

- Suppose we have a collection of jointly distributed random variables

$\underline{X} = [X_1, \dots, X_d]^T$ with joint pdf (or pmf) $f_{\underline{X}}(x_1, \dots, x_d)$.

Questions

- ① How can we compactly represent this joint distribution?
 - ② If the variables are observable, how can we learn such a compact representation from data?
 - ③ How can we use this representation to infer the distribution for one subset of the variables given another in a reasonable amount of time?
- These are the questions to which the theory of graphical models aim to provide a general answer. We will see that the theory of graphical models contains a collection of theorems for answering these questions. As the theorems are model agnostic they can be applied broadly to problems in statistical modeling and machine learning.
 - In this module, we will prove these basic theorems and derive their consequences for the three questions above.

1 Representation

- The first goal of graphical models is to provide an informative representation of the joint distribution $f_{\underline{X}}(x_1, \dots, x_d)$.
- To do this we start with a rather uninformative representation that can be derived for any joint distribution.

Theorem ① (Chain Rule) The pdf (pmf) for $\underline{X} = [X_1, \dots, X_d]^T$ satisfies

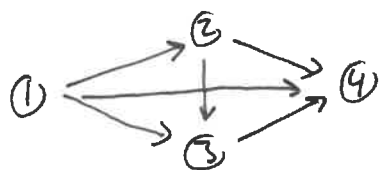
$$f_{\underline{X}}(x_1, \dots, x_d) = \prod_{i=1}^d f_{X_i | \underline{X}_{[i-1]}}(x_i | \underline{x}_{[i-1]})$$

$$([m] = \{1, \dots, m\} \quad \text{and} \quad \underline{X}_S = [X_i : i \in S]^T, \quad \underline{x}_S = [x_i : i \in S]^T)$$

(2)

- To represent the distribution we can draw a graph representing the dependencies specified by the conditional factors: $(f_{\underline{x}_i | \underline{x}_S}) \Rightarrow j \rightarrow i \forall j \in S$

$d = 4$ $f_{\underline{x}}(x_1, x_2, x_3, x_4) = f_{x_1}(x_1) f_{x_2 | x_1}(x_2 | x_1) f_{x_3 | x_1, x_2}(x_3 | x_1, x_2) f_{x_4 | x_1, x_2, x_3}(x_4 | x_1, x_2, x_3)$



- Pretty useless but still captures intuition for classic models:

Example 1 (Beta-Binomial) Consider the joint distribution $\underline{z} = [\underline{x}, \theta]^T$

where

$$\underline{x} | \theta = \theta \sim \text{Bin}(n, \theta)$$

$$\theta \sim \text{Beta}(\alpha, \beta)$$

The joint distribution is

$$f_{\underline{z}}(x, \theta) = f_{\underline{x} | \theta}(x | \theta) f_{\theta}(\theta) = \binom{n}{x} \theta^x (1-\theta)^{n-x} \frac{\Gamma(\alpha+\beta)}{\Gamma(\alpha)\Gamma(\beta)} \theta^{\alpha-1} (1-\theta)^{\beta-1}$$

which we represent as

$$\theta \rightarrow \underline{x}$$

("x depends on θ ").

- The graph is not so informative yet. But it becomes increasingly useful as the model becomes increasingly complex.
- When building models for many observables with many parameters we tend to make use of conditional independence assumptions. The assumptions lead to more interesting/informative graph structure.

Definition 1 For $\underline{x} = [x_1, \dots, x_d]^T$ and subsets $A, B, C \subseteq [d]$ disjoint with $A, B \neq \emptyset$, we say $\underline{x}_A \perp\!\!\!\perp \underline{x}_B | \underline{x}_C$ (~~with~~ \underline{x}_A is independent of \underline{x}_B given \underline{x}_C) if

$$f_{\underline{x}_A | \underline{x}_B | \underline{x}_C}(x_A, x_B | x_C) = f_{\underline{x}_A | \underline{x}_C}(x_A | x_C) f_{\underline{x}_B | \underline{x}_C}(x_B | x_C)$$

for $\forall \underline{x}_A, \underline{x}_B, \underline{x}_C$ with $f_{\underline{x}_C}(\underline{x}_C) > 0$.

