

# Bayes in AI

## 4.1 ('Small') PGMs and (shallow) NNs

David Ríos Insua & Roi Naveiro

[david.rios@icmat.es](mailto:david.rios@icmat.es) [roi.naveiro@icmat.es](mailto:roi.naveiro@icmat.es)

# Objectives

## Earlier sessions

1. Bayes or die
2. MCMC
3. Large scale Bayes: VB and SGMCMC

## Today

- 4.1 Bayes in AI: (small) PGMs and (shallow) NNs [david.rios@icmat.es](mailto:david.rios@icmat.es)
- 4.2 GPs and Bayesian NNs in function spaces [simon.rodriguez@icmat.es](mailto:simon.rodriguez@icmat.es)
- 4.3 Modern research in ML [roi.naveiro@icmat.es](mailto:roi.naveiro@icmat.es)

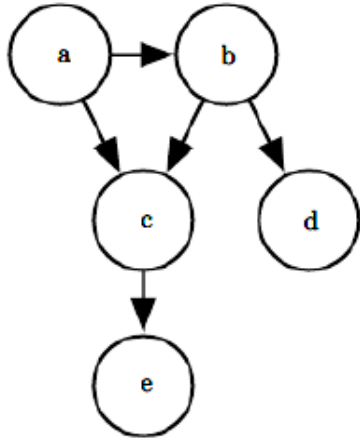
# PGMs. Motivation

# Motivation

- Simple way to visualize structure of probabilistic models
- Designing and motivating new models
- Understanding properties like conditional independence
- Complex computations viewed through simple graphical manipulations
- Explainable and interpretable
- Deep belief nets in deep learning

# Concept

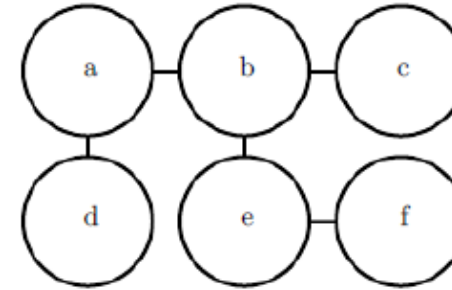
$$p(\mathbf{x}) = \prod_i p(x_i | Pa_{\mathcal{G}}(x_i))$$



$$p(a, b, c, d, e) = p(a)p(b | a)p(c | a, b)p(d | b)p(e | c)$$

Bayesian networks. Directed, Acyclic

$$\tilde{p}(\mathbf{x}) = \prod_{\mathcal{C} \in \mathcal{G}} \phi(\mathcal{C})$$



$$p(a, b, c, d, e, f) = \frac{1}{Z} \phi_{a,b}(a, b) \phi_{b,c}(b, c) \phi_{a,d}(a, d) \phi_{b,e}(b, e) \phi_{e,f}(e, f)$$

Markov fields. Undirected

# Probabilistic graphical models. Directed Bayesian networks

# Directed PGMs

As basic tools for qualitative modelling of uncertainty use probabilistic influence diagrams a.k.a. causal networks, Bayesian networks, Belief networks,... See the excellent

[http://en.wikipedia.org/wiki/Bayesian\\_network](http://en.wikipedia.org/wiki/Bayesian_network)

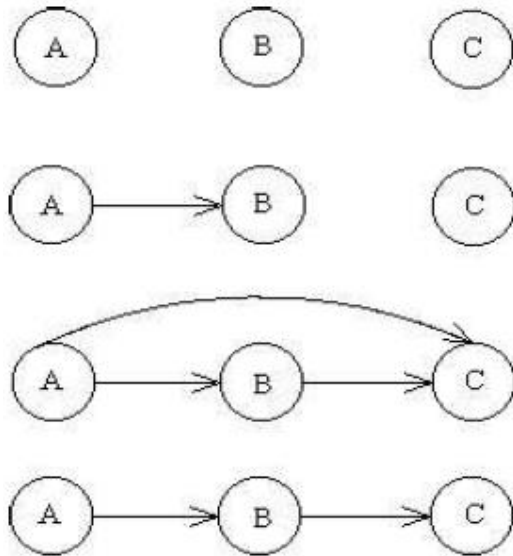
They are **influence diagrams** with chance nodes only. Qualitatively they describe a probabilistic model through

$$P(A_1, A_2, \dots, A_n) = P(A_1 \mid \text{ant}(A_1)) \dots P(A_n \mid \text{ant}(A_n))$$

where  $\text{ant}(A_i)$  are the antecessors of node  $A_i$ .

In what follows we see several PIDs and we need to indicate the entailed probabilistic model

