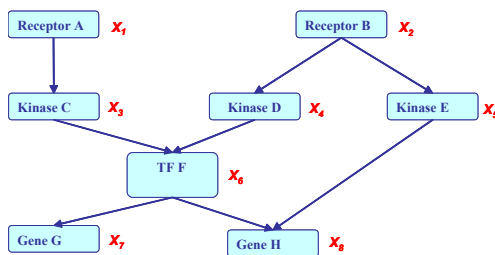
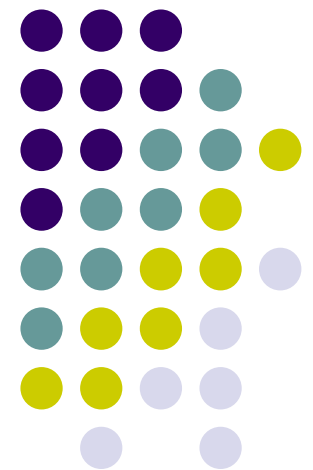


Machine Learning

10-701, Fall 2016

Graphical Models and Exact Inference

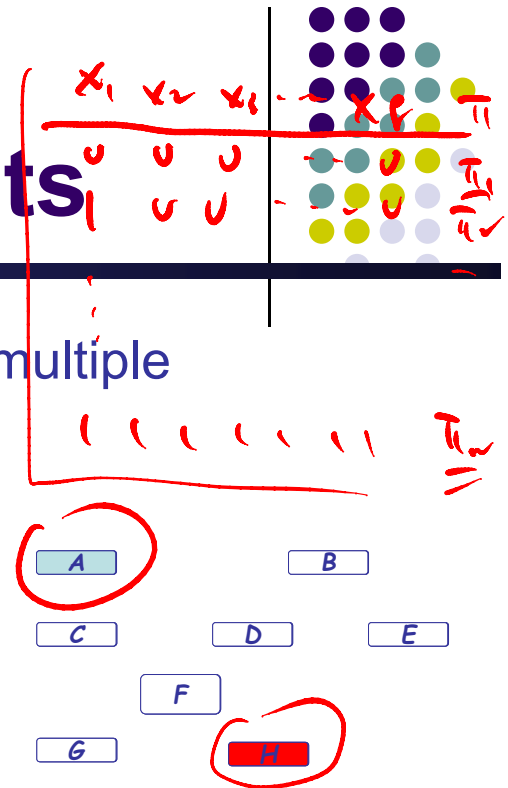
Eric Xing



Lecture 17, November 7, 2016

Reading: Chap. 8, C.B book

Recap of Basic Prob. Concepts



- Representation: what is the joint probability dist. on multiple variables?

$$P(X_1, X_2, X_3, X_4, X_5, X_6, X_7, X_8)$$

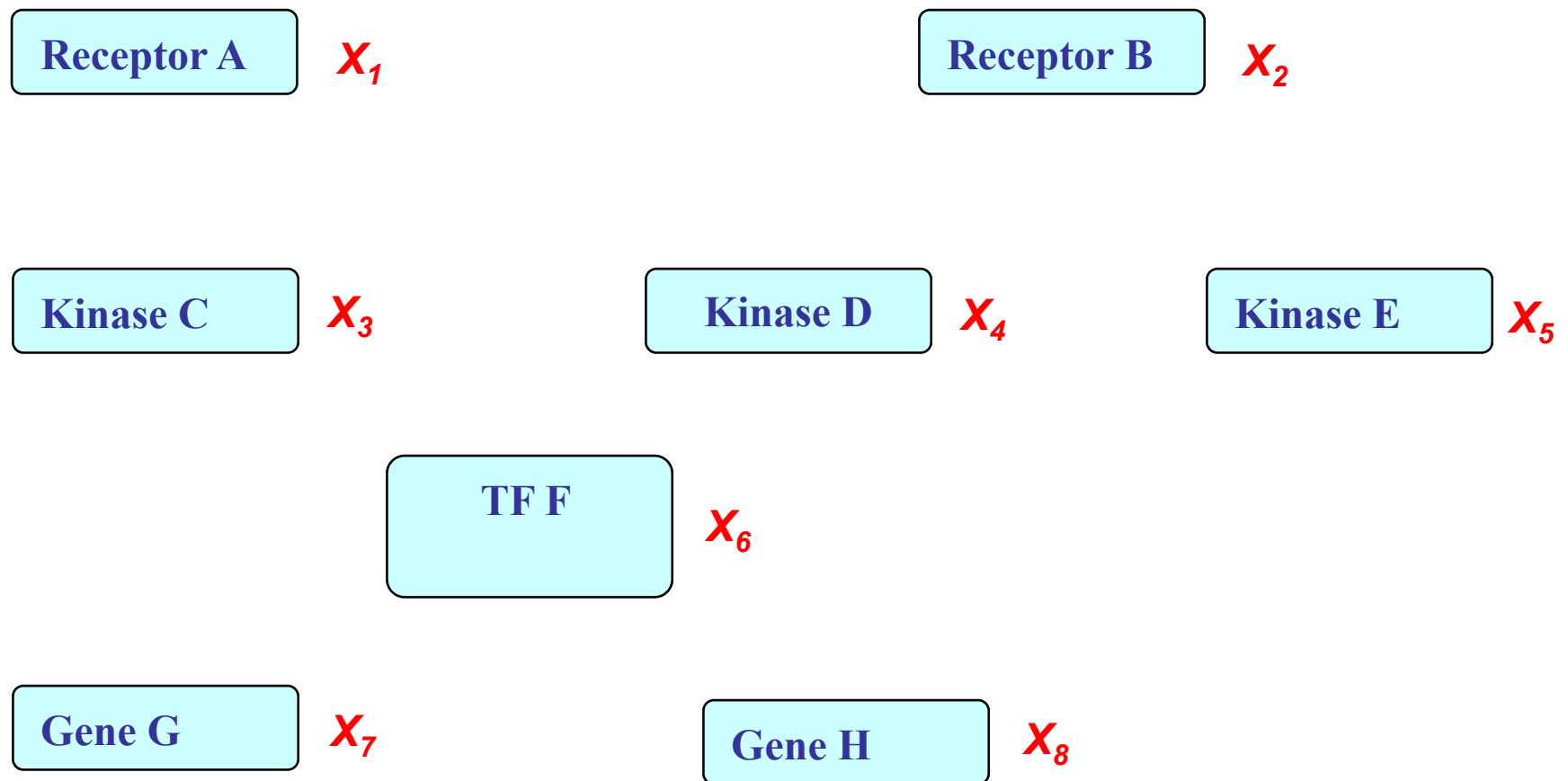
- How many state configurations in total? --- 2^8
 - Are they all needed to be represented?
 - **Do we get any scientific/medical insight?**
- Learning: where do we get all this probabilities?
 - Maximal-likelihood estimation? but how many data do we need?
 - Are there other est. principles?
 - Where do we put domain knowledge in terms of plausible relationships between variables, and plausible values of the probabilities?
 - Inference: If not all variables are observable, how to compute the conditional distribution of latent variables given evidence?
 - Computing $p(H|A)$ would require summing over all 2^6 configurations of the unobserved variables

What is a Graphical Model?

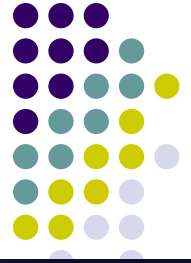
--- Multivariate Distribution in High-D Space



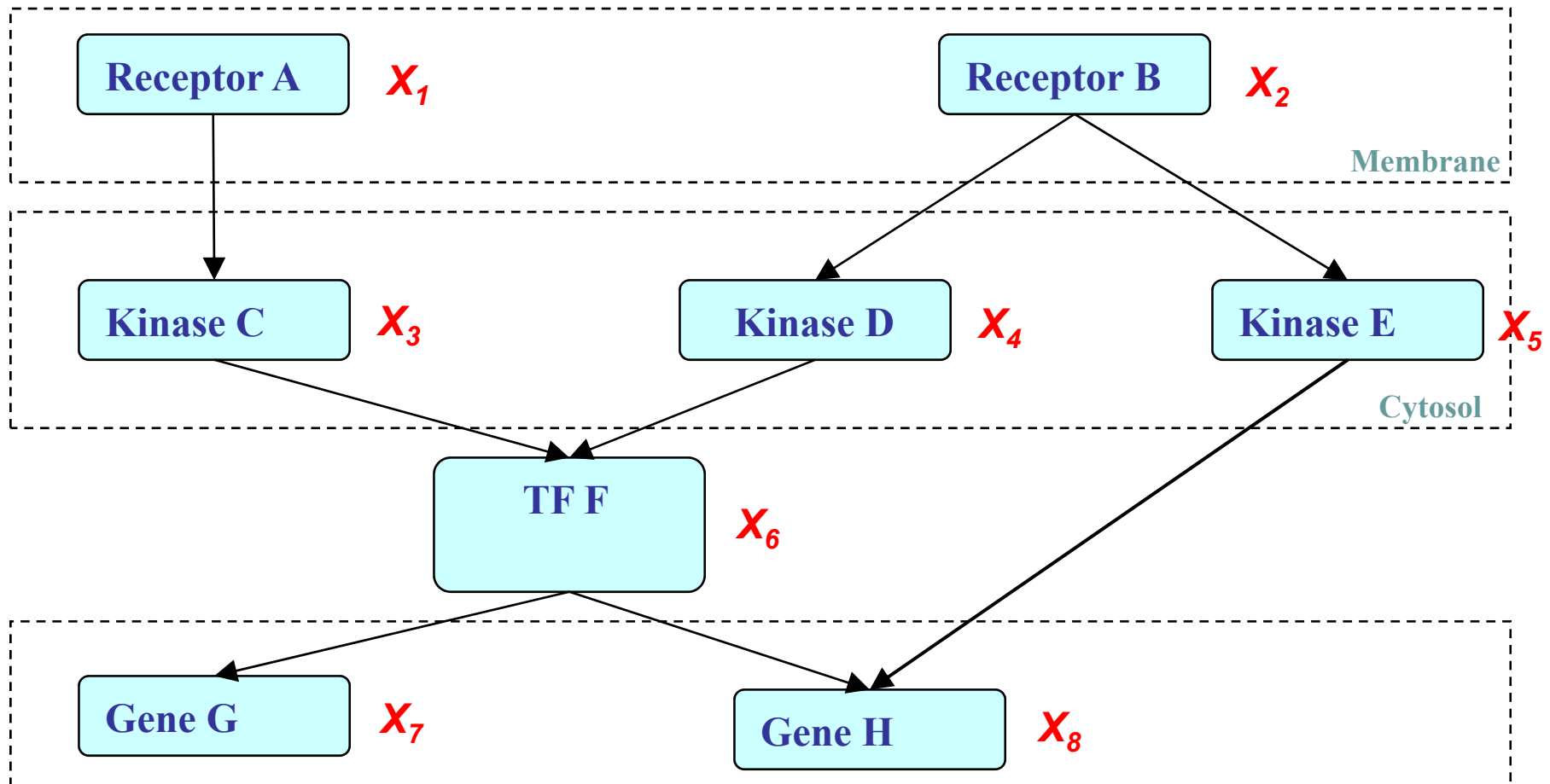
- A possible world for cellular signal transduction:



GM: Structure Simplifies Representation

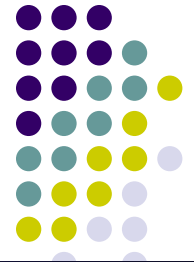


- Dependencies among variables

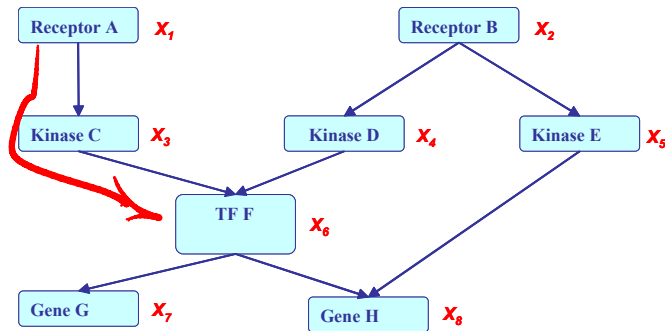


Probabilistic Graphical Models

$P(x_i | \text{Pa}(x_i))$
 $P(x_i)$



- If X_i 's are **conditionally independent** (as described by a **PGM**), the joint can be factored to a product of simpler terms, e.g.,



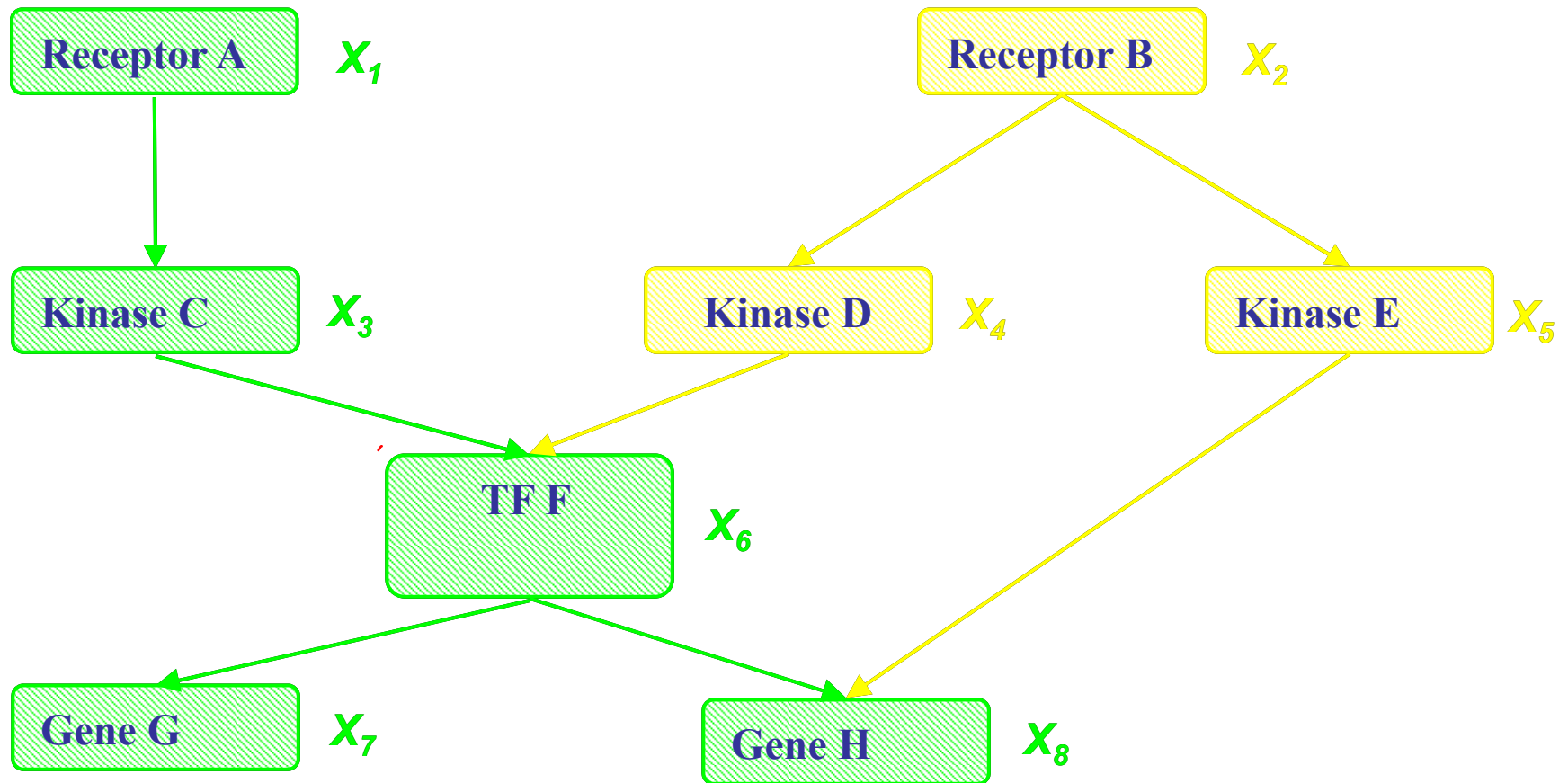
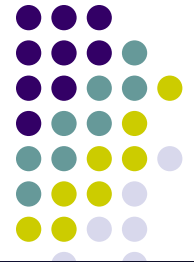
$$\begin{aligned}
 & \cancel{P(X_1, X_2, X_3, X_4, X_5, X_6, X_7, X_8)} \\
 &= P(X_1) P(X_2) P(X_3 | X_1) P(X_4 | X_2) P(X_5 | X_2) \\
 & P(X_6 | X_3, X_4) P(X_7 | X_6) P(X_8 | X_5, X_6) \quad \uparrow
 \end{aligned}$$

Stay tune for what are these independencies!

- Why we may favor a PGM?
 - Incorporation of domain knowledge and causal (logical) structures
 - 1+1+2+2+2+4+2+4=18, a 16-fold reduction from 2^8 in representation cost !

GM: Data Integration

Modularity



More Data Integration



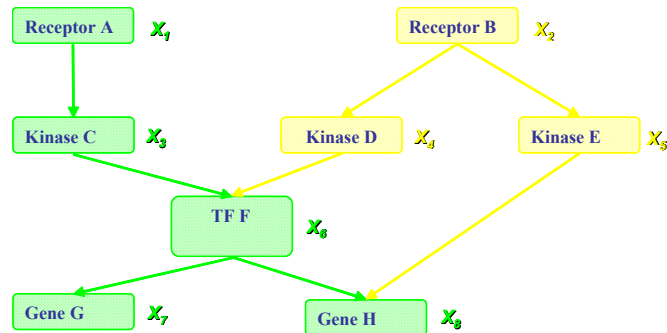
- Text + Image + Network → Holistic Social Media

- Genome + Proteome + Transcriptome + Phenome + ... → PanOmic Biology



Probabilistic Graphical Models

- If X_i 's are **conditionally independent** (as described by a **PGM**), the joint can be factored to a product of simpler terms, e.g.,



$$\begin{aligned} &P(X_1, X_2, X_3, X_4, X_5, X_6, X_7, X_8) \\ &= P(X_2) P(X_4 | X_2) P(X_5 | X_2) P(X_1) P(X_3 | X_1) \\ &P(X_6 | X_3, X_4) P(X_7 | X_6) P(X_8 | X_5, X_6) \end{aligned}$$

- Why we may favor a PGM?
 - Incorporation of domain knowledge and causal (logical) structures
 $2+2+4+4+4+8+4+8=36$, an 8-fold reduction from 2^8 in representation cost !
 - Modular combination of heterogeneous parts – data fusion



Rational Statistical Inference

The Bayes Theorem:

Posterior probability

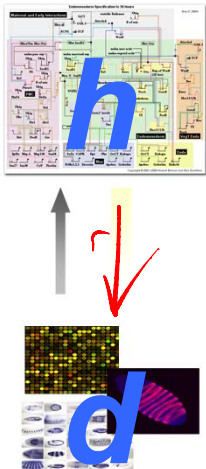
Likelihood

Prior probability

$$p(h | d) = \frac{p(d | h) p(h)}{\sum_{h' \in H} p(d | h') p(h')}$$

Sum over space of hypotheses

$p(d|h)$

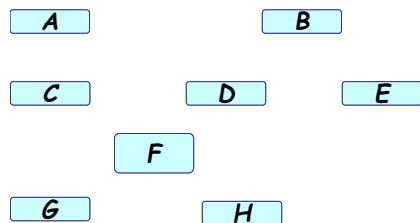


- This allows us to capture uncertainty about the model in a principled way
- But how can we specify and represent a complicated model?
 - Typically the number of genes need to be modeled are in the order of thousands!



GM: MLE and Bayesian Learning

- Probabilistic statements of Θ is conditioned on the values of the observed variables \mathbf{A}_{obs} and prior $p(\cdot | \chi)$

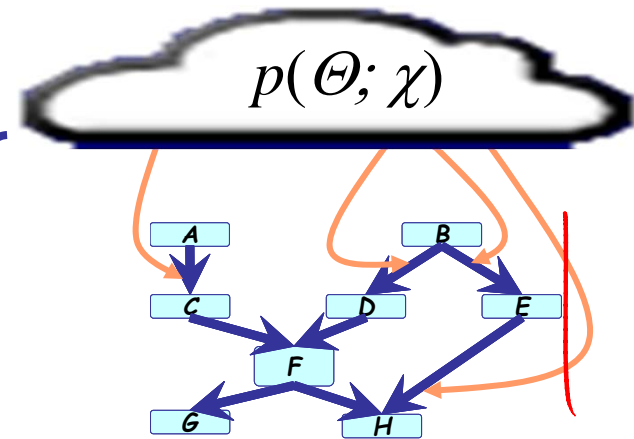
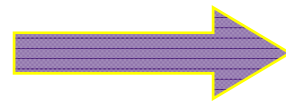


$(A,B,C,D,E,\dots)=(T,F,F,T,F,\dots)$

$\mathbf{A} = (A,B,C,D,E,\dots)=(T,F,T,T,F,\dots)$

.....

