

# Clustering Analysis

Abby Smith

11/10/2020

```
knitr::opts_chunk$set(echo=F, cache=T, results='hide', message = F, warning = F)
```

## Questions of Interest

Edovo is interested in a clustering of the course content. They have over 15,000 courses. Each course has tags attached to it, representing different topics that are covered in the course (such as “CHRISTIANITY”), but also features of the course itself (such as “INTERMEDIATE”).

We have the tags for each course, and the navigation map of courses.

We are interested in:

- What courses have similar tags?
- What courses have similar features?
- What are the most important features in determining similar tags?
- How can we incorporate the navigation map in our clustering?
- What do the users look like that are in a given tag cluster? In terms of time as active user, facility, etc.

(**UPDATE 11/16**): Not only do we want a clustering analysis of users + content; we want to determine an *engagement score* for each user on the platform.

- What clusters of content are “quality”? Can we provide a quantitative metric for that?
- How can we collapse this into a score?

## The Data

We get the data into a **wide** format- each row represents a course, its title, lesson count, total page count, number of items across the pages and then columns representing each tag and a corresponding column with that tag’s ID. Because a course has up to 32 different tags, there are  $32 \times 2 = 64$  different columns.

## Question of Representation: Long vs. Wide?

## Key Decisions in Clustering

- Do we one-hot encode categorical variables?
- How many clusters do we use?
- Should we use a joint dimension reduction + clustering approach?
- How to deal with many, many categorical features (tags)?

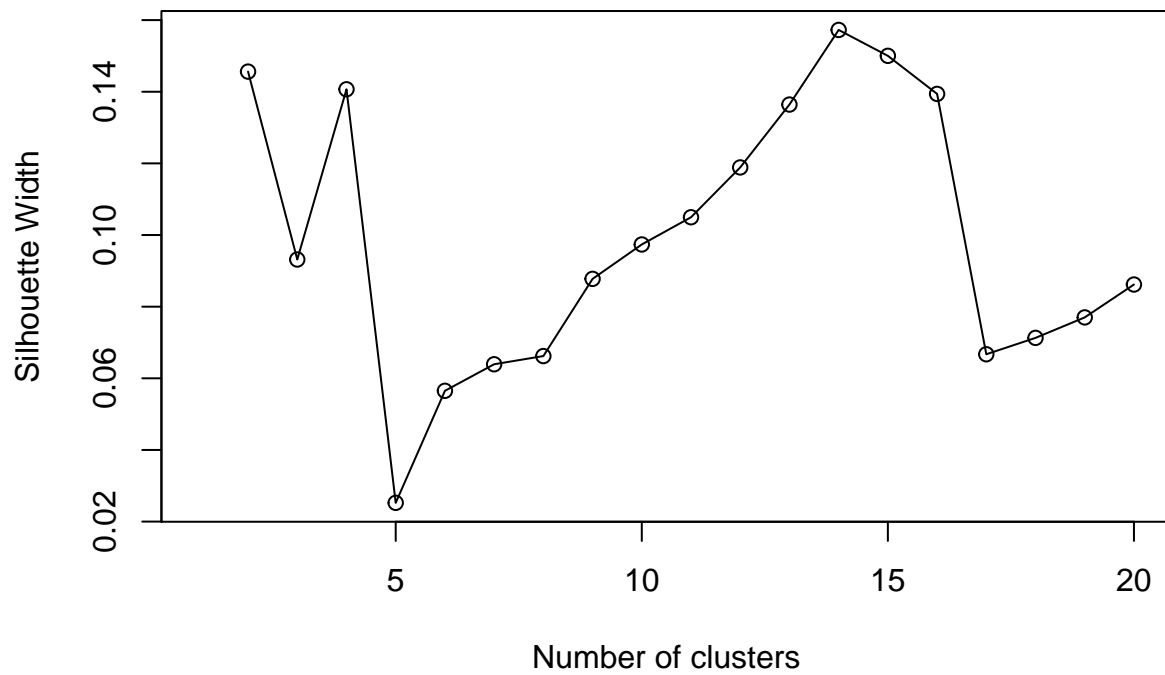
## Clustering Approaches

In our clustering we want to prioritize **numerical** features over **categorical**. More specifically- we want to account for hierarchy of the course in the JSON navigation map. If two users both watched just 2 pages of a very specifically tagged activity- we want to weight that higher than tags that are further up in the hierarchy.

