# Song Recommender

## MathStat Application

**MathStat**

Professor Tillmann Schwörer

Summer Semester 2020

Alexander P. Kleine::          MatNr. 934944

Robert Weber::          MatNr. 936729

# Overview

- The Task

- Obtaining Data from a Music Test Survey

- Project Vision

- Wrangling… Wrangling… Wrangling

- EDA – Explorative Data Analysis (k-Means Clustering)

- Recommenderlab Package (Recommender System)

- RStudio – Guideline Markdown

# The Chosen Task

**Setup:**

Your team collects a selection of jokes (~ 25), ideally jokes that appeal to different humour preferences. Each course participant must rate a mandatory subset of jokes (~15). Additionally each participant rates a subset of the non-mandatory jokes.

**Possible analytic questions:**

- **Distances**: Measure distances of jokes and/or persons in terms of joke ratings, using alternative distance measures

- **Recommendations**: Make an individualised joke recommendation.
  Recommend one of the bonus jokes that the closest neighbour liked

- **Clustering**: Use k-means or hierarchical clustering to identify homogenous clusters of jokes or persons **Principle component analysis/Singular value decomposition**: Find a lower dimensional repre- sentation (latent factors) of persons and/or jokes.

**Songs, movies, or books analysis** The same setup and the analytic questions as above can be addressed using an alternative domain, e.g.: Songs, Movies, Books

Note: It won't be feasible to get a full set of ratings for 15 movies or books from each participant. Hence, you will have a *sparse* rating matrix. This makes some of the analytic questions more difficult or unfeasable (notably clustering and PCA/SVD). However, it is possible to measure distances for each pair of movie or book raters, based on the subset of movies/books that both have rated. And thus, to make personalized recommendations.

# Obtaining Data from a CATI/CAWI Music Test

Choosing the Domain: **Songs**

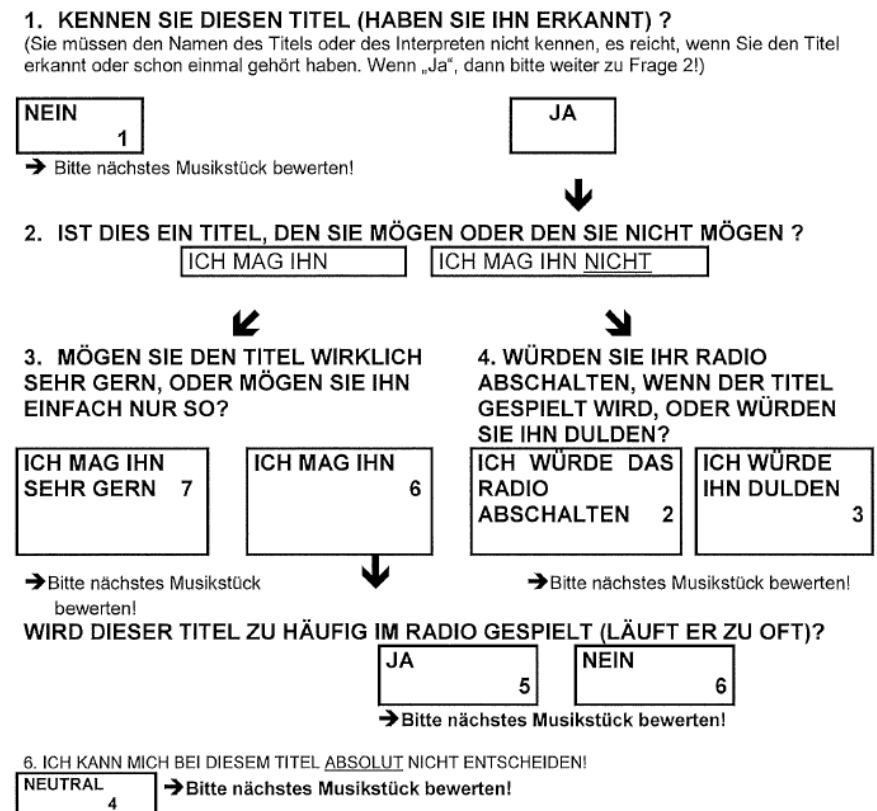Using **current Music Test Data** from a music research unit

