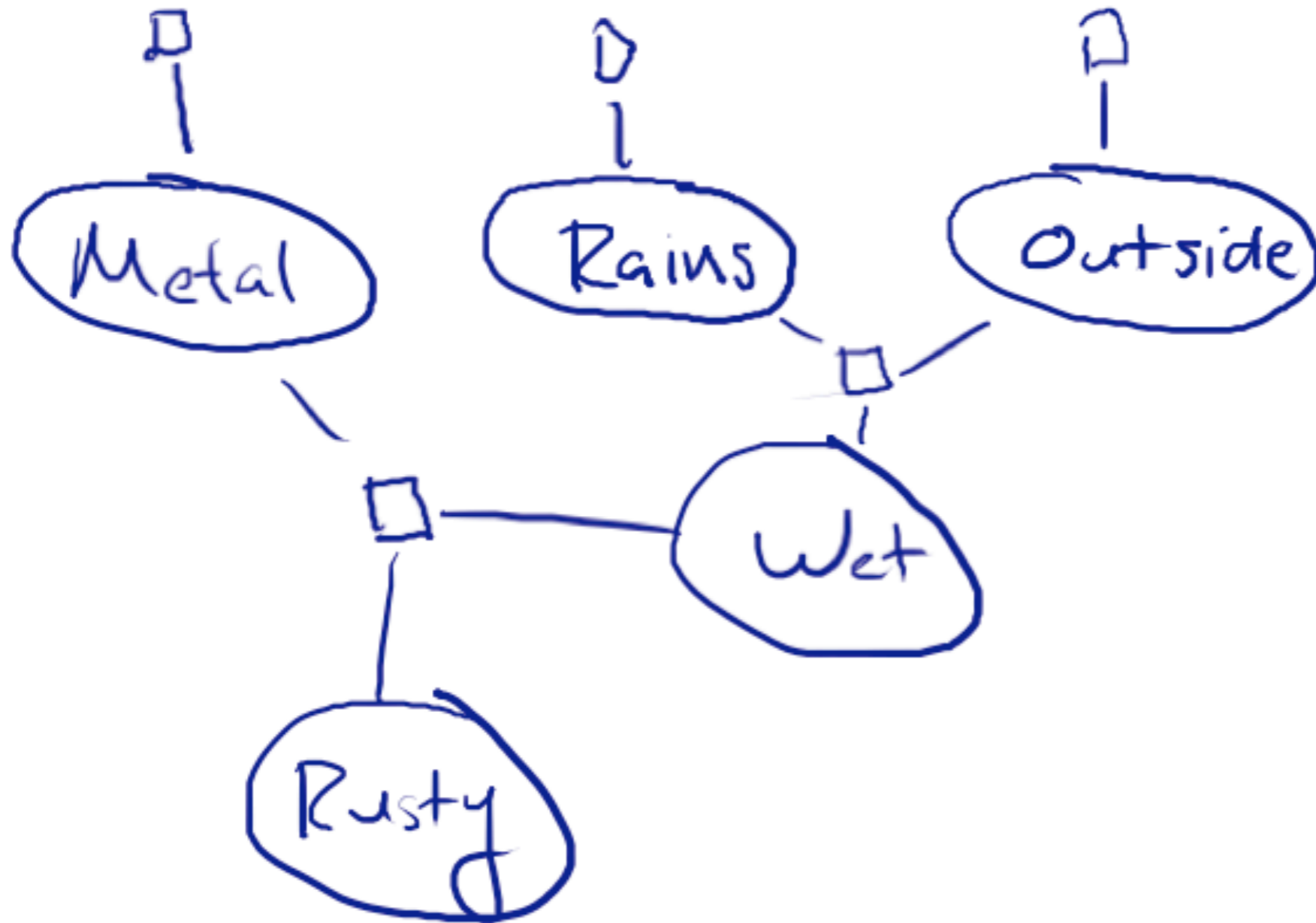
- Conversion possible, but more involved
 - ▶ Each representation can handle **any** distribution
 - ▶ But, size/complexity of graph may differ
- 2 cases for conversion:
 - ▶ without adding nodes:
 - ▶ adding nodes:

Independence



- Just like Bayes nets, there are graphical tests for independence and conditional independence
- Simpler, though:
 - ▶ Cover up all observed nodes
 - ▶ Look for a path

Independence example



Modeling independence

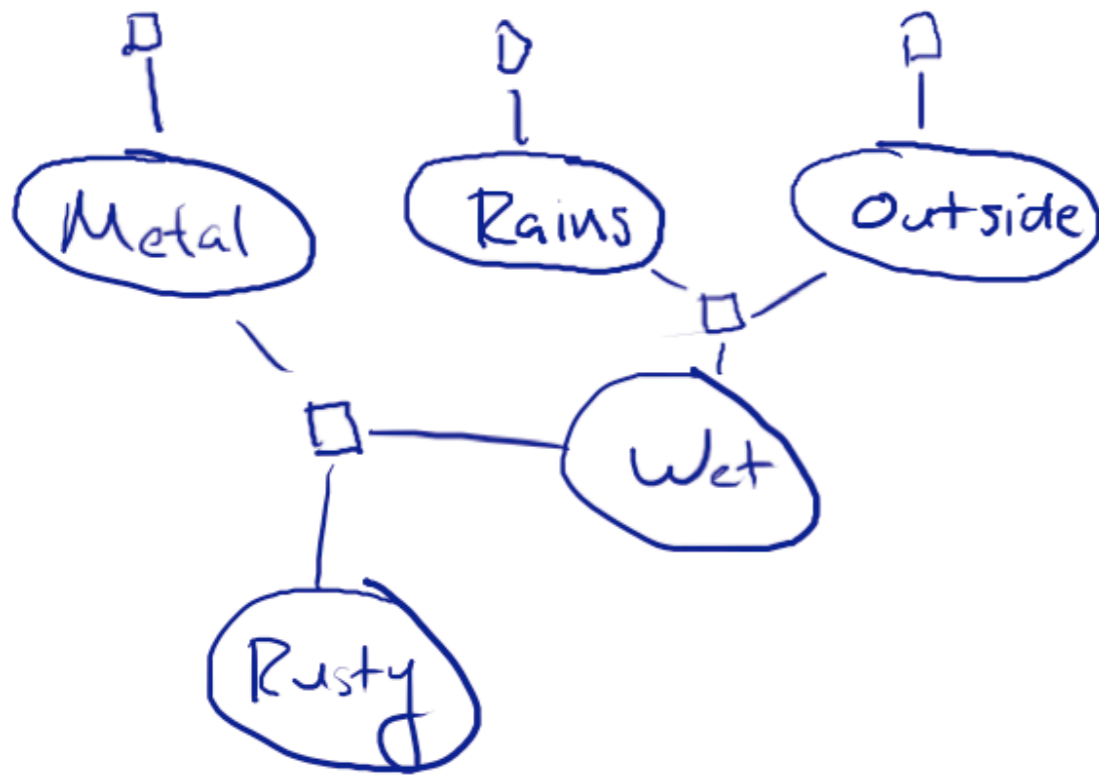
- Take a Bayes net, list the (conditional) independences
- Convert to a factor graph, list the (conditional) independences
- Are they the same list?
- What happened?

Inference

- Inference: prior + evidence \rightarrow posterior
- We gave examples of inference in a Bayes net, but not a general algorithm
- Reason: general algorithm uses factor-graph representation
- Steps: instantiate evidence, eliminate nuisance nodes, normalize, answer query

Inference

$$P(M, R_a, O, W, R_u) = \phi_1(M) \phi_2(R_a) \phi_3(O) \phi_4(R_a, O, W) \phi_5(M, W, R_u) / Z$$



