# COMP90051 **Statistical Machine Learning**

Semester 2, 2015

Lecturer:  Ben Rubinstein

6. PGM Representation

# Next Lectures

- Representation of joint distributions

- Examples

- Probabilistic inference
    * Computing other distributions from joint

- Statistical inference
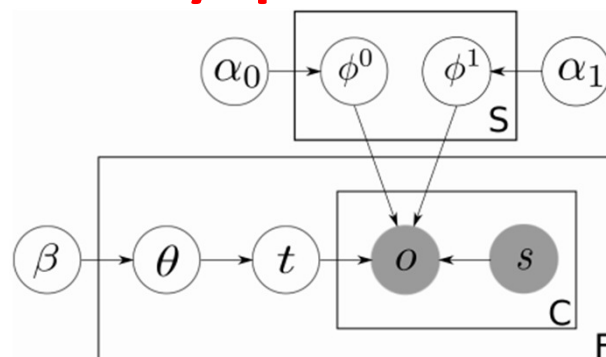    * Learn parameters from (missing) data

# Probabilistic Graphical Models

*Marriage of graph theory and probability theory.*
*Tool of choice for Bayesian statistical learning.*

*We'll stick with easier discrete case,*
*ideas generalise to continuous.*

# Motivation by practical importance



- **Many applications**
  - ∗ Phylogenetic trees
  - ∗ Pedigrees, Linkage analysis
  - ∗ Error-control codes
  - ∗ Speech recognition
  - ∗ Document topic models
  - ∗ Probabilistic parsing
  - ∗ Image segmentation
  - ∗ …

- **Unifies many previously-discovered algorithms**
  - ∗ HMMs
  - ∗ Kalman filters
  - ∗ Mixture models
  - ∗ LDA
  - ∗ MRFs
  - ∗ CRF
  - ∗ Logistic, linear regression
  - ∗ …

4

# Motivation by way of comparison

## Bayesian statistical learning

- Model joint distribution of X's,Y and parameter r.v.'s
    * "Priors": marginals on parameters

- Training: update prior to posterior using observed data

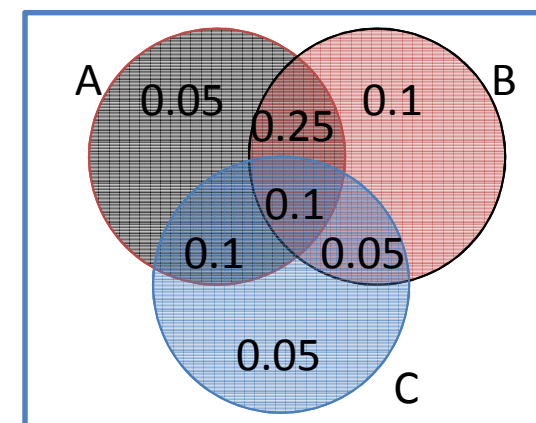- Prediction: output posterior, or some function of it (MAP)

## PGMs aka "Bayes Nets"

- Efficient joint representation
    * Independence made explicit
    * Trade-off between expressiveness and need for data, easy to make
    * Easy for practitioners to model

- Algorithms for fit parameters, compute marginals, posterior

# Everything Starts at the Joint Distribution

- All joint distributions on discrete r.v.'s can be represented as tables

- #rows grows exponentially with #r.v.'s

- Example: Truth Tables
    * $M$ Boolean r.v.'s require $2^M$-1 rows
    * Table assigns probability per row

| A | B | C | Prob |
|---|---|---|------|
| 0 | 0 | 0 | 0.30 |
| 0 | 0 | 1 | 0.05 |
| 0 | 1 | 0 | 0.10 |
| 0 | 1 | 1 | 0.05 |
| 1 | 0 | 0 | 0.05 |
| 1 | 0 | 1 | 0.10 |
| 1 | 1 | 0 | 0.25 |
| 1 | 1 | 1 | ? |