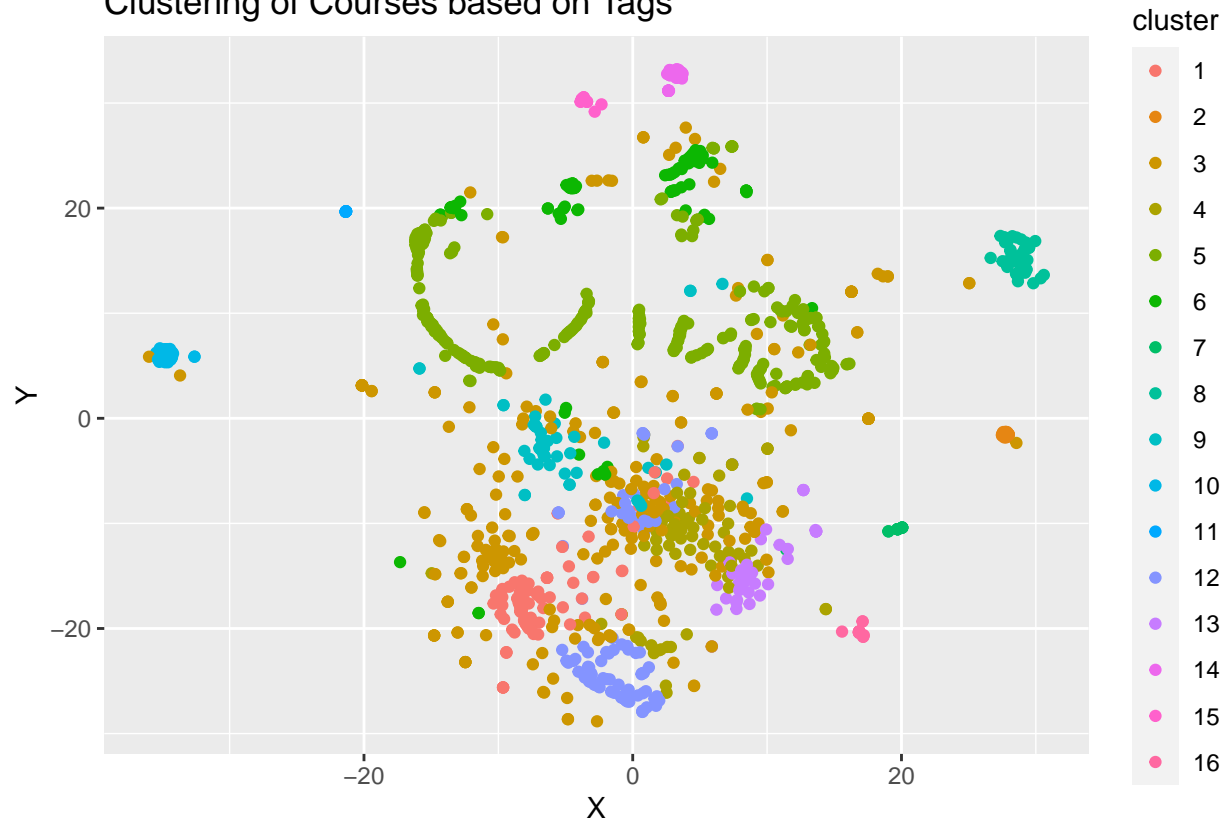
## K-Means with Gower Distance

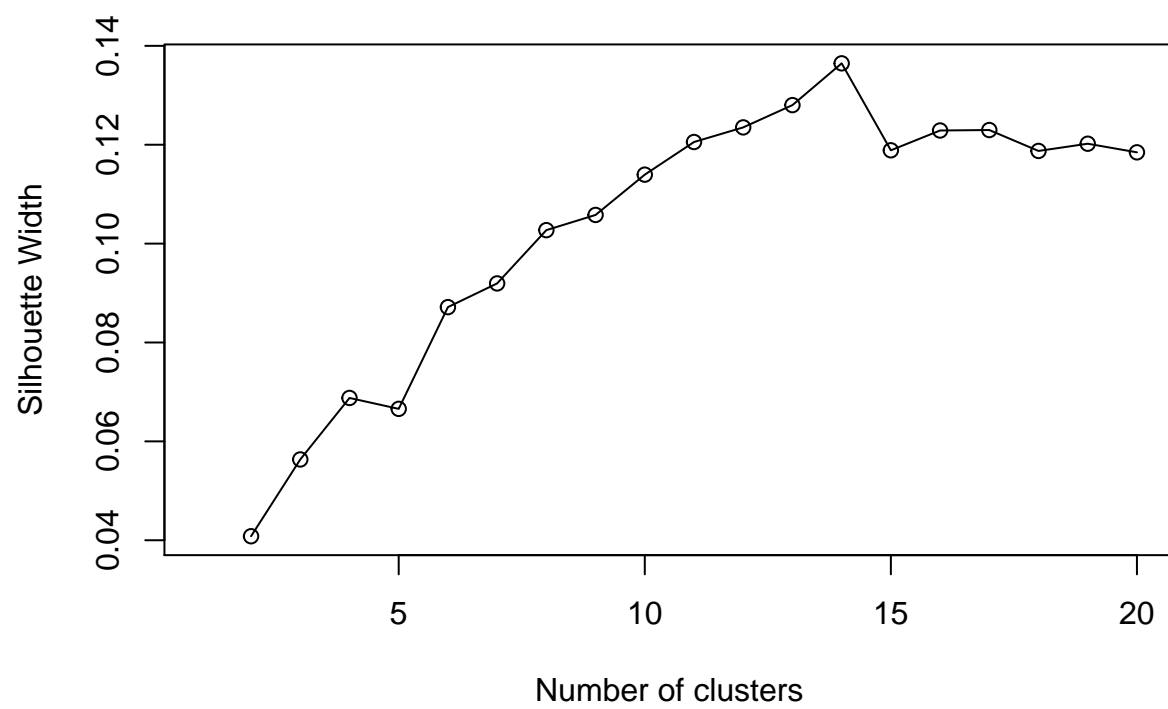
In order for a yet-to-be-chosen algorithm to group observations together, we first need to define some notion of (dis)similarity between observations. A popular choice for clustering is Euclidean distance. However, Euclidean distance is only valid for continuous variables, and thus is not applicable here. In order for a clustering algorithm to yield sensible results, we have to use a distance metric that can handle mixed data types. In this case, we will use something called Gower distance.

We will look at a few different approaches from this clustering. We will plot a low dimensional representation for each of these clusterings

