

Lecture 2

Expectation Maximisation and Dynamic Networks

G. Englebienne
M. Poel

University of Twente

Introduction

kMeans Clustering

Mixtures of Gaussians

The General EM Algorithm

The E-step and M-step

Why it works

Models of data sequences

Hidden Markov Models

Linear Dynamical System

Introduction

kMeans Clustering

Mixtures of Gaussians

The General EM Algorithm

The E-step and M-step

Why it works

Models of data sequences

Hidden Markov Models

Linear Dynamical System

Last week, we have seen how the factorisation of probability distributions could be represented as graphs.

- ▶ We have seen algorithms to efficiently compute marginal probabilities on such graphs.
- ▶ We have seen how this could be used to compute conditional probabilities.
- ▶ However in we needed to assume that the parameters of the distributions were known.

Today, we see how to learn these parameters efficiently.

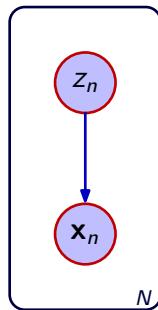
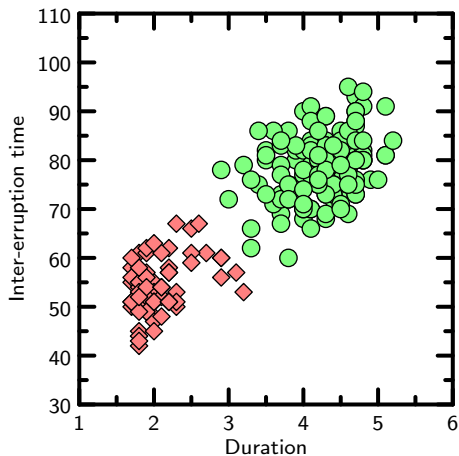
We want to find the parameters that maximise the likelihood

- ▶ MAP treatment is very similar
- ▶ (exact) full Bayesian treatment is often not tractable
- ▶ We'll see how to use sampling instead

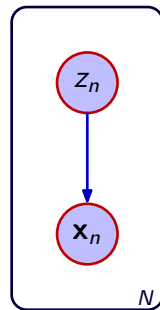
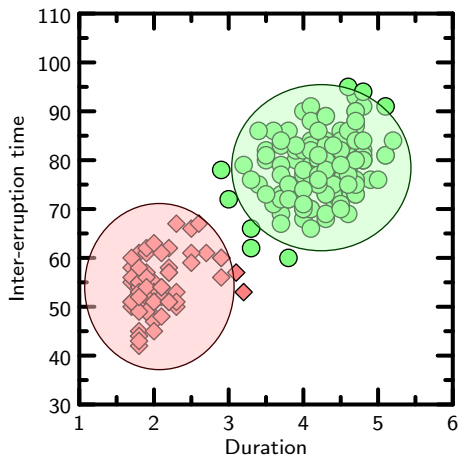
How can we find the maximum of the likelihood?

- ▶ If everything is observed, it's easy
- ▶ If we have latent variables, it's hard

Example



Example



$$p(\mathbf{x}_n, z_n) = \begin{cases} p(\mathcal{C}_1) p(\mathbf{x}_n|\mathcal{C}_1) & \text{if } z_n = 1 \\ p(\mathcal{C}_2) p(\mathbf{x}_n|\mathcal{C}_2) & \text{if } z_n = 0 \end{cases}$$

► Parametrisation:

$$p(\mathcal{C}_1) = \pi$$

$$p(\mathcal{C}_2) = 1 - \pi$$

$$p(\mathbf{x}_n|\mathcal{C}_1) = \mathcal{N}(\mathbf{x}_n; \boldsymbol{\mu}_1, \boldsymbol{\Sigma}_1)$$

$$p(\mathbf{x}_n|\mathcal{C}_2) = \mathcal{N}(\mathbf{x}_n; \boldsymbol{\mu}_2, \boldsymbol{\Sigma}_2)$$

$$\implies \boldsymbol{\theta} = \{\pi, \boldsymbol{\mu}_1, \boldsymbol{\Sigma}_1, \boldsymbol{\mu}_2, \boldsymbol{\Sigma}_2\}$$

(Complete) likelihood:

$$p(\mathbf{x}_n, z_n|\boldsymbol{\theta}) = [p(\mathcal{C}_1) p(\mathbf{x}_n|\mathcal{C}_1)]^{z_n} [p(\mathcal{C}_2) p(\mathbf{x}_n|\mathcal{C}_2)]^{1-z_n}$$

$$p(\mathbf{x}_n, z_n) = \begin{cases} p(\mathcal{C}_1) p(\mathbf{x}_n|\mathcal{C}_1) & \text{if } z_n = 1 \\ p(\mathcal{C}_2) p(\mathbf{x}_n|\mathcal{C}_2) & \text{if } z_n = 0 \end{cases}$$

► Parametrisation:

$$p(\mathcal{C}_1) = \pi$$

$$p(\mathcal{C}_2) = 1 - \pi$$

$$p(\mathbf{x}_n|\mathcal{C}_1) = \mathcal{N}(\mathbf{x}_n; \boldsymbol{\mu}_1, \boldsymbol{\Sigma}_1)$$

$$p(\mathbf{x}_n|\mathcal{C}_2) = \mathcal{N}(\mathbf{x}_n; \boldsymbol{\mu}_2, \boldsymbol{\Sigma}_2)$$

$$\implies \boldsymbol{\theta} = \{\pi, \boldsymbol{\mu}_1, \boldsymbol{\Sigma}_1, \boldsymbol{\mu}_2, \boldsymbol{\Sigma}_2\}$$

(Complete) likelihood:

$$p(\mathbf{x}_n, \mathbf{z}_n|\boldsymbol{\theta}) = [p(\mathcal{C}_1) p(\mathbf{x}_n|\mathcal{C}_1)]^{z_n} [p(\mathcal{C}_2) p(\mathbf{x}_n|\mathcal{C}_2)]^{1-z_n}$$

Complete likelihood:

$$p(\{\mathbf{x}_i, \mathbf{z}_i\}|\boldsymbol{\theta}) = \prod_{i=1}^N [\pi \mathcal{N}(\mathbf{x}_i; \boldsymbol{\mu}_1, \boldsymbol{\Sigma}_1)]^{z_i} [(1 - \pi) \mathcal{N}(\mathbf{x}_i; \boldsymbol{\mu}_2, \boldsymbol{\Sigma}_2)]^{1-z_i}$$

Maximise: take logarithm, set derivative = 0

$$\pi_1 = \frac{\sum_{i=1}^N z_i}{N}$$

$$\boldsymbol{\mu}_1 = \frac{\sum_{i=1}^N z_i \mathbf{x}_i}{\sum_{i=1}^N z_i}$$

$$\boldsymbol{\Sigma}_1 = \frac{\sum_{i=1}^N z_i \mathbf{x}_i \mathbf{x}_i^\top}{\sum_{i=1}^N z_i} - \boldsymbol{\mu}_1 \boldsymbol{\mu}_1^\top$$

$$\pi_2 = \frac{\sum_{i=1}^N (1 - z_i)}{N}$$

$$\boldsymbol{\mu}_2 = \frac{\sum_{i=1}^N (1 - z_i) \mathbf{x}_i}{\sum_{i=1}^N (1 - z_i)}$$

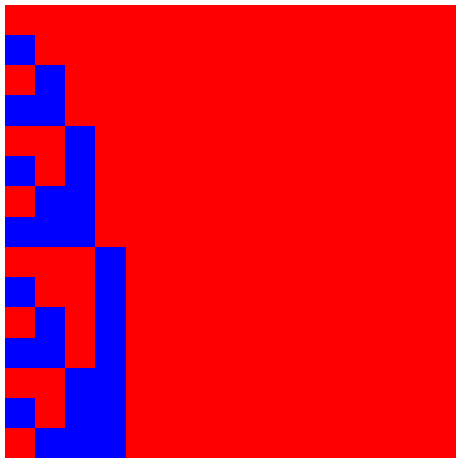
$$\boldsymbol{\Sigma}_2 = \frac{\sum_{i=1}^N (1 - z_i) \mathbf{x}_i \mathbf{x}_i^\top}{\sum_{i=1}^N (1 - z_i)} - \boldsymbol{\mu}_2 \boldsymbol{\mu}_2^\top$$

Hidden class label

Slide 10 of 61

So, what happens when the class label is not observed?

- List all possible assignments, pick best

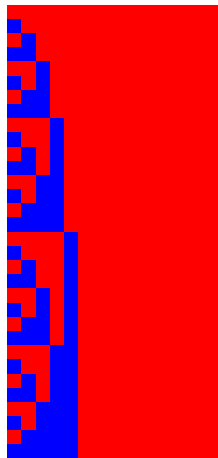


Hidden class label

Slide 10 of 61

So, what happens when the class label is not observed?

- List all possible assignments, pick best

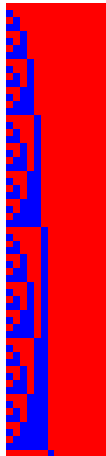


Hidden class label

Slide 10 of 61

So, what happens when the class label is not observed?

- List all possible assignments, pick best

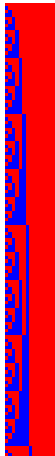


Hidden class label

Slide 10 of 61

So, what happens when the class label is not observed?

- ▶ List all possible assignments, pick best

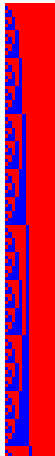


Hidden class label

Slide 10 of 61

So, what happens when the class label is not observed?

- ▶ List all possible assignments, pick best

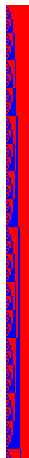


Hidden class label

Slide 10 of 61

So, what happens when the class label is not observed?

- List all possible assignments, pick best



Hidden class label

Slide 10 of 61

So, what happens when the class label is not observed?

- ▶ List all possible assignments, pick best



Hidden class label

Slide 10 of 61

So, what happens when the class label is not observed?

- ▶ List all possible assignments, pick best



Hidden class label

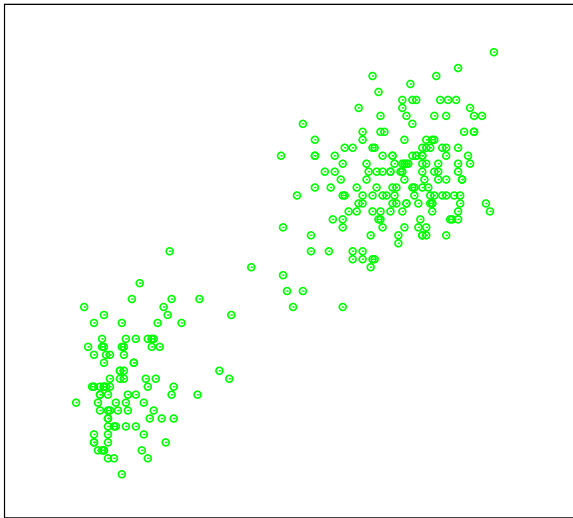
Slide 10 of 61

So, what happens when the class label is not observed?

- ▶ List all possible assignments, pick best

Example: kMeans clustering

Cluster data in two groups in an unsupervised manner:



In order to cluster the data, we need:

- ▶ Some representation of what a cluster looks like
 - ▶ Let's assume for now that each cluster is fully defined by its centre.
- ▶ An assignment of each datapoint to one of the clusters
 - ▶ Let's assume that this is defined by the Euclidean distance to the clusters' centres.

The best configuration is the one where all datapoints are as close as possible to their cluster's centre

In order to cluster the data, we need:

- ▶ Some representation of what a cluster looks like
 - ▶ Let's assume for now that each cluster is fully defined by its centre.
- ▶ An assignment of each datapoint to one of the clusters
 - ▶ Let's assume that this is defined by the Euclidean distance to the clusters' centres.

The best configuration is the one where all datapoints are as close as possible to their cluster's centre

In order to cluster the data, we need:

- ▶ Some representation of what a cluster looks like
 - ▶ Let's assume for now that each cluster is fully defined by its centre.
- ▶ An assignment of each datapoint to one of the clusters
 - ▶ Let's assume that this is defined by the Euclidean distance to the clusters' centres.

The best configuration is the one where all datapoints are as close as possible to their cluster's centre

An exhaustive search for the optimal clustering is intractable and requires C^N operations

- ▶ Where C is the number of clusters
- ▶ and N is the number of datapoints.

How do we find the optimal clustering without exhaustive search? Solve the clustering iteratively:

- ▶ Initialise the cluster means at random
- ▶ Repeat until convergence
 1. Assign each data point to the closest cluster mean
 2. Update each cluster's centre according to the associated data

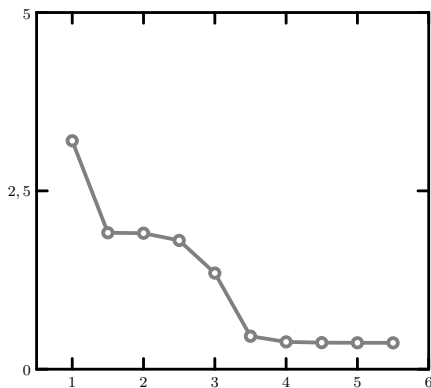
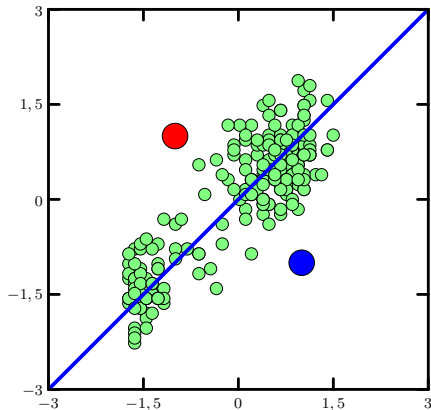
An exhaustive search for the optimal clustering is intractable and requires C^N operations

- ▶ Where C is the number of clusters
- ▶ and N is the number of datapoints.

