

# **Unsupervised Learning**

## **Bayesian Model Comparison**

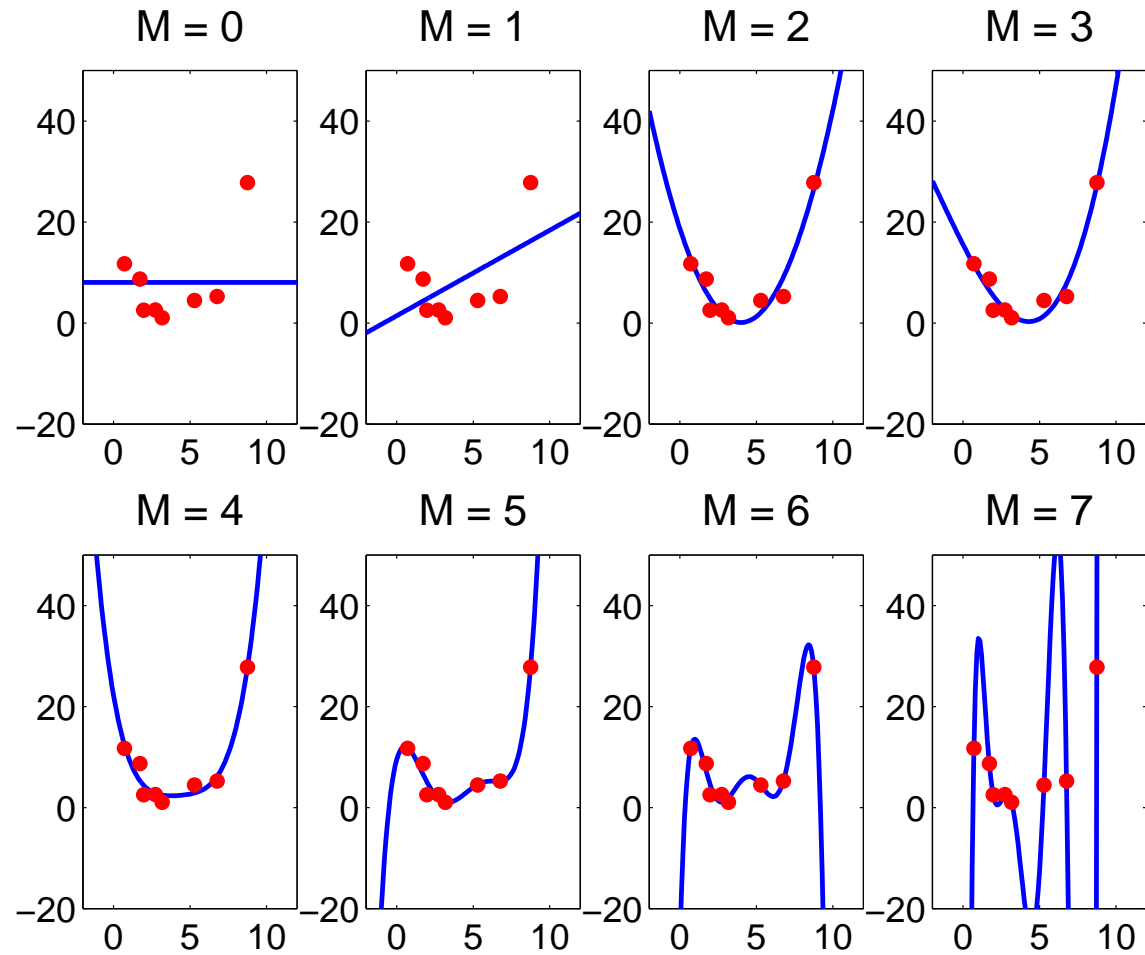
**Zoubin Ghahramani**

`zoubin@gatsby.ucl.ac.uk`

**Gatsby Computational Neuroscience Unit, and  
MSc in Intelligent Systems, Dept Computer Science  
University College London**

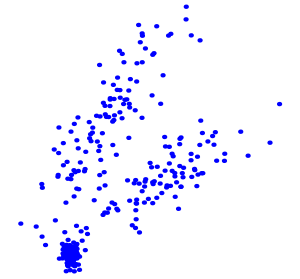
**Term 1, Autumn 2005**

# Model complexity and overfitting: a simple example

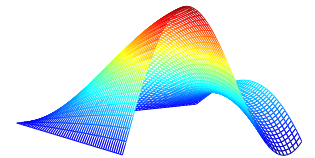


# Learning Model Structure

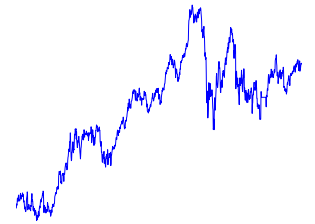
How many clusters in the data?



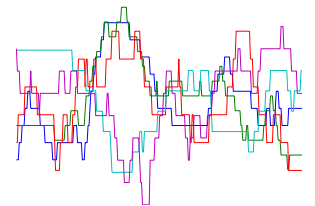
What is the intrinsic dimensionality of the data?



Is this input relevant to predicting that output?



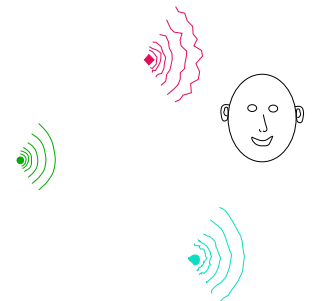
What is the order of a dynamical system?



How many states in a hidden Markov model?

SVYDAAAQLTADVKKDLRDSWKVIGSDKKGNGVALMTTY

How many auditory sources in the input?



