

Methods for Unsupervised Clustering

W. Evan Johnson, Ph.D.
Professor, Division of Infectious Disease
Director, Center for Data Science
Rutgers University – New Jersey Medical School

8/25/2022

Supervised vs. Unsupervised Machine Learning

Machine learning algorithms are generally classified into two categories. Most of the algorithms we have described up to now are examples of a general approach referred to as **supervised** machine learning. The name comes from the fact that we use the outcomes in a training set to **supervise** the creation of our prediction algorithm.

There is another subset of machine learning methods referred to as **unsupervised**. In this subset we do not necessarily know the outcomes and instead are interested in discovering groups. These algorithms are also referred to as **clustering** algorithms since predictors are used to define **clusters**.

Classification vs. Clustering

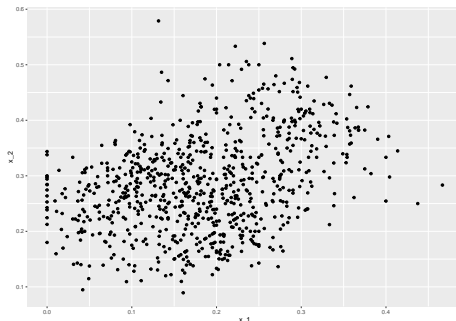
Classification and **clustering** are two methods of pattern identification used in machine learning.

Although both techniques have certain similarities, the difference lies in the fact that *classification* uses predefined classes in which objects are assigned, while *clustering* identifies similarities between objects, which it groups according to those characteristics in common and which differentiate them from other groups of objects. These groups are known as “clusters”.

Supervised vs. Unsupervised Machine Learning

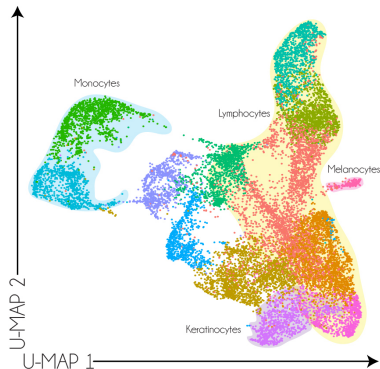
Sometimes clustering is not be very useful. For example, if we are simply given the heights we may not be able to discover two groups, males and females. The digits example would also be also be challenging:

```
library(tidyverse)
library(dslabs)
data("mnist_27")
mnist_27$train %>% qplot(x_1, x_2, data = .)
```



Unsupervised Machine Learning

However, there are applications in which unsupervised learning can be a powerful technique, in particular as an exploratory tool, and can be combined with multiple methods.



Unsupervised Machine Learning

There are many algorithms for unsupervised learning. We have already learned about **PCA** and **UMAP** for dimension reduction. Here we introduce two methods for clustering: **hierarchical clustering** and **k-means**.

Unsupervised Machine Learning and clustering

A first step in any clustering algorithm is defining a distance between observations or groups of observations. Then we need to decide how to join observations into clusters.

Hierarchical clustering starts by defining each observation as a separate group, and distances are calculated between every group (distance matrix). Then the two closest groups are merged into a single group, and this new group (two observations) is represented by its centroid. Distances between this new group and the rest are calculated, and then the next two closest groups are merged. This process is repeated until there is just one group including all the observations.

Hierarchical clustering

Consider the ratings of 50 movies from 139 different critics:

```
library(dslabs); data("movielens")
top <- movielens %>% group_by(movieId) %>%
  summarize(n=n(), title = first(title)) %>%
  top_n(50, n) %>% pull(movieId)

x <- movielens %>% filter(movieId %in% top) %>%
  group_by(userId) %>% filter(n() >= 25) %>%
  ungroup() %>% select(title, userId, rating) %>%
  spread(userId, rating)

row_names <- str_remove(x$title, ": Episode") %>% str_trunc(20)
x <- x[,-1] %>% as.matrix()
x <- sweep(x, 2, colMeans(x, na.rm = TRUE))
x <- sweep(x, 1, rowMeans(x, na.rm = TRUE))
rownames(x) <- row_names
```


Hierarchical clustering

We want to use these data to find out if there are clusters of movies based on the ratings from 139 movie raters. A first step is to find the distance between each pair of movies using the `dist` function:

```
d <- dist(x)
```

Hierarchical clustering

With the distance between each pair of movies computed, we need an algorithm to define groups from these. The `hclust` function implements this algorithm and it takes a distance as input.

```
h <- hclust(d)
```

```
h
```

```
##
```

```
## Call:
```

```
## hclust(d = d)
```

```
##
```