# Probabilistic diagrams with three nodes

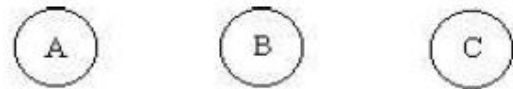


Before moving forward, write the entailed probabilistic model

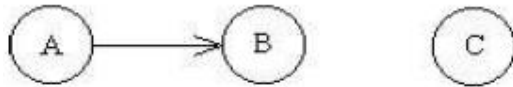


# Probabilistic diagrams with three nodes

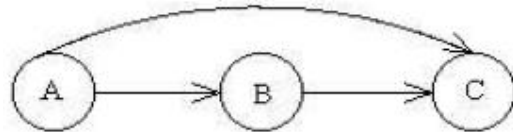
Model  $P(A, B, C)$



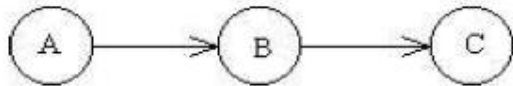
$P(A)P(B)P(C)$



$P(A) P(B|A) P(C)$



$P(A)P(B|A)P(C|A,B)$

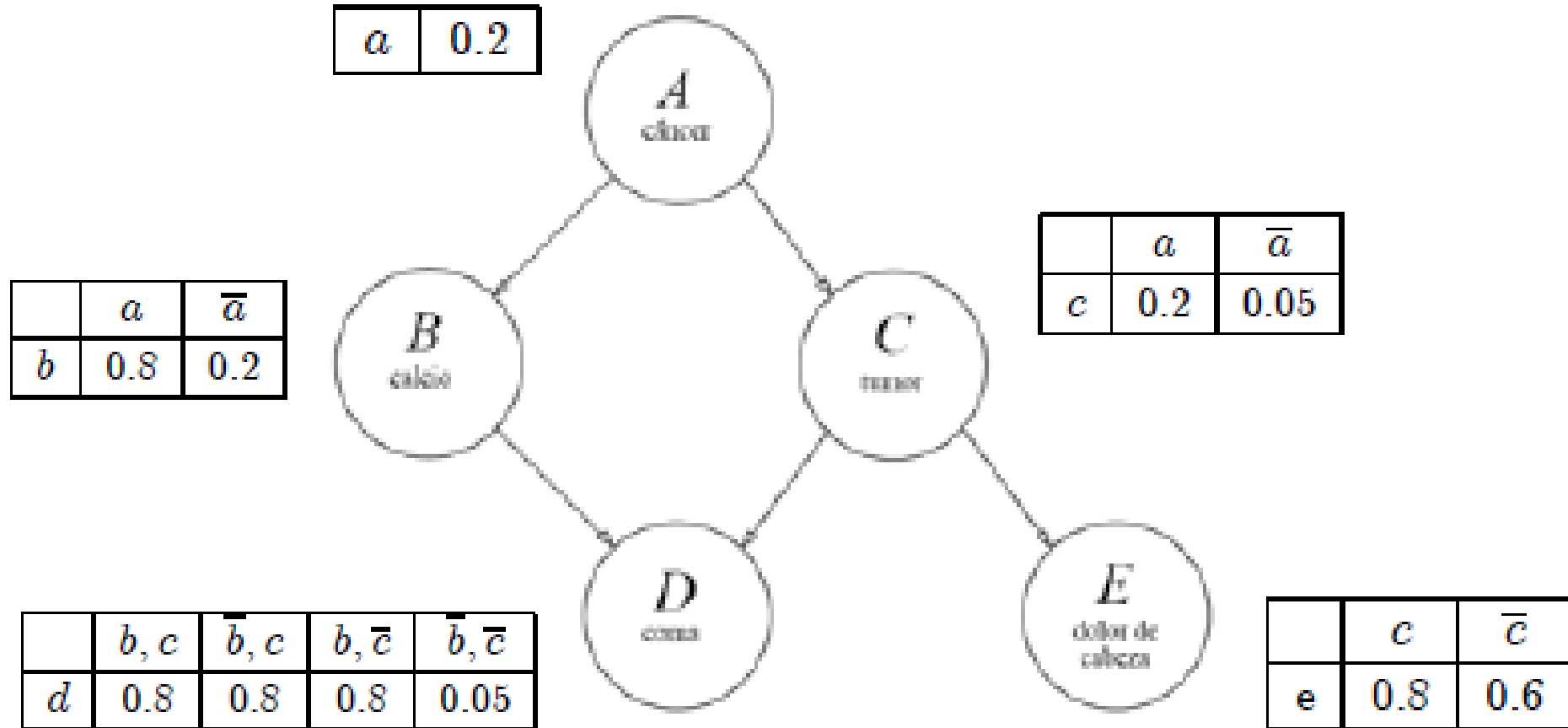


$P(A)P(B|A)P(C|B)$

First case, independence. Third case, A and C are conditionally independent given B.

Read [http://en.wikipedia.org/wiki/Conditional\\_independence](http://en.wikipedia.org/wiki/Conditional_independence)

# The hidden info



$$P(A, B, C, D, E) = P(A)P(B|A)P(C|A)P(D|B, C)P(E|C)$$

# Probabilistic diagrams. Asia

An example referring to lung diseases

A breathing condition (dyspnea) may be due to tuberculosis, lung cancer or bronchitis, none of them or several of them. A recent visit to Asia, increases the chances of tuberculosis, whereas smoking is a risk factor for lung cancer and bronchitis. The results of an X-ray may not discriminate between cancer and tuberculosis, as neither the presence or absence of dyspnea does.

# Probabilistic diagrams

An example referring to lung diseases:

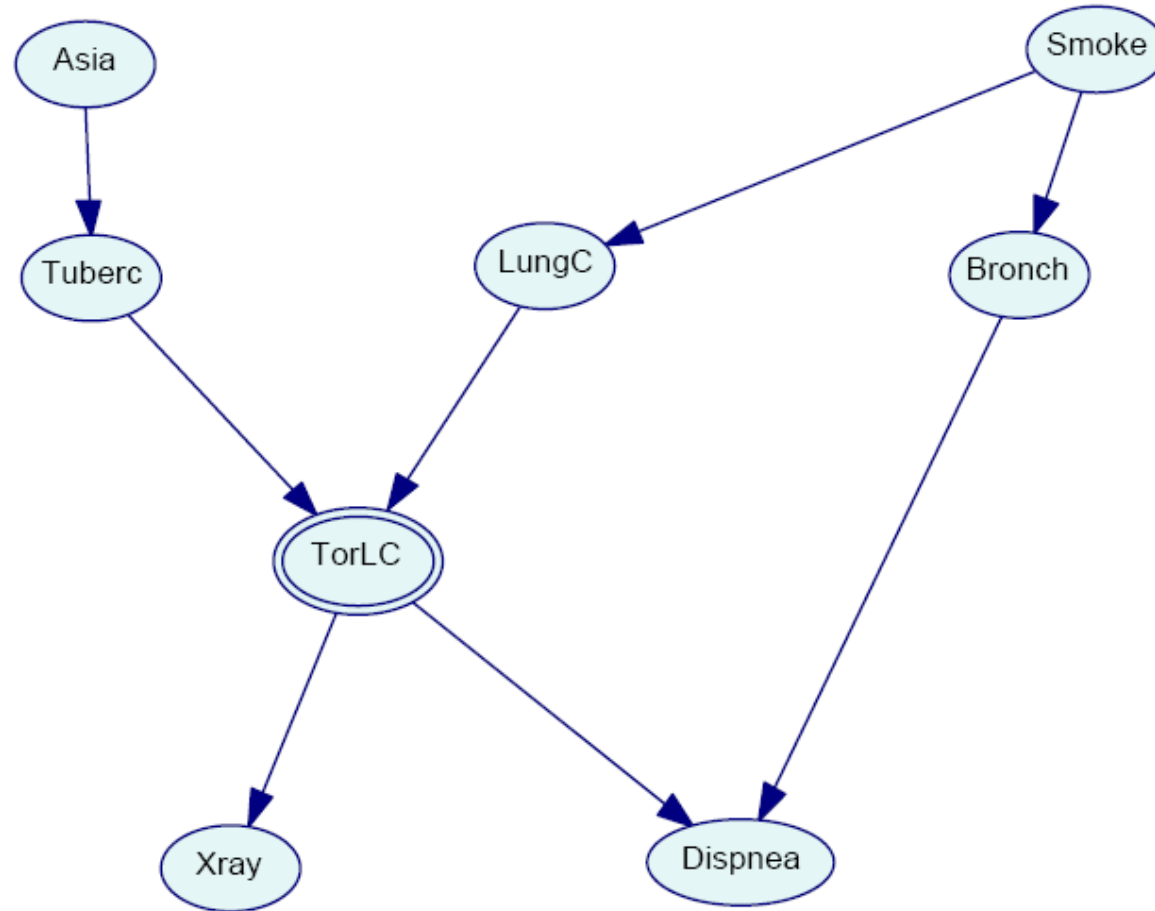
A breathing condition (dyspnea) **may be due** to tuberculosis, lung cancer or bronchitis, none of them or several of them. A recent visit to Asia, **increases the chances** of tuberculosis, whereas smoking is a **risk factor** for lung cancer and bronchitis. The results of an X-ray **may not discriminate** between cancer and tuberculosis, as neither the presence or absence of dyspnea does.

