

Clustering - Badisches Landesmuseum

Arne Nommensen, Mara Lindhorn, Martin Schanze, Nick Jürgensen, Ninik Fatmawati-Lottes, Onur Zere

26.01.2022

1. Import der Daten & Libraries

Libraries laden

Hier werden die benötigten Libraries geladen

```
rm(list = ls())

library(dplyr)
library(psych)
library(factoextra)
library(deepANN)
library(kohonen)
library(tidyverse)
library(readxl)
library(readr)
library(ggplot2)
library(reshape)
```

Datensatz laden

Der Datensatz stammt aus der Umfrage "Creative User Empowerment" (Ende 2021) des Badischen Landesmuseums.

```
rawdata <- read.csv("umfrageonline-2728300.csv", header = TRUE, sep = ";", encoding="UTF-8")
```

2. Pre-Processing der Daten

Die Daten müssen für das Aufstellen des Modells vorverarbeitet werden.

Relevante Daten extrahieren

Es werden aus der CSV Datei einige Spalten nicht übernommen. Bei den nicht übernommenen Spalten handelt es sich zum Beispiel um die Frage "Ich möchte noch folgendes sagen", da Freitext von diesem Modell nicht aufgegriffen werden kann.

```
datanew <- rawdata[c(6L:9L, 11L, 13L:14L, 16L:22L, 27L, 30L:41L, 44L:46L, 49L:56L, 59L:67L, 70L:75L, 79L:84L, 87L:98L, 102L:105L, 107L:110L, 112L:114L) ]  
headers <- colnames(datanew)
```

Funktion zum Anpassen der Überschriften

Die Überschriften der übernommenen Spalten werden angepasst, um die Weiterverarbeitung zu erleichtern. Die Bedeutung der einzelnen Überschriften können Sie der Umfrage entnehmen, die Reihenfolge der Spalten wurde nicht verändert.

```
ueberschriften_setzen <- function(datanew) {  
names(datanew)[1] <- "Datenschutz"  
names(datanew)[2] <- "Teilnahme"  
names(datanew)[3] <- "Geschlecht"  
names(datanew)[4] <- "Alter"  
names(datanew)[5] <- "PLZ"  
names(datanew)[6] <- "Bildungsabschluss"  
names(datanew)[7] <- "Anzahl_Besuche"  
names(datanew)[8] <- "Beurteilung_Sammlung"  
names(datanew)[9] <- "Beurteilung_Sonder"  
names(datanew)[10] <- "Beurteilung_Bildungsangebot"  
names(datanew)[11] <- "Beurteilung_Expo"  
names(datanew)[12] <- "Beurteilung_Online"  
names(datanew)[13] <- "Beurteilung_Apps"  
names(datanew)[14] <- "Beurteilung_Erreichbarkeit"  
names(datanew)[15] <- "Empfehlung"  
names(datanew)[16] <- "Interesse_Archaeologie"  
names(datanew)[17] <- "Interesse_Weltkultur"  
names(datanew)[18] <- "Interesse_Mittelalter"  
names(datanew)[19] <- "Interesse_Design"  
names(datanew)[20] <- "Interesse_Baden"  
names(datanew)[21] <- "Interesse_1900"  
names(datanew)[22] <- "Interesse_Aktuell"  
names(datanew)[23] <- "Interesse_Musik"  
names(datanew)[24] <- "Interesse_Materiell"  
names(datanew)[25] <- "Interesse_Immateriell"  
names(datanew)[26] <- "Interesse_HdK"  
names(datanew)[27] <- "Interesse_3D"  
names(datanew)[28] <- "Praesenz"  
names(datanew)[29] <- "Gesellschaft"  
names(datanew)[30] <- "Digitale Angebote"  
names(datanew)[31] <- "Digital_Fuehrung"  
names(datanew)[32] <- "Digital_Workshops"  
names(datanew)[33] <- "Digital_Ausstellungen"  
names(datanew)[34] <- "Digital_Podcasts"  
names(datanew)[35] <- "Digital_Kataloge"  
names(datanew)[36] <- "Digital_Kurse"  
names(datanew)[37] <- "Digital_Spiele"  
names(datanew)[38] <- "Digital_SocialMedia"  
names(datanew)[39] <- "Grunddigital_Freude"  
names(datanew)[40] <- "Grunddigital_Unterhaltung"  
names(datanew)[41] <- "Grunddigital_Erkenntnisse"  
names(datanew)[42] <- "Grunddigital_Faehigkeiten"  
names(datanew)[43] <- "Grunddigital_Schule"  
names(datanew)[44] <- "Grunddigital_Spezial"  
names(datanew)[45] <- "Grunddigital_Forschung"  
names(datanew)[46] <- "Grunddigital_Beruf"  
names(datanew)[47] <- "Grunddigital_Weltweit"  
names(datanew)[48] <- "Anforderungen_Praxis"  
names(datanew)[49] <- "Anforderungen_Ausstellungen"  
names(datanew)[50] <- "Anforderungen_gleicheInhalte"
```

```
names(datanew)[51] <- "Anforderungen_Bildungsmaterial"  
names(datanew)[52] <- "Anforderungen_MehrInhalte"  
names(datanew)[53] <- "Anforderungen_NeueErlebnisse"  
names(datanew)[54] <- "Erlebnis_Video"  
names(datanew)[55] <- "Erlebnis_Lesen"  
names(datanew)[56] <- "Erlebnis_Zuhoeren"  
names(datanew)[57] <- "Erlebnis_Entdecken"  
names(datanew)[58] <- "Erlebnis_Interaktion"  
names(datanew)[59] <- "Erlebnis_Spielerisch"  
names(datanew)[60] <- "KI_Uebersetzung"  
names(datanew)[61] <- "KI_Vertiefungsinfos"  
names(datanew)[62] <- "KI_indv_Empfehlung"  
names(datanew)[63] <- "KI_Texterstellung"  
names(datanew)[64] <- "KI_Bildererkennung"  
names(datanew)[65] <- "KI_Spracherkennung"  
names(datanew)[66] <- "KI_Chatbot"  
names(datanew)[67] <- "KI_generierte_Kunst"  
names(datanew)[68] <- "KI_Emotionserkennung"  
names(datanew)[69] <- "KI_Geschichten_generieren"  
names(datanew)[70] <- "KI_Zusammenhaenge_sichtbar_machen"  
names(datanew)[71] <- "KI_neue_kreative_Prozesse"  
names(datanew)[72] <- "KIASPEKTE_Verstehen"  
names(datanew)[73] <- "KIASPEKTE_Mitgestalten_koennen"  
names(datanew)[74] <- "KIASPEKTE_NeueGeschichten"  
names(datanew)[75] <- "KIASPEKTE_Barrierefreiheit"  
names(datanew)[76] <- "Helfen_Untertitel_und_AlternativeMedien"  
names(datanew)[77] <- "Helfen_Alternative_fuer_nicht_Textinhalte"  
names(datanew)[78] <- "Helfen_leichte_Sprache"  
names(datanew)[79] <- "Helfen_autom_Uebersetzungen"  
names(datanew)[80] <- "digitale_Affinitaet"  
names(datanew)[81] <- "Freizeit_Std"  
names(datanew)[82] <- "Besonders_wichtig"  
return (datanew)  
}  
datanew <- ueberschriften_setzen(datanew)
```

