

Lecture 13: Gaussian Mixture Models, EM, Model Selection

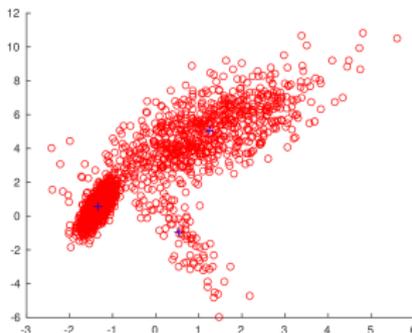
Recommended reading: Bishop, Chapter 1.3, 3.1, 9.2

Duncan Gillies and Marc Deisenroth

Department of Computing
Imperial College London

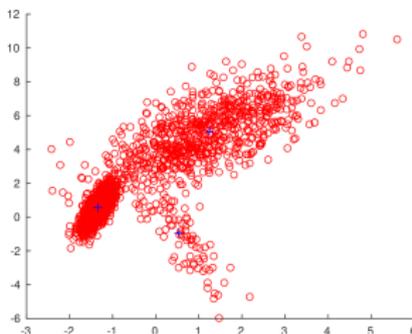
February 17, 2016

Problem Statement



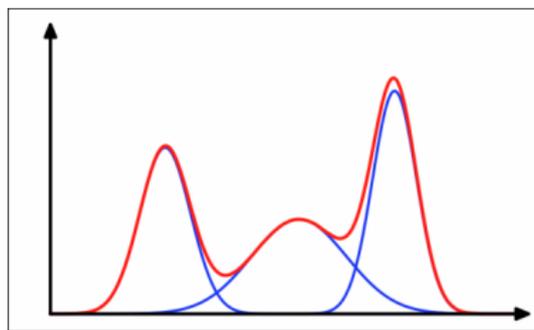
- ▶ Often, we are given a set of points whose density we wish to model
- ▶ Example: Find mean, variance of a Gaussian
 - ▶▶ MLE/MAP estimation

Problem Statement



- ▶ Often, we are given a set of points whose density we wish to model
- ▶ Example: Find mean, variance of a Gaussian
 - ▶▶ MLE/MAP estimation
- ▶ Gaussians (or similarly all other distributions we encountered so far) have very limited modeling capabilities.
 - ▶▶ **Mixture models** are more flexible

Gaussian Mixtures



$$p(\mathbf{x}) = \sum_{k=1}^K \pi_k \mathcal{N}(\mathbf{x} | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)$$

$$0 \leq \pi_k \leq 1$$

$$\sum_k \pi_k = 1$$

- ▶ Individual components are Gaussian distributions
- ▶ Each component is weighted by π_k

Parameter Learning for GMMs

- ▶ Objective: Maximum likelihood estimate of model parameters θ
- ▶ $\theta := \{\pi_k, \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k, k = 1, \dots, K\}$
- ▶ Maximum likelihood estimate:

Parameter Learning for GMMs

- ▶ Objective: Maximum likelihood estimate of model parameters θ
- ▶ $\theta := \{\pi_k, \mu_k, \Sigma_k, k = 1, \dots, K\}$
- ▶ Maximum likelihood estimate:

$$\begin{aligned}\theta^* &= \arg \max_{\theta} p(\mathbf{x}|\theta) = \arg \max_{\theta} \sum_{k=1}^K \pi_k \mathcal{N}(\mathbf{x} | \mu_k, \Sigma_k) \\ &\stackrel{\log}{=} \arg \max_{\theta} \log \sum_{k=1}^K \pi_k \mathcal{N}(\mathbf{x} | \mu_k, \Sigma_k)\end{aligned}$$

Parameter Learning for GMMs

- ▶ Objective: Maximum likelihood estimate of model parameters θ
- ▶ $\theta := \{\pi_k, \mu_k, \Sigma_k, k = 1, \dots, K\}$
- ▶ Maximum likelihood estimate:

$$\begin{aligned}\theta^* &= \arg \max_{\theta} p(x|\theta) = \arg \max_{\theta} \sum_{k=1}^K \pi_k \mathcal{N}(x | \mu_k, \Sigma_k) \\ &\stackrel{\log}{=} \arg \max_{\theta} \log \sum_{k=1}^K \pi_k \mathcal{N}(x | \mu_k, \Sigma_k)\end{aligned}$$

- ▶ **Problem:** We cannot move the log into the sum
▶▶ **Nasty optimization problem**

Parameter Learning for GMMs

- ▶ Objective: Maximum likelihood estimate of model parameters θ
- ▶ $\theta := \{\pi_k, \mu_k, \Sigma_k, k = 1, \dots, K\}$
- ▶ Maximum likelihood estimate:

$$\begin{aligned}\theta^* &= \arg \max_{\theta} p(\mathbf{x}|\theta) = \arg \max_{\theta} \sum_{k=1}^K \pi_k \mathcal{N}(\mathbf{x} | \mu_k, \Sigma_k) \\ &\stackrel{\log}{=} \arg \max_{\theta} \log \sum_{k=1}^K \pi_k \mathcal{N}(\mathbf{x} | \mu_k, \Sigma_k)\end{aligned}$$

- ▶ **Problem:** We cannot move the log into the sum
 - ▶▶ **Nasty optimization problem**
- ▶▶ Iterative scheme (**EM Algorithm**) for learning parameters

GMM Likelihood

Assume an i.i.d. data set $\mathbf{x} = \mathbf{x}_1, \dots, \mathbf{x}_N$ is given, and we want to determine the optimal parameters $\boldsymbol{\theta}^*$ of the GMM via Maximum Likelihood

1. Likelihood:

$$p(\mathbf{x}|\boldsymbol{\theta}) = \prod_{i=1}^N p(\mathbf{x}_i|\boldsymbol{\theta}), \quad p(\mathbf{x}_i|\boldsymbol{\theta}) = \sum_{k=1}^K \pi_k \mathcal{N}(\mathbf{x}_i | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)$$

GMM Likelihood

Assume an i.i.d. data set $\mathbf{x} = \mathbf{x}_1, \dots, \mathbf{x}_N$ is given, and we want to determine the optimal parameters $\boldsymbol{\theta}^*$ of the GMM via Maximum Likelihood