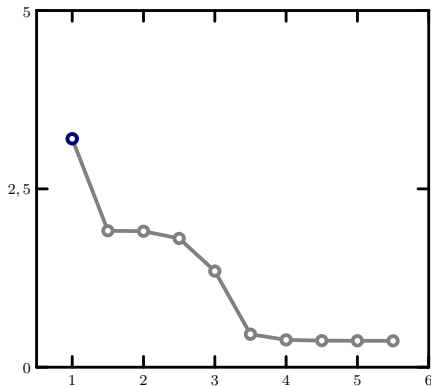
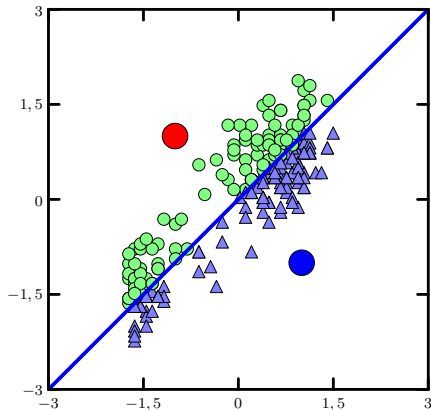
How do we find the optimal clustering without exhaustive search? Solve the clustering iteratively:

- ▶ Initialise the cluster means at random
- ▶ Repeat until convergence
 1. Assign each data point to the closest cluster mean
 2. Update each cluster's centre according to the associated data

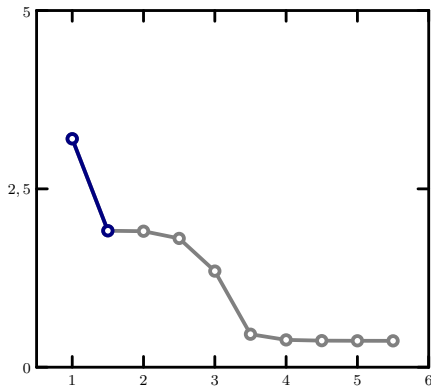
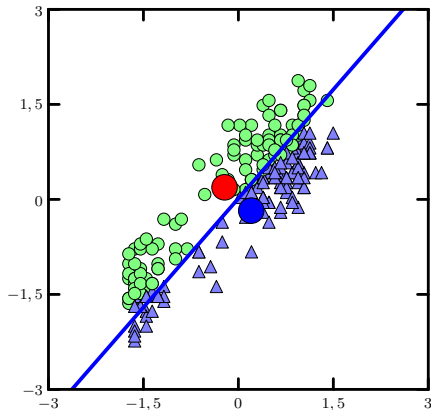
Example: kMeans Clustering



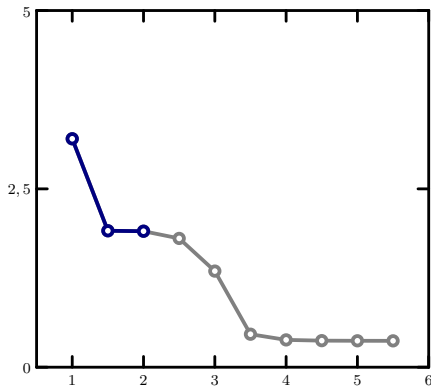
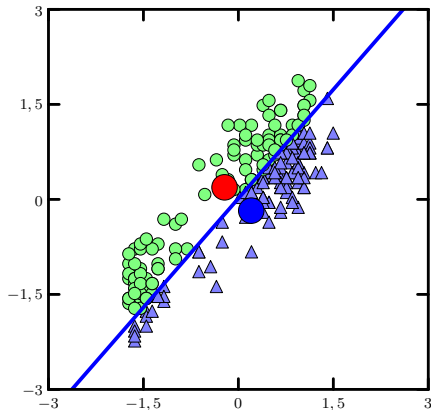
Example: kMeans Clustering



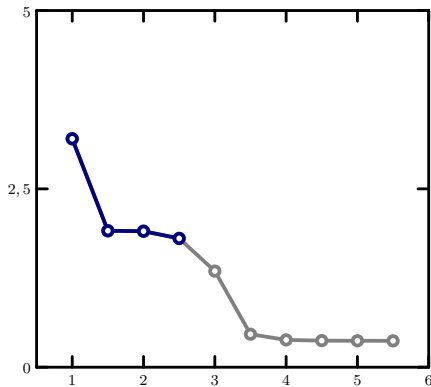
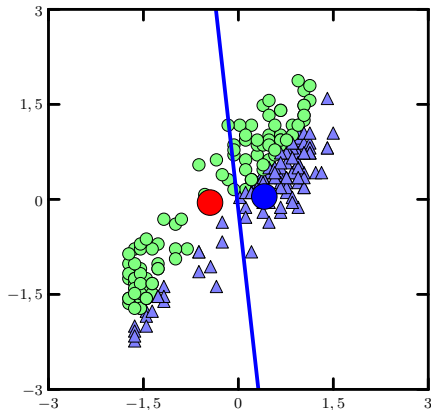
Example: kMeans Clustering



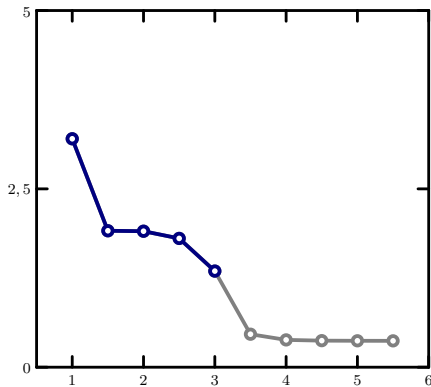
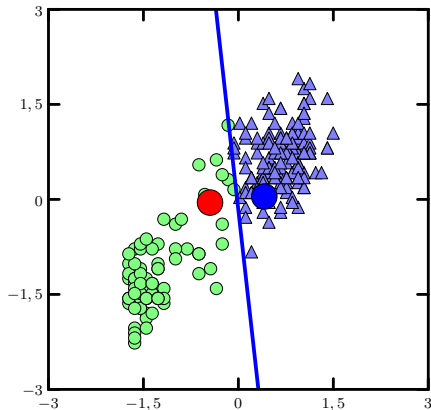
Example: kMeans Clustering



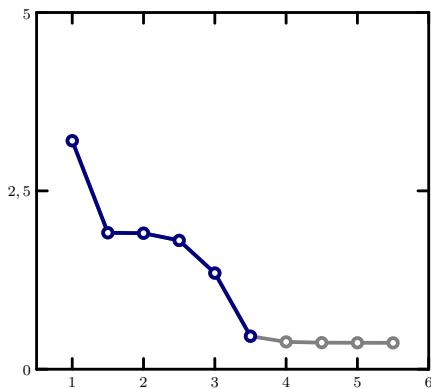
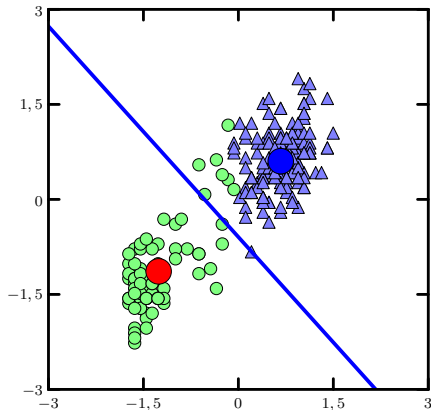
Example: kMeans Clustering



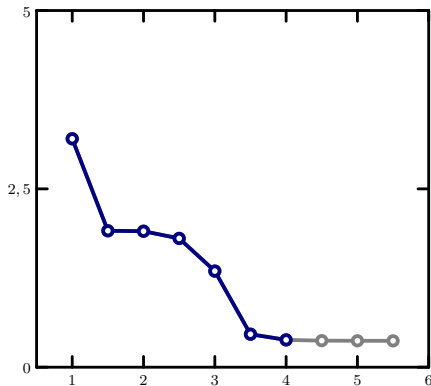
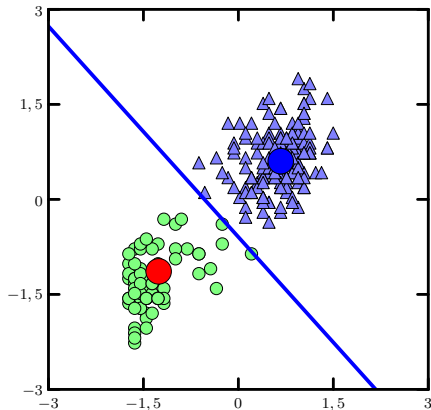
Example: kMeans Clustering



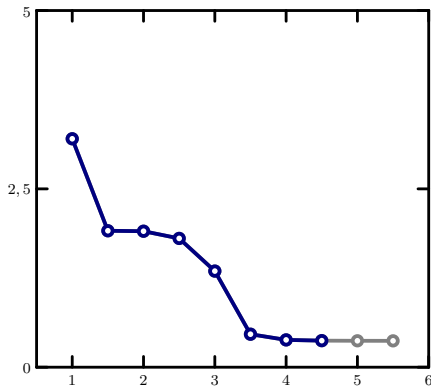
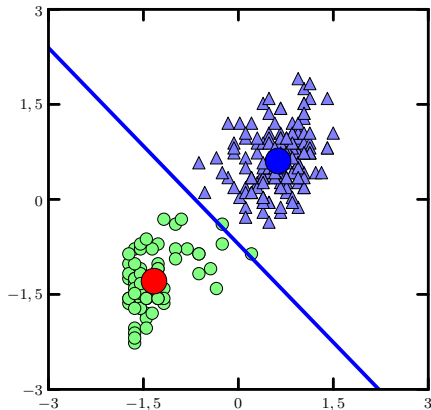
Example: kMeans Clustering



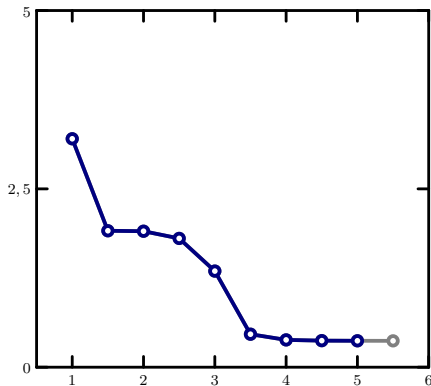
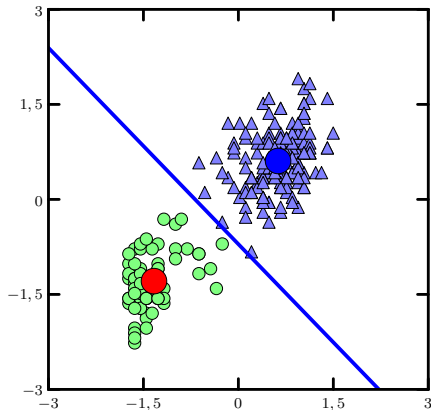
Example: kMeans Clustering



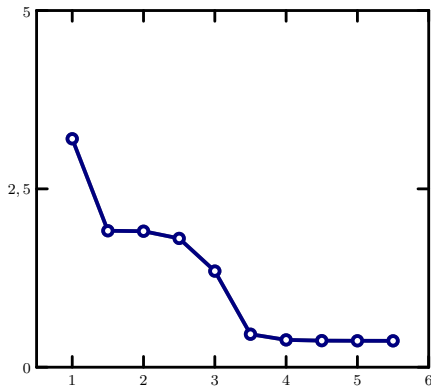
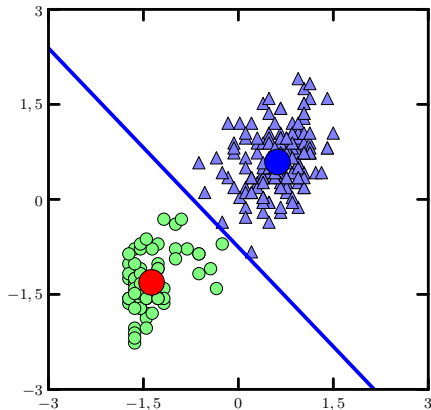
Example: kMeans Clustering



Example: kMeans Clustering



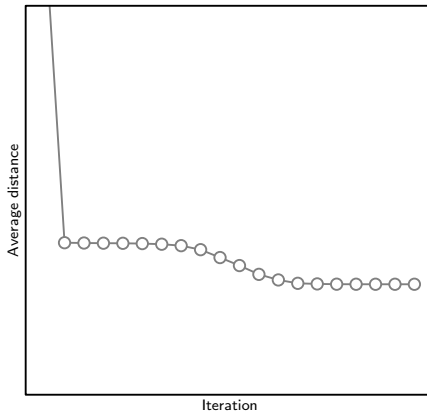
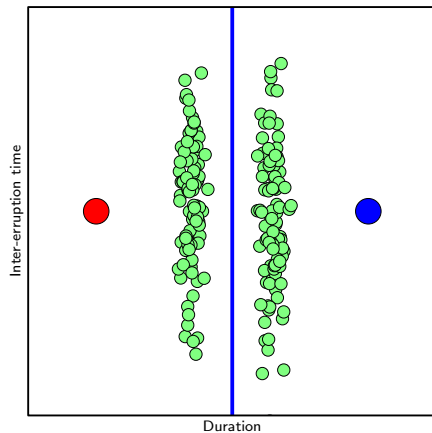
Example: kMeans Clustering



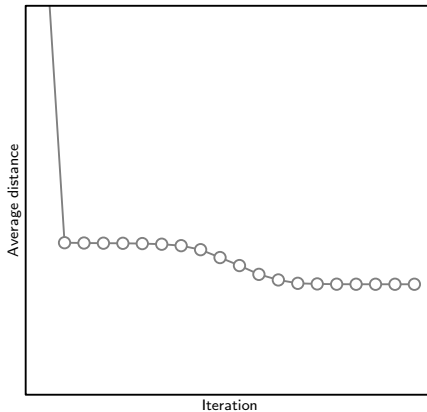
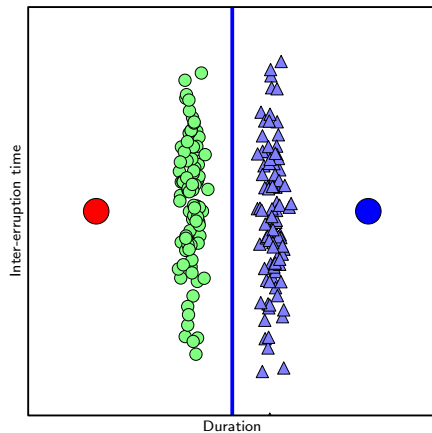
There are some disadvantages to kMeans:

- ▶ **Euclidean distance:**
 - ▶ Only useful for some types of data
 - ▶ Not robust to outliers
 - ▶ Sensitive to scaling of data
 - ▶ Solution: Other distance measures
- ▶ **Hard assignments** At each iteration, each datapoint is assigned to exactly one cluster, even for doubtful cases.

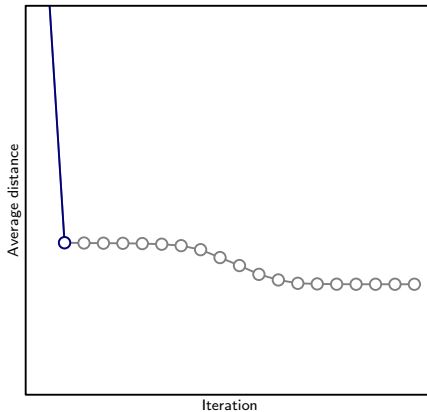
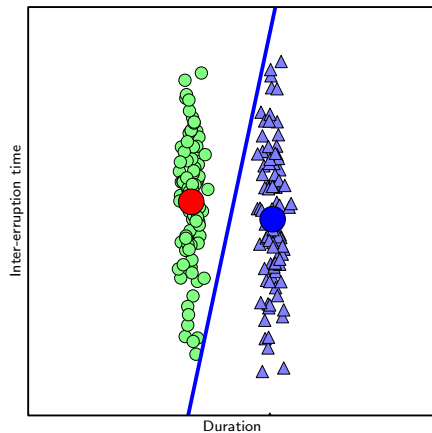
Example: kMeans Clustering



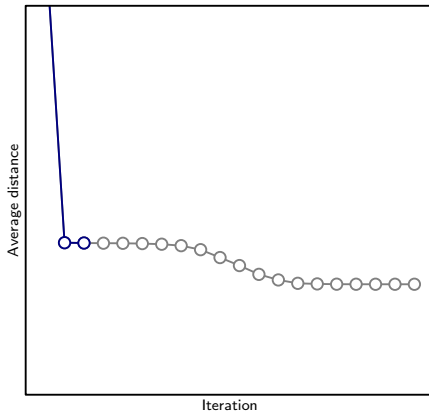
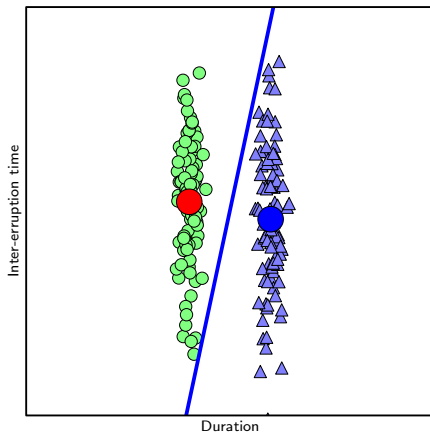
Example: kMeans Clustering



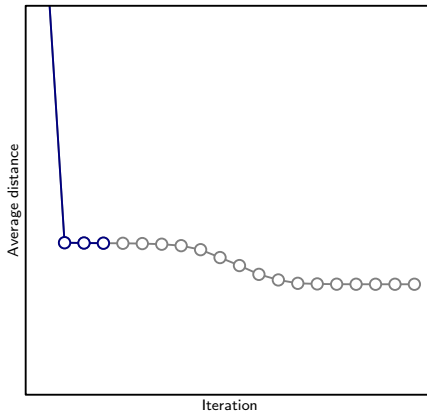
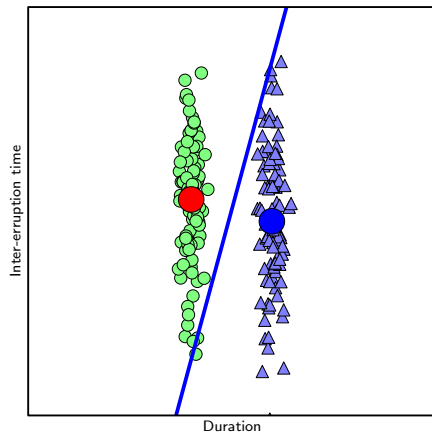
Example: kMeans Clustering