Missing Values entfernen

Datenschutz / Teilnahme nicht zugestimmt - Datensätze entfernt

Die Personen die dem Datenschutz oder der Teilnahme nicht zugestimmt haben werden aus dem Dataframe entfernt. Die Werte in den beiden Spalten werden anschließend auf Null gesetzt, damit diese im nächsten Schritt entfernt werden können.

```
datanew <- subset(datanew, datanew$Datenschutz=="ja")  
datanew <- subset(datanew, datanew$Teilnahme=="ja")  
  
datanew$Datenschutz <- NULL  
datanew$Teilnahme <- NULL
```

True/False durch 0/1 ersetzen

Bei den Fragen mit der Antwortmöglichkeit Ja oder Nein (z.B. "Interesse an Weltkultur?") sind in der CSV datei für Ja = 1 und für Nein = NA gespeichert. Hier werden die NA's mit 0 überschrieben, da diese für "Nein" stehen.

```
print("Ursprüngliche NA's: ")
sum(is.na(datanew[14:25])) # Interesse
sum(is.na(datanew[29:45])) # Digital
sum(is.na(datanew[46:57])) # Anforderungen
sum(is.na(datanew[58:69])) # KI
sum(is.na(datanew[74:77])) # Helfen

datanew[14:25][is.na(datanew[14:25])] <- 0
datanew[29:45][is.na(datanew[29:45])] <- 0
datanew[46:57][is.na(datanew[46:57])] <- 0
datanew[58:69][is.na(datanew[58:69])] <- 0
datanew[74:77][is.na(datanew[74:77])] <- 0

print("Nach der Anpassung: ")
sum(is.na(datanew[14:25])) # Interesse
sum(is.na(datanew[29:45])) # Digital
sum(is.na(datanew[46:57])) # Anforderungen
sum(is.na(datanew[58:69])) # KI
sum(is.na(datanew[74:77])) # Helfen
```

Spalten Beurteilung entfernen, weil zu viele NA's

Teilweise haben Personen nicht alle Fragen beantwortet. Spalten, die zu wenig Informationen enthalten werden entfernt.

```
datanew$Beurteilung_Erreichbarkeit <- NULL
datanew$Beurteilung_Sammlung <- NULL
datanew$Beurteilung_Sonder <- NULL
datanew$Beurteilung_Bildungsangebot <- NULL
datanew$Beurteilung_Expo <- NULL
datanew$Beurteilung_Online <- NULL
datanew$Beurteilung_Apps <- NULL
```