Theorem ② The following are equivalent:

③

(a) $\underline{X}_A \perp\!\!\!\perp \underline{X}_B \mid \underline{X}_C$

(b) $f_{\underline{X}_A \mid \underline{X}_B, \underline{X}_C}(\underline{x}_A \mid \underline{x}_B, \underline{x}_C) = f_{\underline{X}_A \mid \underline{X}_C}(\underline{x}_A \mid \underline{x}_C)$ for all $\underline{x}_A, \underline{x}_B, \underline{x}_C$

(c) $f_{\underline{X}_A \mid \underline{X}_B, \underline{X}_C}(\underline{x}_A \mid \underline{x}_B, \underline{x}_C) = f_{\underline{X}_A \mid \underline{X}_B, \underline{X}_C}(\underline{x}_A \mid \underline{x}'_B, \underline{x}_C)$ for all $\underline{x}_A, \underline{x}_B, \underline{x}'_B, \underline{x}_C$

Proof: (Exercise!)

- When our model makes certain conditional independence assumptions applying Theorem ② (b) to factors in our chain rule factorization can lead to more informative graph structure.

Example ② (A hierarchical model) Consider $\underline{Z} = [\underline{\Lambda}, N, \underline{Y}]^T$ defined by the three-stage hierarchy

$$\underline{Y} \mid N=n \sim \text{Bin}(n, p) \quad (p \text{ fixed})$$

$$N \mid \underline{\Lambda}=\lambda \sim \text{Po}(\lambda)$$

$$\underline{\Lambda} \sim \text{Gamma}(\alpha, \beta)$$

where we assume $\underline{Y} \perp\!\!\!\perp \underline{\Lambda} \mid N$.

By the chain rule,

$$f_{\underline{\Lambda}, N, \underline{Y}}(\lambda, n, \gamma) = f_{\underline{\Lambda}}(\lambda) f_{N \mid \underline{\Lambda}}(n \mid \lambda) f_{\underline{Y} \mid N, \underline{\Lambda}}(\gamma \mid n, \lambda)$$

$$= f_{\underline{\Lambda}}(\lambda) f_{N \mid \underline{\Lambda}}(n \mid \lambda) f_{\underline{Y} \mid N}(\gamma \mid n) \quad \leftarrow \begin{array}{l} \text{Theorem 2.1} \\ + \\ \underline{Y} \perp\!\!\!\perp \underline{\Lambda} \mid N \end{array}$$

Get the graph

$$\underline{\Lambda} \rightarrow N \rightarrow \underline{Y}.$$

General Formula for getting a Graph Representation:

① Apply chain rule

② Apply available CI. relations to shrink conditioning sets:

$$f_{\underline{X}_i \mid \underline{X}_{C_i-1}}(\underline{x}_i \mid \underline{x}_{C_i-1}) = f_{\underline{X}_i \mid \underline{X}_{C_i}}(\underline{x}_i \mid \underline{x}_{C_i}), \quad C_i \subseteq [i-1]$$

③ draw $j \rightarrow i \quad \forall j \in C_i, \forall i \in [d]$.

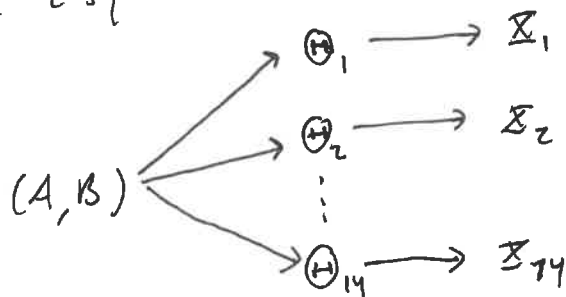
Example ③ (A Bayesian Hierarchical Model)

14 different laboratories have control groups of rats for a cancer study where laboratory i 's control group contains n_i rats. Let θ_i be the probability that a rat in control group i develops tumors. Given that the laboratories have similar practices and resources we expect some correlations between $\theta_1, \dots, \theta_{14}$. So we treat them as a random sample from a $\text{Beta}(\alpha, \beta)$ -population. We have no good guess for the parameters α, β so we put a diffuse (but proper) prior on them with pdf $f_{A,B}(\alpha, \beta)$. We also assume the number of rats x_i in group i that develops tumors depends only on θ_i .

