

(<https://s3-ap-south-1.amazonaws.com/av-blog-media/wp-content/uploads/2019/05/Screenshot-from-2019-05-15-12-34-48.png>).

This technique is generally used for clustering a population into different groups. A few common examples include segmenting customers, clustering similar documents together, recommending similar songs or movies, etc.

There are a LOT more applications of unsupervised learning. If you come across any interesting application, feel free to share them in the comments section below!

Now, there are various algorithms that help us to make these clusters. The most commonly used clustering algorithms are K-means and Hierarchical clustering.

Why Hierarchical Clustering?

We should first know how K-means works before we dive into hierarchical clustering. Trust me, it will make the concept of hierarchical clustering all the more easier.

Here's a brief overview of how K-means works:

1. Decide the number of clusters (k)
2. Select k random points from the data as centroids
3. Assign all the points to the nearest cluster centroid
4. Calculate the centroid of newly formed clusters
5. Repeat steps 3 and 4

Enroll Now

It is an iterative process. It will keep on running until the centroids of newly formed clusters do not change or the maximum number of iterations are reached.

But there are certain challenges with K-means. It always tries to make clusters of the same size. Also, we have to decide the number of clusters at the *beginning* of the algorithm. Ideally, we would not know how many clusters should we have, in the beginning of the algorithm and hence it a challenge with K-means.

This is a gap hierarchical clustering bridges with aplomb. It takes away the problem of having to pre-define the number of clusters. Sounds like a dream! So, let's see what hierarchical clustering is and how it improves on K-means.

What is Hierarchical Clustering?

Let's say we have the below points and we want to cluster them into groups:



<https://s3-ap-south-1.amazonaws.com/av-blog-media/wp-content/uploads/2019/05/Screenshot-from-2019-05-15-13-10-32.png>

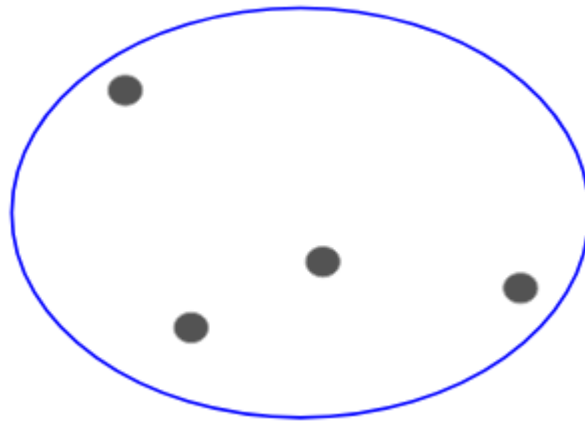
We can assign each of these points to a separate cluster:



<https://s3-ap-south-1.amazonaws.com/av-blog-media/wp-content/uploads/2019/05/Screenshot-from-2019-05-15-13-11-28.png>

Now, based on the similarity of these clusters, we can combine the most similar clusters together and repeat this process until only a single cluster is left:

Enroll Now



<https://s3-ap-south-1.amazonaws.com/av-blog-media/wp-content/uploads/2019/05/Screenshot-from-2019-05-15-13-12-35.png>

We are essentially building a hierarchy of clusters. That's why this algorithm is called hierarchical clustering. I will discuss how to decide the number of clusters in a later section. For now, let's look at the different types of hierarchical clustering.

Types of Hierarchical Clustering

There are mainly two types of hierarchical clustering:

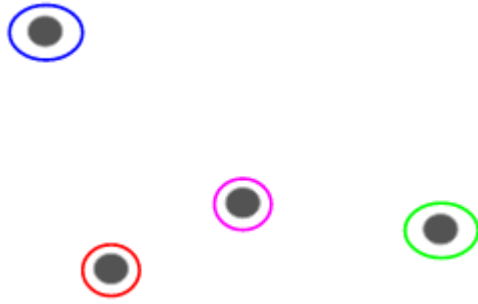
1. Agglomerative hierarchical clustering
2. Divisive Hierarchical clustering

Let's understand each type in detail.