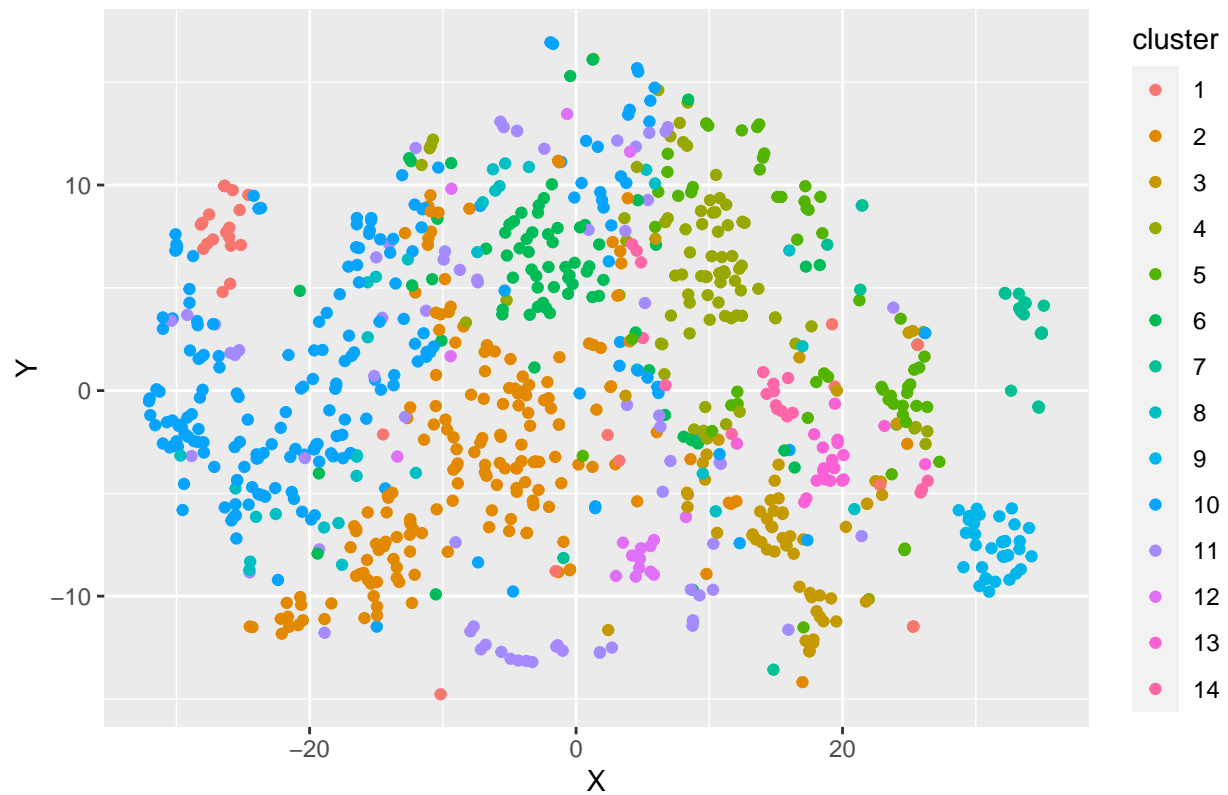


Clustering of Courses based on Tags



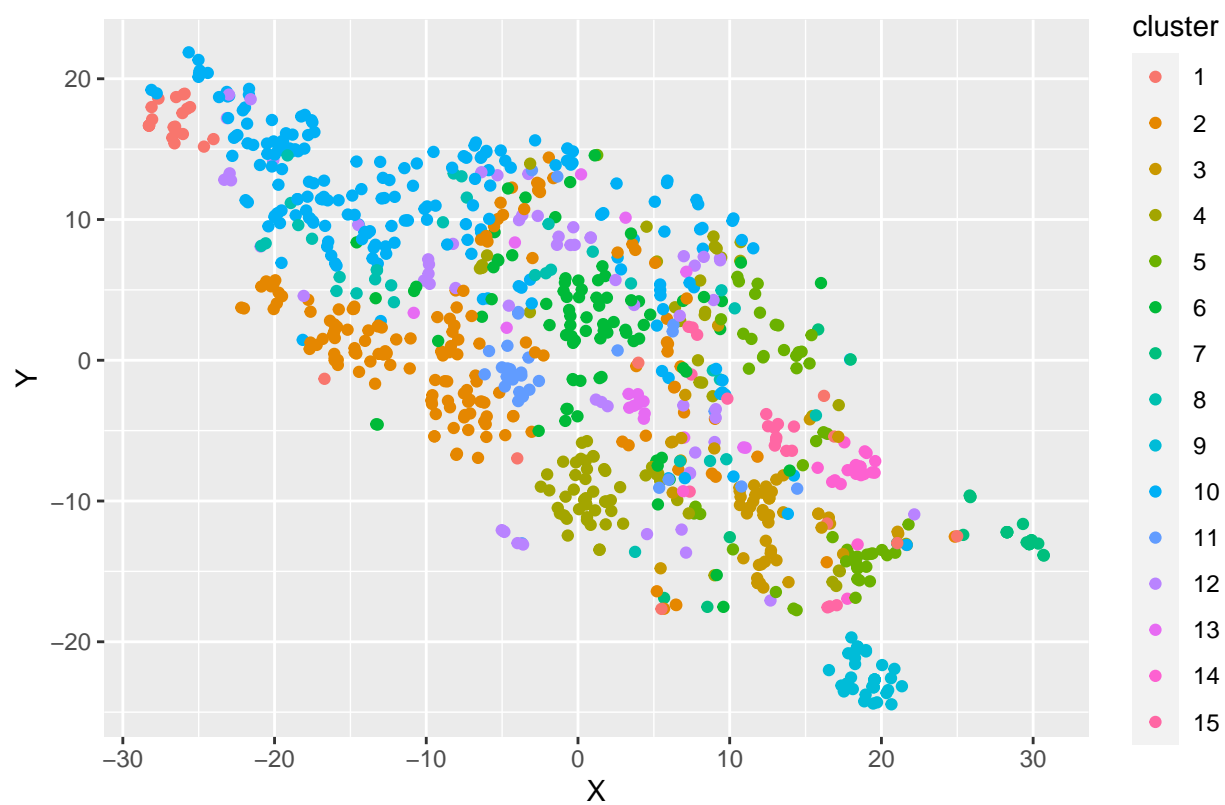


Clustering of All Inmates, Based on Top 5 Tags



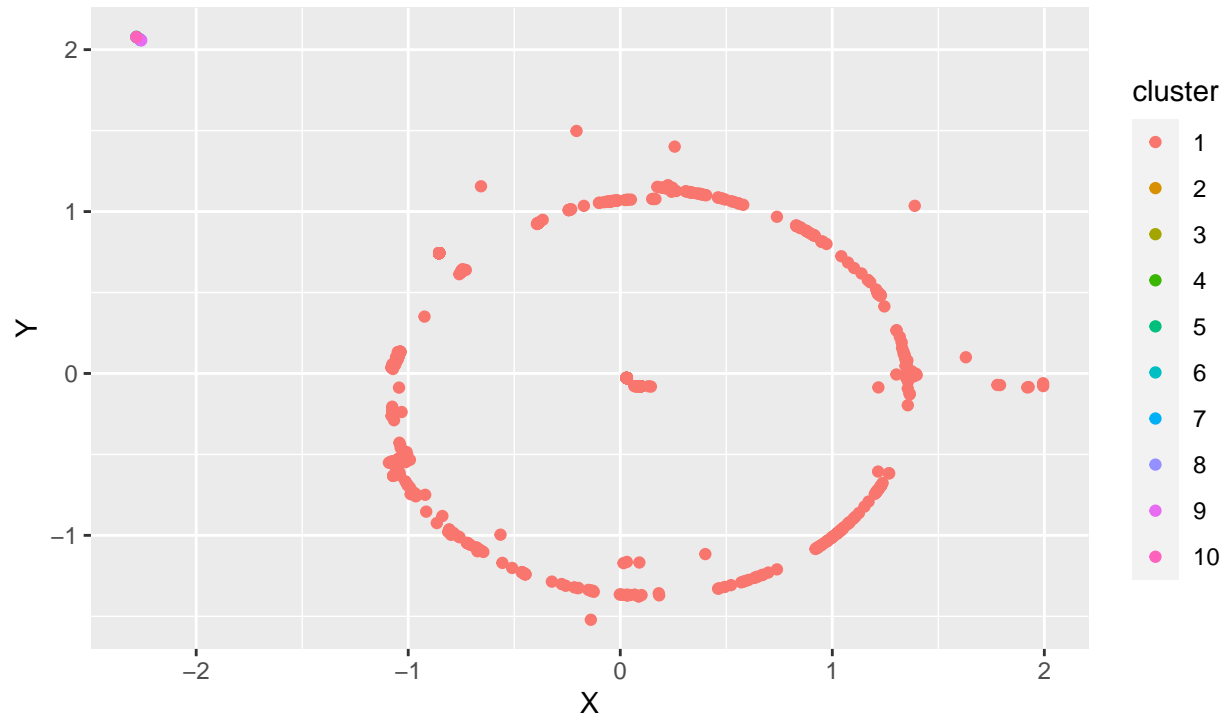
We're also curious, what do inmates in the same cluster have in common?

Because we're interested in the engagement patterns of users, it's probably a good idea for us to account for time as an active user in our clustering.

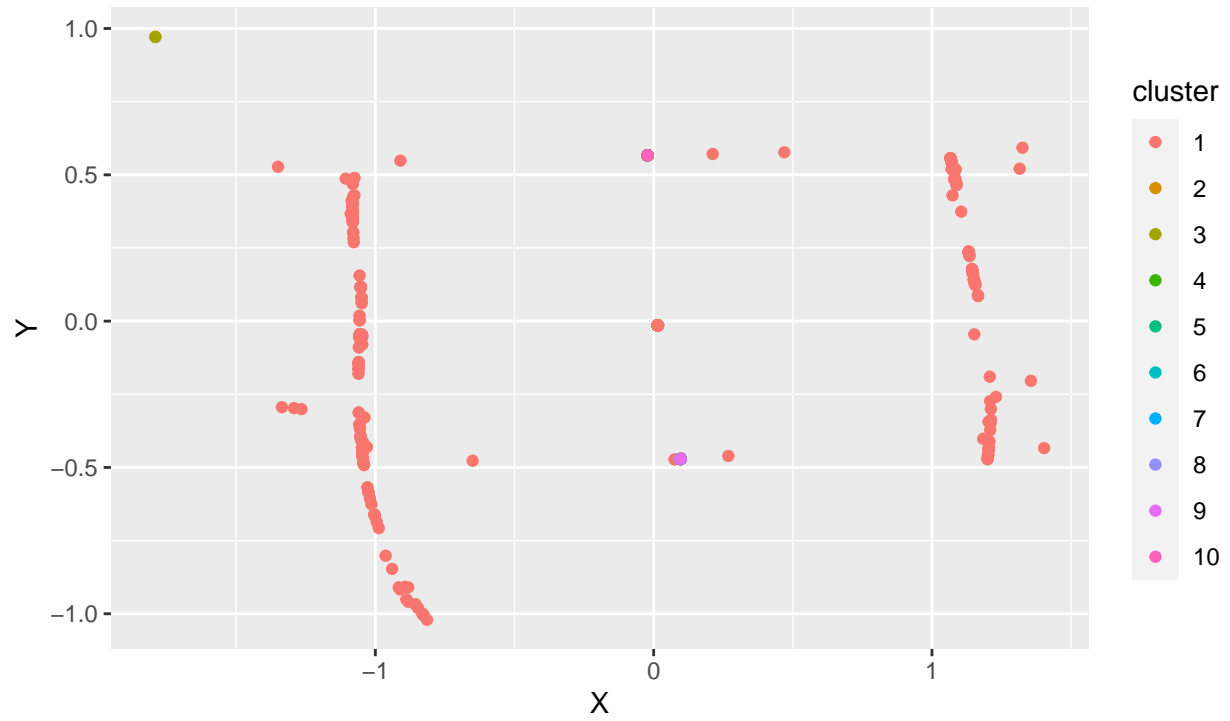


Next, we split the time as an active user up into quantiles and perform the clustering.

Clustering of Users  
Split into Active User Quantiles  
Top 5 Tags in the Top 10 Courses

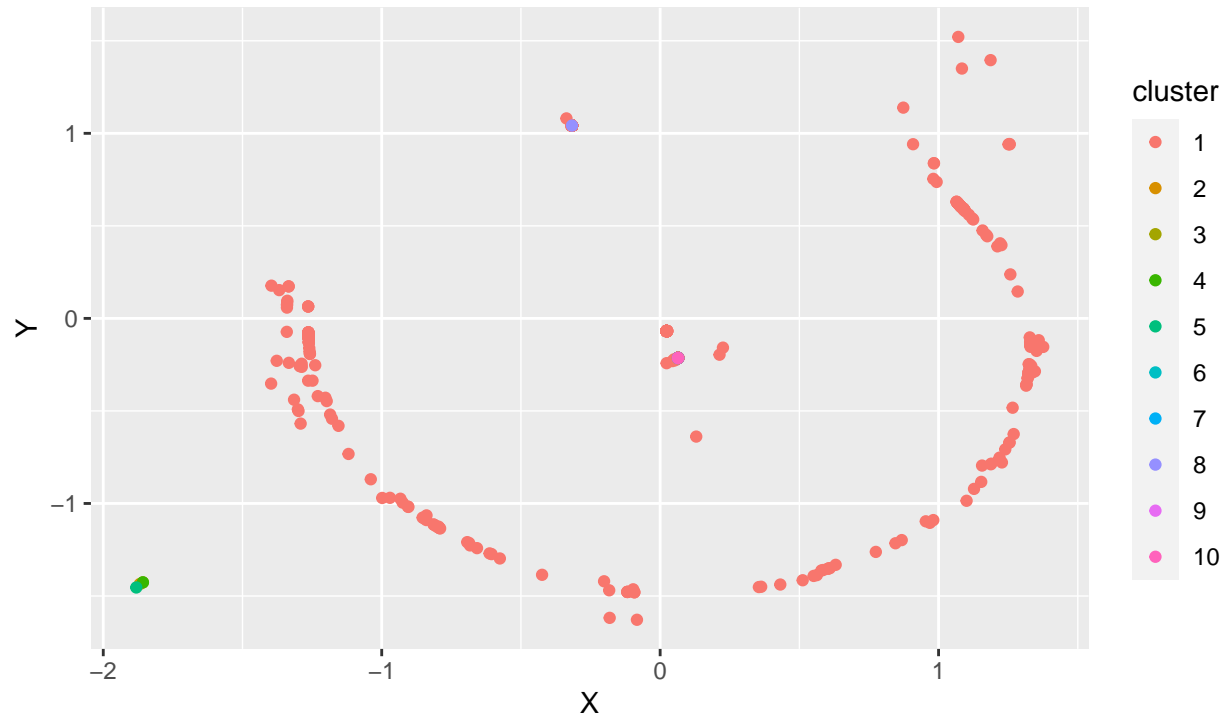


Clustering of Users  
Split into Active User Quantiles  
Top 5 Tags in the Top 10 Courses

