

Introduction to Machine Learning

Brown University CSCI 1950-F, Spring 2011
Prof. Erik Sudderth

Lecture 24: Gibbs Samplers,
Directed Graphical Models

Many figures courtesy Kevin Murphy's textbook,
Machine Learning: A Probabilistic Perspective

Announcements

- May 10 at noon: Homework 9 due
- May 10 at 1pm: Review session for final exam
- May 12 at 1pm: Graduate project presentations
(rumor that free food will be served)
- May 19 at 2pm: Final exam

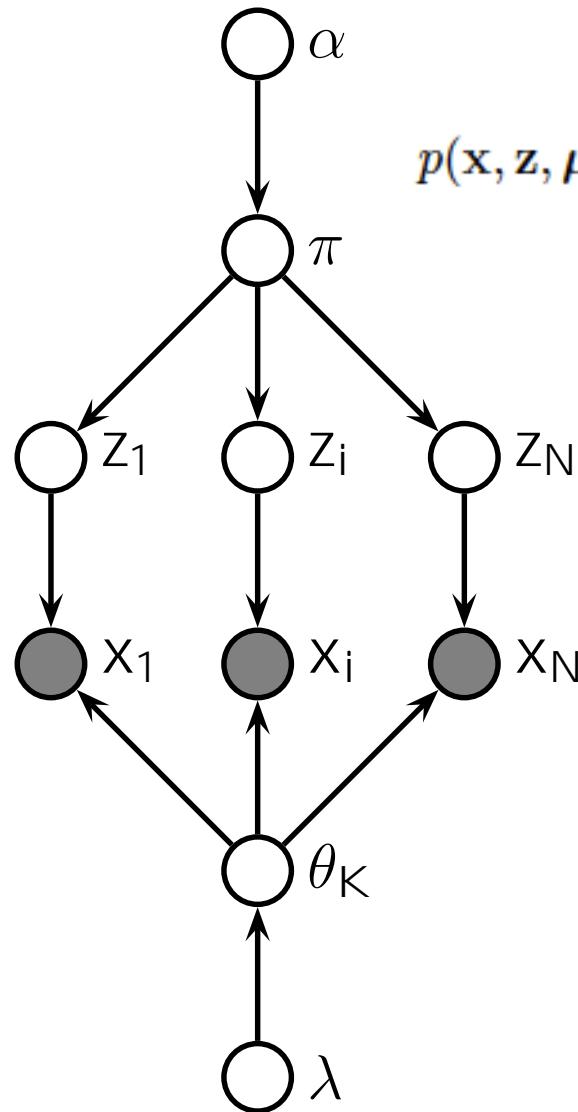
Monte Carlo Estimators

$$\begin{aligned}\mathbb{E}_p[f(x)] &= \int_{\mathcal{X}} f(x)p(x) dx & \{x^{(\ell)}\}_{\ell=1}^L & \text{independent samples} \\ &\approx \frac{1}{L} \sum_{\ell=1}^L f(x^{(\ell)}) = \mathbb{E}_{\tilde{p}}[f(x)] & \tilde{p}(x) &= \frac{1}{L} \sum_{\ell=1}^L \delta(x, x^{(\ell)})\end{aligned}$$

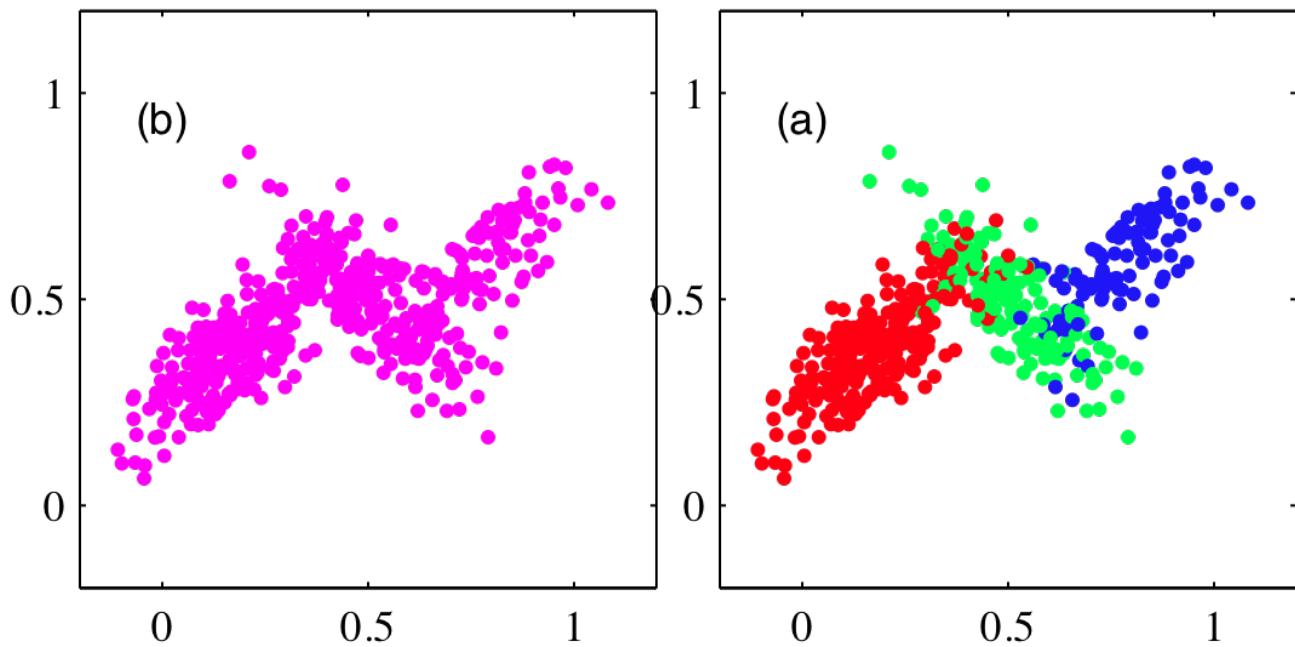
Good properties if L sufficiently large:

- Unbiased for any sample size
- Variance inversely proportional to sample size
(and independent of dimension of space)
- Weak law of large numbers
- Strong law of large numbers
- **Problem:** Drawing samples from complex distributions...

Gibbs Sampler for Gaussian Mixtures



$$\begin{aligned} p(\mathbf{x}, \mathbf{z}, \boldsymbol{\mu}, \boldsymbol{\Sigma}, \boldsymbol{\pi}) &= p(\mathbf{x}|\mathbf{z}, \boldsymbol{\mu}, \boldsymbol{\Sigma}) p(\mathbf{z}|\boldsymbol{\pi}) p(\boldsymbol{\pi}) \prod_{k=1}^K p(\boldsymbol{\mu}_k) p(\boldsymbol{\Sigma}_k) \\ &= \left(\prod_{i=1}^N \prod_{k=1}^K (\pi_k \mathcal{N}(x_i | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k))^{\mathbb{I}(z_i=k)} \right) \\ &\quad \times \text{Dir}(\boldsymbol{\pi} | \boldsymbol{\alpha}) \times \prod_{k=1}^K \mathcal{N}(\boldsymbol{\mu}_k | \mathbf{m}_0, \mathbf{V}_0) \text{IW}(\boldsymbol{\Sigma}_k | \mathbf{S}_0, \nu_0) \end{aligned}$$



Mixture Sampler Pseudocode

Given mixture weights $\pi^{(t-1)}$ and cluster parameters $\{\theta_k^{(t-1)}\}_{k=1}^K$ from the previous iteration, sample a new set of mixture parameters as follows:

1. Independently assign each of the N data points x_i to one of the K clusters by sampling the indicator variables $z = \{z_i\}_{i=1}^N$ from the following multinomial distributions:

$$z_i^{(t)} \sim \frac{1}{Z_i} \sum_{k=1}^K \pi_k^{(t-1)} f(x_i | \theta_k^{(t-1)}) \delta(z_i, k) \quad Z_i = \sum_{k=1}^K \pi_k^{(t-1)} f(x_i | \theta_k^{(t-1)})$$

2. Sample new mixture weights according to the following Dirichlet distribution:

$$\pi^{(t)} \sim \text{Dir}(N_1 + \alpha/K, \dots, N_K + \alpha/K) \quad N_k = \sum_{i=1}^N \delta(z_i^{(t)}, k)$$

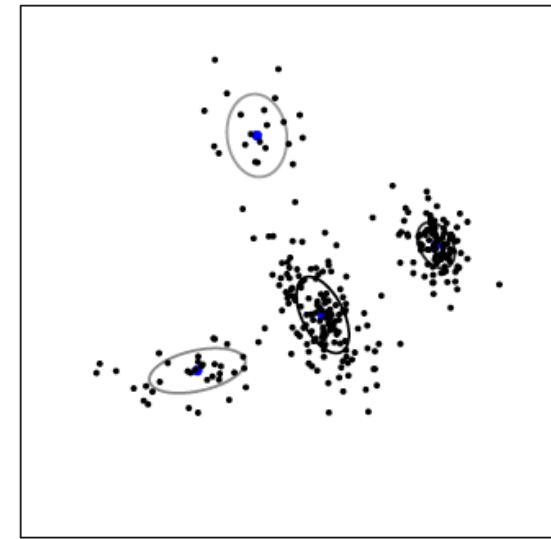
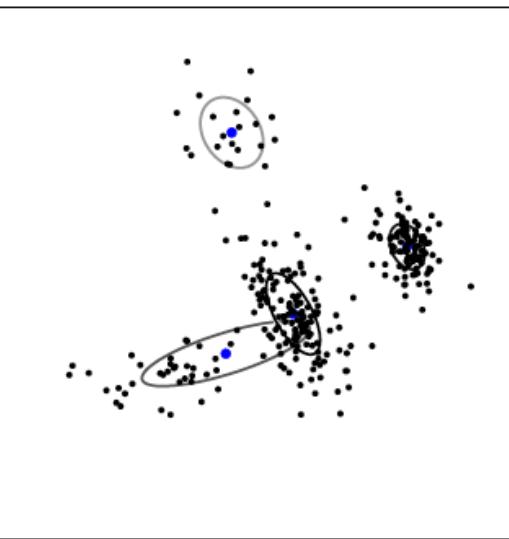
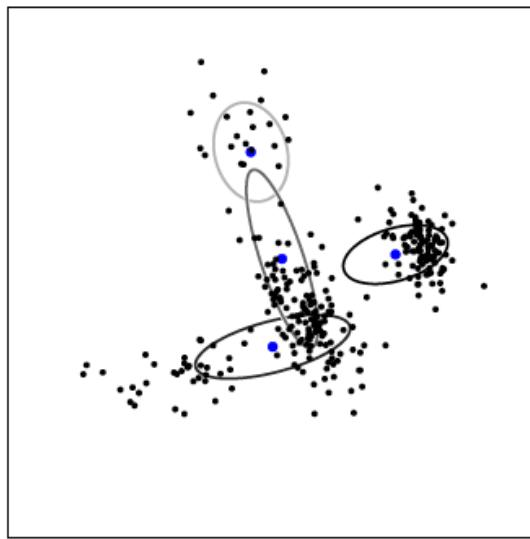
3. For each of the K clusters, independently sample new parameters from the conditional distribution implied by those observations currently assigned to that cluster:

$$\theta_k^{(t)} \sim p(\theta_k | \{x_i | z_i^{(t)} = k\}, \lambda)$$

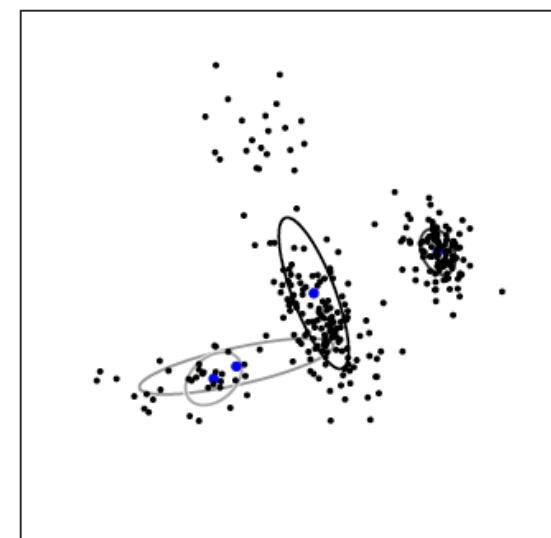
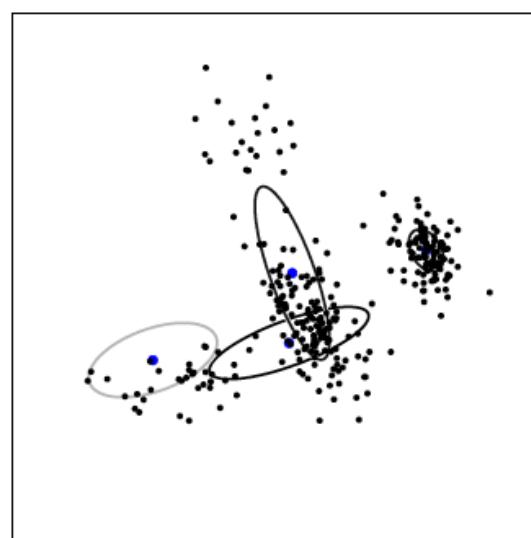
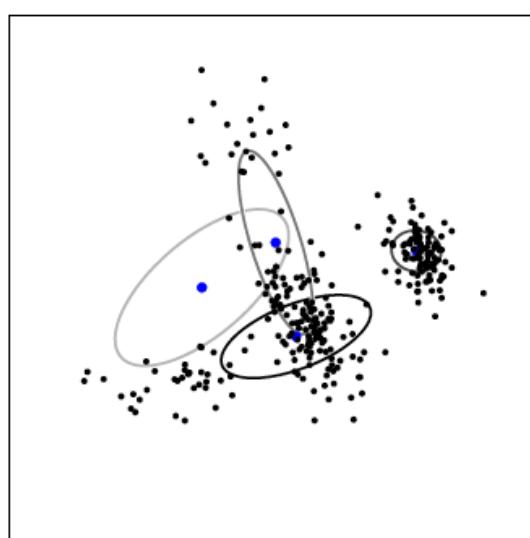
When λ defines a conjugate prior, this posterior distribution is given by Prop. 2.1.4.

Snapshots of Mixture Gibbs Sampler

Initialization A



Initialization B

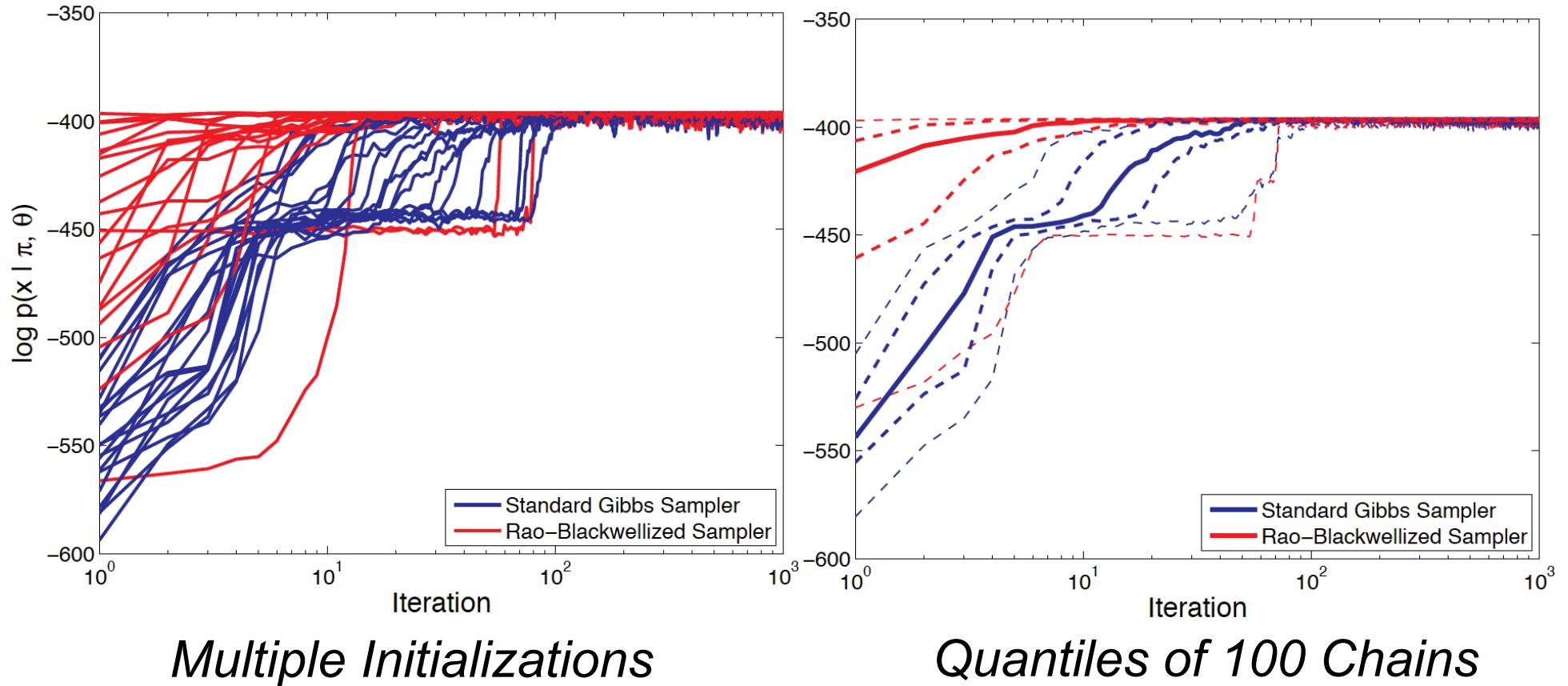


2 Iterations

10 Iterations

50 Iterations

Gibbs: Representation and Mixing

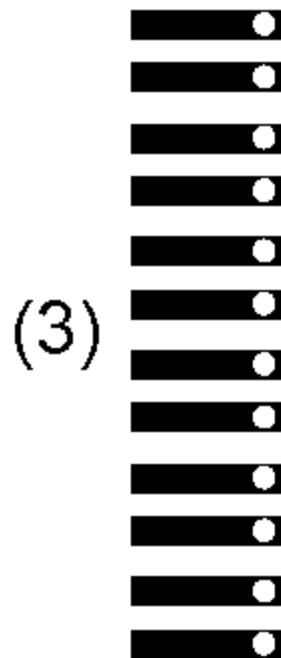


Standard Gibbs: Alternatively sample assignments, parameters
Collapsed Gibbs: Marginalize parameters, sample assignments

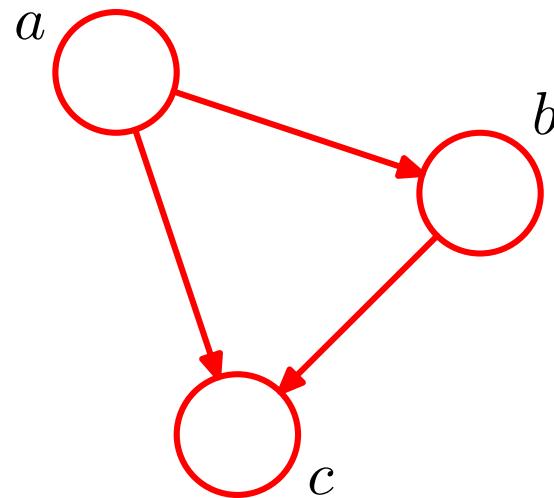
MCMC & Computational Resources



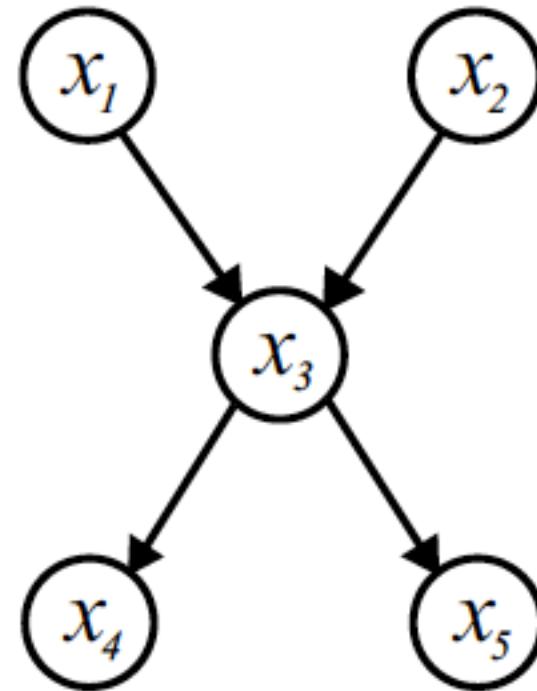
*Best practical option:
A few (> 1) initializations
for as many iterations as possible*



