# Clustering, Mixture Models, and the EM Algorithm

Rebecca C. Steorts

# Agenda

- ▶ Clustering
- ▶ Two Component Mixture Model
- ▶ Latent Variable
- ▶ EM Algorithm

# What will you learn in this lecture

- ▶ Clustering
- ▶ Importance of mixture models
- ▶ Simple illustrations of mixture models
- ▶ We will learn about the two component mixture model (and extension)
- ▶ Will learn about latent variables
- ▶ Two component mixture model (and extension)
- ▶ EM algorithm

# Clustering

Clustering is an **unsupervised method** that divides up data into groups (clusters), so that points in any one group are more "similar'' to each other than to points outside the group

# Clustering methods (that we have covered)

- ▶ Fuzzy clustering (such as deterministic entity resolution or blocking)
- ▶ Overlapping clustering (such as locality sensitive hashing)
- ▶ Fellegi Sunter method (technically a mixture model, stay tuned)

# Clustering methods (that we have not covered)

- Mixture models, such as Gaussian (GMMs)
- K-means
- K-mediods
- Hierarchical clustering
- Among many others

# Application areas

- ▶ Clustering temperatures to identify weather patterns, grouping individuals based on height and weight similarities.
- ▶ Clustering customers based on satisfaction levels (ordinal) or grouping individuals based on gender (categorical).
- ▶ Clustering stock price movements to identify similar trends, grouping temperature readings based on seasonal patterns.
- ▶ Clustering GPS coordinates to identify spatial patterns, grouping locations on a map based on features.

# Application areas

▶ Clustering users in a social network or data based based on their connections (edge structure), grouping academic papers based on citation patterns.
▶ Clustering articles based on topics, grouping emails by subject matter.
▶ Clustering photographs by content, grouping medical images based on visual features.
▶ Clustering music tracks by genre, grouping speech recordings by language.
▶ Clustering surveillance footage to identify similar activities, grouping video clips by visual similarities.

# Importance

- ▶ First proposed by Karl Pearson (1984) and analyzed on crab data.
- ▶ Applications: "agriculture, astronomy, bioinformatics, biology, economics, engineering, genetics, imaging, marketing, medicine, neuroscience, psychiatry, and psychology, among many other fields in the biological, physical, and social sciences". McLachlan et. al (2019).
- ▶ One of the methods in machine learning is **topic modeling**, which identifies "topics" in collections of documents/webpages.
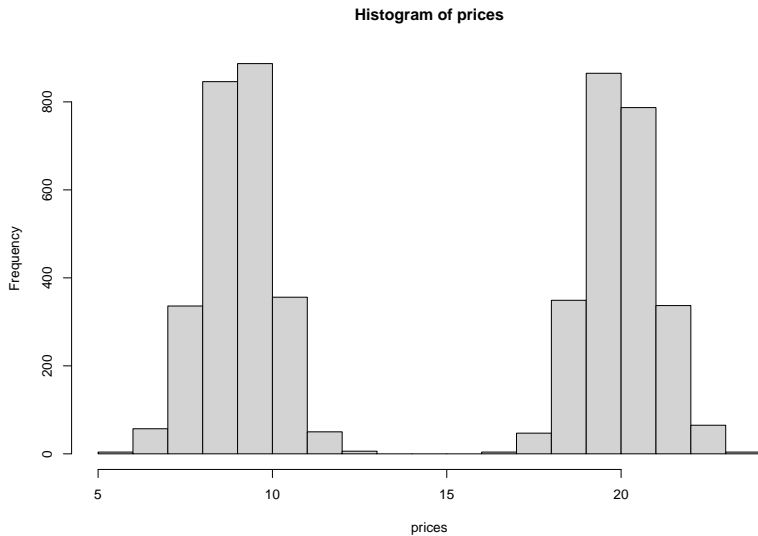- ▶ Topic modeling relies on mixtures models.

# Motivation

- ▶ Suppose we want to simulate the price of a randomly chosen book.

- ▶ Paperbacks are often cheaper than hardbacks, so let's model them separately.

- ▶ Model the price of a book as a mixture model.

- ▶ There will be two components (or clusters) in our model – one for paperbacks and one for hardbacks.

# Model

- Paperback distribution: $N(9, 1)$
- Hardback distribution: $N(20, 2)$
- Assume that there's a there is a 50% chance of choosing a paperback and 50% of choosing hardback.