$(A,B,C,D,E,\dots)=(F,T,T,T,F,\dots)$



C	D	$P(F C,D)$	
c	d	0.9	0.1
c	\bar{d}	0.2	0.8
\bar{c}	d	0.9	0.1
\bar{c}	\bar{d}	0.01	0.99

$$p(\Theta | \mathbf{A}; \chi) \propto p(\mathbf{A} | \Theta) p(\Theta; \chi)$$

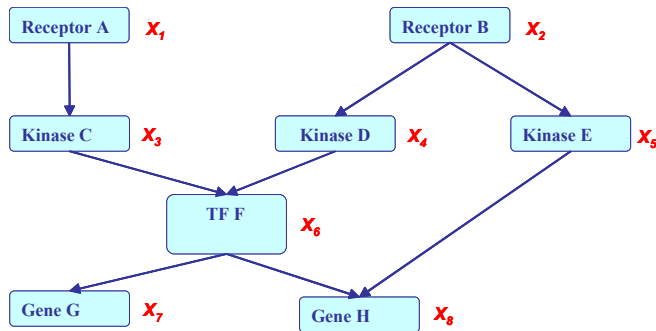
posterior
likelihood
prior

$$\Theta_{\text{Bayes}} = \int \Theta p(\Theta | \mathbf{A}, \chi) d\Theta$$



Probabilistic Graphical Models

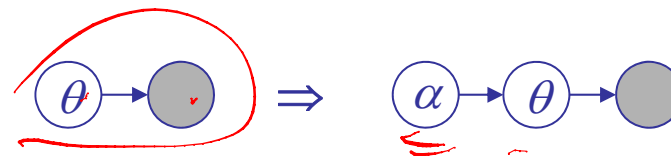
- If X_i 's are **conditionally independent** (as described by a **PGM**), the joint can be factored to a product of simpler terms, e.g.,



$$\begin{aligned}
 &P(X_1, X_2, X_3, X_4, X_5, X_6, X_7, X_8) \\
 &= P(X_1) P(X_2) P(X_3/X_1) P(X_4/X_2) P(X_5/X_2) \\
 &P(X_6/X_3, X_4) P(X_7/X_6) P(X_8/X_5, X_6)
 \end{aligned}$$

- Why we may favor a PGM?
 - Incorporation of domain knowledge and causal (logical) structures
 $2+2+4+4+4+8+4+8=36$, an 8-fold reduction from 2^8 in representation cost !
 - Modular combination of heterogeneous parts – data fusion
 - Bayesian Philosophy

- Knowledge meets data





So What Is a PGM After All?

In a nutshell:

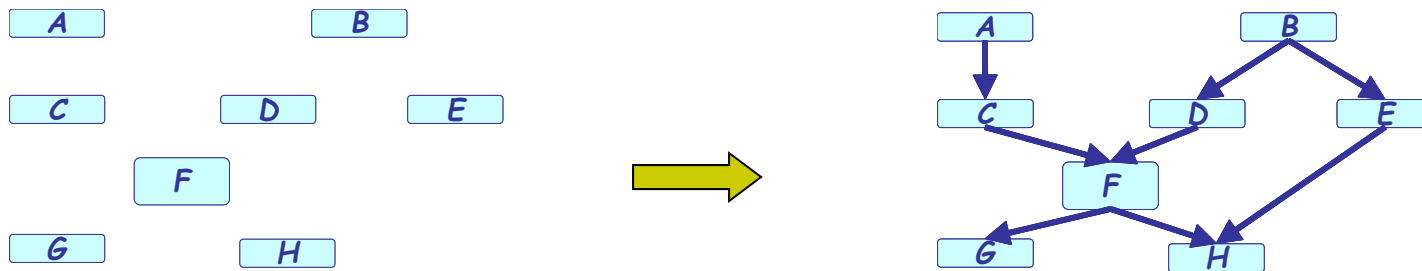
PGM = Multivariate Statistics + Structure

GM = Multivariate Obj. Func. + Structure



So What Is a PGM After All?

- The informal blurb:
 - It is a smart way to **write/specify/compose/design** exponentially-large probability distributions without paying an exponential cost, and at the same time endow the distributions with **structured semantics**



$$P(X_1, X_2, X_3, X_4, X_5, X_6, X_7, X_8)$$

$$P(X_{1:8}) = P(X_1)P(X_2)P(X_3 | X_1 X_2)P(X_4 | X_2)P(X_5 | X_2) \\ P(X_6 | X_3, X_4)P(X_7 | X_6)P(X_8 | X_5, X_6)$$

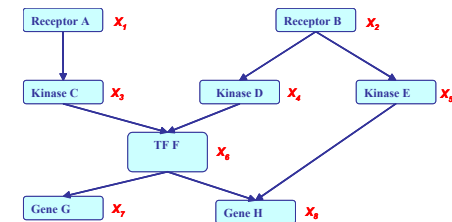
- A more formal description:
 - It refers to a family of distributions on a set of random variables that are compatible with all the probabilistic independence propositions encoded by a graph that connects these variables



Two types of GMs

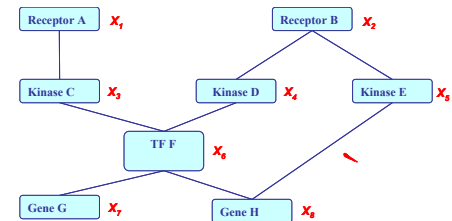
- Directed edges give causality relationships (Bayesian Network or Directed Graphical Model):

$$\begin{aligned}
 &P(X_1, X_2, X_3, X_4, X_5, X_6, X_7, X_8) \\
 &= P(X_1) P(X_2) P(X_3/X_1) P(X_4/X_2) P(X_5/X_2) \\
 &P(X_6/X_3, X_4) P(X_7/X_6) P(X_8/X_5, X_6)
 \end{aligned}$$

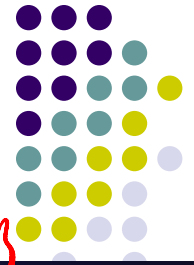


- Undirected edges simply give correlations between variables (Markov Random Field or Undirected Graphical model):

$$\begin{aligned}
 &P(X_1, X_2, X_3, X_4, X_5, X_6, X_7, X_8) \\
 &= \frac{1}{Z} \exp\{E(X_1)+E(X_2)+E(X_3, X_1)+E(X_4, X_2)+E(X_5, X_2) \\
 &+ E(X_6, X_3, X_4)+E(X_7, X_6)+E(X_8, X_5, X_6)\}
 \end{aligned}$$



Towards structural specification of probability distribution



- Separation properties in the graph imply independence properties about the associated variables
- For the graph to be useful, any conditional independence properties we can derive from the graph should hold for the probability distribution that the graph represents

$P(A)$
 $P(B)$

A B

0	6
1	1
0	1
1	6

2.36
7.1
↑

- **The Equivalence Theorem**

For a graph G ,

Let \mathcal{D}_1 denote the family of all distributions that satisfy $I(G)$,

Let \mathcal{D}_2 denote the family of all distributions that factor according to G ,

Then $\mathcal{D}_1 \equiv \mathcal{D}_2$.

(A) (B) \rightarrow A \perp B

$$I(G) \rightarrow \frac{P(A, B)}{= P(A) P(B)}$$

Bayesian Networks

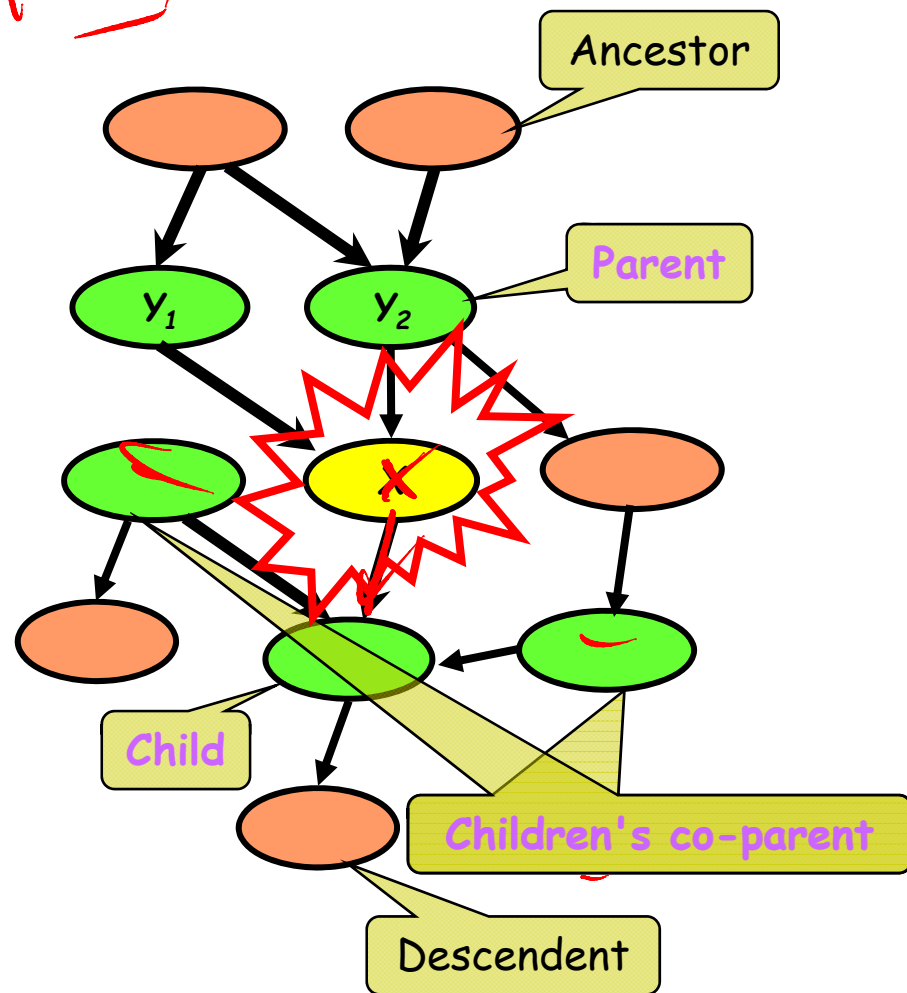


Structure: *DAG*

- Meaning: a node is **conditionally independent** of every other node in the network outside its Markov blanket
- Local conditional distributions (**CPD**) and the **DAG** completely determine the **joint** dist.
- Give **causality** relationships, and facilitate a **generative** process

$$P(X_i | X_{-i}) = P(X_i | X_{M\ddot{a}r})$$

X ⊥ ? | ? ?

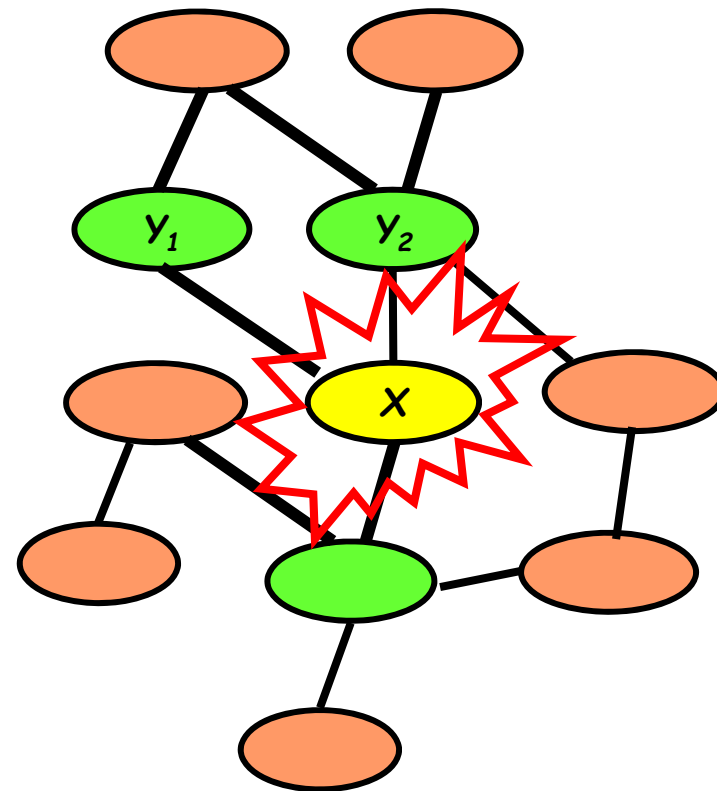




Markov Random Fields

Structure: *undirected graph*

- Meaning: a node is **conditionally independent** of every other node in the network given its **Directed neighbors**
- Local contingency functions (**potentials**) and the **cliques** in the graph completely determine the **joint** dist.
- Give **correlations** between variables, but no explicit way to generate samples