1. Likelihood:

$$p(\mathbf{x}|\boldsymbol{\theta}) = \prod_{i=1}^N p(\mathbf{x}_i|\boldsymbol{\theta}), \quad p(\mathbf{x}_i|\boldsymbol{\theta}) = \sum_{k=1}^K \pi_k \mathcal{N}(\mathbf{x}_i | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)$$

2. Log-likelihood:

$$\log p(\mathbf{x}|\boldsymbol{\theta}) = \sum_{i=1}^N \log p(\mathbf{x}_i|\boldsymbol{\theta}) = \underbrace{\sum_{i=1}^N \log \sum_{k=1}^K \pi_k \mathcal{N}(\mathbf{x}_i | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)}_{=:L}$$

Necessary Optimality Conditions

Learning Objective

Find parameters θ^* that maximize the log-likelihood

$$\log p(\mathbf{x}|\theta) = \sum_{i=1}^N \log p(x_i|\theta) = \underbrace{\sum_{i=1}^N \log \sum_{k=1}^K \pi_k \mathcal{N}(x_i | \mu_k, \Sigma_k)}_{=:L}$$

Necessary Optimality Conditions

Learning Objective

Find parameters θ^* that maximize the log-likelihood

$$\log p(\mathbf{x}|\theta) = \sum_{i=1}^N \log p(\mathbf{x}_i|\theta) = \underbrace{\sum_{i=1}^N \log \sum_{k=1}^K \pi_k \mathcal{N}(\mathbf{x}_i | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)}_{=:L}$$

$$\frac{\partial L}{\partial \boldsymbol{\mu}_k} = \mathbf{0} \Leftrightarrow \sum_{i=1}^N \frac{\partial \log p(\mathbf{x}_i|\theta)}{\partial \boldsymbol{\mu}_k} = 0$$

$$\frac{\partial L}{\partial \boldsymbol{\Sigma}_k} = \mathbf{0} \Leftrightarrow \sum_{i=1}^N \frac{\partial \log p(\mathbf{x}_i|\theta)}{\partial \boldsymbol{\Sigma}_k} = 0$$

$$\frac{\partial L}{\partial \pi_k} = 0 \Leftrightarrow \sum_{i=1}^N \frac{\partial \log p(\mathbf{x}_i|\theta)}{\partial \pi_k} = 0$$

High-Level Gradients

We need to compute gradients of the form

$$\frac{\partial \log p(\mathbf{x}_i | \boldsymbol{\theta})}{\partial \boldsymbol{\theta}}$$

High-Level Gradients

We need to compute gradients of the form

$$\frac{\partial \log p(\mathbf{x}_i|\boldsymbol{\theta})}{\partial \boldsymbol{\theta}}$$

In general,

$$\frac{\partial \log p(\mathbf{x}_i|\boldsymbol{\theta})}{\partial \boldsymbol{\theta}} = \frac{1}{p(\mathbf{x}_i|\boldsymbol{\theta})} \frac{\partial p(\mathbf{x}_i|\boldsymbol{\theta})}{\partial \boldsymbol{\theta}}$$

High-Level Gradients

We need to compute gradients of the form

$$\frac{\partial \log p(\mathbf{x}_i | \boldsymbol{\theta})}{\partial \boldsymbol{\theta}}$$

In general,

$$\frac{\partial \log p(\mathbf{x}_i | \boldsymbol{\theta})}{\partial \boldsymbol{\theta}} = \frac{1}{p(\mathbf{x}_i | \boldsymbol{\theta})} \frac{\partial p(\mathbf{x}_i | \boldsymbol{\theta})}{\partial \boldsymbol{\theta}}$$

With

$$p(\mathbf{x}_i | \boldsymbol{\theta}) = \sum_{j=1}^K \pi_j \mathcal{N}(\mathbf{x}_i | \boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j)$$

we get

$$\frac{\partial p(\mathbf{x}_i | \boldsymbol{\theta})}{\partial \boldsymbol{\mu}_k} = \sum_{j=1}^K \pi_j \frac{\partial \mathcal{N}(\mathbf{x}_i | \boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j)}{\partial \boldsymbol{\mu}_k}$$

In More Detail

$$\frac{\partial p(\mathbf{x}_i|\boldsymbol{\theta})}{\partial \boldsymbol{\mu}_k} = \sum_{j=1}^K \pi_j \frac{\partial \mathcal{N}(\mathbf{x}_i | \boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j)}{\partial \boldsymbol{\mu}_k} =$$

In More Detail

$$\frac{\partial p(\mathbf{x}_i|\boldsymbol{\theta})}{\partial \boldsymbol{\mu}_k} = \sum_{j=1}^K \pi_j \frac{\partial \mathcal{N}(\mathbf{x}_i | \boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j)}{\partial \boldsymbol{\mu}_k} = \pi_k \frac{\partial \mathcal{N}(\mathbf{x}_i | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)}{\partial \boldsymbol{\mu}_k}$$

In More Detail

$$\begin{aligned}\frac{\partial p(\mathbf{x}_i|\boldsymbol{\theta})}{\partial \boldsymbol{\mu}_k} &= \sum_{j=1}^K \pi_j \frac{\partial \mathcal{N}(\mathbf{x}_i | \boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j)}{\partial \boldsymbol{\mu}_k} = \pi_k \frac{\partial \mathcal{N}(\mathbf{x}_i | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)}{\partial \boldsymbol{\mu}_k} \\ &= \pi_k (\mathbf{x}_i - \boldsymbol{\mu}_k)^\top \boldsymbol{\Sigma}_k^{-1} \mathcal{N}(\mathbf{x}_i | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)\end{aligned}$$

In More Detail

$$\begin{aligned}\frac{\partial p(\mathbf{x}_i|\boldsymbol{\theta})}{\partial \boldsymbol{\mu}_k} &= \sum_{j=1}^K \pi_j \frac{\partial \mathcal{N}(\mathbf{x}_i | \boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j)}{\partial \boldsymbol{\mu}_k} = \pi_k \frac{\partial \mathcal{N}(\mathbf{x}_i | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)}{\partial \boldsymbol{\mu}_k} \\ &= \pi_k (\mathbf{x}_i - \boldsymbol{\mu}_k)^\top \boldsymbol{\Sigma}_k^{-1} \mathcal{N}(\mathbf{x}_i | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)\end{aligned}$$