# Motivation



**Histogram of prices**
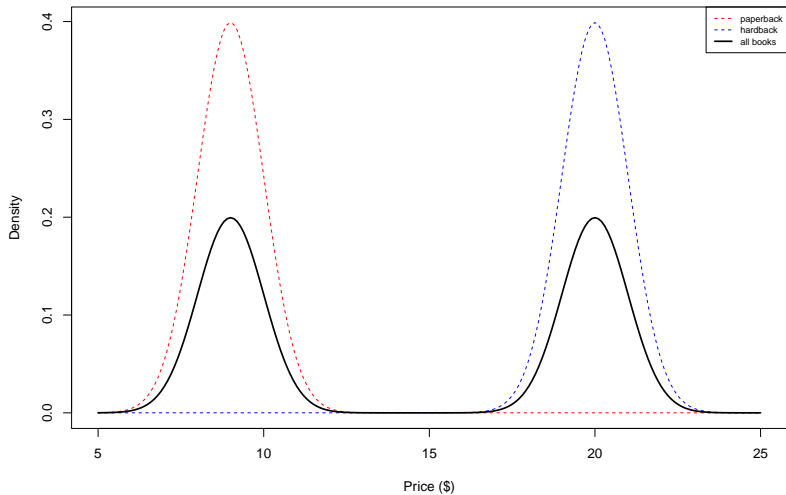
# Motivation

► Are the prices of books unimodal or bimodal?
► Suppose you would want to predict the price of a book. Would its distribution be Normal or something else based on the the histogram.
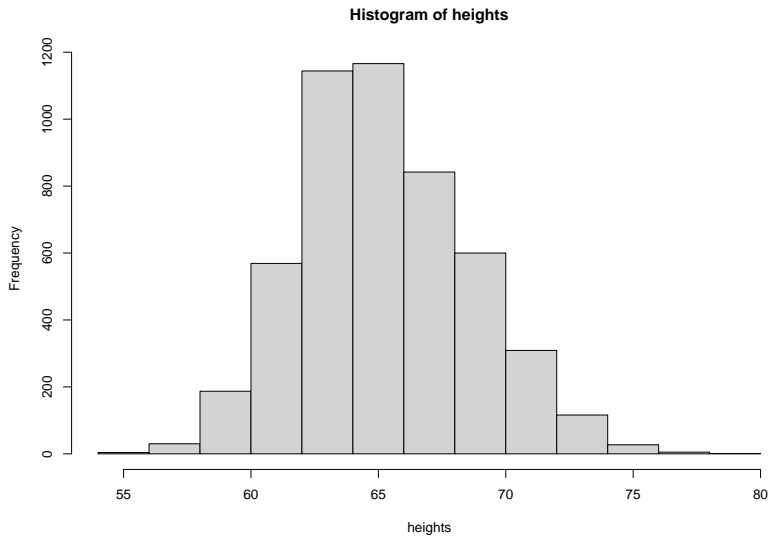
# Motivation

# Motivation

Now assume our data are the heights of students at University X.

Assume the height of a randomly chosen male is normally distributed with a mean equal to 5′9 and a standard deviation of 2.5 inches.

Assume the height of a randomly chosen female is $N(5'4, 2.5)$.

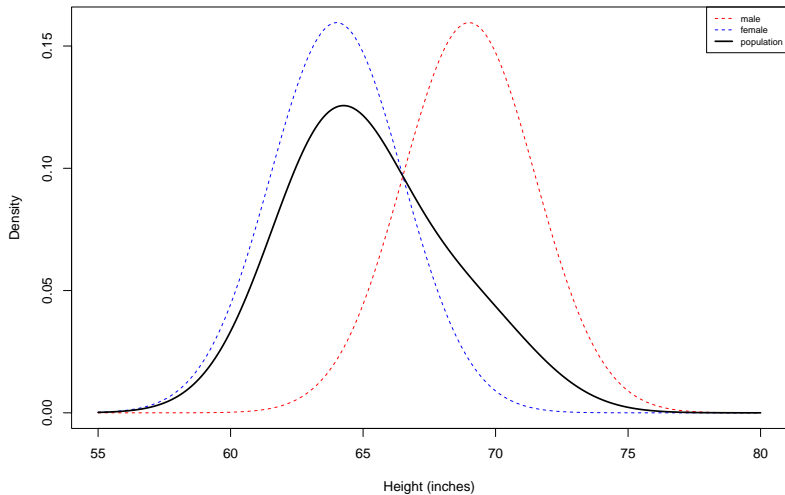Assume that 75% of the population is female and 25% is male.

# Motivation



Histogram of heights

# Motivation

The histogram is now unimodal.

Are heights normally distributed (assuming this model)? Let's invesitgate!

# Motivation

# Motivation

The Gaussian mixture model is unimodal because there is so much overlap between the two densities.