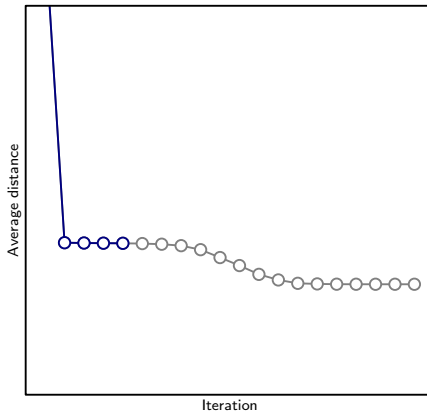
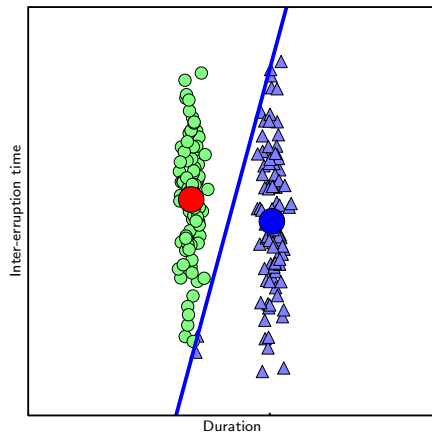
Example: kMeans Clustering



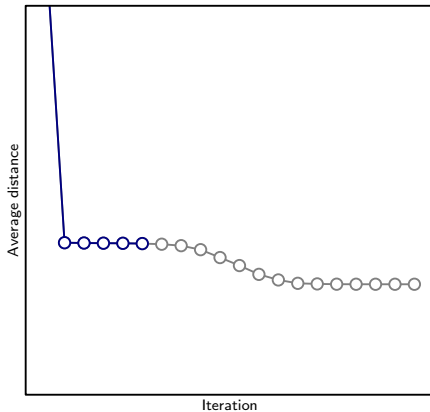
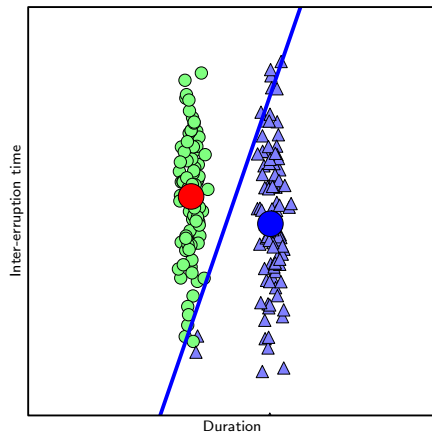
Example: kMeans Clustering



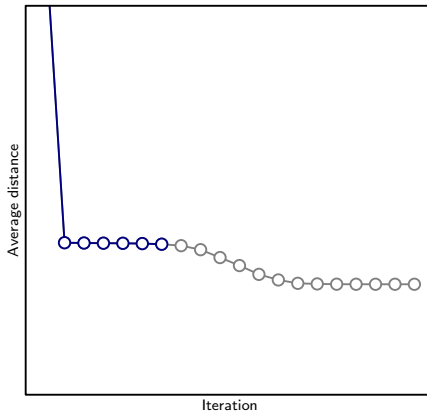
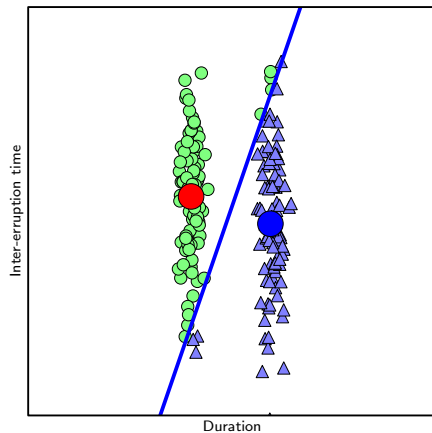
Example: kMeans Clustering



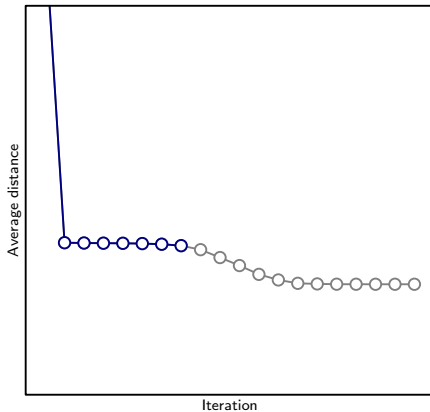
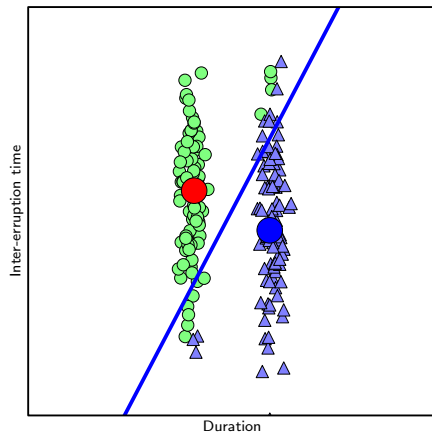
Example: kMeans Clustering



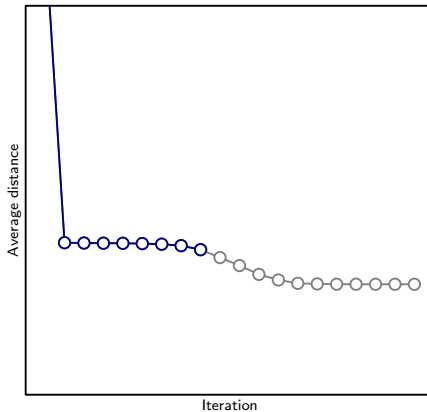
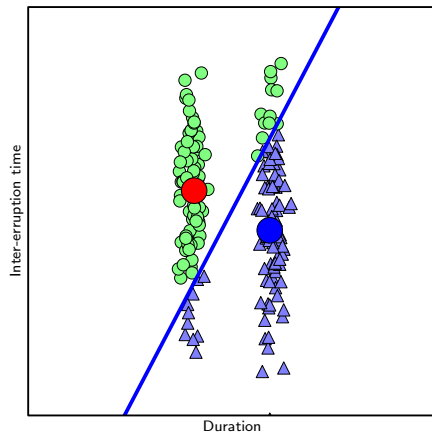
Example: kMeans Clustering



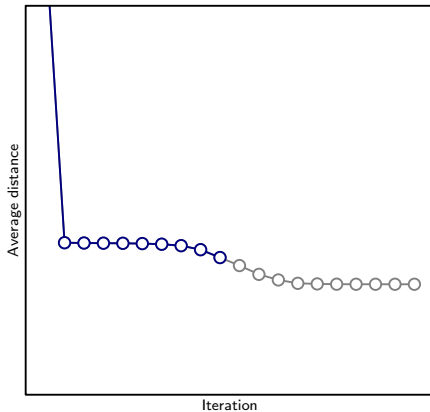
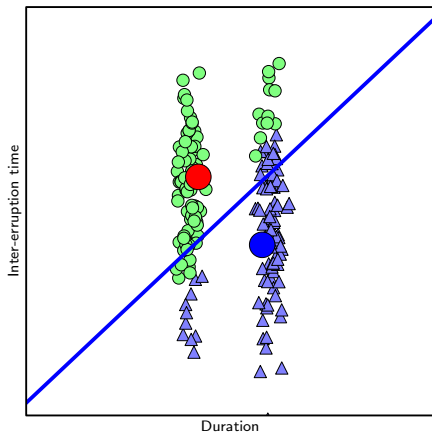
Example: kMeans Clustering



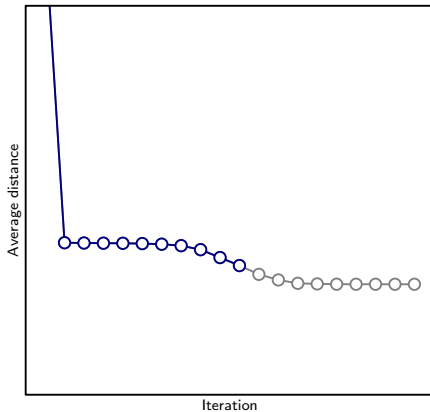
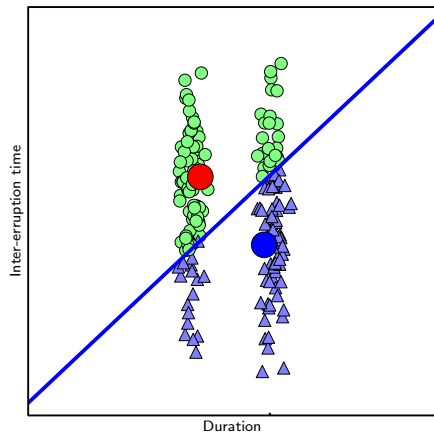
Example: kMeans Clustering



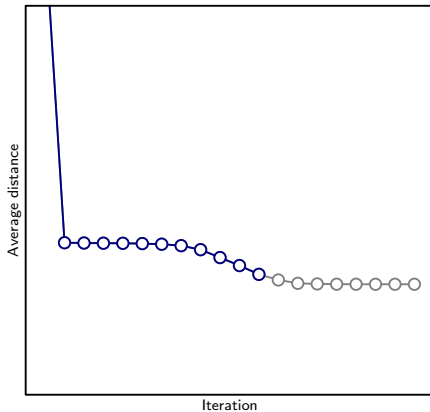
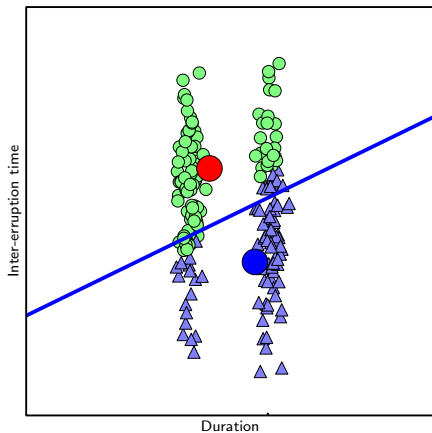
Example: kMeans Clustering



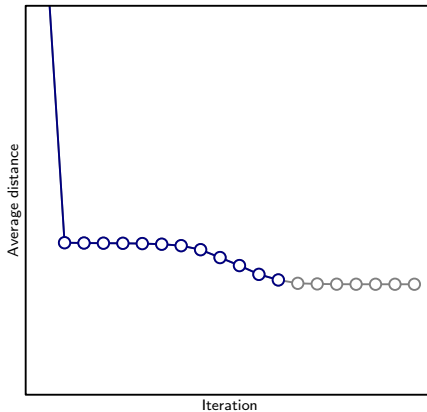
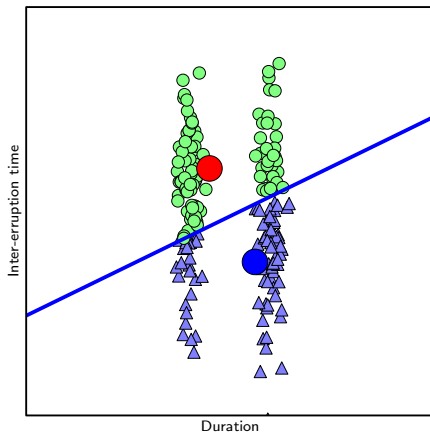
Example: kMeans Clustering



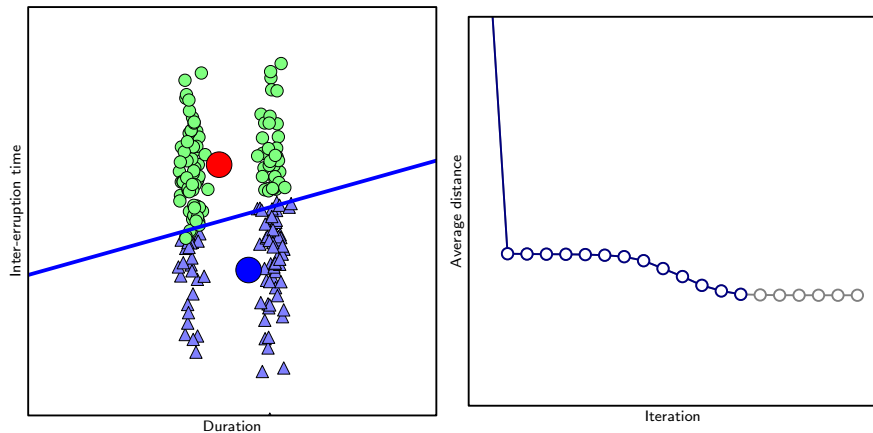
Example: kMeans Clustering



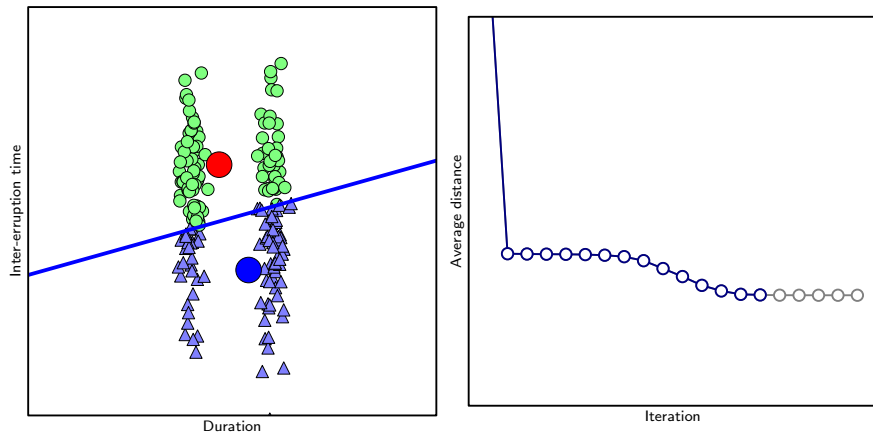
Example: kMeans Clustering



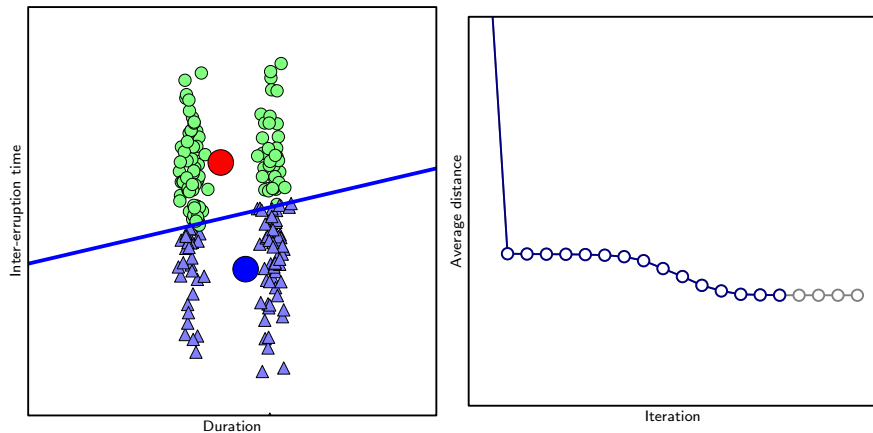
Example: kMeans Clustering



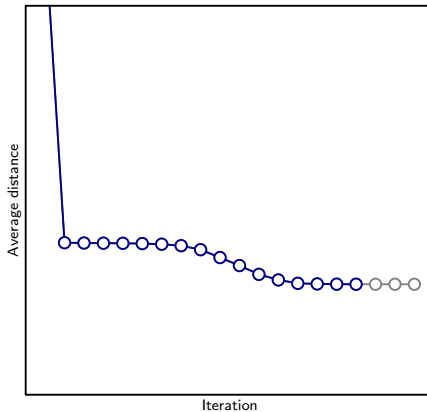
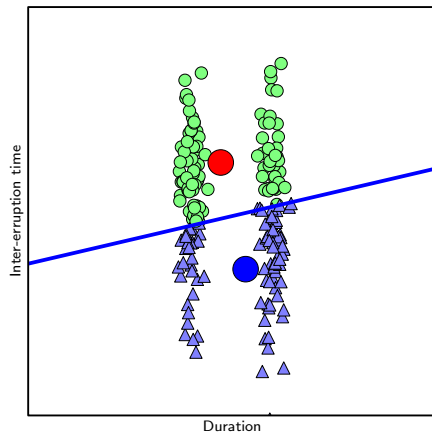
Example: kMeans Clustering