6

# **The Good**: What we can do with the joint

- **Probabilistic inference** from joint on r.v.'s
  - ∗ Computing any other distributions involving our r.v.'s

- Pattern: want a distribution, have joint; use: Bayes rule + marginalisation

- Example: **naïve Bayes classifier**
  - ∗ Predict class $y$ of instance $x$ by maximising

$$\Pr(Y = y | X = x) = \frac{\Pr(Y=y,X=x)}{\Pr(X=x)} = \frac{\Pr(Y=y,X=x)}{\sum_y \Pr(X=x,Y=y)}$$

# **The Bad & Ugly**: Tables *waaaaay* too large*!!*

- **The Bad**: Computational complexity
  - ∗ Tables have exponential number of rows in number of r.v.'s
  - ∗ Therefore → poor space & time to marginalise

- **The Ugly**: Model complexity
  - ∗ Way too flexible
  - ∗ Way too many parameters to fit
    → need lots of data OR will overfit

- Antidote: assume independence!

| A | B | C | Prob |
|---|---|---|------|
| 0 | 0 | 0 | 0.30 |
| 0 | 0 | 1 | 0.05 |
| 0 | 1 | 0 | 0.10 |
| 0 | 1 | 1 | 0.05 |
| 1 | 0 | 0 | 0.05 |
| 1 | 0 | 1 | 0.10 |
| 1 | 1 | 0 | 0.25 |
| 1 | 1 | 1 | ? |

# Example: You're late!

- Modeling a tardy lecturer. Boolean r.v.'s
  - $B$: Ben teaches the class
  - $S$: It is sunny (o.w. bad weather)
  - $L$: The lecturer arrives slightly late (o.w. on time)

- Assume: Ben sometimes delayed by bad weather, Ben more likely late than co-lecturer
  - $\Pr(S|B) = \Pr(S)$, $\Pr(S) = 0.3$, $\Pr(B) = 0.6$

- Lateness not independent on weather, lecturer

  $\Pr(L = true | \dots)$

  |  |  | B | |
  |---|---|---|---|
  |  |  | False | True |
  | S | False | 0.2 | 0.1 |
  |  | True | 0.1 | 0.05 |

  - Need $\Pr(L|B = b, S = s)$ for all combinations

- Need just 6 parameters

# Independence: not a dirty word

| Lazy Lecturer Model | Model details | # params |
|---|---|---|
| Our model with $S, B$ independence | $\Pr(S, B)$ factors to $\Pr(S)\Pr(B)$ | 2 |
| | $\Pr(L|B, S)$ modelled in full | 4 |
| Assumption-free model | $\Pr(L, B, S)$ modelled in full | 7 |

- Independence assumptions
  * Can be reasonable in light of domain expertise
  * Allow us to factor → Key to tractable models

# Factoring Joint Distributions

- Chain Rule: for any ordering of r.v.'s can always factor:

$$\Pr(X_1, X_2, \ldots, X_k) = \prod_{i=1}^{k} \Pr(X_i | X_{i+1}, \ldots, X_k)$$
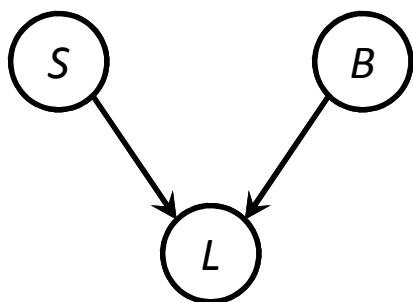
- Model's independence assumptions correspond to
  - Dropping conditioning r.v.'s in the factors!
  - Example unconditional indep.: $\Pr(X_1 | X_2) = \Pr(X_1)$
  - Example conditional indep.: $\Pr(X_1 | X_2, X_3) = \Pr(X_1 | X_2)$

- Example: independent r.v.'s $\Pr(X_1, \ldots, X_k) = \prod_{i=1}^{k} \Pr(X_i)$

- Simpler factors speed inference and avoid overfitting

# Directed PGM

- **Nodes**

- **Edges** (acyclic)

