

EM & Graphical Models

Xuan Zhang

Nov. 18. 2022

I. EM

- EM
- GMM
- EM for a GMM
- K-means as a special case of EM

II. Graphical Models

- Terminology
- D-separation

I.

E

M

The EM Algorithm

X : observed data

Z : hidden variable

θ : model parameters



e.g. $Z \sim \text{Ber}(\theta)$

$$Z = \theta^T X$$

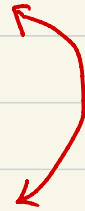
| | k-means | GMM |
|----------|---------------------|--------------------|
| Z | cluster assignments | responsibility |
| θ | μ | μ, Σ, π |

E step:

fix θ , return Z or $q(Z)$

M step:

fix Z , find θ that maximizes likelihood



The EM Algorithm

Goal:

$$\max \sum_{n=1}^N \log p(x_n | \theta) \quad \text{i.e. maximize the log likelihood of the observed data}$$

•
find a θ

Alternatively:

$$\max \sum_{n=1}^N \log \sum_{z_n} p(x_n, z_n | \theta) \quad \text{i.e. maximize the complete data log likelihood}$$

Decomposition of Complete data log Likelihood

Let's introduce a set of arbitrary distributions $q_n(z_n)$ over each hidden variable z_n .

$$\begin{aligned} & \log \sum_{z_n} P(x_n, z_n | \theta) \\ = & \boxed{\sum_{z_n} q_n(z_n) \log \frac{P(x_n, z_n | \theta)}{q_n(z_n)}} - \boxed{\sum_{z_n} q_n(z_n) \log \frac{P(z_n | x_n, \theta)}{q_n(z_n)}} \\ & \mathcal{L}(q, \theta) \qquad \qquad \qquad -KL(q \parallel P) \\ = & \sum_{z_n} q_n(z_n) \left[\log \frac{P(x_n, z_n | \theta)}{q_n(z_n)} - \log \frac{P(z_n | x_n, \theta)}{q_n(z_n)} \right] \\ = & \sum_{z_n} q_n(z_n) \log P(x_n | \theta) = \log P(x_n | \theta) \end{aligned}$$

* KL Divergence: $D_{kl}(q \parallel p) = \sum q \log \frac{q}{p}$

Decomposition of Complete data log likelihood

$$\log P(X_n, \theta) = \log \sum_{z_n} P(X_n, z_n | \theta)$$

$$= \underbrace{\sum_{z_n} q_n(z_n) \log \frac{P(X_n, z_n | \theta)}{q_n(z_n)}}_{\mathcal{L}(q, \theta)} + \underbrace{-\sum_{z_n} q_n(z_n) \log \frac{P(z_n | X_n, \theta)}{q_n(z_n)}}_{KL(q \| P)}$$

$$\mathcal{L}(q, \theta) = \sum_{z_n} q_n(z_n) \log \frac{P(X_n, z_n | \theta)}{q_n(z_n)} \leq \log \sum_{z_n} q_n(z_n) \frac{P(X_n, z_n | \theta)}{q_n(z_n)} = \log P(X_n | \theta)$$

↓
equal when $KL(q \| P) = 0$, i.e. $q_n(z_n) = P(z_n | X_n, \theta)$

* Jensen's Inequality: For a concave function f , $\sum_i \lambda_i f(x_i) \leq f(\sum_i \lambda_i x_i)$

Understand EM with Decomposition

$$\begin{aligned} \log p(X_n, \theta) &= \log \sum_{z_n} p(X_n, z_n | \theta) \\ &= \underbrace{\sum_{z_n} q_n(z_n) \log \frac{p(X_n, z_n | \theta)}{q_n(z_n)}}_{\mathcal{L}(q, \theta)} + \underbrace{-\sum_{z_n} q_n(z_n) \log \frac{p(z_n | X_n, \theta)}{q_n(z_n)}}_{KL(q \| p)} \end{aligned}$$

E step:

fix θ , then $q_n(z_n) = p(z_n | X_n, \theta)$, $KL = 0$, $\mathcal{L} \uparrow$
(Recall: fix θ , return Z or $q(Z)$)

M step:

fix $q_n(z_n)$, find $\arg\max_{\theta} \mathcal{L}(q, \theta)$, $KL \uparrow$, $\mathcal{L} \uparrow$
(Recall: fix Z , find θ that maximizes likelihood)

Another way of Decomposition

$$\begin{aligned}\sum_n \log p(X_n, \theta) &= \sum_n \log \sum_{z_n} p(X_n, z_n | \theta) \\ &= \sum_n \log \left[\sum_{z_n} q_n(z_n) \frac{p(X_n, z_n | \theta)}{q_n(z_n)} \right] \\ &\geq \sum_n \sum_{z_n} q_n(z_n) \log \frac{p(X_n, z_n | \theta)}{q_n(z_n)} = \sum_n \sum_{z_n} q_n(z_n) \log p(X_n, z_n | \theta) + \sum_n \sum_{z_n} q_n(z_n) \log \frac{1}{q_n(z_n)}\end{aligned}$$

$$\boxed{\sum_n \log \sum_{z_n} p(X_n, z_n | \theta)} \geq \boxed{\sum_n \mathbb{E}_{q_n} \log p(X_n, z_n | \theta)} + \sum_n H(q_n)$$

complete data log likelihood expected complete data log likelihood

ELBO
(evidence lower bound)

Another way of Decomposition

$$\sum_n \log \sum_{z_n} p(x_n, z_n | \theta)$$

complete data log likelihood

$$\geq \underbrace{\sum_n \mathbb{E}_{q_n} \log p(x_n, z_n | \theta)}_{\text{expected complete data log likelihood}} + \sum_n H(q_n)$$

expected complete data log likelihood

$$= \sum_n \log p(x_n | \theta) - \text{KL}(q_n(z_n) || p(z_n | x_n, \theta))$$

* see Book 1 8.7.2.2
for derivation

E step:

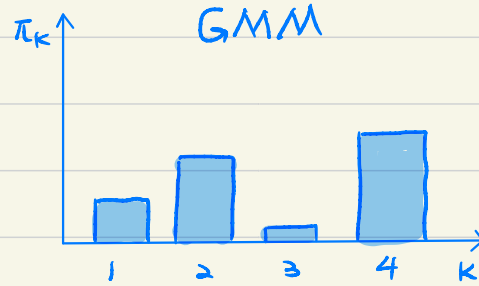
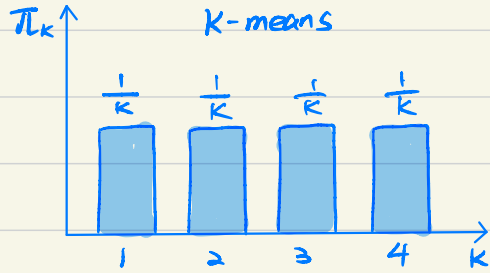
$$\text{fix } \theta, \quad \underbrace{q_n(z_n) = p(z_n | x_n, \theta)}_{\text{the posterior}}, \quad \text{KL} = 0$$

M step:

$$\text{fix } q_n(z_n), \quad H(q_n) \text{ stays constant, } \operatorname{argmax}_{\theta} \sum_n \mathbb{E}_{q_n} \log p(x_n, z_n | \theta)$$

Gaussian Mixture Models (GMM)

π_k : cluster weights, $0 \leq \pi_k \leq 1$, $\sum_{k=1}^K \pi_k = 1$



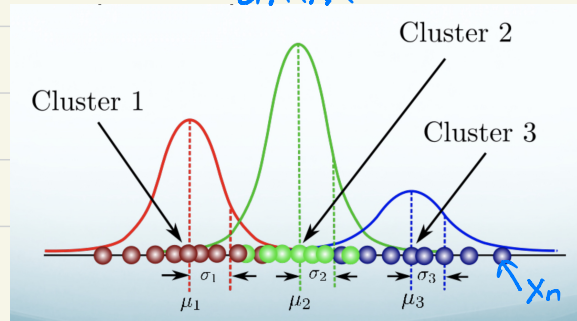
μ_k, Σ_k : probability distribution of points in cluster k

k-means

$$\Sigma_k = I$$

$$\mu_k = \frac{\sum_{n:z_n=k} x_n}{N_k}$$

GMM



GMM

π_k : cluster weights, $0 \leq \pi_k \leq 1$, $\sum_{k=1}^K \pi_k = 1$

μ_k, Σ_k : $N(x_n | \mu_k, \Sigma_k)$

z_n : a selector, $[0, 0, \dots, 1, \dots]$, $P(z_{nk} = 1) = \pi_k$

$$P(x_n | z_{nk} = 1) = N(x_n | \mu_k, \Sigma_k) \quad P(X|Z) = \prod_{n=1}^N \prod_{k=1}^K N(x_n | \mu_k, \Sigma_k)^{z_{nk}}$$

γ_{nk} : responsibility of cluster k for generating example n
determined by $x_n, \pi_k, \mu_k, \Sigma_k$

GMM

$$\pi_k = P(z_{nk} = 1)$$

prior

$$P(x_n | z_{nk} = 1) = \mathcal{N}(x_n | \mu_k, \Sigma_k)$$

likelihood

$$r_{nk} = P(z_{nk} = 1 | x_n)$$

posterior

$$= \frac{P(z_{nk} = 1) P(x_n | z_{nk} = 1)}{P(x_n)} = \frac{P(z_{nk} = 1) P(x_n | z_{nk} = 1)}{\sum_{j=1}^k P(z_{nj} = 1) P(x_n | z_{nj} = 1)}$$

$$= \frac{\pi_k \mathcal{N}(x_n | \mu_k, \Sigma_k)}{\sum_{j=1}^k \pi_j \mathcal{N}(x_n | \mu_j, \Sigma_j)}$$

EM for a GMM

Goal:

$$\begin{aligned} & \max \sum_n \log \sum_k p(X_n, Z_n | \theta) \\ & = \max \sum_n \log \left\{ \sum_k \pi_k \mathcal{N}(X_n | \mu_k, \Sigma_k) \right\} \end{aligned}$$

E step:

$$\text{fix } \pi, \mu, \Sigma, \quad \text{get } r_{nk} = \frac{\pi_k \mathcal{N}(X_n | \mu_k, \Sigma_k)}{\sum_{j=1}^K \pi_j \mathcal{N}(X_n | \mu_j, \Sigma_j)}$$

M step:

fix r_{nk} , maximization (take derivative, set it to 0)
get updated π, μ, Σ

Another way of Decomposition

$$\sum_n \log \sum_{z_n} p(x_n, z_n | \theta) \geq \sum_n \mathbb{E}_{q_n} \log p(x_n, z_n | \theta) + \sum_n H(q_n)$$

complete data log likelihood expected complete data log likelihood

E step:

$$\text{fix } \theta, \quad q_n(z_n) = p(z_n | x_n, \theta)$$

the posterior, rnk!

M step:

$$\text{fix } q_n(z_n), \quad H(q_n) \text{ stays constant, } \operatorname{argmax}_{\theta} \sum_n \mathbb{E}_{q_n} \log p(x_n, z_n | \theta)$$

EM for a GMM

E step:

$$E q_n(Z_{nk}) = \sum_n \sum_k q_n(Z_{nk}) Z_{nk}$$

E step: set $q_n(Z_{nk}) = P(Z_{nk} | X_n, \theta)$
 $= r_{nk}$

$$E q_n(Z_{nk}) = r_{nk}$$

fix π, μ, Σ , get $r_{nk} = \frac{\pi_k N(X_n | \mu_k, \Sigma_k)}{\sum_{j=1}^K \pi_j N(X_n | \mu_j, \Sigma_j)}$

write out *expected complete data log likelihood*:

$$\begin{aligned} E_{q_n} \left[\sum_n \log p(X_n, Z_n | \theta) \right] &= E_{q_n} \left[\sum_n \log p(Z_n | \pi) + \sum_n \log p(X_n | Z_n, \mu, \Sigma) \right] \\ &= E_{q_n} \left[\sum_n \log \left(\prod_k \pi_k^{Z_{nk}} \right) + \sum_n \log \left(\prod_k N(X_n | \mu_k, \Sigma_k)^{Z_{nk}} \right) \right] \\ &= \sum_n \sum_k E_{q_n} [Z_{nk}] \log \pi_k + \sum_n \sum_k E_{q_n} [Z_{nk}] \log N(X_n | \mu_k, \Sigma_k) \\ &= \sum_n \sum_k r_{nk} \log \pi_k + \sum_n \sum_k r_{nk} \log N(X_n | \mu_k, \Sigma_k) \end{aligned}$$