# Using Occam's Razor to Learn Model Structure

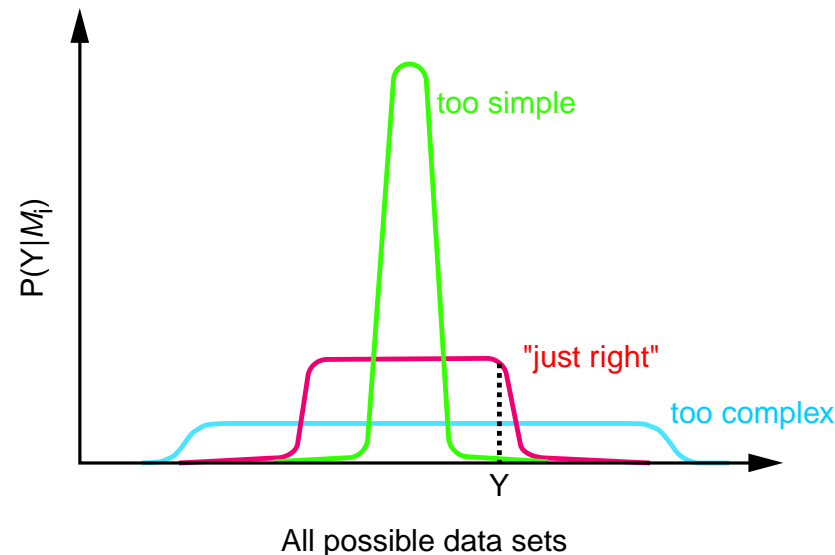
Compare model classes  $m$  using their posterior probability given the data:

$$P(m|\mathbf{y}) = \frac{P(\mathbf{y}|m)P(m)}{P(\mathbf{y})}, \quad P(\mathbf{y}|m) = \int_{\Theta_m} P(\mathbf{y}|\boldsymbol{\theta}_m, m)P(\boldsymbol{\theta}_m|m) d\boldsymbol{\theta}_m$$

**Interpretation of  $P(\mathbf{y}|m)$ :** The probability that *randomly selected* parameter values from the model class would generate data set  $\mathbf{y}$ .

Model classes that are **too simple** are unlikely to generate the data set.

Model classes that are **too complex** can generate many possible data sets, so again, they are unlikely to generate that particular data set at random.



# Bayesian Model Comparison: Terminology

- A **model class**  $m$  is a set of models parameterised by  $\theta_m$ , e.g. the set of all possible mixtures of  $m$  Gaussians.
- The **marginal likelihood** of model class  $m$ :

$$P(\mathbf{y}|m) = \int_{\Theta_m} P(\mathbf{y}|\theta_m, m)P(\theta_m|m) d\theta_m$$

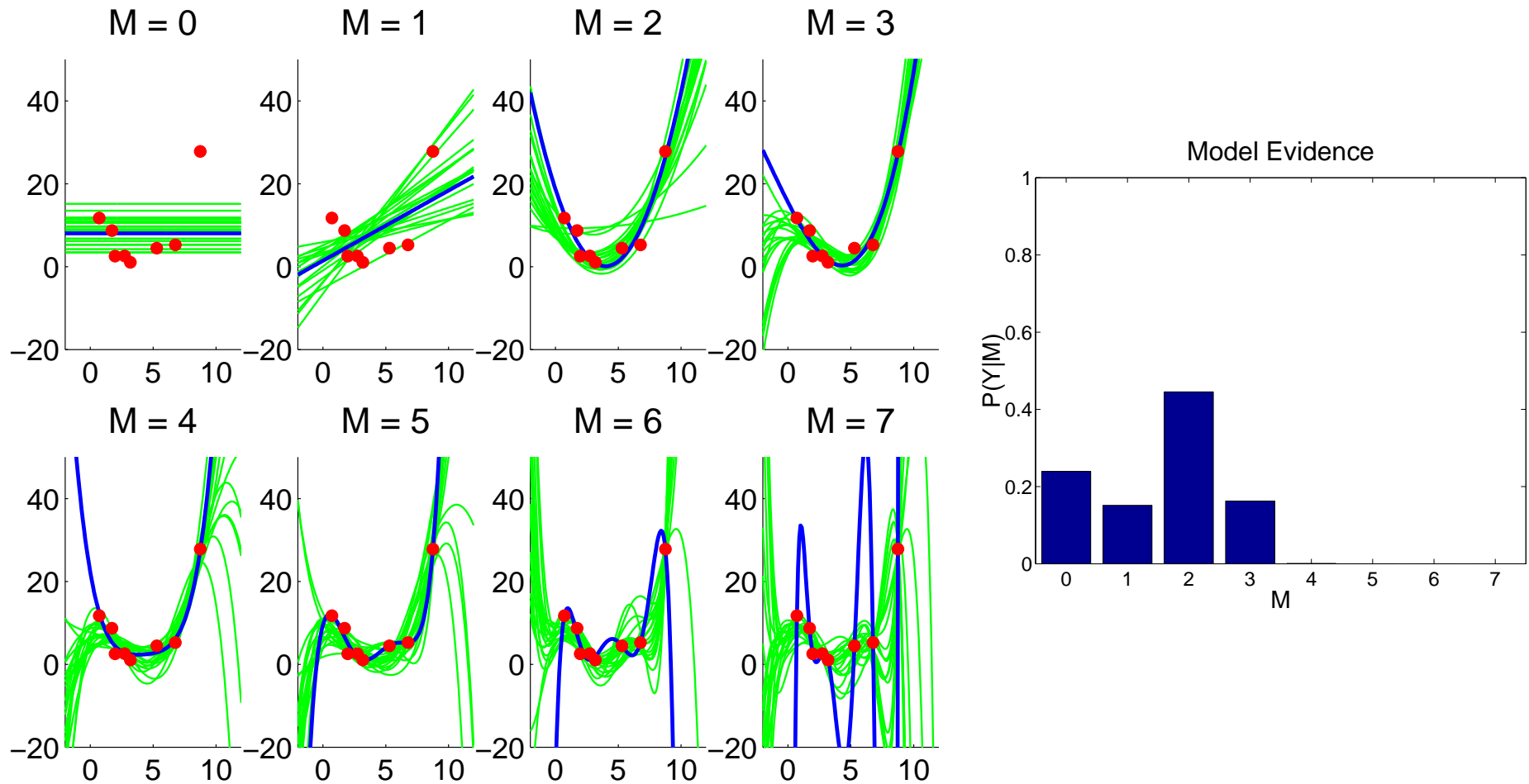
is also known as the **Bayesian evidence** for model  $m$ .