As this data is used to make playlist recommendations or rather **to curate playlists** that are supposed to generate a mixture of maximum number of positive responses while trying to keep negs to a minimum…

We wanted to use this dataset to create a recommender system

# Obtaining Data from a CATI/CAWI Music Test

- Data generated from CATI/CAWI Music Tests

- 120 test persons / 50 song ratings each
  (resulting in 6,000 entries)

- We do have further data on the tested
  person (ZAG, listening type, main station)

- The response: „I do not know this song"
  is present 495 times (8.25 %)



1. **KENNEN SIE DIESEN TITEL (HABEN SIE IHN ERKANNT) ?**
(Sie müssen den Namen des Titels oder des Interpreten nicht kennen, es reicht, wenn Sie den Titel erkannt oder schon einmal gehört haben. Wenn „Ja", dann bitte weiter zu Frage 2!)

NEIN 1 → Bitte nächstes Musikstück bewerten!     JA ↓

2. **IST DIES EIN TITEL, DEN SIE MÖGEN ODER DEN SIE NICHT MÖGEN ?**
ICH MAG IHN     ICH MAG IHN NICHT

3. **MÖGEN SIE DEN TITEL WIRKLICH SEHR GERN, ODER MÖGEN SIE IHN EINFACH NUR SO?**
ICH MAG IHN SEHR GERN 7     ICH MAG IHN 6
→ Bitte nächstes Musikstück bewerten!

4. **WÜRDEN SIE IHR RADIO ABSCHALTEN, WENN DER TITEL GESPIELT WIRD, ODER WÜRDEN SIE IHN DULDEN?**
ICH WÜRDE DAS RADIO ABSCHALTEN 2     ICH WÜRDE IHN DULDEN 3
→ Bitte nächstes Musikstück bewerten!

**WIRD DIESER TITEL ZU HÄUFIG IM RADIO GESPIELT (LÄUFT ER ZU OFT)?**
JA 5     NEIN 6
→ Bitte nächstes Musikstück bewerten!

6. **ICH KANN MICH BEI DIESEM TITEL ABSOLUT NICHT ENTSCHEIDEN!**
NEUTRAL 4 → Bitte nächstes Musikstück bewerten!

# Obtaining Data from a CATI/CAWI Music Test



Data about the tested person

Data about the tested songs

A subset of raw data before wrangling

# Decision concerning the Scope of the Project

Most Common Types of Recommender Systems:

- **Collaborative Filtering Methods**: as in neighbourhood models →

  a ratings matrix includes dependencies between individual items

- **Content-Based Methods:** User's interest can be modeled on the basis of properties

  (or attributes) of the items they have rated in the past.

- **Knowledge-Based Methods**: users interactively specify their interests after which this is

  combined with domain knowledge to provide recommendations

While our data is used in a hybridized manner in practise,

  **we are going to focus on a „Collaborative Filtering Method".**

# Decision concerning the Scope of the Project

- Even though we do have (social) data about the tested person, we are going to assume, that we don't.

- It means that we intentionally take out information and context, which will in turn leave us with a more straightforward approach.

- This takes scope out of this project, while also having practical relevance, since we cannot neccessarily assume to have this information about the audience, we are going to make recommendations to.

- **Collaborative Filtering** refers to the use of ratings from multiple users **to predict missing ratings**.

Attribute information about the users → „content-based rec."

Data about user-item interactions → „collaborative filtering"

A subset of raw data before wrangling

# Wrangling… wrangling… wrangling

- After wrangling, we are left with 50 variables (Songs) and 120 observations (Scores) of these variables.

- The two possible Scores that we have to give special consideration are:  „1" and „4"

- 1: is translated as „I do not know this song" (missing value?); there are 494 of these in the data

- 4: is translated as „I can't say." (neutral), 292 times present

- There seems to be a group of people that either are ignorant of the tested music,  don't like it or can't make up their mind.