EM for a GMM

E step:

expected complete data log likelihood:

$$E_{q_n} \left[\sum_n \log p(x_n, z_n | \theta) \right] = \sum_n \sum_k r_{nk} \log \pi_k + \sum_n \sum_k r_{nk} \log N(x_n | \mu_k, \Sigma_k)$$

M step:

set derivative to 0 w.r.t. π_k, μ_k, Σ_k

$$\pi_k = \frac{1}{N} \sum_n r_{nk} = \frac{r_k}{N} \quad r_k: \text{weighted \# points in cluster } k \quad (r_n = 1)$$

$$\mu_k = \frac{\sum_n r_{nk} x_n}{r_k} \quad \text{weighted average}$$

$$\Sigma_k = \frac{\sum_n r_{nk} x_n x_n^T}{r_k} - \mu_k \mu_k^T \quad \text{also weighted}$$

K-means as a special case of EM

| | GMM | K-means |
|------------|--|--|
| parameters | π, μ, Σ | μ ($\pi_k = \frac{1}{K}, \Sigma_k = I$) |
| E-step | $r_{nk} = \frac{\pi_k \mathcal{N}(x_n \mu_k, \Sigma_k)}{\sum_{j=1}^K \pi_j \mathcal{N}(x_n \mu_j, \Sigma_j)}$ | $z_{nk}^* = 1, \text{ where } k^* = \underset{j}{\operatorname{argmax}} r_{nj}$ ($r_{nk} = \mathbb{I}(z_{nk}^* = 1)$) |
| M-step | $\pi_k = \frac{1}{N} \sum_n r_{nk} = \frac{r_k}{N}$ $\mu_k = \frac{\sum_n r_{nk} x_n}{r_k}$ $\Sigma_k = \frac{\sum_n r_{nk} x_n x_n^T}{r_k} - \mu_k \mu_k^T$ | $\mu_k = \frac{\sum_{n: z_{nk}^* = 1} x_n}{N_k}$ |

II. Graphical Models

Terminology

PGM: probabilistic graphical models

DAG: directed acyclic graph

Bayesian networks (Bayes nets), belief networks: PGMs based on DAGs

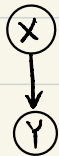
topological ordering: In DAG, nodes are ordered s.t. parents come before children

ordered Markov property: assumption that a node is conditionally independent (CI) of all predecessors given its parents

Markov chain: $(X_1) \rightarrow (X_2) \rightarrow (X_3) \dots \rightarrow (X_T)$

$$P(X_{1:T}) = P(X_1) \prod_{t=2}^T P(X_t | X_{1:t-1}) \quad \text{first-order Markov chain}$$

CPT: conditional probability table, 2d table representing conditional probability distribution (CPD)



| | Y=0 | Y=1 |
|-----|------|------|
| X=0 | 0.05 | 0.95 |
| X=1 | 0.2 | 0.8 |

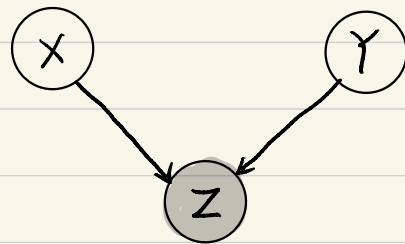
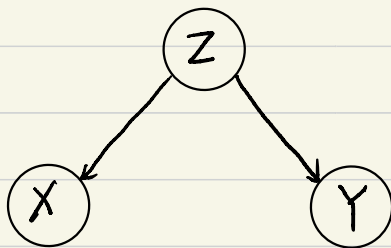
$P(Y|X)$

D-separation

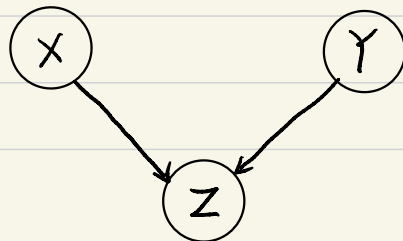
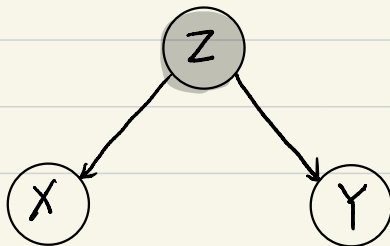
$X_A \perp_G X_B \mid X_C$ A is *conditionally independent* of B given C in graph G.
A is *d-separated* from B given C

We "shade" a node to indicate it's observed.

X, Y are not d-separated



X, Y are d-separated



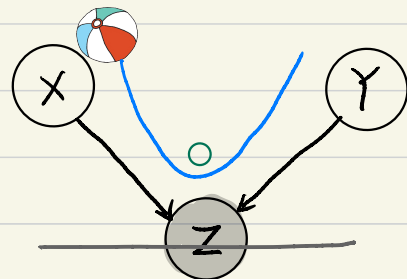
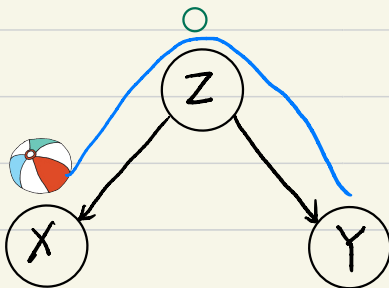
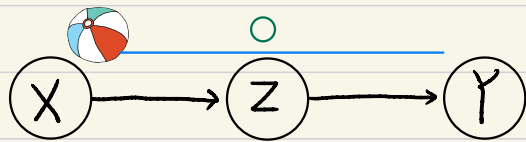
D-separation

Can X pass a ball to Y?