- The ratio of two marginal likelihoods is known as the **Bayes factor**:

$$\frac{P(\mathbf{y}|m)}{P(\mathbf{y}|m')}$$

- The **Occam's Razor** principle is, roughly speaking, that one should prefer simpler explanations than more complex explanations.
- Bayesian inference formalises and *automatically* implements the Occam's Razor principle.

# Bayesian Model Comparison: Occam's Razor at Work



e.g. for quadratic ( $M=2$ ):  $y = a_0 + a_1x + a_2x^2 + \epsilon$ , where  $\epsilon \sim \mathcal{N}(0, \tau)$  and  $\theta_2 = [a_0 \ a_1 \ a_2 \ \tau]$   
demo: polybayes

# Practical Bayesian approaches

- Laplace approximations:
  - Makes a Gaussian approximation about the maximum *a posteriori* parameter estimate.
- Bayesian Information Criterion (BIC)
  - an asymptotic approximation.
- Markov chain Monte Carlo methods (MCMC):
  - In the limit are guaranteed to converge, but:
  - Many samples required to ensure accuracy.
  - Sometimes hard to assess convergence.
- Variational approximations

Note: other deterministic approximations have been developed more recently and can be applied in this context: e.g. Bethe approximations and Expectation Propagation

# Laplace Approximation

data set:  $\mathbf{y}$       models:  $m = 1 \dots, M$       parameter sets:  $\boldsymbol{\theta}_1 \dots, \boldsymbol{\theta}_M$

Model Comparison:  $P(m|\mathbf{y}) \propto P(m)P(\mathbf{y}|m)$

For large amounts of data (relative to number of parameters,  $d$ ) the parameter posterior is approximately Gaussian around the MAP estimate  $\hat{\boldsymbol{\theta}}_m$ :

$$P(\boldsymbol{\theta}_m|\mathbf{y}, m) \approx (2\pi)^{-\frac{d}{2}} |A|^{\frac{1}{2}} \exp \left\{ -\frac{1}{2} (\boldsymbol{\theta}_m - \hat{\boldsymbol{\theta}}_m)^\top A (\boldsymbol{\theta}_m - \hat{\boldsymbol{\theta}}_m) \right\}$$

$$P(\mathbf{y}|m) = \frac{P(\boldsymbol{\theta}_m, \mathbf{y}|m)}{P(\boldsymbol{\theta}_m|\mathbf{y}, m)}$$

Evaluating the above for  $\ln P(\mathbf{y}|m)$  at  $\hat{\boldsymbol{\theta}}_m$  we get the Laplace approximation:

$$\ln P(\mathbf{y}|m) \approx \ln P(\hat{\boldsymbol{\theta}}_m|m) + \ln P(\mathbf{y}|\hat{\boldsymbol{\theta}}_m, m) + \frac{d}{2} \ln 2\pi - \frac{1}{2} \ln |A|$$

$-A$  is the  $d \times d$  Hessian matrix of  $\log P(\boldsymbol{\theta}_m|\mathbf{y}, m)$ :  $A_{kl} = -\frac{\partial^2}{\partial \theta_{mk} \partial \theta_{ml}} \ln P(\boldsymbol{\theta}_m|\mathbf{y}, m) |_{\hat{\boldsymbol{\theta}}_m}$ .

Can also be derived from  $2^{nd}$  order Taylor expansion of log posterior.

The Laplace approximation can be used for model comparison.



# Bayesian Information Criterion (BIC)

BIC can be obtained from the Laplace approximation:

$$\ln P(\mathbf{y}|m) \approx \ln P(\hat{\boldsymbol{\theta}}_m|m) + \ln P(\mathbf{y}|\hat{\boldsymbol{\theta}}_m, m) + \frac{d}{2} \ln 2\pi - \frac{1}{2} \ln |A|$$

in the large sample limit ( $N \rightarrow \infty$ ) where  $N$  is the number of data points,  $A$  grows as  $NA_0$  for some fixed matrix  $A_0$ , so  $\ln |A| \rightarrow \ln |NA_0| = \ln(N^d |A_0|) = d \ln N + \ln |A_0|$ . Retaining only terms that grow in  $N$  we get:

$$\ln P(\mathbf{y}|m) \approx \ln P(\mathbf{y}|\hat{\boldsymbol{\theta}}_m, m) - \frac{d}{2} \ln N$$

Properties:

- Quick and easy to compute
- It does not depend on the prior
- We can use the ML estimate of  $\theta$  instead of the MAP estimate
- It is equivalent to the “Minimum Description Length” (MDL) criterion
- It assumes that in the large sample limit, all the parameters are well-determined (i.e. the model is **identifiable**; otherwise,  $d$  should be the number of **well-determined** parameters)
- **Danger:** counting parameters can be deceiving! (c.f. sinusoid, infinite models)

# Sampling Approximations

Let's consider a non-Markov chain method, **Importance Sampling**:

$$\begin{aligned}\ln P(\mathbf{y}|m) &= \ln \int_{\Theta_m} P(\mathbf{y}|\boldsymbol{\theta}_m, m) P(\boldsymbol{\theta}_m|m) d\boldsymbol{\theta}_m \\ &= \ln \int_{\Theta_m} P(\mathbf{y}|\boldsymbol{\theta}_m, m) \frac{P(\boldsymbol{\theta}_m|m)}{Q(\boldsymbol{\theta}_m)} Q(\boldsymbol{\theta}_m) d\boldsymbol{\theta}_m \\ &\approx \ln \frac{1}{K} \sum_k P(\mathbf{y}|\boldsymbol{\theta}_m^{(k)}, m) \frac{P(\boldsymbol{\theta}_m^{(k)}|m)}{Q(\boldsymbol{\theta}_m^{(k)})}\end{aligned}$$

where  $\boldsymbol{\theta}_m^{(k)}$  are i.i.d. draws from  $Q(\boldsymbol{\theta}_m)$ . Assumes we can **sample from** and **evaluate**  $Q(\boldsymbol{\theta}_m)$  (incl. normalization!) and we can **compute the likelihood**  $P(\mathbf{y}|\boldsymbol{\theta}_m^{(k)}, m)$ .

