

Clustering: K -means and Gaussian Mixture Models

Brian Azizi, Jia Guang Choo

Cavendish Laboratory, Department of Physics, J J Thomson Avenue, Cambridge. CB3 0HE

Clustering

We have an unlabelled training set of N data points, $\{\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(N)}\}$. Each data has D features, i.e. $\mathbf{x}^{(i)} \in \mathbb{R}^D$ for each i . We compactly represent the training data as an $N \times D$ matrix \mathbf{X} .

Our goal is to identify clusters in the data. Data points within the same cluster should be similar to each other and unlike data points in different clusters.

We will discuss two clustering models: K -means and GMM. Both models are usually trained using the *EM-algorithm*.

1. K -means

Defining the Model

Each of the K clusters can be characterized by its centre $\boldsymbol{\mu}_k \in \mathbb{R}^D, k = 1, \dots, K$. Each point $\mathbf{x}^{(i)}$ has a corresponding binary vector $\mathbf{r}^{(i)} \in \mathbb{Z}_2^D$, where $r_k^{(i)} = 1$ if $\mathbf{x}^{(i)}$ belongs to cluster k , 0 otherwise. In general, we aim to minimize the cost function

$$C = \sum_{i=1}^N \sum_{k=1}^K r_k^{(i)} V(\mathbf{x}^{(i)}, \boldsymbol{\mu}_k) \quad (1)$$

where $V(\mathbf{x}^{(i)}, \boldsymbol{\mu}_k)$ measures the ‘distance’ between $\mathbf{x}^{(i)}$ and $\boldsymbol{\mu}_k$. For the purpose of this presentation, we let $V(\mathbf{x}^{(i)}, \boldsymbol{\mu}_k) = \|\mathbf{x}^{(i)} - \boldsymbol{\mu}_k\|^2$.

Minimizing the cost function

A local (and hopefully global) minimum for C is obtained by repeatedly minimizing C with respect to $\mathbf{r}^{(i)}$ while keeping $\boldsymbol{\mu}_k$ fixed, then doing the same with respect to $\boldsymbol{\mu}_k$ while keeping $\mathbf{r}^{(i)}$ fixed, until convergence.

For an initial set of $\boldsymbol{\mu}_k$, we minimize C with respect to $\mathbf{r}^{(i)}$ by setting

$$r_k^{(i)} = \begin{cases} 1 & \text{if } j = \arg \min_j \|\mathbf{x}^{(i)} - \boldsymbol{\mu}_j\|^2, \\ 0 & \text{otherwise.} \end{cases} \quad (2)$$

We then minimize C with respect to $\boldsymbol{\mu}_k$ by differentiating C with respect to $\boldsymbol{\mu}_k$ to get

$$\frac{\partial C}{\partial \boldsymbol{\mu}_k} = 2 \sum_{i=1}^N r_k^{(i)} (\mathbf{x}^{(i)} - \boldsymbol{\mu}_k) = 0 \Rightarrow \boldsymbol{\mu}_k = \frac{\sum_{i=1}^N r_k^{(i)} \mathbf{x}^{(i)}}{\sum_{i=1}^N r_k^{(i)}}. \quad (3)$$

Applications in image compression

Suppose we have an image with N pixels, each comprising of 3 numbers representing the intensity of red, blue and green, and each of these numbers is stored with 8 bits (i.e. values range from 0 to 255). Hence, the total number of bits required to store the image directly would be $24N$.

Now suppose the K -means algorithm is applied, where the features are the values of red, blue and green. If there are K clusters, then only the K vectors corresponding to $\boldsymbol{\mu}_k$ are stored. In addition, the cluster that each pixel belongs to can also be stored in $\log_2 K$ bits. Hence the total number of bits required is $24K + N \log_2 K$, which is much smaller than $24N$.

2. Gaussian Mixture Model

Defining the Model

In latent variable models, we model the probability density function (pdf) of the data as

$$p(\mathbf{x}) = \sum_{\mathbf{z}} p(\mathbf{x}|\mathbf{z})p(\mathbf{z}) \quad (4)$$

The random variable \mathbf{x} represents our data and the random variable \mathbf{z} is a so called *latent* variable, i.e. it is not observed.