# Probabilistic diagrams

An example referring to lung diseases

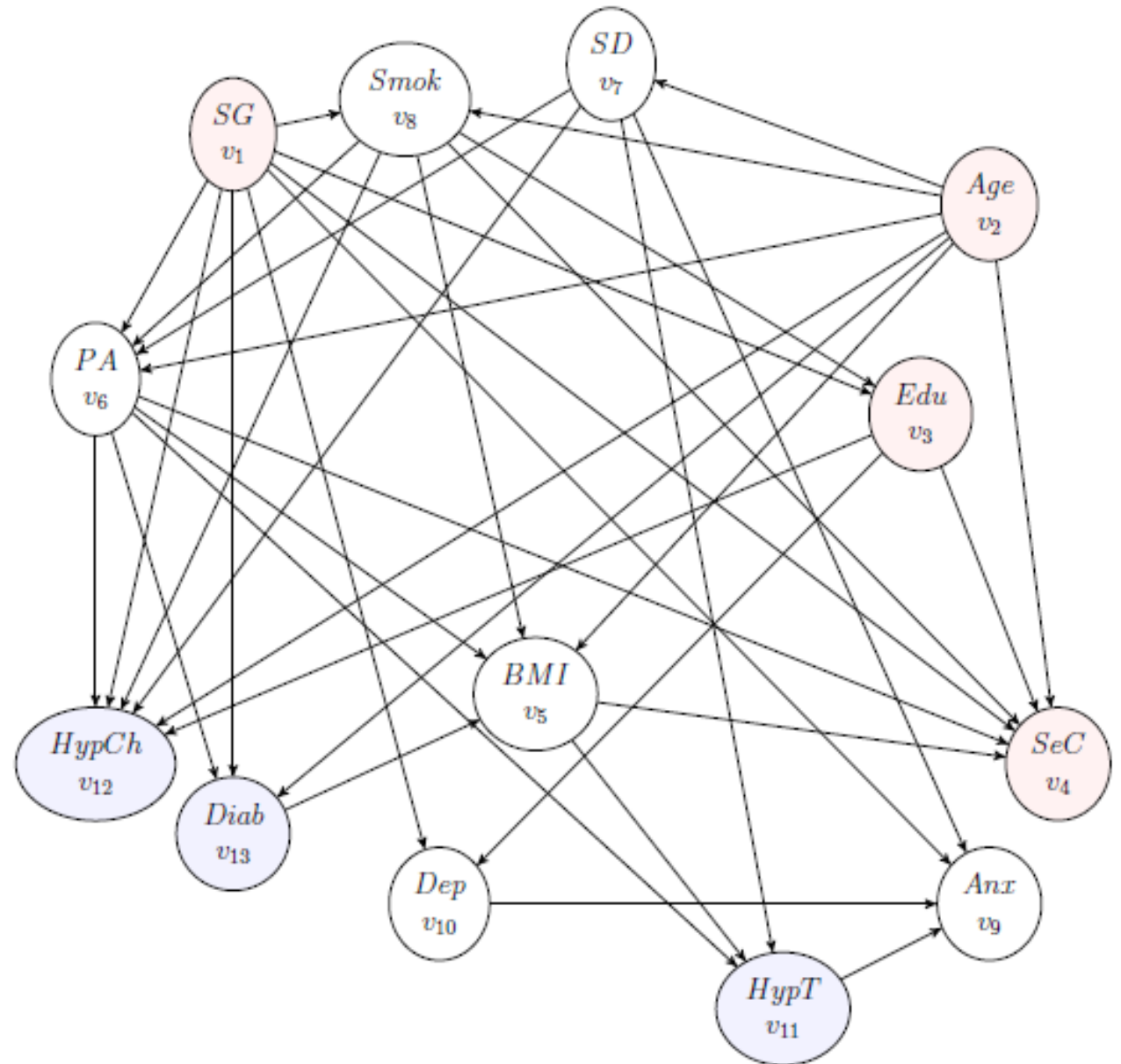
A breathing condition (dyspnea) may be due to tuberculosis, lung cancer or bronchitis, none of them or several of them. A recent visit to Asia, increases the chances of tuberculosis, whereas smoking is a risk factor for lung cancer and bronchitis. The results of an X-ray may not discriminate between cancer and tuberculosis, as neither the presence or absence of dyspnea does.

