

# Lecture 12: Model Comparison

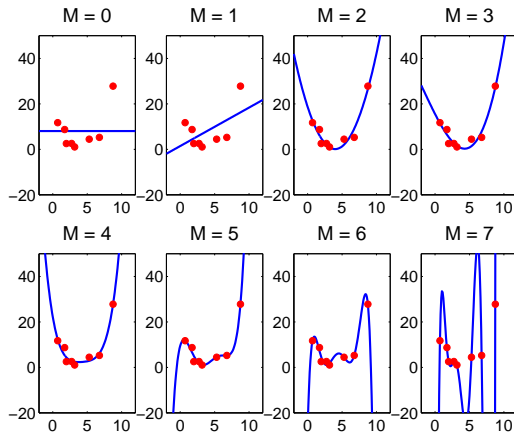
4F13: Machine Learning

Zoubin Ghahramani and Carl Edward Rasmussen

Department of Engineering, University of Cambridge

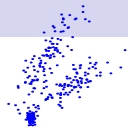
February 24th, 2009

# Model complexity and overfitting: a simple example

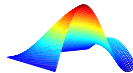


# Learning Model Structure

How many clusters in the data?



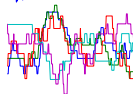
What is the intrinsic dimensionality of the data?



Is this input relevant to predicting that output?



What is the order of a dynamical system?



How many states in a hidden Markov model?

SVYDAAAQLTADVKKDLRDSWKVIGSDKKNGVVALMTTY

How many auditory sources in the input?



# Using Occam's Razor to Learn Model Structure

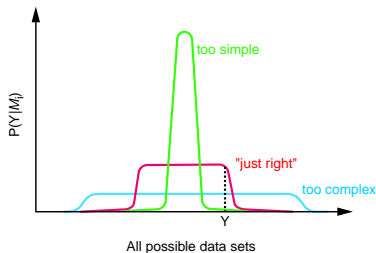
Compare model classes  $m$  using their posterior probability given the data:

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

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

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  $m$ :

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

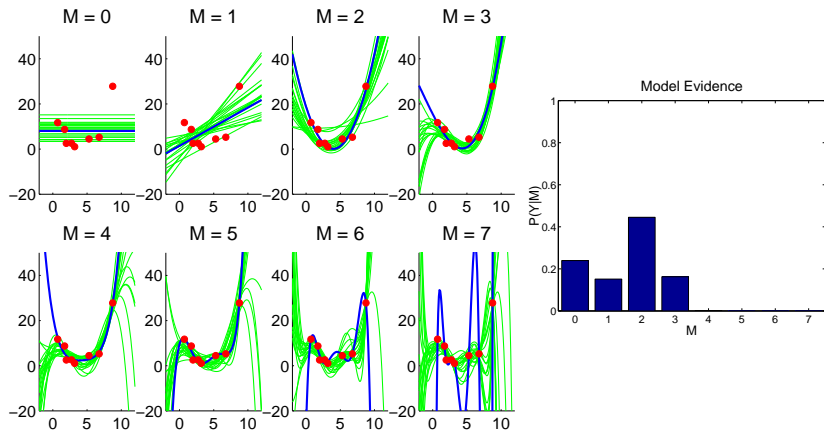
is also known as the **Bayesian evidence** for model  $m$ .

- 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  $\log P(\boldsymbol{\theta}_m|\mathbf{y}, m)$ :

$$A_{k\ell} = -\frac{\partial^2}{\partial \theta_{mk} \partial \theta_{m\ell}} \ln P(\boldsymbol{\theta}_m|\mathbf{y}, m)|_{\hat{\boldsymbol{\theta}}_m}.$$

Can also be derived from 2<sup>nd</sup> 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:

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

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:

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

Properties:

- Quick and easy to compute, and does not depend on the prior
- We can use the ML estimate of  $\theta$  instead of the MAP estimate
- 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)
- It is equivalent to the “Minimum Description Length” (MDL) criterion

# Sampling Approximations

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

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

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**  $P(\mathbf{y}|\boldsymbol{\theta}_m^{(k)}, \mathbf{m})$ .

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

# 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:

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

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:

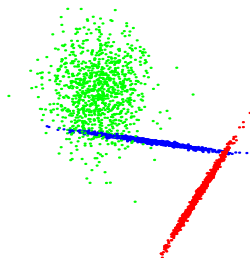
$$Q_{\mathbf{x}}^*(\mathbf{x}) \propto \exp \langle \ln P(\mathbf{x}, \mathbf{y} | \boldsymbol{\theta}) \rangle_{Q_{\boldsymbol{\theta}}(\boldsymbol{\theta})} \quad E\text{-like step}$$

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

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

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

# 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

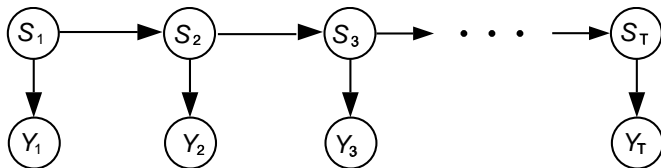
**Inferred structure:**

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,  $s_t$ .

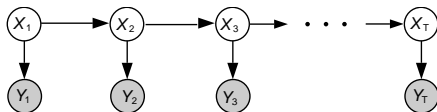
Observations  $y_t$ .

How many hidden states?

What structure state-transition matrix?

demo: `vbhmm_demo`

# 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:
  - Given data, structure and parameters, **Kalman smoothing**  $\rightarrow$  hidden state;
  - Given data and structure, **EM**  $\rightarrow$  hidden state and parameter point estimates;
  - Given data only, **VEM**  $\rightarrow$  **model structure and distributions over parameters and hidden state.**



# Summary & Conclusions

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