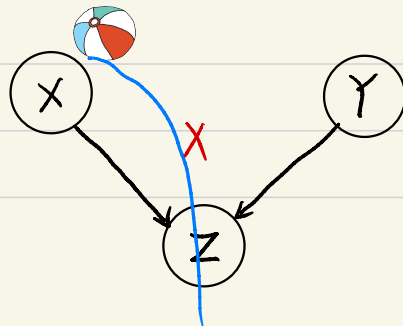
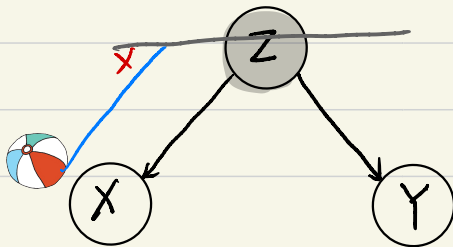
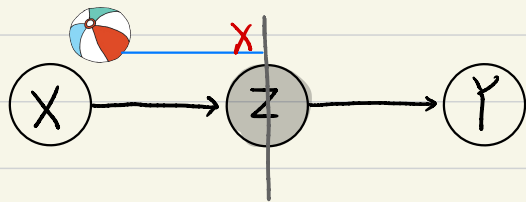
(Balls travel along edges and can travel opposite to edge directions.)

yes: not d-separated
no: d-separated

X, Y are not d-separated

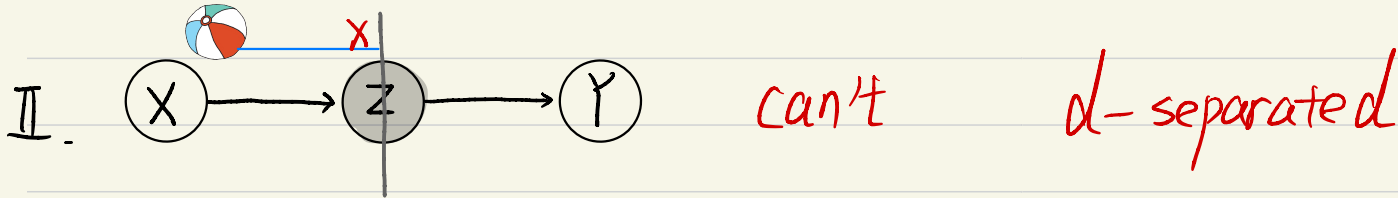
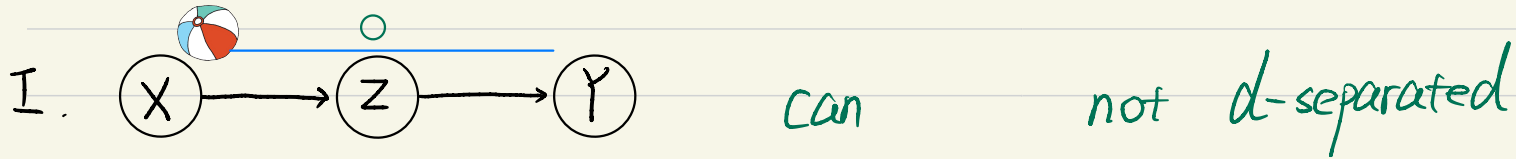


X, Y are d-separated



Can X pass a ball to Y?

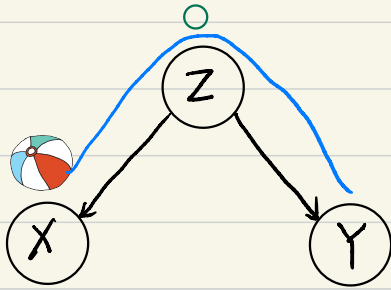
The ball travels along edges and can travel opposite to edge directions.



Can X pass a ball to Y?

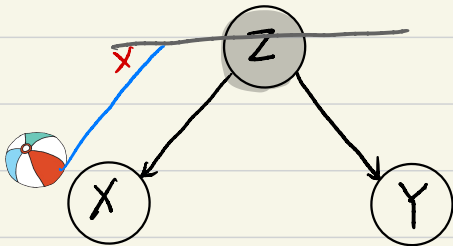
The ball travels along edges and can travel opposite to edge directions.

I.



can not d-separated

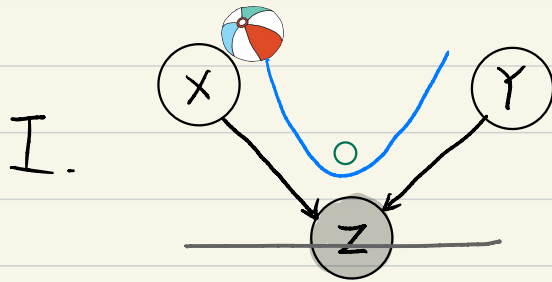
II.



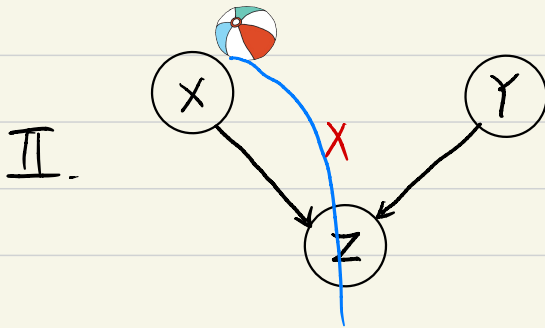
can't d-separated

Can X pass a ball to Y?

The ball travels along edges and can travel opposite to edge directions.



can not d-separated

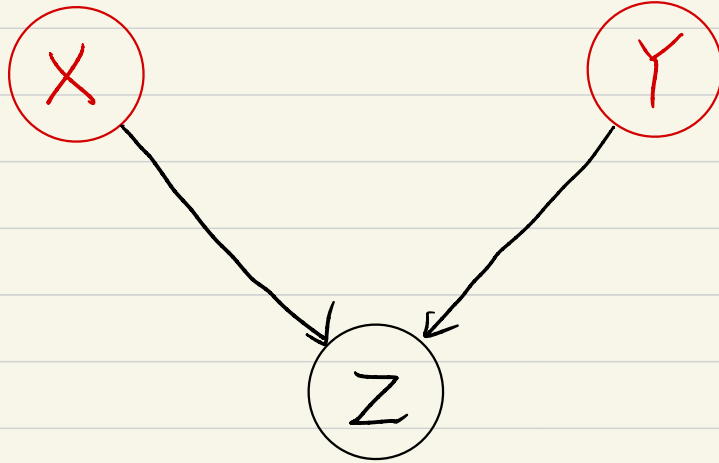


can't d-separated

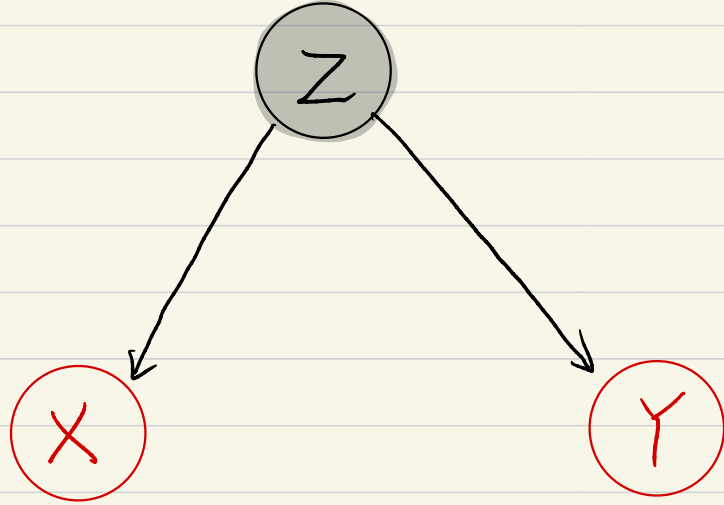


Are \textcircled{X} and \textcircled{Y} d-separated or not?