We could write these CI assumptions explicitly, but it would be difficult to parse given their number. More easily we encode them in the pdf's factorization:

$$f_{\mathbf{x}, \boldsymbol{\theta}, A, B}(\mathbf{x}, \boldsymbol{\theta}, \alpha, \beta) = f_{A,B}(\alpha, \beta) \prod_{i=1}^{14} f_{\theta_i | A, B}(\theta_i | \alpha, \beta) \prod_{i=1}^{14} f_{x_i | \theta_i}(x_i | \theta_i).$$

Even more easy to read is the associated graph:



- Defining a graph as such according to a factorization always results in a DAG (Exercise!).
- The distribution together with its resulting DAG is a pair that we call a DAG model.

Definition (2) $\underline{x} = [x_1, \dots, x_d]^T$ is Markov to a DAG $G = (E_d, E)$ (5)

if

$$f_{\underline{x}}(x_1, \dots, x_d) = \prod_{i=1}^d f_{x_i | \underline{x}_{pa_G(i)}}(x_i | \underline{x}_{pa_G(i)}).$$

- The pair (\underline{x}, G) is called a DAG model.
- The set $M_F(G) = \{ \underline{x} : \underline{x} \text{ is Markov to } G \}$ is also called a DAG model.
- Sometimes (\underline{x}, G) is called a Bayesian network since Bayesian inference works backwards along the edges of the DAG.
- Sometimes (\underline{x}, G) is called a causal model, although a priori there need not be any causal information in the model (there can be though!).
- In the Bayesian setting, using the graph representation can help us quickly identify ways to simplify posterior computations.

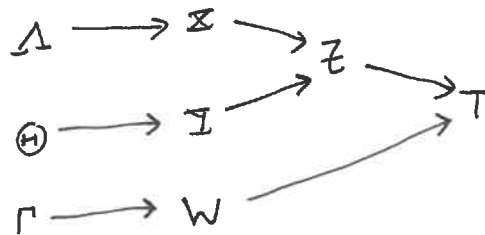
Example (4) Suppose we have a model where

$$T = Z + W, \quad Z = X + Y,$$

$$X | \Lambda = \lambda \sim P_0(\lambda), \quad Y | \Theta = \theta \sim P_0(\theta), \quad W | \Gamma = \gamma \sim P_0(\gamma), \text{ and}$$

Λ, Θ, Γ are independent and Gamma(1, 1)-distributed.

- Rather than specifying a list of CI relations we instead say that $\underline{x} = [\Lambda, \Theta, \Gamma, X, Y, W, Z, T]^T$ is Markov to G :



Given data $[Z, W, T]^T = \mathbf{D}$ we want the posterior

$$f_{\Lambda, \Theta, \Gamma | \mathbf{D}}(\lambda, \theta, \gamma | \mathbf{D}).$$

Since all paths from Λ, Θ to Γ in G are "blocked" by Z, W, T it turns out that $\Lambda, \Theta \perp \Gamma | Z, W, T$.

(6)

- This means we can break up our posterior computation in two

$$f_{\Delta, \Theta, \Gamma | D}(x, \theta, \gamma | D) = f_{\Delta, \Theta | D}(x, \theta | D) f_{\Gamma | D}(\gamma | D)$$

- Note that the CI relation $\Delta, \Theta \perp\!\!\!\perp \Gamma | \mathcal{E}, W, T$ is not a relation that we assume when we assume \underline{X} is Markov to G (Exercise: Check this!)
- Instead this relation is implied by our assumption.
- The fact that we can read CI relations implied by the Markov assumption is one useful property of DAG models!

Question Suppose \underline{X} is Markov to G . What is the complete set of CI relations that we know hold in \underline{X} ? (i.e. What CI relations are implied by the Markov assumption?)