- **Random variables**

- **Conditional dependence**
    - * Node table: $\Pr(child|parents)$
    - * Child directly depends on parents

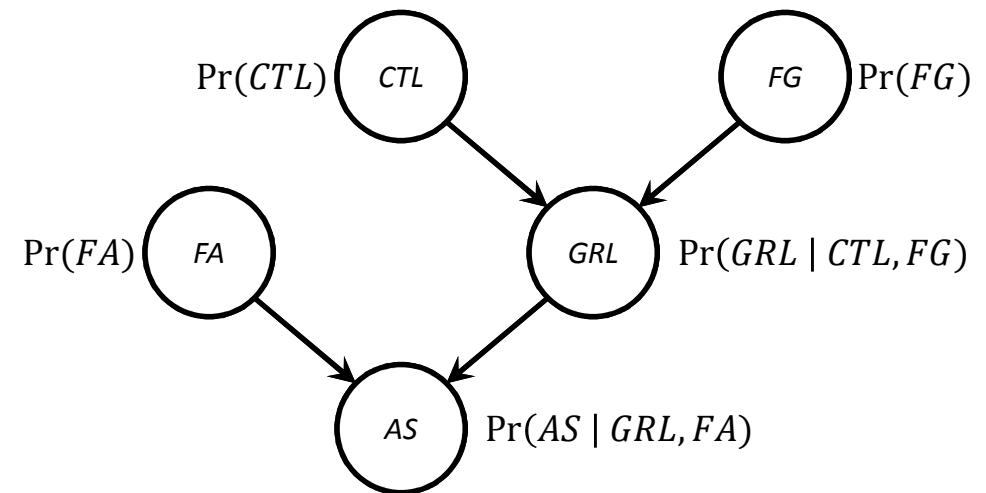- **Joint factorisation**

$$\Pr(X_1, X_2, \ldots, X_k) = \prod_{i=1}^{k} \Pr\big(X_i | X_j \in parents(X_i)\big)$$

*Tardy Lecturer Example*



$\Pr(S)$          $\Pr(B)$
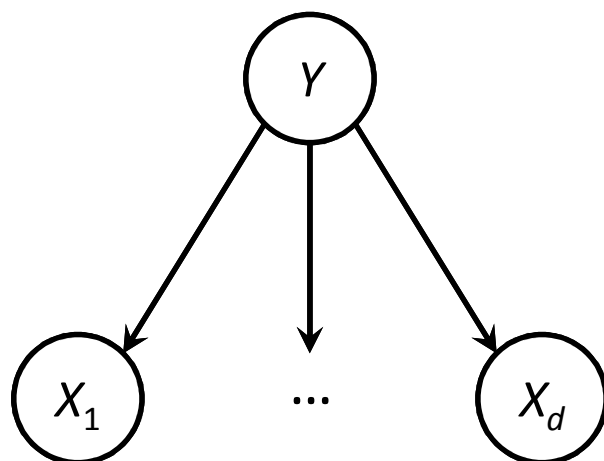
$\Pr(L|S,B)$

# Example: Nuclear power plant

- Core temperature
  → Temperature Gauge
  → Alarm

- Model uncertainty in monitoring failure
  * GRL: gauge reads low
  * CTL: core temperature low
  * FG: faulty gauge
  * FA: faulty alarm
  * AS: alarm sounds

- PGMs to the rescue!