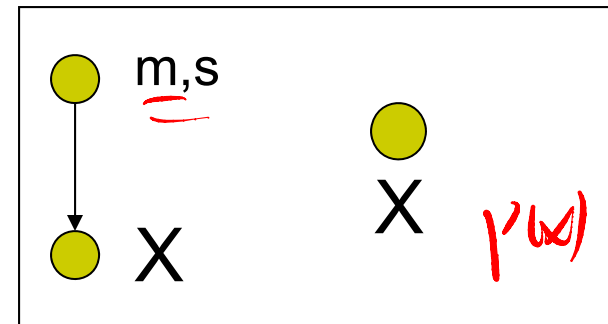




GMs are your old friends

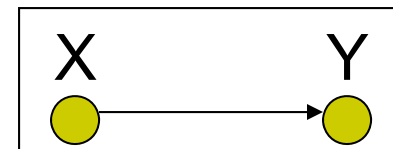
Density estimation

Parametric and nonparametric methods



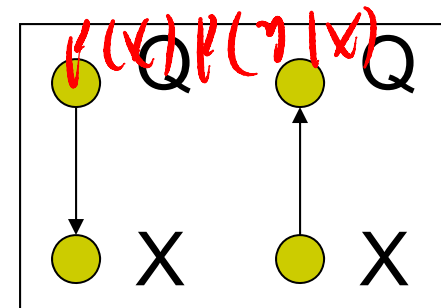
Regression

Linear, conditional mixture, nonparametric



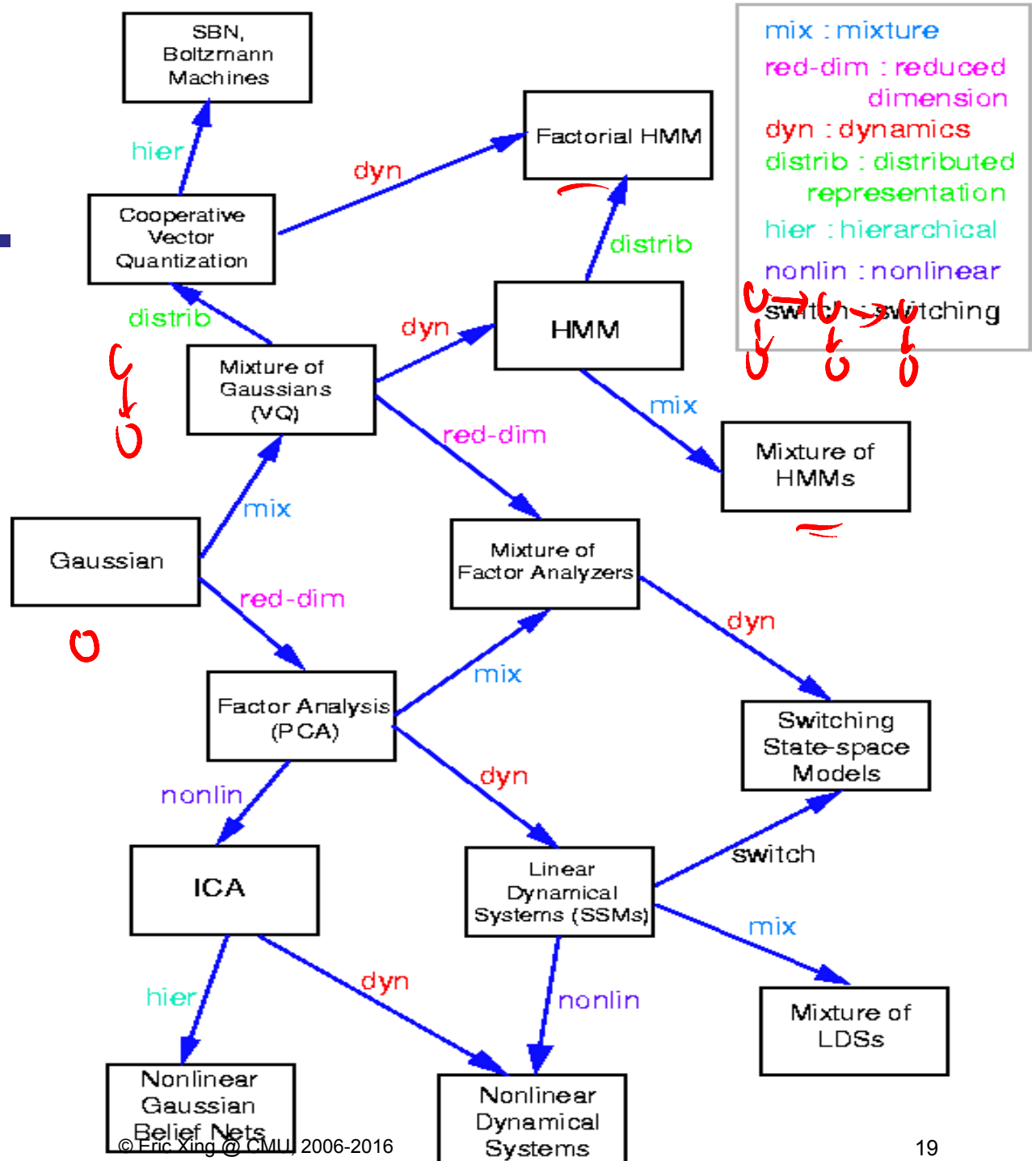
Classification

Generative and discriminative approach



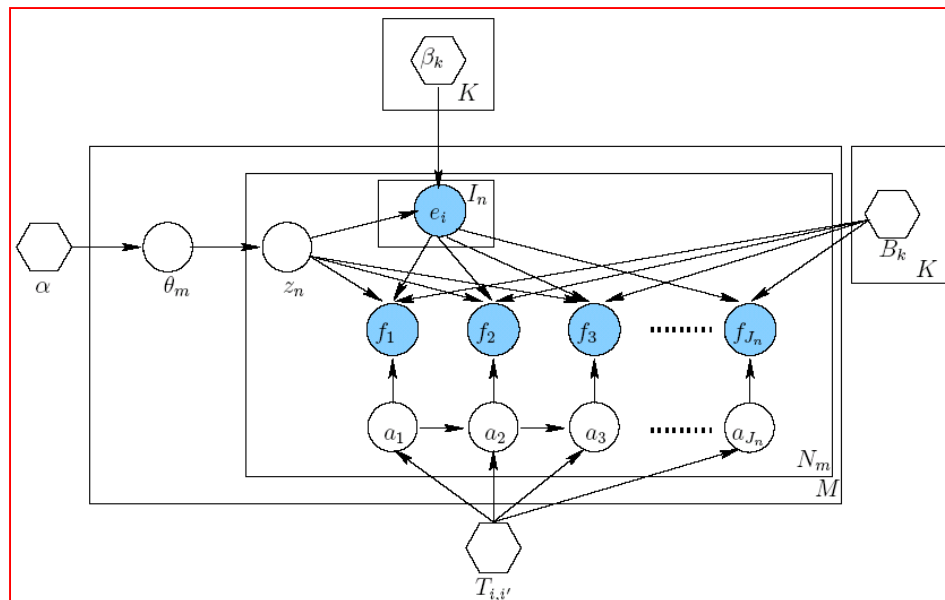
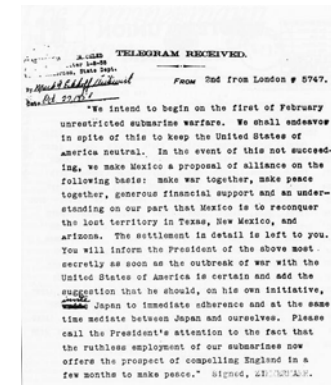
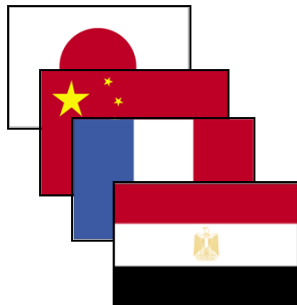
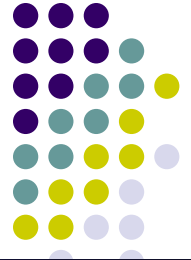
Clustering

An (incomplete) genealogy of graphical models



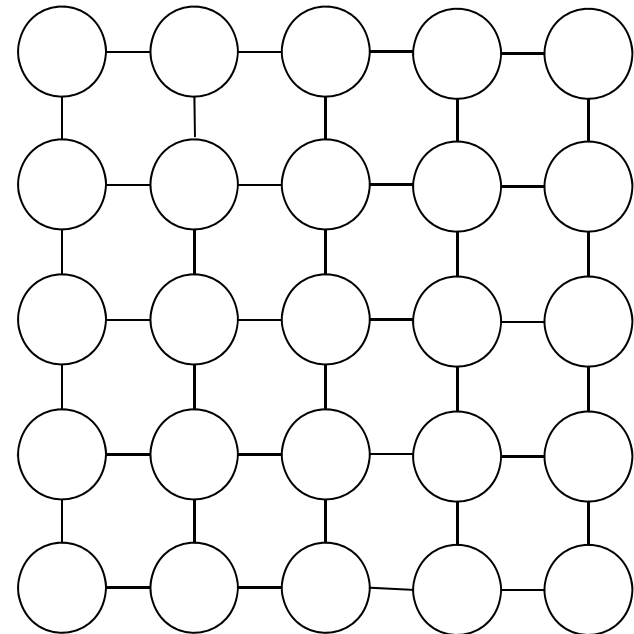
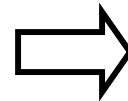
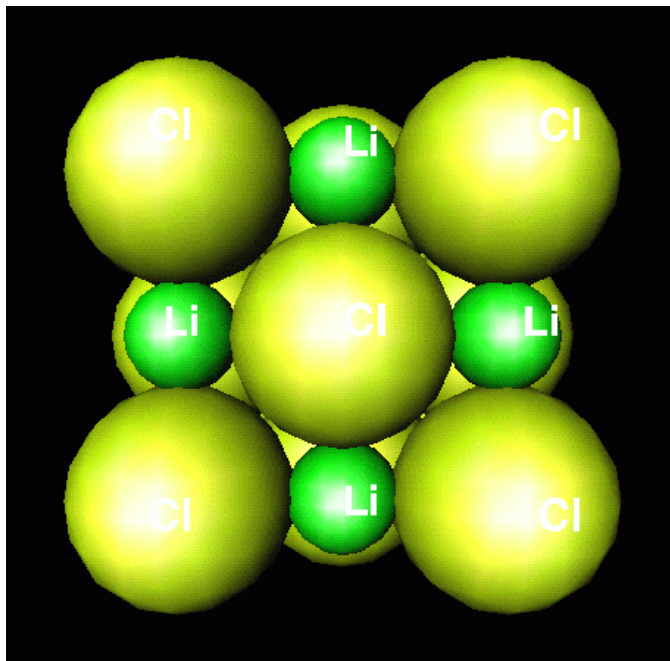
(Picture by Zoubin Ghahramani and Sam Roweis)