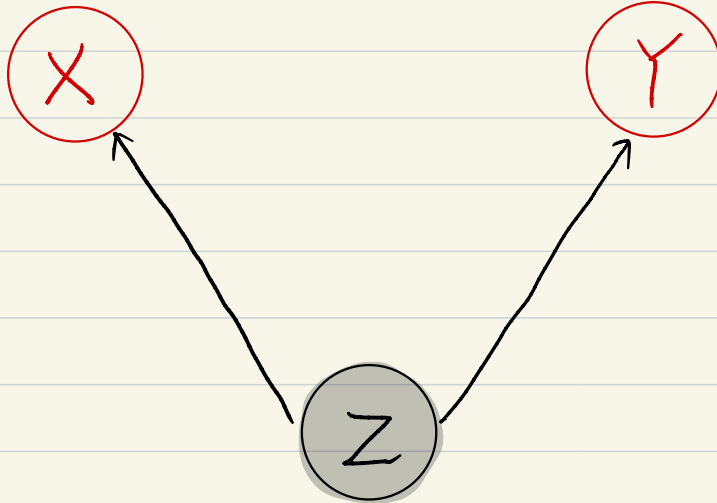
1. d-separated ?



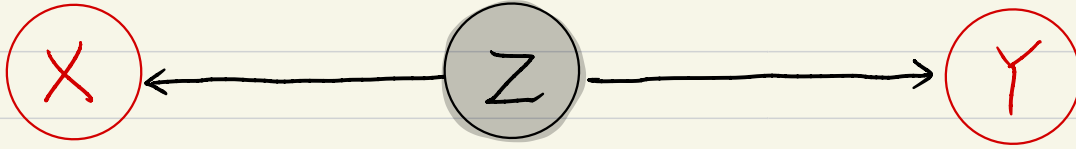
2. d-separated?



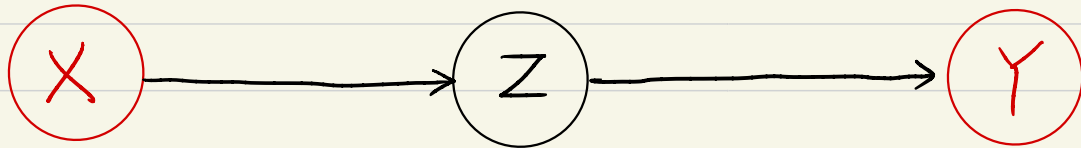
3. d-separated?



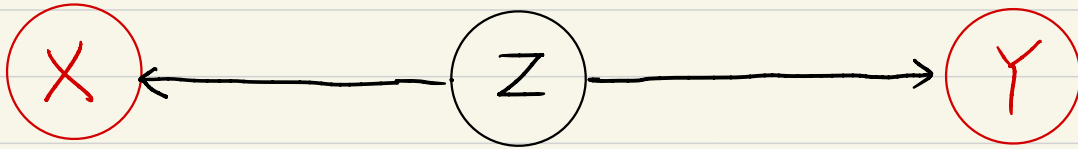
4. d-separated?



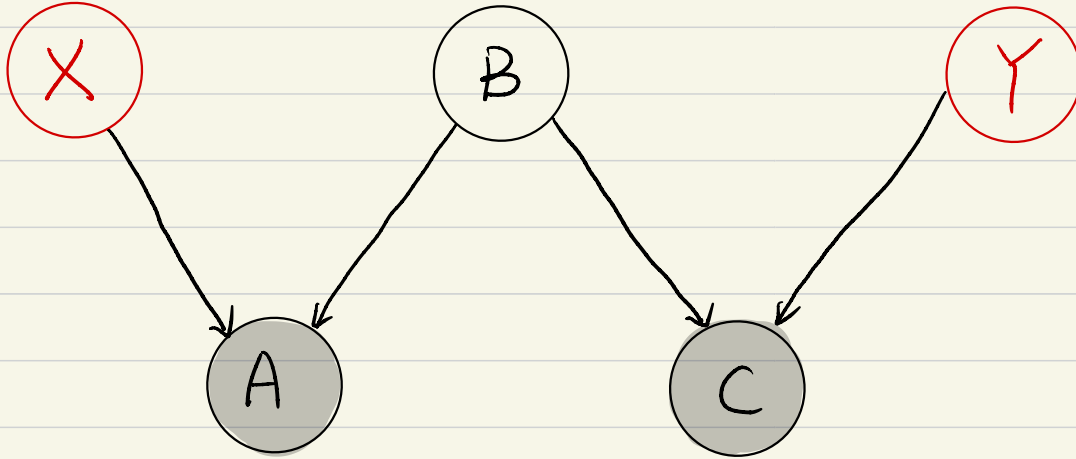
5. d-separated?