Weitere Anpassung des Datensatzes

Im Folgenden werden noch die Postleitzahl und die Spalten mit ordinalen Werten angepasst (dort wo es fehlende Werte gibt). Darüberhinaus werden alle Freitext-Angaben beim Bildungsabschluss mit der Kategorie "Sonstiges" überschrieben.

```

#Weitere NA's anzeigen
apply(datanew, function(column) {sum(is.na(column)) })

#NA der PLZ-Column
datanew <- datanew[complete.cases(datanew[, 3]),]

#NA's der ordinalen Spalten durch Mittelwert ersetzen
for(i in 1:ncol(datanew)) {
  datanew[, i][is.na(datanew[, i])] <- round(x= mean(datanew[, i], na.rm=TRUE), dig
its = 0)
}

#Level Bildungsabschluss einschränken, da Freitext
datanew$Bildungsabschluss[which(datanew$Bildungsabschluss != "Haupt-/Realschulabschlu
ss (Mittlere Reife)" & datanew$Bildungsabschluss != "Abitur/(Fach-)Hochschulreife" &
datanew$Bildungsabschluss != "Studium (Fachhochschul-/Hochschulabschluss)")] <- "Sonst
iges"

```

Faktorisierung kategorialer Größen

Zur Verarbeitung durch das Modell ist eine Faktorisierung der kategorialen Größen notwendig. So wird zum Beispiel der Ausprägung “weiblich” des Merkmals Geschlecht die Zahl 0 zugeordnet, während “männlich” als 1 dargestellt wird. Die Ausprägungen sind besonders für die spätere Interpretation wichtig. Eine veranschaulichte Darstellung finden Sie in der Präsentation unter “Visualisierung”.

```

datanew$Geschlecht <- factor(datanew$Geschlecht, levels = c("weiblich","männlich"),o
rdered = TRUE)
datanew$Alter <- factor(datanew$Alter, levels = c("15-29 Jahre", "30-59 Jahre", "60 J
ahre und älter"), ordered = TRUE)
datanew$Bildungsabschluss <- factor(datanew$Bildungsabschluss, levels = c("Haupt-/Rea
lschulabschluss (Mittlere Reife)","Abitur/(Fach-)Hochschulreife", "Studium (Fachhochs
chul-/Hochschulabschluss)","Sonstiges"), ordered = TRUE)
datanew$Anzahl_Besuche <- factor(datanew$Anzahl_Besuche, levels = c("bisher gar nich
t", "nur digital", "höchstens 1-mal im Jahr","2- bis 3-mal im Jahr","mehr als 3-mal i
m Jahr"), ordered = TRUE)
datanew$Empfehlung <- factor(datanew$Empfehlung, levels = c("nein","ja"), ordered =
TRUE)
datanew$Praesenz <- factor(datanew$Praesenz, levels = c("ein Museum digital zu erlebe
n","ein Museum direkt vor Ort zu besuchen"), ordered = TRUE)
datanew$Gesellschaft <- factor(datanew$Gesellschaft, levels = c("allein","mit anderen
zusammen"), ordered = TRUE)
datanew$`Digitale Angebote` <- factor(datanew$`Digitale Angebote`, levels = c("nei
n", "ja"), ordered = TRUE)
datanew$Freizeit_Std <- factor(datanew$Freizeit_Std, levels = c("höchstens 10 Stunde
n","11 bis 20 Stunden","21 bis 30 Stunden","mehr als 30 Stunden"), ordered = TRUE)
datanew$Besonders_wichtig <- factor(datanew$Besonders_wichtig, levels = c("Tradition
& Ordnung","Modernisierung & Selbstverwirklichung","Neuorientierung & Sac
hlichkeit"), ordered = TRUE)

```

Neu entstandene Missing Values entfernen

Durch die Faktorisierung sind einige neue NA's entstanden, diese werden entfernt.

```
row.has.na <- apply(datanew, 1, function(x){any(is.na(x))})
sum(row.has.na)
sapply(datanew, function(column) {sum(is.na(column)) })
datanew <- na.omit(datanew)
row.has.na <- apply(datanew, 1, function(x){any(is.na(x))})
sum(row.has.na)
```