$$\phi_1(M) = \begin{matrix} T & 0.9 \\ F & 0.1 \end{matrix}$$

$$\phi_2(R_a) = \begin{matrix} T & 0.7 \\ F & 0.3 \end{matrix}$$

$$\phi_3(O) = \begin{matrix} T & 0.2 \\ F & 0.8 \end{matrix}$$

$$\phi_4(R_a, O, W) =$$

$$TTT \quad 0.9$$

$$TTF \quad 0.1$$

$$TFT \quad 0.1$$

$$TFF \quad 0.9$$

$$FTT \quad 0.1$$

$$FTF \quad 0.9$$

$$FFT \quad 0.1$$

$$FFF \quad 0.9$$

$$\phi_5(M, W, R_u) =$$

$$TTT \quad 0.8$$

$$TTF \quad 0.2$$

$$TFT \quad 0.1$$

$$TFF \quad 0.9$$

$$FTT \quad 0$$

$$FTF \quad 1$$

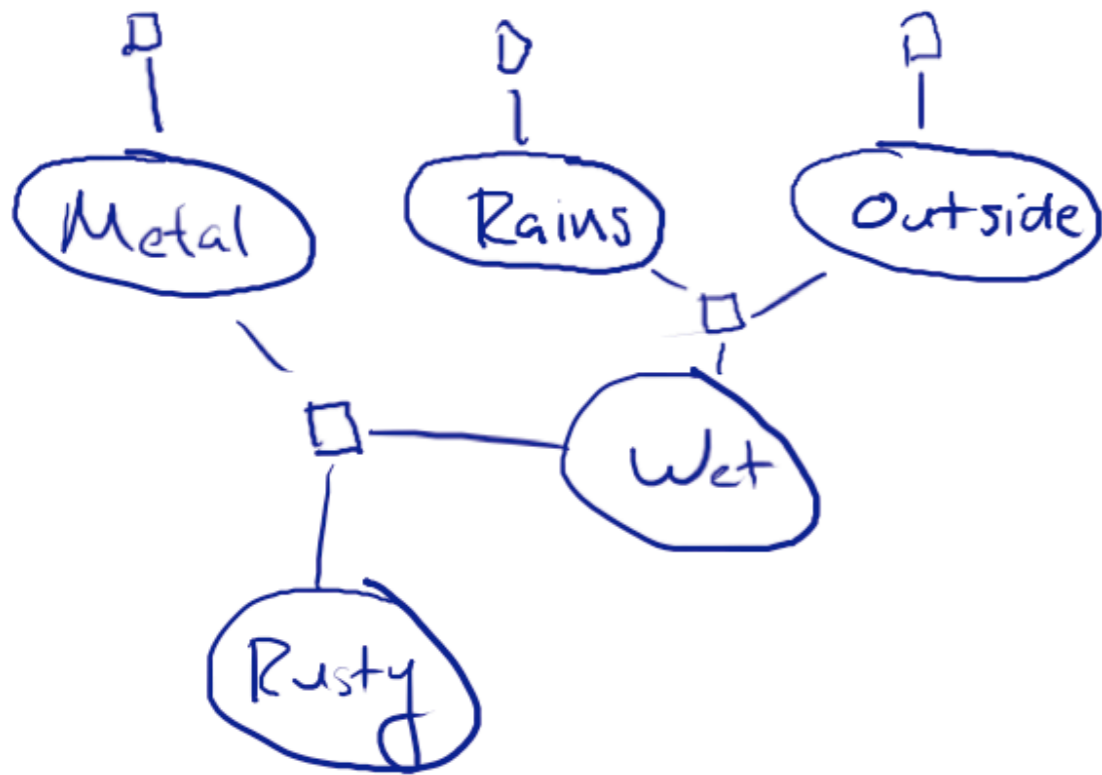
$$FFT \quad 0$$

$$FFF \quad 1$$

- Typical Q: given $R_a=F$, $R_u=T$, what is $P(W)$?

Incorporate evidence

$$P(M, R_a, O, W, R_u) = \phi_1(M) \phi_2(R_a) \phi_3(O) \phi_4(R_a, O, W) \phi_5(M, W, R_u) / Z$$



$$\phi_1(M) = \begin{matrix} T & 0.9 \\ F & 0.1 \end{matrix}$$

$$\phi_2(R_a) = \begin{matrix} T & 0.7 \\ F & 0.3 \end{matrix}$$

$$\phi_3(O) = \begin{matrix} T & 0.2 \\ F & 0.8 \end{matrix}$$

$$\phi_4(R_a, O, W) =$$

$$TTT \quad 0.9$$

$$TFE \quad 0.1$$

$$TFT \quad 0.1$$

$$TFE \quad 0.9$$

$$FTT \quad 0.1$$

$$FTF \quad 0.9$$

$$FFT \quad 0.1$$

$$FFF \quad 0.9$$

$$\phi_5(M, W, R_u) =$$

$$TTT \quad 0.8$$

$$TFE \quad 0.2$$

$$TFT \quad 0.1$$

$$TFE \quad 0.9$$

$$FTT \quad 0$$

$$FTF \quad 1$$

$$FFT \quad 0$$

$$FFF \quad 1$$

Condition on $R_a=F, R_u=T$

Eliminate nuisance nodes

$$P(M, R, O, W, P) = \phi_1(M) \phi_2(R) \phi_3(O) \phi_4(R, O, W) \phi_5(M, W, P) / Z$$

- Remaining nodes: M, O, W
- Query: P(W)
- So, O&M are nuisance—marginalize away
- Marginal =

Elimination order

$$\sum_M \sum_O \phi_1(\mu) \phi_3(O) \phi_4(O, \omega) \phi_5(\mu, \omega) / Z$$

- Sum out the nuisance variables in turn
- Can do it in any order, but some orders may be easier than others
- Let's do O, then M

$$\phi_3(O) = \begin{matrix} T & 0.2 \\ F & 0.8 \end{matrix}$$

$$\phi_4(\cancel{O}, O, \omega) =$$

T	TT	0.1
F	TF	0.9
T	FT	0.1
F	FF	0.9

One last elimination

$$\phi_1(M) = \begin{array}{l} T \ 0.9 \\ F \ 0.1 \end{array}$$

$$\phi_6(\omega) = \begin{array}{l} T \ 0.1 \\ F \ 0.9 \end{array}$$

$$\phi_5(M, \omega, \omega) =$$

T	T	T	0.8
T	F	T	0.1
F	T	T	0
F	F	T	0

Checking our work



- <http://www.aispace.org/bayes/version5.1.6/bayes.jnlp>

Discussion

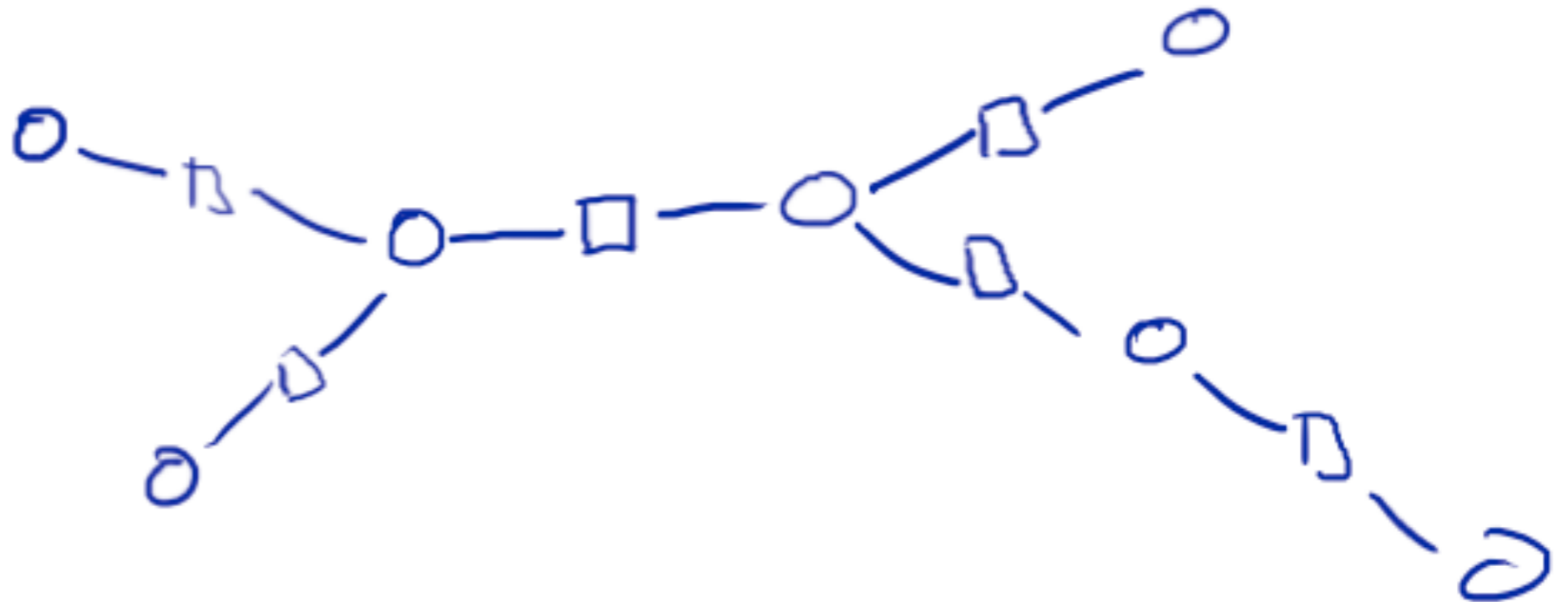
- Steps: instantiate evidence, eliminate nuisance nodes, normalize, answer query
 - ▶ each elimination introduces a new table, makes some old tables irrelevant
- Normalization
- Each elim. order introduces different tables
 - ▶ some tables bigger than others
- FLOP count; treewidth