Directed Graphical Models



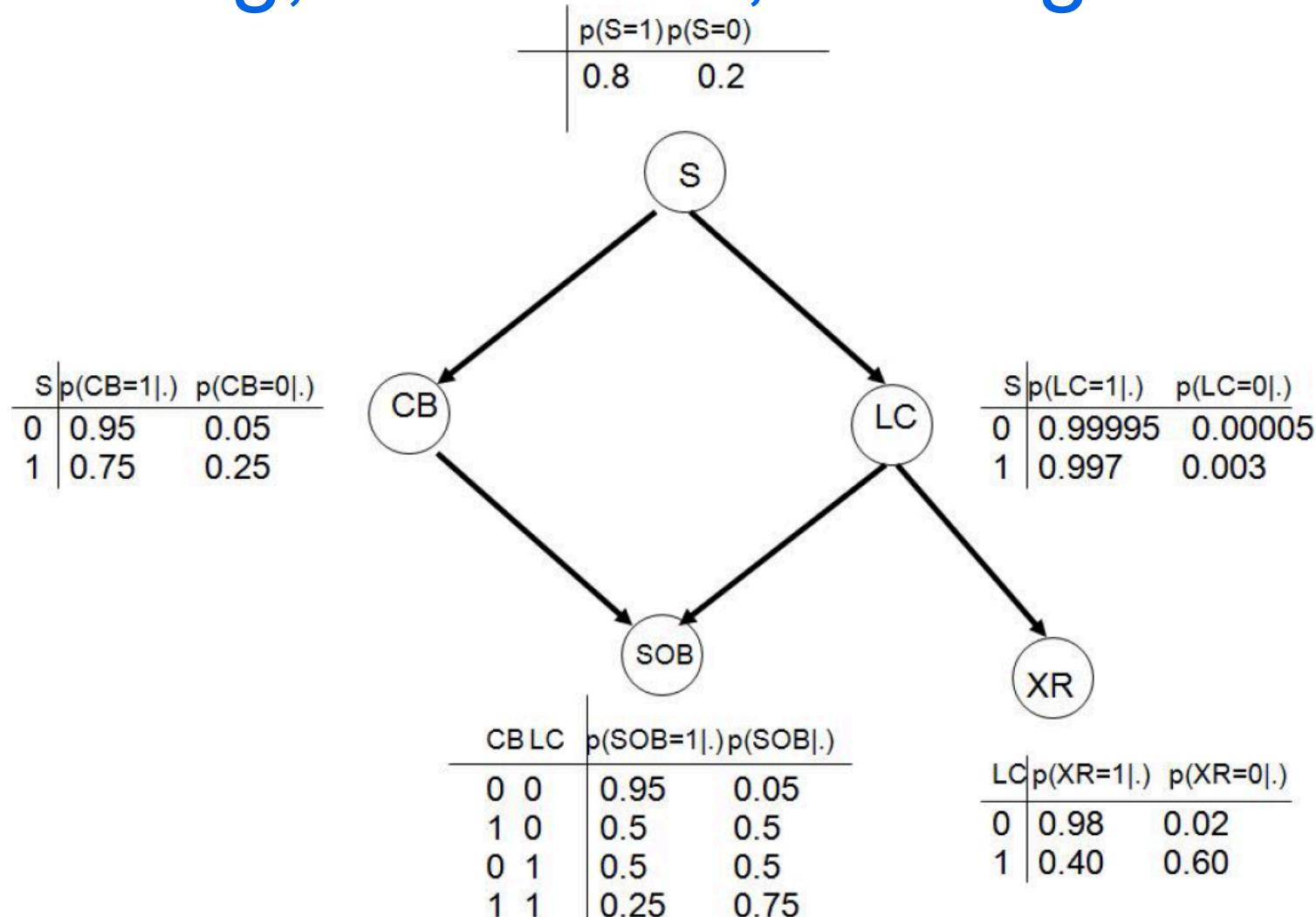
$$p(\mathbf{x}_{1:D}|G) = \prod_{s=1}^D p(x_s|\mathbf{x}_{\text{pa}(s)})$$



$$p(\mathbf{x}) = p(x_1) p(x_2) p(x_3 \mid x_1, x_2) p(x_4 \mid x_3) p(x_5 \mid x_3)$$

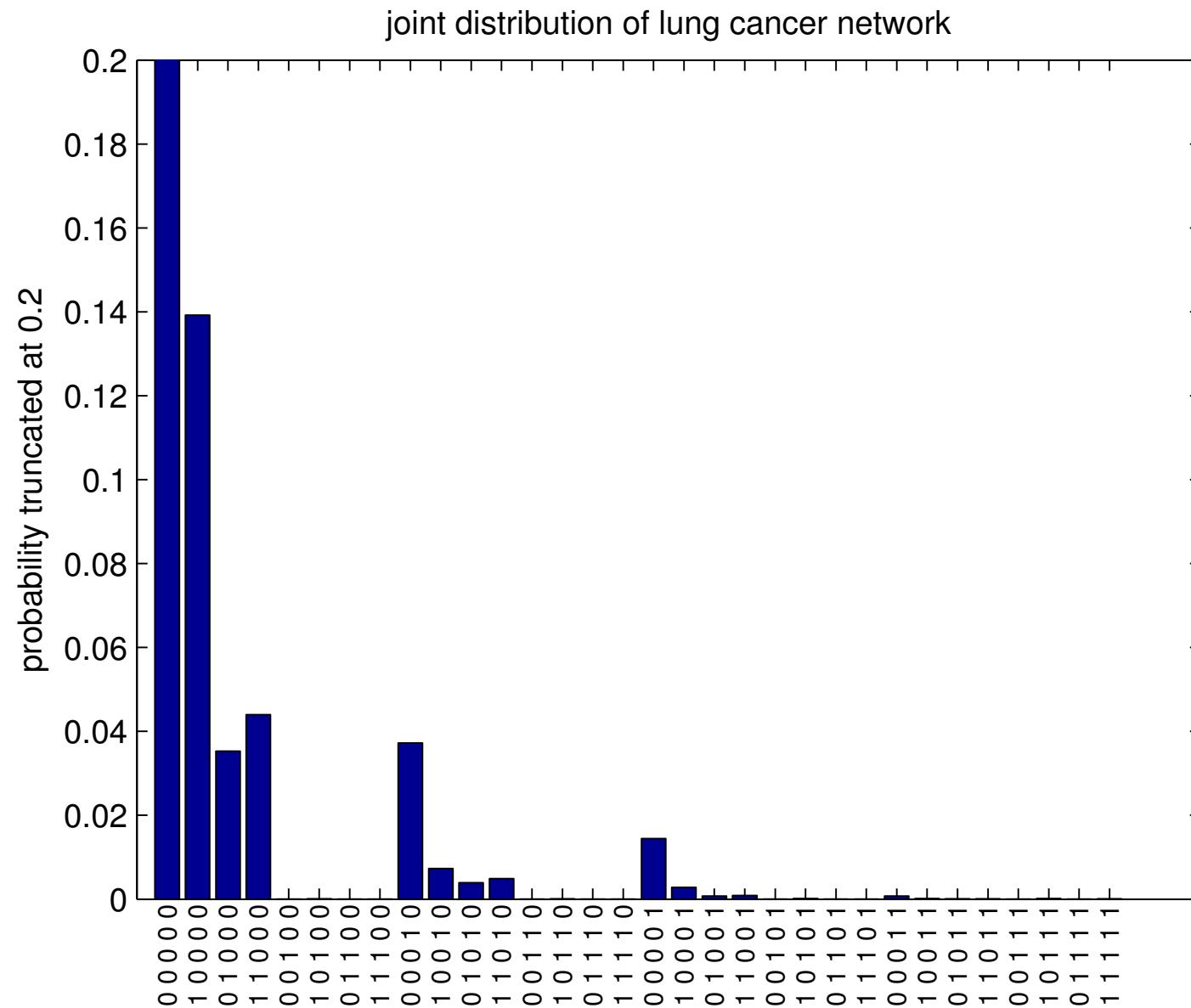
- Allows easy visualization of model similarities and differences
- Suggests ways of generalizing and combining models
- General-purpose learning and inference algorithms

Smoking, Bronchitis, & Lung Cancer

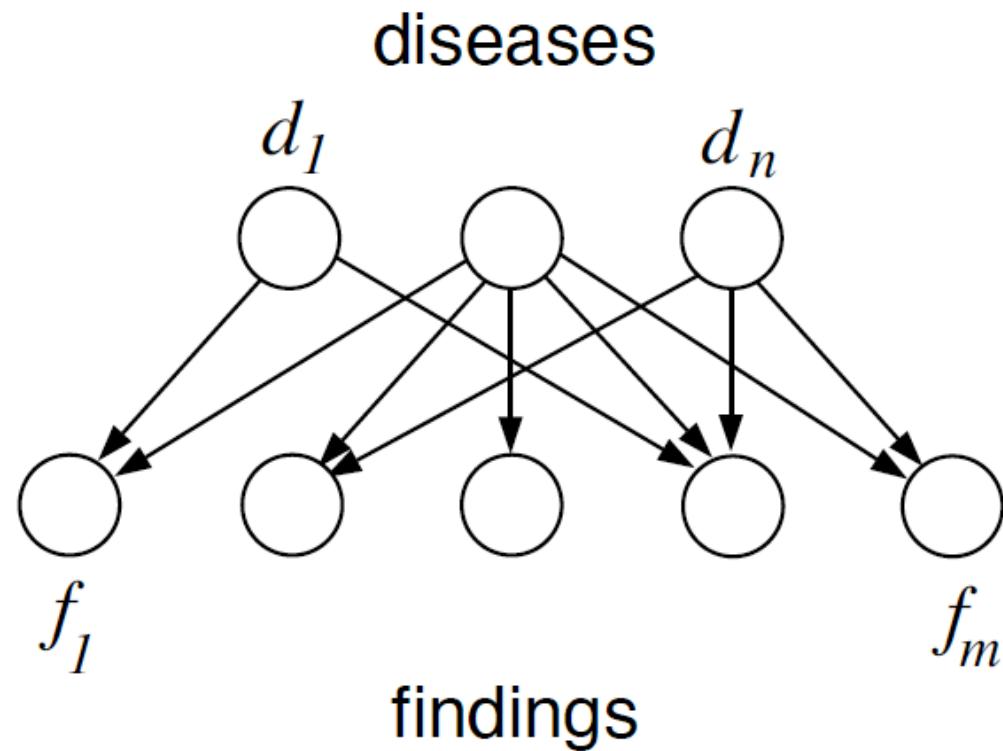


$$\begin{aligned}
 p(X_S, X_{CB}, X_{LC}, X_{SOB}, X_{XR}) &= p(X_S)p(X_{CB}|X_S)p(X_{LC}|X_{CB}, X_S)p(X_{SOB}|X_{LC}, X_{CB}, X_S) \\
 &\quad \times p(X_{XR}|X_{SOB}, X_{LC}, X_{CB}, X_S) \\
 &= p(X_S)p(X_{CB}|X_S)p(X_{LC}|X_S)p(X_{SOB}|X_{CB}, X_{LC})p(X_{XR}|X_{LC})
 \end{aligned}$$