Overall:

$$\frac{\partial L}{\partial \boldsymbol{\mu}_k} = \sum_{i=1}^N \frac{\partial \log p(\mathbf{x}_i|\boldsymbol{\theta})}{\partial \boldsymbol{\mu}_k} = \sum_{i=1}^N \frac{1}{p(\mathbf{x}_i|\boldsymbol{\theta})} \frac{\partial p(\mathbf{x}_i|\boldsymbol{\theta})}{\partial \boldsymbol{\mu}_k}$$

In More Detail

$$\begin{aligned}\frac{\partial p(\mathbf{x}_i|\boldsymbol{\theta})}{\partial \boldsymbol{\mu}_k} &= \sum_{j=1}^K \pi_j \frac{\partial \mathcal{N}(\mathbf{x}_i | \boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j)}{\partial \boldsymbol{\mu}_k} = \pi_k \frac{\partial \mathcal{N}(\mathbf{x}_i | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)}{\partial \boldsymbol{\mu}_k} \\ &= \pi_k (\mathbf{x}_i - \boldsymbol{\mu}_k)^\top \boldsymbol{\Sigma}_k^{-1} \mathcal{N}(\mathbf{x}_i | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)\end{aligned}$$

Overall:

$$\begin{aligned}\frac{\partial L}{\partial \boldsymbol{\mu}_k} &= \sum_{i=1}^N \frac{\partial \log p(\mathbf{x}_i|\boldsymbol{\theta})}{\partial \boldsymbol{\mu}_k} = \sum_{i=1}^N \frac{1}{p(\mathbf{x}_i|\boldsymbol{\theta})} \frac{\partial p(\mathbf{x}_i|\boldsymbol{\theta})}{\partial \boldsymbol{\mu}_k} \\ &= \sum_{i=1}^N (\mathbf{x}_i - \boldsymbol{\mu}_k)^\top \boldsymbol{\Sigma}_k^{-1} \underbrace{\frac{\pi_k \mathcal{N}(\mathbf{x}_i | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)}{\sum_j \pi_j \mathcal{N}(\mathbf{x}_i | \boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j)}}_{:=r_{ik}}\end{aligned}$$

In More Detail

$$\begin{aligned}\frac{\partial p(\mathbf{x}_i|\boldsymbol{\theta})}{\partial \boldsymbol{\mu}_k} &= \sum_{j=1}^K \pi_j \frac{\partial \mathcal{N}(\mathbf{x}_i | \boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j)}{\partial \boldsymbol{\mu}_k} = \pi_k \frac{\partial \mathcal{N}(\mathbf{x}_i | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)}{\partial \boldsymbol{\mu}_k} \\ &= \pi_k (\mathbf{x}_i - \boldsymbol{\mu}_k)^\top \boldsymbol{\Sigma}_k^{-1} \mathcal{N}(\mathbf{x}_i | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)\end{aligned}$$

Overall:

$$\begin{aligned}\frac{\partial L}{\partial \boldsymbol{\mu}_k} &= \sum_{i=1}^N \frac{\partial \log p(\mathbf{x}_i|\boldsymbol{\theta})}{\partial \boldsymbol{\mu}_k} = \sum_{i=1}^N \frac{1}{p(\mathbf{x}_i|\boldsymbol{\theta})} \frac{\partial p(\mathbf{x}_i|\boldsymbol{\theta})}{\partial \boldsymbol{\mu}_k} \\ &= \sum_{i=1}^N (\mathbf{x}_i - \boldsymbol{\mu}_k)^\top \boldsymbol{\Sigma}_k^{-1} \underbrace{\frac{\pi_k \mathcal{N}(\mathbf{x}_i | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)}{\sum_j \pi_j \mathcal{N}(\mathbf{x}_i | \boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j)}}_{:=r_{ik}} = \sum_{i=1}^N r_{ik} (\mathbf{x}_i - \boldsymbol{\mu}_k)^\top \boldsymbol{\Sigma}_k^{-1} = \mathbf{0}\end{aligned}$$

In More Detail

$$\begin{aligned}\frac{\partial p(\mathbf{x}_i|\boldsymbol{\theta})}{\partial \boldsymbol{\mu}_k} &= \sum_{j=1}^K \pi_j \frac{\partial \mathcal{N}(\mathbf{x}_i | \boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j)}{\partial \boldsymbol{\mu}_k} = \pi_k \frac{\partial \mathcal{N}(\mathbf{x}_i | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)}{\partial \boldsymbol{\mu}_k} \\ &= \pi_k (\mathbf{x}_i - \boldsymbol{\mu}_k)^\top \boldsymbol{\Sigma}_k^{-1} \mathcal{N}(\mathbf{x}_i | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)\end{aligned}$$

Overall:

$$\begin{aligned}\frac{\partial L}{\partial \boldsymbol{\mu}_k} &= \sum_{i=1}^N \frac{\partial \log p(\mathbf{x}_i|\boldsymbol{\theta})}{\partial \boldsymbol{\mu}_k} = \sum_{i=1}^N \frac{1}{p(\mathbf{x}_i|\boldsymbol{\theta})} \frac{\partial p(\mathbf{x}_i|\boldsymbol{\theta})}{\partial \boldsymbol{\mu}_k} \\ &= \sum_{i=1}^N (\mathbf{x}_i - \boldsymbol{\mu}_k)^\top \boldsymbol{\Sigma}_k^{-1} \underbrace{\frac{\pi_k \mathcal{N}(\mathbf{x}_i | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)}{\sum_j \pi_j \mathcal{N}(\mathbf{x}_i | \boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j)}}_{:=r_{ik}} = \sum_{i=1}^N r_{ik} (\mathbf{x}_i - \boldsymbol{\mu}_k)^\top \boldsymbol{\Sigma}_k^{-1} = \mathbf{0} \\ &\Leftrightarrow \sum_{i=1}^N r_{ik} \mathbf{x}_i = \sum_{i=1}^N r_{ik} \boldsymbol{\mu}_k\end{aligned}$$