In this example, observe that the population density is not symmetric, and therefore not normally distributed.

# Goal

The goal of this module is to introduce **mixture models**, which are commonly used in applications in classical and modern machine learning.

We will do this using a **latent variable**.

# Background

A **latent variable** is the true version of the state of a random variable that is unknown and not directly observed.[1]

---

[1]We will not delve into the properties of latent variables in this course.

# Mixture models can be viewed as probabilistic clustering

▶ Mixture models put similar data points into "clusters".

▶ This is appealing as we can potentially compare different probabilistic clustering methods by how well they predict (under cross-validation). We will not explore this in this particular lecture.

▶ This contrasts other methods such as k-means and hierarchical clustering as they produce clusters (and not predictions), so it's difficult to test if they are correct/incorrect.[2]

---

[2]Explore looking at these on your own and see if you can determine their limitations practically, compared to other machine learning models.

# Two-component mixture model

Assume that both mixture components have the same precision, $\lambda = 1/\sigma^2$, which is fixed and known.

Let $\pi$ be the mixture proportion for the first component.

Then the two-component Normal mixture model is:

$$X_1, \ldots, X_n \mid \mu, \pi \sim F(\mu, \pi) \tag{1}$$

where $F(\mu, \pi)$ is the distribution with p.d.f.

$$f(x|\mu, \pi) = (1 - \pi)\mathcal{N}(x \mid \mu_0, \lambda^{-1}) + \pi\mathcal{N}(x \mid \mu_1, \lambda^{-1}).$$

# Likelihood

The likelihood is

$$p(x_{1:n}|\mu, \pi) = \prod_{i=1}^{n} f(x_i|\mu, \pi)$$
$$= \prod_{i=1}^{n} \left[ (1 - \pi)\mathcal{N}(x_i \mid \mu_0, \lambda^{-1}) + \pi\mathcal{N}(x_i \mid \mu_1, \lambda^{-1}) \right].$$

# Likelihood

What do you notice about the likelihood function?

$$p(x_{1:n}|\mu, \pi) = \prod_{i=1}^{n} f(x_i|\mu, \pi)$$
$$= \prod_{i=1}^{n} \left[ (1-\pi)\mathcal{N}(x_i \mid \mu_0, \lambda^{-1}) + \pi\mathcal{N}(x_i \mid \mu_1, \lambda^{-1}) \right].$$

## Likelihood

The **likelihood** is very complicated function of $\mu$ and $\pi$.

This makes working with it directly to find the MLE (or other estimates) difficult.

Thus, we will rewrite the likelihood using **latent variables**.

# Latent allocation variables to the rescue!

Define an equivalent model that includes latent "allocation''
variables $Z_1, \ldots, Z_n$.

These indicate which mixture component each data point comes
from–that is, $Z_i$ indicates whether subject $i$ is from component 1 or
2.

$$X_i \mid \mu, Z \sim \mathcal{N}(\mu_{Z_i}, \lambda^{-1}) \text{ independently for } i = 1, \ldots, n. \tag{2}$$

$$Z_1, \ldots, Z_n \mid \mu, \pi \overset{iid}{\sim} \text{Bernoulli}(\pi) \tag{3}$$

How can we check that the latent allocation model is equivalent to
our original model? (Exercise to complete outside of class).

# Extension to k-components

Assume we observe $X_1, \ldots, X_n$ and that each $X_i$ is sampled from one of $K$ **mixture components**.

Associated with each random variable $X_i$ is a latent variable (or label) $Z_i \in \{1, \ldots, K\}$ which indicates which component $X_i$ came from.

# Notation

Let $\pi_k$ be called **mixture proportions** or **mixture weights**, which represent the probability that $X_i$ belongs to the $k$-th mixture component.

The mixture proportions are non-negative and they sum to one, $\sum_{k=1}^{K} \pi_k = 1$.

We call $P(X_i \mid Z_i = k)$ the **mixture component**, and it represents the distribution of $X_i$ assuming it came from component $k$.

# Law of Total Probability

From the law of total probability, it follows that

$$P(X_i = x) = \sum_{k=1}^{K} P(X_i = x | Z_i = k) \underbrace{P(Z_i = k)}_{\pi_k} \qquad (4)$$

$$= \sum_{k=1}^{K} P(X_i = x | Z_i = k) \pi_k \qquad (5)$$

$$= \sum_{k=1}^{K} \pi_k P(X_i = x | Z_i = k) \qquad (6)$$

# Mixture Models

For discrete random variables these mixture components can be any probability mass function $p(\cdot \mid Z_k)$.