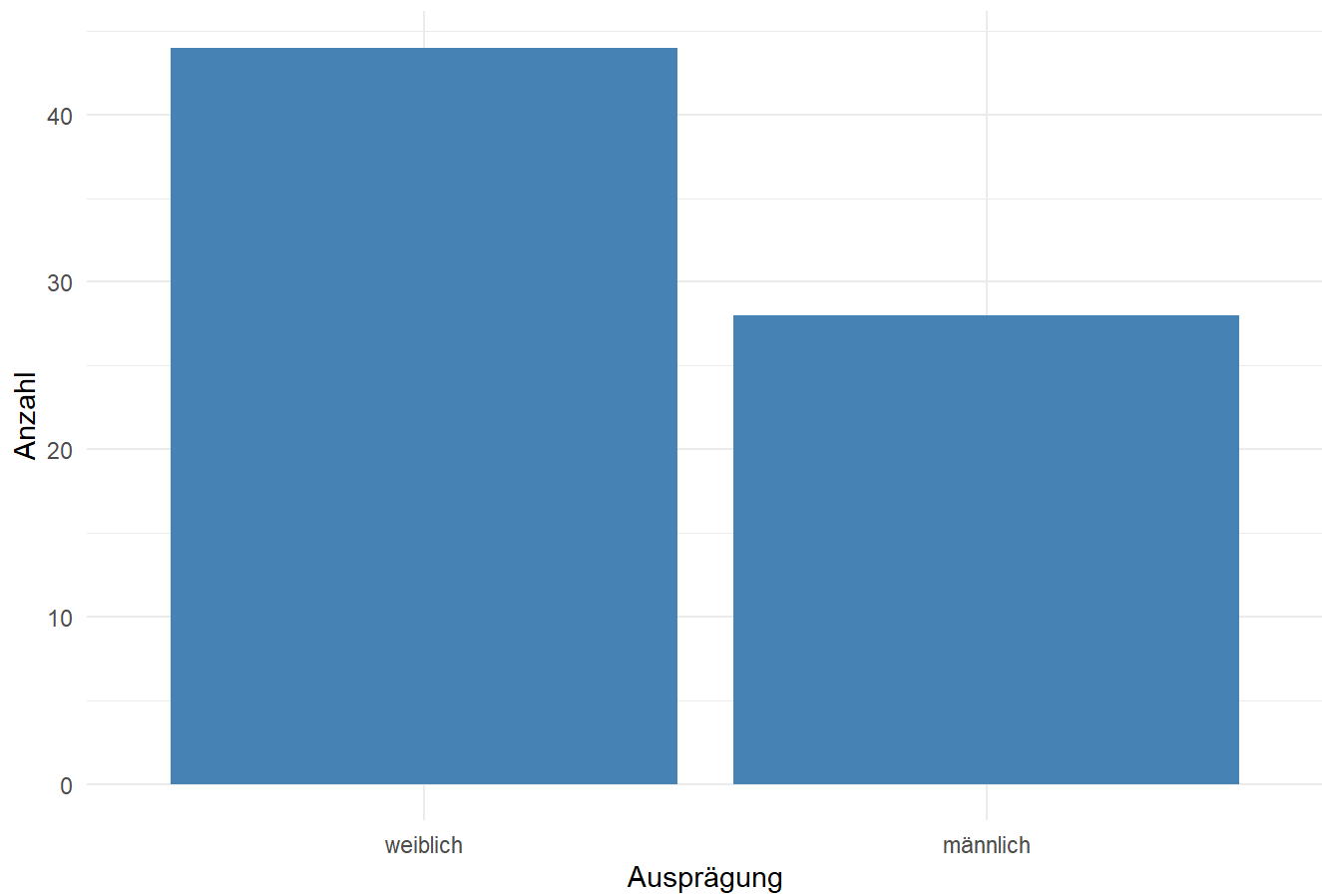
3. Explorative Datenanalyse

Für einen Überblick über den vorhandenen Datensatz werden im Folgenden verschiedene Häufigkeitsdiagramme erstellt.

```
#summary(datanew)

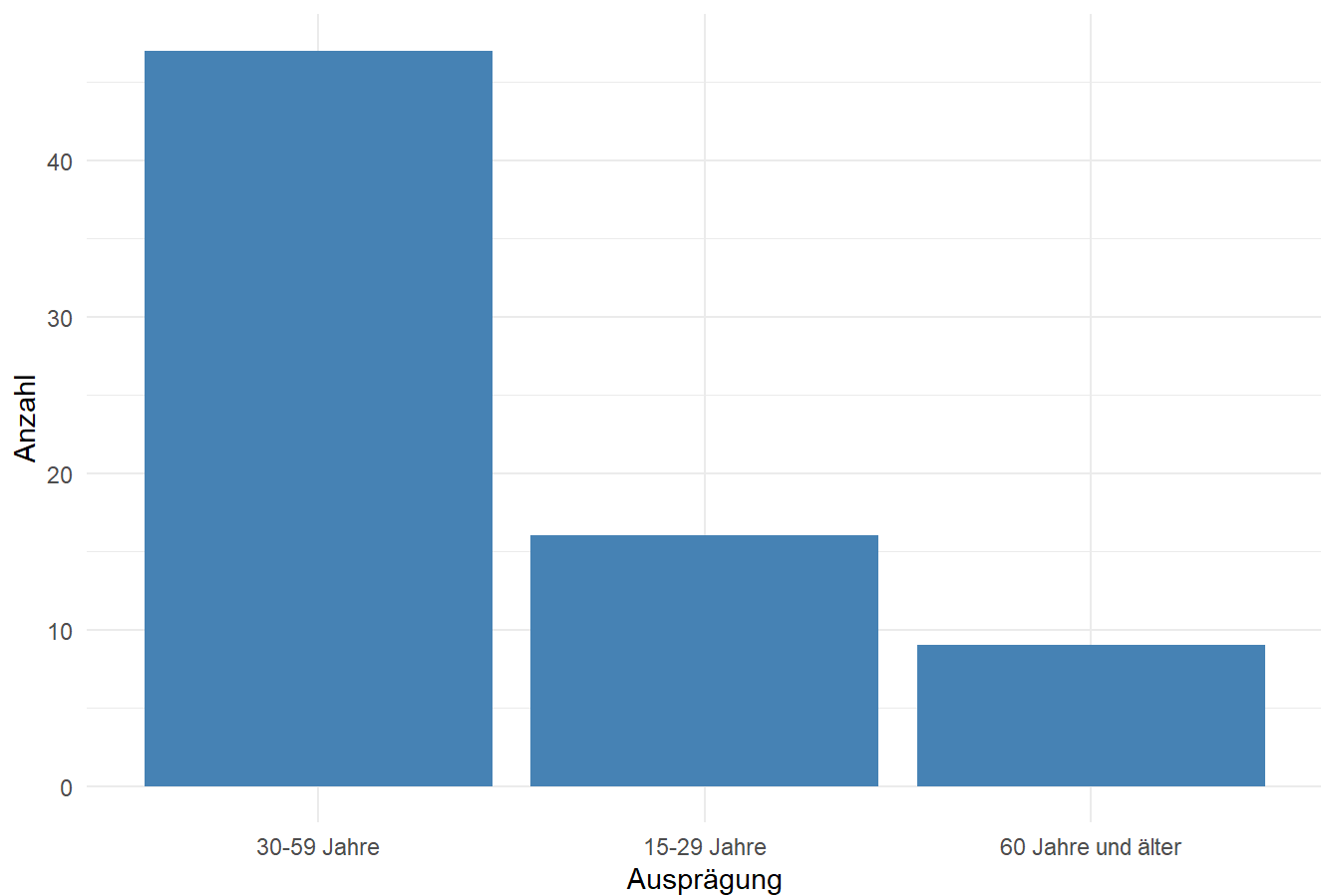
datanew %>% count(Geschlecht) %>%
  ggplot(aes(x=reorder(Geschlecht, -n), y=n))+
  geom_col(fill="steelblue")+