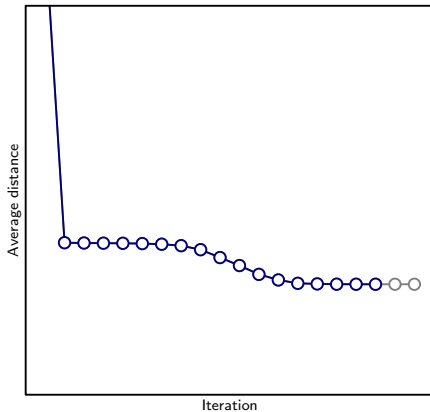
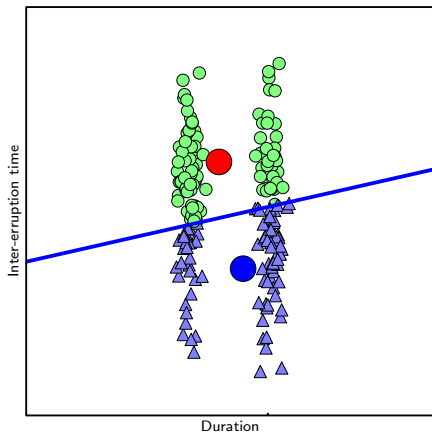
Example: kMeans Clustering



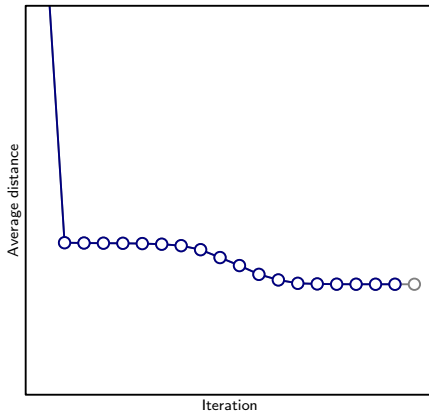
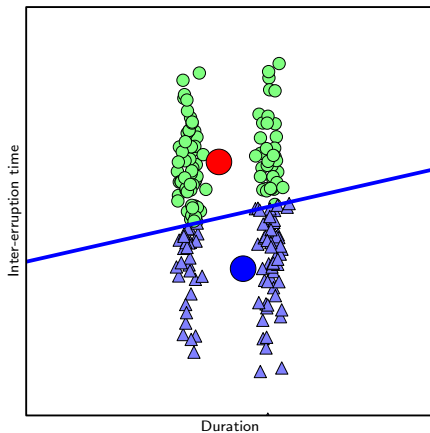
Example: kMeans Clustering



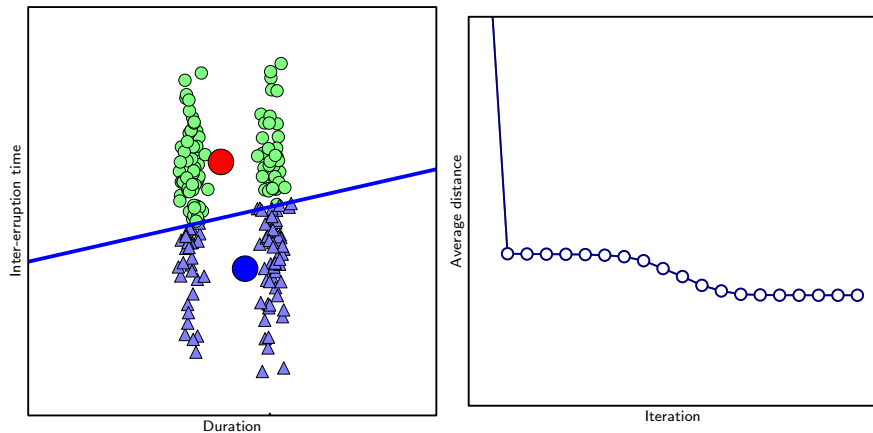
Example: kMeans Clustering



Example: kMeans Clustering



Example: kMeans Clustering

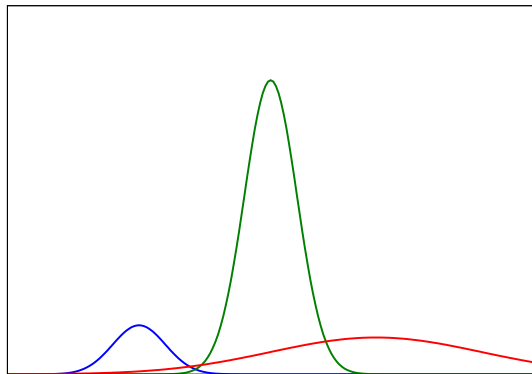


A mixture of Gaussians is a linear combination of Gaussians:

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

where $0 \leq \pi_k \leq 1$ and

$$\sum_{k=1}^K \pi_k = 1 \quad (2)$$

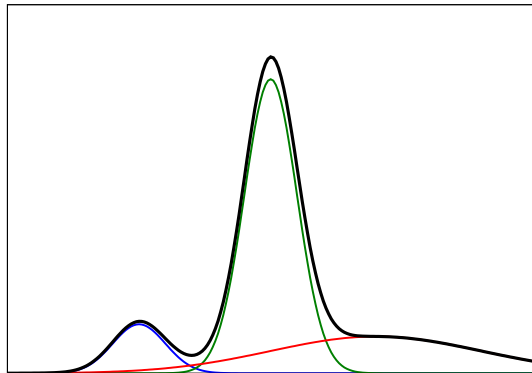


A mixture of Gaussians is a linear combination of Gaussians:

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

where $0 \leq \pi_k \leq 1$ and

$$\sum_{k=1}^K \pi_k = 1 \quad (2)$$

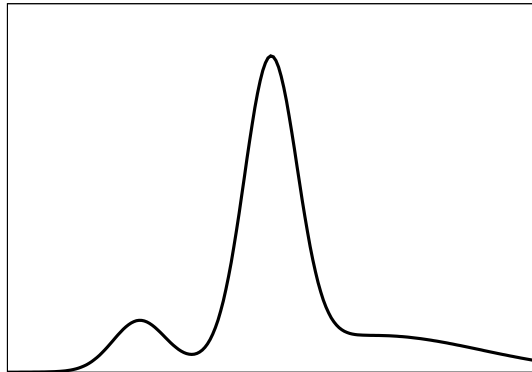


A mixture of Gaussians is a linear combination of Gaussians:

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

where $0 \leq \pi_k \leq 1$ and

$$\sum_{k=1}^K \pi_k = 1 \quad (2)$$



We introduce a binary random variable \mathbf{z} in 1-of-K encoding, the probability of \mathbf{z} can be written as

$$p(z_k = 1) = \pi_k \text{ so that } p(\mathbf{z}) = \prod_{k=1}^K \pi_k^{z_k}$$

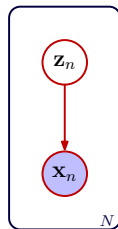
We choose the conditional distribution of \mathbf{x} given a z_k as

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

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

Then we have:

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



Similar to kMeans, finding the optimal parameters for $p(\mathbf{x})$ in a mixture of Gaussians is hard.

However, with our new representation, we can now work with $p(\mathbf{x}, \mathbf{z})$ rather than $p(\mathbf{x})$. In particular, consider the log-likelihood of N datapoints \mathbf{X} :

$$\ln p(\mathbf{X}|\pi, \boldsymbol{\mu}, \boldsymbol{\Sigma}) = \sum_{n=1}^N \ln \left[\sum_{k=1}^K \pi_k \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k) \right] \quad (3)$$

Setting the first derivative with respect to μ_k equal to zero, gives:

$$0 = - \sum_{n=1}^N \frac{\pi_k \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)}{\sum_{\ell} \pi_{\ell} \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_{\ell}, \boldsymbol{\Sigma}_{\ell})} \boldsymbol{\Sigma}_k^{-1} (\mathbf{x}_n - \boldsymbol{\mu}_k) \quad (4)$$

Similar to kMeans, finding the optimal parameters for $p(\mathbf{x})$ in a mixture of Gaussians is hard.

However, with our new representation, we can now work with $p(\mathbf{x}, \mathbf{z})$ rather than $p(\mathbf{x})$. In particular, consider the log-likelihood of N datapoints \mathbf{X} :

$$\ln p(\mathbf{X}|\pi, \boldsymbol{\mu}, \boldsymbol{\Sigma}) = \sum_{n=1}^N \ln \left[\sum_{k=1}^K \pi_k \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k) \right] \quad (3)$$

Setting the first derivative with respect to μ_k equal to zero, gives:

$$0 = - \sum_{n=1}^N \underbrace{\frac{\pi_k \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)}{\sum_{\ell} \pi_{\ell} \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_{\ell}, \boldsymbol{\Sigma}_{\ell})}}_{p(z_k | \mathbf{x}_n)} \boldsymbol{\Sigma}_k^{-1} (\mathbf{x}_n - \boldsymbol{\mu}_k) \quad (4)$$

Using the notation $N_k = \sum_{n=1}^N p(z_k|\mathbf{x}_n)$, we obtain that the values for $\boldsymbol{\mu}_k$, $\boldsymbol{\Sigma}_k$ and π_k that maximise the likelihood are:

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

$$\boldsymbol{\mu}_k = \frac{1}{N_k} \sum_{n=1}^N p(z_k|\mathbf{x}_n) \mathbf{x}_n$$

$$\boldsymbol{\Sigma}_k = \frac{1}{N_k} \sum_{n=1}^N p(z_k|\mathbf{x}_n) (\mathbf{x}_n - \boldsymbol{\mu}_k)(\mathbf{x}_n - \boldsymbol{\mu}_k)^\top$$

However $p(z_k|\mathbf{x}_n)$ is a function of $\pi_1, \dots, \pi_K, \boldsymbol{\mu}_1, \dots, \boldsymbol{\mu}_K, \boldsymbol{\Sigma}_1, \dots, \boldsymbol{\Sigma}_K$

Using the notation $N_k = \sum_{n=1}^N p(z_k|\mathbf{x}_n)$, we obtain that the values for $\boldsymbol{\mu}_k$, $\boldsymbol{\Sigma}_k$ and π_k that maximise the likelihood are:

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

$$\boldsymbol{\mu}_k = \frac{1}{N_k} \sum_{n=1}^N p(z_k|\mathbf{x}_n) \mathbf{x}_n$$

$$\boldsymbol{\Sigma}_k = \frac{1}{N_k} \sum_{n=1}^N p(z_k|\mathbf{x}_n) (\mathbf{x}_n - \boldsymbol{\mu}_k)(\mathbf{x}_n - \boldsymbol{\mu}_k)^\top$$

However $p(z_k|\mathbf{x}_n)$ is a function of $\pi_1, \dots, \pi_K, \boldsymbol{\mu}_1, \dots, \boldsymbol{\mu}_K, \boldsymbol{\Sigma}_1, \dots, \boldsymbol{\Sigma}_K$

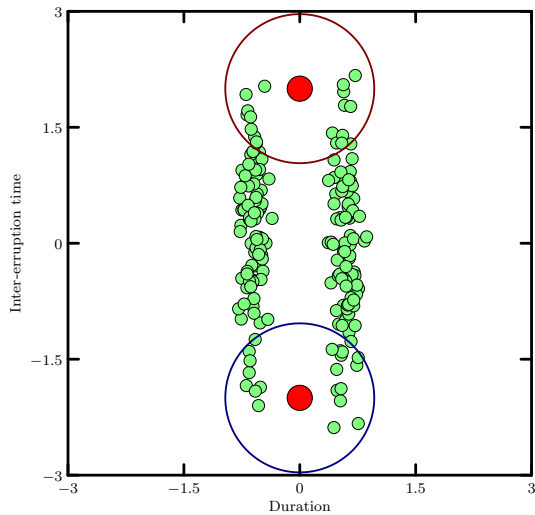
The Expectation-Maximisation algorithm is an iterative update where:

E-step Compute the posterior probabilities of the latent variables given the data and the current parameters (also called *responsibilities*), $p(z_k | \mathbf{x}_n, \boldsymbol{\theta})$

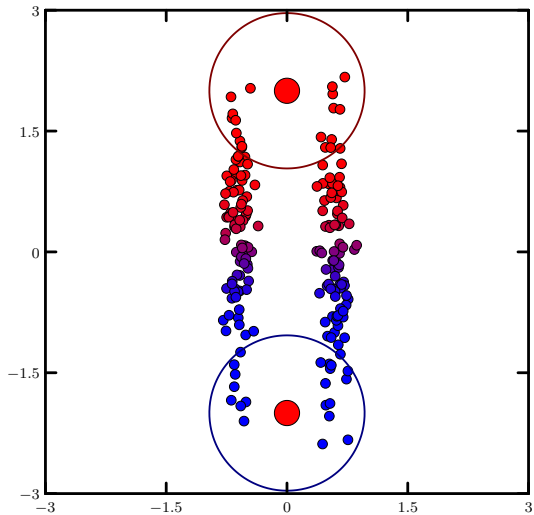
M-step Optimise the expectation of the complete log-likelihood with respect to the parameters

In practice, stop when the increase in likelihood falls below a certain threshold.