# Wrangling... wrangling... wrangling

- After wrangling, we are left with 50 variables (Songs) and 120 observations (Scores of tested people) of these variables.

- For the next steps it is assumed that the market for attention to popular music is efficient, i.e. that the score „1" can be interpreted as a maximum disinterest in that particular song.

- Scores of 4 then represent true indifference and we have a true scale of 1:7.

- First glances at the test scores let us to postulate at least three patterns: Very positive, very negative and mixed scores.

- The question, whether the first two are essentially similar and just shifted patterns will have to be adressed later.

- **Let's have a look at a particular song's test pattern  (the one with the highest mean score) on the next slide before going on to distances**
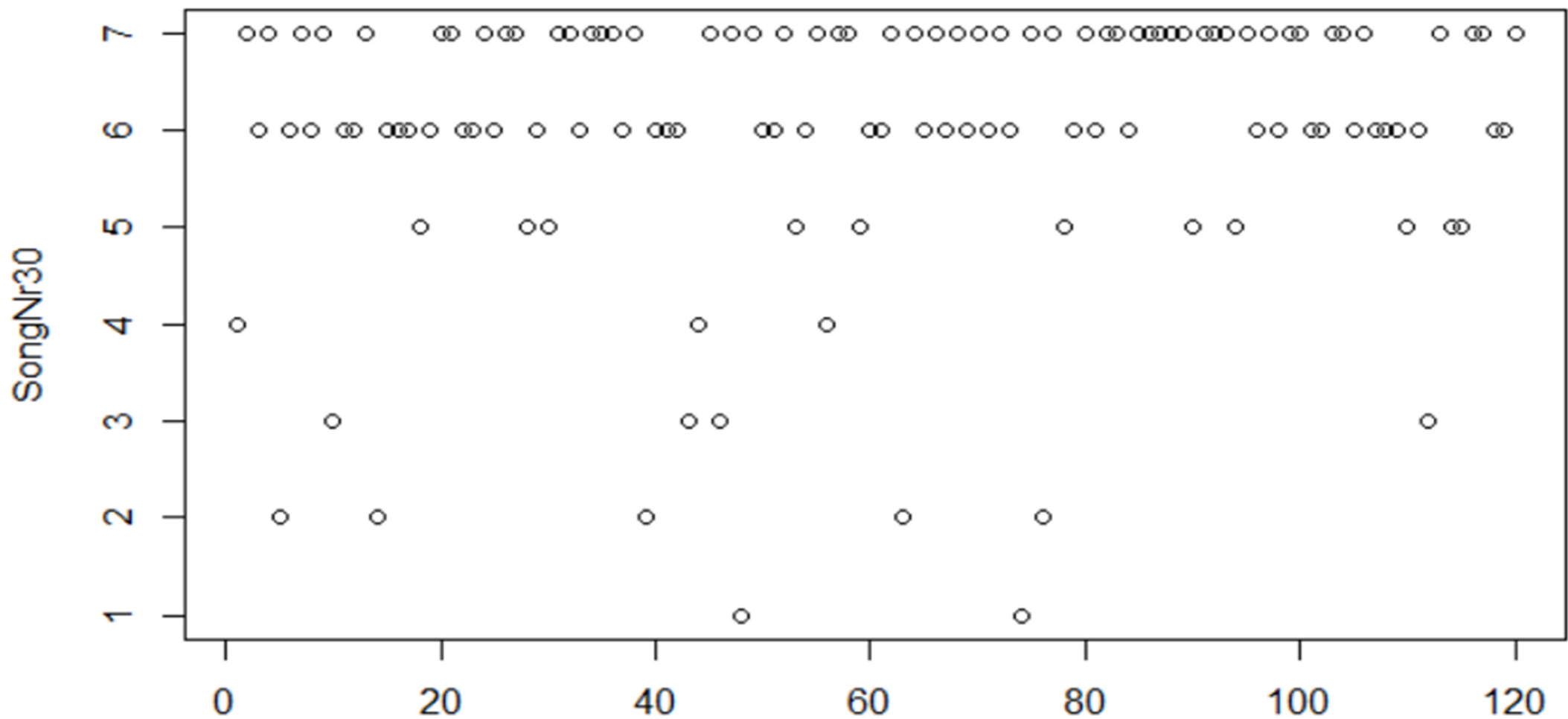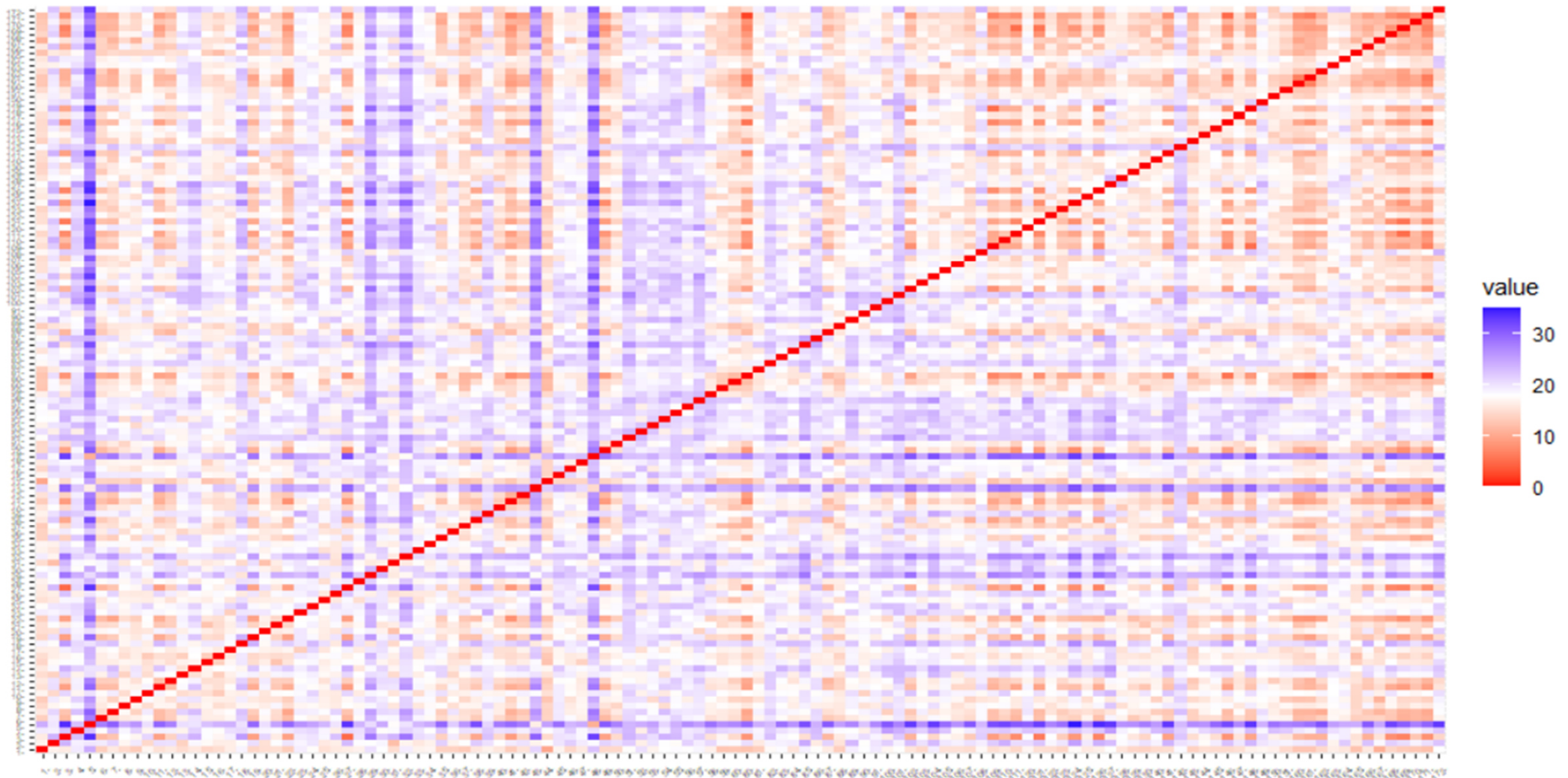


