

STA 414/2104

Statistical Methods for Machine Learning and Data Mining

Radford M. Neal, University of Toronto, 2013

Week 7

# Clustering and Mixture Models

# Unsupervised Learning, Clustering, and Mixture Distributions

Recall that *unsupervised learning* does not focus on predicting anything in particular, but rather tries to find “interesting” aspects of the data.

One possible informal objective is to find *clusters* of similar items. One possible formal objective is modeling the *probability distribution* of all observed variables.

It can be useful to model a probability distribution in which variables are dependent using *latent* (also called *hidden*) variables to express some or all of the dependencies.

When there is one discrete latent variable, the model will express the distribution as a *mixture* of distributions. The latent variable can also be seen as identifying the cluster an item belongs to.

**Example:** We have data on symptoms of patients (body temperature, blood pressure, etc.). We could cluster the patients, hoping the clusters will correspond to “diseases”. We could model the distribution of symptoms using a discrete latent variable, hoping it will represent the disease a patient has. The distribution of symptoms will be a mixture of distributions for each disease.

# K-Means Clustering

Suppose we have data  $x_1, \dots, x_n$ , with each  $x_i$  being a vector of  $p$  variables.

We aim to divide the data into  $K$  clusters of similar items, measuring similarity by Euclidean distance (perhaps rescaling variables to all have standard deviation one).

The *K-means algorithm* iteratively finds  $K$  centres for the clusters, and assigns each item to the cluster whose centre it is nearest to:

- 1) Initialize  $\mu_1, \dots, \mu_K$  somehow (eg, set them to randomly selected data items).
- 2) Repeat the following:
  - a) For  $i = 1, \dots, n$ , assign data item  $i$  to the cluster,  $k$ , for which  $\|x_i - \mu_k\|$  is smallest. (Prefer the current assignment if it is tied for the best.)
  - b) For  $k = 1, \dots, K$ , set  $\mu_k$  to the mean of data items assigned to cluster  $k$ .until there is no change in the cluster assignments from the previous iteration.

It's easy to see that this process converges, since each step reduces the value of

$$J = \sum_{i=1}^n \|x_i - \mu_{c_i}\|^2$$

where  $c_i$  is the cluster assigned to data item  $i$ . The algorithm may not find the global minimum of  $J$ , however.

## From Clustering to Mixture Models

K-means clustering assigns a definite cluster to each data item. But often this is unrealistic — eg, some patients have combinations of the symptoms we have measured that are consistent with more than one disease.

The clusters found by the K-means algorithm are described only by the mean of the data items in the cluster. We might like a more complete description of what items in a cluster are like.

Both of these issues can be addressed by modeling the data as coming from a *mixture* distribution, with mixture components corresponding to clusters.