  ggtitle("Häufigkeitenverteilung Geschlecht") + xlab("Ausprägung") + ylab("Anzahl")
+
  theme_minimal()
```

Häufigkeitenverteilung Geschlecht

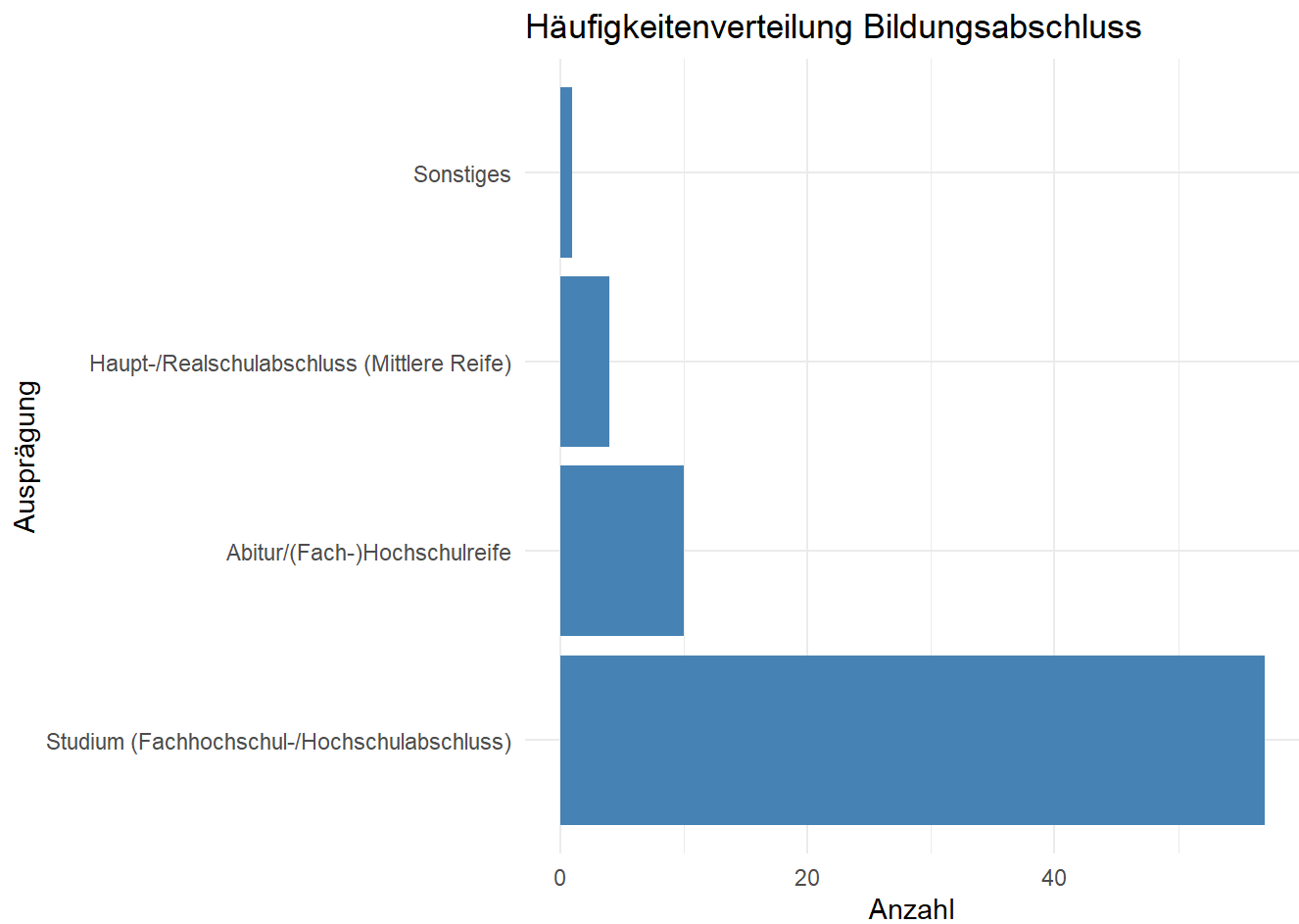


```
datanew %>% count(Alter) %>%  
  ggplot(aes(x=reorder(Alter, -n), y=n))+  
  geom_col(fill="steelblue")+  
  ggtitle("Häufigkeitenverteilung Alter") +  
  xlab("Ausprägung") + ylab("Anzahl") +  
  theme_minimal()
```