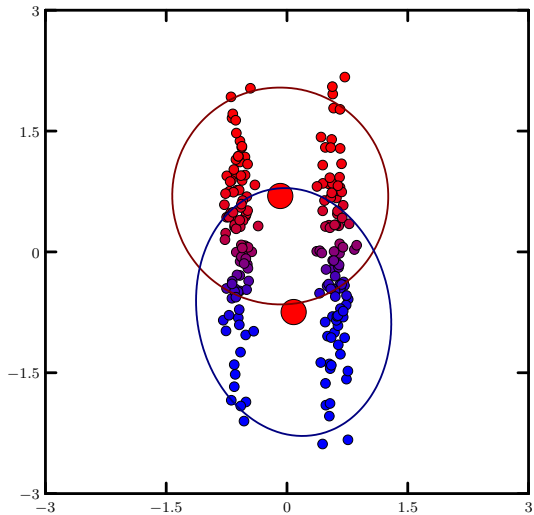
Example: Mixtures of Gaussians



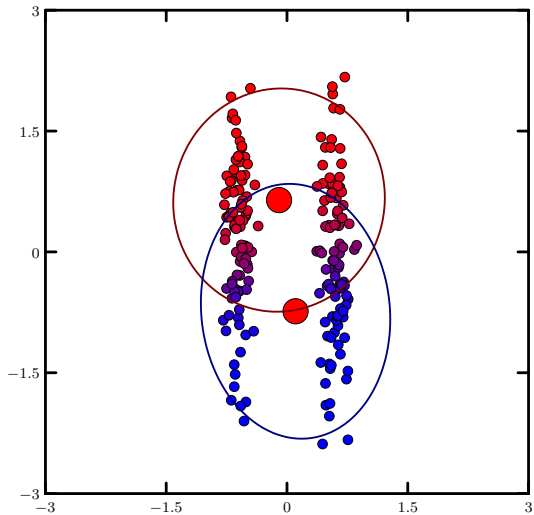
Example: Mixtures of Gaussians



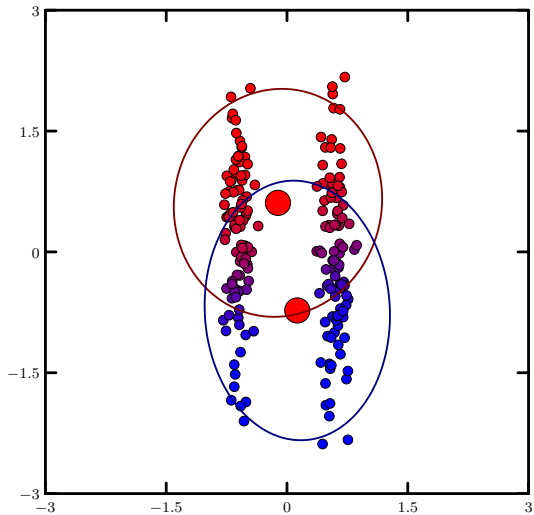
Example: Mixtures of Gaussians



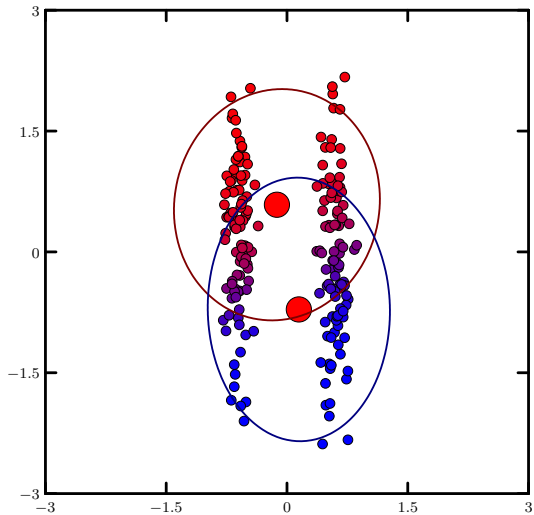
Example: Mixtures of Gaussians



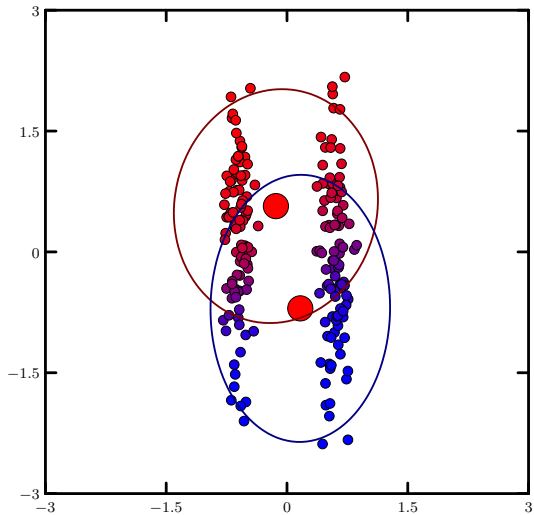
Example: Mixtures of Gaussians



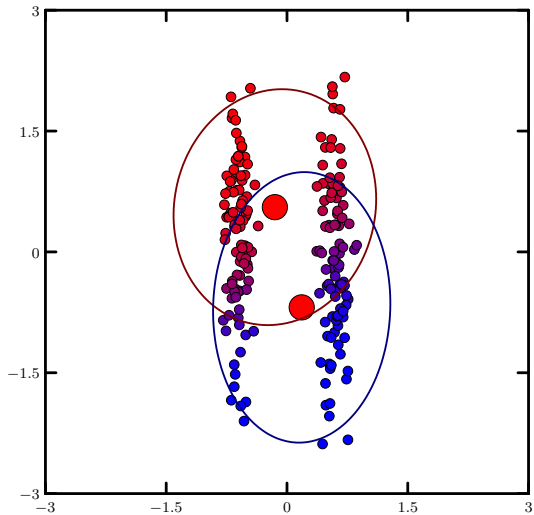
Example: Mixtures of Gaussians



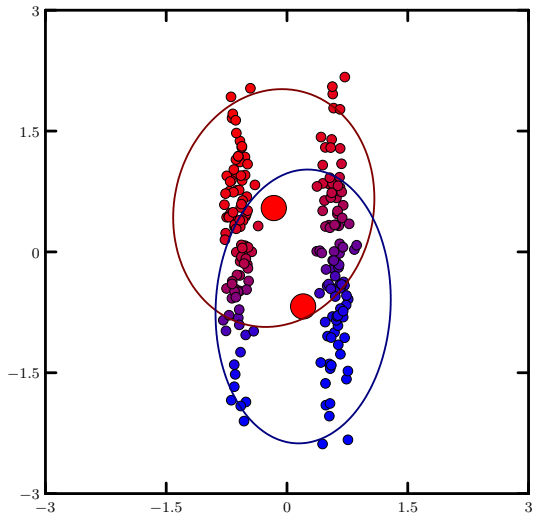
Example: Mixtures of Gaussians



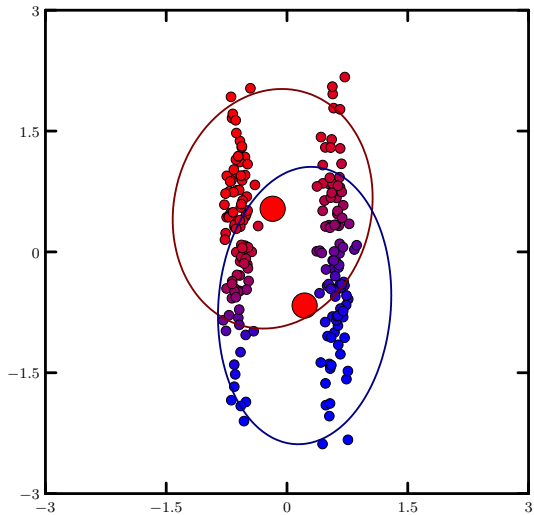
Example: Mixtures of Gaussians



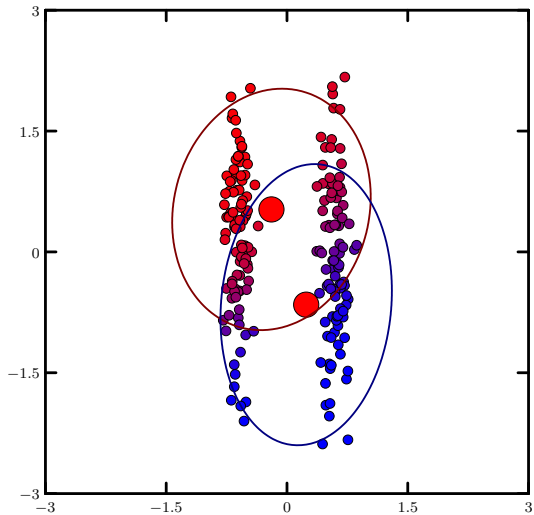
Example: Mixtures of Gaussians



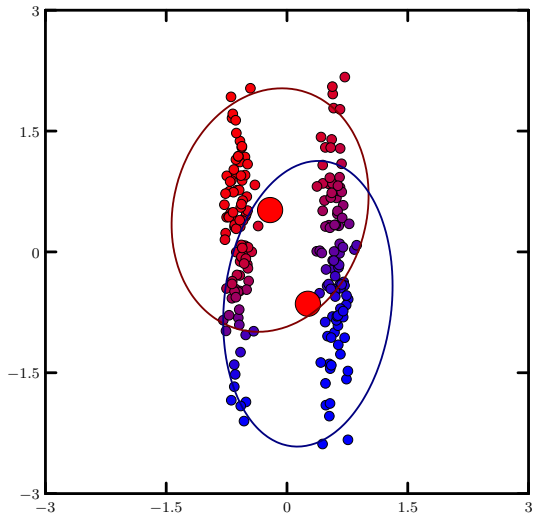
Example: Mixtures of Gaussians



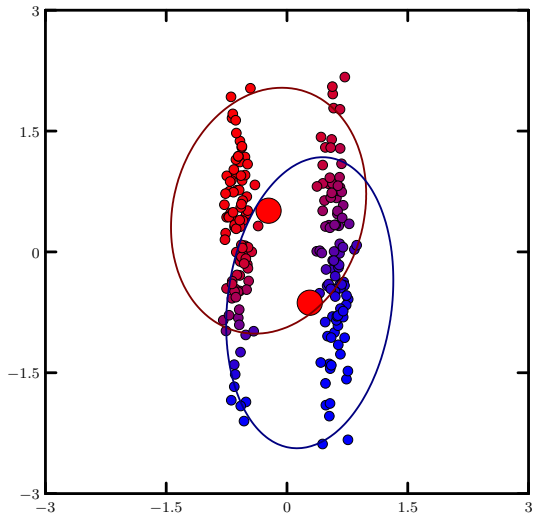
Example: Mixtures of Gaussians



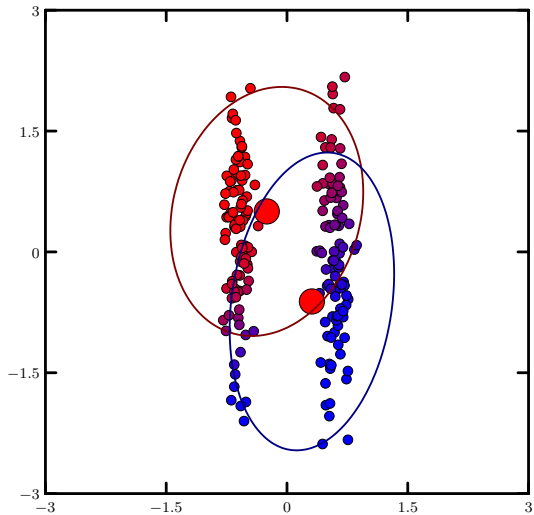
Example: Mixtures of Gaussians



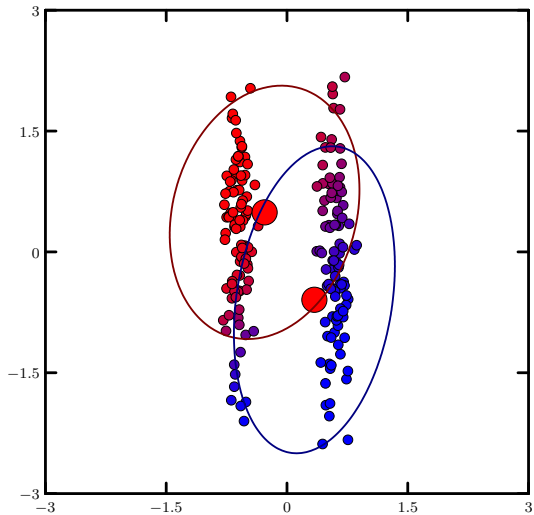
Example: Mixtures of Gaussians



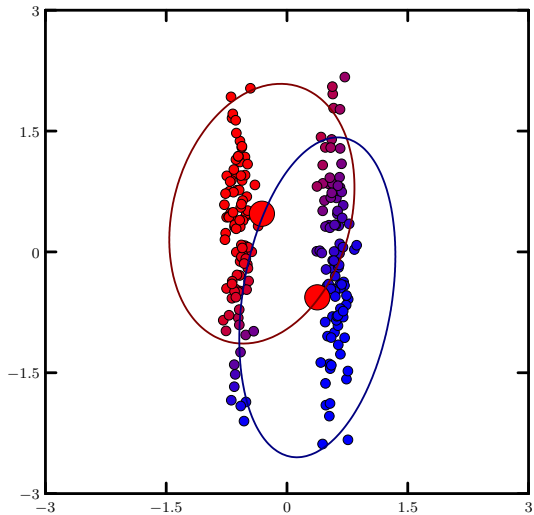
Example: Mixtures of Gaussians



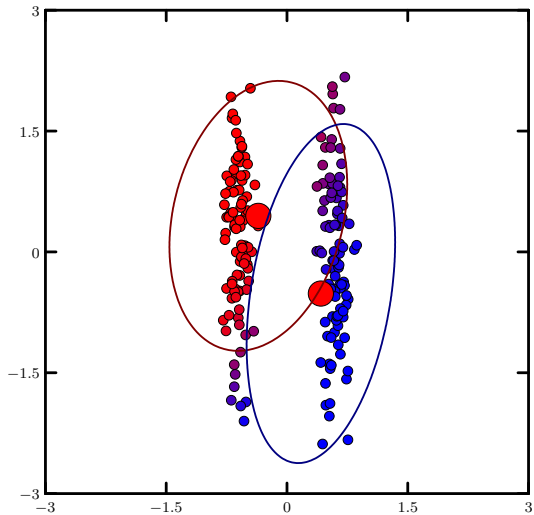
Example: Mixtures of Gaussians



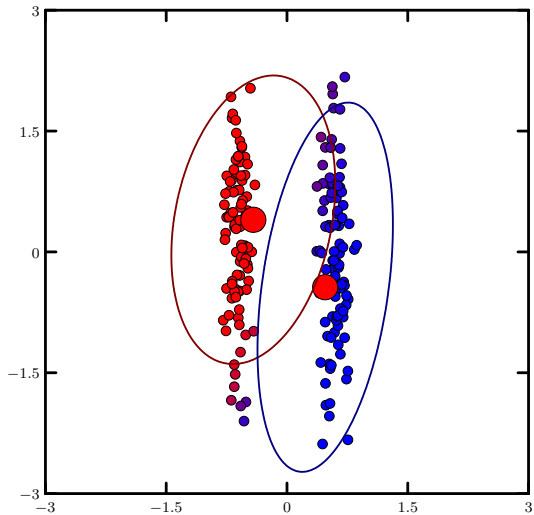
Example: Mixtures of Gaussians



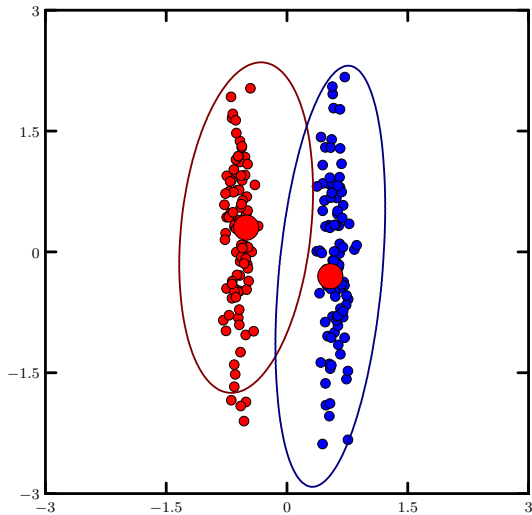
Example: Mixtures of Gaussians



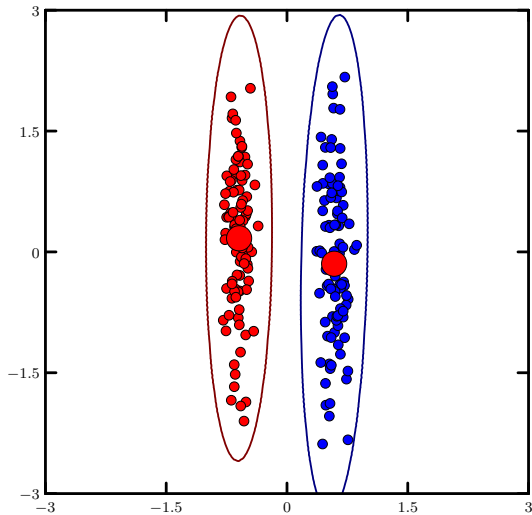
Example: Mixtures of Gaussians



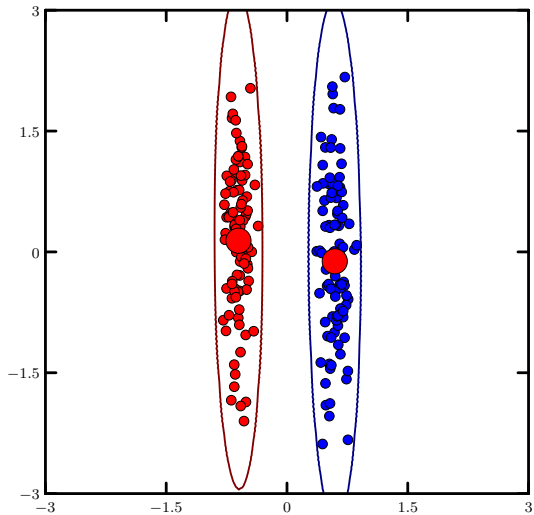
Example: Mixtures of Gaussians



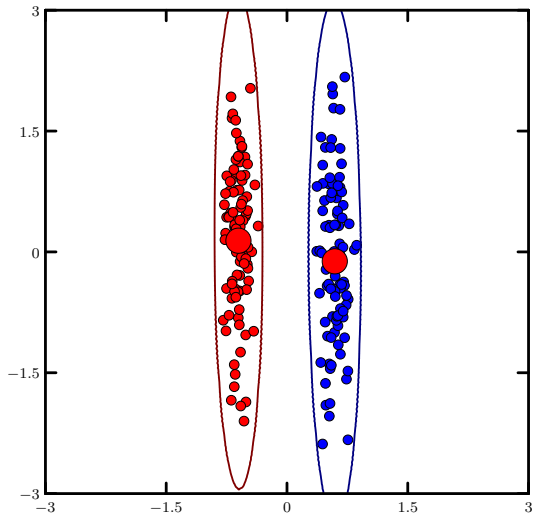
Example: Mixtures of Gaussians



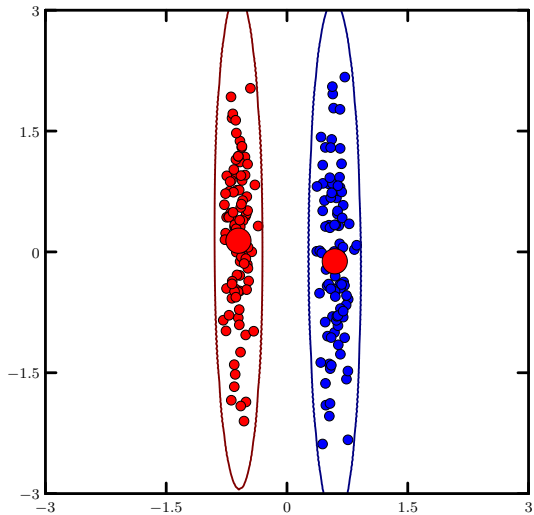
Example: Mixtures of Gaussians



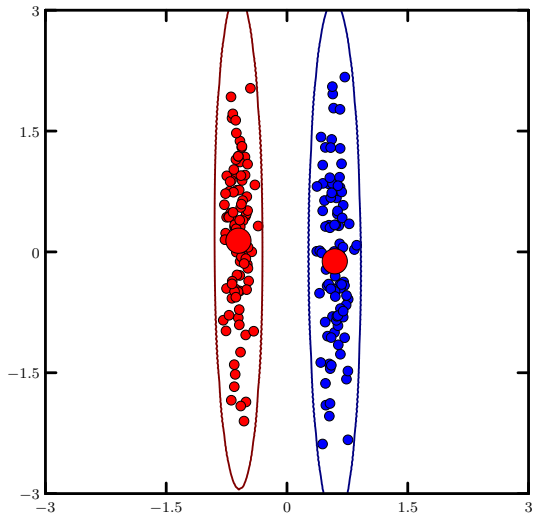
Example: Mixtures of Gaussians



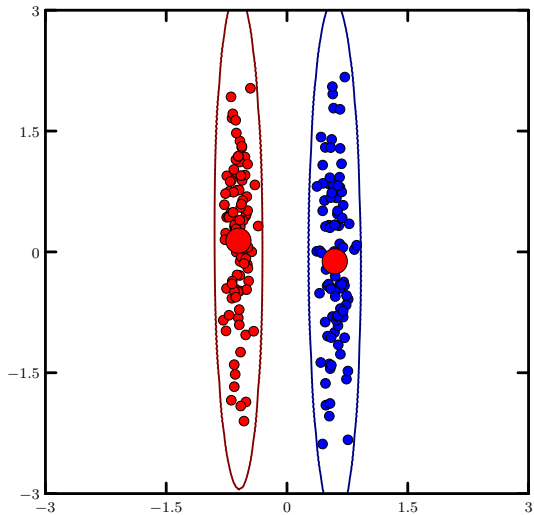
Example: Mixtures of Gaussians



Example: Mixtures of Gaussians



Example: Mixtures of Gaussians



Some things to notice

- ▶ The problem is ill-posed: consider a component k with covariance σI . If the mean of the component falls exactly on a data point, its contribution to the likelihood is

$$\ln p(\mathbf{x}_n | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k) = \frac{1}{\sqrt{2\pi\sigma}} \quad (5)$$

which, in the limit of $\sigma \rightarrow 0$ goes to infinity.

- ▶ A suitable prior on $\boldsymbol{\theta}$ avoids this problem
- ▶ The kMeans algorithm is equivalent with EM for a Gaussian mixture model, where the covariance is σI for *all* mixture components in the limit $\sigma \rightarrow 0$.

Introduction

kMeans Clustering

Mixtures of Gaussians

The General EM Algorithm

The E-step and M-step

Why it works

Models of data sequences

Hidden Markov Models

Linear Dynamical System

In general, the EM algorithm is defined as follows. Optimising the *complete* log-likelihood

$$p(\mathbf{X}, \mathbf{Z} | \theta) \tag{6}$$

would be easy, but we only observe \mathbf{X} . So let's optimise our best estimate of the complete log-likelihood: the expectation of the complete log-likelihood under our current parameter estimates θ^{old} :

$$\mathcal{Q}(\theta, \theta^{old}) = \mathbb{E}_{p(\mathbf{Z} | \mathbf{X}, \theta^{old})} [\ln p(\mathbf{X}, \mathbf{Z} | \theta)] \tag{7}$$

$$= \sum_{\mathbf{z}} p(\mathbf{Z} | \mathbf{X}, \theta) \ln p(\mathbf{X}, \mathbf{Z} | \theta) \tag{8}$$

- In the E-step we evaluate the distribution of the latent variables

$$p(\mathbf{Z}|\mathbf{X}, \boldsymbol{\theta}^{old}) \quad (9)$$

so that we can compute the expectation of the complete log-likelihood

- In the M-step we maximise the complete log-likelihood with respect to $\boldsymbol{\theta}$

$$\boldsymbol{\theta}^{new} \leftarrow \arg \max_{\boldsymbol{\theta}} Q(\boldsymbol{\theta}, \boldsymbol{\theta}^{old}) \quad (10)$$

If we incorporate a prior over θ , $p(\theta)$, the EM algorithm changes slightly:

E-Step : $p(\mathbf{Z}|\mathbf{X}, \theta^{old})$ does not depend on $p(\theta^{old})$ since θ^{old} are given, so the E-step remains identical.

M-Step : We now optimise

$$\mathbb{E} [\ln(p(\mathbf{X}, \mathbf{Z}|\theta)p(\theta))] \quad (11)$$