Although importance sampling does not work well in high dimensions, it inspires the following approach: Create a **Markov chain**,  $Q_k \rightarrow Q_{k+1} \dots$  for which:

- $Q_k(\boldsymbol{\theta})$  can be evaluated including normalization
- $\lim_{k \rightarrow \infty} Q_k(\boldsymbol{\theta}) = P(\boldsymbol{\theta}|\mathbf{y}, m)$

# Variational Bayesian Learning

## Lower Bounding the Marginal Likelihood

Let the hidden latent variables be  $\mathbf{x}$ , data  $\mathbf{y}$  and the parameters  $\boldsymbol{\theta}$ .

Lower bound the marginal likelihood (Bayesian model evidence) using Jensen's inequality:

$$\begin{aligned}\ln P(\mathbf{y}) &= \ln \int d\mathbf{x} d\boldsymbol{\theta} P(\mathbf{y}, \mathbf{x}, \boldsymbol{\theta}) && \text{||}m \\ &= \ln \int d\mathbf{x} d\boldsymbol{\theta} Q(\mathbf{x}, \boldsymbol{\theta}) \frac{P(\mathbf{y}, \mathbf{x}, \boldsymbol{\theta})}{Q(\mathbf{x}, \boldsymbol{\theta})} \\ &\geq \int d\mathbf{x} d\boldsymbol{\theta} Q(\mathbf{x}, \boldsymbol{\theta}) \ln \frac{P(\mathbf{y}, \mathbf{x}, \boldsymbol{\theta})}{Q(\mathbf{x}, \boldsymbol{\theta})}.\end{aligned}$$

Use a simpler, factorised approximation to  $Q(\mathbf{x}, \boldsymbol{\theta})$ :

$$\begin{aligned}\ln P(\mathbf{y}) &\geq \int d\mathbf{x} d\boldsymbol{\theta} Q_{\mathbf{x}}(\mathbf{x}) Q_{\boldsymbol{\theta}}(\boldsymbol{\theta}) \ln \frac{P(\mathbf{y}, \mathbf{x}, \boldsymbol{\theta})}{Q_{\mathbf{x}}(\mathbf{x}) Q_{\boldsymbol{\theta}}(\boldsymbol{\theta})} \\ &= \mathcal{F}(Q_{\mathbf{x}}(\mathbf{x}), Q_{\boldsymbol{\theta}}(\boldsymbol{\theta}), \mathbf{y}).\end{aligned}$$

Maximize this lower bound.

# Variational Bayesian Learning . . .

Maximizing this **lower bound**,  $\mathcal{F}$ , leads to **EM-like** updates:

$$Q_{\mathbf{x}}^*(\mathbf{x}) \propto \exp \langle \ln P(\mathbf{x}, \mathbf{y} | \boldsymbol{\theta}) \rangle_{Q_{\boldsymbol{\theta}}(\boldsymbol{\theta})} \quad E\text{-like step}$$

$$Q_{\boldsymbol{\theta}}^*(\boldsymbol{\theta}) \propto P(\boldsymbol{\theta}) \exp \langle \ln P(\mathbf{x}, \mathbf{y} | \boldsymbol{\theta}) \rangle_{Q_{\mathbf{x}}(\mathbf{x})} \quad M\text{-like step}$$

Maximizing  $\mathcal{F}$  is equivalent to minimizing KL-divergence between the *approximate posterior*,  $Q(\boldsymbol{\theta})Q(\mathbf{x})$  and the *true posterior*,  $P(\boldsymbol{\theta}, \mathbf{x} | \mathbf{y})$ .

$$\begin{aligned} \ln P(\mathbf{y}) - \mathcal{F}(Q_{\mathbf{x}}(\mathbf{x}), Q_{\boldsymbol{\theta}}(\boldsymbol{\theta}), \mathbf{y}) &= \\ \ln P(\mathbf{y}) - \int d\mathbf{x} d\boldsymbol{\theta} Q_{\mathbf{x}}(\mathbf{x}) Q_{\boldsymbol{\theta}}(\boldsymbol{\theta}) \ln \frac{P(\mathbf{y}, \mathbf{x}, \boldsymbol{\theta})}{Q_{\mathbf{x}}(\mathbf{x}) Q_{\boldsymbol{\theta}}(\boldsymbol{\theta})} &= \\ \int d\mathbf{x} d\boldsymbol{\theta} Q_{\mathbf{x}}(\mathbf{x}) Q_{\boldsymbol{\theta}}(\boldsymbol{\theta}) \ln \frac{Q_{\mathbf{x}}(\mathbf{x}) Q_{\boldsymbol{\theta}}(\boldsymbol{\theta})}{P(\mathbf{x}, \boldsymbol{\theta} | \mathbf{y})} &= KL(Q || P) \end{aligned}$$

# Conjugate-Exponential models

Let's focus on *conjugate-exponential* (CE) models, which satisfy (1) and (2):

- **Condition (1)**. The joint probability over variables is in the exponential family:

$$P(\mathbf{x}, \mathbf{y}|\boldsymbol{\theta}) = f(\mathbf{x}, \mathbf{y}) g(\boldsymbol{\theta}) \exp \{ \boldsymbol{\phi}(\boldsymbol{\theta})^\top \mathbf{u}(\mathbf{x}, \mathbf{y}) \}$$

where  $\boldsymbol{\phi}(\boldsymbol{\theta})$  is the vector of *natural parameters*,  $\mathbf{u}$  are *sufficient statistics*

- **Condition (2)**. The prior over parameters is conjugate to this joint probability:

$$P(\boldsymbol{\theta}|\eta, \boldsymbol{\nu}) = h(\eta, \boldsymbol{\nu}) g(\boldsymbol{\theta})^\eta \exp \{ \boldsymbol{\phi}(\boldsymbol{\theta})^\top \boldsymbol{\nu} \}$$

where  $\eta$  and  $\boldsymbol{\nu}$  are hyperparameters of the prior.