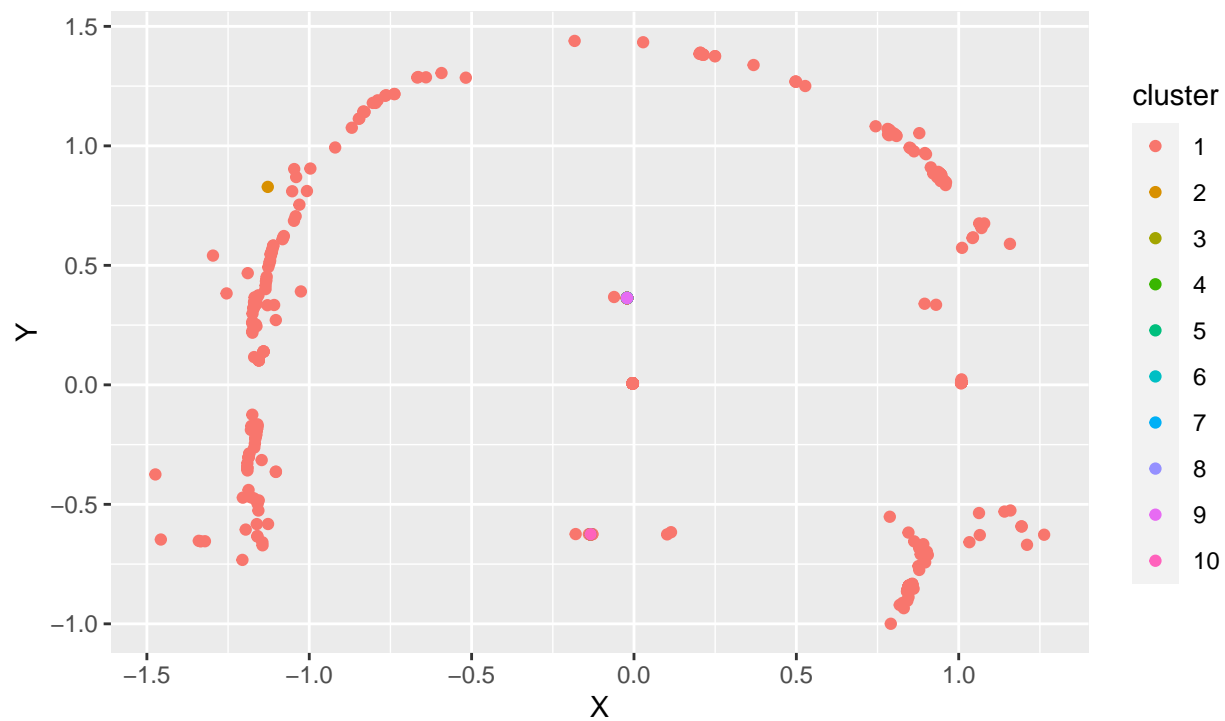




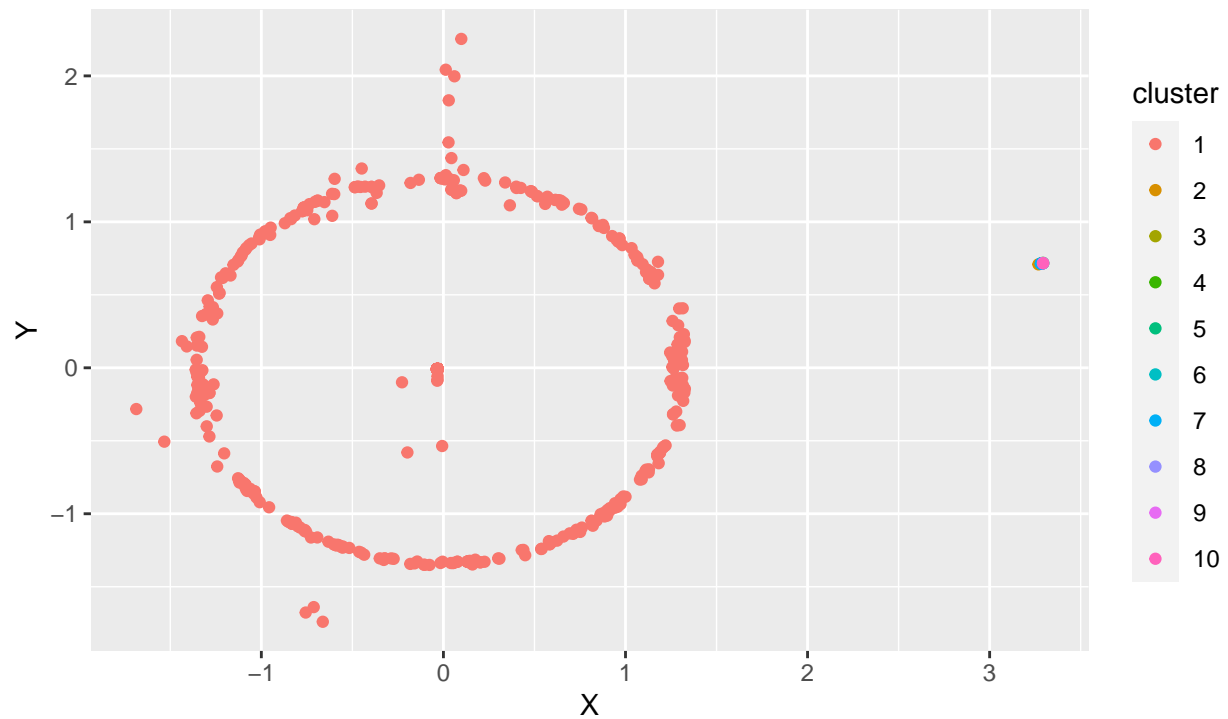
Clustering of Users  
Split into Active User Quantiles  
Top 5 Tags in the Top 10 Courses



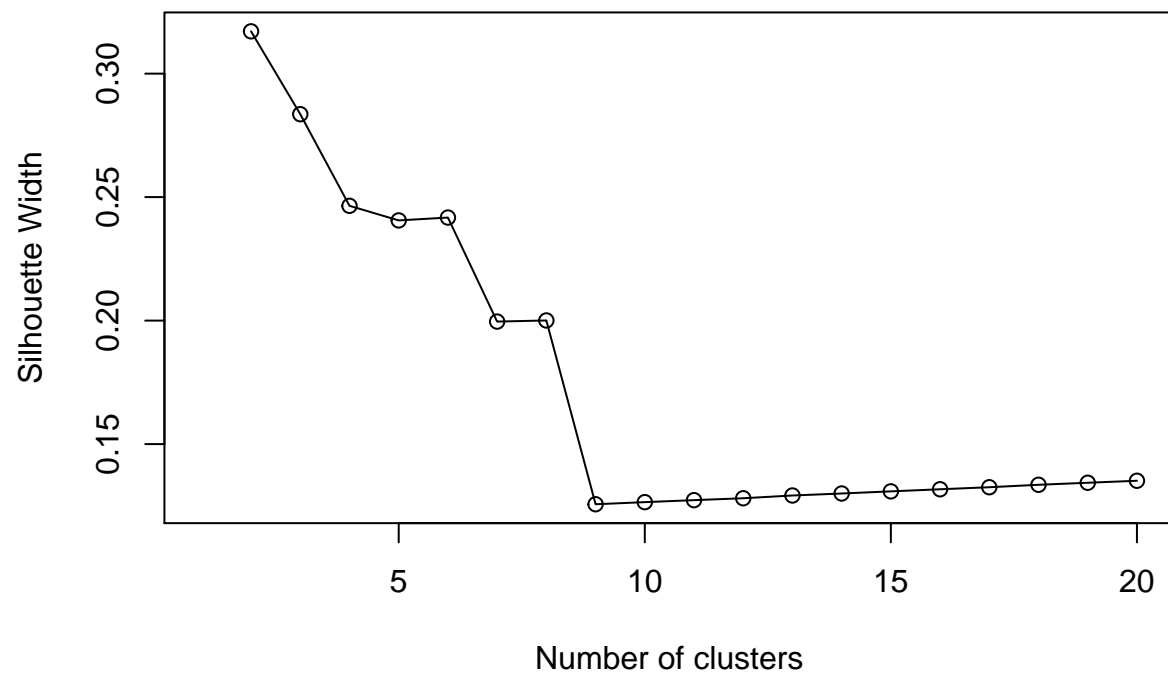
Clustering of Users  
Split into Active User Quantiles  
Top 5 Tags in the Top 10 Courses