Fancier GMs: machine translation



The HM-BiTAM model
(B. Zhao and E.P Xing,
ACL 2006)

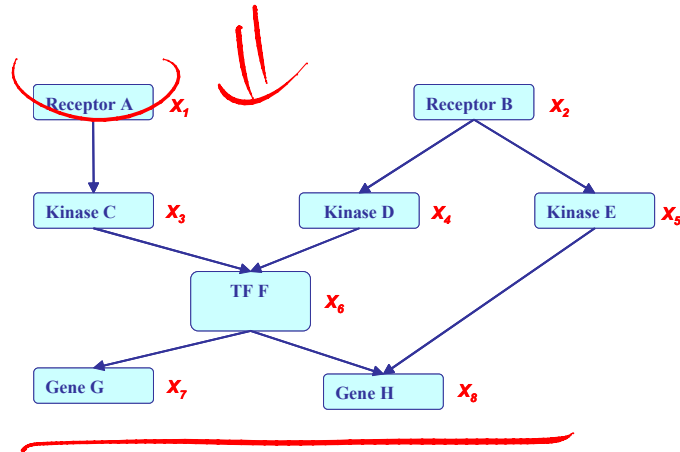
Fancier GMs: solid state physics



Ising/Potts model



Bayesian Network: Factorization Theorem



$$P(X_1, X_2, X_3, X_4, X_5, X_6, X_7, X_8) \\ = P(X_1) P(X_2) P(X_3 | X_1) P(X_4 | X_2) P(X_5 | X_2) \\ P(X_6 | X_3, X_4) P(X_7 | X_6) P(X_8 | X_5, X_6)$$

local cond. dist
 $P(X_i | X_{\pi_i})$

- Theorem:**

Given a DAG, The most general form of the probability distribution that is consistent with the (probabilistic independence properties encoded in the) graph factors according to “node given its parents”:

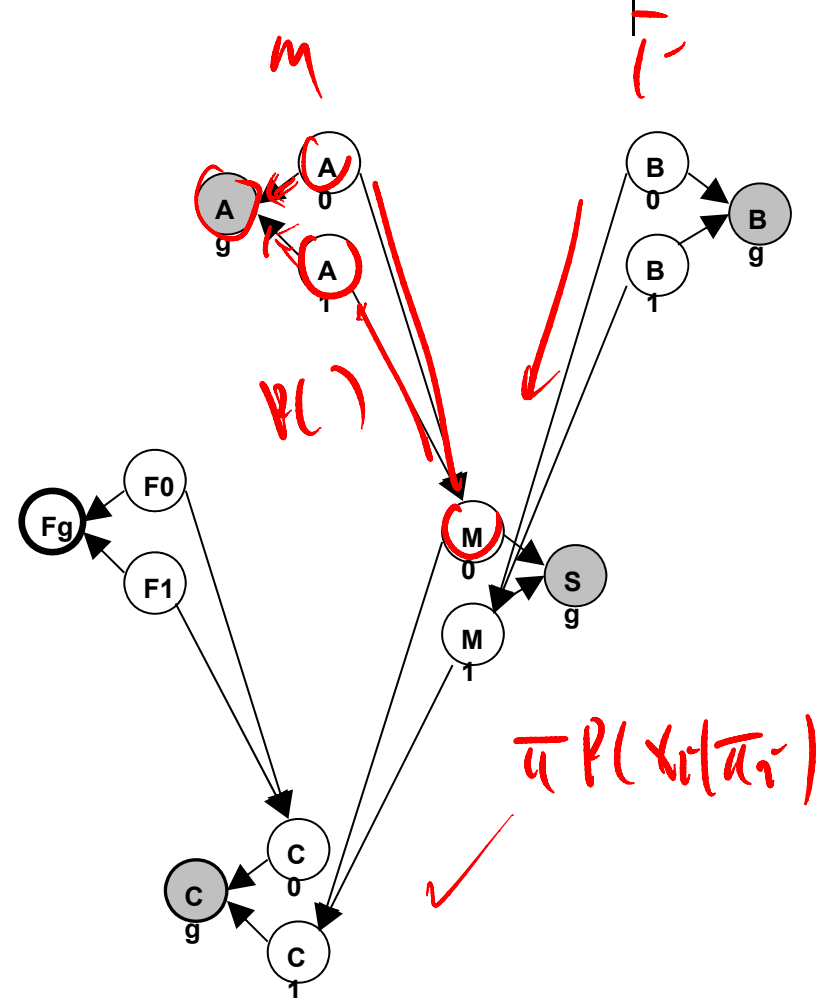
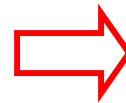
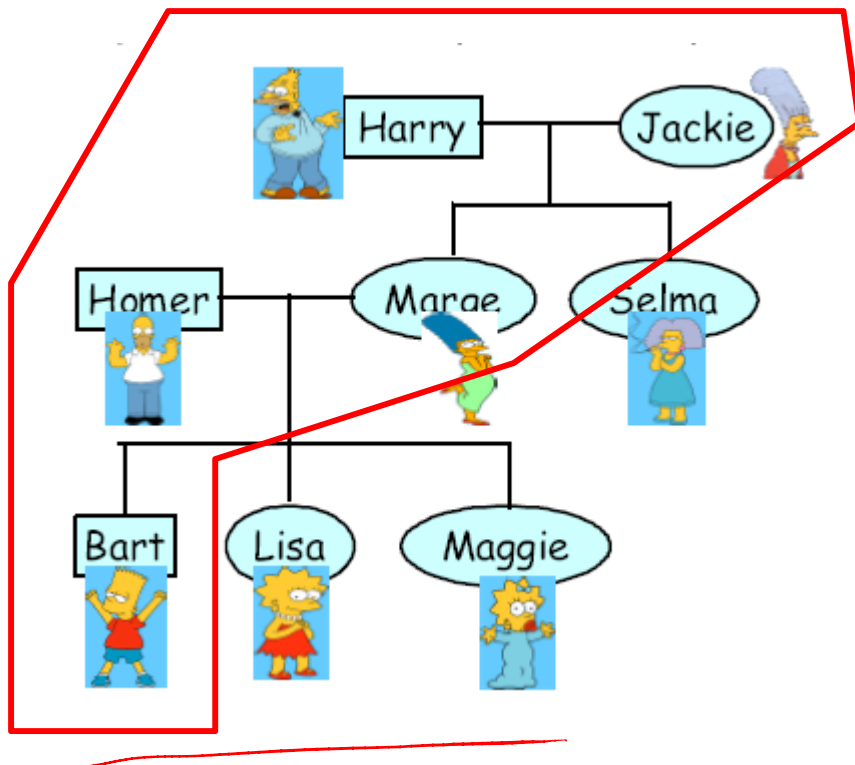
$$P(\mathbf{X}) = \prod_i P(X_i | \mathbf{X}_{\pi_i})$$

where \mathbf{X}_{π_i} is the set of parents of x_i . d is the number of nodes (variables) in the graph.



Example: a pedigree of people

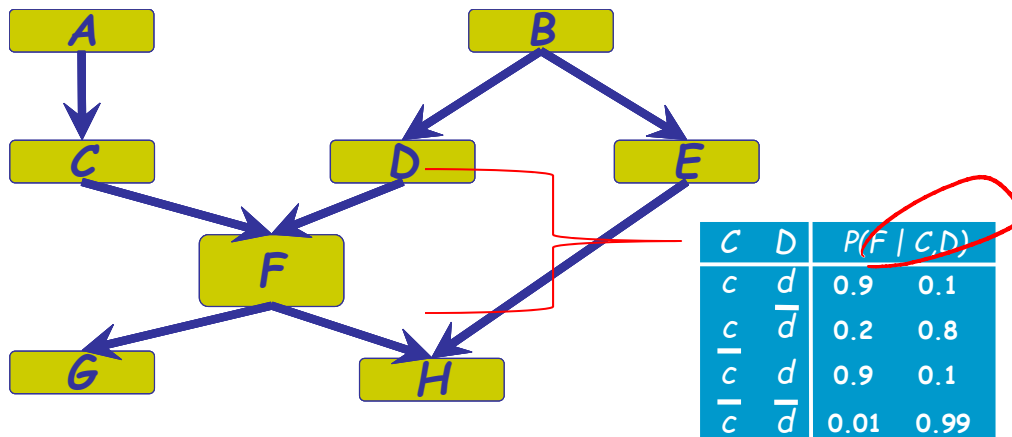
- Genetic Pedigree



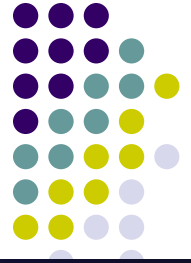


Specification of a BN

- There are two components to any GM:
 - the qualitative specification
 - the quantitative specification



Qualitative Specification

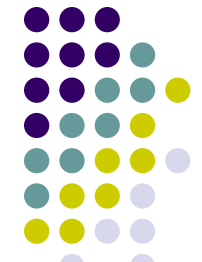


- Where does the qualitative specification come from?
 - Prior knowledge of causal relationships
 - Prior knowledge of modular relationships
 - Assessment from experts
 - Learning from data
 - We simply link a certain architecture (e.g. a layered graph)
 - ...

Local Structures & Independencies

$I(C|A)$

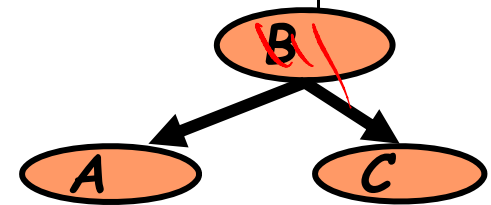
$A \perp C | B$



- Common parent

- Fixing B decouples A and C

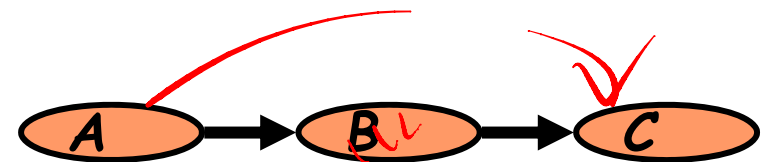
"given the level of gene B, the levels of A and C are independent"



- Cascade

- Knowing B decouples A and C

"given the level of gene B, the level gene A provides no extra prediction value for the level of gene C"



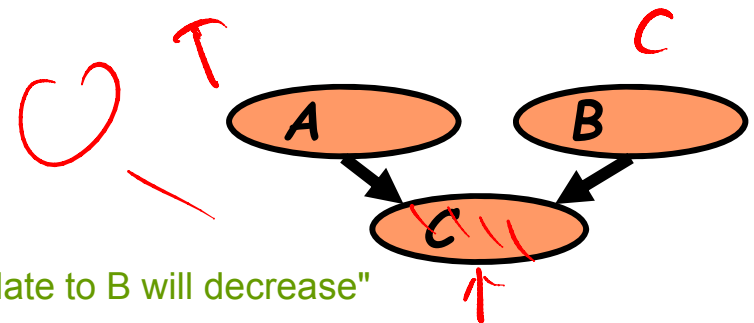
$A \perp C | B$

- V-structure

- Knowing C couples A and B

because A can "explain away" B w.r.t. C

"If A correlates to C, then chance for B to also correlate to B will decrease"



$A \perp B$
 $A \perp B | C$

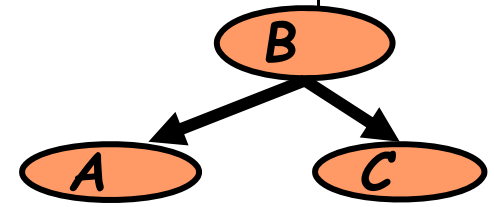
- The language is compact, the concepts are rich!