Corresponding Joint Distribution

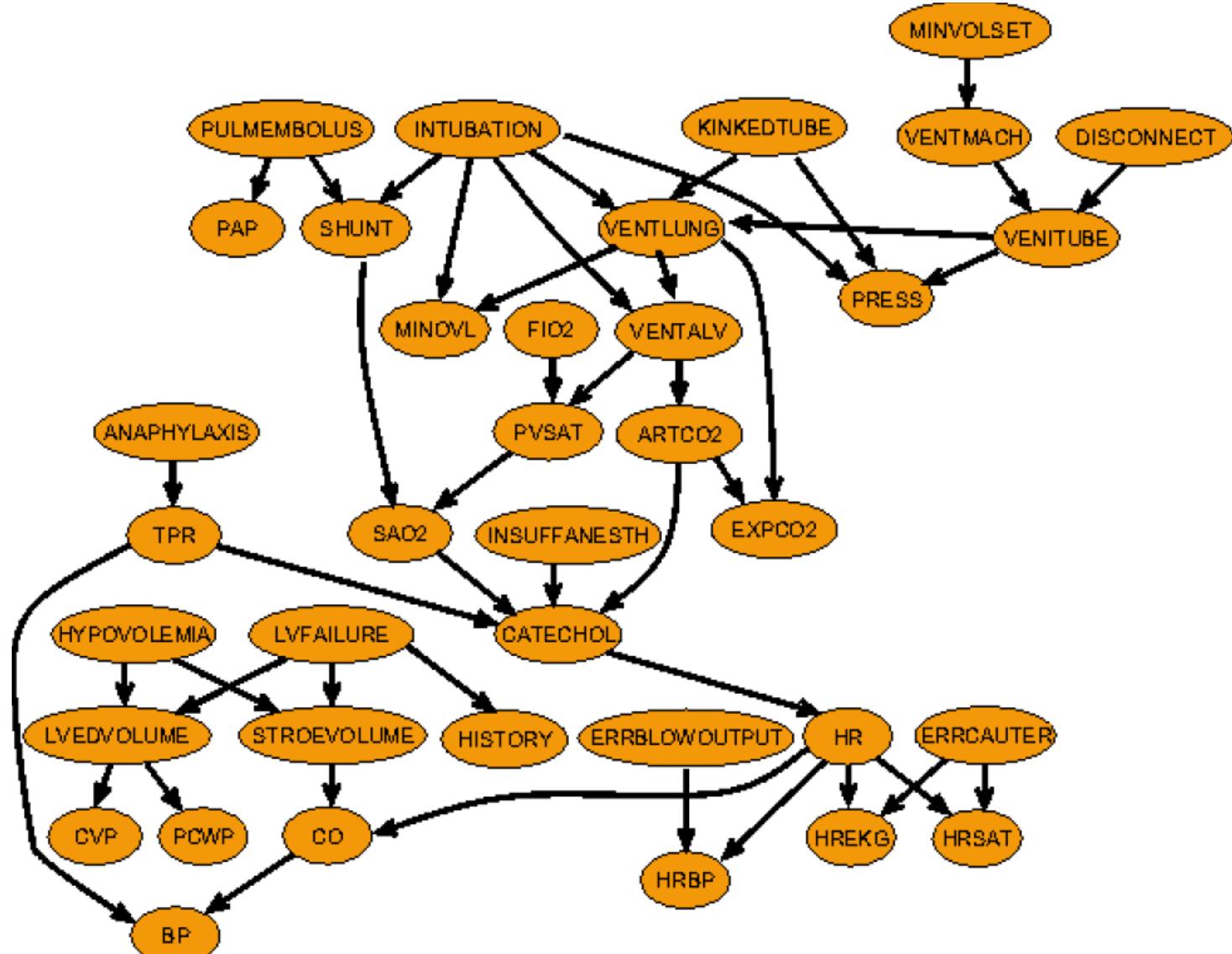


Medical Diagnosis



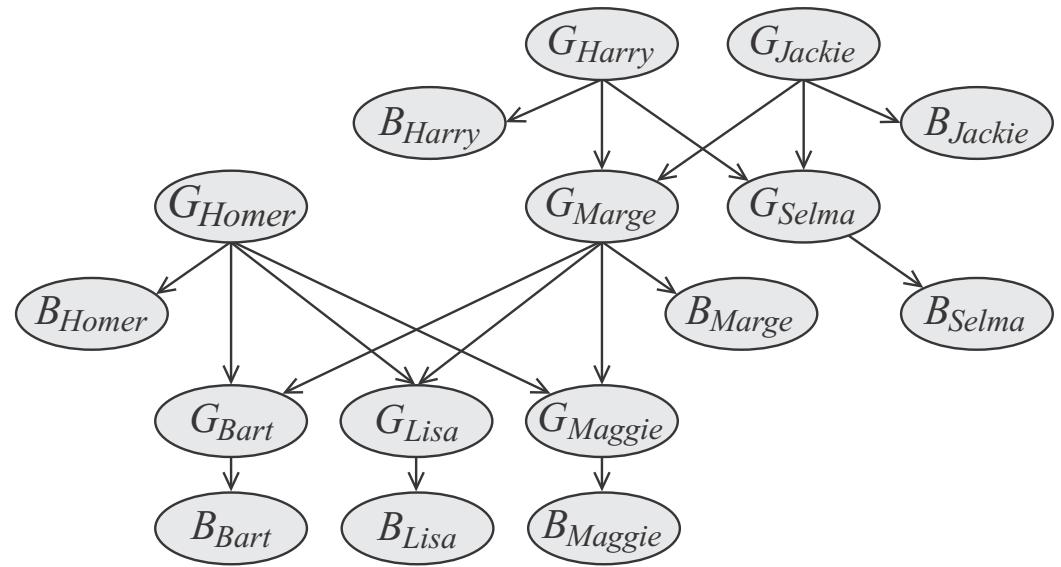
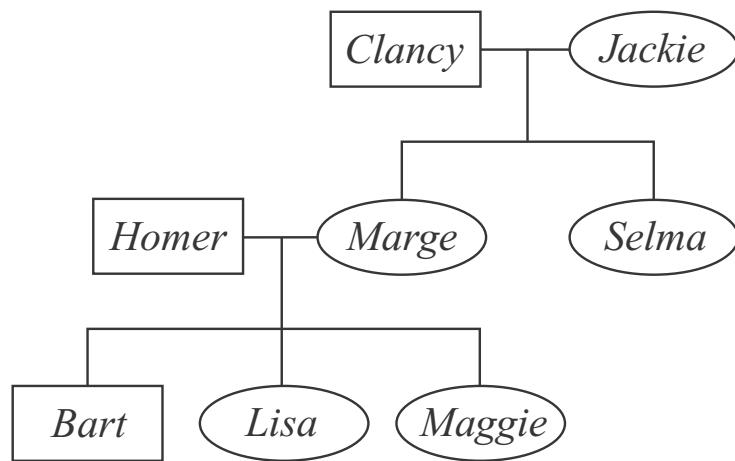
Parameterization: Noisy-OR, logistic regression,
generalized linear models...

Expert Systems



Alarm Network for ICU Monitoring

Bloodtypes and Genotypes



Learning for Graphical Models

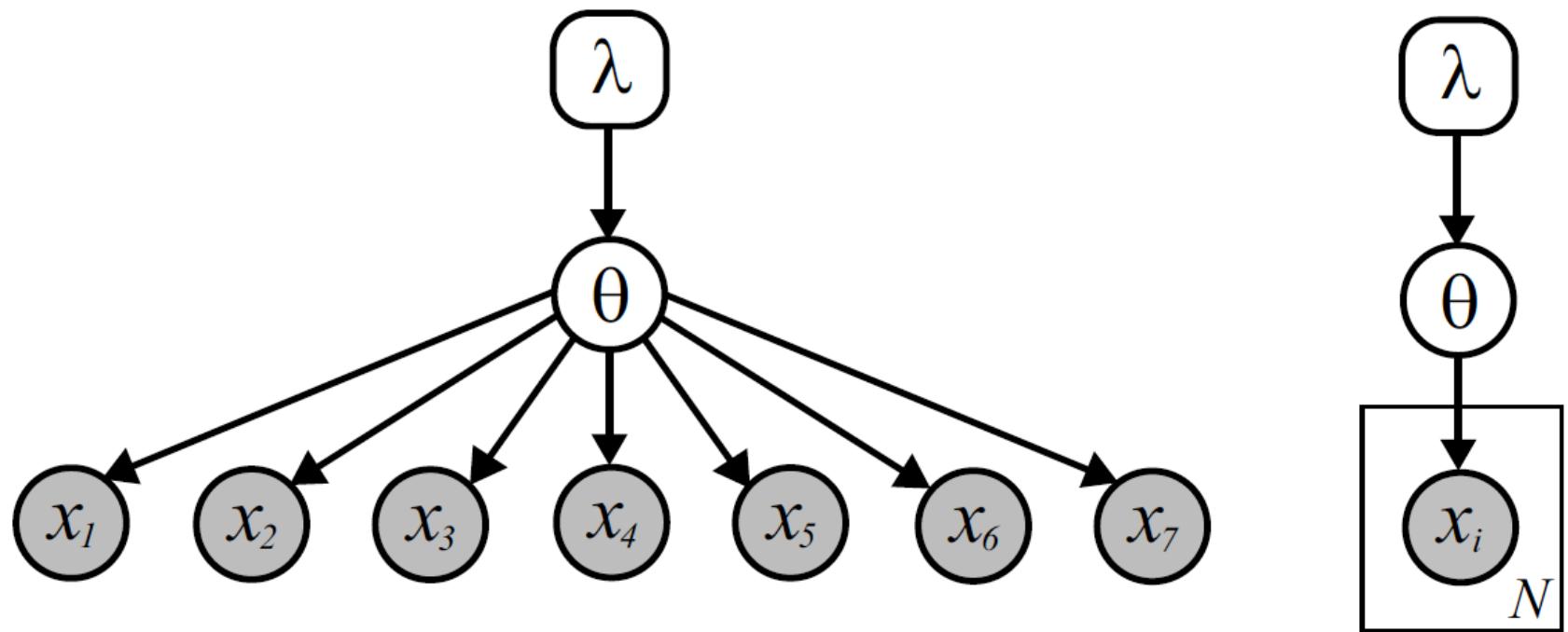
$$p(\mathbf{x}_{1:D}|G) = \prod_{s=1}^D p(x_s|\mathbf{x}_{\text{pa}(s)})$$