# Probabilistic diagrams



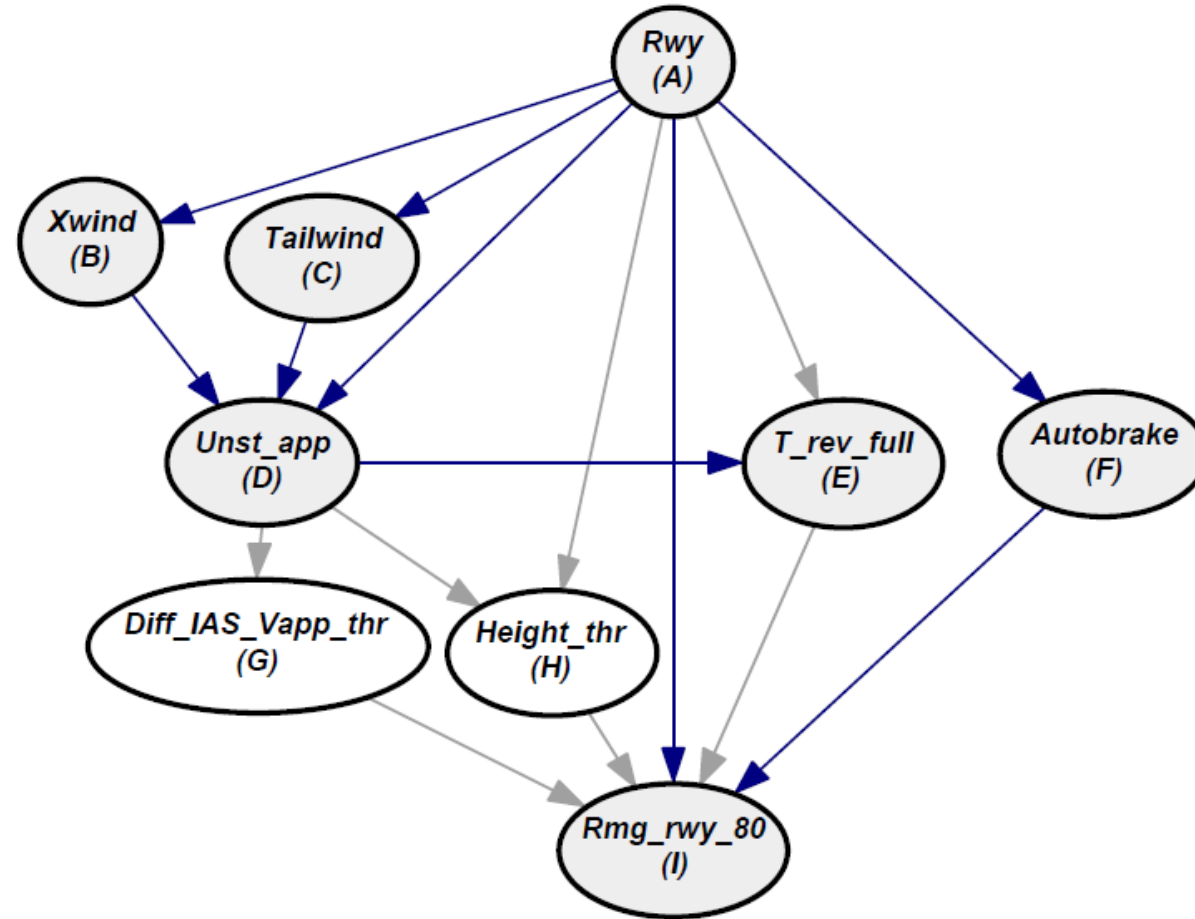
$$P(A,T,S,L,B,O,X,D) = P(A)P(T|A)P(S)P(L|S)P(B|S)P(O|T,L)P(X|O)P(D|O,B)$$