A simple justification

$$\begin{aligned} & \underline{P(A, B, C)} \\ &= \underline{P(B)} \underline{P(A|B)} \underline{P(C|B)} \end{aligned}$$

$$\begin{aligned} & \frac{P(A, B, C)}{P(B)} \\ &= P(A|B) P(C|B) \end{aligned}$$



$$\begin{aligned} & \underline{A \perp C | B} \\ & \Downarrow \\ & P(A, C | B) \\ &= \underline{P(A|B) P(C|B)} \end{aligned}$$

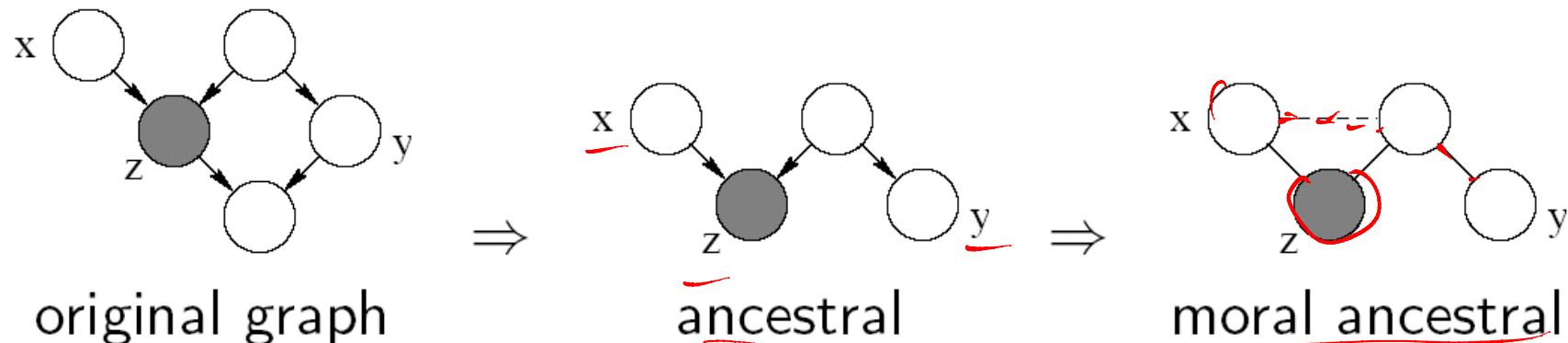


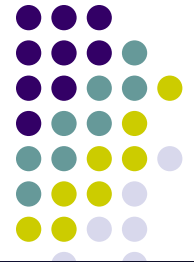
Graph separation criterion

- D-separation criterion for Bayesian networks (D for Directed edges):

Definition: variables x and y are *D-separated* (conditionally independent) given z if they are separated in the *moralized* ancestral graph

- Example:

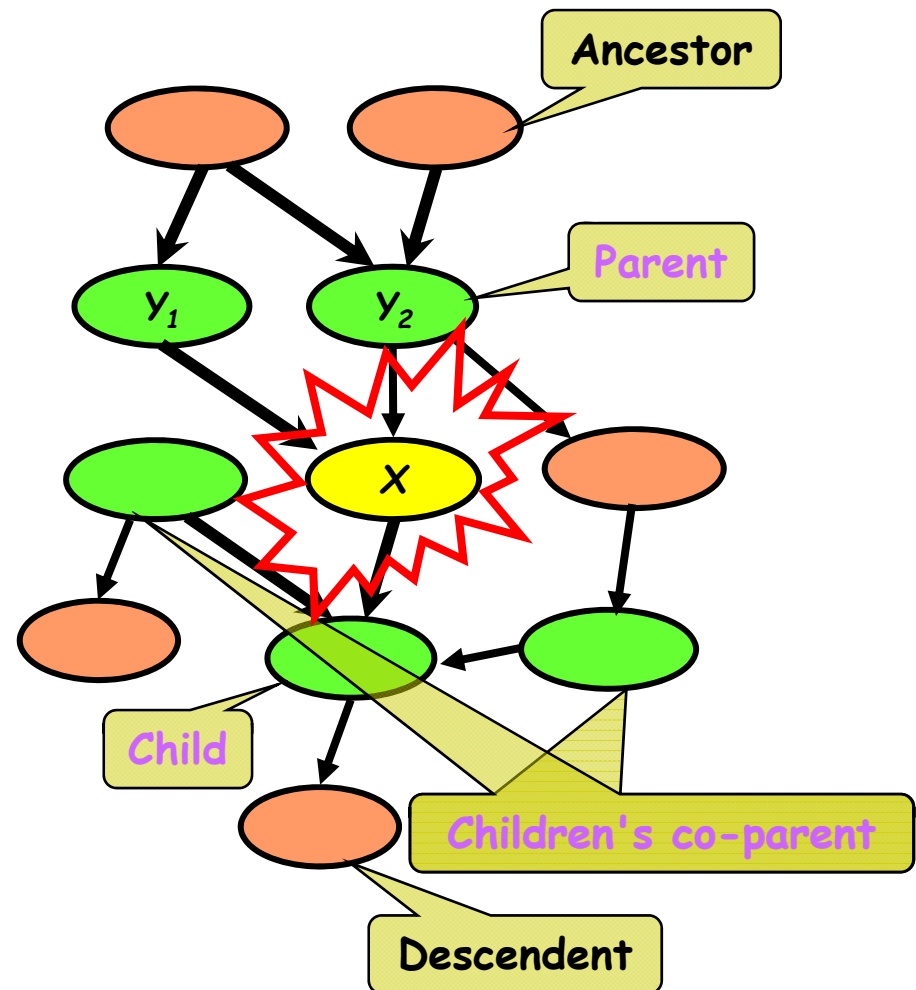




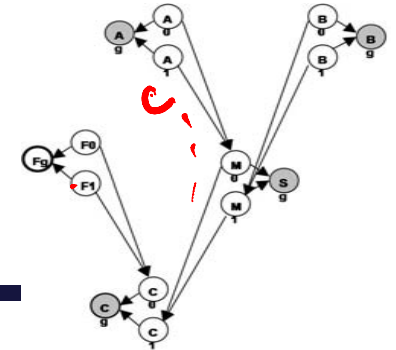
Local Markov properties of DAGs

Structure: DAG

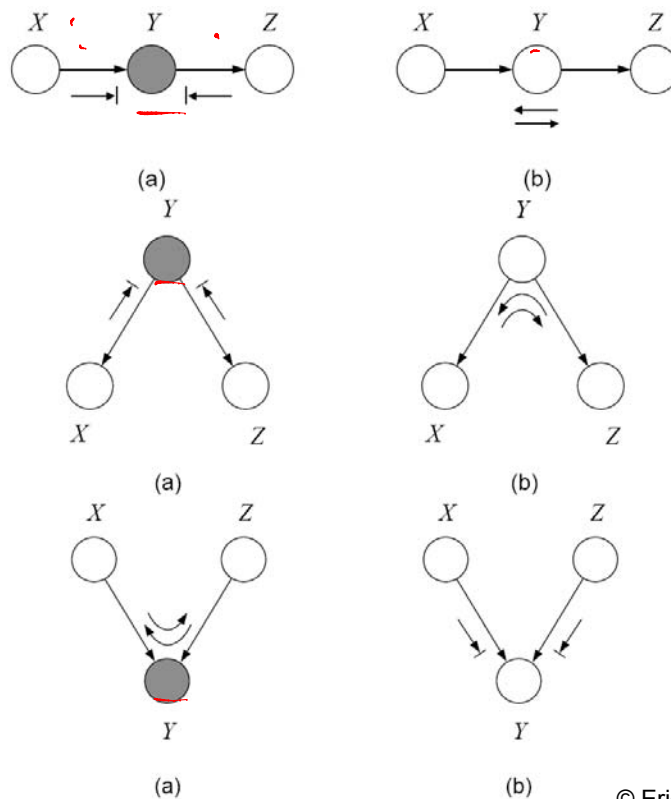
- Meaning: a node is **conditionally independent** of every other node in the network outside its **Markov blanket**
- Local conditional distributions (**CPD**) and the **DAG** completely determine the **joint** dist.
- Give **causality** relationships, and facilitate a **generative** process



Global Markov properties of DAGs



- X is **d-separated** (directed-separated) from Z given Y if we can't send a ball from any node in X to any node in Z using the "Bayes-ball" algorithm illustrated below (and plus some boundary conditions):



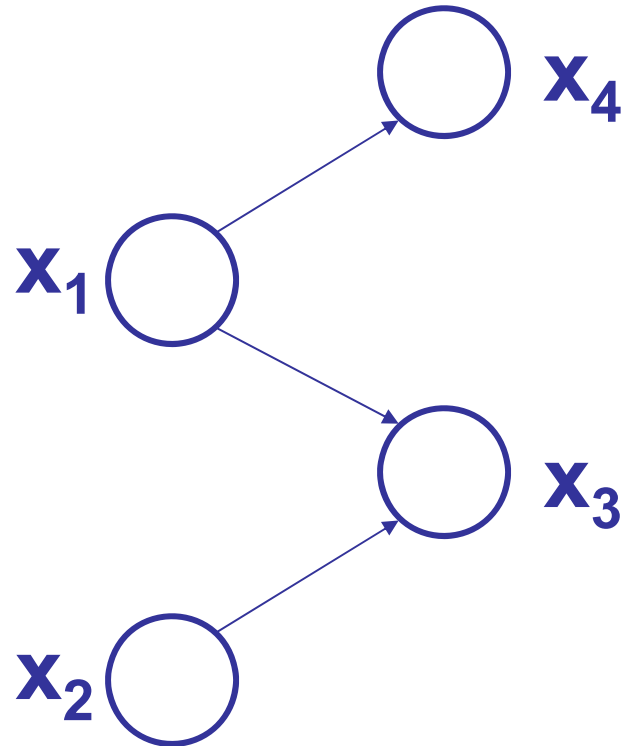
- Defn: $I(G)$ = all independence properties that correspond to d-separation:

$$I(G) = \{X \perp Z | Y : \text{dsep}_G(X; Z | Y)\}$$

- D-separation is sound and complete



Example:



- Complete the I(G) of this graph:

Essentially: A BN is a database of Pr. Independence statements among variables.

Towards quantitative specification of probability distribution



- Separation properties in the graph imply independence properties about the associated variables
- For the graph to be useful, any conditional independence properties we can derive from the graph should hold for the probability distribution that the graph represents

$P(I(G))$
 $P(I(G))$

$P(x_i \dots x_n)$

- **The Equivalence Theorem**

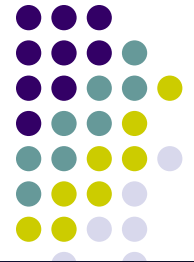
For a graph G ,

Let \mathcal{D}_1 denote the family of all distributions that satisfy $I(G)$,

Let \mathcal{D}_2 denote the family of all distributions that factor according to G ,

Then $\mathcal{D}_1 \equiv \mathcal{D}_2$.