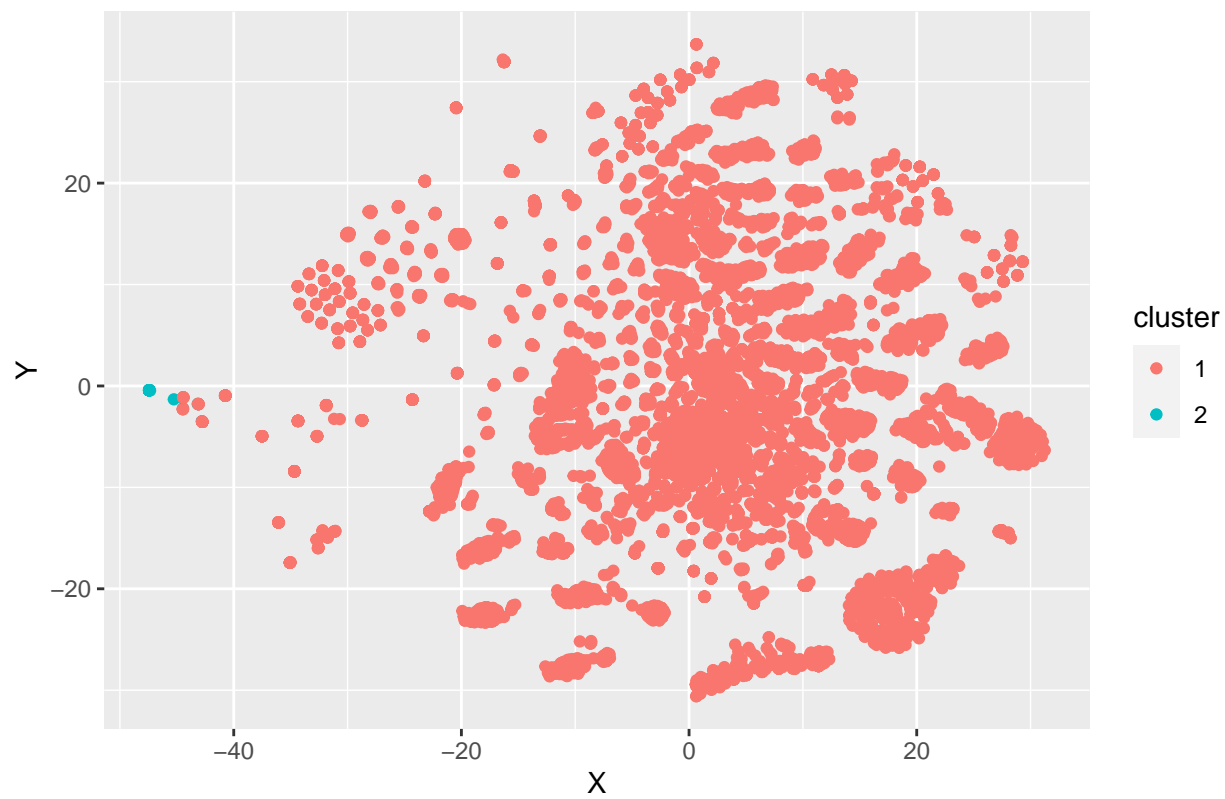
Clustering of Users  
Split into Active User Quantiles  
Top 5 Tags in the Top 10 Courses



**THIS IS A BIG TO-DO:** Incorporating nav-map information. Here we have features like `total_pages` and `total_page_items` but we need to play around with weighting!