Agglomerative Hierarchical Clustering

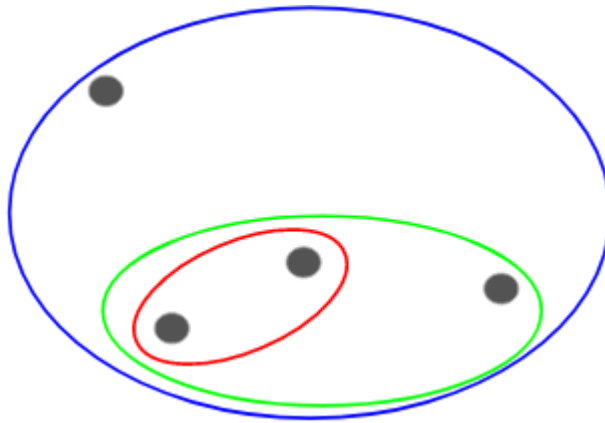
We assign each point to an individual cluster in this technique. Suppose there are 4 data points. We will assign each of these points to a cluster and hence will have 4 clusters in the beginning:

Enroll Now



<https://s3-ap-south-1.amazonaws.com/av-blog-media/wp-content/uploads/2019/05/Screenshot-from-2019-05-15-13-11-28.png>

Then, at each iteration, we merge the closest pair of clusters and repeat this step until only a single cluster is left:



<https://s3-ap-south-1.amazonaws.com/av-blog-media/wp-content/uploads/2019/05/Screenshot-from-2019-05-15-13-31-06.png>

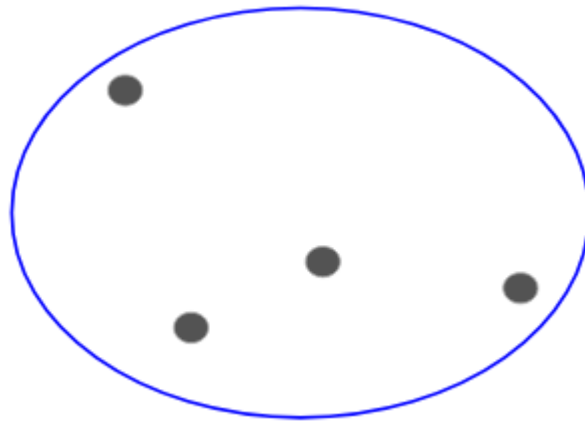
We are merging (or adding) the clusters at each step, right? Hence, this type of clustering is also known as **additive hierarchical clustering**.

Divisive Hierarchical Clustering

Divisive hierarchical clustering works in the opposite way. Instead of starting with n clusters (in case of n observations), we start with a single cluster and assign all the points to that cluster.

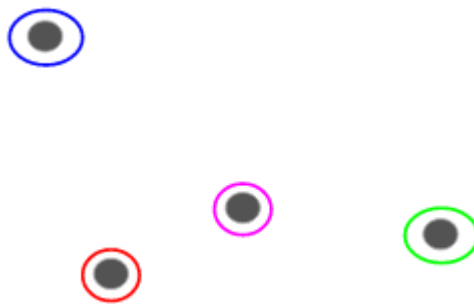
So, it doesn't matter if we have 10 or 1000 data points. All these points will belong to the same cluster at the beginning:

Enroll Now



<https://s3-ap-south-1.amazonaws.com/av-blog-media/wp-content/uploads/2019/05/Screenshot-from-2019-05-15-13-12-35.png>

Now, at each iteration, we split the farthest point in the cluster and repeat this process until each cluster only contains a single point:



We are splitting (or dividing) the clusters at each step, hence the name divisive hierarchical clustering.

Agglomerative Clustering is widely used in the industry and that will be the focus in this article. Divisive hierarchical clustering will be a piece of cake once we have a handle on the agglomerative type.

Steps to Perform Hierarchical Clustering

We merge the most similar points or clusters in hierarchical clustering – we know this. Now the question is – how do we decide which points are similar and which are not? It's one of the most important questions in clustering!

Here's one way to calculate similarity – Take the distance between the centroids of these clusters. The points having the least distance are referred to as similar points and we can merge them. We can refer to this as a **distance-based algorithm** as well (since we are calculating the distances between the clusters).

Enroll Now

In hierarchical clustering, we have a concept called a **proximity matrix**. This stores the distances between each point. Let's take an example to understand this matrix as well as the steps to perform hierarchical clustering.

Setting up the Example



(<https://s3-ap-south-1.amazonaws.com/av-blog-media/wp-content/uploads/2019/05/teacher-student.jpg>).