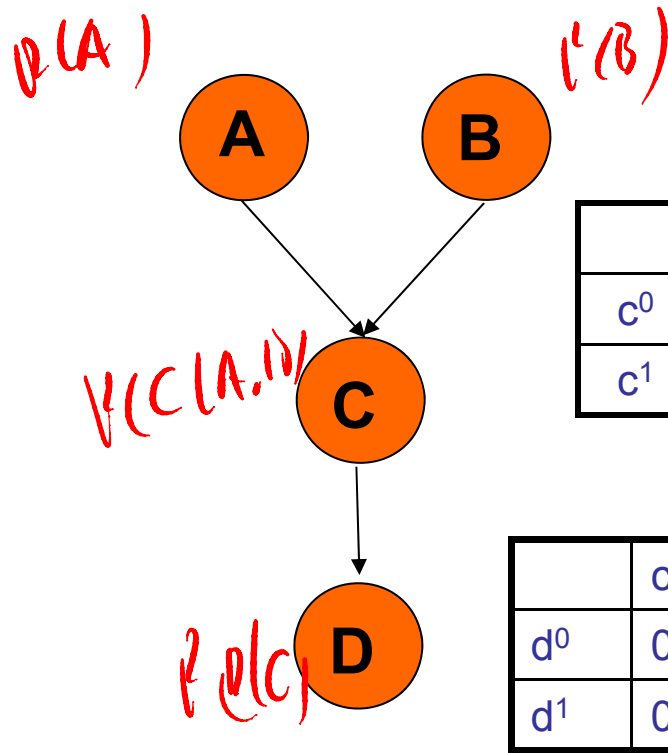
Conditional probability tables (CPTs)



a^0	0.75
a^1	0.25

b^0	0.33
b^1	0.67

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



	a^0b^0	a^0b^1	a^1b^0	a^1b^1
c^0	0.45	1	0.9	0.7
c^1	0.55	0	0.1	0.3

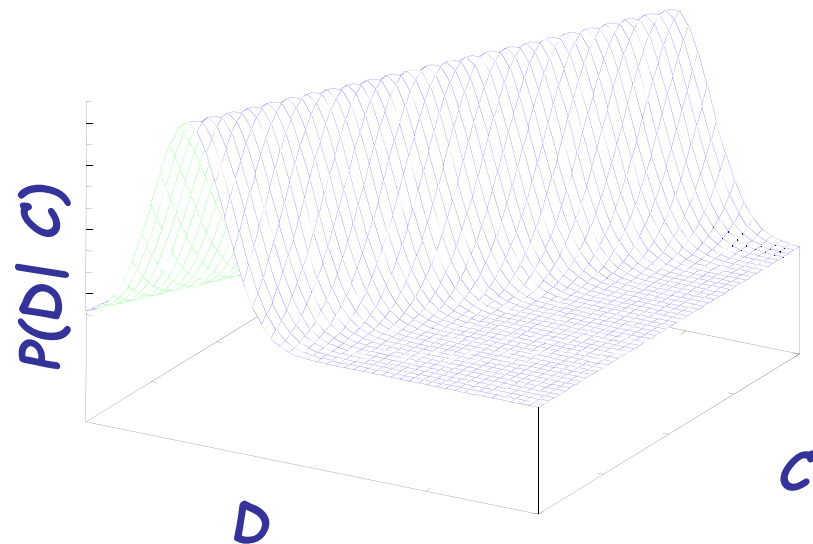
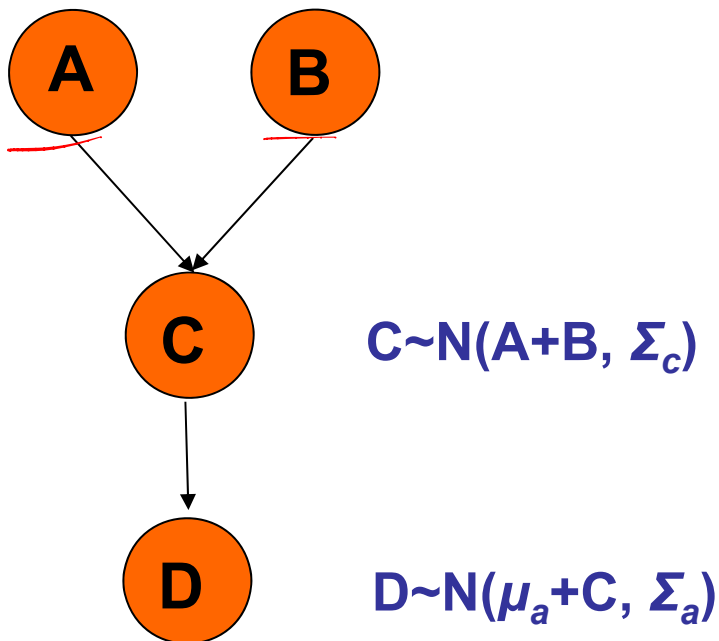
	c^0	c^1
d^0	0.3	0.5
d^1	0.7	0.5

Conditional probability density func. (CPDs)

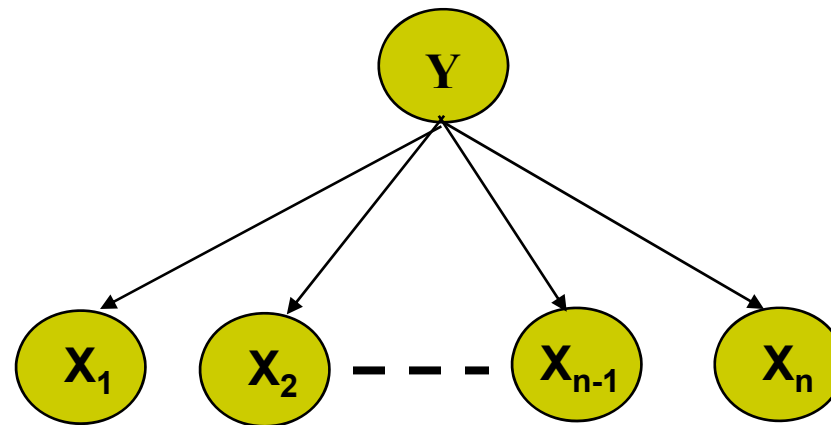
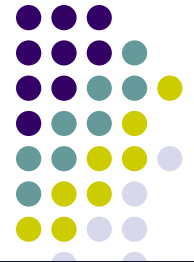


$$A \sim N(\mu_a, \Sigma_a) \quad B \sim N(\mu_b, \Sigma_b)$$

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



Conditional Independencies



$$P(Y, X)$$

Label

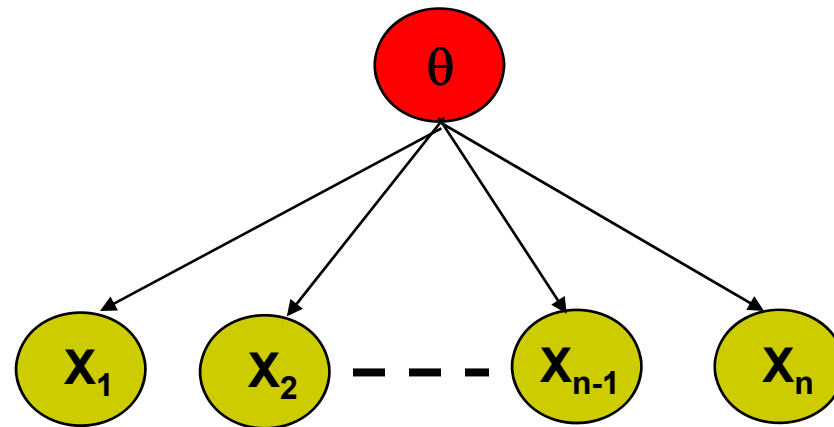
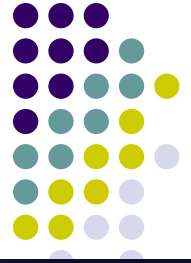
$$= \prod_i P(X_i | y) P(y)$$

Features

What is this model

1. When Y is observed?
2. When Y is unobserved?

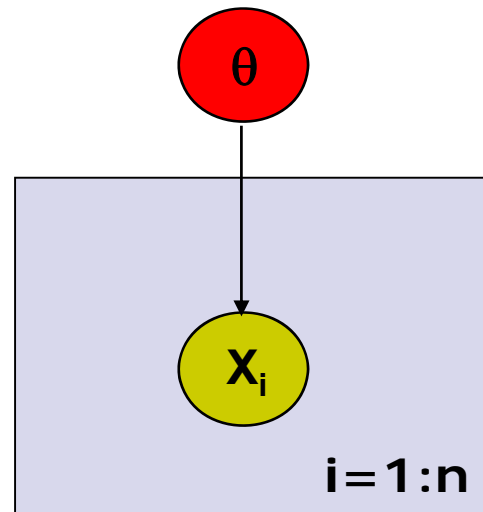
Conditionally Independent Observations



Model parameters

Data = $\{y_1, \dots, y_n\}$

“Plate” Notation



Model parameters

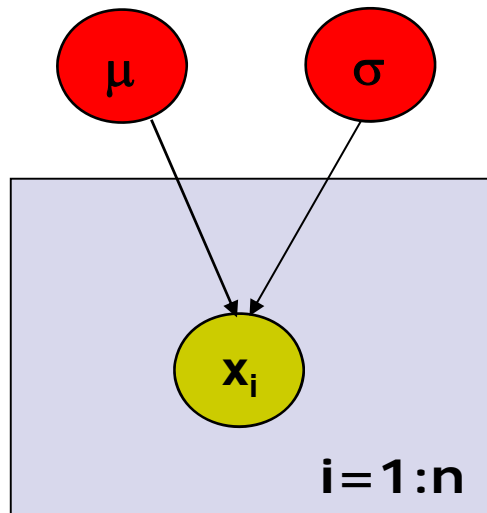
Data = $\{x_1, \dots, x_n\}$

Plate = rectangle in graphical model

**variables within a plate are replicated
in a conditionally independent manner**



Example: Gaussian Model



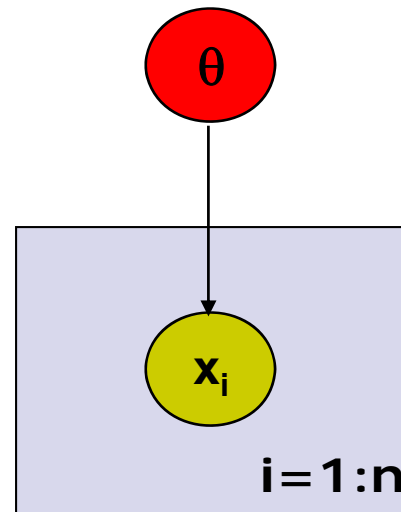
Generative model:

$$\begin{aligned} p(x_1, \dots, x_n \mid \mu, \sigma) &= \prod p(x_i \mid \mu, \sigma) \\ &= p(\text{data} \mid \text{parameters}) \\ &= p(D \mid \theta) \end{aligned}$$

where $\theta = \{\mu, \sigma\}$

- Likelihood $= p(\text{data} \mid \text{parameters})$
 $= p(D \mid \theta)$
 $= L(\theta)$
- Likelihood tells us how likely the observed data are conditioned on a particular setting of the parameters
 - Often easier to work with $\log L(\theta)$

Bayesian models



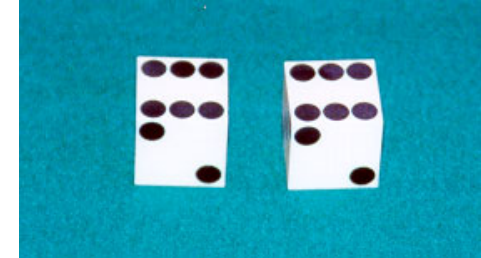
A Generative Scheme for model design



Suppose you were told about the following story before heading to Vegas...