For continuous random variables they can be any probability density function $f(\cdot \mid Z_k)$.

The corresponding pmf and pdf for the mixture model are:

$$p(x) = \sum_{k=1}^{K} \pi_k p(x \mid Z_k)$$

$$f_x(x) = \sum_{k=1}^{K} \pi_k f_{x \mid Z_k}(x \mid Z_k)$$

## Likelihood

The likelihood of observing independent samples $X_1, \ldots, X_n$ with mixture proportion vector $\pi = (\pi_1, \pi_2, \ldots, \pi_K)$ is therefore:

$$L(X_1, \ldots, X_n \mid \pi) = \prod_{i=1}^{n} P(X_i|\pi) = \prod_{i=1}^{n} \sum_{k=1}^{K} P(X_i|Z_i = k)\pi_k$$

# Estimation

Now assume we are in the Gaussian mixture model setting where the $k$-th component is $N(\mu_k, \sigma^2)$ and the mixture proportions are $\pi_k$.

How can we estimate $\{\mu_k, \sigma^2, \pi_k\}$ from the observed data $X_1, \ldots, X_n$?

Solution: EM Algorithm.

Why? The MLE is not possible to find in closed form. Think about why this is the case.

# Conditional and marginal distributions

Recall that the conditional distribution $X_i | Z_i = k \sim N(\mu_k, \sigma_k^2)$, where $\pi_k = P(Z_i = k)$.

The marginal distribution of $X_i$ is:

$$P(X_i = x) = \sum_{k=1}^{K} P(Z_i = k) P(X_i = x | Z_i = k) \qquad (7)$$

$$= \sum_{k=1}^{K} \pi_k N(x \mid \mu_k, \sigma_k^2) \qquad (8)$$

Note: $\sigma_k^2 = \sigma^2$ moving forward.

# Joint distribution

The joint probability of observations $X_1, \ldots, X_n$ is

$$P(X_1 = x_1, \ldots, X_n = x_n) = \prod_{i=1}^{n} \sum_{k=1}^{K} \pi_k N(x_i \mid \mu_k, \sigma_k^2)$$

## Exercise

Show that

$$\log P(X_1, \ldots, X_n \mid \mu_1, \ldots, \mu_K) \tag{9}$$

$$= \log \prod_{i=1}^{n} P(x_i \mid \mu_1, \ldots, \mu_K) \tag{10}$$

$$= \sum_{i=1}^{n} \log\Big[\sum_{k=1}^{K} P(x_i \mid \pi_k, \mu_1, \ldots, \mu_K)\pi_k\Big] \tag{11}$$

# Background

Recall that

$$\frac{\partial \log f(x)}{\partial dx} = \frac{1}{f(x)}\frac{\partial f(x)}{\partial dx}.$$

## Exercise

Show that

$$\frac{\partial \log P(X_1, \ldots, X_n \mid \mu_1, \ldots, \mu_K)}{\partial \mu_k} \tag{12}$$

$$= \sum_{i=1}^{n} P(\pi_k \mid x_i, \mu_1, \ldots, \mu_K) \frac{(x_i - \mu_k)}{\sigma} \tag{13}$$

This implies that

$$\mu_k = \frac{\sum_{i=1}^{n} P(\pi_k \mid x_i, \mu_1, \ldots, \mu_K) x_i}{\sum_{i=1}^{n} P(\pi_k \mid x_i, \mu_1, \ldots, \mu_K)},$$

which is a non-linear equation of the $\mu_k$'s.

# Intuition of EM

$$\mu_k = \frac{\sum_{i=1}^n P(\pi_k \mid x_i, \mu_1, \ldots, \mu_K) x_i}{\sum_{i=1}^n P(\pi_k \mid x_i, \mu_1, \ldots, \mu_K)},$$

▶ E-step: If for each $x_i$ we knew that for each $\pi_k$ the prob. that $\mu_k$ was in component $\pi_k$ is $P(\pi_k \mid x_i, \mu_1, \ldots, \mu_K)$. Then we could compute $\mu_k$.