Suppose a teacher wants to divide her students into different groups. She has the marks scored by each student in an assignment and based on these marks, she wants to segment them into groups. There's no fixed target here as to how many groups to have. Since the teacher does not know what type of students should be assigned to which group, it cannot be solved as a supervised learning problem. So, we will try to apply hierarchical clustering here and segment the students into different groups.

Let's take a sample of 5 students:

Enroll Now

Student_ID	Marks
1	10
2	7
3	28
4	20
5	35

<https://s3-ap-south-1.amazonaws.com/av-blog-media/wp-content/uploads/2019/05/Screenshot-from-2019-05-17-16-12-33.png>

Creating a Proximity Matrix

First, we will create a proximity matrix which will tell us the distance between each of these points. Since we are calculating the distance of each point from each of the other points, we will get a square matrix of shape $n \times n$ (where n is the number of observations).

Let's make the 5×5 proximity matrix for our example:

ID	1	2	3	4	5
1	0	3	18	10	25
2	3	0	21	13	28
3	18	21	0	8	7
4	10	13	8	0	15
5	25	28	7	15	0

<https://s3-ap-south-1.amazonaws.com/av-blog-media/wp-content/uploads/2019/05/Screenshot-from-2019-05-15-14-46-36.png>

Enroll Now

The diagonal elements of this matrix will always be 0 as the distance of a point with itself is always 0. We will use the Euclidean distance formula to calculate the rest of the distances. So, let's say we want to calculate the distance between point 1 and 2:

$$\sqrt{(10-7)^2} = \sqrt{9} = 3$$

Similarly, we can calculate all the distances and fill the proximity matrix.

Steps to Perform Hierarchical Clustering

Step 1: First, we assign all the points to an individual cluster:



(<https://s3-ap-south-1.amazonaws.com/av-blog-media/wp-content/uploads/2019/05/Screenshot-from-2019-05-15-14-54-20.png>).

Different colors here represent different clusters. You can see that we have 5 different clusters for the 5 points in our data.

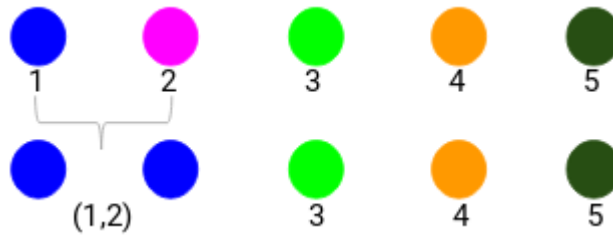
Step 2: Next, we will look at the smallest distance in the proximity matrix and merge the points with the smallest distance. We then update the proximity matrix:

ID	1	2	3	4	5
1	0	3	18	10	25
2	3	0	21	13	28
3	18	21	0	8	7
4	10	13	8	0	15
5	25	28	7	15	0

(<https://s3-ap-south-1.amazonaws.com/av-blog-media/wp-content/uploads/2019/05/Screenshot-from-2019-05-15-14-57-30.png>).

Here, the smallest distance is 3 and hence we will merge point 1 and 2:

Enroll Now



Let's look at the updated clusters and accordingly update the proximity matrix:

Student_ID	Marks
(1,2)	10
3	28
4	20
5	35

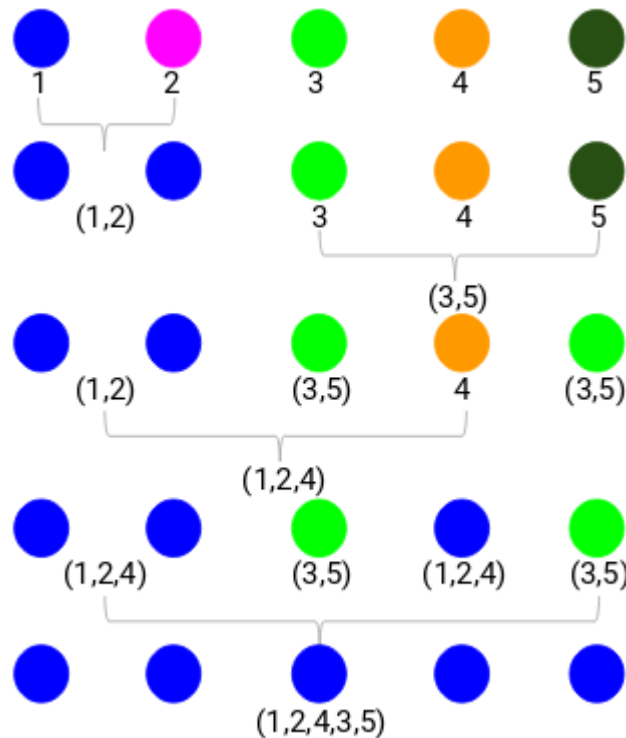
Here, we have taken the maximum of the two marks (7, 10) to replace the marks for this cluster. Instead of the maximum, we can also take the minimum value or the average values as well. Now, we will again calculate the proximity matrix for these clusters:

ID	(1,2)	3	4	5
(1,2)	0	18	10	25
3	18	0	8	7
4	10	8	0	15
5	25	7	15	0

Step 3: We will repeat step 2 until only a single cluster is left.

So, we will first look at the minimum distance in the proximity matrix and then merge the closest pair of clusters. We will get the merged clusters as shown below after repeating these steps:

Enroll Now



(<https://s3-ap-south-1.amazonaws.com/av-blog-media/wp-content/uploads/2019/05/Screenshot-from-2019-05-15-15-03-17.png>).

We started with 5 clusters and finally have a single cluster. **This is how agglomerative hierarchical clustering works.** But the burning question still remains – how do we decide the number of clusters? Let's understand that in the next section.

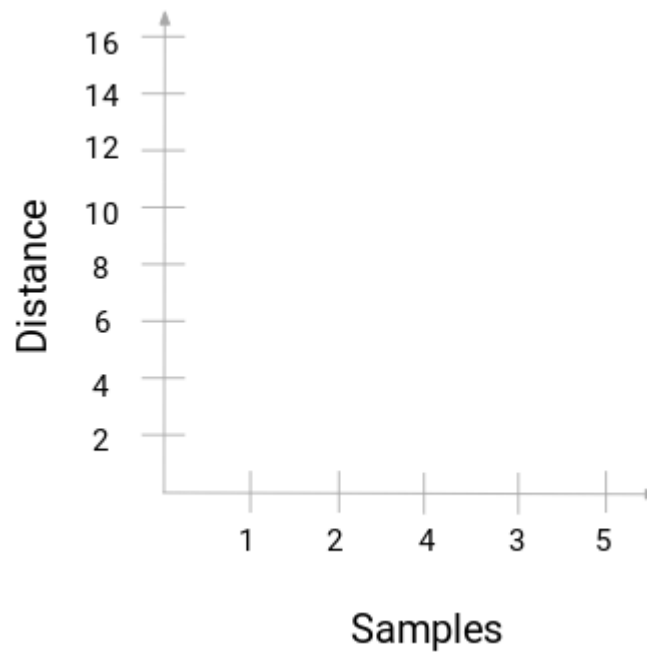
How should we Choose the Number of Clusters in Hierarchical Clustering?

Ready to finally answer this question that's been hanging around since we started learning? To get the number of clusters for hierarchical clustering, we make use of an awesome concept called a **Dendrogram**.

A dendrogram is a tree-like diagram that records the sequences of merges or splits.

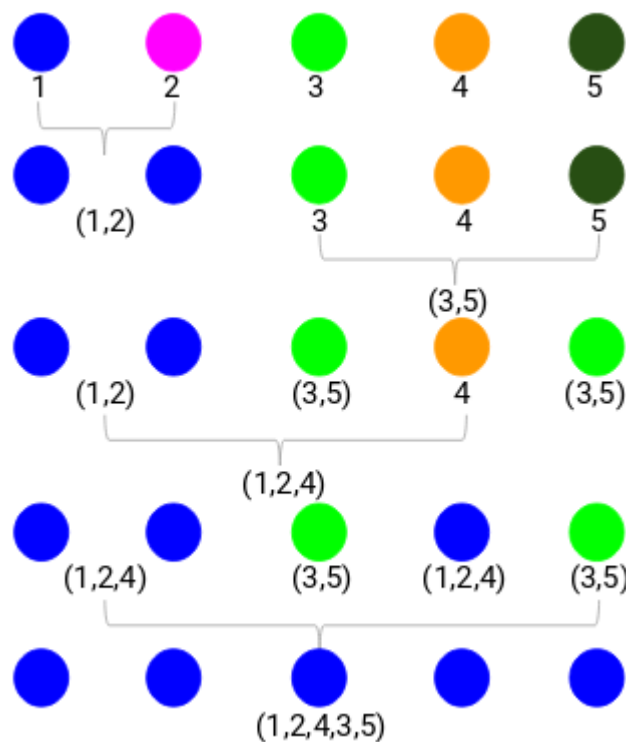
Let's get back to our teacher-student example. Whenever we merge two clusters, a dendrogram will record the distance between these clusters and represent it in graph form. Let's see how a dendrogram looks like:

Enroll Now



(<https://s3-ap-south-1.amazonaws.com/av-blog-media/wp-content/uploads/2019/05/Screenshot-from-2019-05-16-17-11-39.png>).