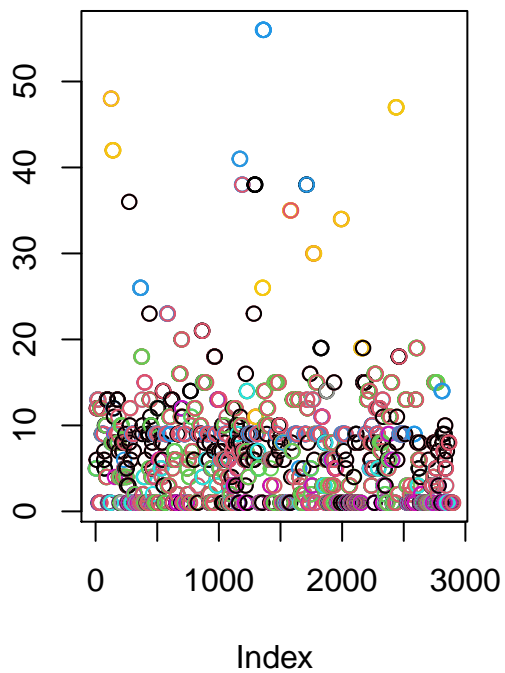
Clustering with Nav Map Info



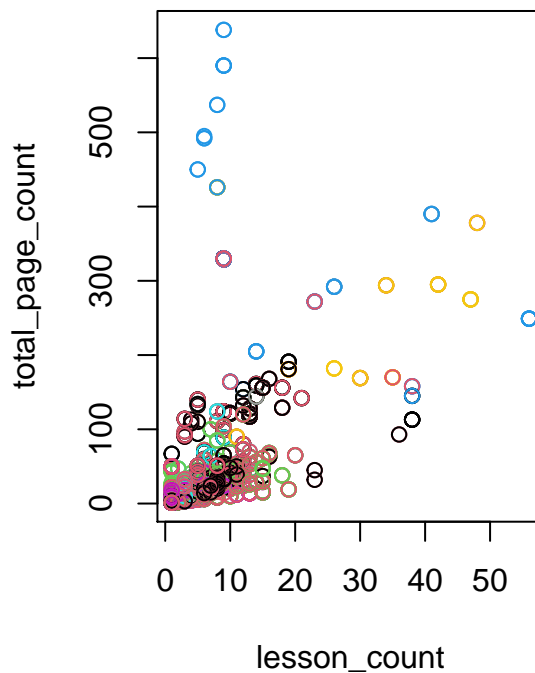
## K Prototypes

data.frame(tags\_with\_parsed\_json2 %>% select\_if(is.numeric

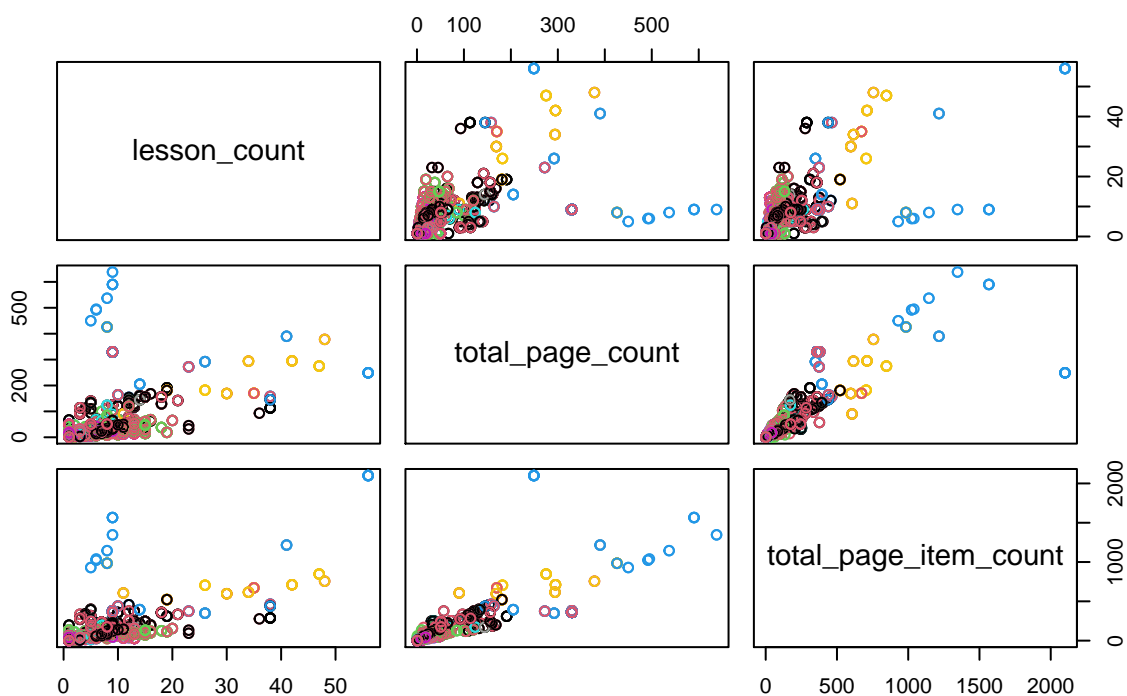
### K-prototypes



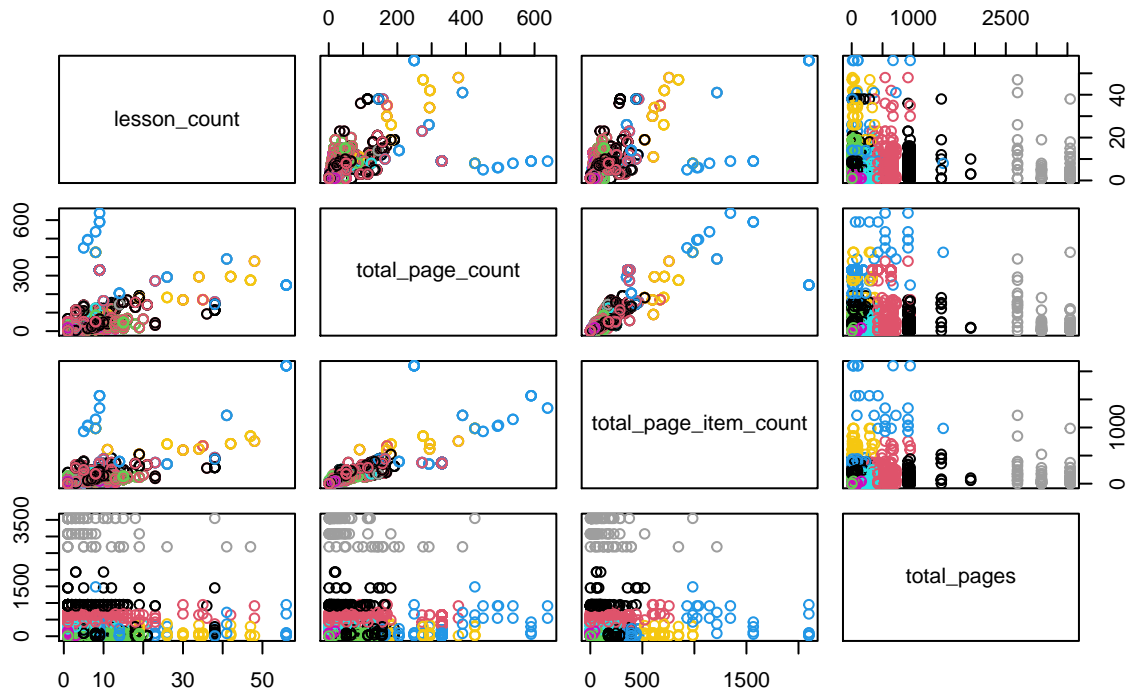
### K-prototypes

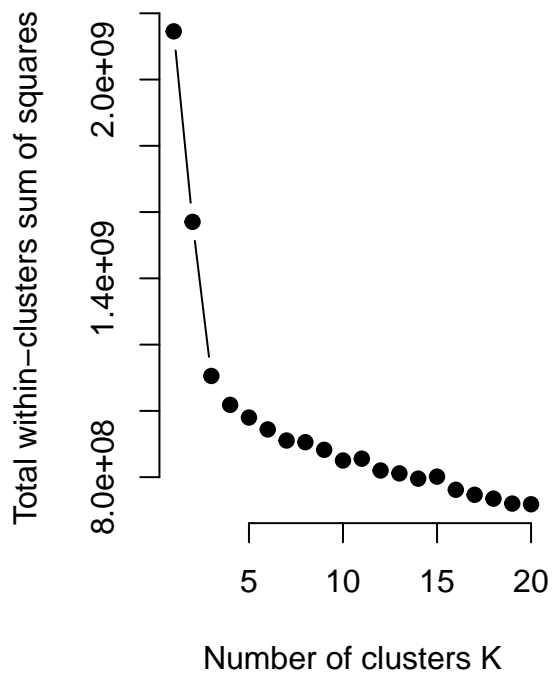


### K-prototypes



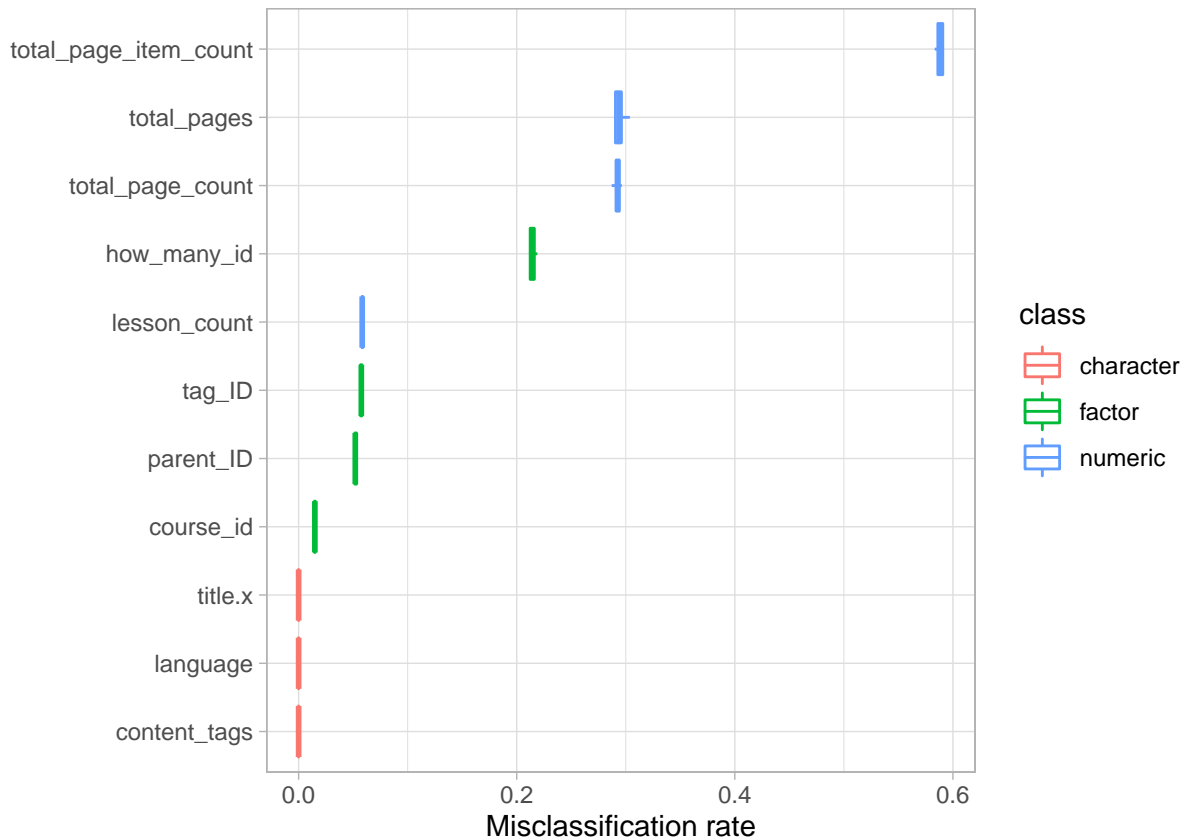
## K-prototypes







What features are most important?