▶ M-step: If we knew each $\mu_k$, then we could compute $P(\pi_k \mid x_i, \mu_1, \ldots, \mu_K)$ for each $\mu_k$ and $x_i$

# EM Algorithm

Initalize all the unknown parameters. On iteration $t$, let the estimates be $\lambda^{(t)} = \{\mu_1^{(t)}, \ldots, \mu_k^{(t)}\}$

1. E-Step:

$$P(\pi_k \mid x_i, \lambda^{(t)}) = \frac{P(\pi_k \mid x_i, \lambda^{(t)})x_i}{P(\pi_k \mid x_i, \lambda^{(t)})} \tag{14}$$

2. M-Step:

$$\mu_k^{(t+1)} = \frac{\sum_{i=1}^n P(\pi_k \mid x_i, \lambda^{(t)})x_i}{\sum_{i=1}^n P(\pi_k \mid x_i, \lambda^{(t)})} \tag{15}$$

# Exercise

Assume our mixture components are fully specified Gaussian distributions with $K = 2$ components:
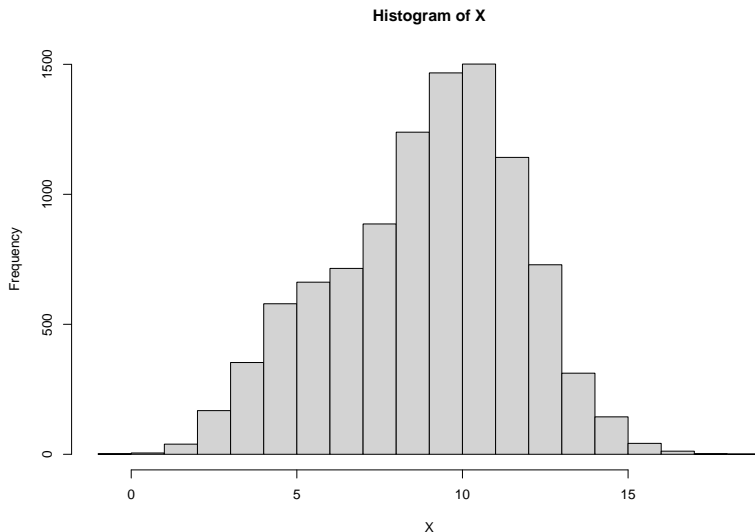
$$X_i \mid Z_i = 0 \sim N(5, 1.5) \qquad (16)$$
$$X_i \mid Z_i = 1 \sim N(10, 2) \qquad (17)$$

Let the true mixture proportions be $P(Z_i = 0) = 0.25$ and $P(Z_i = 1) = 0.75$, respectively.

# Exercise

1. Simulate data from the mixture model on the previous slide, which should produce the following histogram.



**Histogram of X**

# Exercise

Compute the likelihood $P(X_i|Z_i = 0)$ and $P(X_i|Z_i = 1)$ and store these in a matrix $L$.

# Exercise

Implement the E and M step in a function called `emIteration`.
Then evaluate the EM and verify that your estimates are 0.29 and
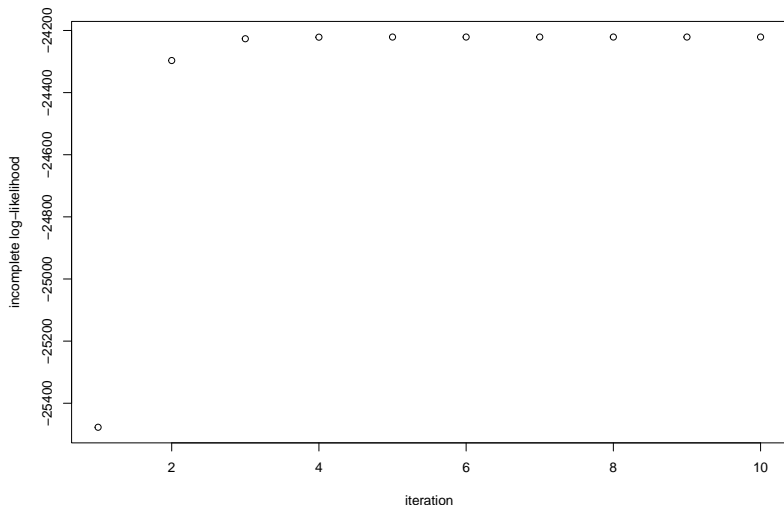0.71, respectively.

Plot the incomplete log-likelihood versus the iteration. What do you
observe regarding its behavior.

# Perform EM

```
#perform EM
ee <- mixture.EM(w.init=c(0.5,0.5), L)
print(paste("Estimate = (", round(ee[[1]][1],2), ",",
            round(ee[[1]][2],2), ")", sep=""))
```

```
## [1] "Estimate = (0.22,0.78)"
```

# Plot



The log-likelihood is strictly increasing, meaning that we have reached a local maxima.

# R packages for mixture models

▶ The `mclust` package
(http://www.stat.washington.edu/mclust/) is standard for
Gaussian mixtures.

▶ The `mixtools` considers classic parametric densities, mixtures
of regressions, and some non-parametric mixtures.