where $p(\theta)$ is constant with respect to $p(\mathbf{Z}|\mathbf{X}, \theta^{old})$, so that we get:

$$\theta^{new} = \arg \max_{\theta} \mathcal{Q}(\theta, \theta^{old}) + \ln p(\theta) \quad (12)$$

Lower bound: why EM works

Slide 28 of 61

Consider the quantity

$$\begin{aligned}\sum_{\mathbf{z}} q(\mathbf{z}) \ln \left[\frac{p(\mathbf{x}, \mathbf{z} | \boldsymbol{\theta})}{q(\mathbf{z})} \right] \\&= \sum_{\mathbf{z}} q(\mathbf{z}) [\ln p(\mathbf{z} | \mathbf{x}, \boldsymbol{\theta}) + \ln p(\mathbf{x} | \boldsymbol{\theta}) - \ln q(\mathbf{z})] \\&= \sum_{\mathbf{z}} q(\mathbf{z}) \ln \frac{p(\mathbf{z} | \mathbf{x}, \boldsymbol{\theta})}{q(\mathbf{z})} + \ln p(\mathbf{x} | \boldsymbol{\theta}) \underbrace{\sum_{\mathbf{z}} q(\mathbf{z})}_{=1}\end{aligned}$$

so that

$$\ln p(\mathbf{x} | \boldsymbol{\theta}) = \underbrace{\sum_{\mathbf{z}} q(\mathbf{z}) \ln \left[\frac{p(\mathbf{x}, \mathbf{z} | \boldsymbol{\theta})}{q(\mathbf{z})} \right]}_{\mathcal{L}(q, \boldsymbol{\theta})} + \underbrace{\left(- \sum_{\mathbf{z}} q(\mathbf{z}) \ln \frac{p(\mathbf{z} | \mathbf{x}, \boldsymbol{\theta})}{q(\mathbf{z})} \right)}_{KL(q || p)}$$

The KL-divergence is positive

Consider the log likelihood

Slide 29 of 61

$$1 = \sum_{\mathbf{z}} p(\mathbf{z}) = \sum_{\mathbf{z}} q(\mathbf{z}) \frac{p(\mathbf{z})}{q(\mathbf{z})}$$

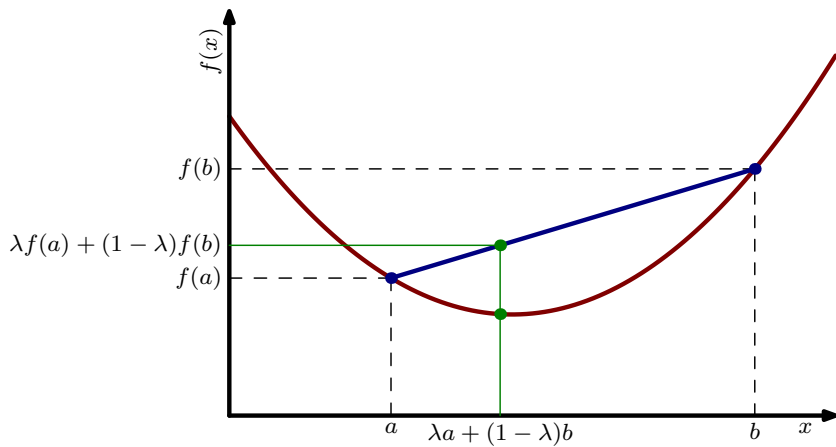
$$\ln 1 = \ln \sum_{\mathbf{z}} q(\mathbf{z}) \frac{p(\mathbf{z})}{q(\mathbf{z})}$$

$$0 \geq \sum_{\mathbf{z}} q(\mathbf{z}) \ln \frac{p(\mathbf{z})}{q(\mathbf{z})}$$

$$0 \leq - \sum_{\mathbf{z}} q(\mathbf{z}) \ln \frac{p(\mathbf{z})}{q(\mathbf{z})}$$

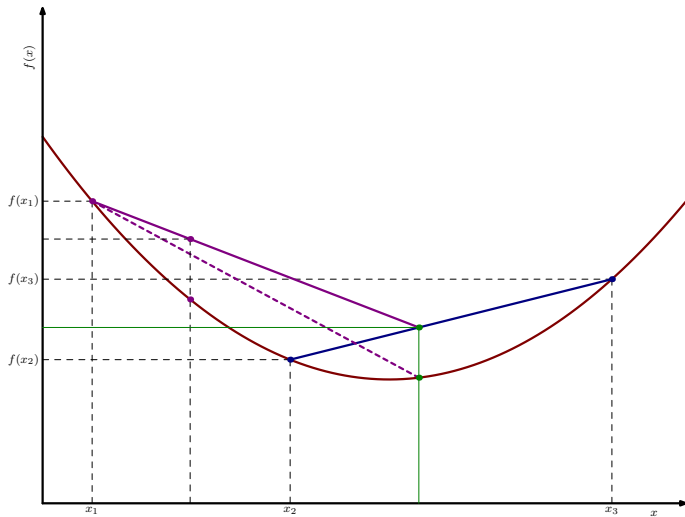
By *Jensen's inequality*, if $f(x)$ is concave:

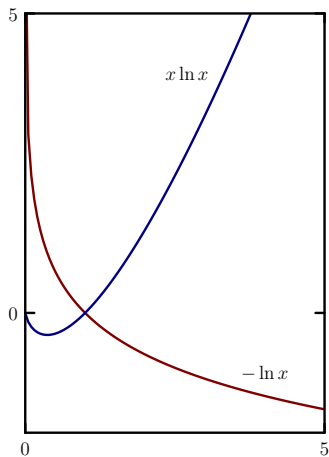
$$\mathbb{E}[f(x)] \geq f(\mathbb{E}[x])$$



Jensen's inequality

Slide 30 of 61





During the E-step, we maximise $\mathcal{L}(q, \theta)$ with respect to $q(\mathbf{z})$, leaving θ^{old} untouched. Since $\ell(\theta^{old})$ does not depend on $q(\mathbf{z})$, this can only be achieved by setting

$$KL(q||p) = 0$$

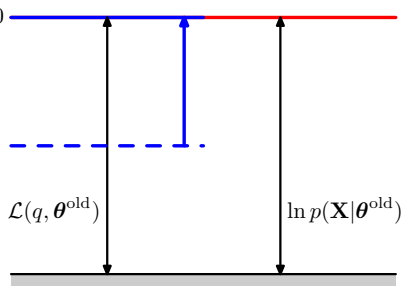
$$KL(q||p) = 0$$

in other words

$$q(\mathbf{z}) = p(\mathbf{z}|\mathbf{x}, \theta)$$

so that

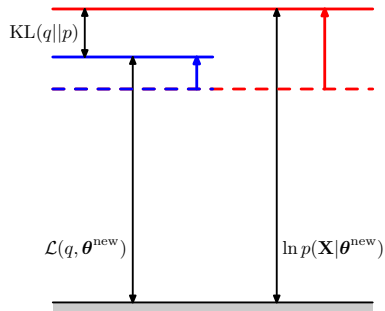
$$\mathcal{L}(q, \theta^{old}) = \ln p(\mathbf{x}|\theta^{old})$$

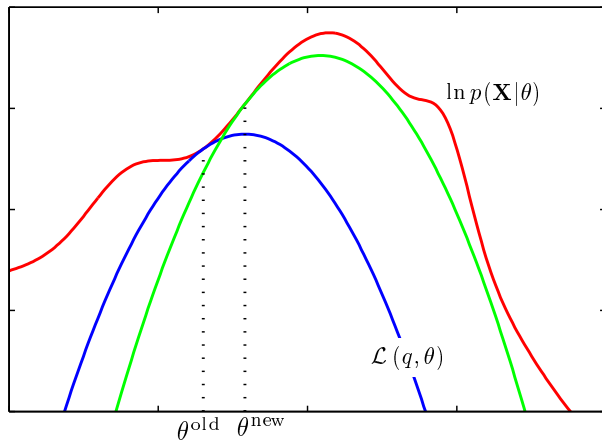


During the M-step, we maximise $\mathcal{L}(q, \theta)$ with respect to θ , leaving $q(\mathbf{z})$ untouched. Since $q(\mathbf{z}) = p(\mathbf{z}|\mathbf{x}, \theta^{old})$,

$$\begin{aligned}\mathcal{L}(q, \theta) &= \sum_{\mathbf{z}} p(\mathbf{z}|\mathbf{x}, \theta^{old}) \ln p(\mathbf{z}, \mathbf{x}|\theta) - \sum_{\mathbf{z}} p(\mathbf{z}|\mathbf{x}, \theta^{old}) \ln p(\mathbf{z}|\mathbf{x}, \theta^{old}) \\ &= \mathcal{Q}(\theta, \theta^{old}) + \text{const}\end{aligned}$$

- ▶ Maximising $\mathcal{L}(q, \theta)$ with respect to θ changes $\ln p(\mathbf{z}|\mathbf{x}, \theta)$, so that $KL(q||p)$ increases.
- ▶ $p(\mathbf{x}|\theta)$ therefore increases at least as much as $\mathcal{L}(q, \theta)$





Sometimes, the following extensions of EM are used:

- ▶ When the datapoints are independent, the responsibilities z_n depend on x_n and θ only, so that the E and M step can be computed online rather than in batch. This can converge faster than the batch version.
- ▶ Sometimes the E-step or M-step (or both) remain intractable. Increasing the likelihood (rather than maximising it) still guarantees increasing the likelihood. This is called Generalised EM (GEM)

Introduction

kMeans Clustering

Mixtures of Gaussians

The General EM Algorithm

The E-step and M-step

Why it works

Models of data sequences

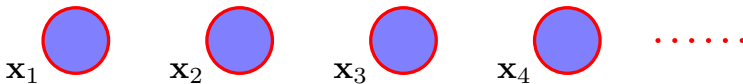
Hidden Markov Models

Linear Dynamical System

Simplest solution:

- ▶ Assume that all observations are independent
- ▶ Problem: fails to exploit sequential patterns in the data

Example: Rainy days



The only information we can obtain from this model is the frequency of rainy days.

