

compact representations of distributions  
Directed Acyclic Graphical models  
Bayes-ball algorithm



A6: - typo in Fenchel dual:

$$D(y) = -f^*(\underline{-y^*}) - g^*(A^T y)$$

- extra hints now included

Midterm (a modest proposal) change?

- Cancel class Monday? (you can go to Bertsekas talk)

- midterm November 12, 17, 19?

- A6 due Monday?

- tutorials moved?

3: code rec

2: <sup>major</sup> inconvenience

1: minor inconvenience



## Compact representations of joint distributions

Notation for today:

- we'll use  $x_j$  for  $(x_i)_j$

Consider  $x \in \{0,1\}^d$  (binary vectors)

We want to model  $p(x)$  (joint distribution)

- why? scientific discovery, outlier detection,  
 $p(x_i | y_i)$  in generative model,  $p(x_A | x_B)$

- Previous: Naive Bayes, Gaussian, mixture models

- Today: "graphical" models.



(acyclic)

Basic idea behind directed graphical models:

- use product rule repeatedly:

$$p(x) = p(x_1) p(x_2 | x_1) p(x_3 | x_2, x_1) p(x_4 | x_3, x_2, x_1) \dots p(x_d | x_{1:(d-1)})$$
$$= \prod_{j=1}^d p(x_j | x_{1:(j-1)})$$

$2^{d-1}$

- too many parameters!

Solutions:

- "parsimonious" parameterization

$$p(x_j | x_{1:(j-1)}) = f(W^T x_{1:(j-1)})$$

$d$  parameters

- Conditional independence

eg. naive Bayes:  $x_j \perp x_{1:(j-1)} | y_i \Rightarrow p(x_j | x_{1:(j-1)} | y_i) = p(x_j | y_i)$

Markov chain:  $x_j \perp x_{1:(j-2)} | x_{j-1} \Rightarrow p(x_j | x_{1:(j-1)}) = p(x_j | x_{j-1})$

General:  $x_j \perp x_{1:(j-1) \setminus \pi(j)} | \pi(j)$   
"parents"

- can also do both.