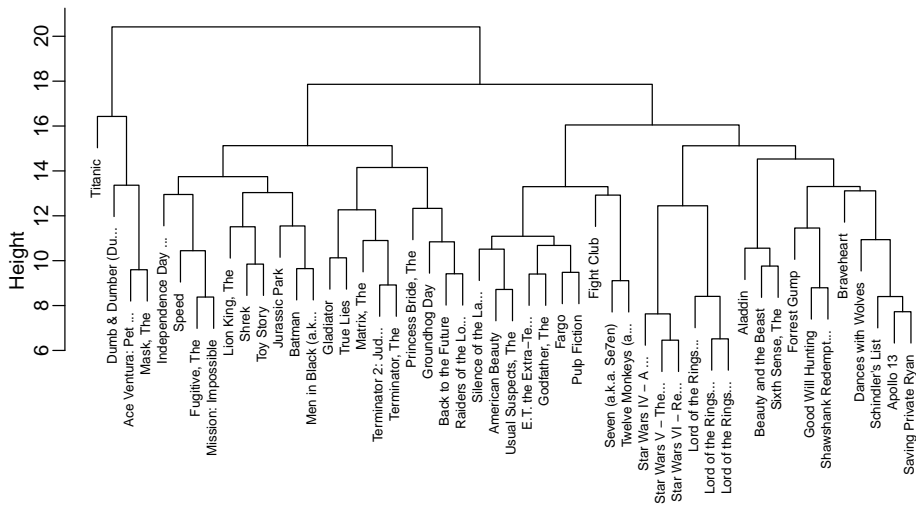
```
## Cluster method      : complete
```

```
## Distance             : euclidean
```

```
## Number of objects: 50
```

Hierarchical clustering

We can see the resulting groups using a **dendrogram**.

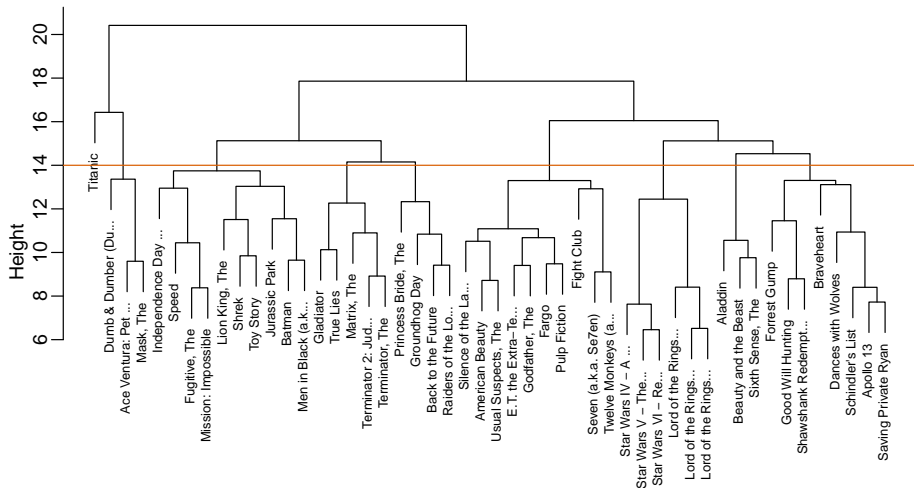


Hierarchical clustering

This graph gives us an approximation between the distance between any two movies. To find this distance we find the first location, from top to bottom, where these movies split into two different groups. The height of this location is the distance between these two groups. So, for example, the distance between the three *Star Wars* movies is 8 or less, while the distance between *Raiders of the Lost of Ark* and *Silence of the Lambs* is about 17.

Hierarchical clustering

To generate actual groups: 1) decide on a maximum distance to be in the same group or 2) decide on the number of groups. For example:



Hierarchical clustering

The function `cutree` can be applied to the output of `hclust` to perform either of these two operations and generate groups.

```
# Maximum Distance
```

```
groups <- cutree(h, h = 14)  
table(groups)
```

```
## groups  
##  1  2  3  4  5  6  7  8  9  
##  3  3 10  8  4 10  5  6  1
```

```
#Number of groups
```

```
groups <- cutree(h, k = 10)  
table(groups)
```

```
## groups  
##  1  2  3  4  5  6  7  8  9 10  
##  3  3 10  8  4  6  4  5  6  1
```

Hierarchical clustering

The clustering provides some insights, e.g., Group 4 appears to be blockbusters:

```
names(groups)[groups==4]
```

```
## [1] "Apollo 13"           "Braveheart"          "Dances with Wolves"  
## [4] "Forrest Gump"        "Good Will Hunting"   "Saving Private Ryan"  
## [7] "Schindler's List"    "Shawshank Redempt..."
```

And group 9 appears to be fantasy/nerd movies:

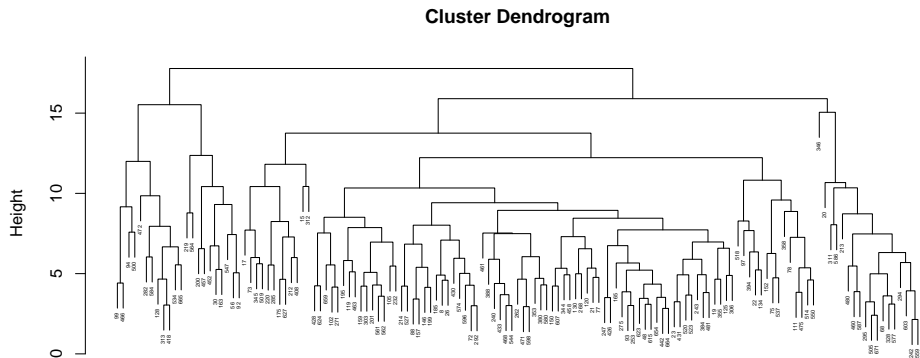
```
names(groups)[groups==9]
```

```
## [1] "Lord of the Rings..." "Lord of the Rings..." "Lord of the Rings..."  
## [4] "Star Wars IV - A ..." "Star Wars V - The..." "Star Wars VI - Re..."
```

Hierarchical clustering

We can also explore the data to see if there are clusters of movie raters.

```
h_2 <- dist(t(x)) %>% hclust()  
plot(h_2, cex = 0.35)
```



k-means

We can also use the k-means algorithm. Here, we have to pre-define k , the number of clusters we want to define.

The k-means algorithm is iterative. The first step is to define k centers. Then each observation is assigned to the cluster with the closest center to that observation. In a second step the centers are redefined using the observation in each cluster: the column means are used to define a **centroid**. We repeat these two steps until the centers converge.

k-means

The `kmeans` function included in R-base does not handle NAs. For illustrative purposes we will fill out the NAs with 0s. In general, the choice of how to fill in missing data, or if one should do it at all, should be made with care.

```
x_0 <- x; x_0[is.na(x_0)] <- 0
set.seed(0)
k <- kmeans(x_0, centers = 10)
```

k-means

The cluster assignments are in the cluster component:

```
groups <- k$cluster  
table(groups)
```

```
## groups  
##  1  2  3  4  5  6  7  8  9 10  
##  1  9  3  3  4  3  3  5  6 13
```

k-means

This yields some interesting groups:

```
names(groups)[groups==4]
```

```
## [1] "Ace Ventura: Pet ..." "Dumb & Dumber (Du..." "Mask, The"
```

```
names(groups)[groups==6]
```

```
## [1] "Lord of the Rings..." "Lord of the Rings..." "Lord of the Rings..."
```

```
names(groups)[groups==7]
```

```
## [1] "Star Wars IV - A ..." "Star Wars V - The..." "Star Wars VI - Re..."
```

```
names(groups)[groups==9]
```

```
## [1] "Braveheart" "Dances with Wolves" "Godfather, The"
## [4] "Good Will Hunting" "Schindler's List" "Shawshank Redempt..."
```

```
names(groups)[groups==10]
```

```
## [1] "Aladdin" "Apollo 13" "Beauty and the Beast"
## [4] "E.T. the Extra-Te..." "Forrest Gump" "Gladiator"
## [7] "Lion King, The" "Princess Bride, The" "Raiders of the Lo..."
## [10] "Saving Private Ryan" "Shrek" "Sixth Sense, The"
## [13] "Toy Story"
```

k-means

Note that because the first center is chosen at random, the final clusters are random. We impose some stability by repeating the entire function several times and averaging the results. The number of random starting values to use can be assigned through the `nstart` argument.

```
k <- kmeans(x_0, centers = 10, nstart = 25)
```

Heatmaps

A powerful visualization tool for discovering clusters or patterns in your data is the heatmap. The idea is simple: plot an image of your data matrix with colors used as the visual cue and both the columns and rows ordered according to the results of a clustering algorithm. We will demonstrate this with the `tissue_gene_expression` dataset. We will scale the rows of the gene expression matrix.