If we do not assume any independence, we can factorise the distribution as:

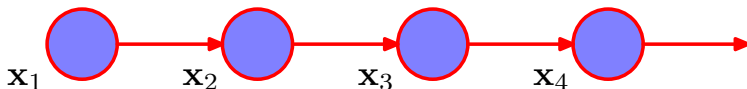
$$p(\mathbf{x}_1, \dots, \mathbf{x}_N) = \prod_{n=1}^N p(\mathbf{x}_n | \mathbf{x}_1, \dots, \mathbf{x}_{n-1}) \quad (13)$$

However:

- ▶ In general the most recent observations are more likely to affect the current observation.
- ▶ If we assume that only the single most recent observation affects the current observation, we obtain the **first order Markov chain**

$$p(\mathbf{x}_1, \dots, \mathbf{x}_N) = p(\mathbf{x}_1) \prod_{n=2}^N p(\mathbf{x}_n | \mathbf{x}_{n-1}) \quad (14)$$

First order Markov Chain

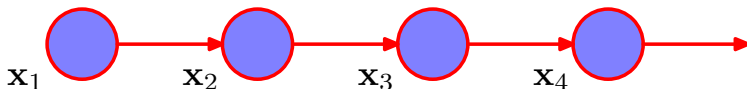


- In most applications, we constrain $p(x_n|x_{n-1})$ to be equal over the length of the chain: **homogeneous markov chain**

Example

This was used to model character occurrences in text

First order Markov Chain



- In most applications, we constrain $p(x_n|x_{n-1})$ to be equal over the length of the chain: **homogeneous markov chain**

Example

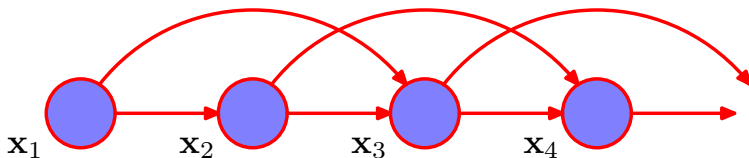
This was used to model character occurrences in text

First order Markov chains are still very restrictive:

- ▶ Modelling trends requires more information from the past
- ▶ Second order Markov chain:

$$p(\mathbf{x}_1, \dots, \mathbf{x}_N) = p(\mathbf{x}_1)p(\mathbf{x}_2|\mathbf{x}_1) \prod_{n=3} p(\mathbf{x}_n|\mathbf{x}_{n-1}, \mathbf{x}_{n-2}) \quad (15)$$

Second order Markov chain



There is a high penalty for the extra flexibility:

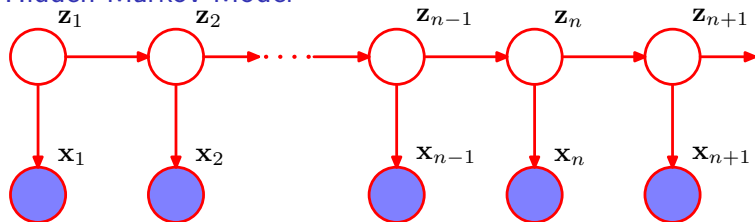
- ▶ Imagine K discrete states
 - ▶ First order chain: $K \times (K - 1)$ independent parameters
 - ▶ Second order chain: $K^2 \times (K - 1)$ parameters
 - ▶ In general, M th order: $K^{M-1} \times (K - 1)$ parameters
- ▶ The complexity grows exponentially with M

This becomes impractical very quickly.

How can we keep a longer memory and still restrict the number of parameters?

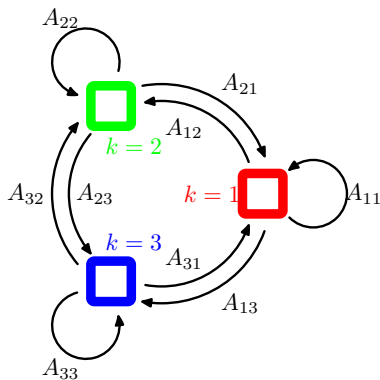
- ▶ Create a latent variable model: the probability of observations depends on the latent variable (see Mixtures of Gaussians)
- ▶ The latent variables are part of a latent Markov chain
- ▶ D-separation shows that all observations are now dependent

Hidden Markov Model

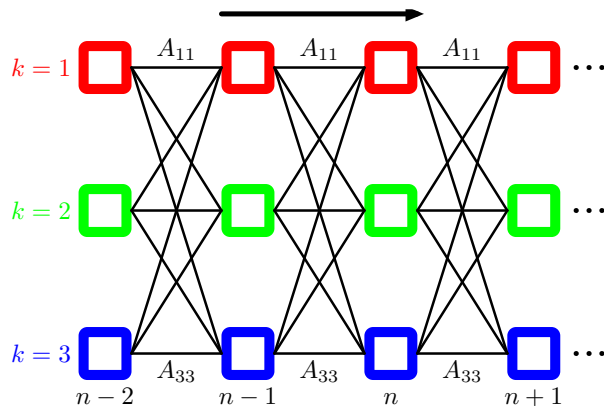


Hidden Markov Models are often viewed as Finite State Machines

- ▶ The arrows indicate the *transition probabilities*, typically denoted by the matrix **A**
- ▶ *Emission probabilities* are associated with each state
- ▶ Not depicted here, is the probability of starting a sequence in a given state k , typically denoted by π_k



We can then “unroll” the FSA over time to form a lattice



Inference allows us to answer the questions:

- ▶ What is the probability of each state, for each observation

Example

In speech recognition, what is the probability that a certain phoneme corresponds to a certain audio observation?

- ▶ What is the probability of seeing an observation sequence, according to the model?

Example

In elderly care, it is useful to monitor how well patients perform activities at home. It is very hard to have a model of the problems they may have, but it is possible to detect when their activities become very unlikely according to the model.

Inference allows us to answer the questions:

- ▶ What is the probability of each state, for each observation

Example

In speech recognition, what is the probability that a certain phoneme corresponds to a certain audio observation?

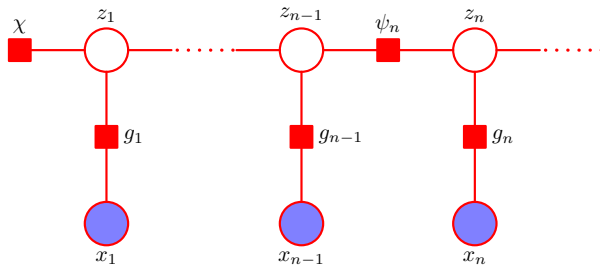
- ▶ What is the probability of seeing an observation sequence, according to the model?

Example

In elderly care, it is useful to monitor how well patients perform activities at home. It is very hard to have a model of the problems they may have, but it is possible to detect when their activities become very unlikely according to the model.

HMM have been around for much longer than graphical models

- ▶ Specific versions of the sum-product and max-product were developed for the HMM
- ▶ The sum-product rule for HMM is known as the Forward-Backward or Baum-Welch algorithm
- ▶ It consists of a single forward and a single backward pass through the chain



The message leaving the first node is

$$\begin{aligned}\mu_{z_1 \rightarrow \psi_1} &= p(x_1 | z_1) p(z_1) \\ &= p(x_1, z_1)\end{aligned}$$

The message leaving the second node is

$$\begin{aligned}\mu_{z_2 \rightarrow \psi_3} &= p(x_2 | z_2) \sum_{z_1} p(x_1, z_1) p(z_2 | z_1) \\ &= p(x_1, x_2, z_2)\end{aligned}$$

In general:

$$\begin{aligned}\mu_{\psi_{n-1} \rightarrow z_n} &= p(x_1, \dots, x_n, z_n) \\ &= \alpha(z_n)\end{aligned}$$

The Backward pass

The message entering the penultimate node is

$$\begin{aligned}\mu_{\psi_{N-1} \rightarrow z_{N-1}} &= \sum_{z_N} p(z_N | z_{N-1}) p(x_N | z_N) \\ &= p(x_N | z_{N-1})\end{aligned}$$

The message entering the antepenultimate node is

$$\begin{aligned}\mu_{\psi_{N-2} \rightarrow z_{N-2}} &= \sum_{z_{N-1}} p(z_{N-1} | z_{N-2}) p(x_{N-1} | z_{N-1}) p(x_N | z_{N-1}) \\ &= p(x_{N-1}, x_N | z_{N-2})\end{aligned}$$

In general, the “backward” message is:

$$\begin{aligned}\mu_{\psi_n \rightarrow z_n} &= p(x_{n+1}, \dots, x_N | z_n) \\ &= \beta(z_n)\end{aligned}$$

The sequence of hidden states often has a real, physical explanation

Example

In speech recognition, the sequence of phonemes for a sequence of sound observations. It is therefore useful to have the most likely single sequence of latent states, rather than the sequence of most likely latent states

- ▶ We obtain this using the max-sum algorithm
- ▶ For the HMM, this is called the Viterbi Algorithm

We can train the HMM using the EM algorithm, optimising

$$\mathcal{Q}(\boldsymbol{\theta}, \boldsymbol{\theta}^{\text{old}}) = \sum_{\mathbf{z}} p(\mathbf{z}|\mathbf{X}, \boldsymbol{\theta}^{\text{old}}) \ln p(\mathbf{X}, \mathbf{z}|\boldsymbol{\theta}) \quad (16)$$

We introduce

$$\gamma(z_{nk}) = p(z_{nk}|\mathbf{X}, \boldsymbol{\theta}^{\text{old}}) \quad (17)$$

$$\xi(z_{n-1j}, z_{n,k}) = p(z_{n-1j}, z_{n,k}|\mathbf{X}, \boldsymbol{\theta}^{\text{old}}) \quad (18)$$

so that the expectation of the complete log-likelihood becomes

$$\begin{aligned} \mathcal{Q}(\boldsymbol{\theta}, \boldsymbol{\theta}^{\text{old}}) = & \sum_k \gamma(z_{1k}) \ln \pi_k + \sum_{n=2}^N \sum_{j=1}^K \sum_{k=1}^K \xi(z_{n-1j}, z_{n,k}) \ln A_{jk} + \\ & \sum_{n=1}^N \sum_{k=1}^K \gamma(z_{nk}) \ln p(\mathbf{x}_n|z_k) \end{aligned}$$

The *responsibilities* $\gamma(z_{nk})$ can be computed from the forward and backward probabilities:

$$\gamma(\mathbf{z}_n) = p(\mathbf{z}_n | \mathbf{X}, \theta^{\text{old}}) \quad (19)$$

$$= \frac{\alpha(\mathbf{z}_n)\beta(\mathbf{z}_n)}{p(\mathbf{X})} \quad (20)$$

and similarly for the pairwise responsibilities:

$$\xi(\mathbf{z}_{n-1}, \mathbf{z}_n) = p(\mathbf{z}_{n-1}, \mathbf{z}_n | \mathbf{X}) \quad (21)$$

$$= \frac{\alpha(\mathbf{z}_{n-1})p(\mathbf{z}_n | \mathbf{z}_{n-1})p(\mathbf{x}_n | \mathbf{z}_n)\beta(\mathbf{z}_n)}{p(\mathbf{X})} \quad (22)$$

- ▶ Taking the first derivative of $Q(\theta, \theta^{\text{old}})$ with respect to the parameters and setting it to zero gives us the update rules:

$$\pi_k = \frac{\gamma(z_{1k})}{\sum_{j=1}^K \gamma(z_{1j})} \quad (23)$$

$$A_{jk} = \frac{\sum_{n=2}^N \xi(z_{n-1,j}, z_{nk})}{\sum_{l=2}^K \sum_{n=2}^N \xi(z_{n-1,j}, z_{nl})} \quad (24)$$

- ▶ optimising the *emission* probabilities depends on the specific form of these, and is identical to a mixture model.

Example: Gaussian observations

If the observations have a Gaussian distribution, the maximum likelihood parameters are given by:

$$\boldsymbol{\mu}_k = \frac{\sum_{n=1}^N \gamma(z_{nk}) \mathbf{x}_n}{\sum_{n=1}^N \gamma(z_{nk})} \quad (25)$$

$$\boldsymbol{\Sigma}_k = \frac{\sum_{n=1}^N \gamma(z_{nk}) (\mathbf{x}_n - \boldsymbol{\mu}_k)(\mathbf{x}_n - \boldsymbol{\mu}_k)^\top}{\sum_{n=1}^N \gamma(z_{nk})} \quad (26)$$

(See Gaussian Mixture Models)

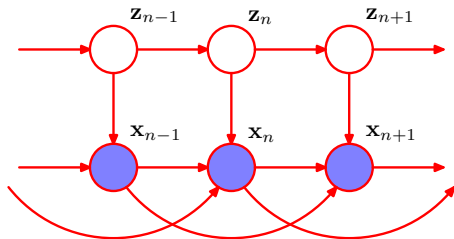
The basic HMM framework can be extended in many ways

- ▶ Autoregressive HMM

- ▶ Input-output HMM

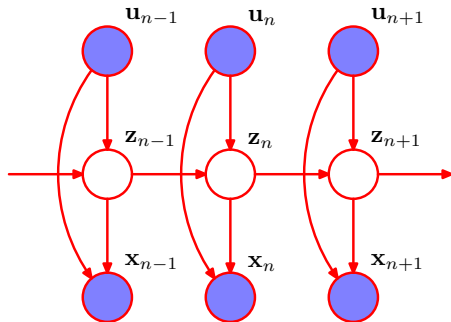
- ▶ Factorial HMM

- ▶ Explicit time model
(Hidden semi-Markov
model)



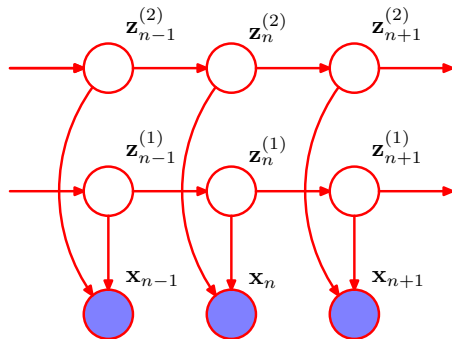
The basic HMM framework can be extended in many ways

- ▶ Autoregressive HMM
- ▶ Input-output HMM
- ▶ Factorial HMM
- ▶ Explicit time model
(Hidden semi-Markov model)



The basic HMM framework can be extended in many ways

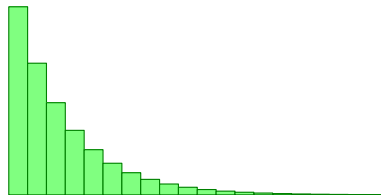
- ▶ Autoregressive HMM
- ▶ Input-output HMM
- ▶ Factorial HMM
- ▶ Explicit time model
(Hidden semi-Markov model)



The basic HMM framework can be extended in many ways

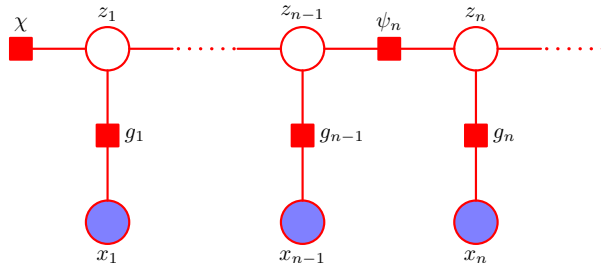
- ▶ Autoregressive HMM
- ▶ Input-output HMM
- ▶ Factorial HMM
- ▶ Explicit time model
(Hidden semi-Markov model)

$$p(z_n, \dots, z_{n+t} = z) = A_{zz}^t (1 - A_{zz})$$

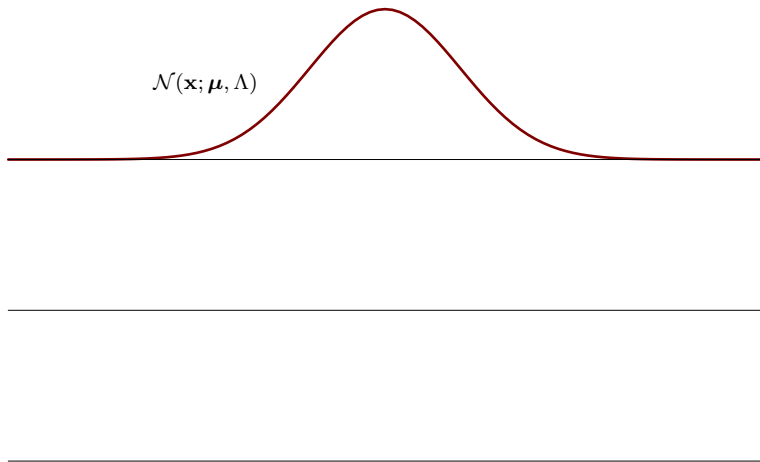


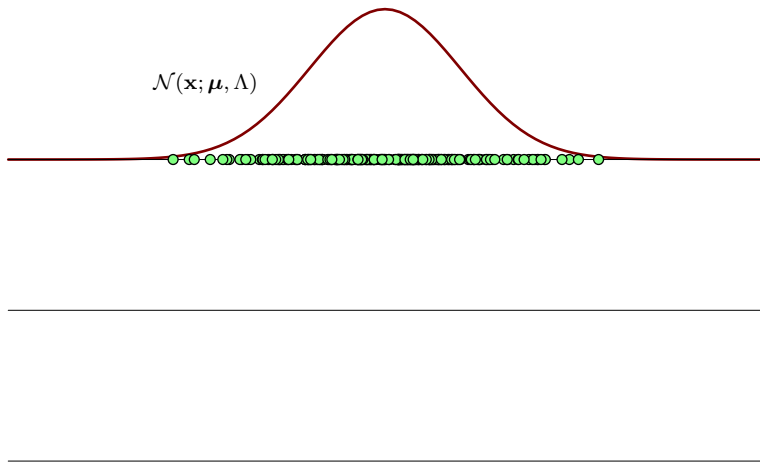
In the HMM, we assume that the latent variables are discrete.