- What does "implied" mean?

The Conditional Independence Axioms

(a) (Symmetry) $\underline{X}_A \perp\!\!\!\perp \underline{X}_B | \underline{X}_C \Rightarrow \underline{X}_B \perp\!\!\!\perp \underline{X}_A | \underline{X}_C$

(b) (Decomposition) $\underline{X}_A \perp\!\!\!\perp \underline{X}_{B \cup D} | \underline{X}_C \Rightarrow \underline{X}_A \perp\!\!\!\perp \underline{X}_B | \underline{X}_C$

(c) (Weak Union) $\underline{X}_A \perp\!\!\!\perp \underline{X}_{B \cup D} | \underline{X}_C \Rightarrow \underline{X}_A \perp\!\!\!\perp \underline{X}_B | \underline{X}_{C \cup D}$

(d) (Contraction) $\underline{X}_A \perp\!\!\!\perp \underline{X}_B | \underline{X}_{C \cup D}$ and $\underline{X}_A \perp\!\!\!\perp \underline{X}_D | \underline{X}_C \Rightarrow \underline{X}_A \perp\!\!\!\perp \underline{X}_{B \cup D} | \underline{X}_C$

(e) (Intersection) If $f_{\underline{X}}(x) > 0 \forall x$ then if $\underline{X}_A \perp\!\!\!\perp \underline{X}_B | \underline{X}_{C \cup D}$ and $\underline{X}_A \perp\!\!\!\perp \underline{X}_D | \underline{X}_{C \cup B} \Rightarrow \underline{X}_A \perp\!\!\!\perp \underline{X}_{B \cup D} | \underline{X}_C$.

Example ⑤ If $\underline{X} = [x_1, x_2, x_3, x_4]^T$ is Markov to

$$\textcircled{1} \rightarrow \textcircled{2} \rightarrow \textcircled{3} \rightarrow \textcircled{4}$$

then by definition we know $x_4 \perp\!\!\!\perp x_1, x_2 | x_3$ and $x_3 \perp\!\!\!\perp x_1 | x_2$.

Using the CI axioms we can show $x_1 \perp\!\!\!\perp x_4 | x_2$ (Exercise!)

Theorem ③ (Verma, Pearl) The complete set of CI relations implied by the Markov condition is

$$C(G) = \{ \underline{X}_A \perp\!\!\!\perp \underline{X}_B \mid \underline{X}_C : A, B \text{ d-separated given } C \text{ in } G \}.$$

• By the end of the first two lectures we will have proven

Theorem ④ Let $\underline{X} = [X_1, \dots, X_d]^T$ a distribution and $G = ([d], E)$ a DAG.

The following are equivalent:

(a) \underline{X} is Markov to G

(b) \forall triples $A, B, C \subseteq [d]$ where A, B are d-separated given C in G we have $\underline{X}_A \perp\!\!\!\perp \underline{X}_B \mid \underline{X}_C$.

2 Structure Learning

• We may find ourselves having jointly distributed observed variables $\underline{X} = [X_1, \dots, X_d]^T$, no idea how to parametrize the model, but would still like to learn a graphical structure that the distribution follows.