In mixture models, \mathbf{x} is generated from one of K possible base distributions and \mathbf{z} tells us from which. Here, K is analogous to number of clusters K in K -means,. So, to get the required behaviour, \mathbf{z} has a categorical nature. We generally use the so-called 1-of- K representation of \mathbf{z} , where

$$\mathbf{z} \in \{0, 1\}^K, \quad \sum_{k=1}^K z_k = 1 \quad (5)$$

For example, if $K = 3$ and for a particular point $\mathbf{x}^{(i)}$ we have $\mathbf{z}^{(i)} = (0, 1, 0)$, then this means that $\mathbf{x}^{(i)}$ came from the 2nd base distribution (think: 2nd cluster).

Of course, we don't know \mathbf{z} , since it is a latent variable. So we define a distribution for it:

$$p(z_k = 1) = \pi_k \quad (6)$$

or equivalently

$$p(\mathbf{z}) = \prod_{k=1}^K \pi_k^{z_k} \quad (7)$$

Equation (4) then gives:

$$\begin{aligned} p(\mathbf{x}) &= \sum_{\mathbf{z}} p(\mathbf{x}|\mathbf{z})p(\mathbf{z}) \\ &= \sum_{\mathbf{z}} p(\mathbf{x}|\mathbf{z}) \prod_{k=1}^K \pi_k^{z_k} \\ &= \sum_{k=1}^K \pi_k p(\mathbf{x}|z_k = 1) \end{aligned} \quad (8)$$

Now we have general mixture models of K base distribution. In *Gaussian Mixture Models*, each base distribution is a multivariate Gaussian. The k th base distribution has parameters $(\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)$ and has pdf

$$\mathcal{N}(\mathbf{x} | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k) = \frac{1}{(2\pi)^{D/2}} \frac{1}{|\boldsymbol{\Sigma}_k|^{1/2}} \exp \left\{ -\frac{1}{2} (\mathbf{x} - \boldsymbol{\mu}_k)^T \boldsymbol{\Sigma}_k^{-1} (\mathbf{x} - \boldsymbol{\mu}_k) \right\} \quad (9)$$

The final Gaussian mixture model is then

$$p(\mathbf{x}) = \sum_{k=1}^K \pi_k \mathcal{N}(\mathbf{x} | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k) \quad (10)$$

Training the model

To train the model, i.e. find the parameters μ_k and Σ_k for $k = 1, \dots, K$, we need a set of training data $\mathbf{X} \in \mathbb{R}^{N \times D}$. We assume that the samples were idenpendently generated from our model, which allows us to write the likelihood function as

$$p(\mathbf{X} | \pi, \mu, \Sigma) = \prod_{i=1}^N \sum_{k=1}^K \pi_k \mathcal{N}(\mathbf{x}^{(i)} | \mu_k, \Sigma_k). \quad (11)$$

We maximize the (log) likelihood to get the parameters. The solution can be expressed as

$$\begin{aligned} \mu_k &= \frac{1}{N_k} \sum_{i=1}^N \gamma(z_k^{(i)}) \mathbf{x}^{(i)}, \\ \Sigma_k &= \frac{1}{N_k} \sum_{i=1}^N \gamma(z_k^{(i)}) (\mathbf{x}^{(i)} - \mu_k)(\mathbf{x}^{(i)} - \mu_k)^T, \\ \pi_k &= \frac{N_k}{N}, \\ N_k &= \sum_{i=1}^N \gamma(z_k^{(i)}), \\ \gamma(z_k^{(i)}) &= \frac{\pi_k \mathcal{N}(\mathbf{x}^{(i)} | \mu_k, \Sigma_k)}{\sum_{j=1}^K \pi_j \mathcal{N}(\mathbf{x}^{(i)} | \mu_j, \Sigma_j)}. \end{aligned} \quad (12)$$

The quantity $\gamma(z_k^{(i)})$ is referred to as the *responsibility* that cluster k takes in explaining sample $\mathbf{x}^{(i)}$, while N_k is the effective number of data points in cluster k .

Note that this is not a closed-form solution since the quantities depend on each other. However, we can try to find the solution iteratively as follows: First, randomly initialize the quantities μ_k , Σ_k and π_k . Then calculate $\gamma(z_k^{(i)})$ and N_k . Then calculate μ_k , Σ_k and π_k and so on until convergence. This iterative scheme is known as the *Expectation-Maximization algorithm* for the Gaussian mixture model.

For more information and details on the derivation, see [1] [2].

References

- [1] Georgios Pilikos, Unsupervised Learning, Lecture Notes in Machine Learning, Chapter 2, 2015
- [2] Christopher Bishop, Pattern Recognition and Machine Learning, Chapter 9, 2015, Springer