Treewidth examples

Chain

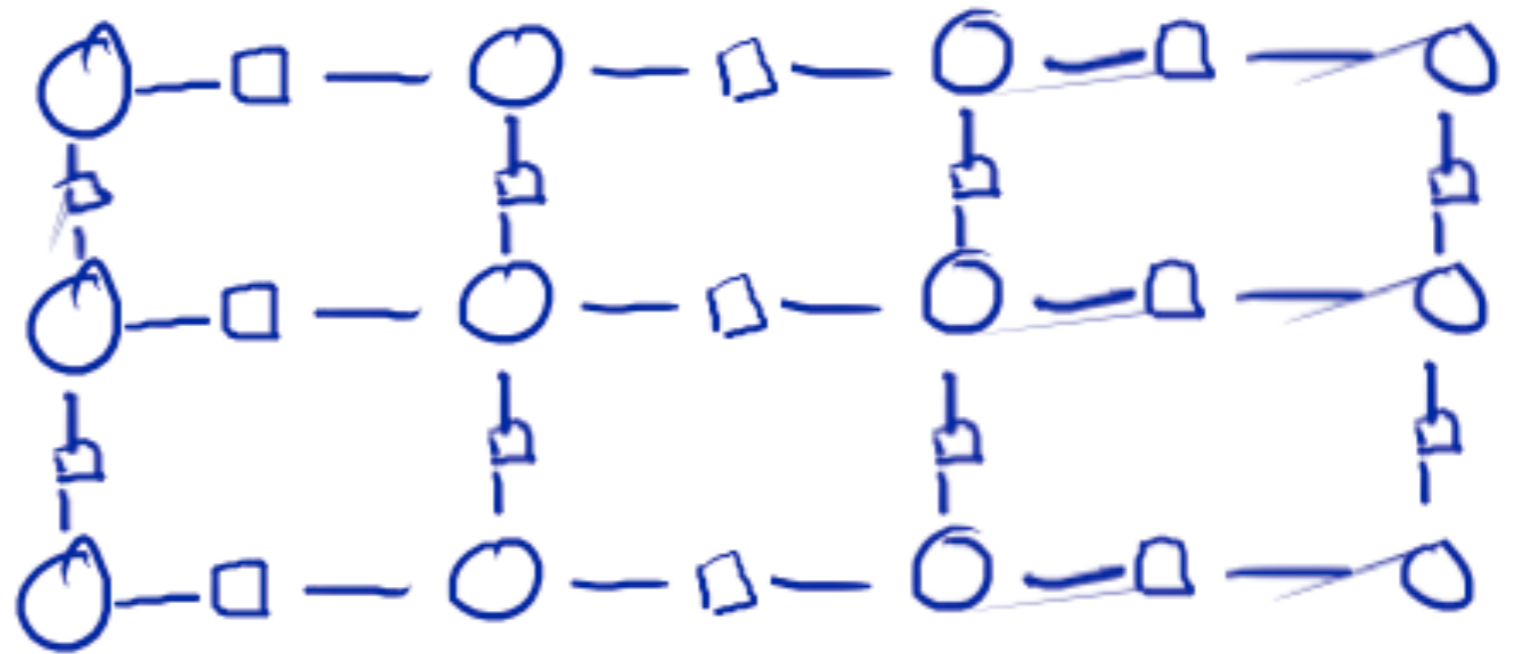


Tree

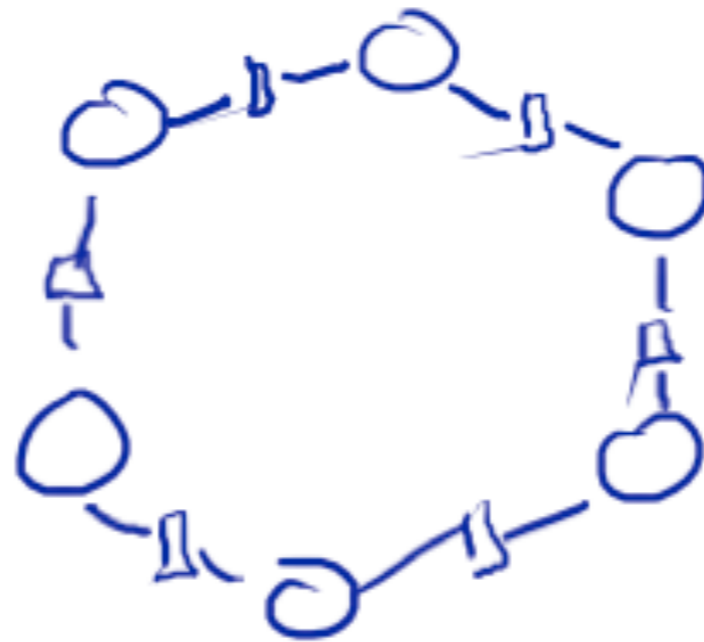


Treewidth examples

Parallel chains



Cycle



Discussion

- Several relationships between GMs and logic (similar DP algorithm, use of independent choices + logical consequences to represent a GM, factor graph with 0-1 potentials = CSP, MAP assignment = ILP)
- Directed v. undirected: advantages to both
- Lifted reasoning
 - ▶ Propositional logic + objects = FOL
 - ▶ FO GMs are a current hot topic of research (plate models, MLNs, ICL)—not solved yet!

Discussion: belief propagation

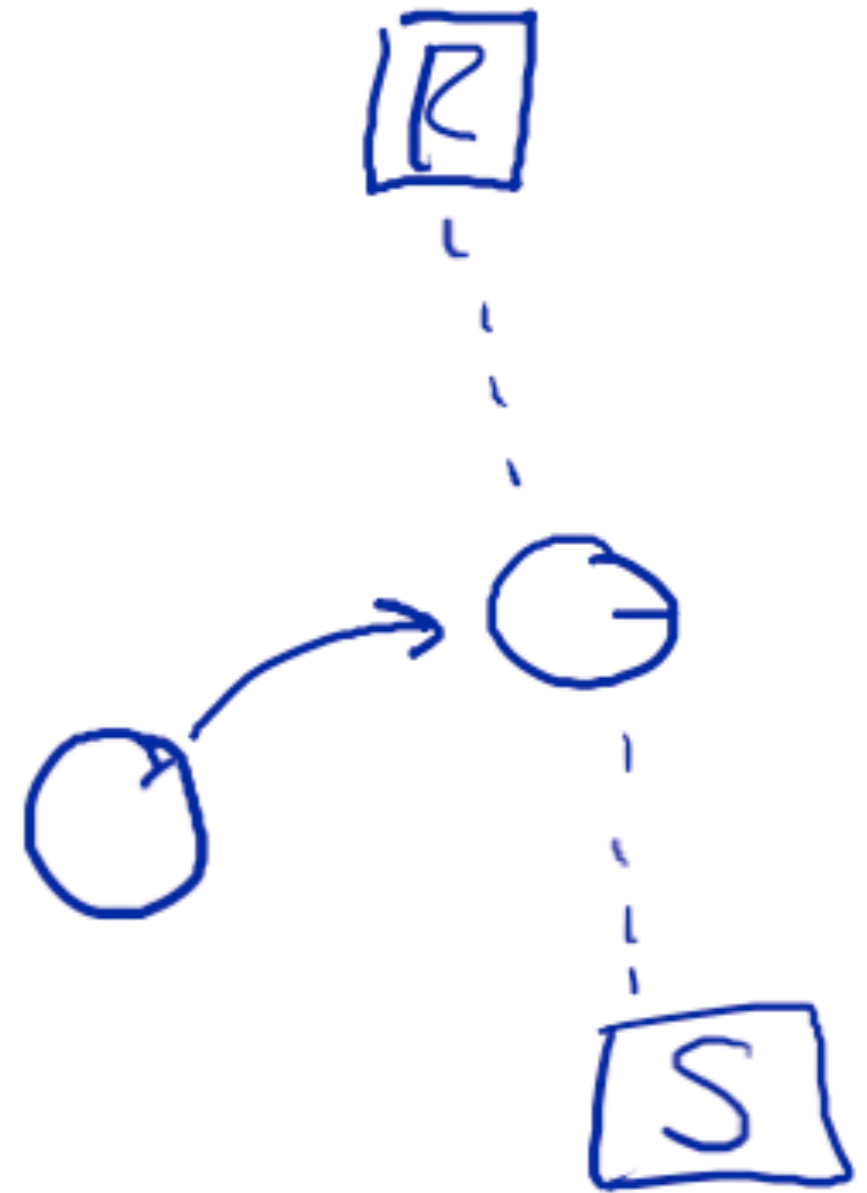
- Suppose we want all 1-variable marginals
- Could do N runs of variable elimination
- Or: the BP algorithm simulates N runs for the price of 2
- For details: Kschischang et al. reading



HMMs and DBNs

Inference over time

- Consider a robot:
 - ▶ true state (x, y, θ)
 - ▶ controls (v, w)
 - ▶ N range sensors (here $N=2: r, s$)