- Given a sequence of complete observations of all variables in graph
 - Log-likelihood decomposes: Predict each child node given its parent nodes
 - Discrete child node: classification problem
 - Continuous child node: regression problem
- What if some variables are unobserved?
 - Expectation-Maximization (EM) algorithm
 - MCMC algorithms
 - Alternate between “filling in” and re-estimating parameters

Machine Learning Problems

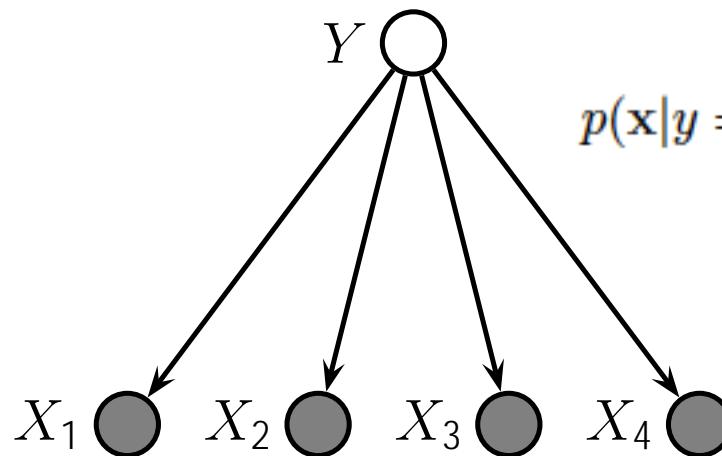
	<i>Supervised Learning</i>	<i>Unsupervised Learning</i>
<i>Discrete</i>	classification or categorization	clustering
<i>Continuous</i>	regression	dimensionality reduction

Learning with a Prior and Independent Observations



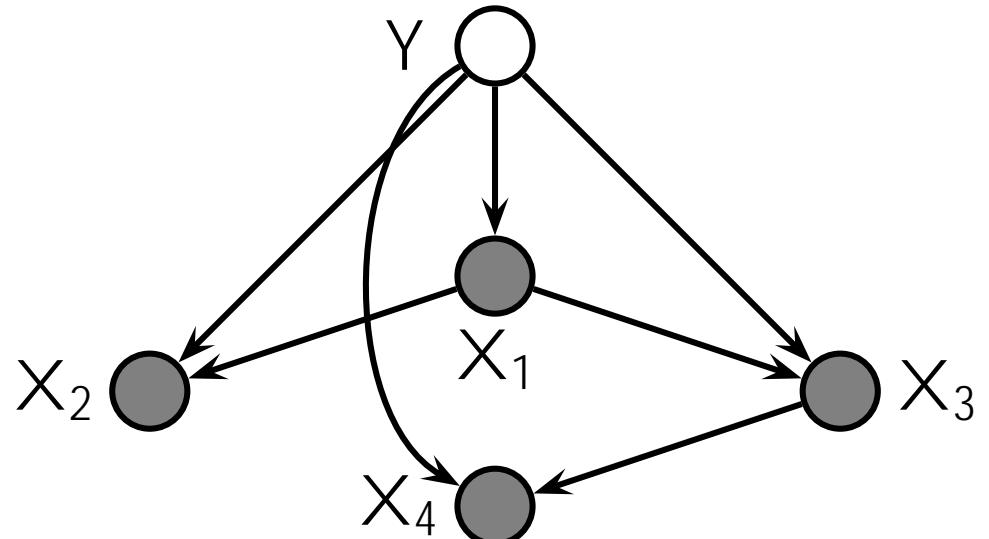
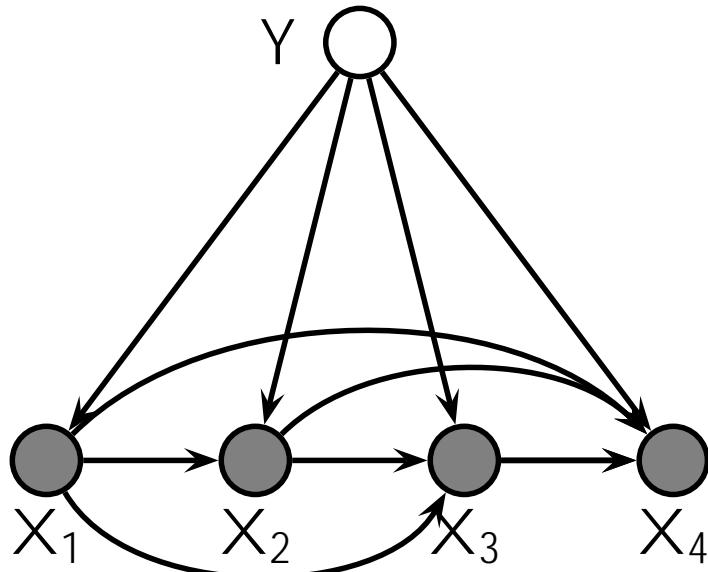
Sometimes-used convention: Shaded nodes are observed