In More Detail

$$\begin{aligned}\frac{\partial p(\mathbf{x}_i|\boldsymbol{\theta})}{\partial \boldsymbol{\mu}_k} &= \sum_{j=1}^K \pi_j \frac{\partial \mathcal{N}(\mathbf{x}_i | \boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j)}{\partial \boldsymbol{\mu}_k} = \pi_k \frac{\partial \mathcal{N}(\mathbf{x}_i | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)}{\partial \boldsymbol{\mu}_k} \\ &= \pi_k (\mathbf{x}_i - \boldsymbol{\mu}_k)^\top \boldsymbol{\Sigma}_k^{-1} \mathcal{N}(\mathbf{x}_i | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)\end{aligned}$$

Overall:

$$\begin{aligned}\frac{\partial L}{\partial \boldsymbol{\mu}_k} &= \sum_{i=1}^N \frac{\partial \log p(\mathbf{x}_i|\boldsymbol{\theta})}{\partial \boldsymbol{\mu}_k} = \sum_{i=1}^N \frac{1}{p(\mathbf{x}_i|\boldsymbol{\theta})} \frac{\partial p(\mathbf{x}_i|\boldsymbol{\theta})}{\partial \boldsymbol{\mu}_k} \\ &= \sum_{i=1}^N (\mathbf{x}_i - \boldsymbol{\mu}_k)^\top \boldsymbol{\Sigma}_k^{-1} \underbrace{\frac{\pi_k \mathcal{N}(\mathbf{x}_i | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)}{\sum_j \pi_j \mathcal{N}(\mathbf{x}_i | \boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j)}}_{:=r_{ik}} = \sum_{i=1}^N r_{ik} (\mathbf{x}_i - \boldsymbol{\mu}_k)^\top \boldsymbol{\Sigma}_k^{-1} = \mathbf{0} \\ \Leftrightarrow \sum_{i=1}^N r_{ik} \mathbf{x}_i &= \sum_{i=1}^N r_{ik} \boldsymbol{\mu}_k \Leftrightarrow \boldsymbol{\mu}_k = \frac{\sum_{i=1}^N r_{ik} \mathbf{x}_i}{\sum_{i=1}^N r_{ik}} = \frac{1}{N_k} \sum_{i=1}^N r_{ik} \mathbf{x}_i\end{aligned}$$

Similarly...

$$\frac{\partial L}{\partial \Sigma_k} = \mathbf{0} \Leftrightarrow \Sigma_k = \frac{1}{N_k} \sum_{i=1}^N r_{ik} (\mathbf{x}_n - \boldsymbol{\mu}_k)(\mathbf{x}_n - \boldsymbol{\mu}_k)^\top$$

$$\frac{\partial L}{\partial \pi_k} = \mathbf{0} \Leftrightarrow \pi_k = \frac{N_k}{N} \quad \blacktriangleright \text{Requires Lagrange multipliers}$$

Similarly...

$$\frac{\partial L}{\partial \Sigma_k} = \mathbf{0} \Leftrightarrow \Sigma_k = \frac{1}{N_k} \sum_{i=1}^N r_{ik} (\mathbf{x}_n - \boldsymbol{\mu}_k) (\mathbf{x}_n - \boldsymbol{\mu}_k)^\top$$

$$\frac{\partial L}{\partial \pi_k} = \mathbf{0} \Leftrightarrow \pi_k = \frac{N_k}{N} \quad \blacktriangleright \text{Requires Lagrange multipliers}$$

- ▶ **Bad news:** These results do not constitute a closed-form solution of the parameters $\boldsymbol{\mu}_k, \Sigma_k, \pi_k$ of the mixture model because the responsibilities r_{ik} depend on those parameters in a complex way.
- ▶ **Good news:** Results suggest a simple **iterative** scheme for finding a solution to the MLE problem.

EM Algorithm

- ▶ Iterative scheme for learning parameters in mixture models and latent-variable models
 1. Choose initial values for $\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k, \pi_k$
 2. Until convergence, alternate between
 - ▶ **E-step:** Evaluate the responsibilities r_{ik} (posterior probability of data point i belonging to mixture component k)
 - ▶ **M-step:** Use the updated responsibilities to re-estimate the parameters $\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k, \pi_k$
- ▶ Every step in the EM algorithm increases the likelihood function
- ▶ Convergence: Check log-likelihood or the parameters

Implementation

1. Initialize μ_k, Σ_k, π_k

Implementation

1. Initialize $\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k, \pi_k$
2. **E-step:** Evaluate responsibilities for every data point \mathbf{x}_i using current parameters $\pi_k, \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k$:

$$r_{ik} = \frac{\pi_k \mathcal{N}(\mathbf{x}_i | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)}{\sum_j \pi_j \mathcal{N}(\mathbf{x}_i | \boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j)}$$

Implementation

1. Initialize $\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k, \pi_k$
2. **E-step:** Evaluate responsibilities for every data point \mathbf{x}_i using current parameters $\pi_k, \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k$:

$$r_{ik} = \frac{\pi_k \mathcal{N}(\mathbf{x}_i | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)}{\sum_j \pi_j \mathcal{N}(\mathbf{x}_i | \boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j)}$$