# Hypertension



Build the probabilistic model

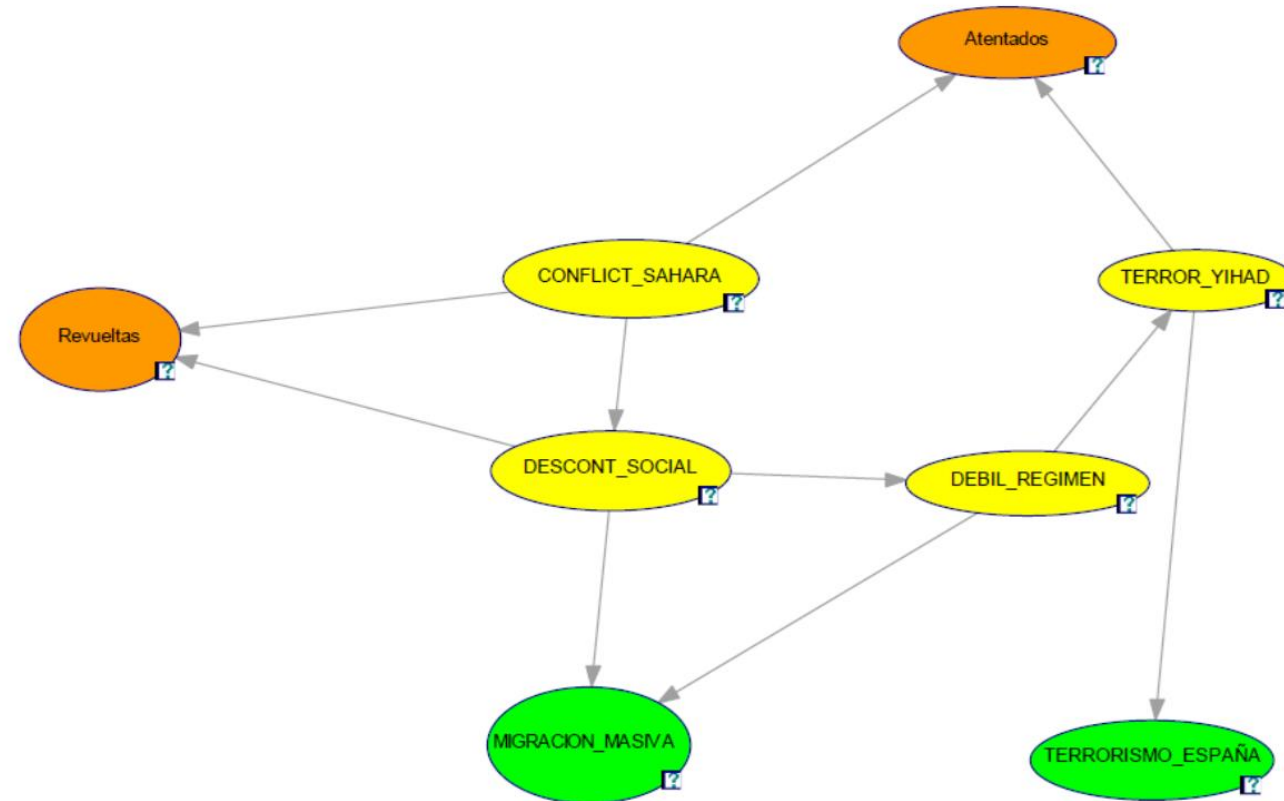
# Runway excursions at airports



Build the probabilistic model

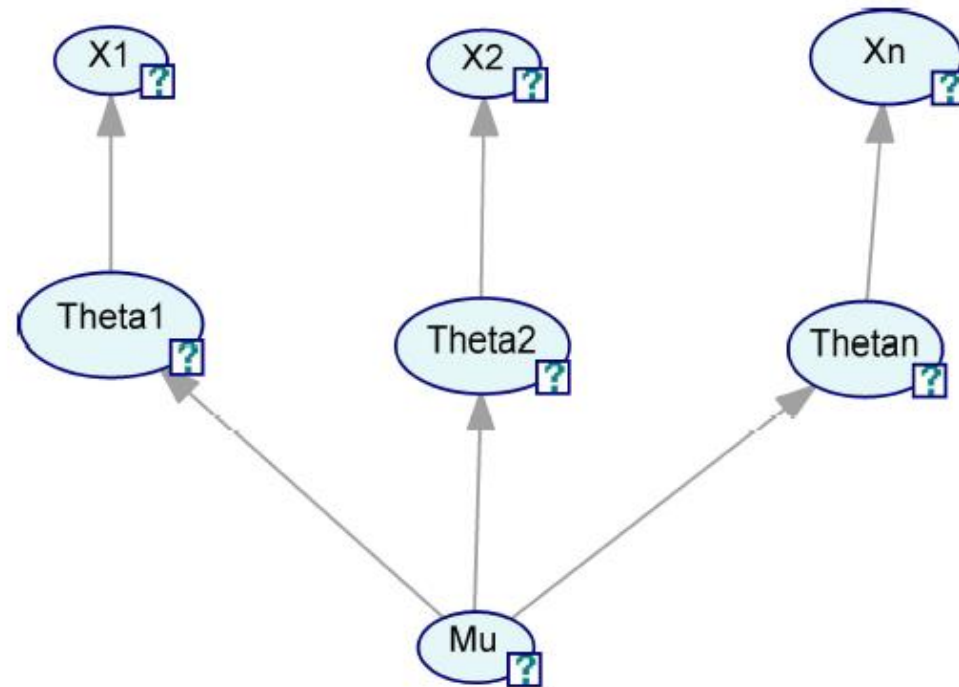


# National security

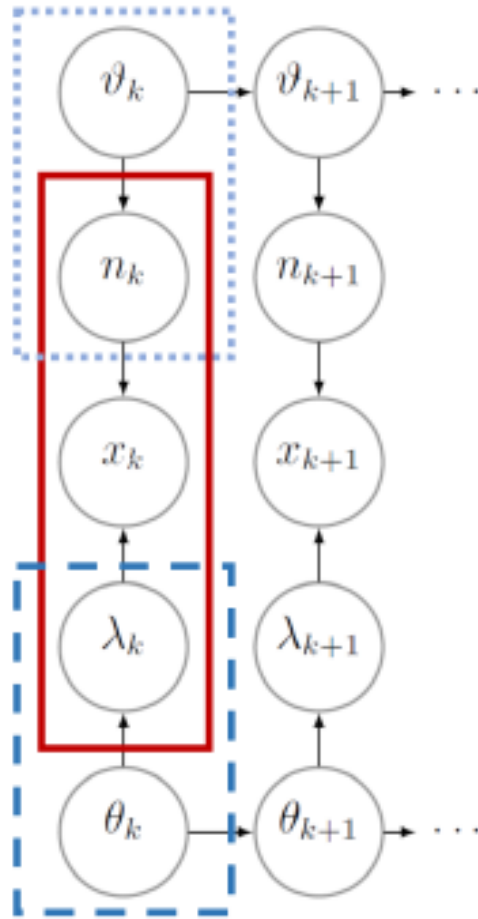


Build the probabilistic model

# Statistical models as PGMs. Hierarchical models



# National aviation safety plan



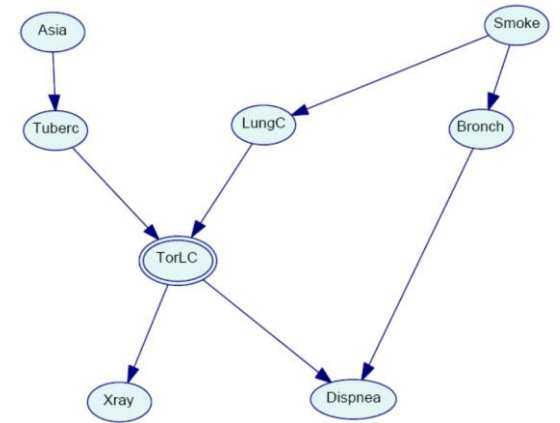
$$\left\{ \begin{array}{l} \left\{ \begin{array}{l} n_k = \mathbf{H}_k \vartheta_k + z_k, \quad z_k \sim N(\mathbf{0}, \Sigma_k) \\ \vartheta_k = \mathbf{J}_k \vartheta_{k-1} + \xi_k, \quad \xi_k \sim N(\mathbf{0}, \mathbf{S}_k) \end{array} \right. \\ \vartheta_0 \sim N(\eta_0, \mathbf{S}_0) \\ x_k | \lambda_k, n_k \sim Po(\lambda_k n_k), \quad \lambda_k = \exp(u_k) \\ \left\{ \begin{array}{l} u_k = \mathbf{F}_k \theta_k + v_k, \quad v_k \sim N(0, V_k) \\ \theta_k = \mathbf{G}_k \theta_{k-1} + w_k, \quad w_k \sim N(\mathbf{0}, \mathbf{W}_k) \end{array} \right. \\ \theta_0 \sim N(m_0, \mathbf{C}_0), \end{array} \right.$$