# Directed acyclic graphical (DAG) models

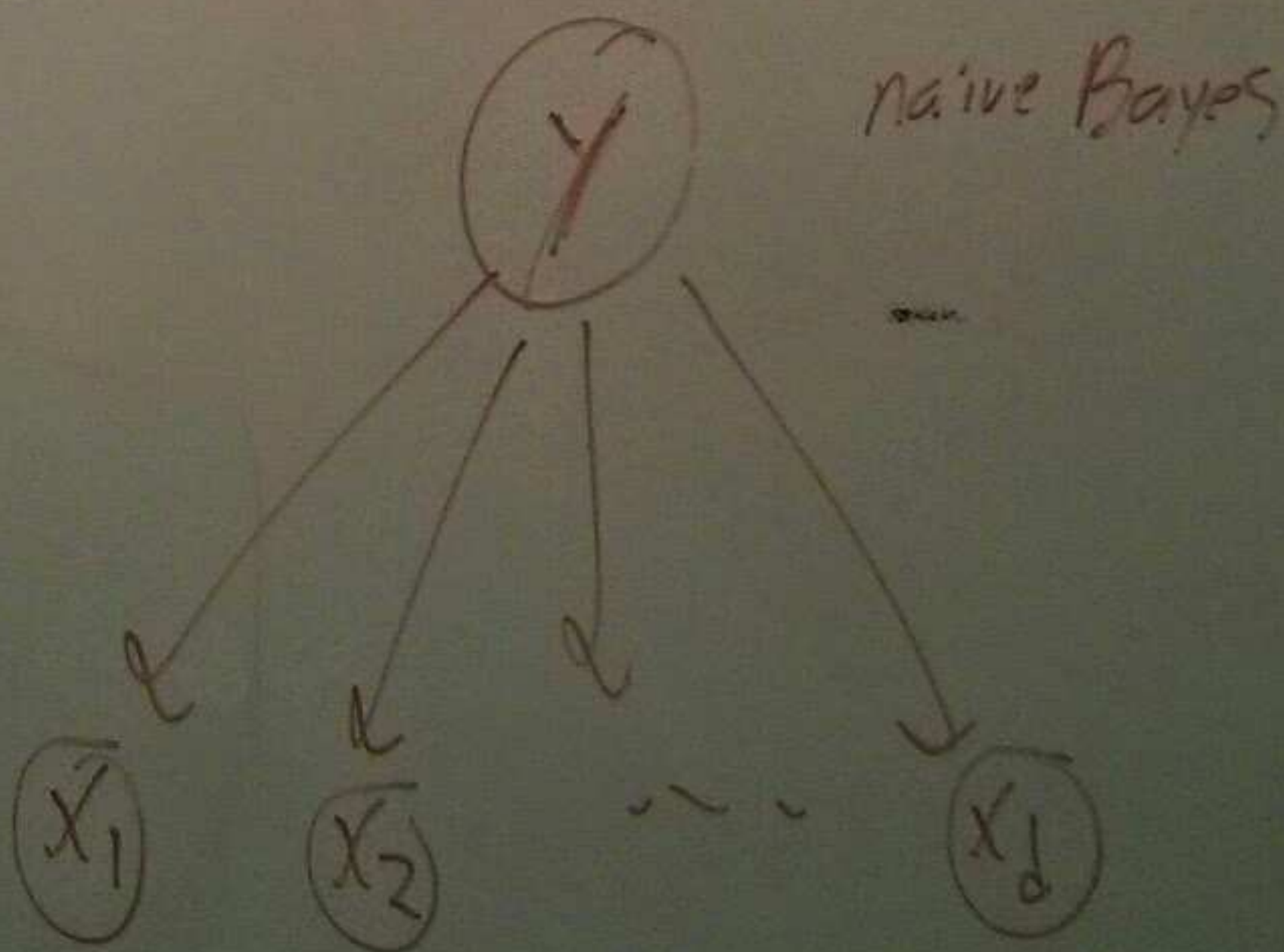
$$\text{Graph } G = (V, E)$$

"vertices" "edges"

$V$ : random variables  $x_j$

$E$ :  $x_{\pi(j)} \rightarrow x_j$

$$p(x) = \prod_{j=1}^d p(x_j | x_{\pi(j)})$$



$(d-1)$   
 $d-1$

$x_{1:(j-1)}$   
parameters

$$p(x_j | y_i) \\ = p(x_j | x_{j-1})$$

Markov

$x_1$

2nd

$x_2$



models

"Bayesian networks"

"belief networks"

"causal networks"

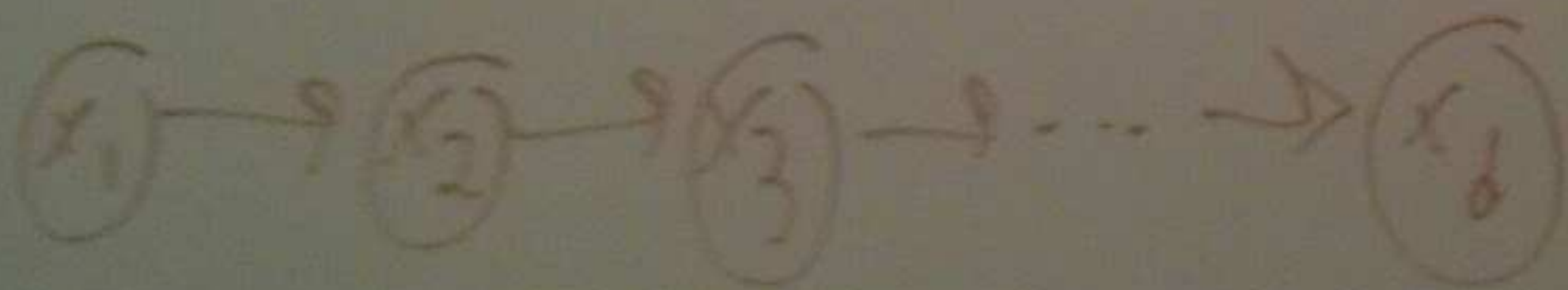


(you need a justification  
to interpret edges  
causally)