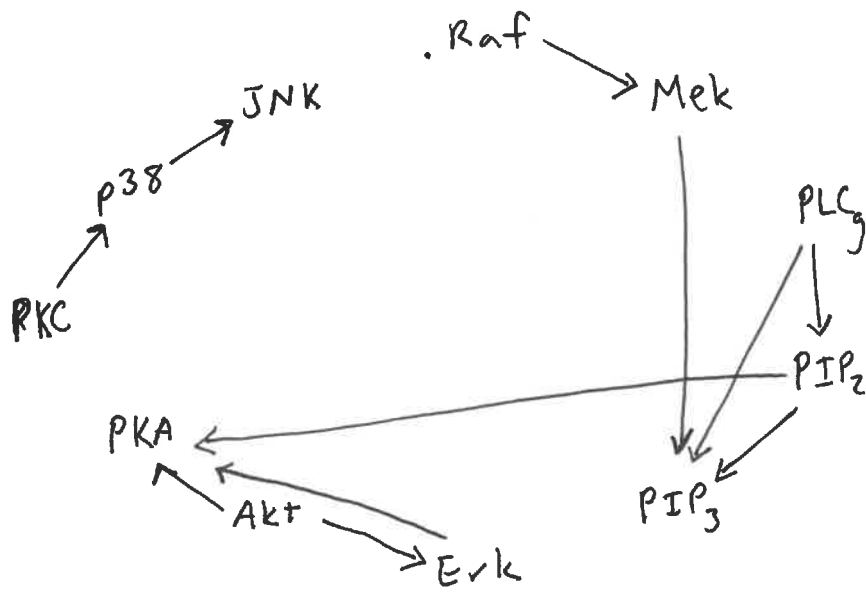
• This is the following unsupervised learning problem:

Problem Given a random sample (data) from $\underline{X} = [X_1, \dots, X_d]^T$ can we find a "good" or "the best" DAG to which \underline{X} is Markov?

• This problem is central in the field of causal inference, where it is called DAG structure learning or causal discovery.

• Applications are numerous and include fields like bioinformatics:

Example ⑥ (Protein Signaling Network) The abundances of 14 different phosphoproteins and phospholipids in primary human immune system cells were measured in 1755 individuals cells. Treating the joint distribution of these 14 different molecules as a multivariate normal distribution we can apply one of the causal discovery algorithms that we will see in this class to estimate a DAG to which the data-generating distribution is Markov:

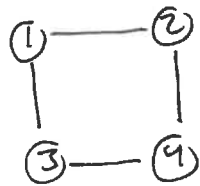


(PC Algorithm, $\alpha = 0.05$)

- The causal discovery algorithm used here relies on Theorem 7, using the characterization of the Markov property given in (b).
- To prove this theorem, and develop such algorithms, we will first prove similar theorems for undirected graphical models, which similarly provide useful representations of joint distributions:

Definition 3 A distribution $\underline{x} = [x_1, \dots, x_d]^T$ is Markov to an undirected graph $G = ([d], E)$ if $x_i \perp x_j \mid \underline{x}_{[d] \setminus \{i,j\}} \forall \{i,j\} \notin E$.

Example 7 $\underline{x} = [x_1, x_2, x_3, x_4]^T$ is Markov to



if $x_1 \perp x_4 \mid x_2, x_3$ and $x_2 \perp x_3 \mid x_1, x_4$.

- We can similarly learn undirected graph representations of a distribution.

Example 7 Returning to our molecules $\underline{x} = [x_1, \dots, x_{14}]^T$ from Example 6, we assume again that $\underline{x} \sim N(\mu, \Sigma)$ for some unknown mean vector $\mu \in \mathbb{R}^{14}$ and positive definite covariance matrix $\Sigma \in \mathbb{R}^{d \times d}$.

- Since we assume a normal model we know that $x_i \perp x_j \mid \underline{x}_{[14] \setminus \{i,j\}} \iff \text{Cov}[x_i, x_j \mid \underline{x}_{[14] \setminus \{i,j\}}] = 0$.

• To compute $\text{Cov}[\mathbf{z}_i, \mathbf{z}_j | \mathbf{z}_{[d] \setminus \{i,j\}}]$ we consider the covariance matrix for $[\mathbf{z}_i, \mathbf{z}_j]^T | \mathbf{z}_{[14] \setminus \{i,j\}} = \mathbf{z}_{[14] \setminus \{i,j\}} \sim \mathcal{N}(\mu_{[i,j]}, \Sigma'_{[i,j]})$ where