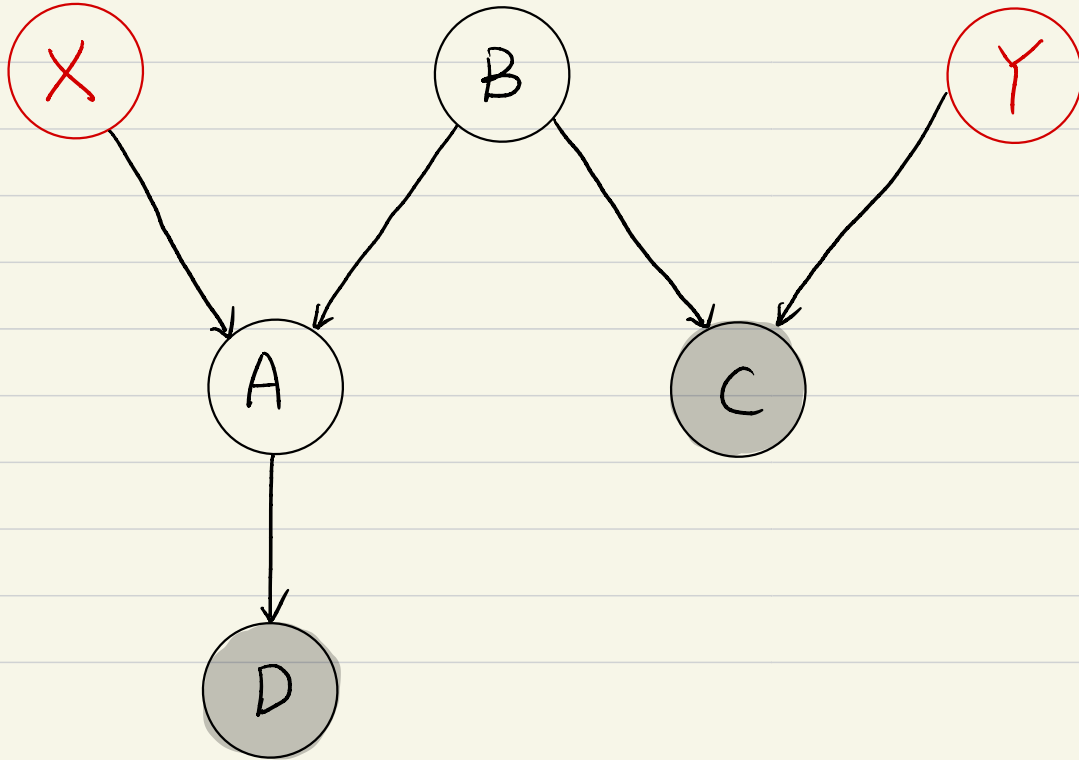
6. d-separated?



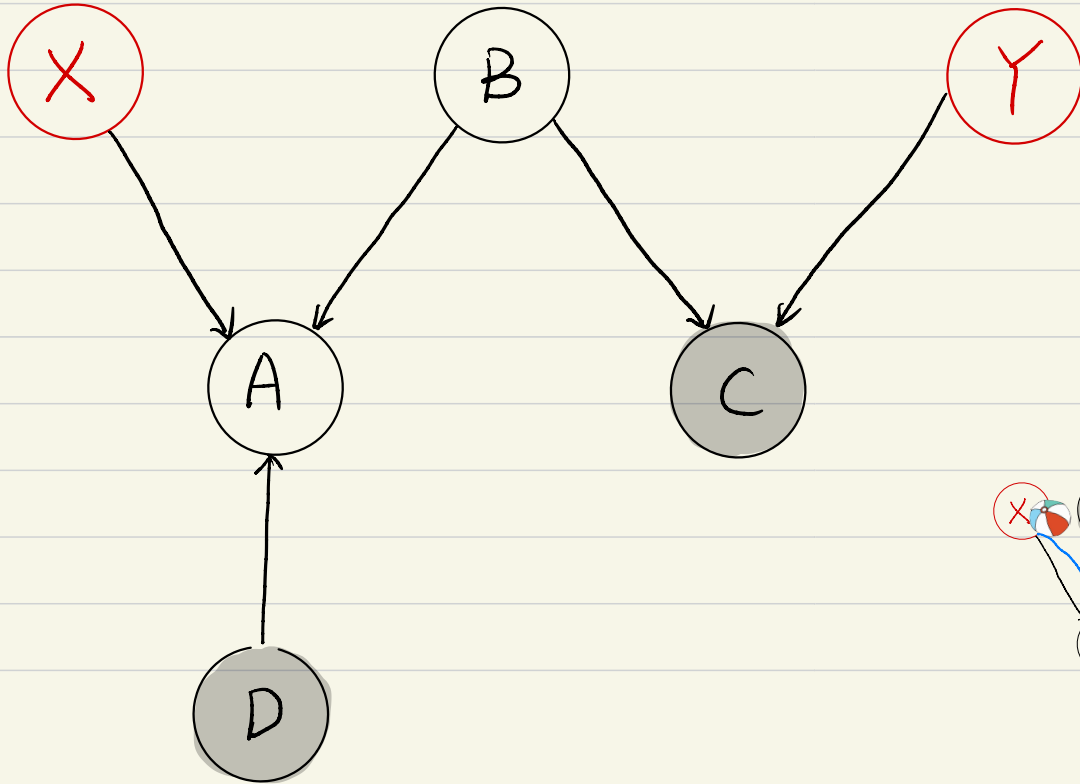
7. d-separated?



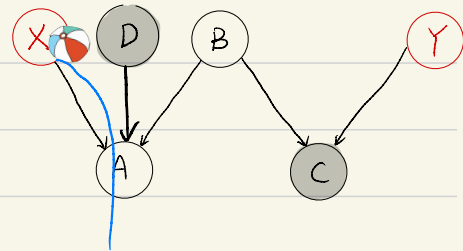
8. d-separated?