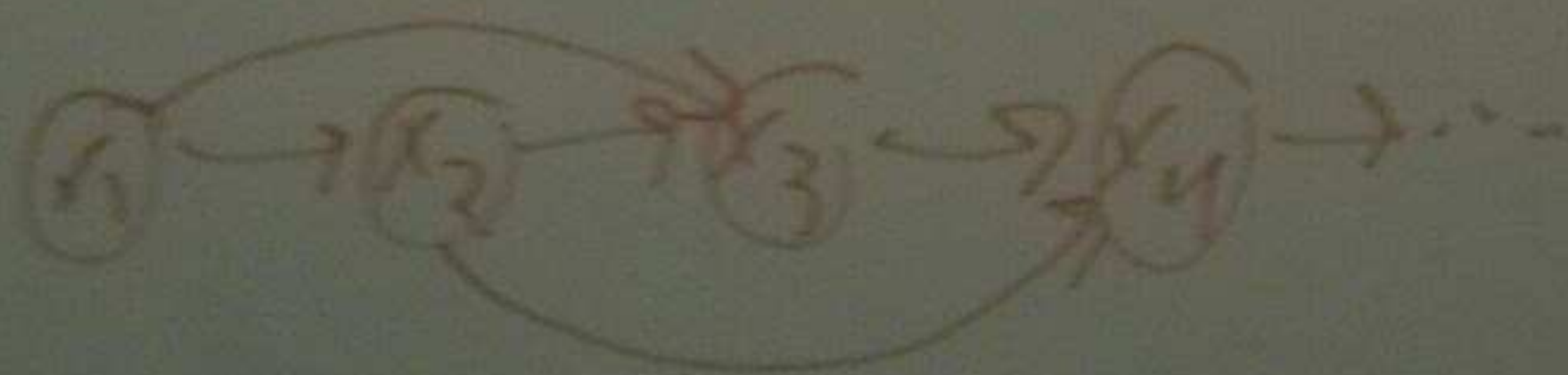
Hidden Markov Models

$x_1$

Markov chain



2nd order Markov chain

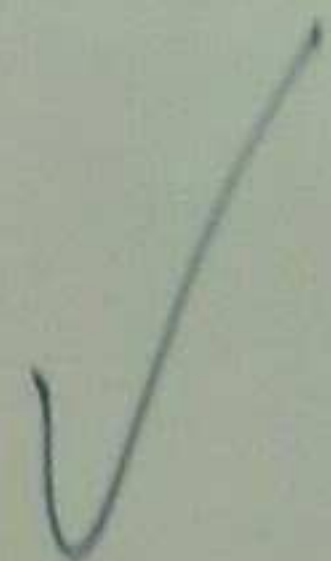


$$p(x_j | x_{j-1}, x_{j-2})$$



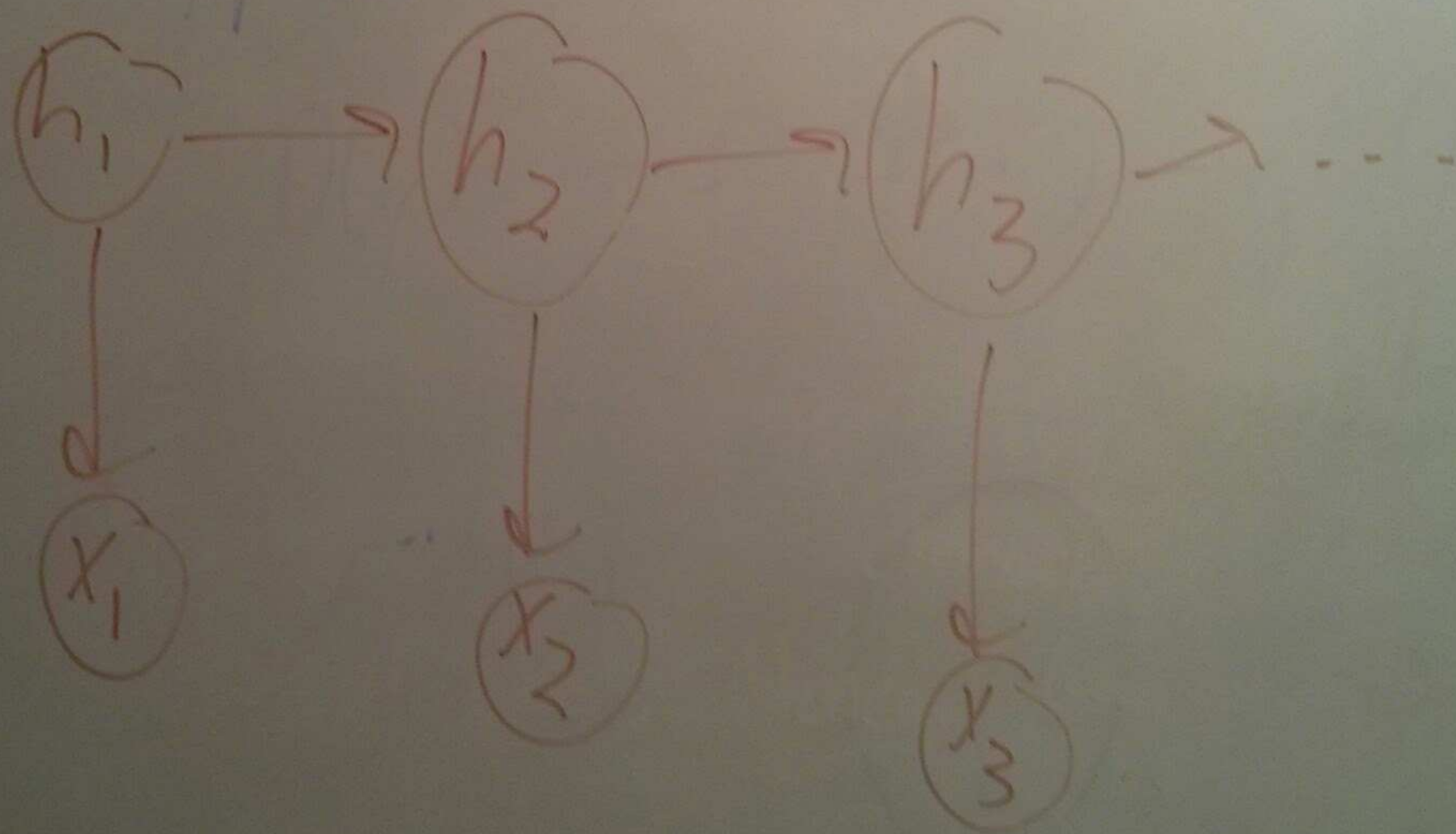
Directed Acyclic Graphical models

Bayes-ball algorithm

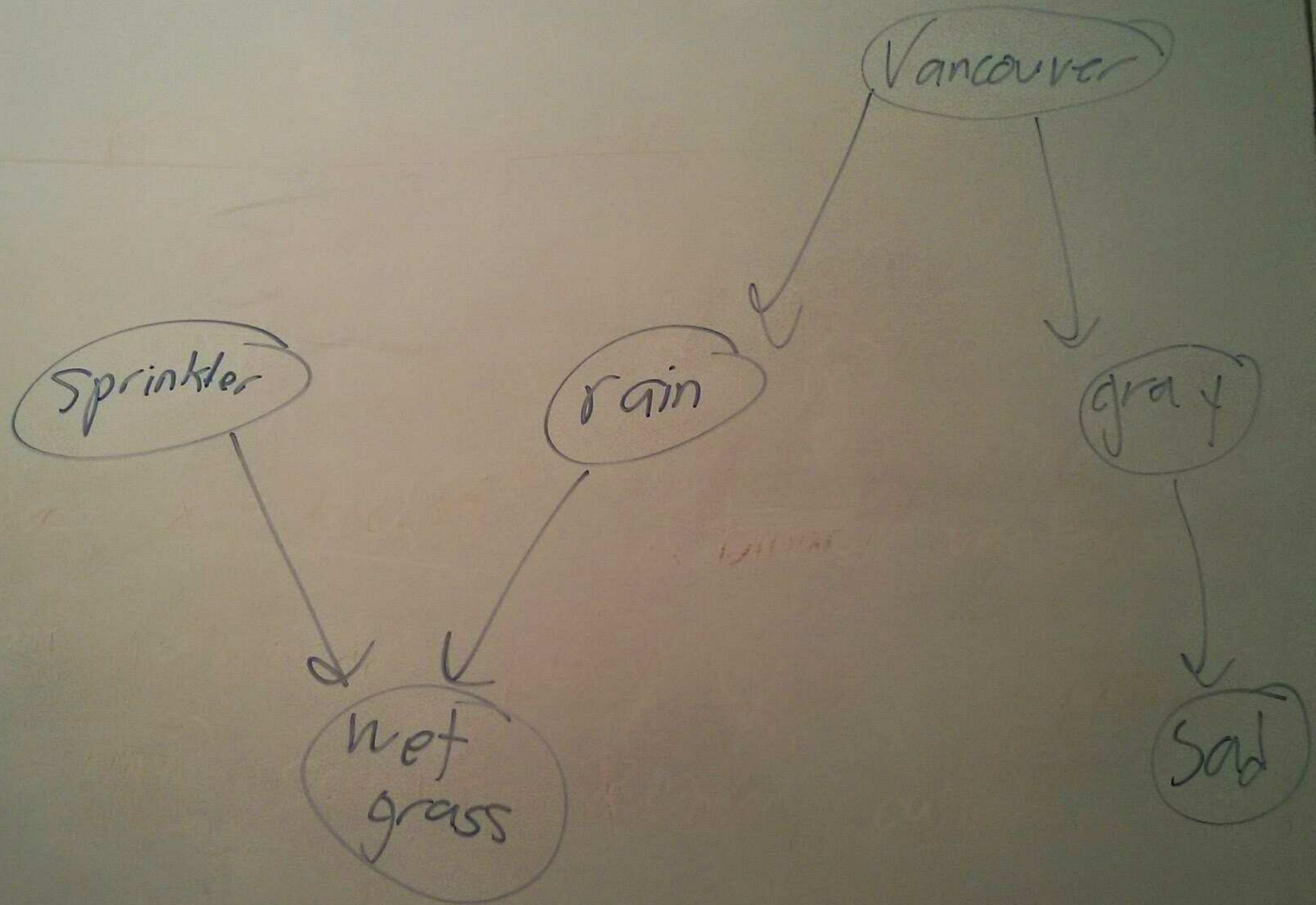


HMM

(Kalman filtering)









## Learning:

- fit each  $p(x_j | x_{\pi(j)})$  independently

- "inference":

computing $p(x)$	}	easy
computing $p(x_j   x_{\pi(j)})$		
computing $p(x_j   x_{j+1})$	}	#P-hard











Earthquake

Burglary

Alarm

Stuff is missing

Burglary  $\perp$  Call

Neighbour calls police

Call  $\perp$

Stuff is missing

Earthquake  $\perp$  Call

Earthquake  $\perp$  Call | Alarm

Alarm  $\perp$  Stuff is missing

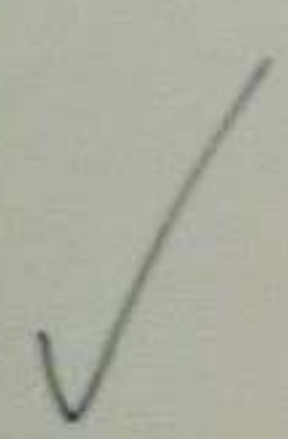
Alarm  $\perp$  Stuff is missing | Burglary