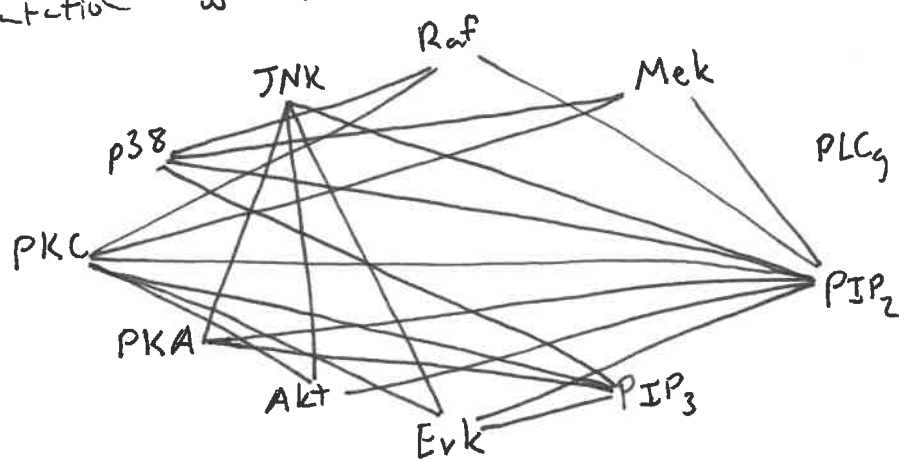
$$\mu'_{[i,j]} = \mu_{[i,j]} + \Sigma_{[i,j], [14] \setminus \{i,j\}} \Sigma_{[14] \setminus \{i,j\}, [14] \setminus \{i,j\}}^{-1} (\mathbf{z}_{[14] \setminus \{i,j\}} - \mu_{[14]})$$

$$\Sigma'_{[i,j]} = \Sigma_{[i,j], [i,j]} - \Sigma_{[i,j], [14] \setminus \{i,j\}} \Sigma_{[14] \setminus \{i,j\}, [14] \setminus \{i,j\}}^{-1} \Sigma_{[14] \setminus \{i,j\}, [i,j]}$$

• $\Sigma'_{[i,j]}$ is a 2×2 matrix whose off-diagonal entry is equal to $\text{Cov}[\mathbf{z}_i, \mathbf{z}_j | \mathbf{z}_{[d] \setminus \{i,j\}}]$.

• For each pair $\{i,j\} \in [14] \times [14]$ we can do a hypothesis test to check if this off-diagonal entry is (believably) zero.

• The result of these tests gives us the following undirected representation of relations in the data-generating distribution:

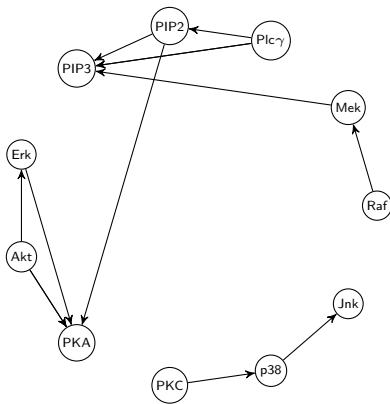


(Pearson correlation, $\alpha=0.0$)

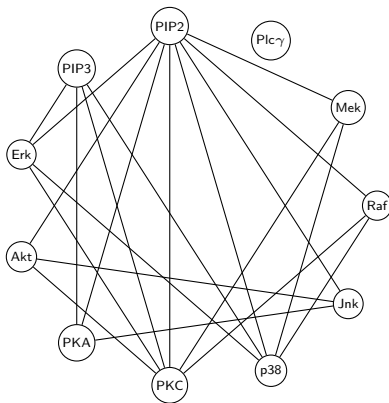
[3] Inference with Graph Representations

- Lectures (5) and (6) will begin the discussion of how we can use a graph representation for inference (either a DAG or an undirected graph)
- We already saw a little example of how a graph may help us compute posterior distributions more efficiently.
- In these lectures we will look at exact inference algorithms for posterior computations that are motivated by the graph structure

- (10)
- We will see that these inference algorithms have complexity bounds given by the structure of the graph.
 - In module (2), you will study approximate inference algorithms that may be used even in cases where the complexity bounds we see in this module suggest that using the exact inference methods is infeasible.



Sachs data set estimated DAG;
PC algorithm ($\alpha = 0.05$).



Sachs data set estimated undirected graph;
Pairwise CI testing with Pearson Correlation Tests ($\alpha = 0.05$).