3. **M-step:** Re-estimate parameters $\pi_k, \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k$ using the current responsibilities r_{ik} (from E-step):

$$\boldsymbol{\mu}_k = \frac{1}{N_k} \sum_{i=1}^N r_{ik} \mathbf{x}_i$$

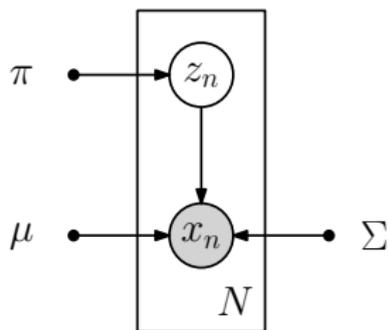
$$\boldsymbol{\Sigma}_k = \frac{1}{N_k} \sum_{i=1}^N r_{ik} (\mathbf{x}_i - \boldsymbol{\mu}_k)(\mathbf{x}_i - \boldsymbol{\mu}_k)^\top$$

$$\pi_k = \frac{N_k}{N}$$

Example

Demo

The Latent-Variable Perspective



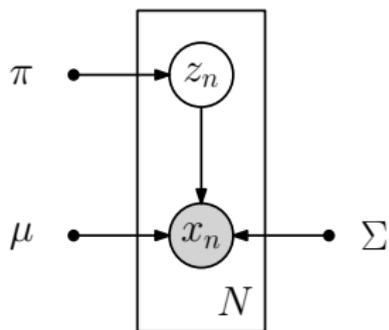
$$p(z_k = 1) = \pi_k, \quad 0 \leq \pi_k \leq 1, \quad \sum_k \pi_k = 1$$

$$p(\mathbf{x}|z_k = 1) = \mathcal{N}(\mathbf{x} | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)$$

$$p(\mathbf{x}) = \sum_z p(\mathbf{z})p(\mathbf{x}|\mathbf{z}) = \sum_{k=1}^K \pi_k \mathcal{N}(\mathbf{x} | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)$$

- ▶ $\mathbf{z}_n = (z_1, \dots, z_K)$ is a discrete latent variable. Exactly one entry of \mathbf{z}_n is 1, all others are 0 ►► 1-of- K code

The Latent-Variable Perspective



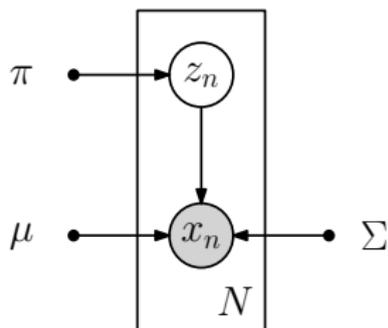
$$p(z_k = 1) = \pi_k, \quad 0 \leq \pi_k \leq 1, \quad \sum_k \pi_k = 1$$

$$p(\mathbf{x}|z_k = 1) = \mathcal{N}(\mathbf{x} | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)$$

$$p(\mathbf{x}) = \sum_z p(\mathbf{z})p(\mathbf{x}|\mathbf{z}) = \sum_{k=1}^K \pi_k \mathcal{N}(\mathbf{x} | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)$$

- ▶ $\mathbf{z}_n = (z_1, \dots, z_K)$ is a discrete latent variable. Exactly one entry of \mathbf{z}_n is 1, all others are 0 \blacktriangleright 1-of- K code
- ▶ For every observed data point \mathbf{x}_n there is a corresponding latent variable \mathbf{z}_n , which indicates which mixture component generated \mathbf{x}_n

The Latent-Variable Perspective



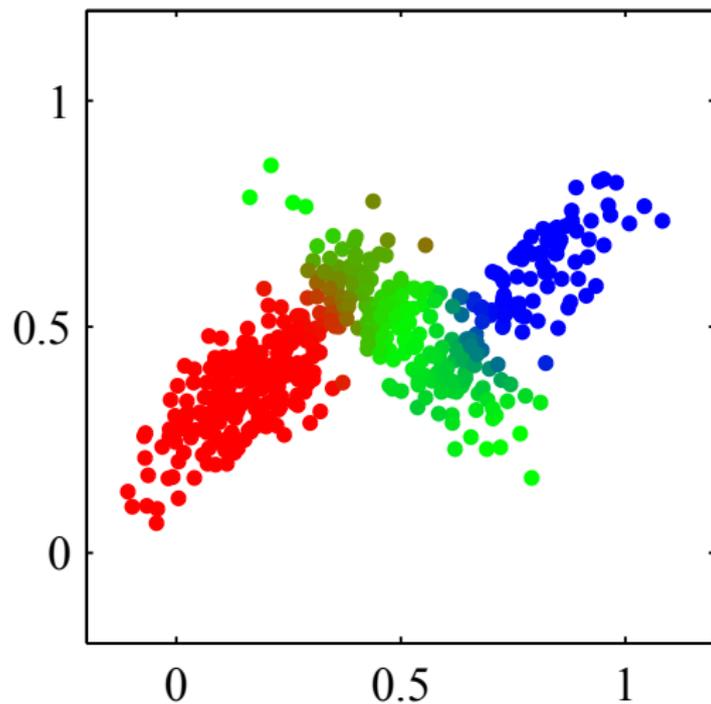
$$p(z_k = 1) = \pi_k, \quad 0 \leq \pi_k \leq 1, \quad \sum_k \pi_k = 1$$

$$p(\mathbf{x}|z_k = 1) = \mathcal{N}(\mathbf{x} | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)$$

$$p(\mathbf{x}) = \sum_z p(\mathbf{z})p(\mathbf{x}|\mathbf{z}) = \sum_{k=1}^K \pi_k \mathcal{N}(\mathbf{x} | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)$$

- ▶ $\mathbf{z}_n = (z_1, \dots, z_K)$ is a discrete latent variable. Exactly one entry of \mathbf{z}_n is 1, all others are 0 \blacktriangleright 1-of- K code
- ▶ For every observed data point \mathbf{x}_n there is a corresponding latent variable \mathbf{z}_n , which indicates which mixture component generated \mathbf{x}_n
- ▶ Posterior $p(z_k = 1 | \mathbf{x}_i) = r_{ik}$ corresponds to the “responsibility” (see earlier) that mixture component k generated data point i .

