

Lecture 12 and 13: Model Comparison

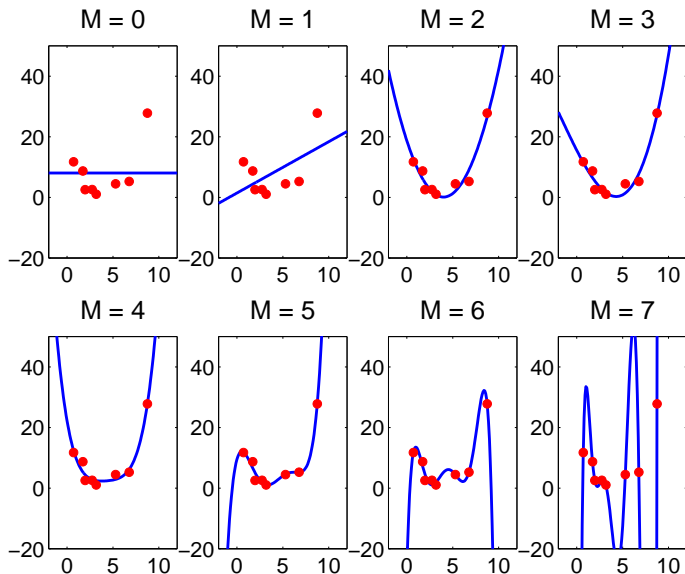
4F13: Machine Learning

Zoubin Ghahramani and Carl Edward Rasmussen

Department of Engineering, University of Cambridge

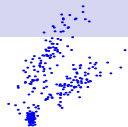
February 24th and 25th, 2010

Model complexity and overfitting: a simple example

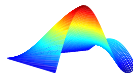


Learning Model Structure

How many clusters in the data?



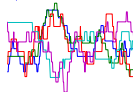
What is the intrinsic dimensionality of the data?



Is this input relevant to predicting that output?



What is the order of a dynamical system?



How many states in a hidden Markov model?

SVYDAAAQLTADVKKDLRDSWKVIGSDKKGNVALMTTY

How many auditory sources in the input?



Using Occam's Razor to Learn Model Structure

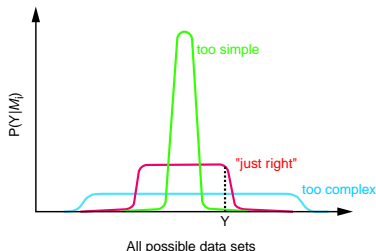
Compare model classes m using their posterior probability given the data:

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

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 θ_m , e.g. the set of all possible mixtures of m Gaussians.
- The **marginal likelihood** of model class m :

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

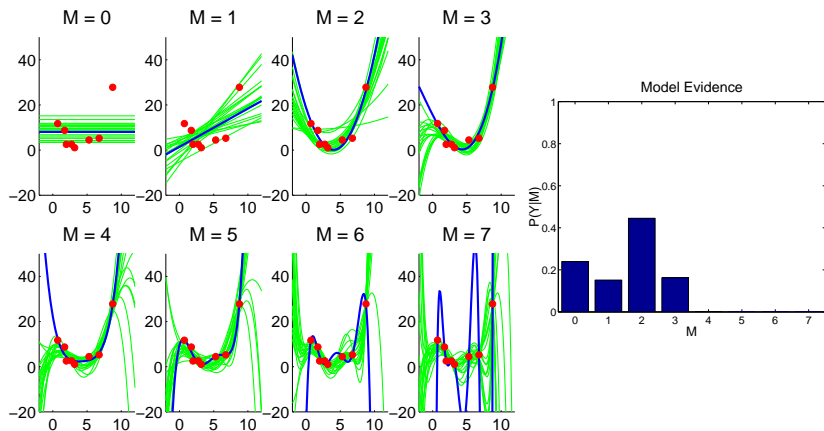
is also known as the **Bayesian evidence** for model m .