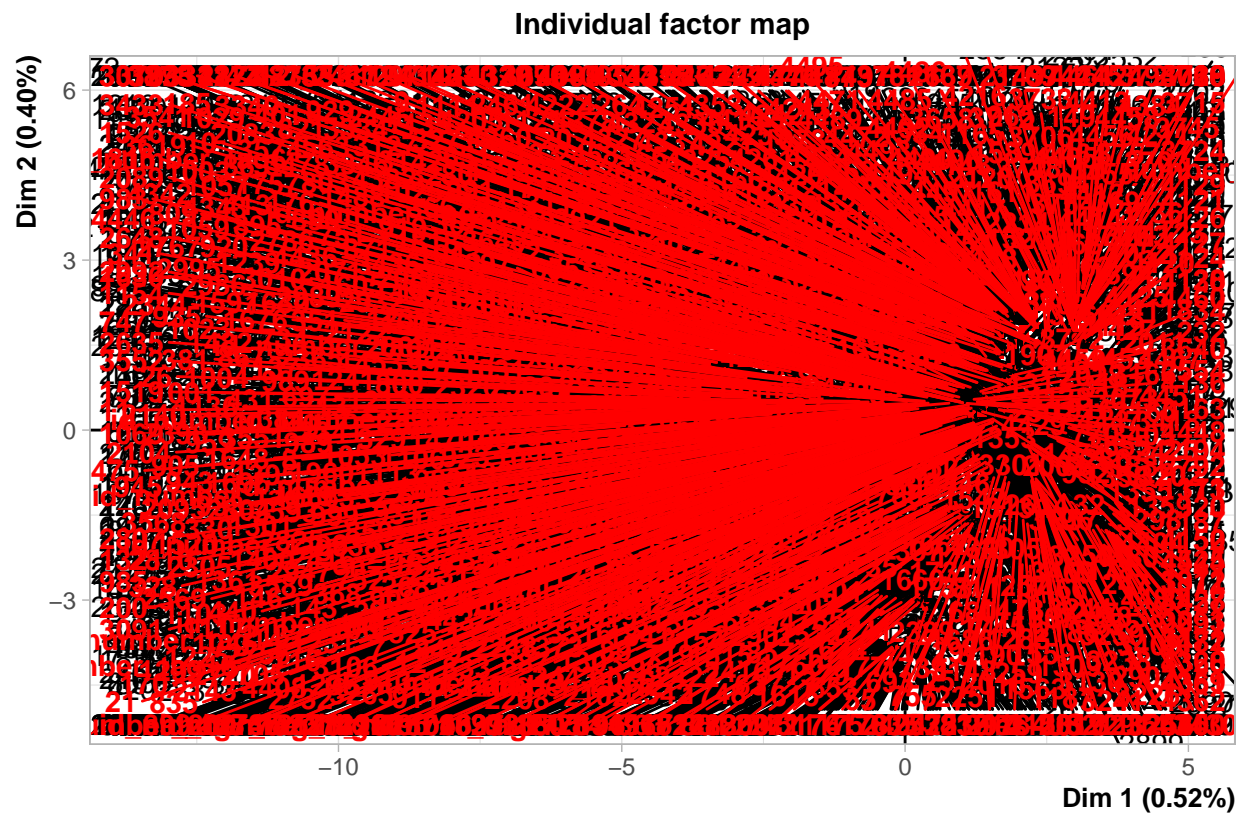


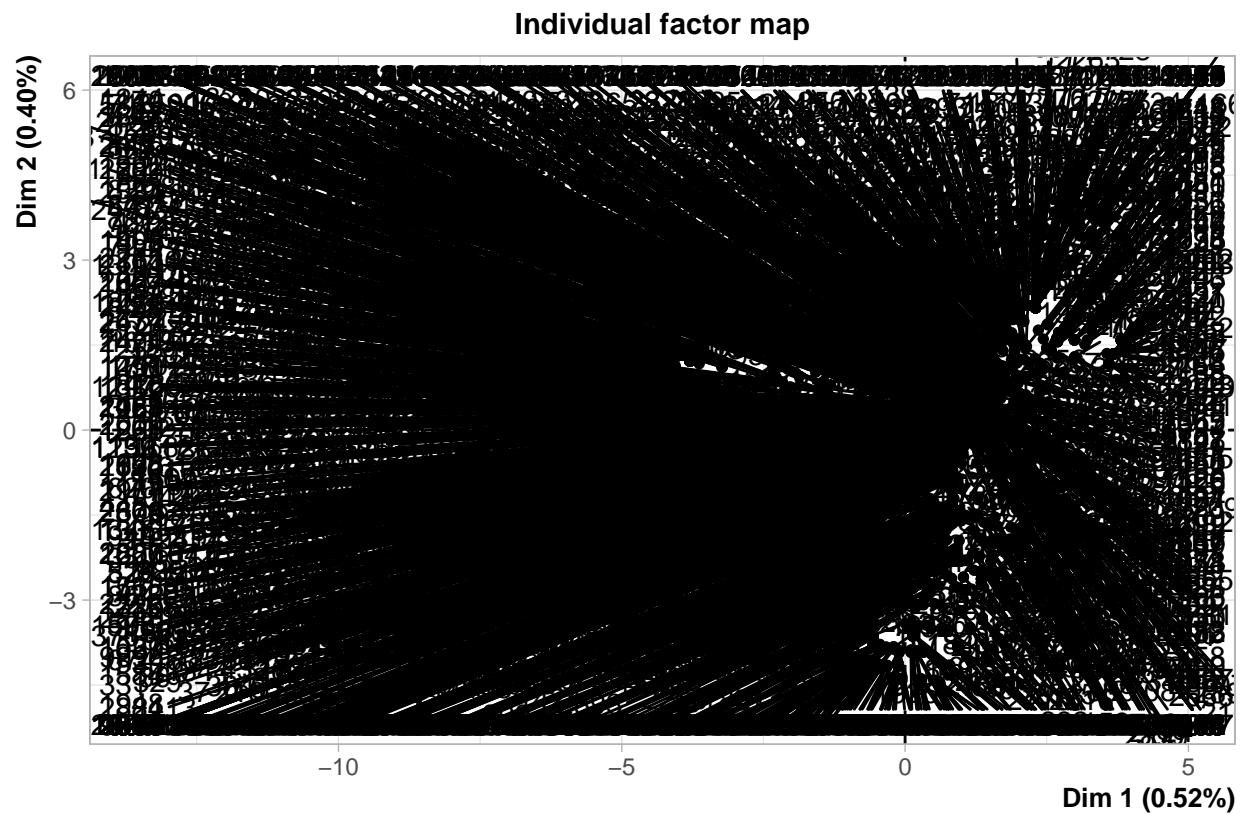
### Dimension Reduction? Or at least PCA/ FAMD

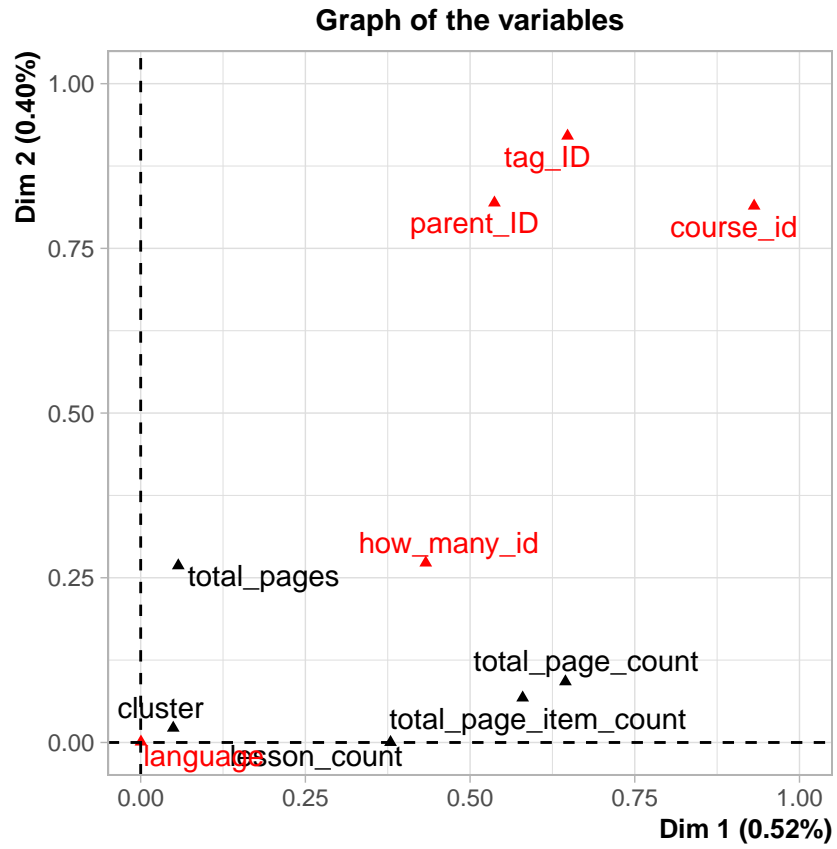
Assuming we represent the data in **wide** format– we may want to look into joint dimension reduction + clustering algorithms. UPDATE: this ay not be a good idea unless we have exclusively numerical features.

We note that the first two components when running these things capture **very little variance**. So how useful is this for anything? :(

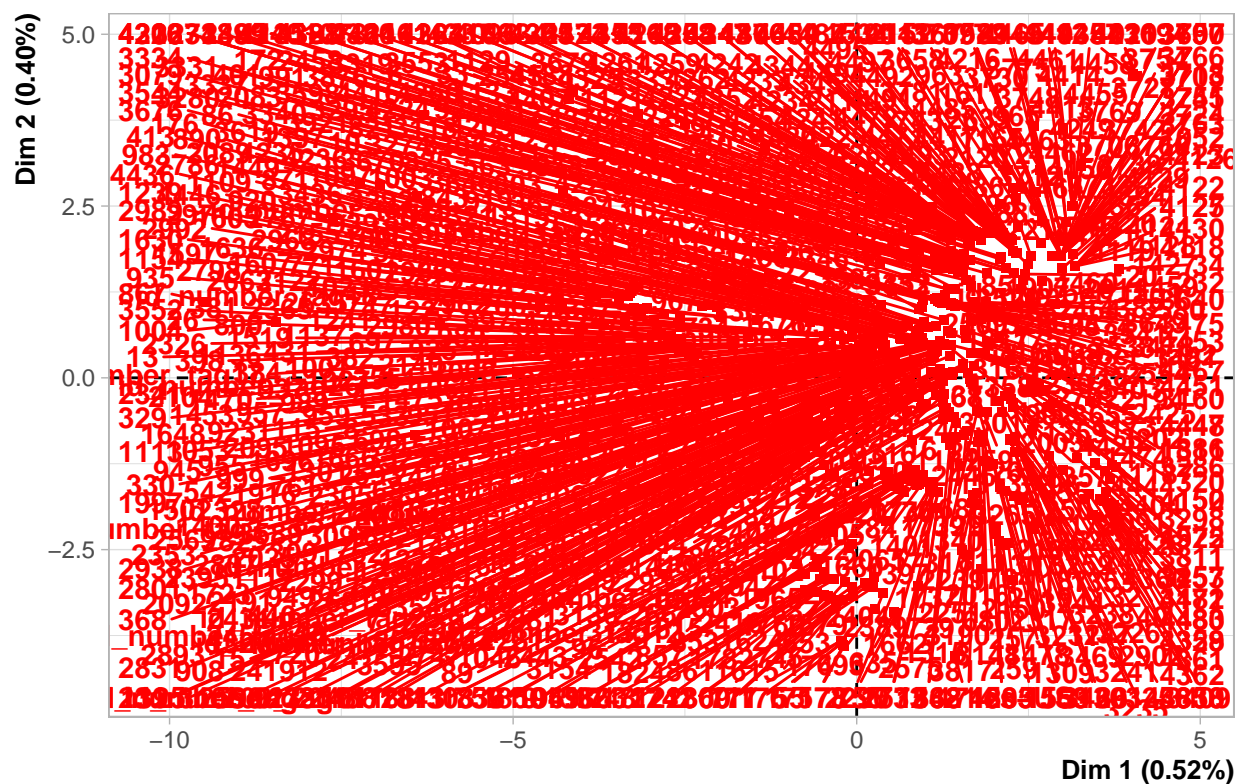
I think this is necessary if we want to make the clustering unit **INMATES**–the data frame would have many columns if we had each tag ID for each course associated with the inmate, so I think this definately has merits.