Conjugate priors are computationally convenient and have an intuitive interpretation:

- $\eta$ : number of pseudo-observations
- $\boldsymbol{\nu}$ : values of pseudo-observations

# Conjugate-Exponential examples

In the **CE** family:

- Gaussian mixtures
- factor analysis, probabilistic PCA
- hidden Markov models and factorial HMMs
- linear dynamical systems and switching models
- discrete-variable belief networks

Other as yet undreamt-of models can combine Gaussian, Gamma, Poisson, Dirichlet, Wishart, Multinomial and others.

Not in the **CE** family:

- Boltzmann machines, MRFs (no simple conjugacy)
- logistic regression (no simple conjugacy)
- sigmoid belief networks (not exponential)
- independent components analysis (not exponential)

Note: one can often approximate these models with models in the **CE** family.

## A Useful Result

Given an iid data set  $\mathbf{y} = (\mathbf{y}_1, \dots, \mathbf{y}_n)$ , if the model is **CE** then:

(a)  $Q_{\boldsymbol{\theta}}(\boldsymbol{\theta})$  is also **conjugate**, *i.e.*

$$Q_{\boldsymbol{\theta}}(\boldsymbol{\theta}) = h(\tilde{\eta}, \tilde{\boldsymbol{\nu}}) g(\boldsymbol{\theta})^{\tilde{\eta}} \exp \{ \boldsymbol{\phi}(\boldsymbol{\theta})^\top \tilde{\boldsymbol{\nu}} \}$$

where  $\tilde{\eta} = \eta + n$  and  $\tilde{\boldsymbol{\nu}} = \boldsymbol{\nu} + \sum_i \bar{\mathbf{u}}(\mathbf{x}_i, \mathbf{y}_i)$ .

(b)  $Q_{\mathbf{x}}(\mathbf{x}) = \prod_{i=1}^n Q_{\mathbf{x}_i}(\mathbf{x}_i)$  is of the **same form** as in the E step of regular EM, but using **pseudo parameters** computed by averaging over  $Q_{\boldsymbol{\theta}}(\boldsymbol{\theta})$

$$Q_{\mathbf{x}_i}(\mathbf{x}_i) \propto f(\mathbf{x}_i, \mathbf{y}_i) \exp \{ \bar{\boldsymbol{\phi}}(\boldsymbol{\theta})^\top \mathbf{u}(\mathbf{x}_i, \mathbf{y}_i) \} = P(\mathbf{x}_i | \mathbf{y}_i, \bar{\boldsymbol{\phi}}(\boldsymbol{\theta}))$$

### KEY points:

(a) the approximate parameter posterior is of the same form as the prior, so it is **easily summarized** in terms of two sets of hyperparameters,  $\tilde{\eta}$  and  $\tilde{\boldsymbol{\nu}}$ ;

(b) the approximate hidden variable posterior, *averaging over all parameters*, is of the same form as the hidden variable posterior for a *single setting of the parameters*, so again, it is **easily computed** using the usual methods.

# The Variational Bayesian EM algorithm

## EM for MAP estimation

**Goal:** maximize  $p(\boldsymbol{\theta}|\mathbf{y}, m)$  w.r.t.  $\boldsymbol{\theta}$

**E Step:** compute

$$q_{\mathbf{x}}^{(t+1)}(\mathbf{x}) = p(\mathbf{x}|\mathbf{y}, \boldsymbol{\theta}^{(t)})$$

**M Step:**

$$\boldsymbol{\theta}^{(t+1)} = \operatorname{argmax}_{\boldsymbol{\theta}} \int q_{\mathbf{x}}^{(t+1)}(\mathbf{x}) \ln p(\mathbf{x}, \mathbf{y}, \boldsymbol{\theta}) d\mathbf{x}$$

## Variational Bayesian EM

**Goal:** lower bound  $p(\mathbf{y}|m)$

**VB-E Step:** compute

$$q_{\mathbf{x}}^{(t+1)}(\mathbf{x}) = p(\mathbf{x}|\mathbf{y}, \bar{\boldsymbol{\phi}}^{(t)})$$

**VB-M Step:**

$$q_{\boldsymbol{\theta}}^{(t+1)}(\boldsymbol{\theta}) = \exp \left[ \int q_{\mathbf{x}}^{(t+1)}(\mathbf{x}) \ln p(\mathbf{x}, \mathbf{y}, \boldsymbol{\theta}) d\mathbf{x} \right]$$

## Properties:

- Reduces to the EM algorithm if  $q_{\boldsymbol{\theta}}(\boldsymbol{\theta}) = \delta(\boldsymbol{\theta} - \boldsymbol{\theta}^*)$ .
- $\mathcal{F}_m$  increases monotonically, and incorporates the model complexity penalty.
- Analytical parameter distributions (but not constrained to be Gaussian).
- VB-E step has same complexity as corresponding E step.
- We can use the junction tree, belief propagation, Kalman filter, etc, algorithms in the VB-E step of VB-EM, but **using expected natural parameters**,  $\bar{\boldsymbol{\phi}}$ .