The Dishonest Casino !!!



A casino has two dice:

- **Fair die**

$$P(1) = P(2) = P(3) = P(5) = P(6) = 1/6$$

- **Loaded die**

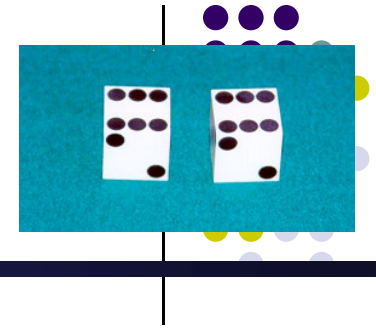
$$P(1) = P(2) = P(3) = P(5) = 1/10$$

$$P(6) = 1/2$$

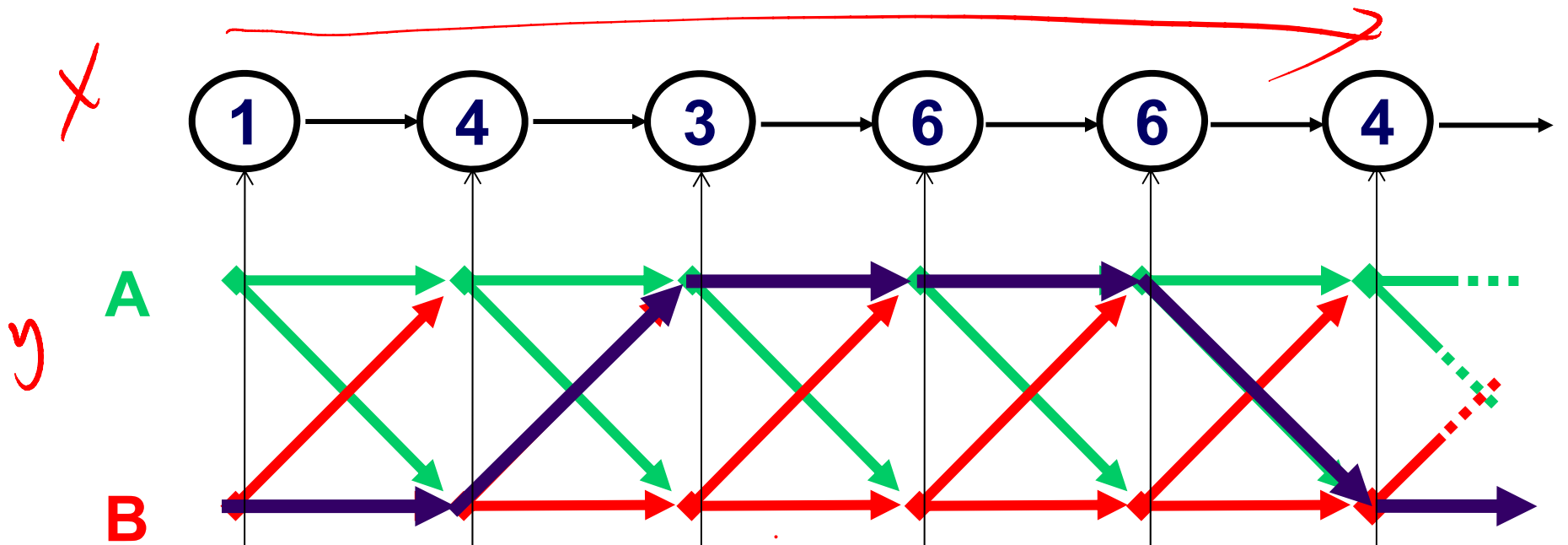
Casino player switches back-&-forth between fair and loaded die once every 20 turns



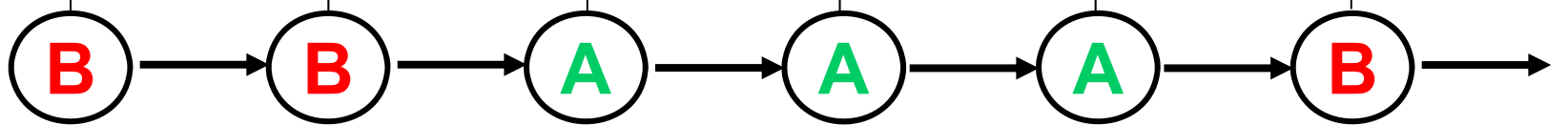
An HMM is a Stochastic Generative Model



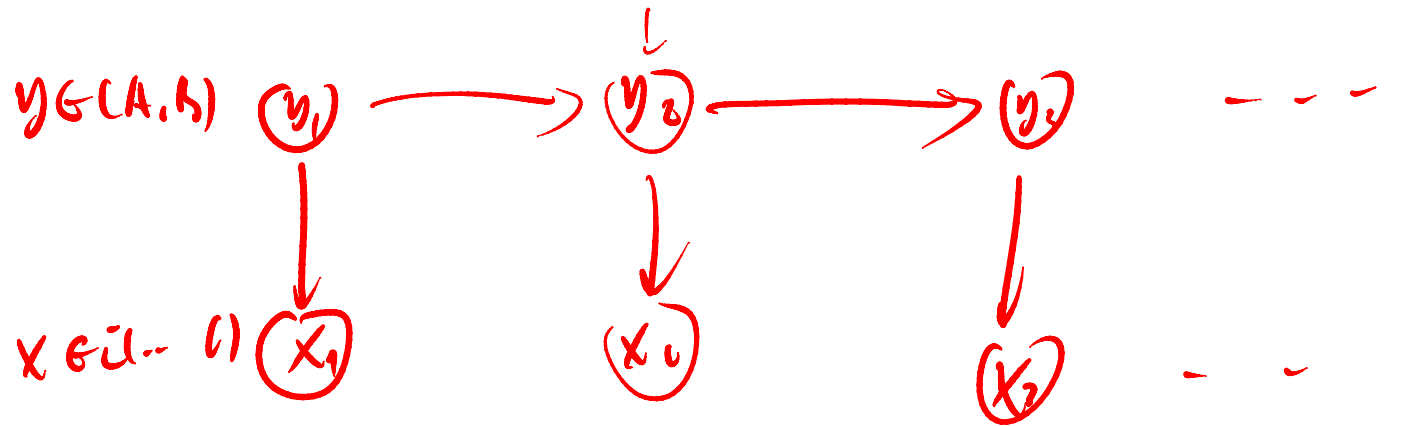
- Observed sequence:



- Hidden sequence (a parse or segmentation):



A Generative Scheme for model design



$$p(x_0, x_T, y_1, \dots, y_T)$$



Definition (of HMM)

- **Observation space**

Alphabetic set:

$$\mathcal{C} = \{c_1, c_2, \dots, c_K\}$$

Euclidean space:

$$\mathbb{R}^d$$

- **Index set of hidden states**

$$\mathcal{I} = \{1, 2, \dots, M\}$$

- **Transition probabilities** between any two states

$$p(y_t^j = 1 | y_{t-1}^i = 1) = a_{i,j},$$

or $p(y_t | y_{t-1}^i = 1) \sim \text{Multinomial}(a_{i,1}, a_{i,2}, \dots, a_{i,M}), \forall i \in \mathcal{I}.$

- **Start probabilities**

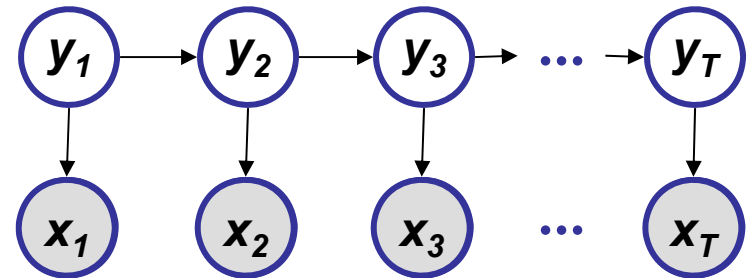
$$p(y_1) \sim \text{Multinomial}(\pi_1, \pi_2, \dots, \pi_M).$$

- **Emission probabilities** associated with each state

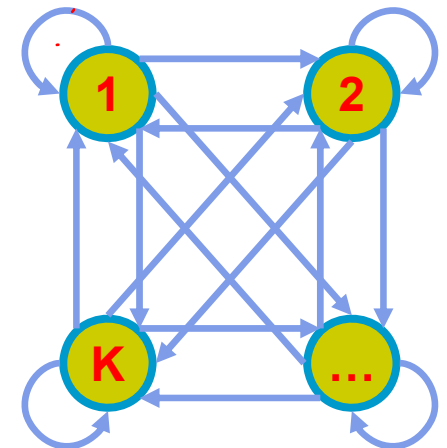
$$p(x_t | y_t^i = 1) \sim \text{Multinomial}(b_{i,1}, b_{i,2}, \dots, b_{i,K}), \forall i \in \mathcal{I}.$$

or in general:

$$p(x_t | y_t^i = 1) \sim f(\cdot | \theta_i), \forall i \in \mathcal{I}.$$



Graphical model



State automata



Why graphical models

- A language for communication
- A language for computation
- A language for development

- **Origins:**

- Wright 1920's
- Independently developed by Spiegelhalter and Lauritzen in statistics and Pearl in computer science in the late 1980's



Why graphical models

- **Probability theory** provides the **glue** whereby the parts are combined, ensuring that the system as a whole is consistent, and providing ways to interface models to data.
- The **graph theoretic** side of graphical models provides both an intuitively appealing interface by which humans can model highly-interacting sets of variables as well as a data structure that lends itself naturally to the design of efficient general-purpose algorithms.
- **Many of the classical multivariate probabilistic systems** studied in fields such as statistics, systems engineering, information theory, pattern recognition and statistical mechanics **are special cases of the general graphical model formalism**
- The graphical model framework provides a way to view all of these systems as instances of a **common underlying formalism**.

Summary



- Represent dependency structure with a directed acyclic graph
 - Node \leftrightarrow random variable
 - Edges encode dependencies
 - Absence of edge \rightarrow conditional independence
 - Plate representation
 - A GM is a database of prob. Independence statement on variables
- The factorization theorem of the joint probability
 - Local specification \rightarrow globally consistent distribution
 - Local representation for exponentially complex state-space
 - It is a smart way to **write/specify/compose/design** exponentially-large probability distributions without paying an exponential cost, and at the same time endow the distributions with structured semantics
- Support efficient inference and learning