Model

$$x_{t+1} = x_t + v_t \cos \theta_t + \text{noise}$$

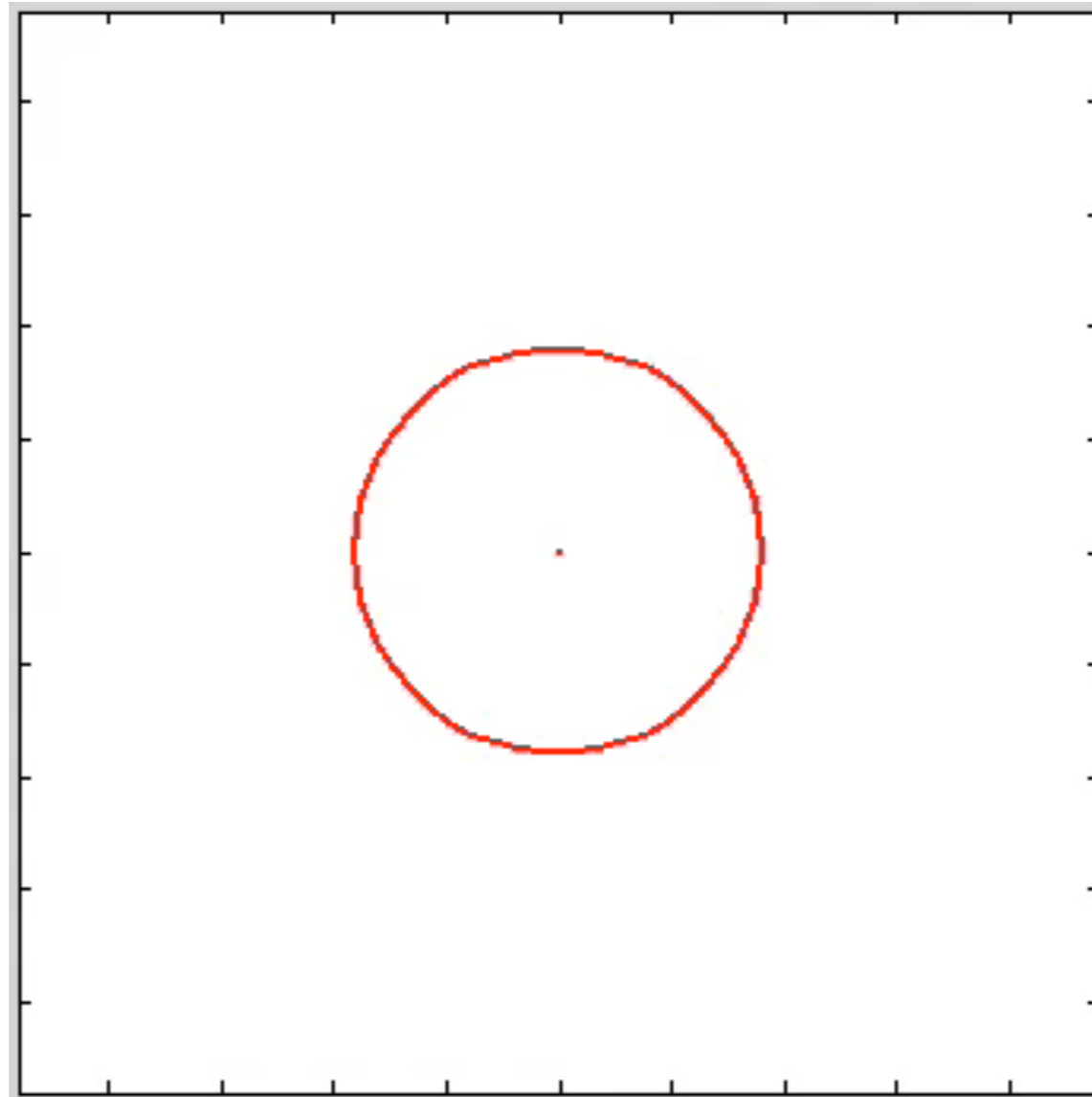
$$y_{t+1} = y_t + v_t \sin \theta_t + \text{noise}$$

$$\theta_{t+1} = \theta_t + w_t + \text{noise}$$

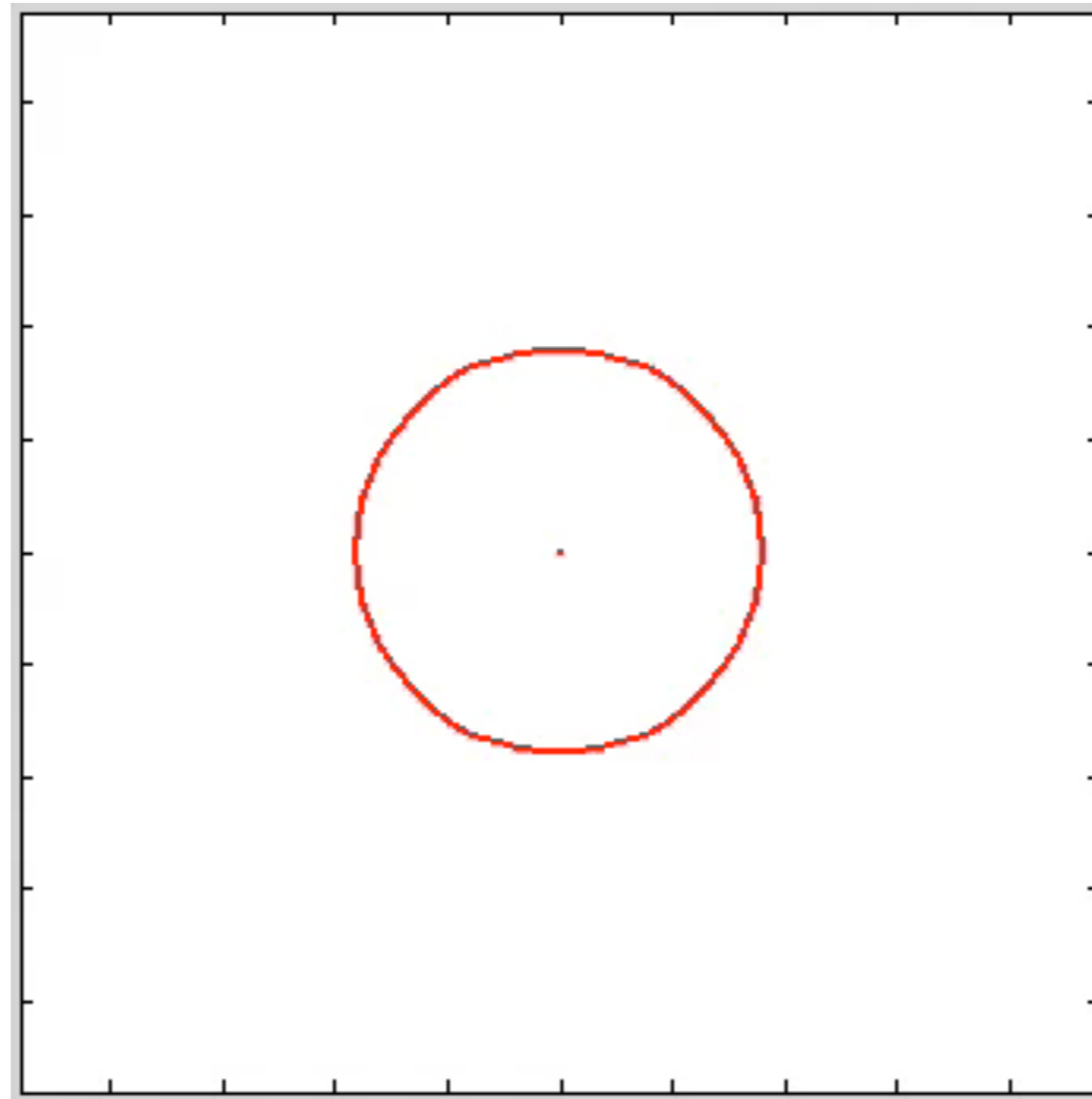
$$r_t = \sqrt{(x_t - x^R)^2 + (y_t - y^R)^2} + \text{noise}$$

$$s_t = \sqrt{(x_t - x^S)^2 + (y_t - y^S)^2} + \text{noise}$$

Model of x, y, θ (r, s unobserved)