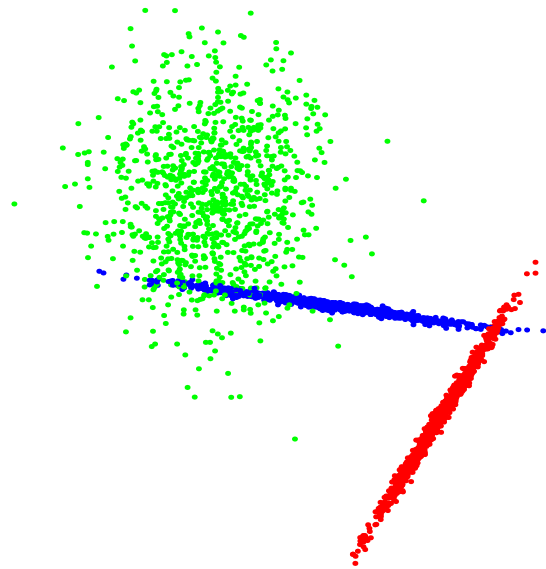
# Variational Bayes: History of Models Treated

- multilayer perceptrons (Hinton & van Camp, 1993)
- mixture of experts (Waterhouse, MacKay & Robinson, 1996)
- hidden Markov models (MacKay, 1995)
- other work by Jaakkola, Jordan, Barber, Bishop, Tipping, etc

## Examples of Variational Learning of Model Structure

- mixtures of factor analysers (Ghahramani & Beal, 1999)
- mixtures of Gaussians (Attias, 1999)
- independent components analysis (Attias, 1999; Miskin & MacKay, 2000; Valpola 2000)
- principal components analysis (Bishop, 1999)
- linear dynamical systems (Ghahramani & Beal, 2000)
- mixture of experts (Ueda & Ghahramani, 2000)
- discrete graphical models (Beal & Ghahramani, 2002)
- **VIBES software** for conjugate-exponential graphs (Winn, 2003)

# Mixture of Factor Analysers



Goal:

- Infer number of clusters
- Infer intrinsic dimensionality of each cluster

Under the assumption that each cluster is Gaussian

embed\_demo

# Mixture of Factor Analysers

True data: 6 Gaussian clusters with dimensions: (1 7 4 3 2 2) embedded in 10-D

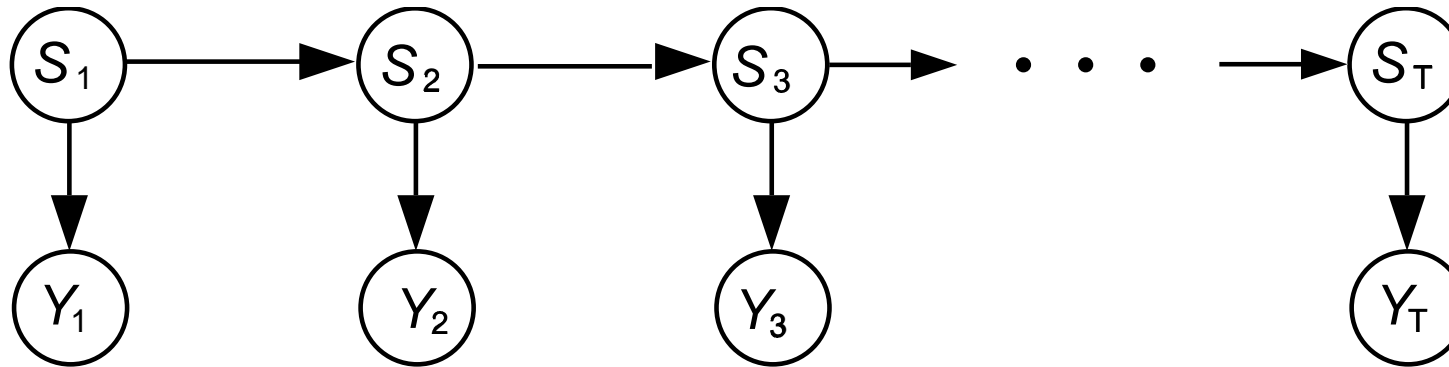
Inferred structure:

number of points per cluster	intrinsic dimensionalities					
	1	7	4	3	2	2
8	2				1	
8	1	2				
16	1	4				2
32	1	6	3	3	2	2
64	1	7	4	3	2	2
128	1	7	4	3	2	2

- Finds the clusters and dimensionalities efficiently.
- The model complexity reduces in line with the lack of data support.

demos: `run_simple` and `ueda_demo`

# Hidden Markov Models



Discrete hidden states,  $s_t$ .

Observations  $y_t$ .