Earthquake  $\perp$  Burglary

Earthquake  $\perp$  Burglary | Alarm

Earthquake  $\perp$  Burglary | Call

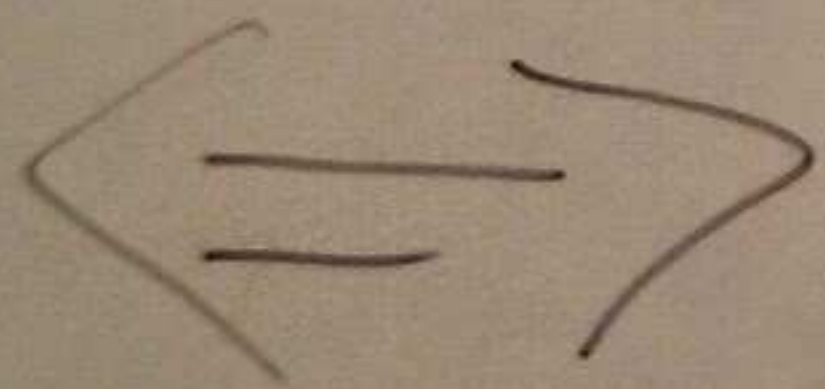
"explaining away"





## ⊛ Gaussian DAGs

$$p(x_j | x_{\pi(j)}) \sim N(\mu_j + \bar{w}_j^T x_{\pi(j)}, \sigma_j^2)$$



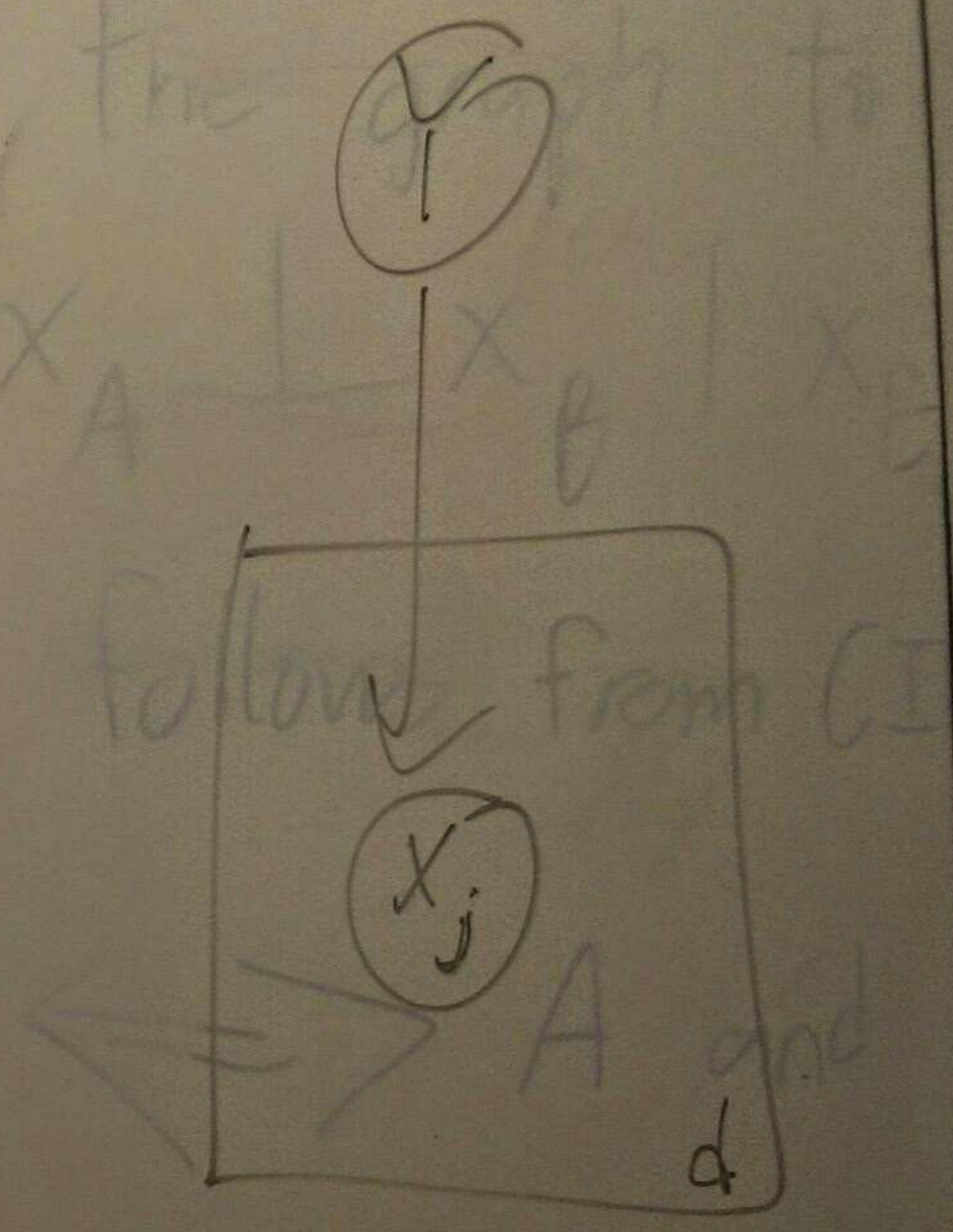
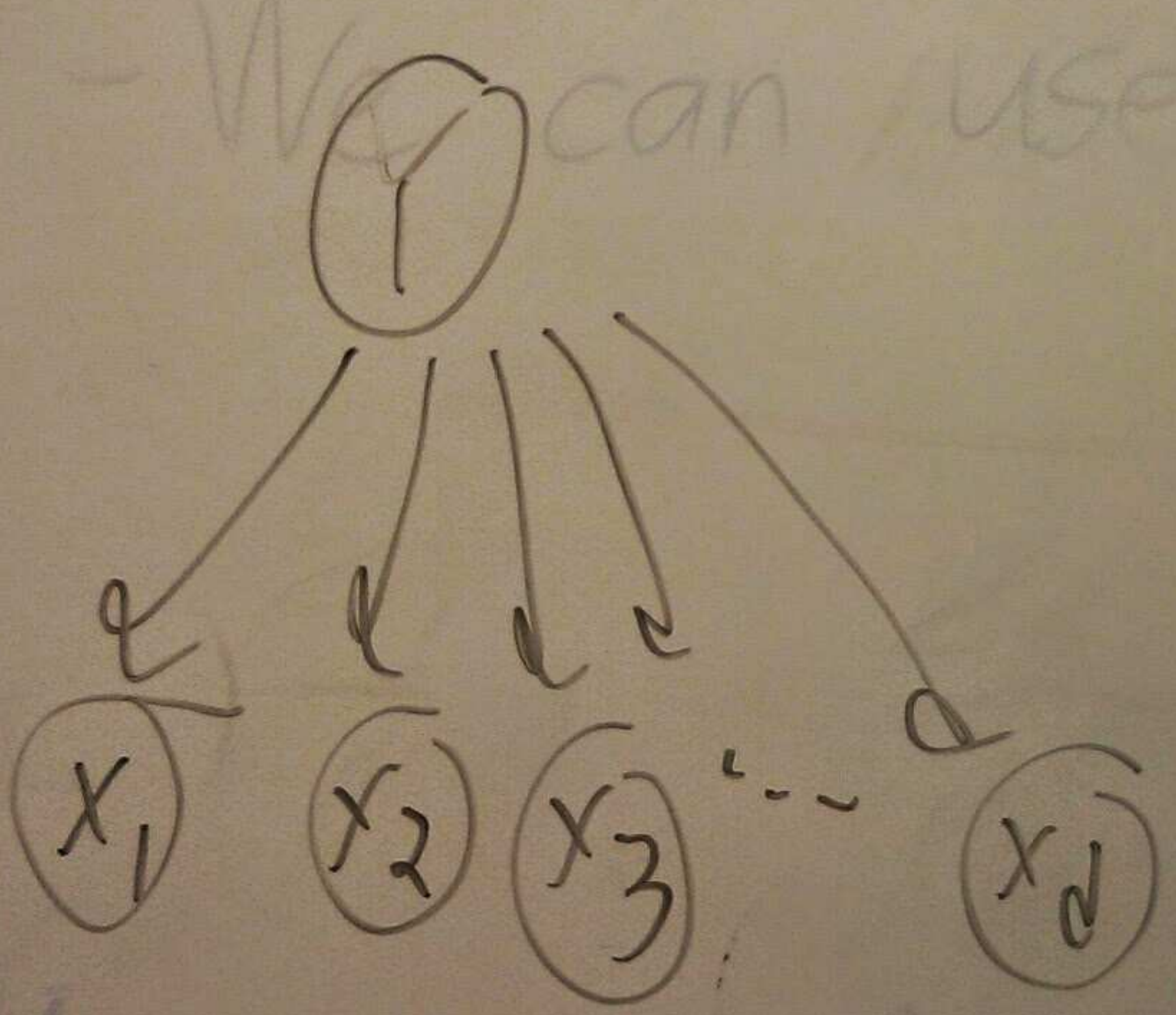
$$x \sim N(\mu, \Sigma)$$

$$\Sigma^{-1} = LL^T$$

∴ non-zero pattern of  $L$   
comes from  $G$ .