Naïve Bayes Classifier

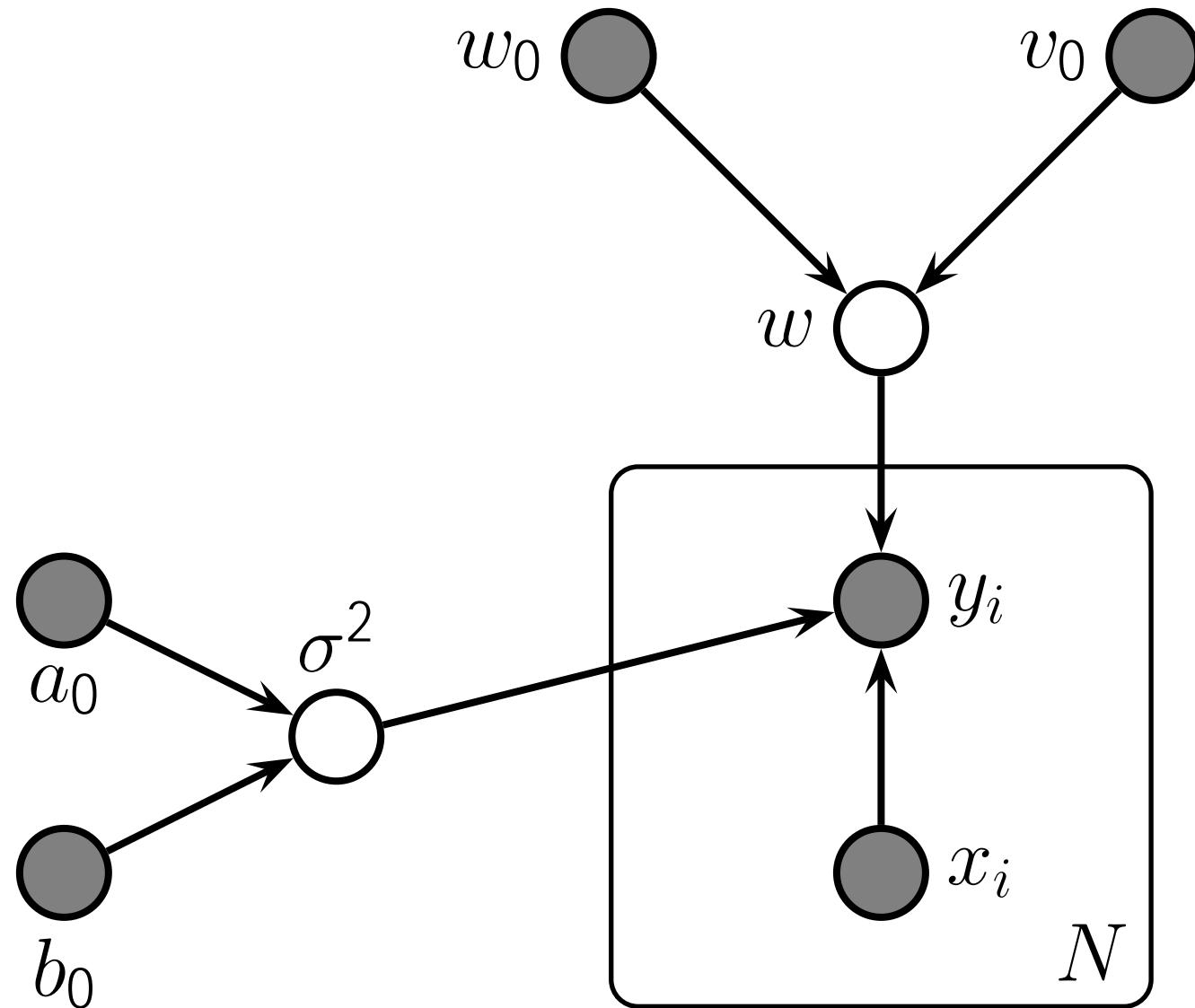


$$p(\mathbf{x}|y = c, \theta) = \prod_{j=1}^D p(x_j|y = c, \theta_j)$$

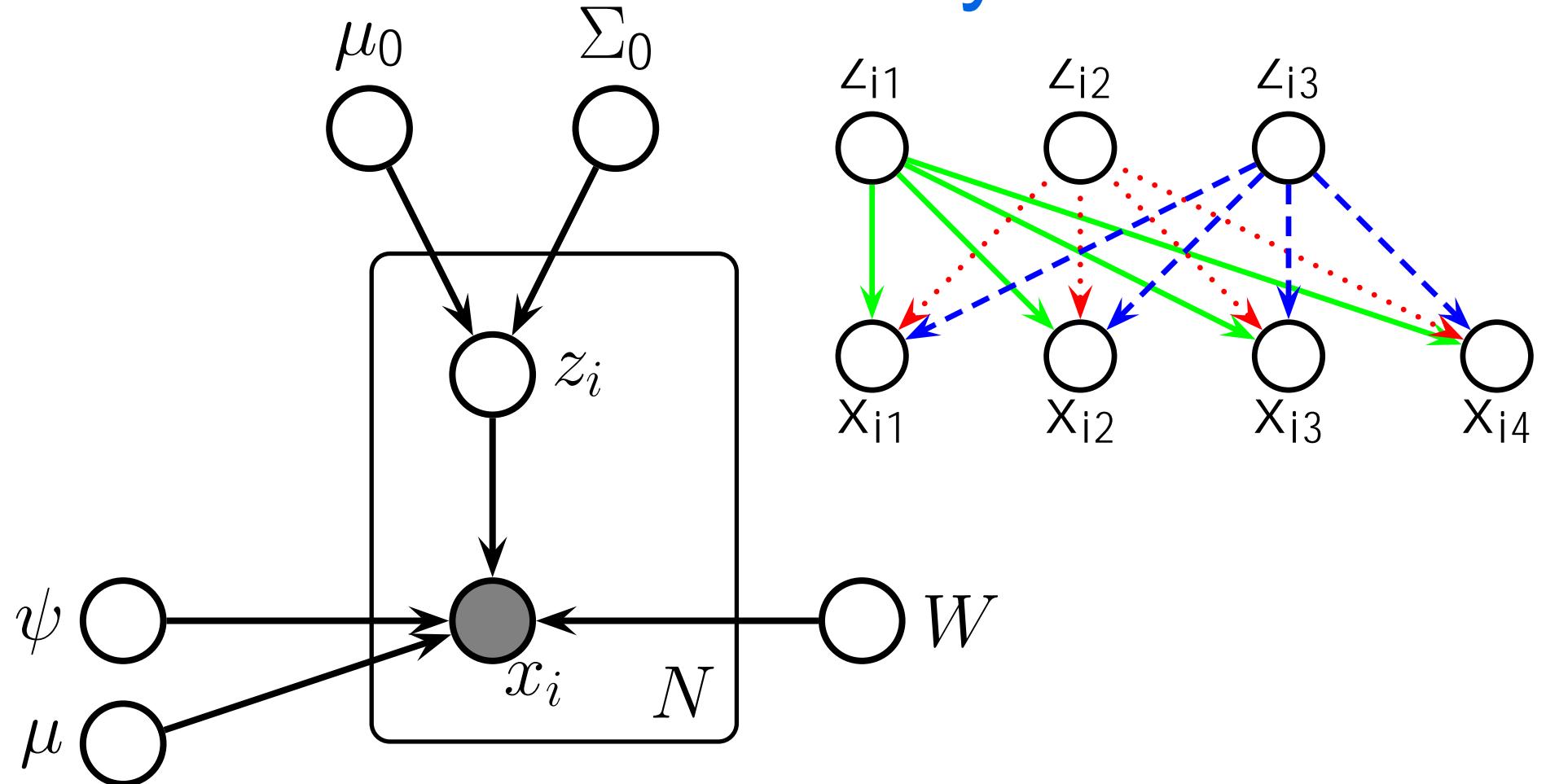
Allowing more dependence among features:



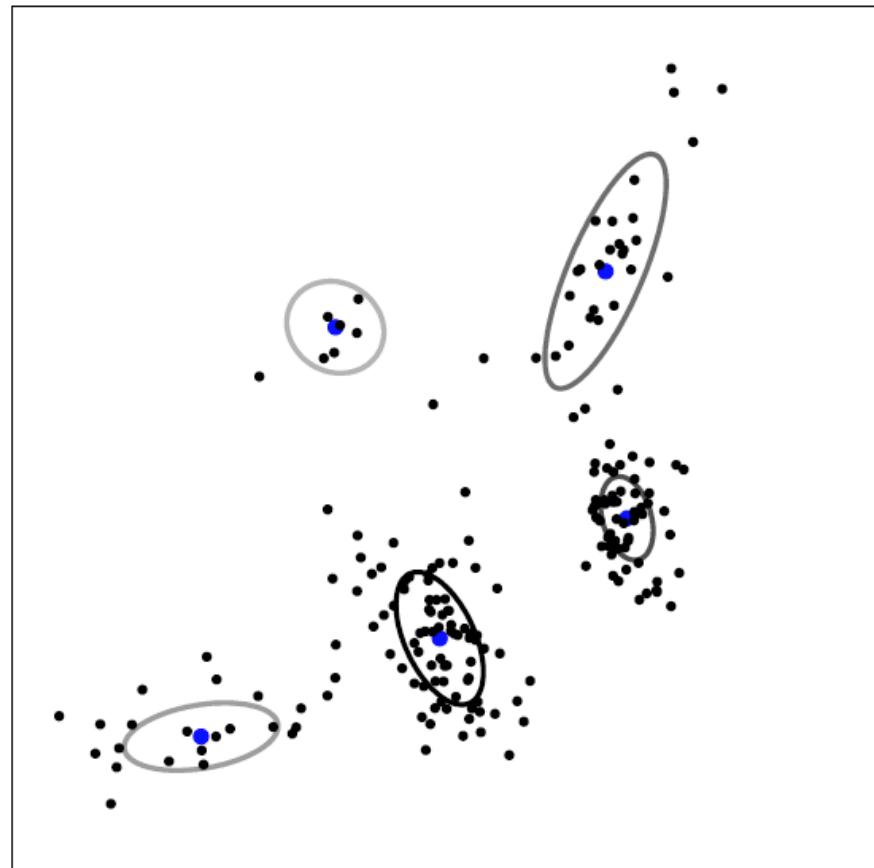
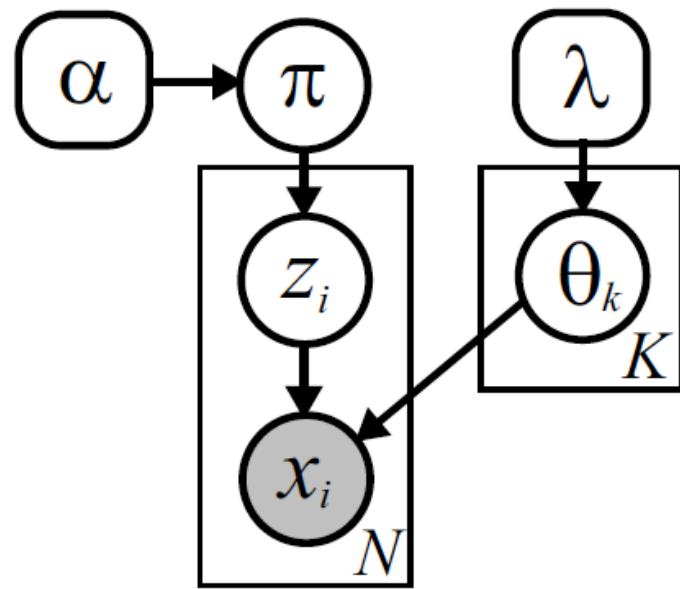
Bayesian Linear Regression



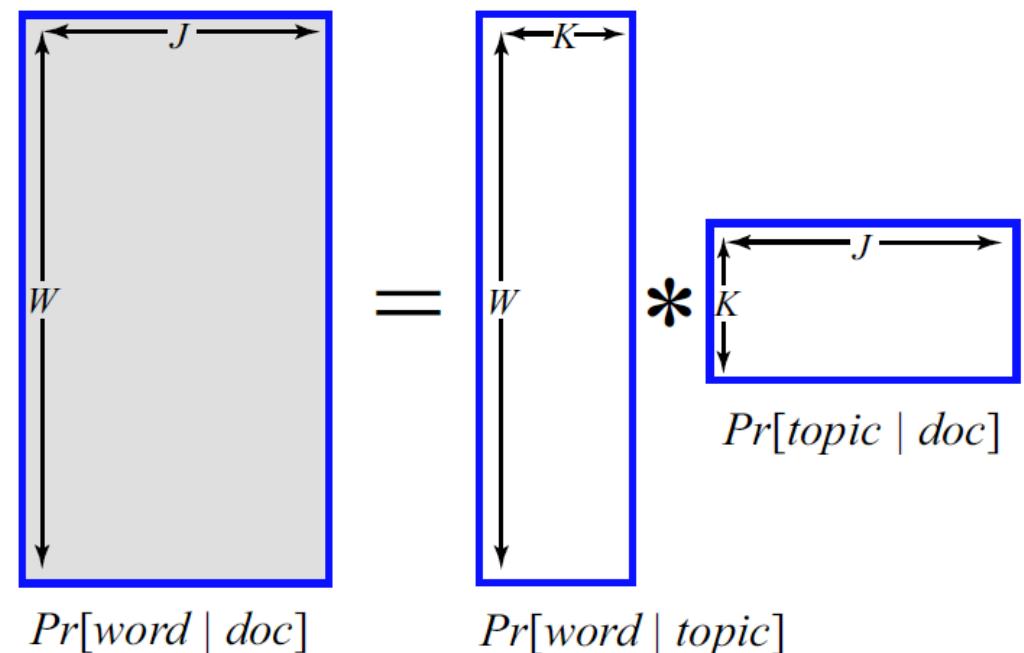
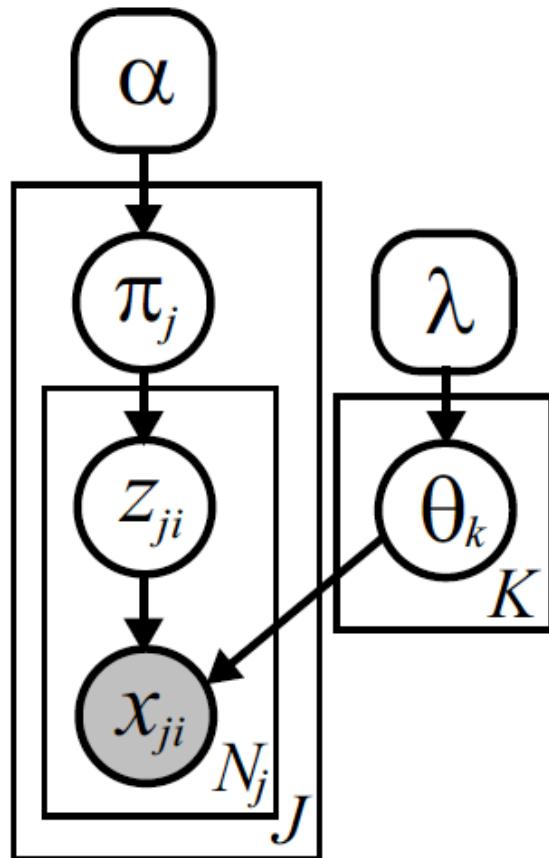
Factor Analysis