User #3 „loves" music in general? (7)

User #5 doesn't know most of the songs... or can't make up his/her mind about music (1, 4)

# Visual:  (120) Scores of SongNr_30



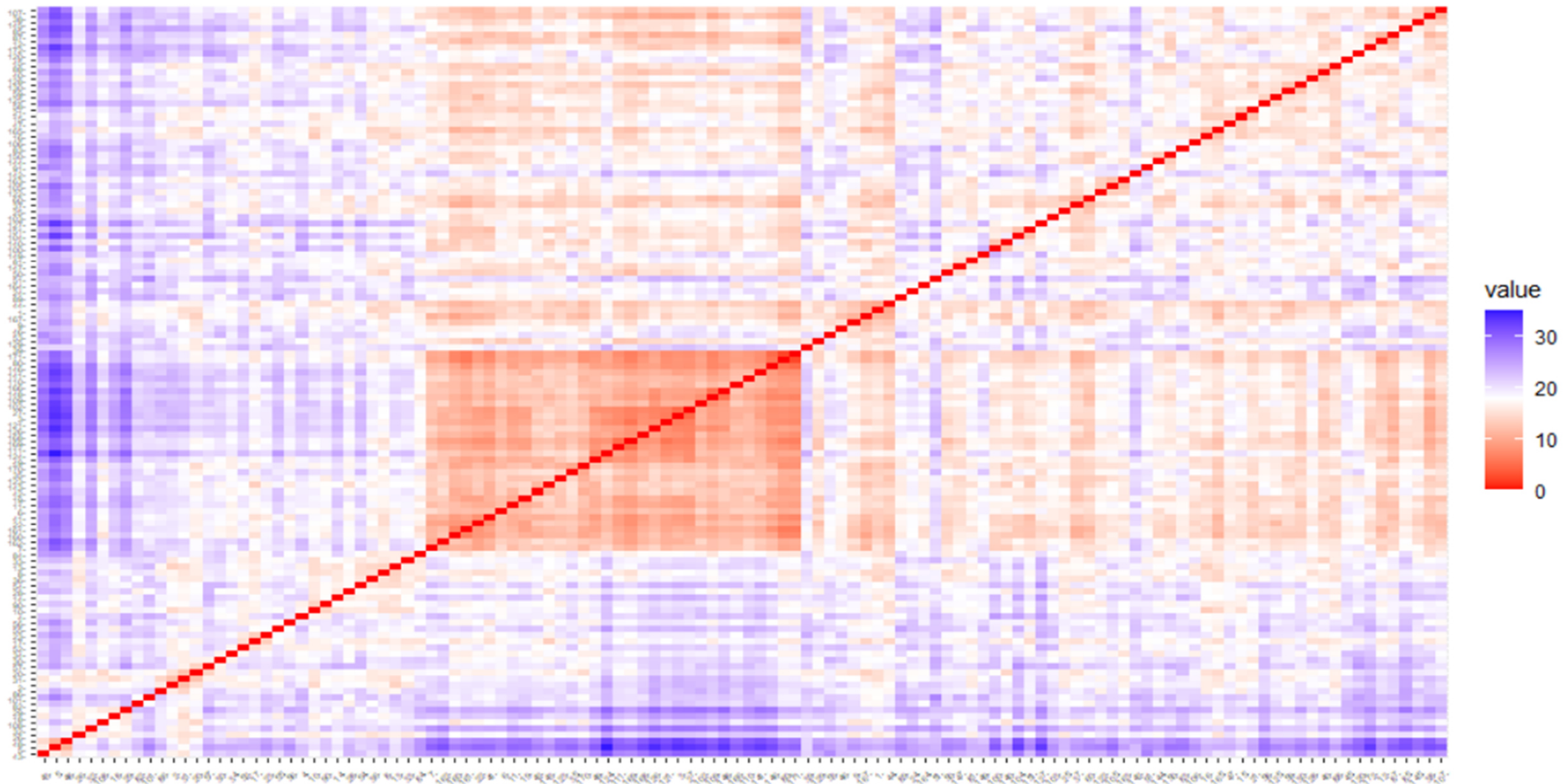(Dynoro, In My Mind: https://open.spotify.com/track/sDtfiOiYSSu01r78p6bZLw)