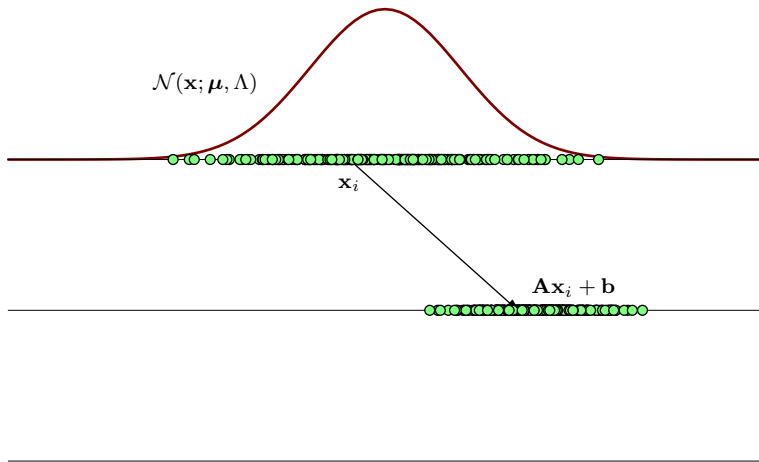
- ▶ The observations (emission probabilities) can have any form (Bernoulli, Gaussian, Mixture of Gaussians, *etc.*)
- ▶ This is not a limitation of the graphical model per se.

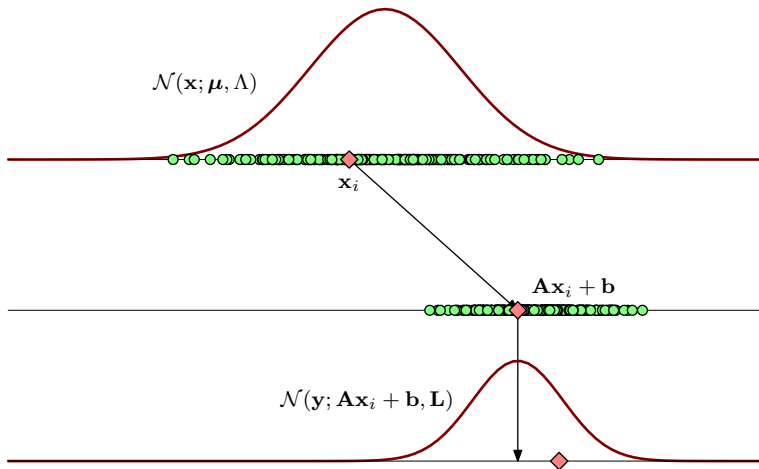


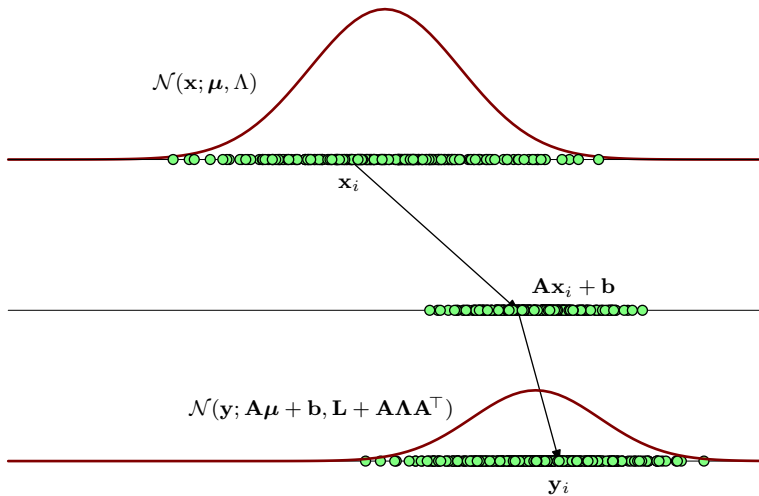
- ▶ One well-known model that allows for continuous latent variables is the Linear Dynamical System

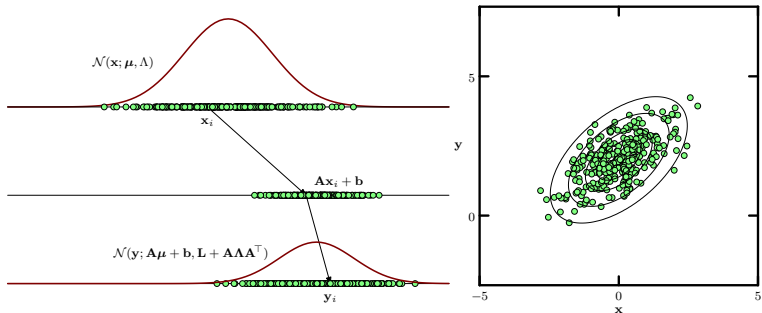


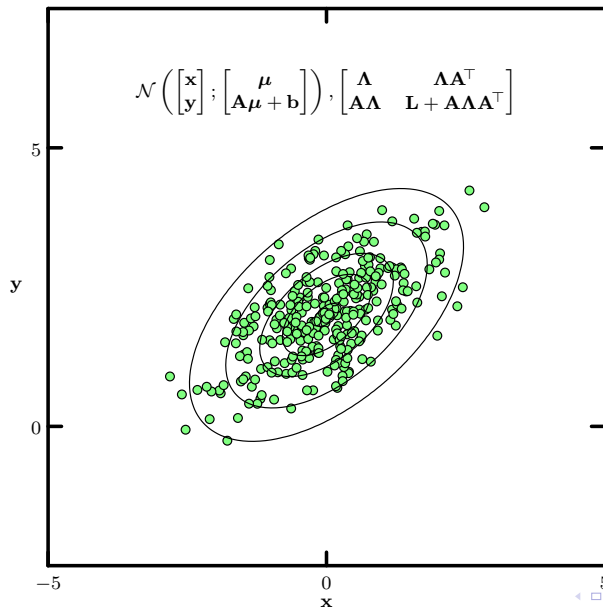


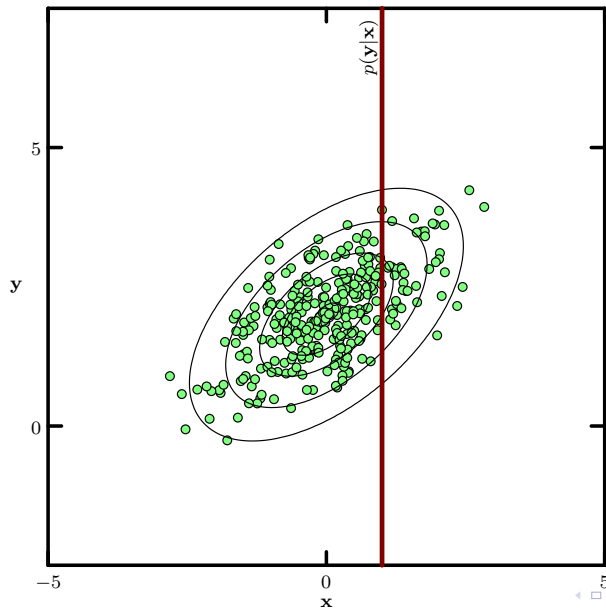


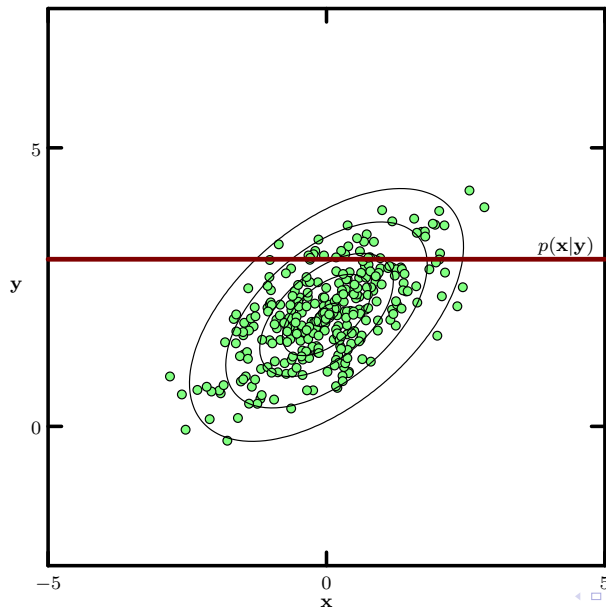












The LDS assumes Gaussian distributions for all variables, where

- ▶ The mean of the next latent variable is a linear function of the previous latent variable:

$$p(\mathbf{z}_n | \mathbf{z}_{n-1}) = \mathcal{N}(\mathbf{z}_n | \mathbf{A}\mathbf{z}_{n-1}, \mathbf{\Gamma}) \quad (27)$$

- ▶ The current observation has a Gaussian distribution centred around the current latent variable

$$p(\mathbf{x}_n | \mathbf{z}_n) = \mathcal{N}(\mathbf{x}_n | \mathbf{C}\mathbf{z}_n, \mathbf{\Sigma}) \quad (28)$$

- ▶ Assuming Gaussian distributions ensures that the complexity of the messages does not increase along the chain
- ▶ The prior $p(\mathbf{z}_0)$ can be a mixture of K Gaussians, however

The LDS assumes Gaussian distributions for all variables, where

- ▶ The mean of the next latent variable is a linear function of the previous latent variable:

$$p(\mathbf{z}_n | \mathbf{z}_{n-1}) = \mathcal{N}(\mathbf{z}_n | \mathbf{A}\mathbf{z}_{n-1} + \mathbf{B}\mathbf{u}, \mathbf{\Gamma}) \quad (27)$$

- ▶ The current observation has a Gaussian distribution centred around the current latent variable

$$p(\mathbf{x}_n | \mathbf{z}_n) = \mathcal{N}(\mathbf{x}_n | \mathbf{C}\mathbf{z}_n, \mathbf{\Sigma}) \quad (28)$$

- ▶ Assuming Gaussian distributions ensures that the complexity of the messages does not increase along the chain
- ▶ The prior $p(\mathbf{z}_0)$ can be a mixture of K Gaussians, however

The LDS assumes Gaussian distributions for all variables, where

- ▶ The mean of the next latent variable is a linear function of the previous latent variable:

$$p(\mathbf{z}_n|\mathbf{z}_{n-1}) = \mathcal{N}(\mathbf{z}_n|\mathbf{A}\mathbf{z}_{n-1} + \mathbf{B}\mathbf{u}, \mathbf{\Gamma}) \quad (27)$$

- ▶ The current observation has a Gaussian distribution centred around the current latent variable

$$p(\mathbf{x}_n|\mathbf{z}_n) = \mathcal{N}(\mathbf{x}_n|\mathbf{C}\mathbf{z}_n, \mathbf{\Sigma}) \quad (28)$$

- ▶ Assuming Gaussian distributions ensures that the complexity of the messages does not increase along the chain
- ▶ The prior $p(\mathbf{z}_0)$ can be a mixture of K Gaussians, however

The LDS assumes Gaussian distributions for all variables, where

- ▶ The mean of the next latent variable is a linear function of the previous latent variable:

$$p(\mathbf{z}_n|\mathbf{z}_{n-1}) = \mathcal{N}(\mathbf{z}_n|\mathbf{A}\mathbf{z}_{n-1} + \mathbf{B}\mathbf{u}, \mathbf{\Gamma}) \quad (27)$$

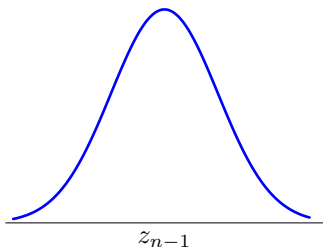
- ▶ The current observation has a Gaussian distribution centred around the current latent variable

$$p(\mathbf{x}_n|\mathbf{z}_n) = \mathcal{N}(\mathbf{x}_n|\mathbf{C}\mathbf{z}_n, \mathbf{\Sigma}) \quad (28)$$

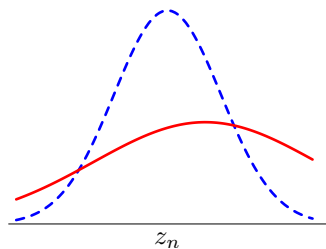
- ▶ Assuming Gaussian distributions ensures that the complexity of the messages does not increase along the chain
- ▶ The prior $p(\mathbf{z}_0)$ can be a mixture of K Gaussians, however

Again, this is a model that was developed long before graphical models

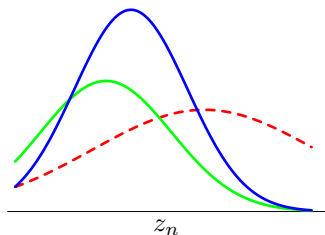
- ▶ If only the forward pass is done, it is called the *Kalman Filter*
- ▶ It was developed for optimal tracking of rockets and satellites
- ▶ If the backwards pass is performed as well, it is called the *Kalman Smoother*. The corresponding equations are called the *Rauch-Tung-Striebel* (RTS) equations.



- ▶ Distribution over \mathbf{z}_{n-1}
- ▶ Prediction over \mathbf{z}_n
- ▶ Probability of the observation \mathbf{x}_n given \mathbf{z}_n , and updated distribution over \mathbf{z}_n



- ▶ Distribution over \mathbf{z}_{n-1}
- ▶ Prediction over \mathbf{z}_n
- ▶ Probability of the observation \mathbf{x}_n given \mathbf{z}_n , and updated distribution over \mathbf{z}_n



- ▶ Distribution over \mathbf{z}_{n-1}
- ▶ Prediction over \mathbf{z}_n
- ▶ Probability of the observation \mathbf{x}_n given \mathbf{z}_n , and updated distribution over \mathbf{z}_n

The Kalman filter is often used for tracking

- ▶ The matrices **A** and **C** are then known *a priori* — e.g.

$$\mathbf{z} = \begin{bmatrix} x \\ y \\ v_x \\ v_y \end{bmatrix} \qquad \mathbf{A} = \begin{bmatrix} 1 & 0 & 1 & 0 \\ 0 & 1 & 0 & 1 \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 \end{bmatrix} \qquad (29)$$

$$\mathbf{x} = \begin{bmatrix} x \\ y \end{bmatrix} \qquad \mathbf{C} = \begin{bmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \end{bmatrix} \qquad (30)$$

- ▶ Covariance matrices **Γ** and **Σ** must be estimated

... Sometimes process must be learnt as well

As always, we can use EM to learn the model parameters

- ▶ The latent variables are continuous, posterior PDF is a Gaussian
- ▶ The expectation of the complete log-likelihood can be optimised in closed form
 - ▶ With respect to the distribution parameters $(\boldsymbol{\mu}, \boldsymbol{\Sigma}, \boldsymbol{\Gamma})$
 - ▶ With respect to the deterministic model parameters (\mathbf{A}, \mathbf{C})

Today, we have seen:

- ▶ Learning when we can do inference (Bishop, p. 439–441)
- ▶ Examples of the EM algorithm (Bishop, p. 423–439)
- ▶ A formal analysis of EM (Bishop, p. 450–453)
- ▶ Compared k-means with mixtures of Gaussians
- ▶ Markov Models (Bishop, p. 605–610)
- ▶ HMM (Bishop, p. 610–630)
- ▶ Introduction to the LDS (Bishop, p. 635–637)

Lab:

- ▶ The EM algorithm for GMM