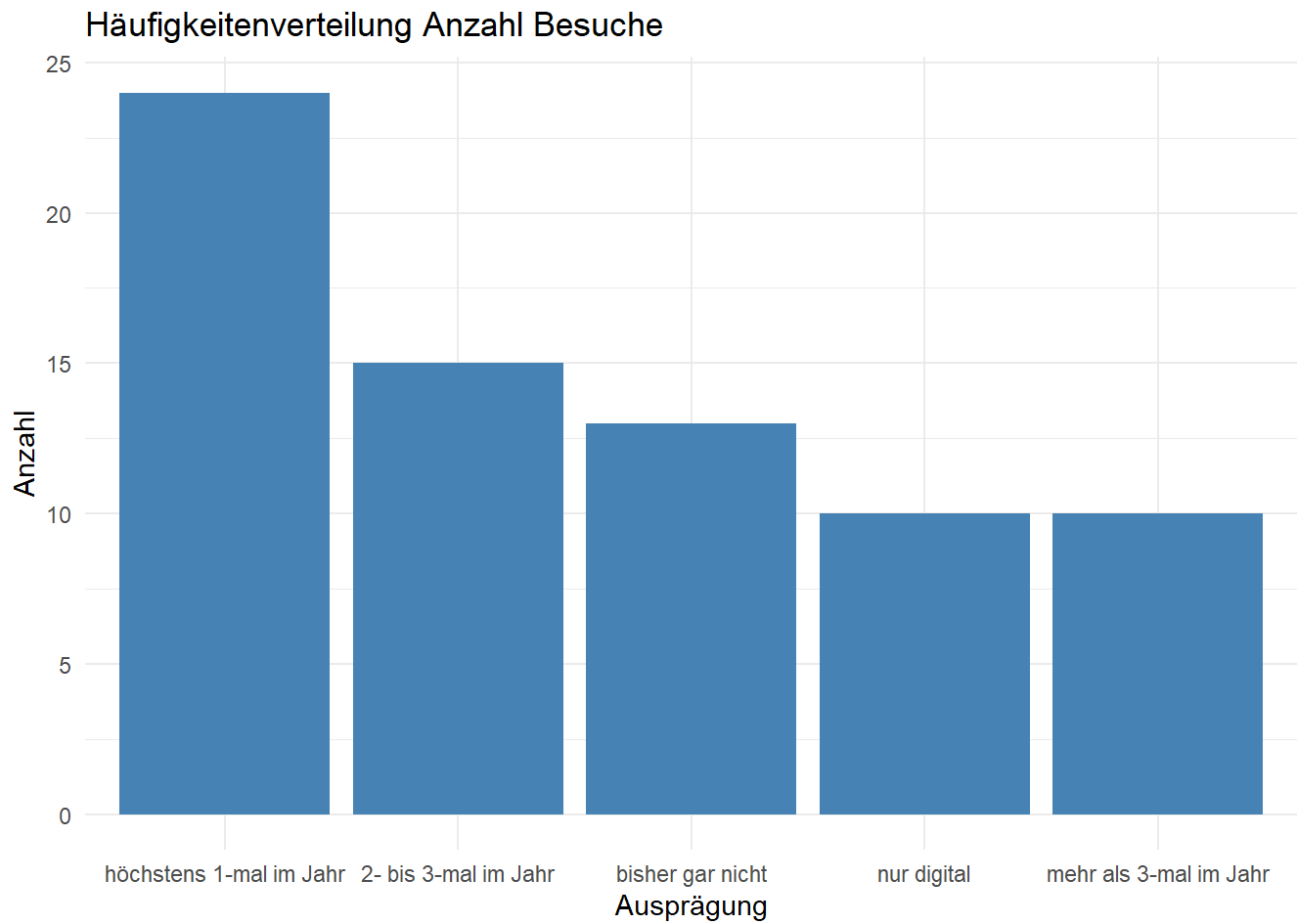

Häufigkeitenverteilung Alter



```
datanew %>% count(Bildungsabschluss) %>%  
  ggplot(aes(x=n, y=reorder(Bildungsabschluss, -n)))+  
  geom_col(fill="steelblue")+  
  ggtitle("Häufigkeitenverteilung Bildungsabschluss") + xlab("Anzahl") + ylab("Ausprä  
gung") +  