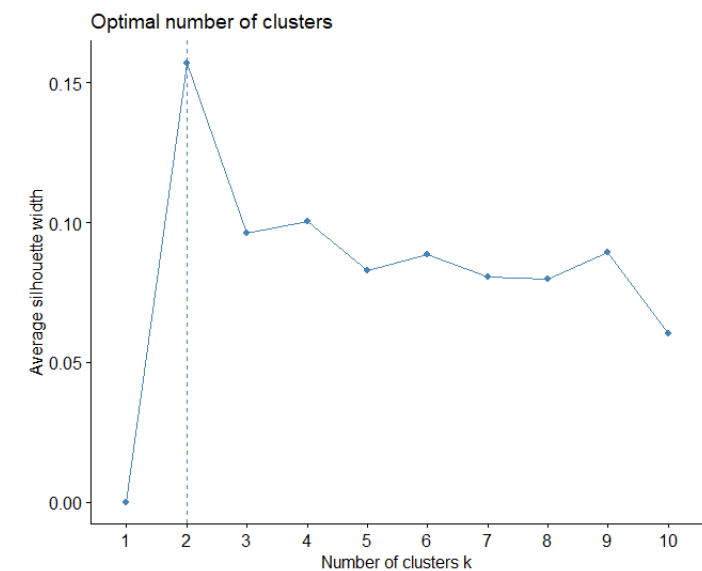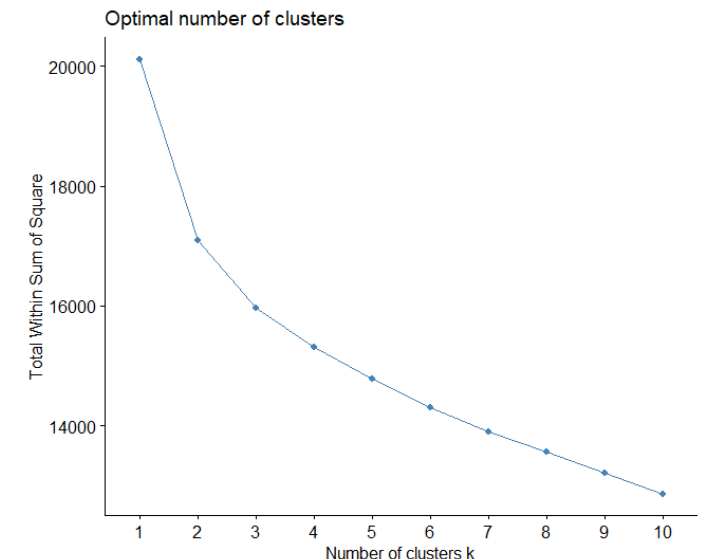# Distance Visualisation (fviz_dist, „Euclidean", unordered) (invariant to scaling)

# Ordered Dissimilarity Image (ODI) (not completely invariant to scaling)

# Clustering Approaches

- We do not have to standardize as all data is on the same scale

- k-Means Clustering uses Euclidean distances, which should give a first idea about clustering possibilities

- Using the Elbow Approach (via fviz_nbclust, method = „wss"), the optimal number of clusters could be 3. This would split the data into three clusters of sizes 25, 42, 53.

- Using the „Silhouette Method" (againg via fviz_nbclust), the optimal number of clusters should be 2. This splits the data into two groups of almost equal size (58, 62).

# Overview

- The Task

- Obtaining Data from a Music Test Survey

- Project Vision

- Wrangling… Wrangling… Wrangling

- EDA – Explorative Data Analysis (k-Means Clustering)

- **Recommenderlab Package (Recommender System)**

- **RStudio – Guideline Markdown**