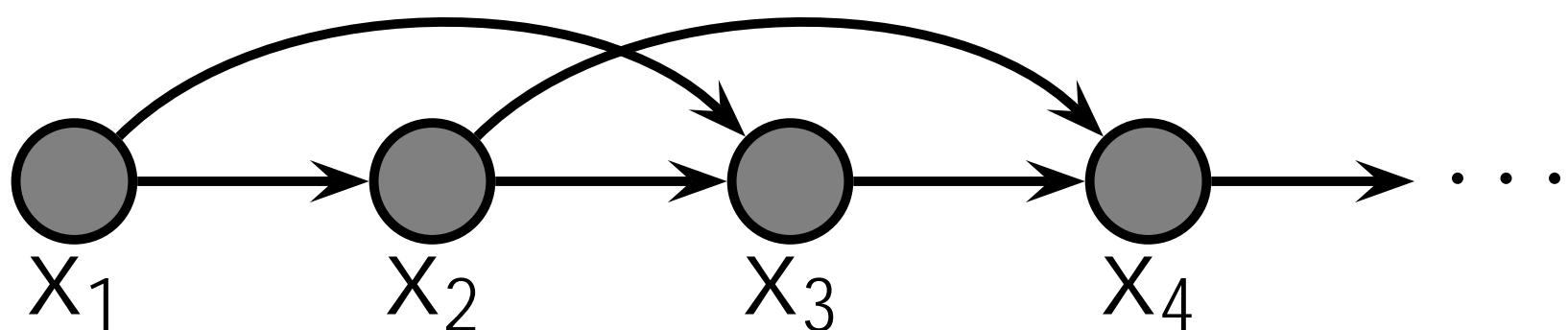
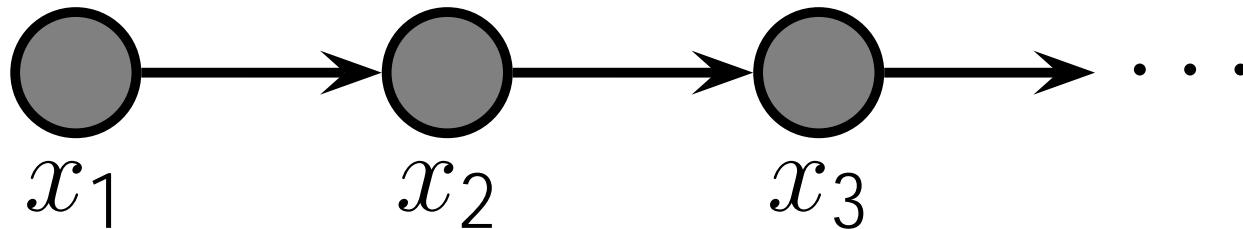
Mixture Models



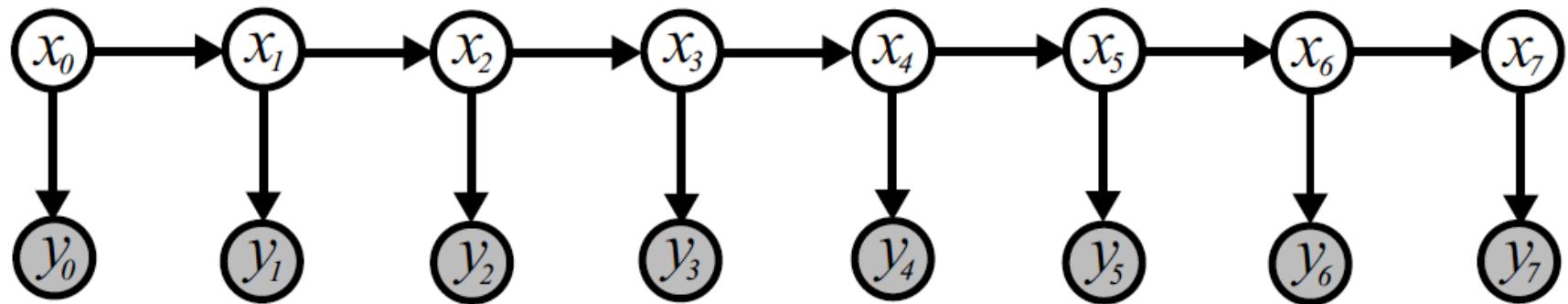
Admixture or Topic Models (Latent Dirichlet Allocation)



Markov Chains of Varying Orders

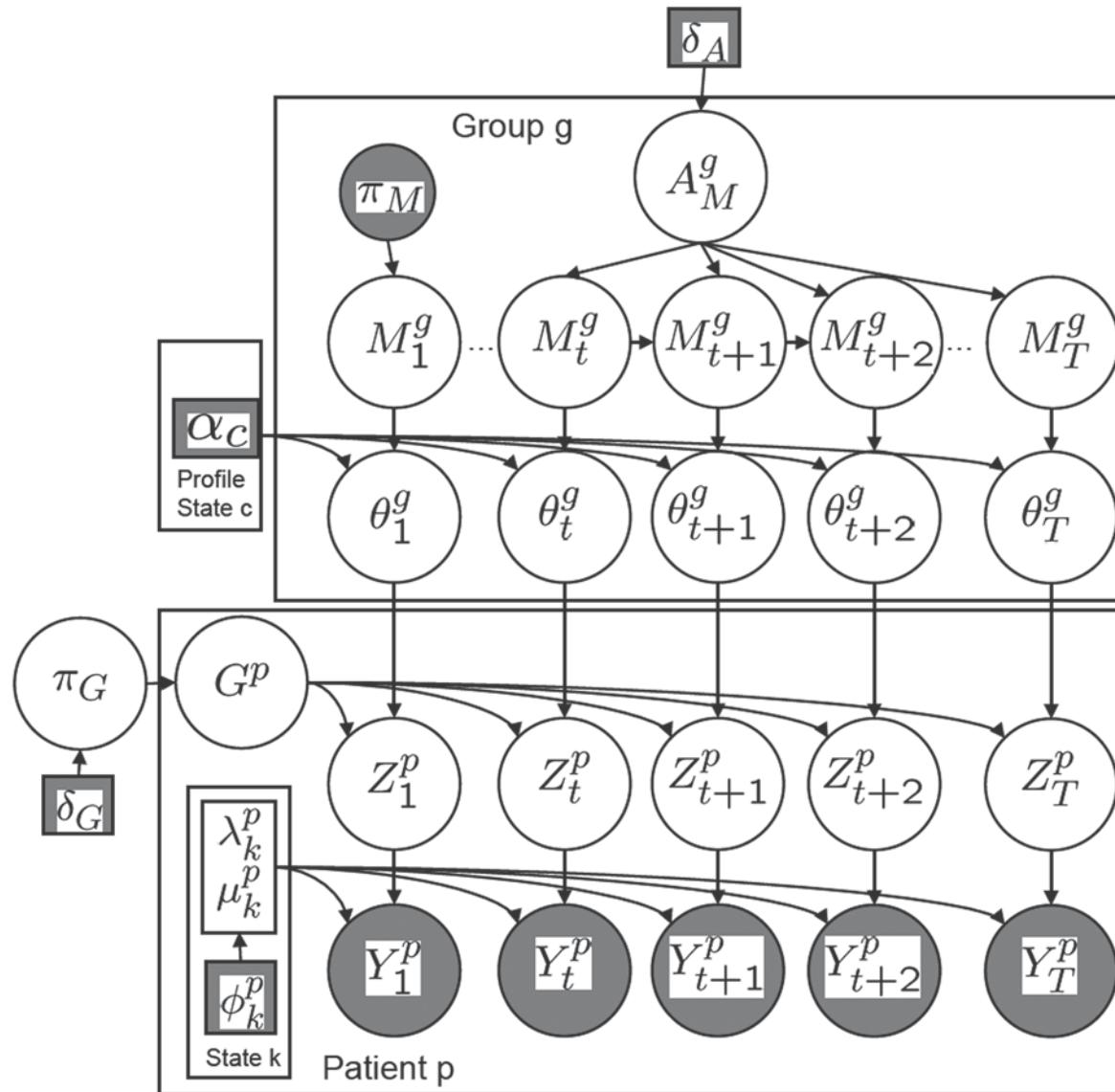


Hidden Markov Model (& linear dynamical system)

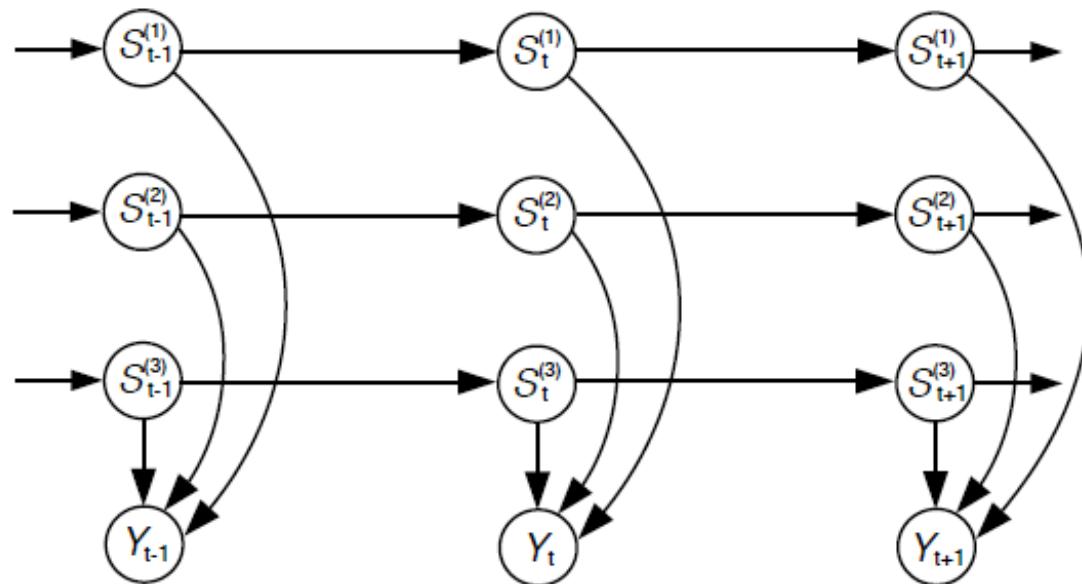


Advanced Models (presented for your amusement)