  theme_minimal()
```

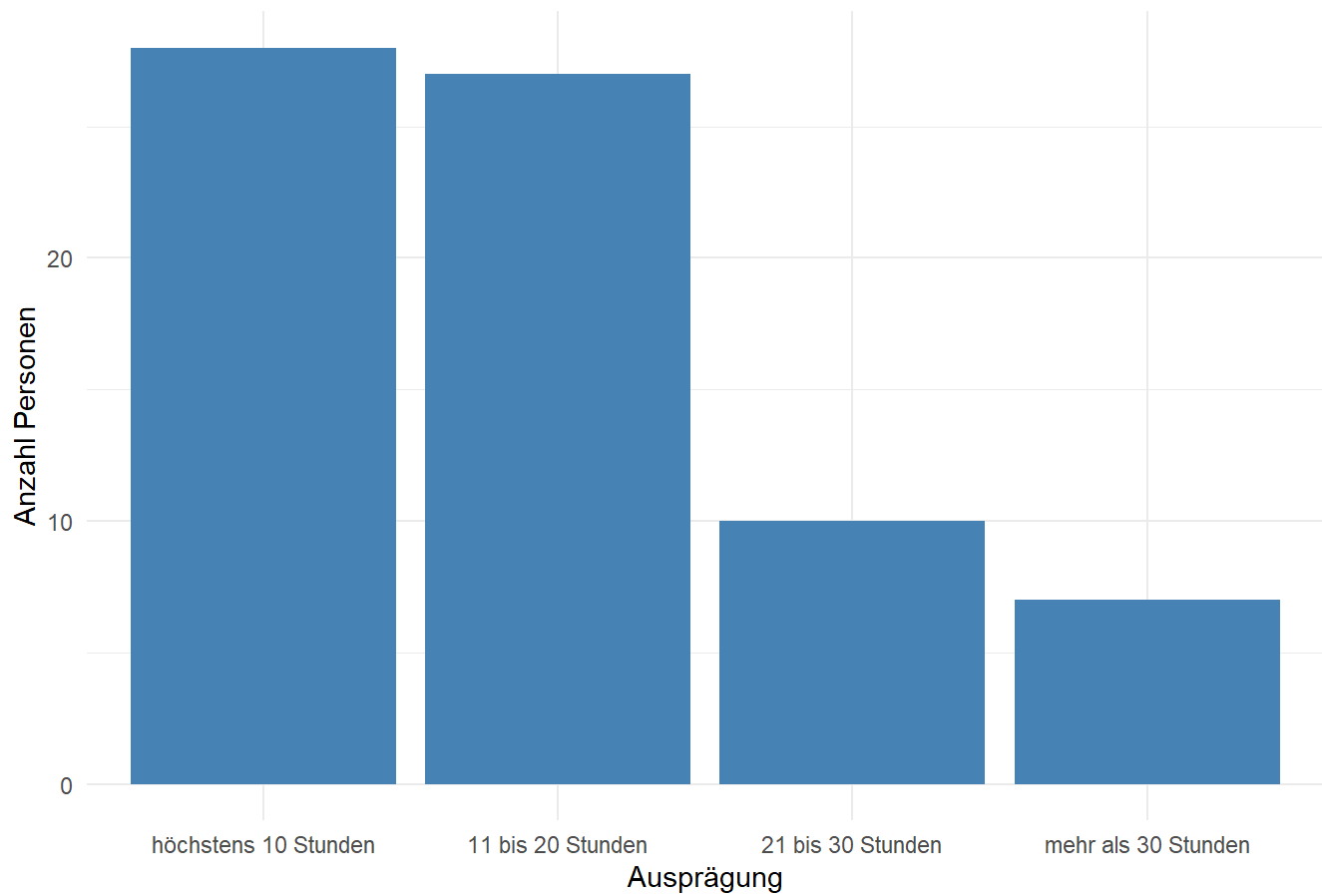


```
datanew %>% count(Anzahl_Besuche) %>%  
  ggplot(aes(x=reorder(Anzahl_Besuche, -n), y=n))+  
  geom_col(fill="steelblue")+  
  ggtitle("Häufigkeitenverteilung Anzahl Besuche") + xlab("Ausprägung") + ylab("Anzahl") +  
  theme_minimal()
```