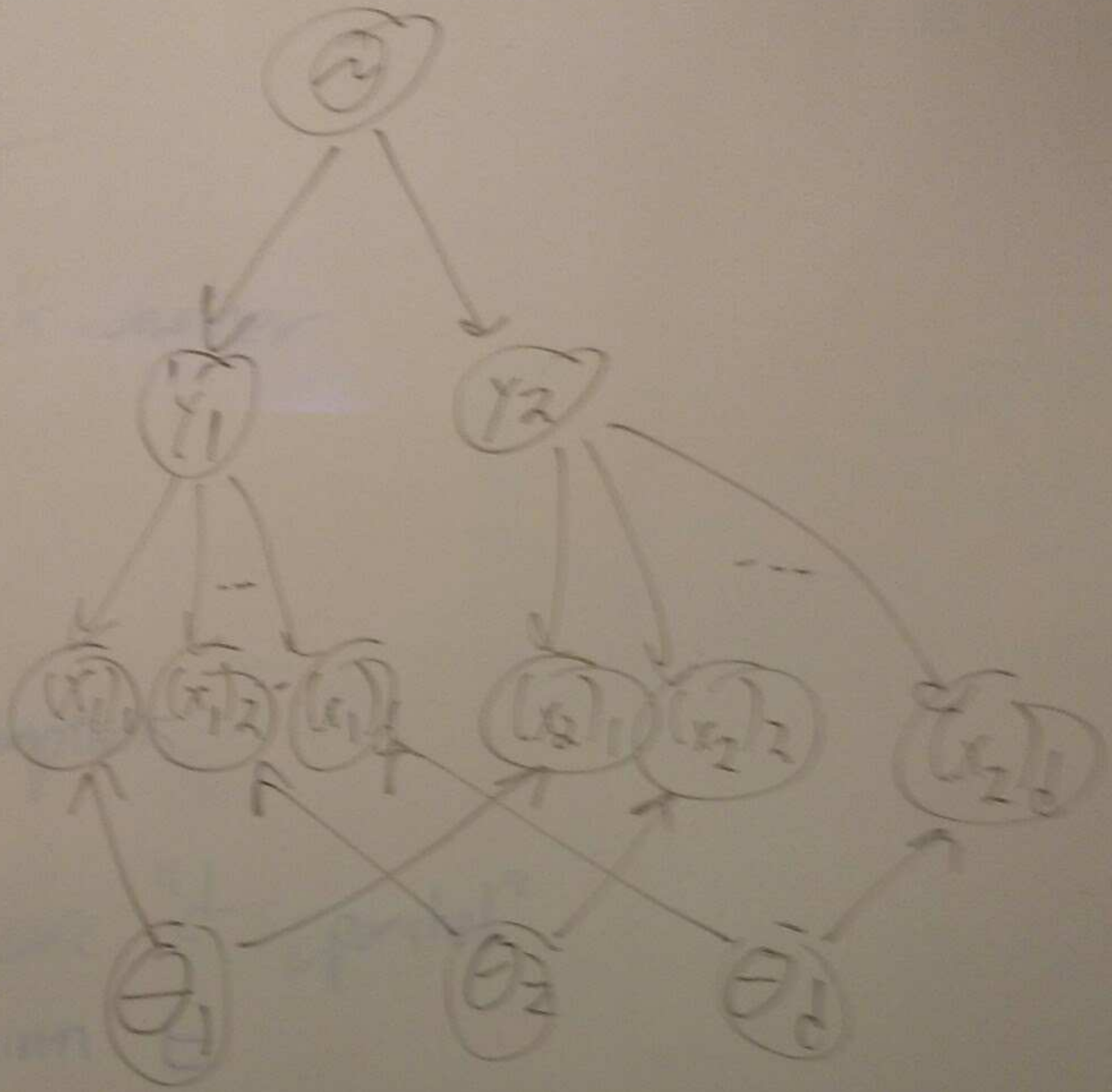


# Plate Notation





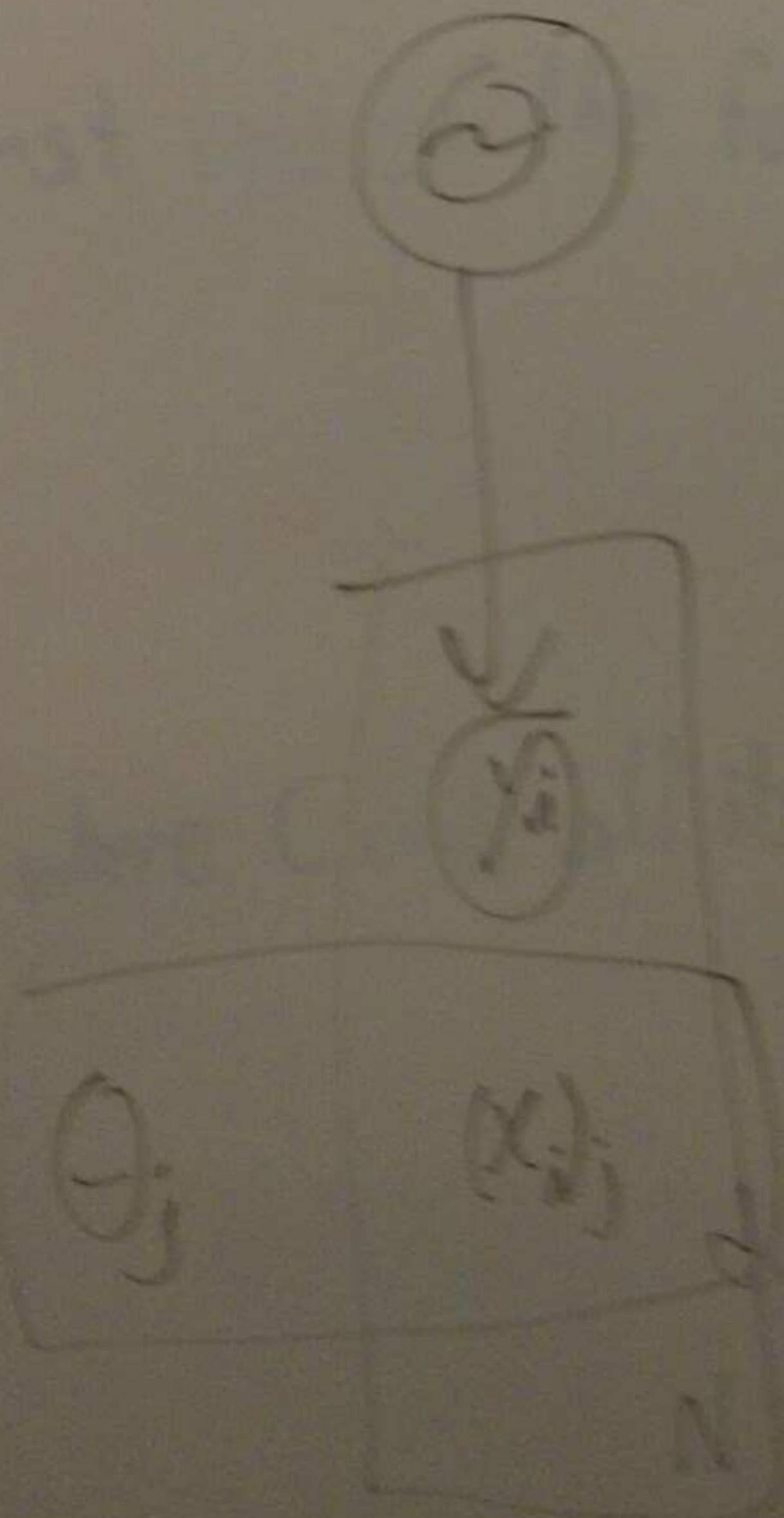