Visualizing the Responsibilities



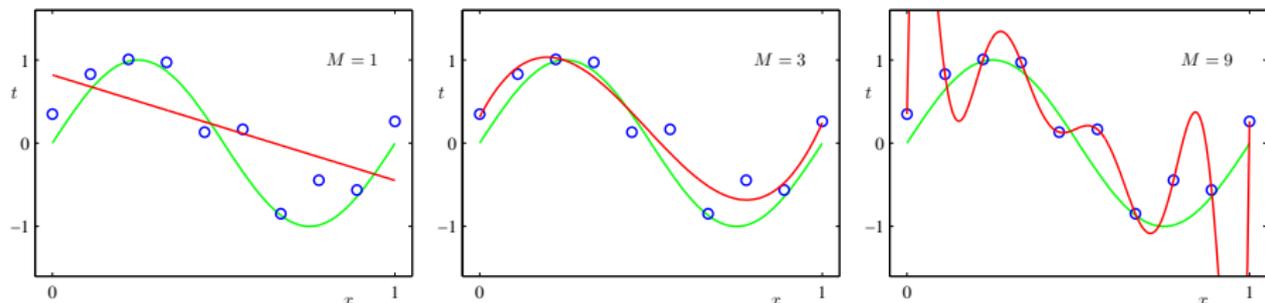
From PRML (Bishop, 2006)

EM with Latent Variables (see CO-495)

- ▶ Latent-variable perspective gives rise to a **general EM algorithm for maximum likelihood parameter estimation** (regression, classification, dimensionality reduction, density estimation, ...), see Dempster et al., (1977)
- ▶ EM iteratively maximizes a lower bound on the log likelihood $\log p(\mathbf{X}|\boldsymbol{\theta})$
- ▶ At the same time, EM iteratively minimizes the KL divergence $\text{KL}(q(\mathbf{Z})||p(\mathbf{Z}|\mathbf{X}, \boldsymbol{\theta}))$ between an approximate posterior $q(\mathbf{Z})$ and the true (but unknown) posterior distribution $p(\mathbf{Z}|\mathbf{X}, \boldsymbol{\theta})$

Model Selection

Model Selection

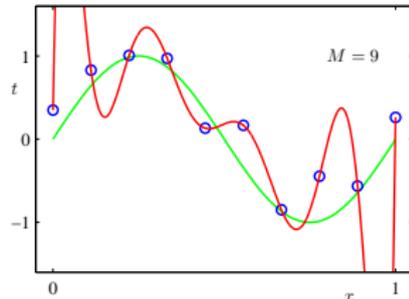
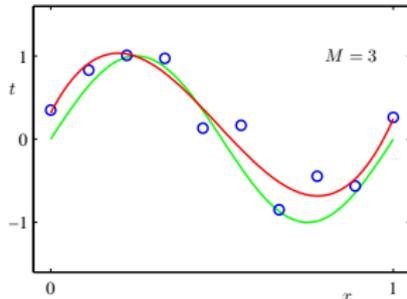
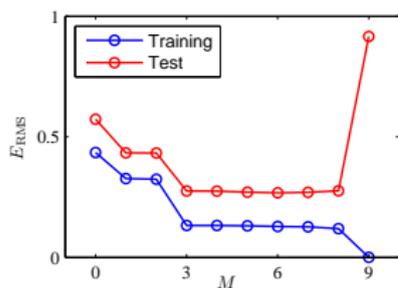


From PRML (Bishop, 2006)

Sometimes, we have to make high-level decisions about the model we want to use:

- ▶ Number of components in a mixture model
- ▶ Network architecture of (deep) neural networks
- ▶ Type of kernel in a support vector machine
- ▶ Degree of a polynomial in a regression problem

Test vs Training Error

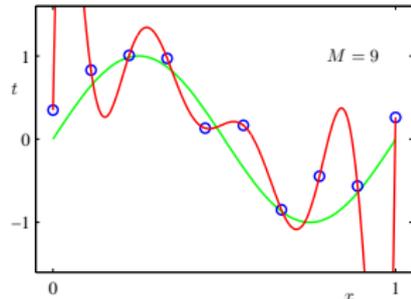
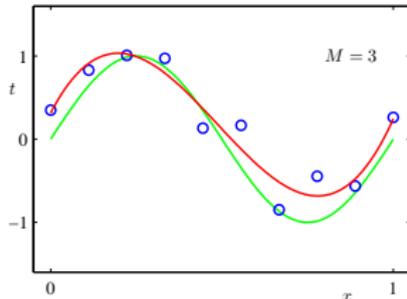
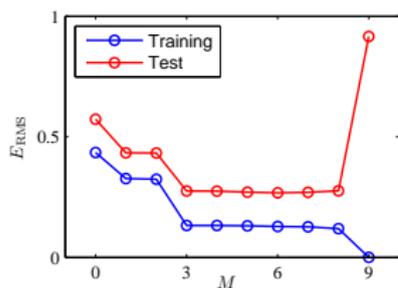


From PRML (Bishop, 2006)

General problem:

- ▶ Model fits training data perfectly, but may not do well on test data ► **Overfitting** (especially with MLE)

Test vs Training Error

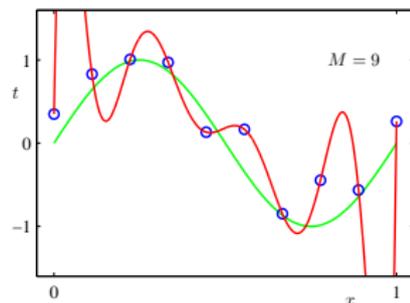
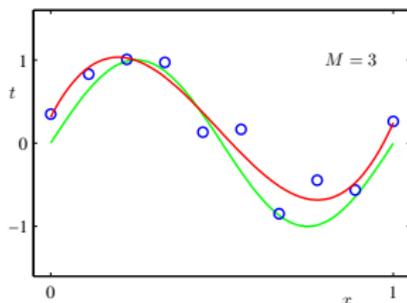
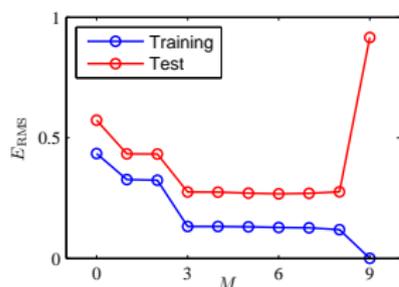


From PRML (Bishop, 2006)

General problem:

- ▶ Model fits training data perfectly, but may not do well on test data ► **Overfitting** (especially with MLE)
- ▶ Training performance \neq test performance, but we are largely interested in test performance