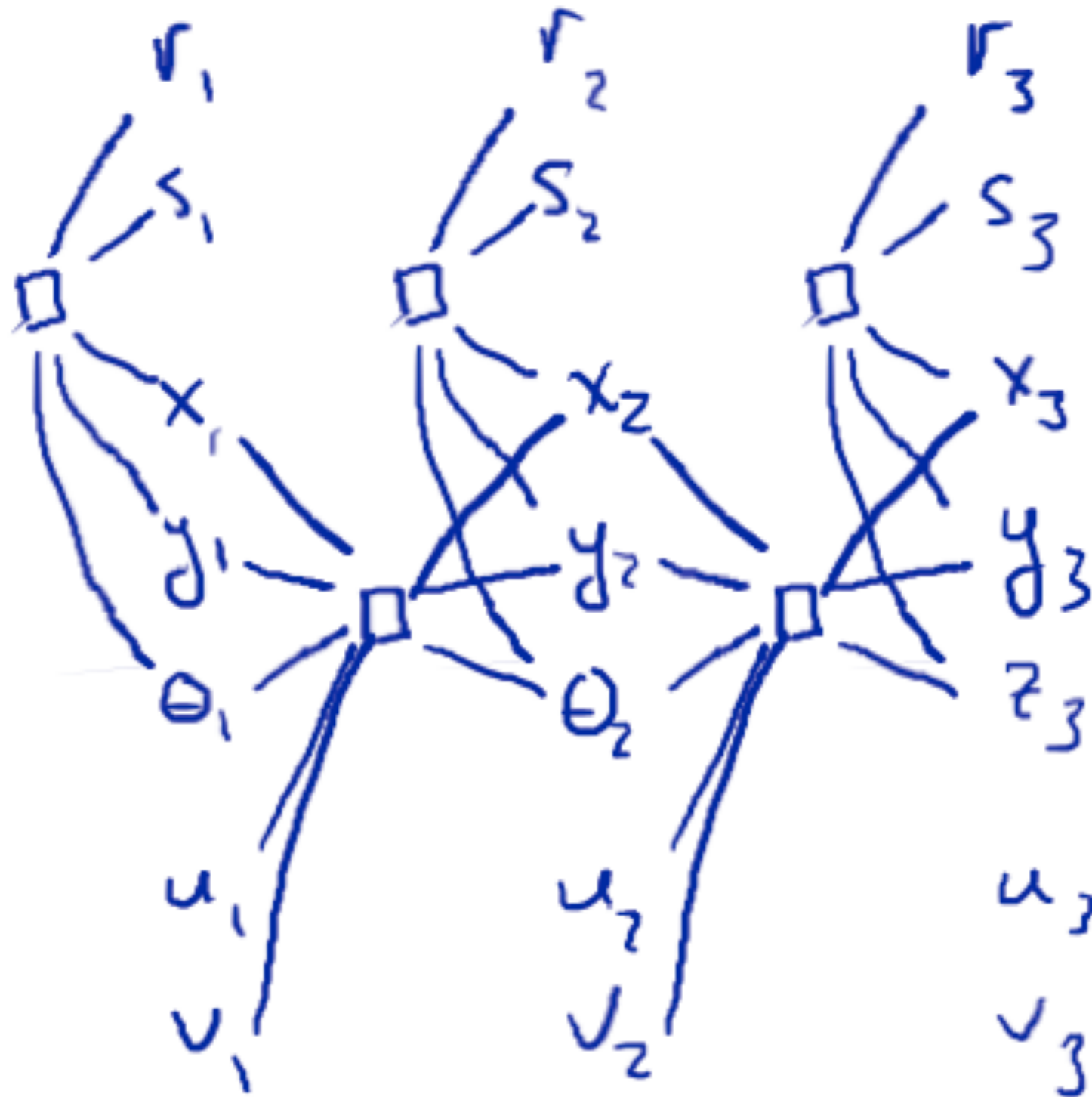


Goal: inference over time



- $N=1$ sensor, repeatedly observe range = $l_m + \text{noise}$

Factor graph



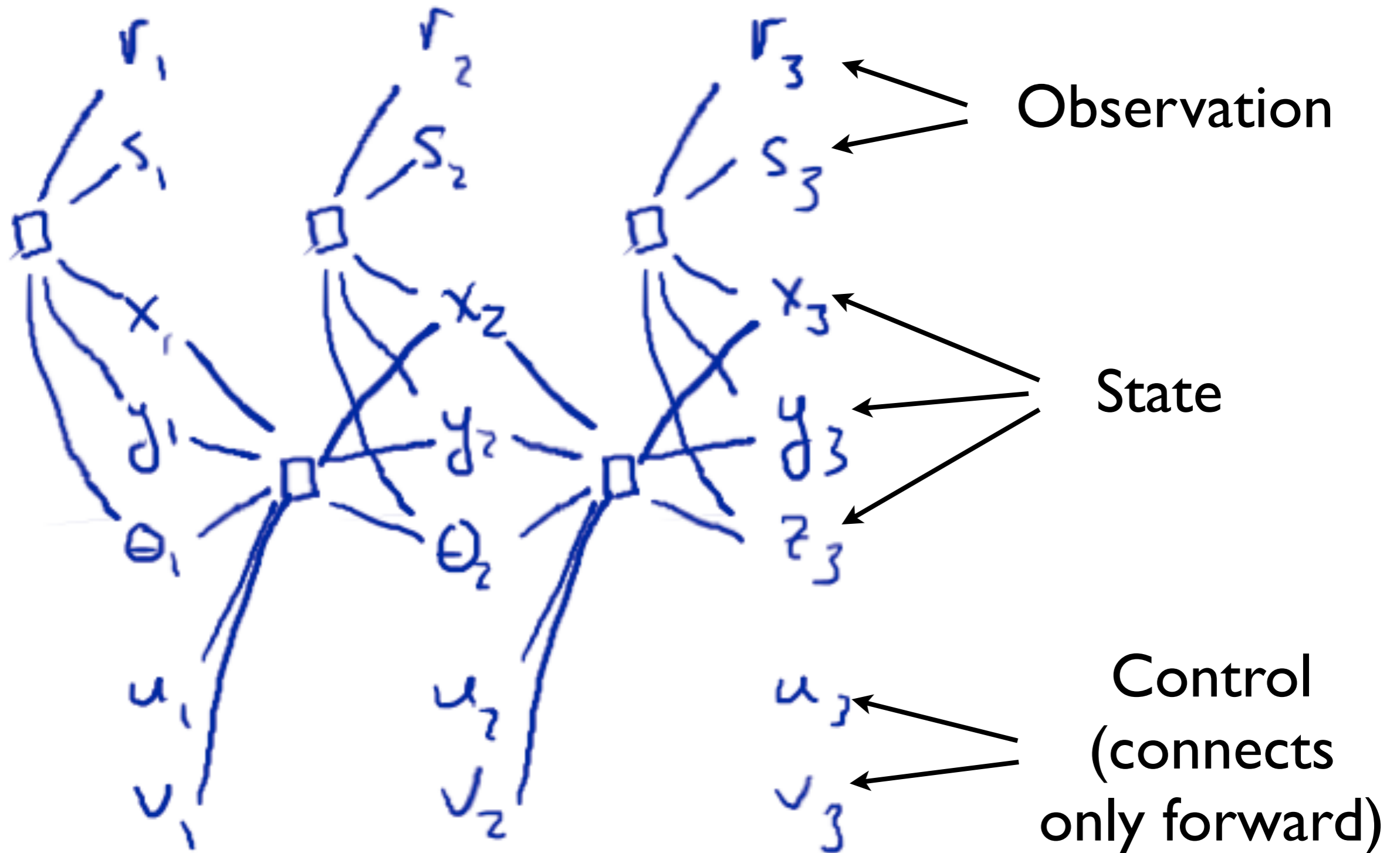
Dynamic Bayes Network

- DBN: factor graph composed of a single structural unit repeated over time
 - ▶ conceptually infinite to right, but in practice cut off at some maximum T
- Factors **must** be conditional distributions