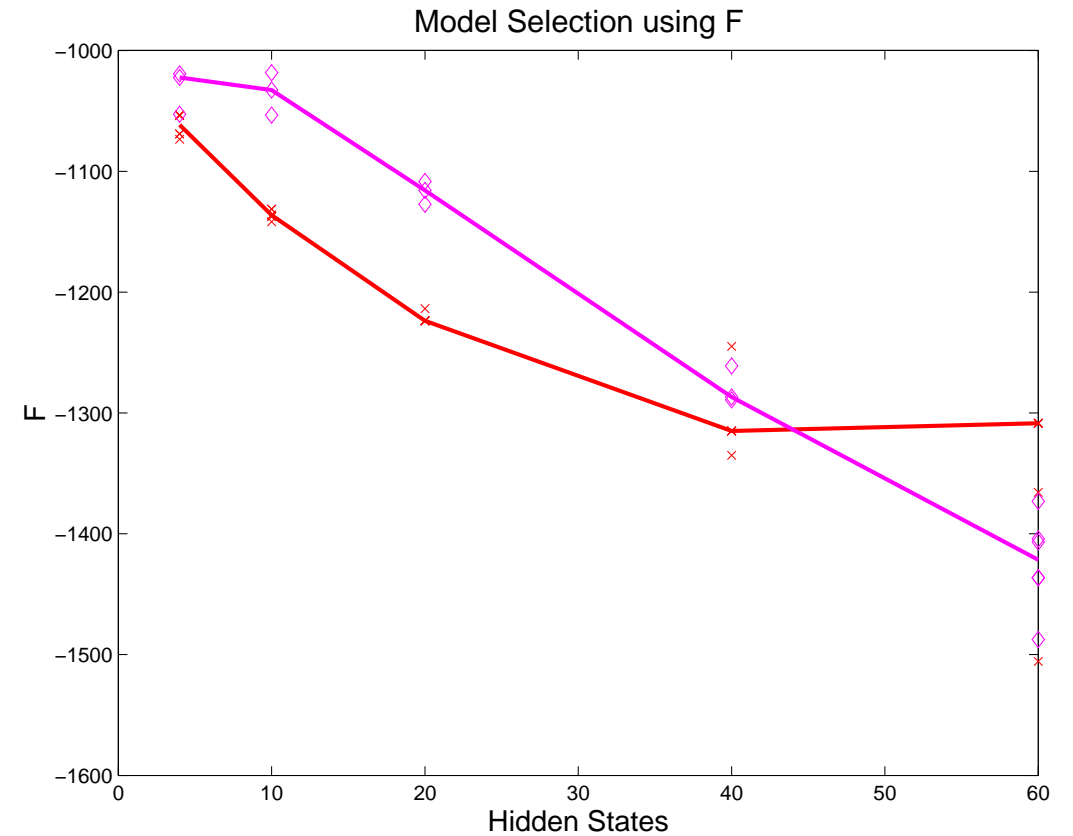
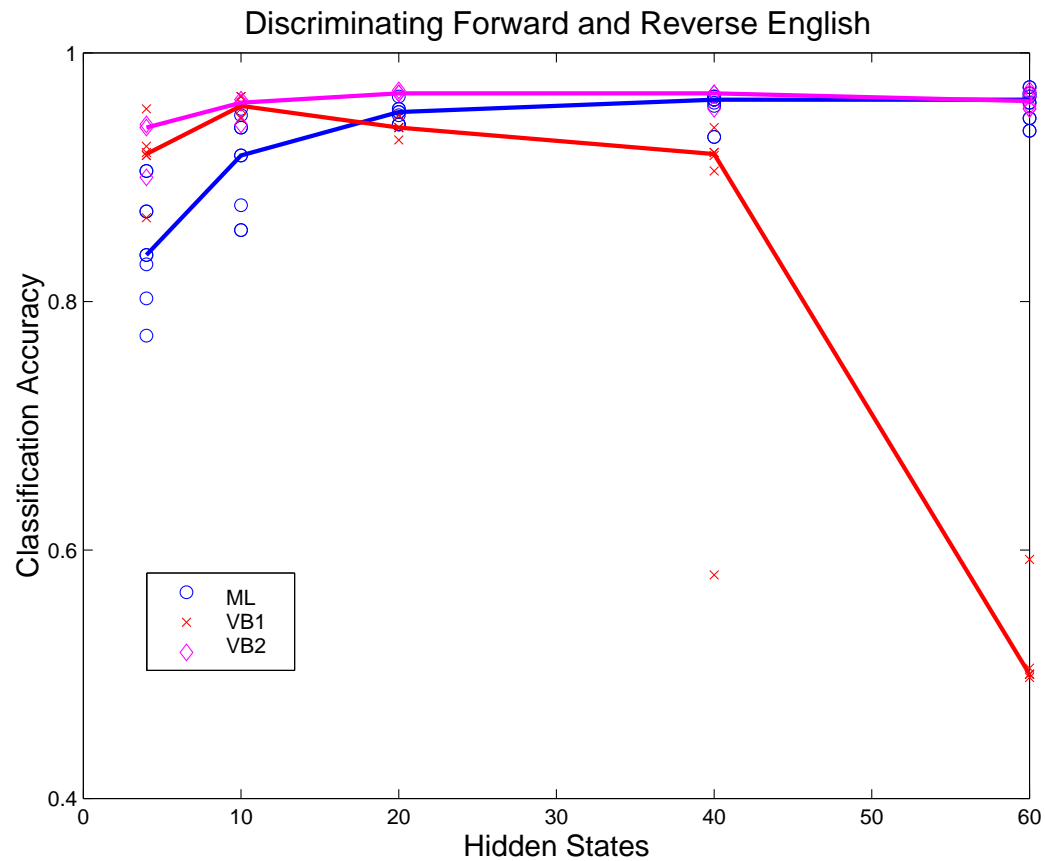
How many hidden states?

What structure state-transition matrix?

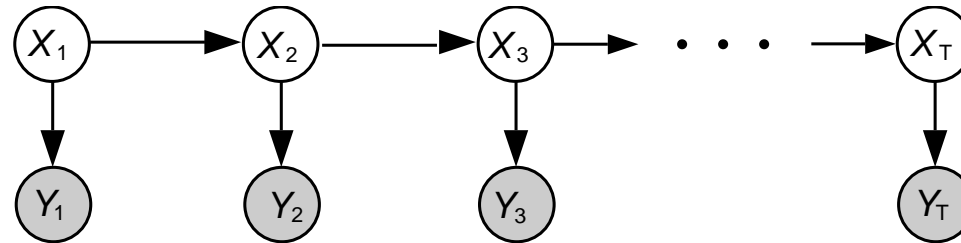
demo: vbhmm\_demo

# Hidden Markov Models: Discriminating Forward from Reverse English

First 8 sentences from Alice in Wonderland.  
Compare VB-HMM with ML-HMM.



# Linear Dynamical Systems



- Assumes  $y_t$  generated from a sequence of Markov *hidden* state variables  $x_t$
- If transition and output functions are linear, time-invariant, and noise distributions are Gaussian, this is a **linear-Gaussian state-space model**:

$$\mathbf{x}_t = A\mathbf{x}_{t-1} + \mathbf{w}_t, \quad \mathbf{y}_t = C\mathbf{x}_t + \mathbf{v}_t$$

- Three levels of inference:
  - I Given data, structure and parameters, **Kalman smoothing**  $\rightarrow$  hidden state;
  - II Given data and structure, **EM**  $\rightarrow$  hidden state and parameter point estimates;
  - III Given data only, **VEM**  $\rightarrow$  **model structure and distributions over parameters and hidden state.**

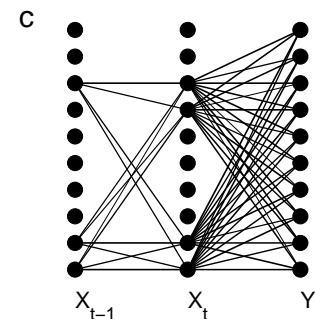
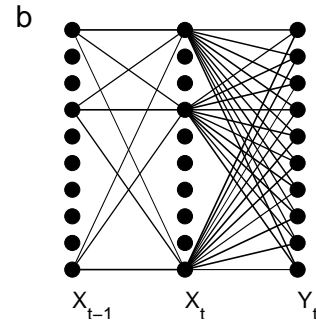
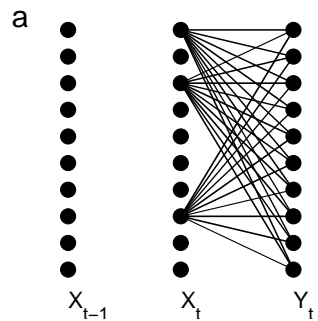
# Linear Dynamical System Results