Test vs Training Error

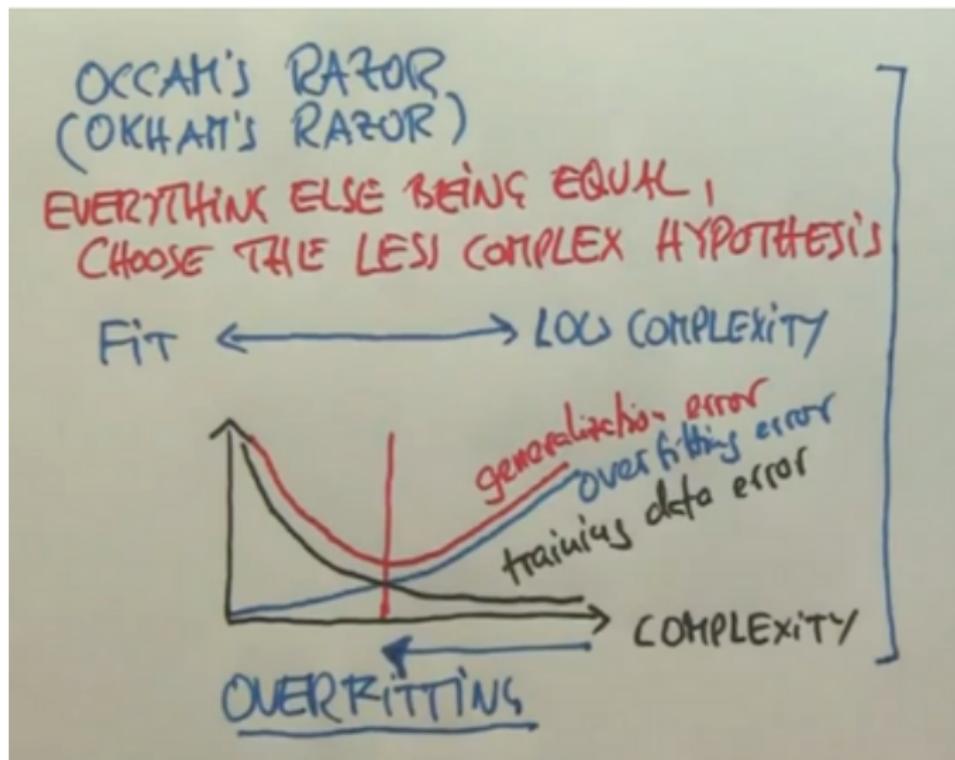


From PRML (Bishop, 2006)

General problem:

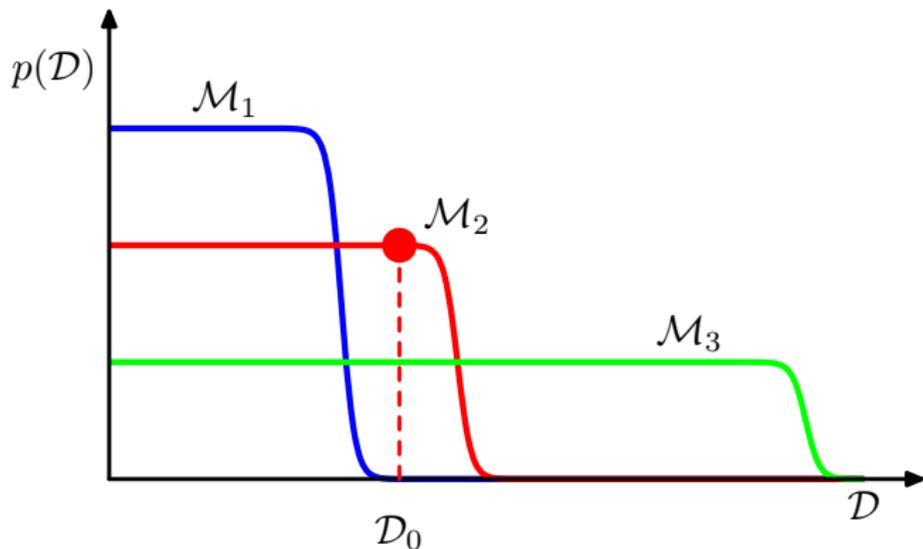
- ▶ Model fits training data perfectly, but may not do well on test data ▶ **Overfitting** (especially with MLE)
- ▶ Training performance \neq test performance, but we are largely interested in test performance
- ▶ Need mechanisms for assessing how a model generalizes to unseen test data ▶ **Model selection**

Occam's Razor



From crowfly.net

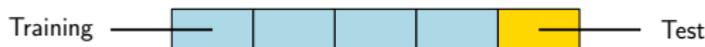
Occam's Razor (2)



From PRML (Bishop, 2006)

- Choose the simplest model that explains the data reasonably well

Cross Validation



- ▶ Partition your training data into L subsets
- ▶ Train the model on $L - 1$ subsets
- ▶ Evaluate the model on the other subset

Cross Validation



- ▶ Partition your training data into L subsets
- ▶ Train the model on $L - 1$ subsets
- ▶ Evaluate the model on the other subset
- ▶ To reduce variability, multiple rounds of cross-validation are performed using different partitions, and the validation results are averaged over the rounds.

Cross Validation



- ▶ Partition your training data into L subsets
- ▶ Train the model on $L - 1$ subsets
- ▶ Evaluate the model on the other subset
- ▶ To reduce variability, multiple rounds of cross-validation are performed using different partitions, and the validation results are averaged over the rounds.
- ▶ Train many models, compare test error

Cross Validation



- ▶ Partition your training data into L subsets
- ▶ Train the model on $L - 1$ subsets
- ▶ Evaluate the model on the other subset
- ▶ To reduce variability, multiple rounds of cross-validation are performed using different partitions, and the validation results are averaged over the rounds.
- ▶ Train many models, compare test error

Number of training runs increases with the number of partitions

Information Criteria

- ▶ Correct for the bias of MLE by addition of a penalty term to compensate for the overfitting of more complex models (with lots of parameters)

Information Criteria

- ▶ Correct for the bias of MLE by addition of a penalty term to compensate for the overfitting of more complex models (with lots of parameters)
- ▶ Akaike Information Criterion (Akaike 1974):

$$\ln p(\mathbf{x}|\boldsymbol{\theta}_{\text{ML}}) - M$$

where M is the number of model parameters

Information Criteria

- ▶ Correct for the bias of MLE by addition of a penalty term to compensate for the overfitting of more complex models (with lots of parameters)
- ▶ Akaike Information Criterion (Akaike 1974):

$$\ln p(\mathbf{x}|\boldsymbol{\theta}_{\text{ML}}) - M$$

where M is the number of model parameters

- ▶ Bayesian Information Criterion/MDL (Schwarz 1978) (for exponential family distributions):

$$\ln p(\mathbf{x}) = \ln \int p(\mathbf{x}|\boldsymbol{\theta})p(\boldsymbol{\theta})d\boldsymbol{\theta} \approx \ln p(\mathbf{x}|\boldsymbol{\theta}_{\text{ML}}) - \frac{1}{2}M \ln N$$

where N is the number of data points and M is the number of parameters.