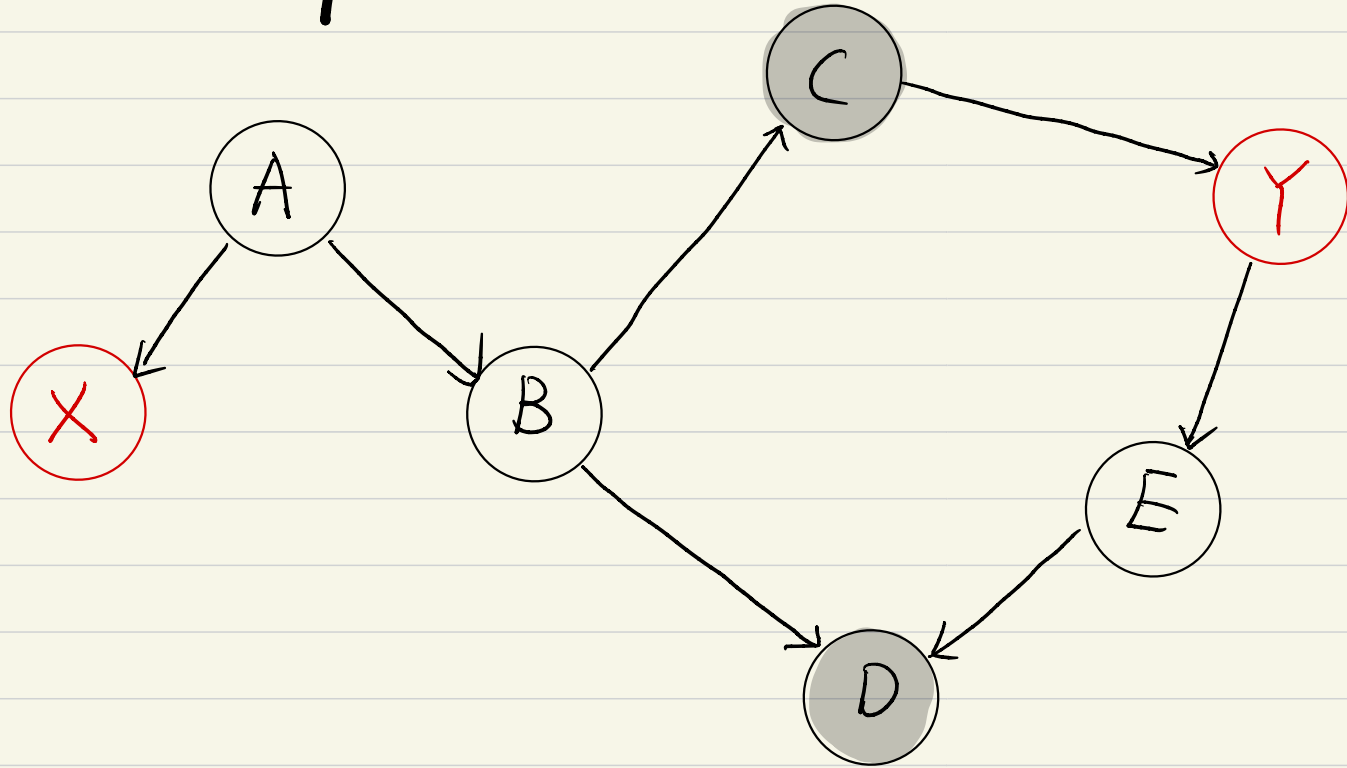
9. d-separated?



Yes!

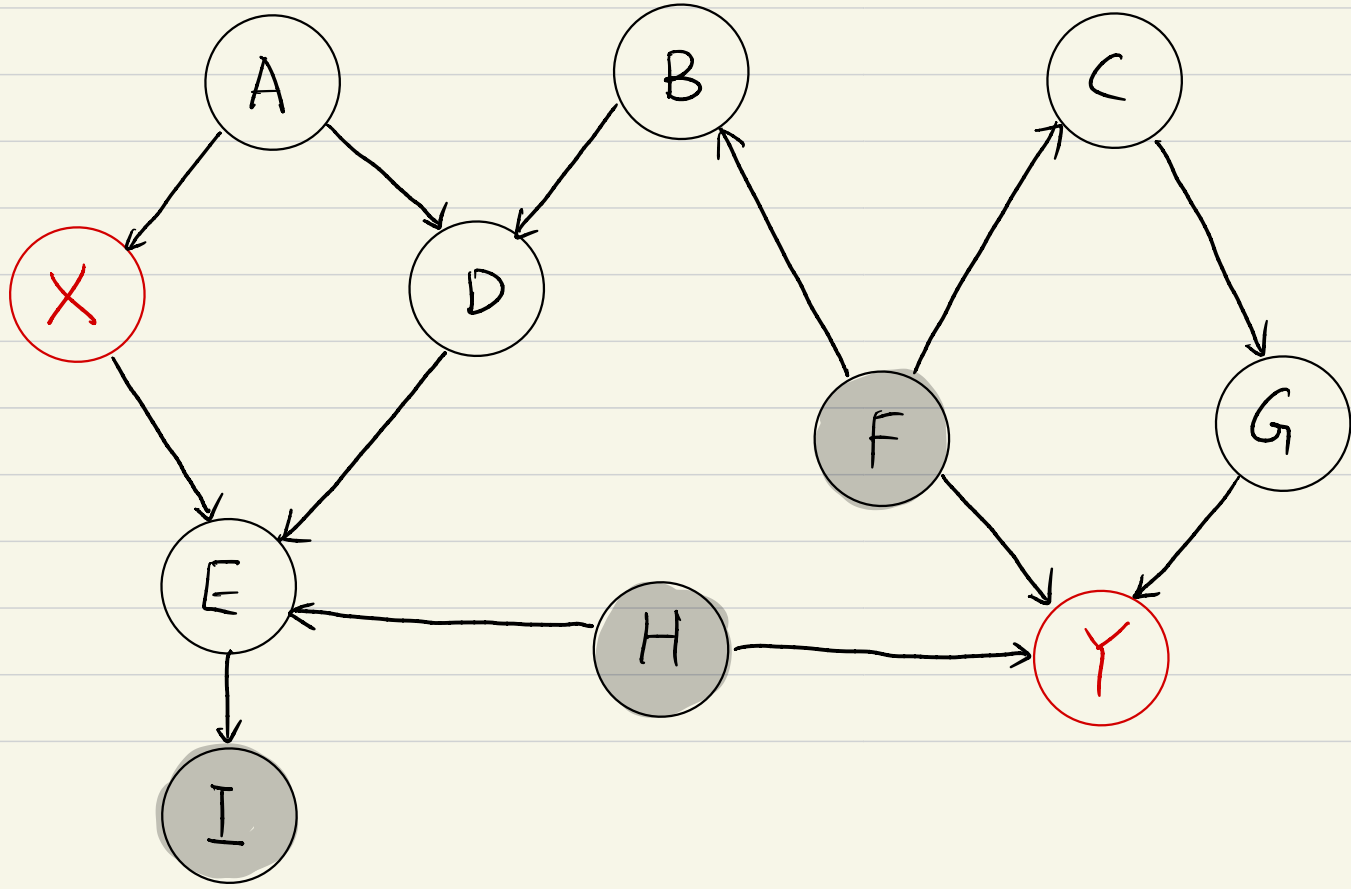


10. d-separated?



11. d-separated?

The last one!