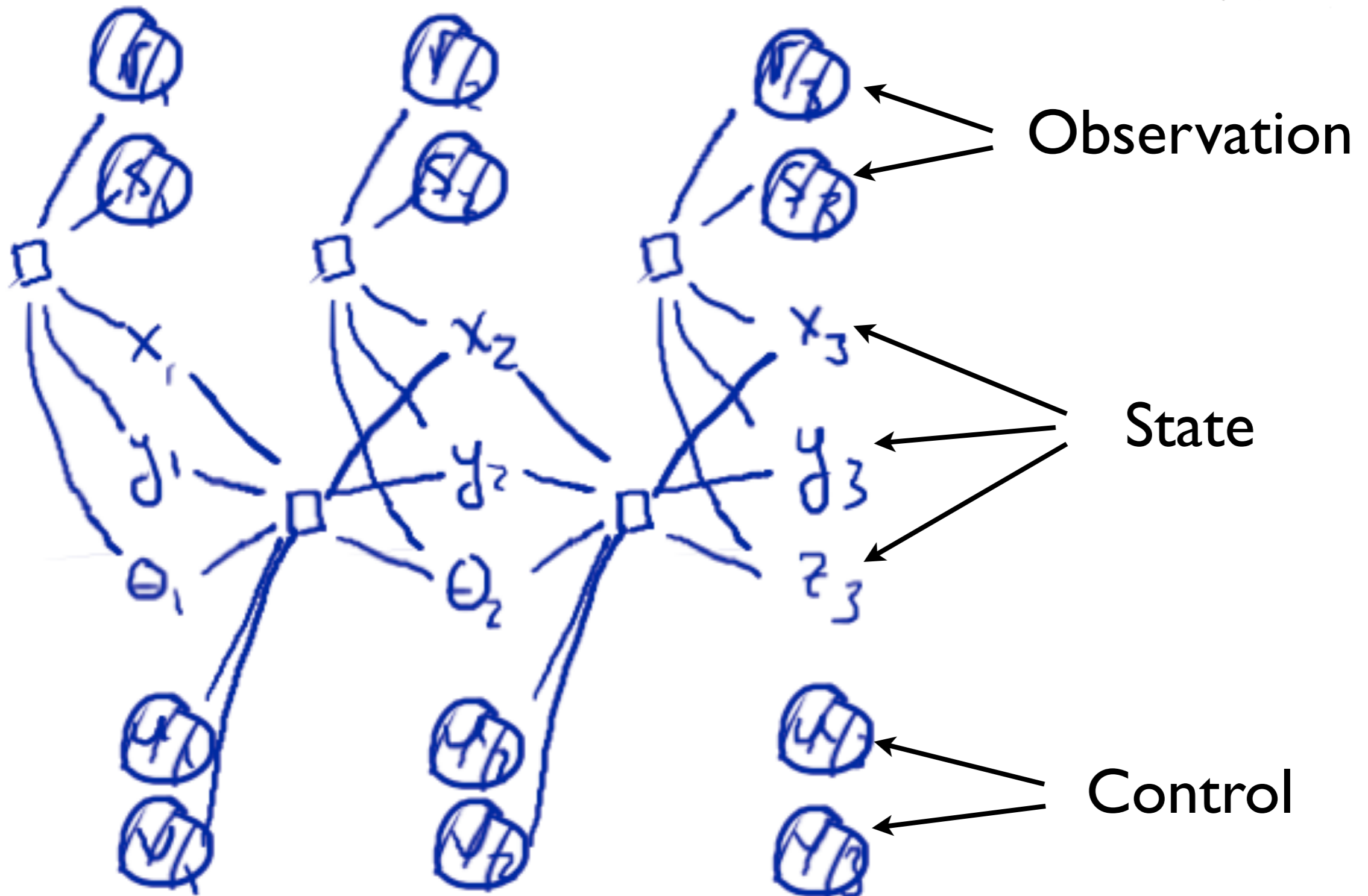
$$\forall x_t, y_t, \theta_t, u_t, v_t \quad \sum_{x_{t+1}, y_{t+1}, \theta_{t+1}} \phi(x_t, y_t, \theta_t, u_t, v_t, x_{t+1}, y_{t+1}, \theta_{t+1}) = 1$$

$$\forall x_t, y_t, \theta_t \quad \sum_{r_t, s_t} \phi(x_t, y_t, \theta_t, r_t, s_t) = 1$$

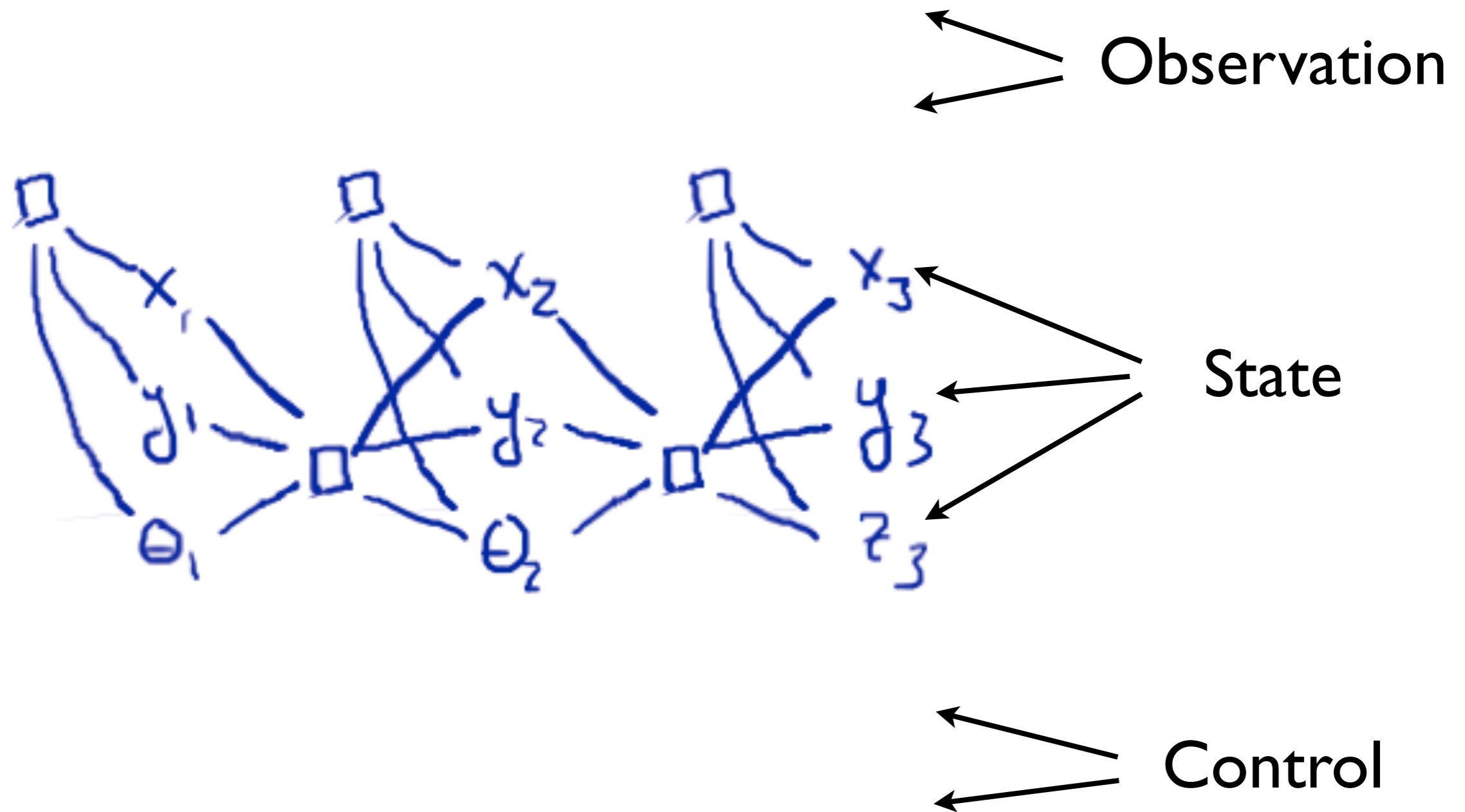
Three kinds of variable



Condition on obs, do(control)



Condition on obs, do(control)



Simplified version

- State: $x_t \in \{1, 2, 3\}$
- Observation: $y_t \in \{L, H\}$
- Control: just one (i.e., no choice)—“keep going”

Hidden Markov Models



- This is an HMM—a DBN with:
 - ▶ one state variable
 - ▶ one observation variable

Potentials

		X_{t+1}		
		1	2	3
X_t	1	0.7	0.3	0
	2	0.3	0.3	0.3
	3	0	0.3	0.7

		Y_t	
		L	H
X_t	1	0.67	0.33
	2	0.5	0.5
	3	0.33	0.67

HMM inference

- Condition on $y_1 = H, y_2 = H, y_3 = L$
- What is $P(X_2 \mid HHL)$?