Inferring model structure:

a) SSM(0,3) i.e. FA

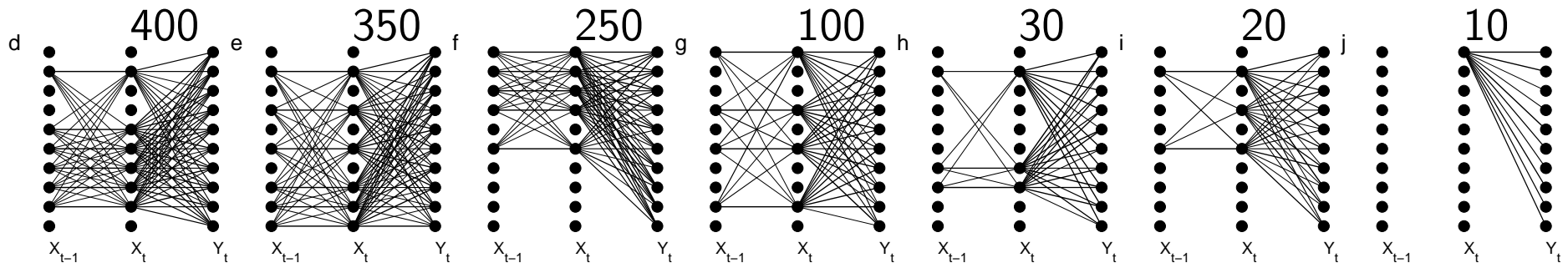
b) SSM(3,3)

c) SSM(3,4)



Inferred model complexity reduces with less data:

True model: SSM(6,6) ● 10-dim observation vector.



demo: bayeslds

# Independent Components Analysis

Blind Source Separation:  $5 \times 100$  msec speech and music sources linearly mixed to produce 11 signals (microphones)

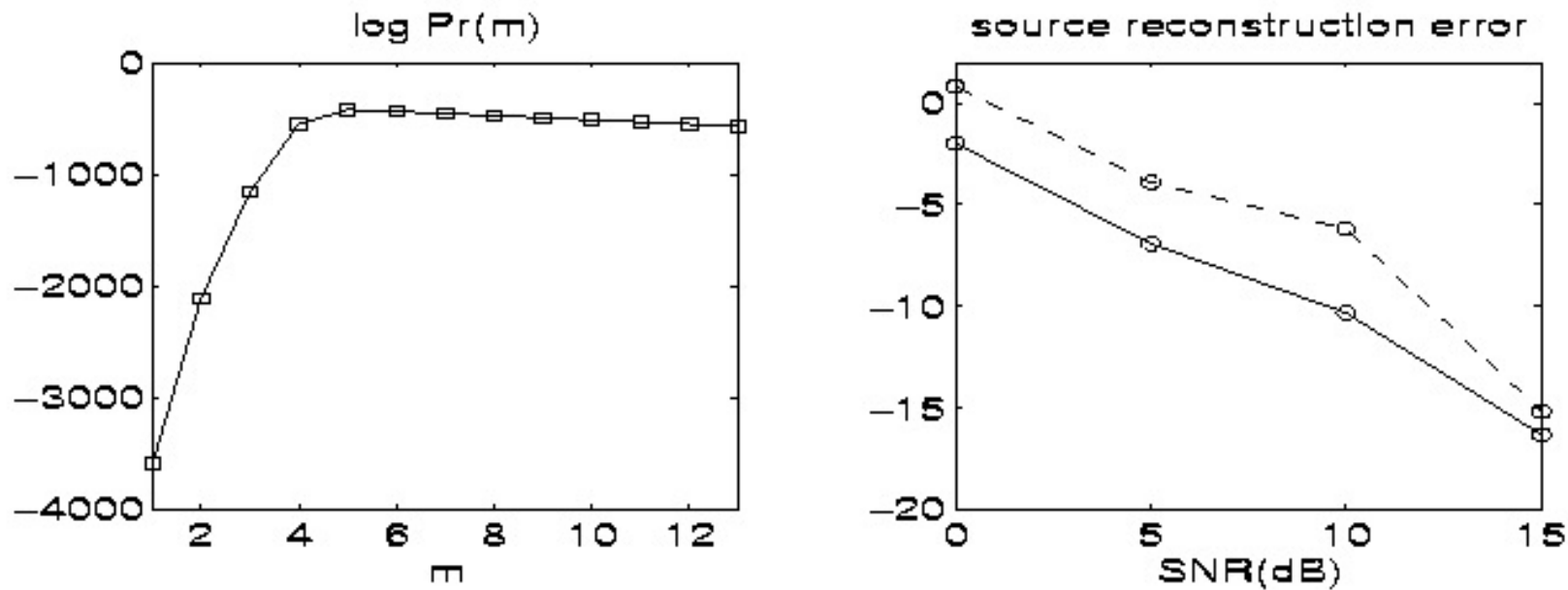


Figure 2. Application of VB to blind source separation algorithm (see text).

from Attias (2000)



# Summary & Conclusions

- Bayesian learning avoids overfitting and can be used to do model comparison / selection.
- But we need approximations:
  - Laplace
  - BIC
  - Sampling
  - Variational

## Other topics I would have liked to cover in the course but did not have time to...

- Nonparametric Bayesian learning, infinite models and Dirichlet Processes
- Bayes net structure learning
- Nonlinear dimensionality reduction
- More on loopy belief propagation, Bethe free energies, and Expectation Propagation
- Exact sampling
- Particle filters
- Semi-supervised learning
- Gaussian processes and SVMs (supervised!)
- Reinforcement Learning