$\Pr(CTL)$ CTL          FG $\Pr(FG)$

$\Pr(FA)$ FA          GRL $\Pr(GRL \mid CTL, FG)$

AS  $\Pr(AS \mid GRL, FA)$

Joint $\Pr(CTL, FG, FA, GRL, AS)$ given by

$\Pr(AS|FA, GRL) \Pr(FA) \Pr(GRL|CTL, FG) \Pr(CTL) \Pr(FG)$

# Naïve Bayes



$Y \sim \text{Bernoulli}(\theta)$
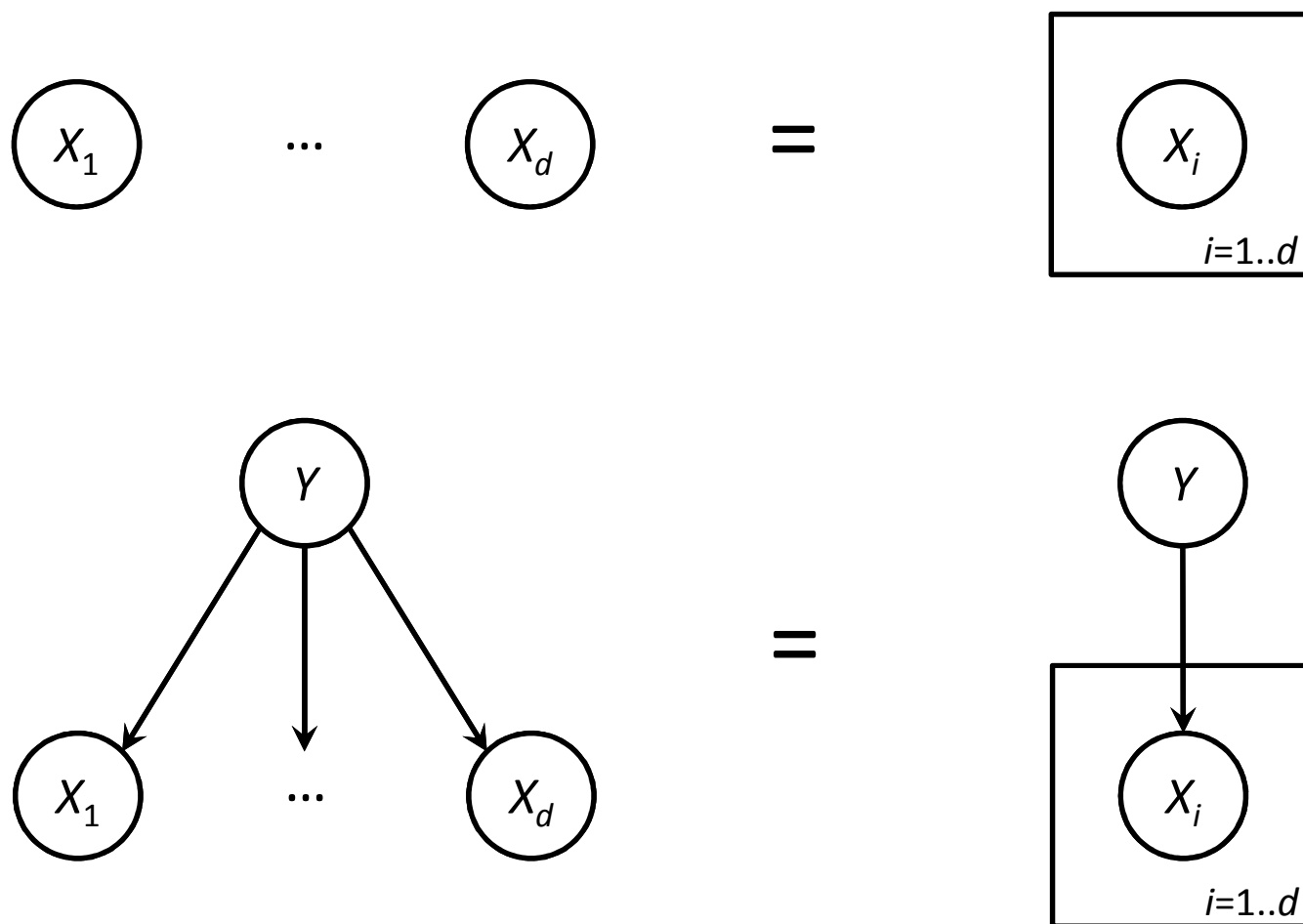
$X_j | Y \sim \text{Bernoulli}(\theta_{j,Y})$