Information Criteria

- ▶ Correct for the bias of MLE by addition of a penalty term to compensate for the overfitting of more complex models (with lots of parameters)
- ▶ Akaike Information Criterion (Akaike 1974):

$$\ln p(\mathbf{x}|\boldsymbol{\theta}_{\text{ML}}) - M$$

where M is the number of model parameters

- ▶ Bayesian Information Criterion/MDL (Schwarz 1978) (for exponential family distributions):

$$\ln p(\mathbf{x}) = \ln \int p(\mathbf{x}|\boldsymbol{\theta})p(\boldsymbol{\theta})d\boldsymbol{\theta} \approx \ln p(\mathbf{x}|\boldsymbol{\theta}_{\text{ML}}) - \frac{1}{2}M \ln N$$

where N is the number of data points and M is the number of parameters.

- ▶ BIC penalizes model complexity more heavily than AIC.

Bayesian Model Comparison

- ▶ Place a prior $p(M)$ on the class of models

¹When would the integral be tractable?

Bayesian Model Comparison

- ▶ Place a prior $p(M)$ on the class of models
- ▶ Given a training set \mathcal{D} , we compute the **posterior distribution over models** as

$$p(M_i|\mathcal{D}) \propto p(M_i)p(\mathcal{D}|M_i)$$

which allows us to express a preference for different models

¹When would the integral be tractable?

Bayesian Model Comparison

- ▶ Place a prior $p(M)$ on the class of models
- ▶ Given a training set \mathcal{D} , we compute the **posterior distribution over models** as

$$p(M_i|\mathcal{D}) \propto p(M_i)p(\mathcal{D}|M_i)$$

which allows us to express a preference for different models

- ▶ **Model evidence (marginal likelihood):**

$$p(\mathcal{D}|M_i) = \int p(\mathcal{D}|\boldsymbol{\theta}_{M_i})p(\boldsymbol{\theta}_{M_i}|M_i)d\boldsymbol{\theta}_{M_i}$$

¹When would the integral be tractable?

Bayesian Model Comparison

- ▶ Place a prior $p(M)$ on the class of models
- ▶ Given a training set \mathcal{D} , we compute the **posterior distribution over models** as

$$p(M_i|\mathcal{D}) \propto p(M_i)p(\mathcal{D}|M_i)$$

which allows us to express a preference for different models

- ▶ **Model evidence (marginal likelihood):**

$$p(\mathcal{D}|M_i) = \int p(\mathcal{D}|\boldsymbol{\theta}_{M_i})p(\boldsymbol{\theta}_{M_i}|M_i)d\boldsymbol{\theta}_{M_i}$$

- ▶ **Bayes factor** for comparing two models: $p(\mathcal{D}|M_1)/p(\mathcal{D}|M_2)$

¹When would the integral be tractable?

Bayesian Model Comparison

- ▶ Place a prior $p(M)$ on the class of models
- ▶ Given a training set \mathcal{D} , we compute the **posterior distribution over models** as

$$p(M_i|\mathcal{D}) \propto p(M_i)p(\mathcal{D}|M_i)$$

which allows us to express a preference for different models

- ▶ **Model evidence (marginal likelihood):**

$$p(\mathcal{D}|M_i) = \int p(\mathcal{D}|\boldsymbol{\theta}_{M_i})p(\boldsymbol{\theta}_{M_i}|M_i)d\boldsymbol{\theta}_{M_i}$$

- ▶ **Bayes factor** for comparing two models: $p(\mathcal{D}|M_1)/p(\mathcal{D}|M_2)$
- ▶ **Integral often intractable**¹

¹When would the integral be tractable?

Bayesian Model Averaging

- ▶ Place a prior $p(M)$ on the class of models
- ▶ Instead of selecting the “best” model, **integrate out** the corresponding **model parameters** θ_M and **average over all models** $M_i, i = 1, \dots, L$

$$p(\mathcal{D}) = \sum_{i=1}^L p(M_i) \underbrace{\int p(\mathcal{D}|\theta_{M_i})p(\theta_{M_i}|M_i)d\theta_{M_i}}_{=p(\mathcal{D}|M_i)}$$

Bayesian Model Averaging

- ▶ Place a prior $p(M)$ on the class of models
- ▶ Instead of selecting the “best” model, **integrate out** the corresponding **model parameters** θ_M and **average over all models** $M_i, i = 1, \dots, L$

$$p(\mathcal{D}) = \sum_{i=1}^L p(M_i) \underbrace{\int p(\mathcal{D}|\theta_{M_i})p(\theta_{M_i}|M_i)d\theta_{M_i}}_{=p(\mathcal{D}|M_i)}$$

- ▶ Computationally expensive
- ▶ Integral often intractable (still...)

References I

- [1] H. Akaike. A New Look at the Statistical Model Identification. *IEEE Transactions on Automatic Control*, 19(6):716–723, 1974.
- [2] C. M. Bishop. *Pattern Recognition and Machine Learning*. Information Science and Statistics. Springer-Verlag, 2006.
- [3] A. P. Dempster, N. M. Laird, and D. B. Rubin. Maximum Likelihood from Incomplete Data via the EM Algorithm. *Journal of the Royal Statistical Society*, 39(1):1–38, 1977.
- [4] R. M. Neal and G. E. Hinton. A View of the EM Algorithm that Justifies Incremental, Sparse, and Other Variants. In M. I. Jordan, editor, *Learning in Graphical Models*, pages 355–368. MIT Press, 1999.
- [5] G. E. Schwarz. Estimating the Dimension of a Model. *Annals of Statistics*, 6(2):461–464, 1978.