# Inference in graphical models

# General problem

Assuming DAG (arcs and distributions at nodes):

1. Initialisation
2. Absorption of evidence
3. Global propagation of evidence
4. Hypothesising and propagating single pieces of evidence
5. Planning
6. Influential findings

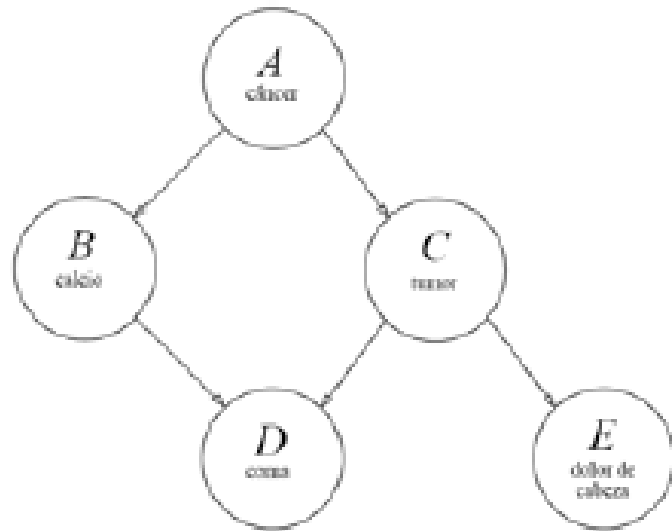


# Gibbs sampler for belief nets

## Conditionals

$$P(X_j = x_j | X_{-j} = x_{-j}) = \alpha P(X_j = x_j | \Pi_{X_j}(x_{-j})) \prod_{Y_j \in \mathcal{S}_j} P(Y_j = y_j | \Pi_{Y_j}(x_j))$$

# Back to example



$\alpha$	0.2
----------	-----

	$a$	$\bar{a}$
$b$	0.8	0.2

	$a$	$\bar{a}$
$c$	0.2	0.05

	$b, c$	$\bar{b}, c$	$b, \bar{c}$	$\bar{b}, \bar{c}$
$d$	0.8	0.8	0.8	0.05

	$c$	$\bar{c}$
$e$	0.8	0.6

$$P(c|\bar{d}, e) = \frac{P(c, \bar{d}, e)}{P(\bar{d}, e)}$$

$$\begin{aligned}
 P(c, \bar{d}, e) &= \sum_{\alpha, \beta} P(\alpha, \beta, c, \bar{d}, e) = \sum_{\alpha, \beta} P(\alpha)P(\beta|\alpha)P(c|\alpha)P(\bar{d}|\beta, c)P(e|c) \\
 &= P(a)P(b|a)P(c|a)P(\bar{d}|b, c)P(e|c) + P(a)P(\bar{b}|a)P(c|a)P(\bar{d}|\bar{b}, c)P(e|c) + \\
 &\quad P(\bar{a})P(b|\bar{a})P(c|\bar{a})P(\bar{d}|b, c)P(e|c) + P(\bar{a})P(\bar{b}|\bar{a})P(c|\bar{a})P(\bar{d}|\bar{b}, c)P(e|c) \\
 &= 0.0118
 \end{aligned}$$

$$P(\bar{d}, e) = \sum_{\alpha, \beta, \gamma} P(\alpha, \beta, \gamma, \bar{d}, e) = 0.410$$

$$P(A, B, C, D, E) = P(A)P(B|A)P(C|A)P(D|B, C)P(E|C)$$

$$P(c|\bar{d}, e) = 0.0287$$

# Back to example

$$P(A|B, C, \bar{d}, e) = P(A|x_{-A}) = \alpha_1 P(A)P(B|A)P(C|A)$$

$$P(B|A, C, \bar{d}, e) = P(B|x_{-B}) = \alpha_2 P(B|A)P(\bar{d}|B, C)$$

$$P(C|A, B, \bar{d}, e) = P(C|x_{-C}) = \alpha_3 P(C|A)P(\bar{d}|B, C)P(e|C)$$

Seleccionar  $B = b_0, C = c_0$  arbitrariamente

Hacer  $j = 1$

Mientras no se juzgue convergencia,

Generar  $A_j = a_j \sim P(A|x_{-A}) = \alpha_{1j} P(A)P(b_{j-1}|A)P(c_{j-1}|A)$

Generar  $B_j = b_j \sim P(B|x_{-B}) = \alpha_{2j} P(B|a_j)P(\bar{d}|B, c_{j-1})$

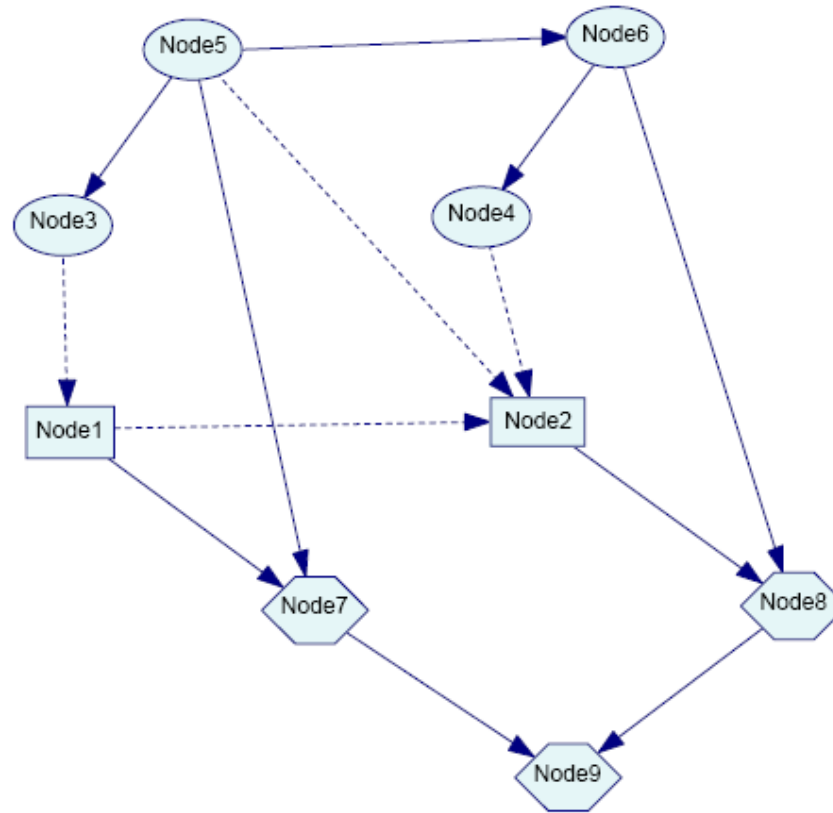
Generar  $C_j = c_j \sim P(C|x_{-C}) = \alpha_{3j} P(C|a_j)P(\bar{d}|b_j, C)P(e|C)$