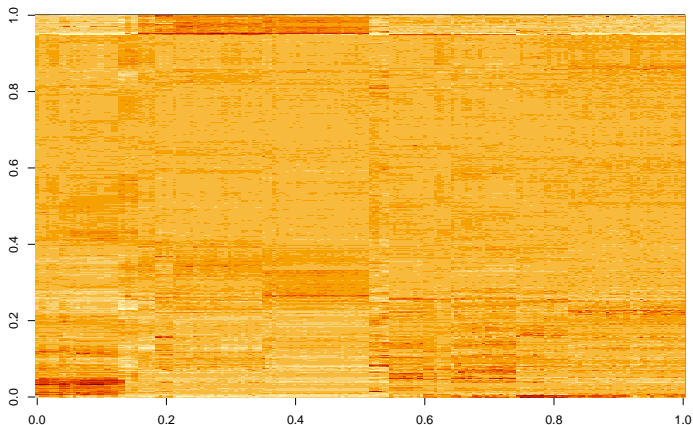
The first step is compute:

```
data("tissue_gene_expression")
x <- sweep(tissue_gene_expression$x, 2,
           colMeans(tissue_gene_expression$x))
h_1 <- hclust(dist(x))
h_2 <- hclust(dist(t(x)))
```

Heatmaps

Now we can use the results of this clustering to order the rows and columns.

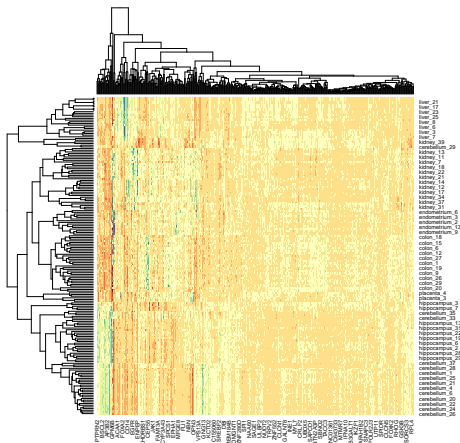
```
image(x[h_1$order, h_2$order])
```



Heatmaps

But there is `heatmap` function that does it for us:

```
heatmap(x, col = RColorBrewer::brewer.pal(11, "Spectral"))
```



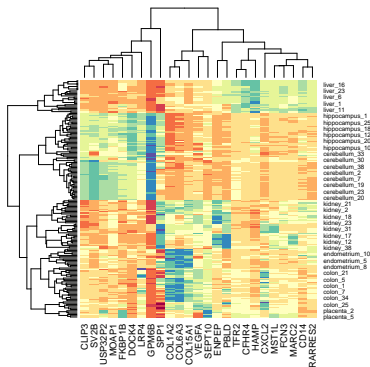
Filtering features

If the information about clusters is included in just a few features, including all the features can add enough noise that detecting clusters becomes challenging. One simple approach to try to remove features with no information is to only include those with high variance. In the movie example, a user with low variance in their ratings is not really informative: all the movies seem about the same to them.

Filtering features

For example, if we include only features (genes) with highest variance.

```
library(matrixStats)
sds <- colSds(x, na.rm = TRUE)
o <- order(sds, decreasing = TRUE)[1:25]
heatmap(x[,o], col = RColorBrewer::brewer.pal(11, "Spectral"))
```



Session Info

```
sessionInfo()
```

```
## R version 4.2.3 (2023-03-15)
## Platform: aarch64-apple-darwin20 (64-bit)
## Running under: macOS Ventura 13.5.1
##
## Matrix products: default
## BLAS: /Library/Frameworks/R.framework/Versions/4.2-arm64/Resources/lib/libRblas.0.dylib
## LAPACK: /Library/Frameworks/R.framework/Versions/4.2-arm64/Resources/lib/libRlapack.dylib
##
## locale:
## [1] en_US.UTF-8/en_US.UTF-8/en_US.UTF-8/C/en_US.UTF-8/en_US.UTF-8
##
## attached base packages:
## [1] stats      graphics  grDevices  utils      datasets  methods   base
##
## other attached packages:
## [1] matrixStats_1.0.0 dslabs_0.7.6      lubridate_1.9.3   forcats_1.0.0
## [5] stringr_1.5.0      dplyr_1.1.3       purrr_1.0.2       readr_2.1.4
## [9] tidyr_1.3.0        tibble_3.2.1      tidyverse_2.0.0   caret_6.0-94
## [13] lattice_0.21-9     ggplot2_3.4.3
##
## loaded via a namespace (and not attached):
## [1] Rcpp_1.0.11        listenv_0.9.0      class_7.3-22
## [4] digest_0.6.33      ipred_0.9-14       foreach_1.5.2
## [7] utf8_1.2.3         parallelly_1.36.0  R6_2.5.1
## [10] plyr_1.8.9         hardhat_1.3.0      stats4_4.2.3
## [13] evaluate_0.22      pillar_1.9.0       rlang_1.1.1
## [16] rstudioapi_0.15.0  data.table_1.14.8  rpart_4.1.19
## [19] Matrix_1.5-4.1     rafalib_1.0.0      rmarkdown_2.25
## [22] labeling_0.4.3     splines_4.2.3      gower_1.0.1
## [25] munsell_0.5.0      compiler_4.2.3     xfun_0.40
```