HMM factors after conditioning

x_1	ϕ_1	x_2	ϕ_2	x_3	ϕ_3
1	.33	1	.33	1	.67
2	.5	2	.5	2	.5
3	.67	3	.67	3	.33

ϕ_4	x_2		
	.67	.33	0
x_1	.33	.33	.33
	0	.33	.67

ϕ_5	x_3		
	.67	.33	0
x_2	.33	.33	.33
	0	.33	.67

Eliminate x_1 and x_3

ϕ_1, ϕ_4

$$\begin{array}{c} x_2 \\ \hline \begin{array}{ccc} 2/9 & 1/9 & 0 \\ x_1 & 1/6 & 1/6 & 1/6 \\ 0 & 2/9 & 4/9 \end{array} \end{array} \rightarrow$$

α_{12}

$$\frac{x_2}{7/18}$$

$$1/2$$

$$11/18$$

ϕ_3, ϕ_5

$$\begin{array}{c} x_3 \\ \hline \begin{array}{ccc} 4/9 & 1/6 & 0 \\ x_2 & 2/9 & 1/6 & 1/9 \\ 0 & 1/6 & 2/9 \end{array} \end{array} \rightarrow$$

β_{23}

$$\frac{x_2}{14/18}$$

$$14/18$$

$$1/2$$

$$-7/18$$

Multiply remaining potentials and renormalize

$$\alpha_{12}$$

$$\frac{x_2}{7/18}$$

$$1/2$$

$$11/18$$

$$\frac{\alpha_{12} \beta_{23} \phi_2}{.079}$$

$$.125$$


$$.158$$

$$\beta_{23}$$

$$\frac{x_2}{14/18}$$

$$1/2$$

$$7/18$$




$$1/2$$

$$.22$$

$$.34$$

$$.44$$

Forward-backward



- You may recognize the above as the forward-backward algorithm
- Special case of dynamic programming / variable elimination / belief propagation



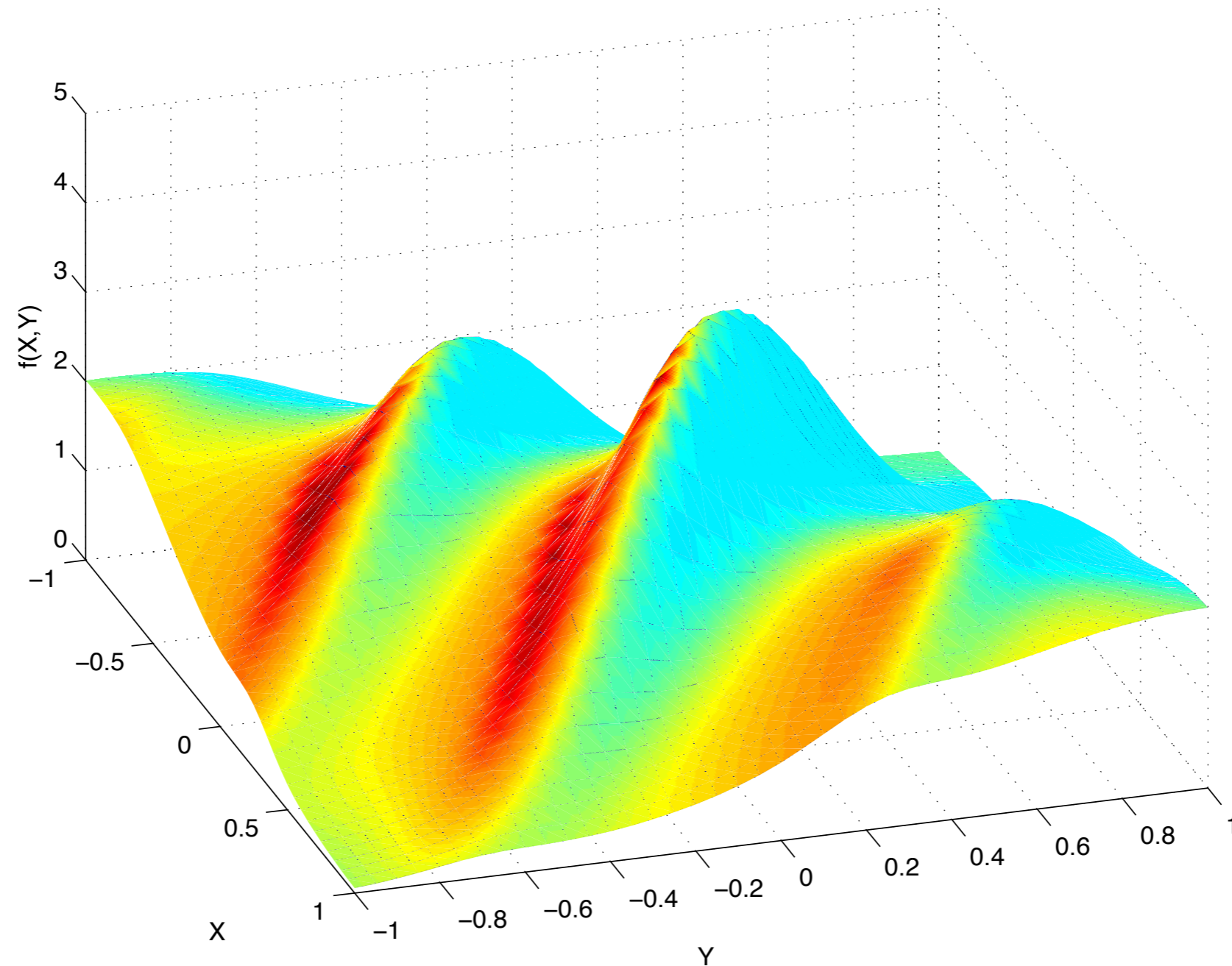
Approximate Inference

Most of the time...



- Treewidth is big
- Variables are high-arity or continuous
- Can't afford exact inference
- Need numerical integration (and/or summation)
- We'll look at randomized algorithms