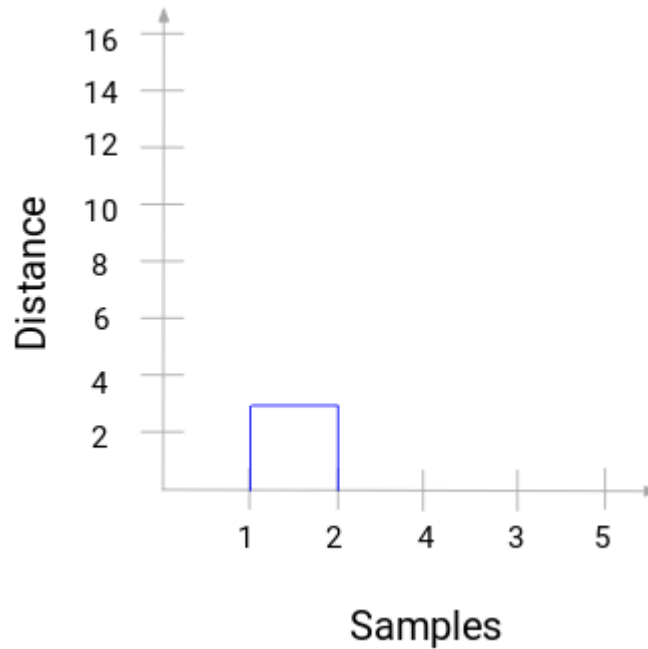
We have the samples of the dataset on the x-axis and the distance on the y-axis. **Whenever two clusters are merged, we will join them in this dendrogram and the height of the join will be the distance between these points.** Let's build the dendrogram for our example:



(<https://s3-ap-south-1.amazonaws.com/av-blog-media/wp-content/uploads/2019/05/Screenshot-from-2019-05-15-15-03-17.png>).

Enroll Now

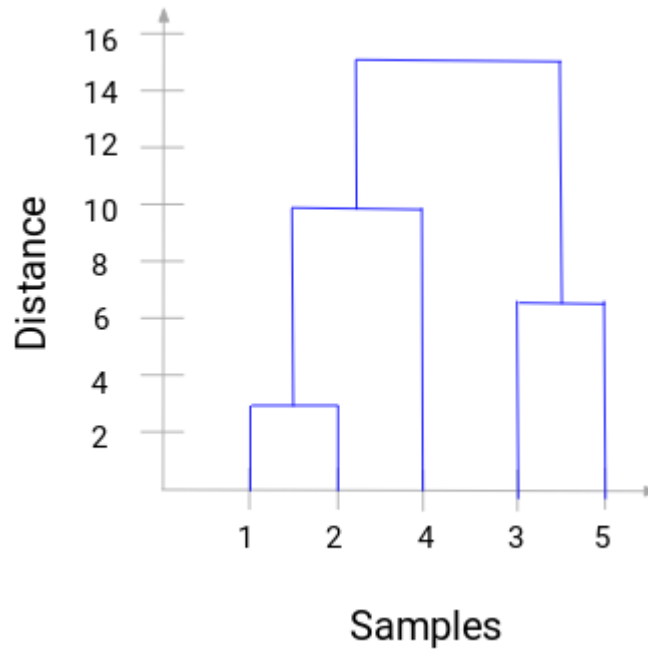
Take a moment to process the above image. We started by merging sample 1 and 2 and the distance between these two samples was 3 (refer to the first proximity matrix in the previous section). Let's plot this in the dendrogram:



<https://s3-ap-south-1.amazonaws.com/av-blog-media/wp-content/uploads/2019/05/Screenshot-from-2019-05-16-18-06-37.png>

Here, we can see that we have merged sample 1 and 2. The vertical line represents the distance between these samples. Similarly, we plot all the steps where we merged the clusters and finally, we get a dendrogram like this:

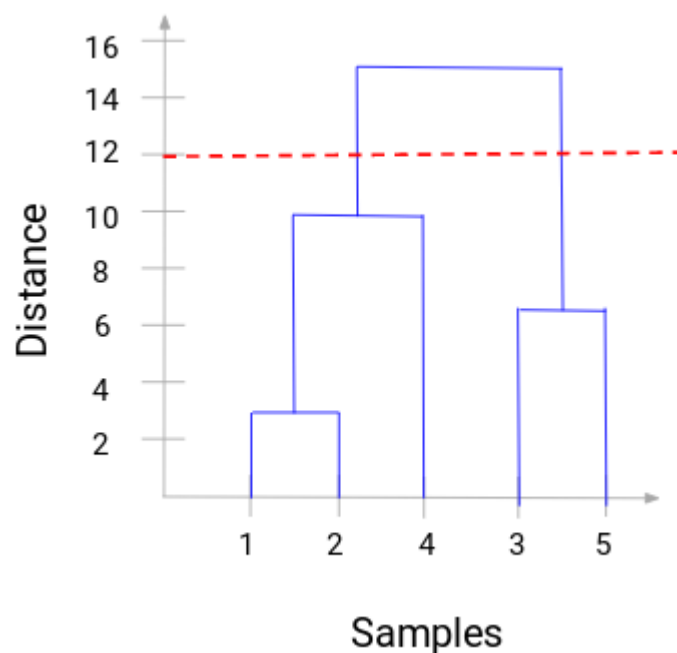
Enroll Now



<https://s3-ap-south-1.amazonaws.com/av-blog-media/wp-content/uploads/2019/05/Screenshot-from-2019-05-16-18-10-06.png>

We can clearly visualize the steps of hierarchical clustering. **More the distance of the vertical lines in the dendrogram, more the distance between those clusters.**

Now, we can set a threshold distance and draw a horizontal line (*Generally, we try to set the threshold in such a way that it cuts the tallest vertical line*). Let's set this threshold as 12 and draw a horizontal line:



Enroll Now

(<https://s3-ap-south-1.amazonaws.com/av-blog-media/wp-content/uploads/2019/05/Screenshot-from-2019-05-16-18-12-55.png>).

The number of clusters will be the number of vertical lines which are being intersected by the line drawn using the threshold. In the above example, since the red line intersects 2 vertical lines, we will have 2 clusters. One cluster will have a sample (1,2,4) and the other will have a sample (3,5). Pretty straightforward, right?

This is how we can decide the number of clusters using a dendrogram in Hierarchical Clustering. In the next section, we will implement hierarchical clustering which will help you to understand all the concepts that we have learned in this article.

Solving the Wholesale Customer Segmentation problem using Hierarchical Clustering

Time to get our hands dirty in Python!

We will be working on a wholesale customer segmentation problem. You can download the dataset using [this link \(https://archive.ics.uci.edu/ml/machine-learning-databases/00292/Wholesale%20customers%20data.csv\)](https://archive.ics.uci.edu/ml/machine-learning-databases/00292/Wholesale%20customers%20data.csv). The data is hosted on the UCI Machine Learning repository. The aim of this problem is to segment the clients of a wholesale distributor based on their annual spending on diverse product categories, like milk, grocery, region, etc.

Let's explore the data first and then apply Hierarchical Clustering to segment the clients.

We will first import the required libraries:

```
1 import pandas as pd
2 import numpy as np
3 import matplotlib.pyplot as plt
4 %matplotlib inline
```

[om/PulkitS01/8ac9bf3b54eb59b4e1d4eaa21d3d774e/raw/6cea281dc4dea42bbcb2160e6cef1535cad765e7/importing_libraries.py](https://gist.github.com/PulkitS01/8ac9bf3b54eb59b4e1d4eaa21d3d774e/raw/6cea281dc4dea42bbcb2160e6cef1535cad765e7/importing_libraries.py)
importing_libraries.py (https://gist.github.com/PulkitS01/8ac9bf3b54eb59b4e1d4eaa21d3d774e#file-importing_libraries-py)
hosted with ❤ by GitHub (<https://github.com>)

Load the data and look at the first few rows:

```
1 data = pd.read_csv('Wholesale customers data.csv')
2 data.head()
```

Enroll Now