Hacer  $j = j + 1$

$$\frac{\#\{C_j = c\}}{M}$$



# Sequential Decisions



Learning structure from data: **Structure learning**. Greedy search based on a scoring function based on an information measure

Learning node distributions....

Deep belief nets

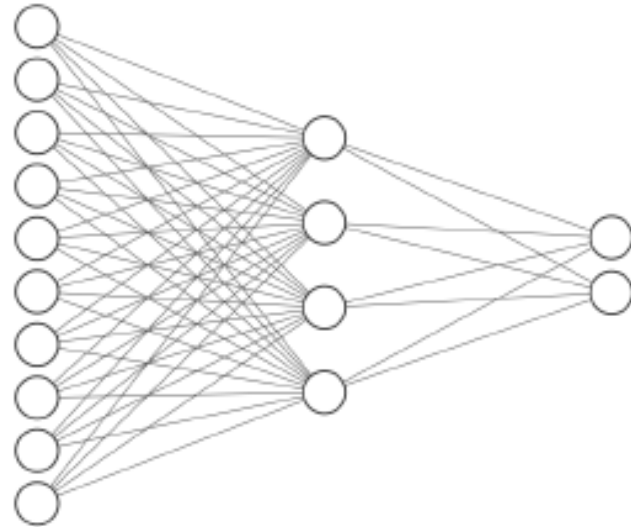
GeNIe

<https://www.bayesfusion.com/influence-diagrams/>

<https://download.bayesfusion.com/files.html?category=Academia>

# Shallow neural nets

# Formulation



Input Layer  $\in \mathbb{R}^{10}$

Hidden Layer  $\in \mathbb{R}^4$

Output Layer  $\in \mathbb{R}^2$

$$y = \sum_{j=1}^m \beta_j \psi(x' \gamma_j) + \epsilon$$

$$\epsilon \sim N(0, \sigma^2),$$

$$\psi(\eta) = \exp(\eta) / (1 + \exp(\eta))$$

Linear in beta's, nonlinear in gamma's

# Training

Given training data, maximise log-likelihood

$$\min_{\beta, \gamma} f(\beta, \gamma) = \sum_{i=1}^n f_i(\beta, \gamma) = \sum_{i=1}^n \left( y_i - \sum_{j=1}^m \beta_j \psi(x_i' \gamma_j) \right)^2$$

Gradient descent

Backpropagation to estimate gradient

# Training with regularisation

$$\min_{\beta, \gamma} f(\beta, \gamma) = \sum_{i=1}^n f_i(\beta, \gamma) = \sum_{i=1}^n \left( y_i - \sum_{j=1}^m \beta_j \psi(x_i' \gamma_j) \right)^2$$

$$\min g(\beta, \gamma) = f(\beta, \gamma) + h(\beta, \gamma),$$

Weight decay

$$h(\beta, \gamma) = \lambda_1 \sum \beta_i^2 + \lambda_2 \sum \sum \gamma_{ji}^2$$

Ridge

# Bayesian analysis of shallow neural nets (fixed arch)

$$y = \sum_{j=1}^m \beta_j \psi(\mathbf{x}'\gamma_j) + \epsilon$$

$$\epsilon \sim N(0, \sigma^2),$$

$$\psi(\eta) = \exp(\eta)/(1 + \exp(\eta))$$

$$\beta_i \sim N(\mu_\beta, \sigma_\beta^2) \text{ and } \gamma_i \sim N(\mu_\gamma, S_\gamma^2)$$

$$\mu_\beta \sim N(a_\beta, A_\beta), \mu_\gamma \sim N(a_\gamma, A_\gamma), \sigma_\beta^{-2} \sim \text{Gamma}(c_b/2, c_b C_b/2)$$

$$S_\gamma^{-1} \sim \text{Wish}(c_\gamma, (c_\gamma C_\gamma)^{-1}) \text{ and } \sigma^{-2} \sim \text{Gamma}(s/2, sS/2)$$

# Bayesian analysis of shallow neural nets (fixed arch)

```
1 Start with arbitrary  $(\beta, \gamma, \nu)$ .
2 while not convergence do
3   Given current  $(\gamma, \nu)$ , draw  $\beta$  from  $p(\beta|\gamma, \nu, y)$  (a multivariate normal).
4   for  $j = 1, \dots, m$ , marginalizing in  $\beta$  and given  $\nu$  do
5     Generate a candidate  $\tilde{\gamma}_j \sim g_j(\gamma_j)$ .
6     Compute  $a(\gamma_j, \tilde{\gamma}_j) = \min\left(1, \frac{p(D|\tilde{\gamma}, \nu)}{p(D|\gamma, \nu)}\right)$  with  $\tilde{\gamma} = (\gamma_1, \gamma_2, \dots, \tilde{\gamma}_j, \dots, \gamma_m)$ .
7     With probability  $a(\gamma_j, \tilde{\gamma}_j)$  replace  $\gamma_j$  by  $\tilde{\gamma}_j$ . If not, preserve  $\gamma_j$ .
8   end
9   Given  $\beta$  and  $\gamma$ , replace  $\nu$  based on their posterior conditionals:
10   $p(\mu_\beta|\beta, \sigma_\beta)$  is normal;  $p(\mu_\gamma|\gamma, S_\gamma)$ , multivariate normal;  $p(\sigma_\beta^{-2}|\beta, \mu_\beta)$ ,
    Gamma;  $p(S_\gamma^{-1}|\gamma, \mu_\gamma)$ , Wishart;  $p(\sigma^{-2}|\beta, \gamma, y)$ , Gamma.
11 end
```