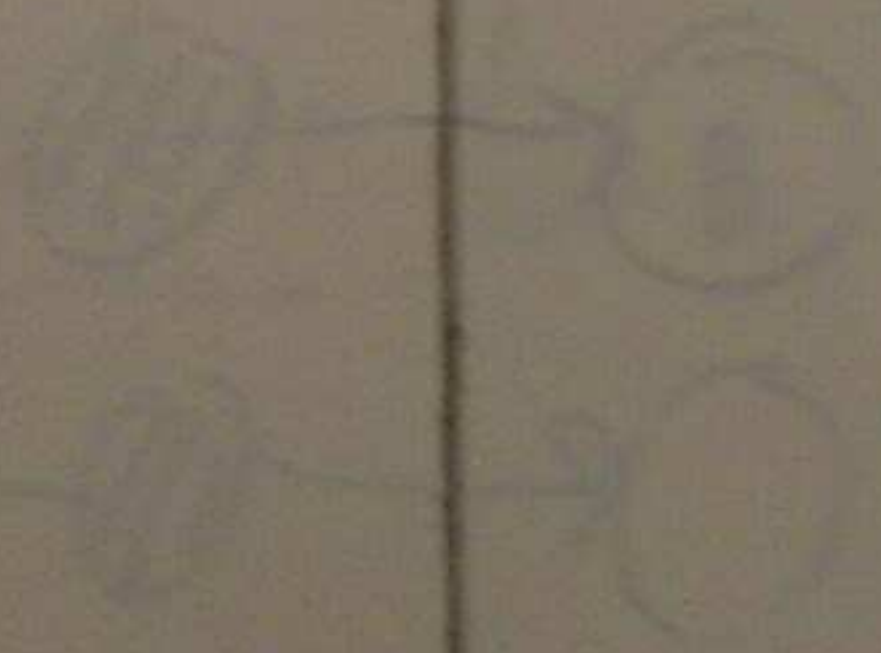
Proof



d

given

least following



is decided  
are values