- 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 $\hat{\boldsymbol{\theta}}_m$:

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

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

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

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

$-\mathbf{A}$ is the $d \times d$ Hessian matrix of $\log P(\boldsymbol{\theta}_m|\mathbf{y}, m)$:

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

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:

$$\log P(\mathbf{y}|\mathbf{m}) \simeq \log P(\hat{\boldsymbol{\theta}}_{\mathbf{m}}|\mathbf{m}) + \log P(\mathbf{y}|\hat{\boldsymbol{\theta}}_{\mathbf{m}}, \mathbf{m}) + \frac{d}{2} \log 2\pi - \frac{1}{2} \log |\mathbf{A}|$$

in the large sample limit ($N \rightarrow \infty$) where N is the number of data points, \mathbf{A} grows as $N\mathbf{A}_0$ for some fixed matrix \mathbf{A}_0 , so

$\log |\mathbf{A}| \rightarrow \log |N\mathbf{A}_0| = \log(N^d |\mathbf{A}_0|) = d \log N + \log |\mathbf{A}_0|$. Retaining only terms that grow in N we get:

$$\log P(\mathbf{y}|\mathbf{m}) \simeq \log P(\mathbf{y}|\hat{\boldsymbol{\theta}}_{\mathbf{m}}, \mathbf{m}) - \frac{d}{2} \log N$$

Properties:

- Quick and easy to compute, and does not depend on the prior
- We can use the ML estimate of $\boldsymbol{\theta}$ 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**:

$$\begin{aligned}\log P(\mathbf{y}|\mathbf{m}) &= \log \int_{\Theta_{\mathbf{m}}} P(\mathbf{y}|\boldsymbol{\theta}_{\mathbf{m}}, \mathbf{m})P(\boldsymbol{\theta}_{\mathbf{m}}|\mathbf{m}) d\boldsymbol{\theta}_{\mathbf{m}} \\ &= \log \int_{\Theta_{\mathbf{m}}} P(\mathbf{y}|\boldsymbol{\theta}_{\mathbf{m}}, \mathbf{m}) \frac{P(\boldsymbol{\theta}_{\mathbf{m}}|\mathbf{m})}{Q(\boldsymbol{\theta}_{\mathbf{m}})} Q(\boldsymbol{\theta}_{\mathbf{m}}) d\boldsymbol{\theta}_{\mathbf{m}} \\ &\simeq \log \frac{1}{K} \sum_k P(\mathbf{y}|\boldsymbol{\theta}_{\mathbf{m}}^{(k)}, \mathbf{m}) \frac{P(\boldsymbol{\theta}_{\mathbf{m}}^{(k)}|\mathbf{m})}{Q(\boldsymbol{\theta}_{\mathbf{m}}^{(k)})}\end{aligned}$$

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

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 θ .

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

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

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

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

Maximize this lower bound.

Variational Bayesian Learning ...

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

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

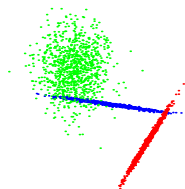
$$Q_\theta^*(\boldsymbol{\theta}) \propto P(\boldsymbol{\theta}) \exp \langle \log P(\mathbf{x}, \mathbf{y} | \boldsymbol{\theta}) \rangle_{Q_x(\mathbf{x})} \quad M\text{-like step}$$

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

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

Mixture of Factor Analysers

Goal: find the *number of clusters* and *their dimensions*.



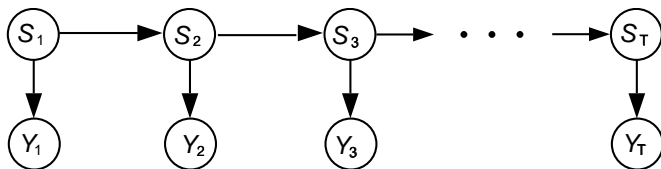
number of points per cluster	intrinsic dimensionalities					
	1	7	4	3	2	2
8	2		1			
8	1	2				
16	1	4				2
32	1	6	3	3	2	2
64	1	7	4	3	2	2
128	1	7	4	3	2	2

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

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

demos: `embed_demo`, `run_simple` and `ueda_demo`

Hidden Markov Models



Discrete hidden states, s_t .

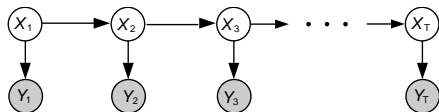
Observations y_t .

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**:

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

- 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.**

Summary & Conclusions

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