Mixture of HMMs

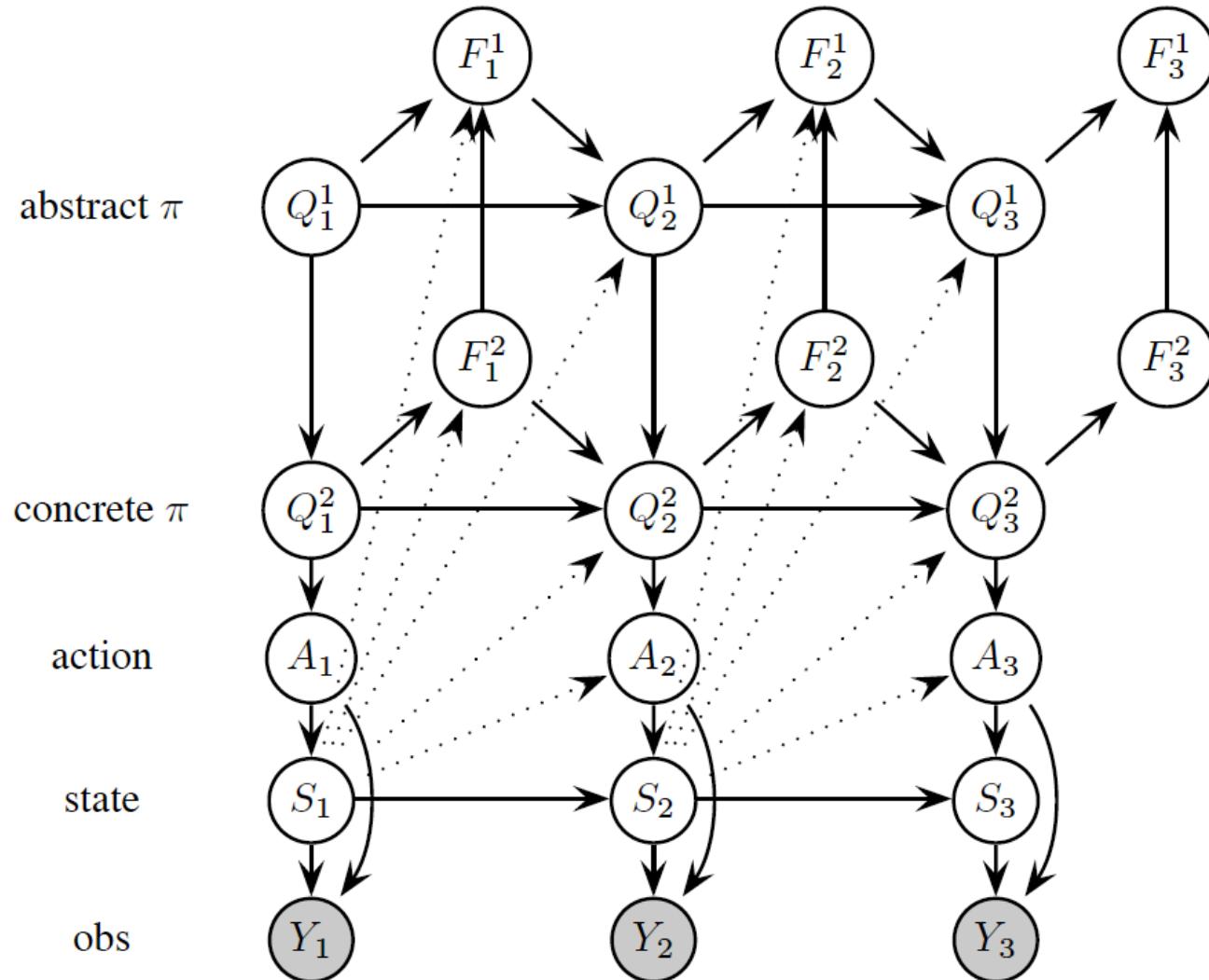


Factorial HMM



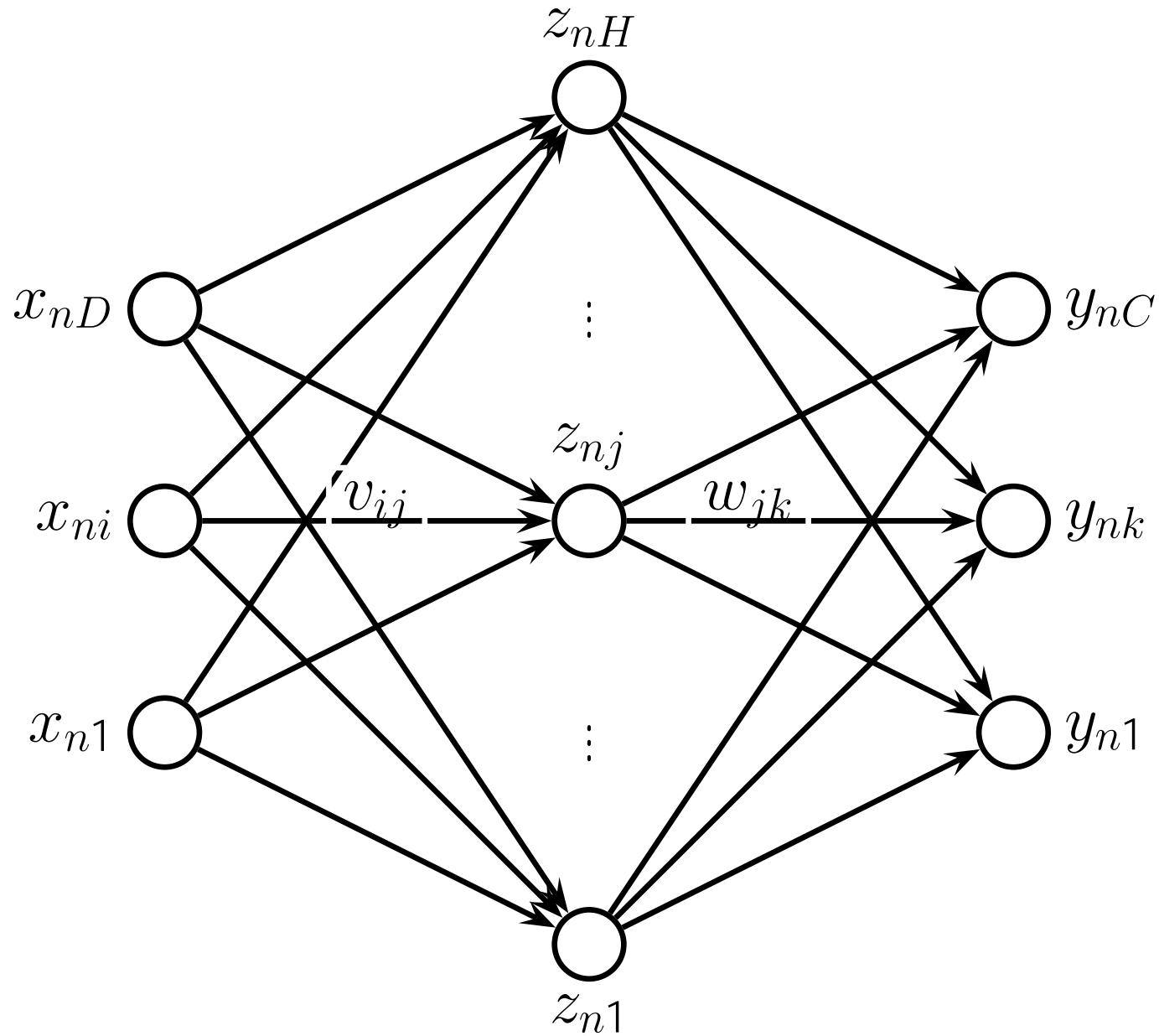
Ghahramani & Jordan, 1997

Dynamic Bayesian Networks

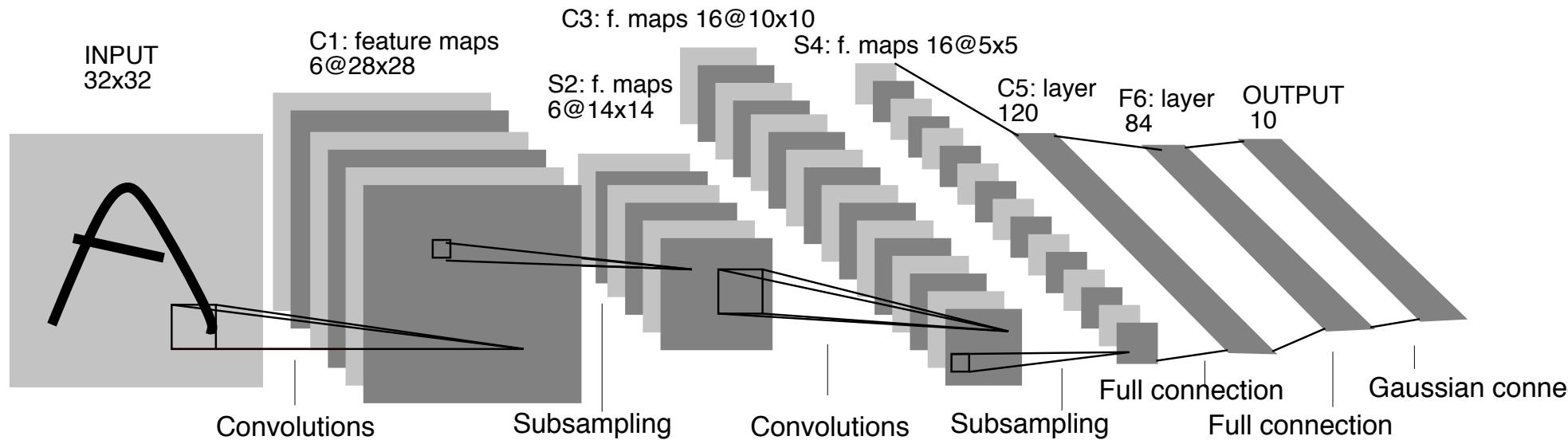


Murphy, 2002

Neural Network



Convolutional Neural Network



4 4->6	5 3->5	3 8->2	1 2->1	5 5->3	4 4->8	2 2->8	3 3->5	6 6->5	7 7->3
4 9->4	8 8->0	7 7->8	5 5->3	7 8->7	6 0->6	3 3->7	2 2->7	3 8->3	4 9->4
8 8->2	5 5->3	4 4->8	3 3->9	0 6->0	9 9->8	9 4->9	6 6->1	4 9->4	1 9->1
9 9->4	2 2->0	1 6->1	3 3->5	3 3->2	9 9->5	0 6->0	6 6->0	5 6->8	8 6->8
4 4->6	7 7->3	9 9->4	4 4->6	2 2->7	7 9->7	4 4->3	9 9->4	9 9->4	9 9->4
7 8->7	4 4->2	8 8->4	5 3->5	6 8->4	6 6->5	8 8->5	3 3->8	3 3->8	9 9->8
1 1->5	9 9->8	6 6->3	0 0->2	6 6->5	9 9->5	0 0->7	6 1->6	4 4->9	1 2->1
2 2->8	8 8->5	4 4->9	7 7->2	7 7->2	6 6->5	9 9->7	1 6->1	6 5->6	5 5->0
4 4->9	2 2->8								

LeNet 5