Answers

1 - 4 : Yes

5 - 8 : No

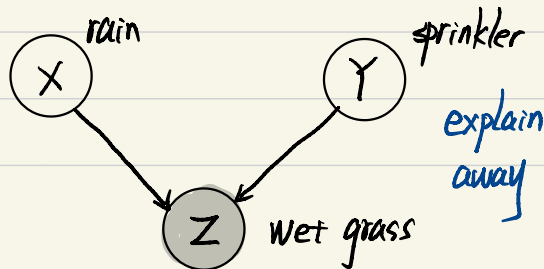
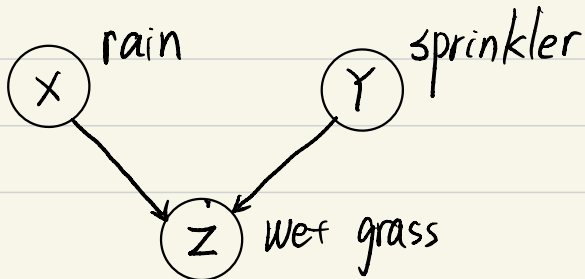
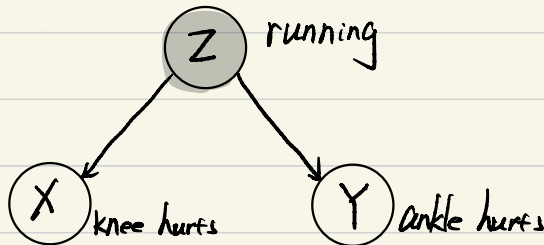
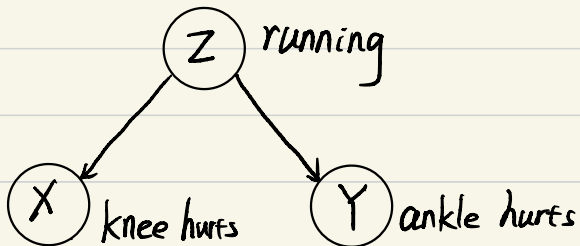
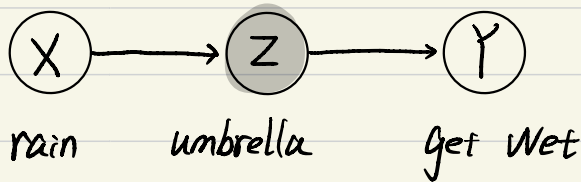
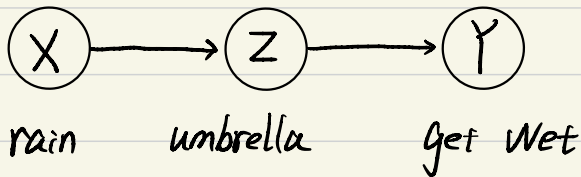
9 : Yes

10 : No

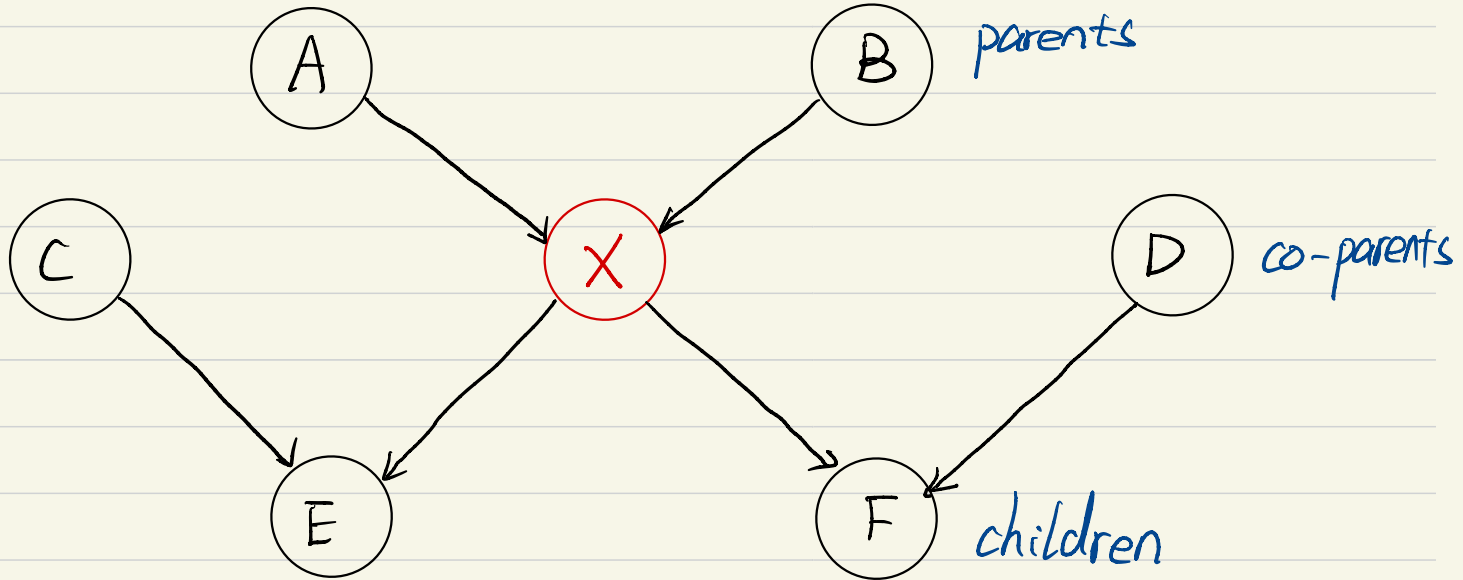
11 : Yes



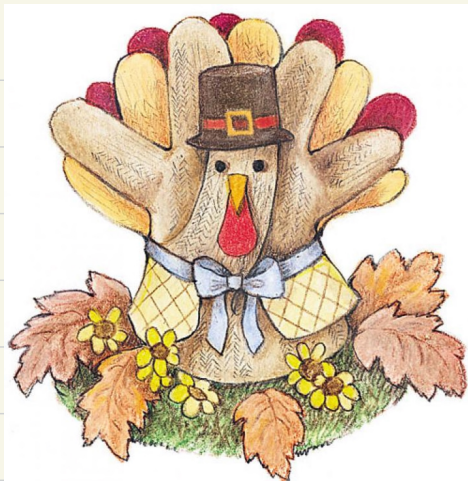
What does it mean?



Markov Blanket



A node conditioned on its Markov blanket is independent from all other nodes in the graph.



HAPPY
 Thanksgiving