# Mixture Distributions

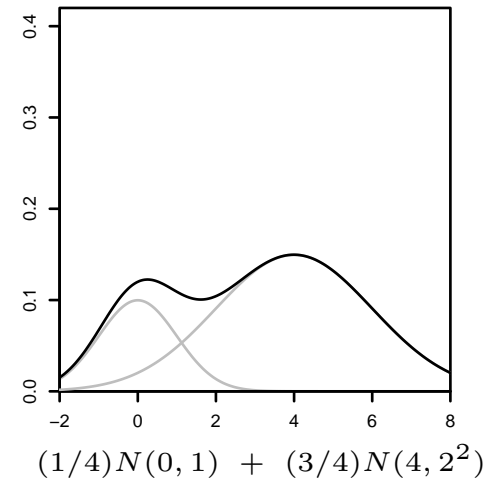
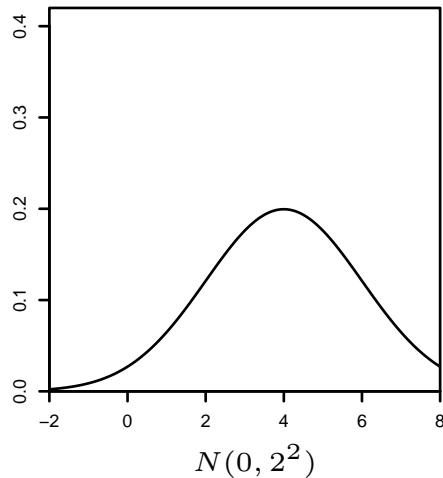
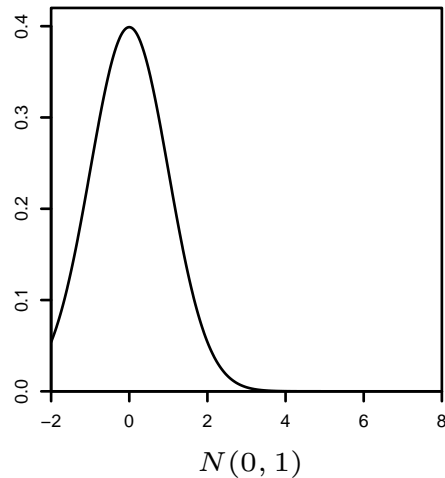
$K$  distributions with probability/density functions  $P_1(x), \dots, P_K(x)$  can be mixed in proportions  $\pi_1, \dots, \pi_K$  to give a mixture distribution with probability/density function

$$P(x) = \sum_{k=1}^K \pi_k P_k(x)$$

For example, when  $x$  is one dimensional, we can mix  $N(0, 1^2)$  and  $N(4, 2^2)$  with proportions  $1/4$  and  $3/4$ , giving the density function

$$P(x) = \frac{1}{4} \frac{1}{\sqrt{2\pi}} \exp(-(1/2)x^2) + \frac{3}{4} \frac{1}{2\sqrt{2\pi}} \exp(-(1/2)(x-4)^2/2^2)$$

as pictured below:



# Gaussian Mixture Models

When data items are real vectors, it may be reasonable to model the data as a mixture of Gaussian distributions, using the data to estimate the mixing proportion, the mean vector, and covariance matrix for each component.

The density function for this model, with  $K$  components, can be written as

$$P(x|\pi, \mu, \Sigma) = \sum_{k=1}^K \pi_k N(x|\mu_k, \Sigma_k)$$

We might allow  $\Sigma_k$  to be any valid covariance matrix, or we might constrain  $\Sigma_k$  to be a diagonal matrix. If the  $\Sigma_k$  are diagonal, *all* the dependence between the variables in  $x$  is a consequence of the distribution being a mixture.

A natural idea is to estimate  $\pi_k$ ,  $\mu_k$ , and  $\Sigma_k$  for  $k = 1, \dots, K$  by maximizing the likelihood. Assuming that data items are independent, the log likelihood is

$$L(\pi, \mu, \Sigma) = \sum_{i=1}^n \log P(x_i|\pi, \mu, \Sigma)$$

where  $x_i$  is the data vector for item  $i$ .

# Issues with Maximum Likelihood for Gaussian Mixture Models

**Non-identifiability:** The global maximum of the likelihood for a mixture model is not unique, since permuting the labels of the mixture components will produce a different set of parameter values that fits the data just as well.

**Other local maxima:** Even aside from re-labellings, there is often more than one local maximum of the likelihood. Finding one of the global maxima (or at least a good local maximum) may require searching for the maximum from many different starting points.

**We don't want the global maximum anyway:** When  $K > 1$ , the actual global maximum will be at a point with infinite likelihood, in which for some component,  $k$ , and some data item,  $i$ , we have  $0 < \pi_k < 1$ ,  $\mu_k = x_i$ , and  $\Sigma_k = 0$ . This gives an infinite spike of probability density at one point, while other points have non-zero probability density from other components.

Because of this problem, we need to try as many starting points as needed to find a good local maximum that *isn't* one where some  $\Sigma_k \rightarrow 0$ .