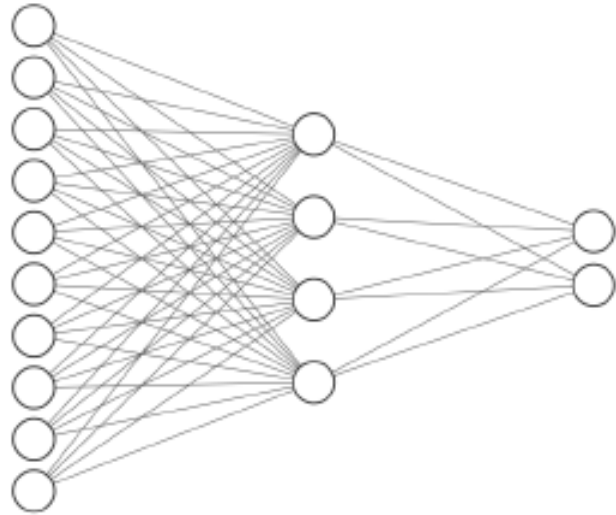


# Bayesian analysis of shallow neural nets (var arch)

$$y = x_i' a + \sum_{j=1}^{m^*} d_j \beta_j \psi(x' \gamma_j) + \epsilon$$
$$\epsilon \sim N(0, \sigma^2),$$
$$\psi(\eta) = \exp(\eta) / (1 + \exp(\eta)),$$
$$Pr(d_j = k | d_{j-1} = 1) = (1 - \alpha)^{1-k} \times \alpha^k, k \in \{0, 1\}$$
$$\beta_i \sim N(\mu_b, \sigma_\beta^2), a \sim N(\mu_a, \sigma_a^2), \gamma_i \sim N(\mu_\gamma, \Sigma_\gamma).$$

Reversible jump algo

# Concept



Input Layer  $\in \mathbb{R}^8$

Hidden Layer  $\in \mathbb{R}^4$

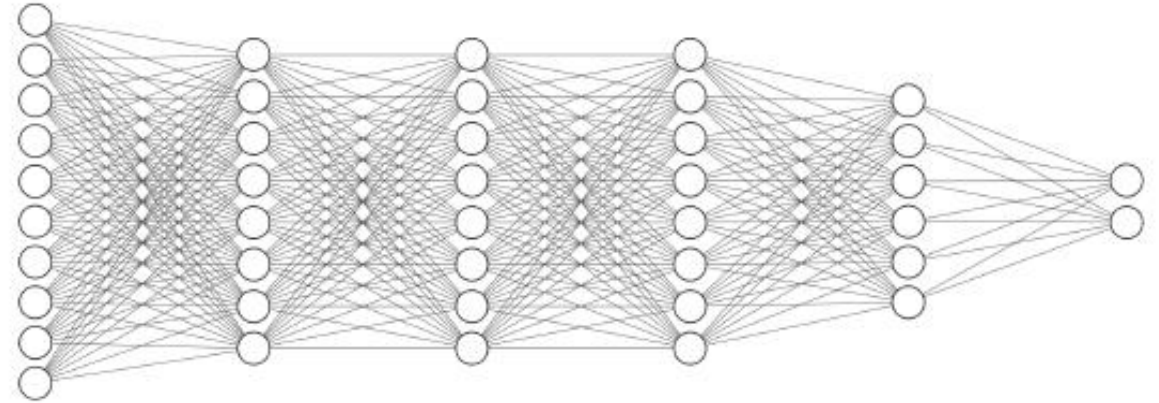
Output Layer  $\in \mathbb{R}^2$

$$y = \sum_{j=1}^m \beta_j \psi(x' \gamma_j) + \epsilon$$

$$\epsilon \sim N(0, \sigma^2),$$

$$\psi(\eta) = \exp(\eta) / (1 + \exp(\eta))$$

(Shallow) Neural nets



$$\{f_0, f_1, \dots, f_{L-1}\}$$

$$z_{l+1} = f_l(z_l, \gamma_l).$$

$$y = \sum_{j=1}^{m_L} \beta_j z_{L,j} + \epsilon$$

$$\epsilon \sim N(0, \sigma^2),$$

Deep neural nets

# Bayesian analysis of deep neural nets

```
1 Start with arbitrary  $(\beta, \gamma, \nu)$ .
2 while not convergence do
3   Given current  $(\gamma, \nu)$ , draw  $\beta$  from  $p(\beta|\gamma, \nu, y)$  (a multivariate normal).
4   for  $j = 1, \dots, m$ , marginalizing in  $\beta$  and given  $\nu$  do
5     Generate a candidate  $\tilde{\gamma}_j \sim g_j(\gamma_j)$ .
6     Compute  $a(\gamma_j, \tilde{\gamma}_j) = \min\left(1, \frac{p(D|\tilde{\gamma}, \nu)}{p(D|\gamma, \nu)}\right)$  with  $\tilde{\gamma} = (\gamma_1, \gamma_2, \dots, \tilde{\gamma}_j, \dots, \gamma_m)$ .
7     With probability  $a(\gamma_j, \tilde{\gamma}_j)$  replace  $\gamma_j$  by  $\tilde{\gamma}_j$ . If not, preserve  $\gamma_j$ .
8   end
9   Given  $\beta$  and  $\gamma$ , replace  $\nu$  based on their posterior conditionals:
10   $p(\mu_\beta|\beta, \sigma_\beta)$  is normal;  $p(\mu_\gamma|\gamma, S_\gamma)$ , multivariate normal;  $p(\sigma_\beta^{-2}|\beta, \mu_\beta)$ ,
    Gamma;  $p(S_\gamma^{-1}|\gamma, \mu_\gamma)$ , Wishart;  $p(\sigma^{-2}|\beta, \gamma, y)$ , Gamma.
11 end
```