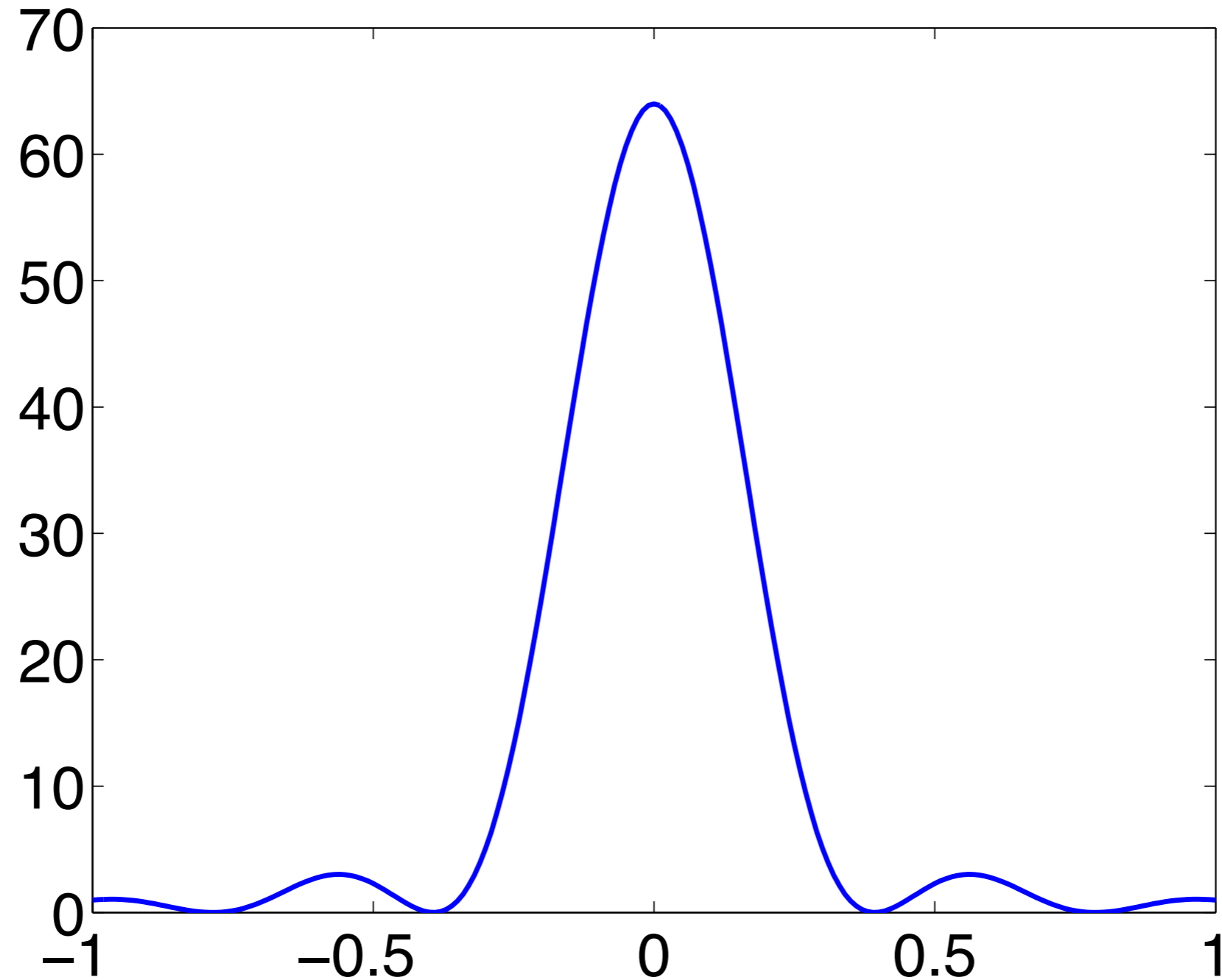
Numerical integration



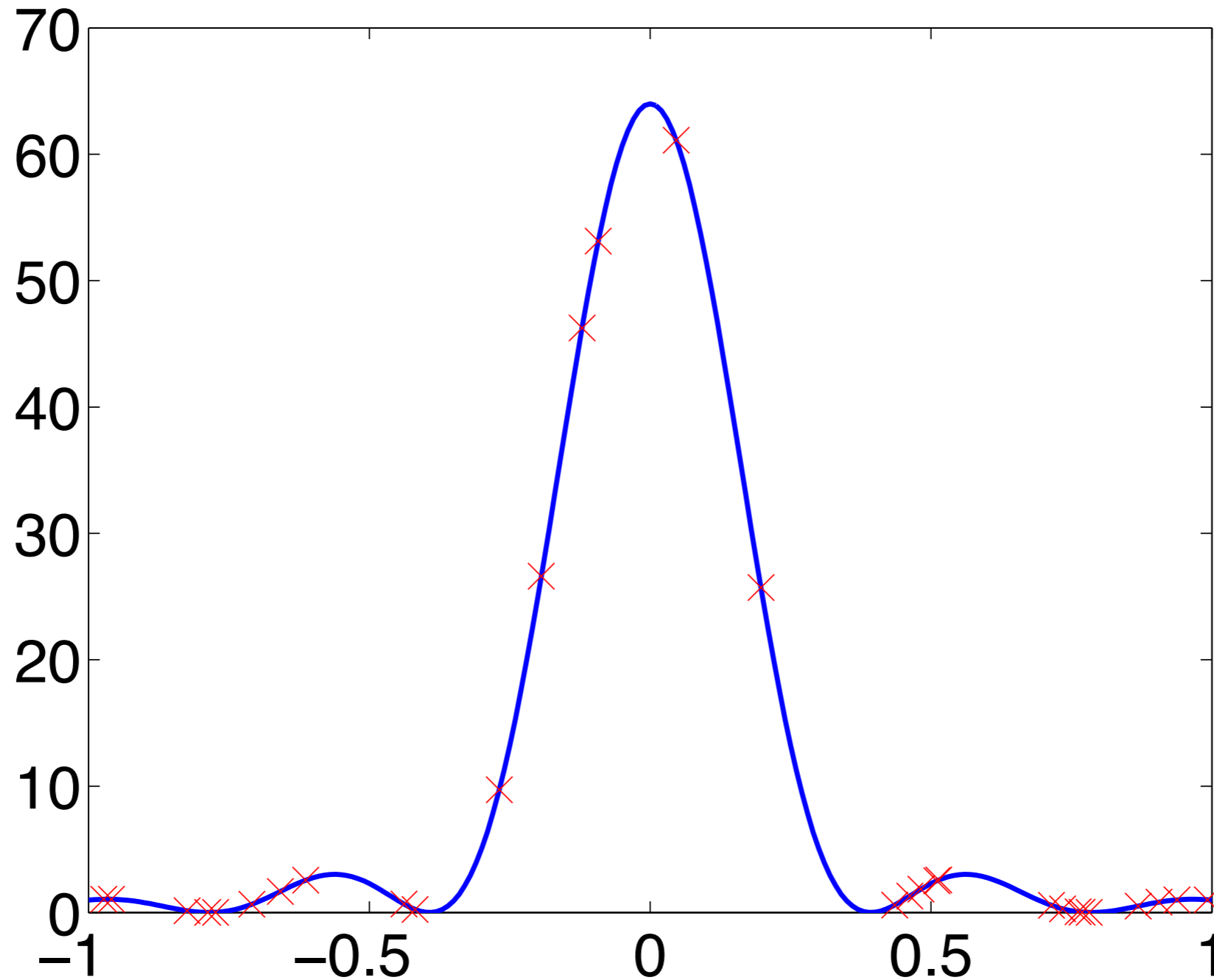
Integration in 1000s of dims



Simple ID problem



Uniform sampling



$$\frac{2}{N} \sum_i f(x_i)$$

Uniform sampling

$$\begin{aligned} E(f(X)) &= \int P(x) f(x) dx \\ &= \frac{1}{V} \int f(x) dx \end{aligned}$$

- So, $V E(f(X))$ is desired integral
- But standard deviation can be big
- Can reduce it by averaging many samples
- But only at rate $1/\sqrt{N}$

Importance sampling

- Instead of $x \sim \text{uniform}$, use $x \sim Q(x)$
- Q = importance distribution
- Should have $Q(x)$ large where $f(x)$ is large
- Problem:

$$E_Q(f(X)) = \int Q(x) f(x) dx$$

Importance sampling

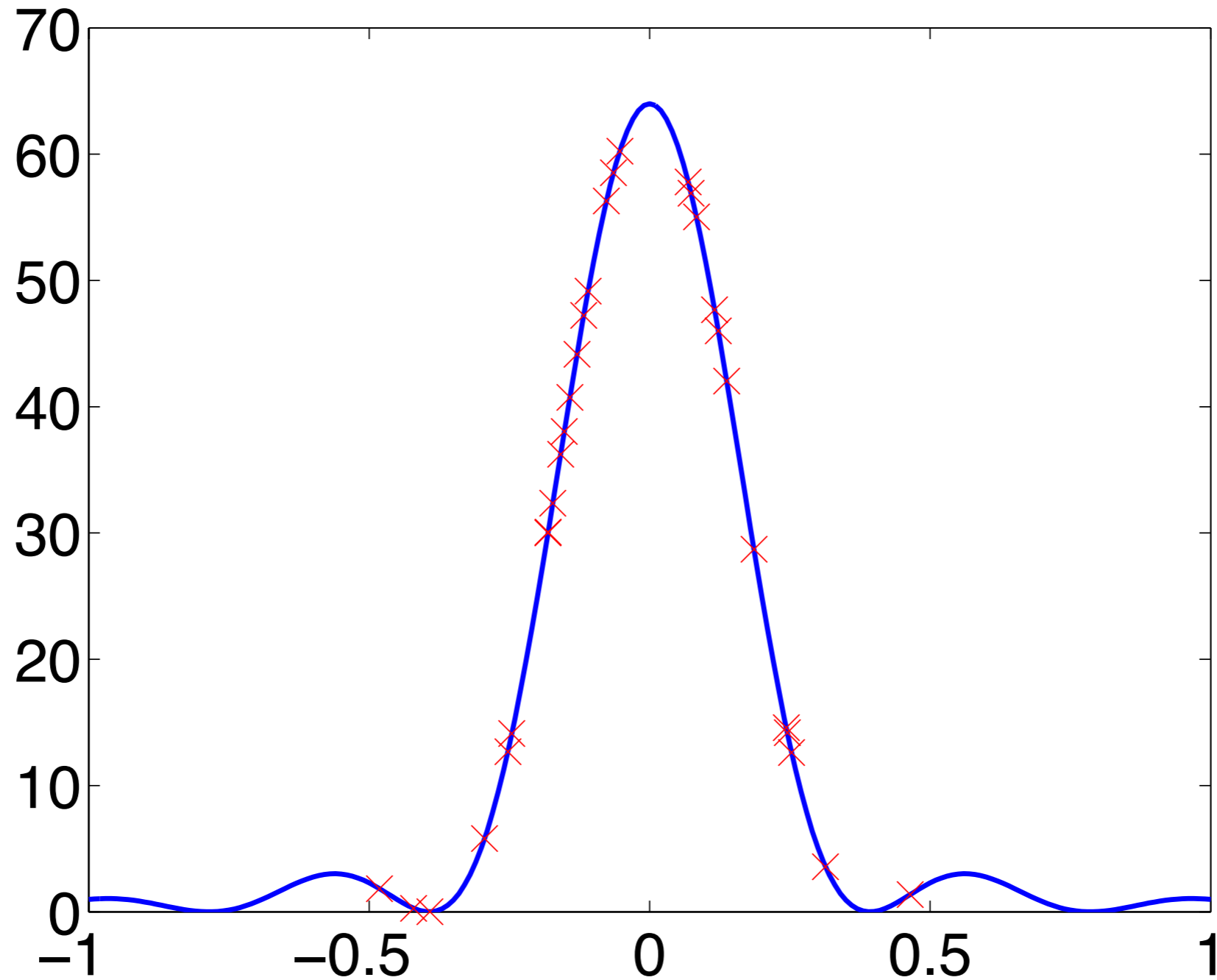
$$h(x) \equiv f(x)/Q(x)$$

$$\begin{aligned} E_Q(h(X)) &= \int Q(x)h(x)dx \\ &= \int Q(x)f(x)/Q(x)dx \\ &= \int f(x)dx \end{aligned}$$

Importance sampling

- So, take samples of $h(X)$ instead of $f(X)$
- $w_i = 1/Q(x_i)$ is **importance weight**
- $Q = 1/V$ yields uniform sampling

Importance sampling



Variance

- How does this help us control variance?
- Suppose f big $\implies Q$ big
- And Q small $\implies f$ small
- Then $h = f/Q$ never gets too big
- Variance of each sample is lower \implies need fewer samples
- A good Q makes a good IS