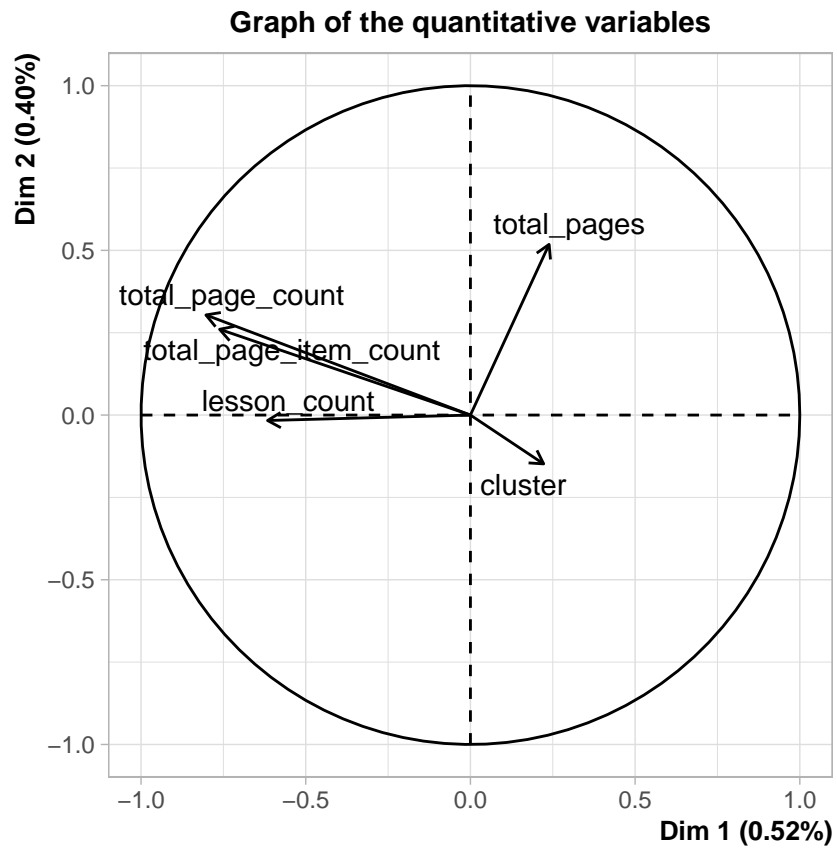


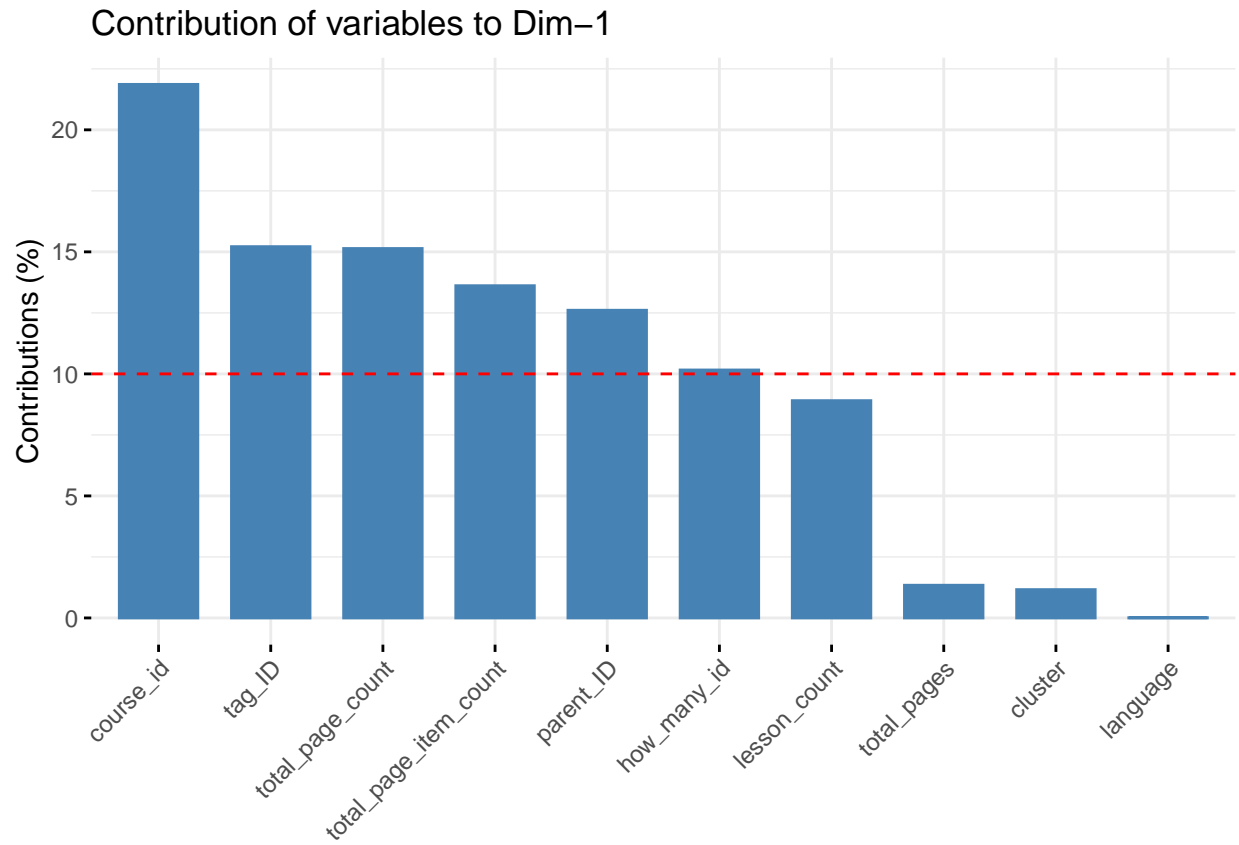


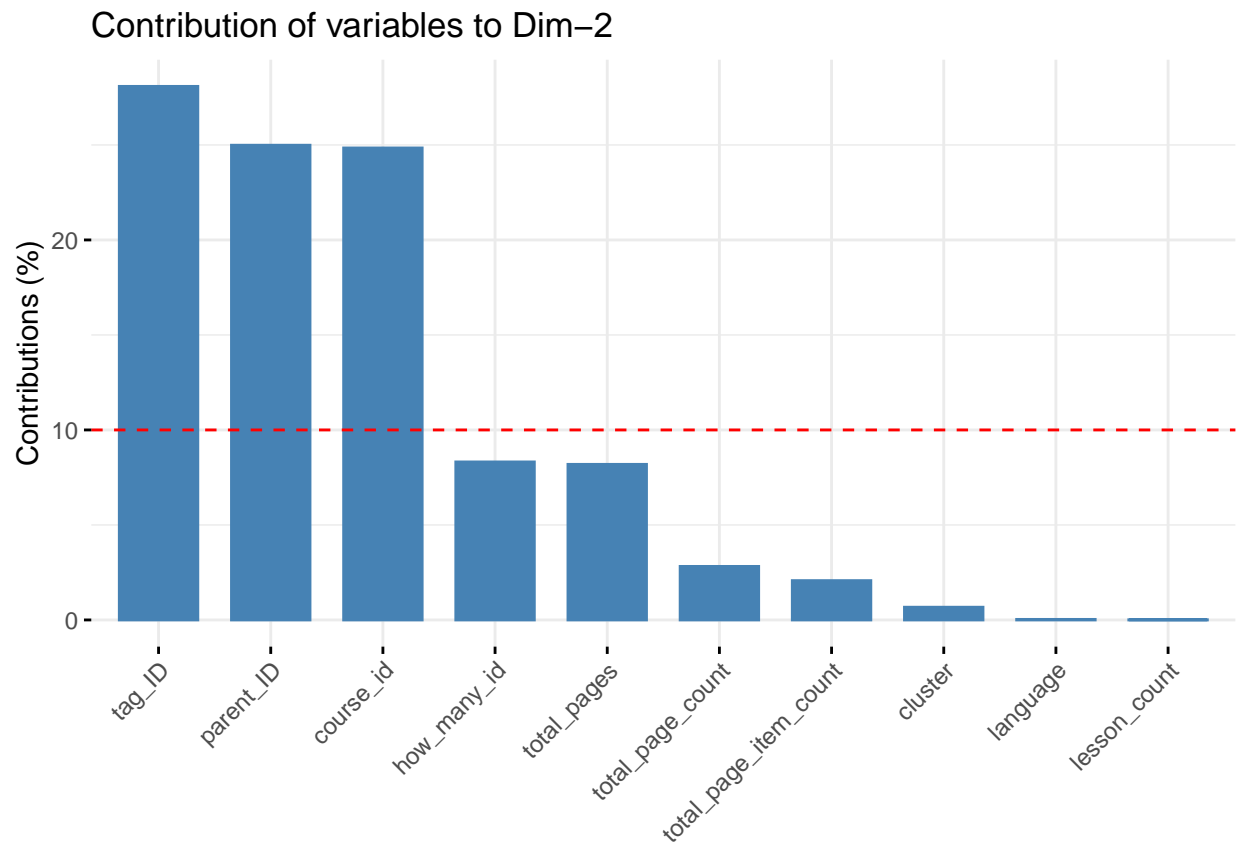


### Graph of the categories









)