$\text{Pr}(Y, X_1, \ldots, X_d)$
$\quad = \text{Pr}(X_1, \ldots, X_d, Y)$
$\quad = \text{Pr}(X_1 | X_2, \ldots, X_d, Y) \, \text{Pr}(X_2 | X_3, \ldots, X_d, Y) \ldots \text{Pr}(X_{d+1} | X_d, Y) \, \text{Pr}(X_d | Y) \, \text{Pr}(Y)$
$\quad = \text{Pr}(X_1 | Y) \, \text{Pr}(X_2 | Y) \ldots \text{Pr}(X_d | Y) \, \text{Pr}(Y)$
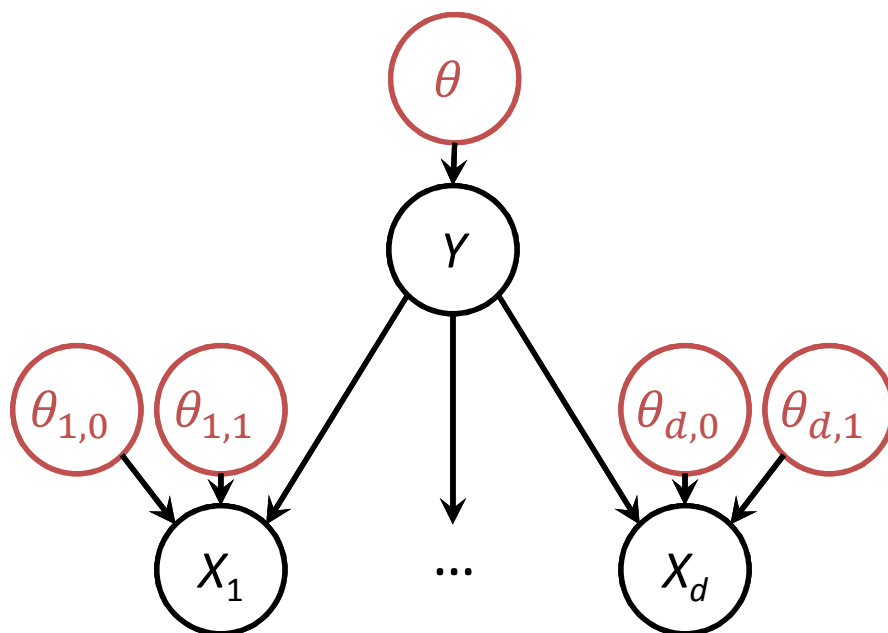
Prediction: predict label maximising $\text{Pr}(Y | X_1, \ldots, X_d)$

# Short-hand for repeats: Plate notation

# PGMs frequentist OR Bayesian…

- PGMs represent joints, which are central to Bayesian

- Catch is that Bayesians add: node per parameters, with table being the parameter's prior



$Y \sim \text{Bernoulli}(\theta)$

$X_j | Y \sim \text{Bernoulli}(\theta_{j,Y})$

$\theta's \sim Beta$

16

# Other types of PGMs…

## Undirected PGM

- Graph
  - * Edges undirected

- Probability
  - * Each node a r.v.
  - * Each clique $C$ has "factor"
    $$f_C(X_j : j \in C) \geq 0$$
  - * Joint $\propto$ product of factors

## Directed PGM

- Graph
  - * Edged directed

- Probability
  - * Each node a r.v.
  - * Each node has conditional
    $$p(X_i | X_j \in parents(X_i))$$
  - * Joint $=$ product of cond'ls

# Summary

- Probabilistic graphical models
  - ∗ Motivation: applications, unifies algorithms
  - ∗ Motivation: ideal tool for Bayesians
  - ∗ Independence lowers computational/model complexity
  - ∗ PGMs: compact representation of factorised joints
  - ∗ Directed and Undirected PGMs