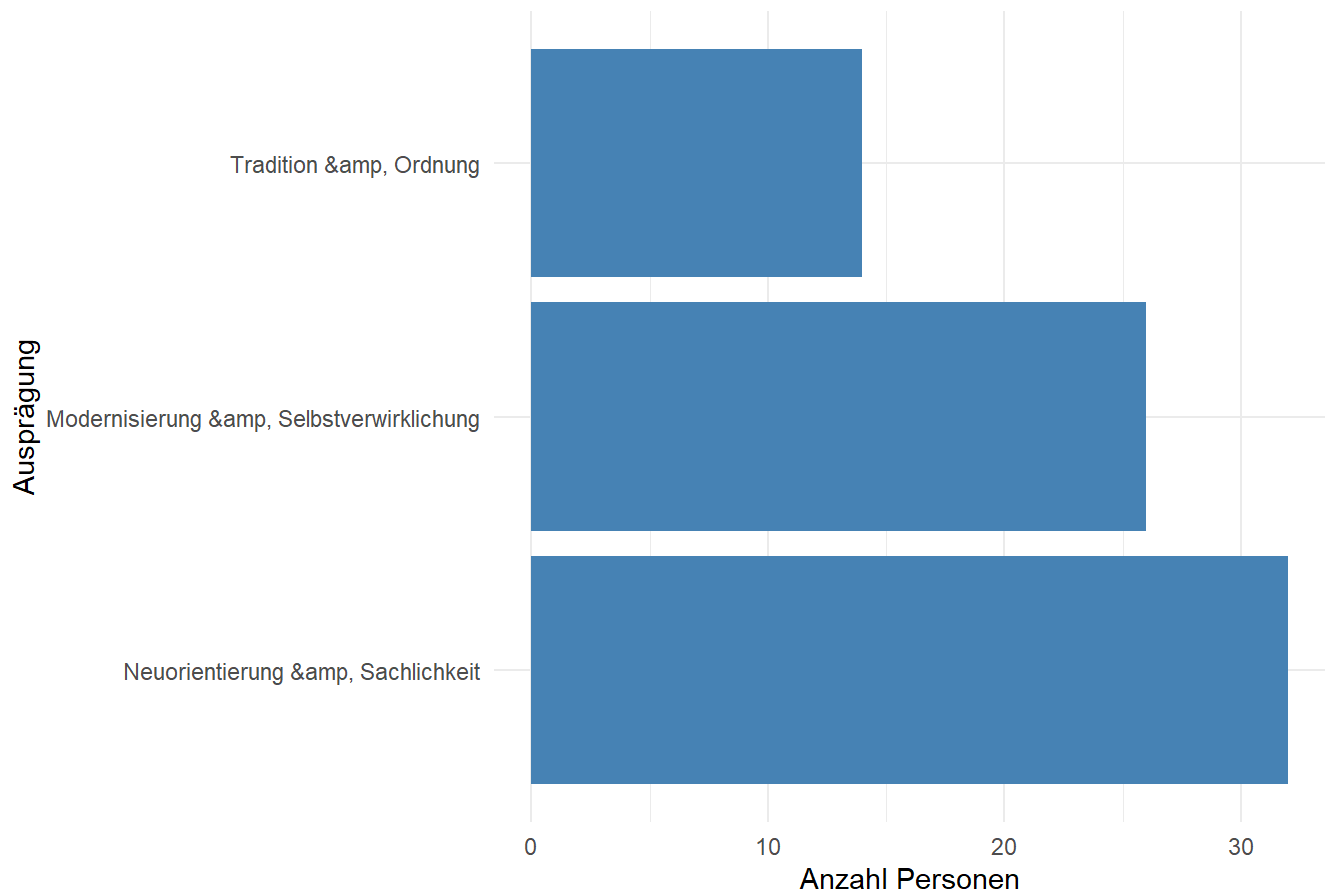
```
datanew %>% count(Freizeit_Std) %>%  
  ggplot(aes(x=reorder(Freizeit_Std, -n), y=n))+  
  geom_col(fill="steelblue")+  
  ggtitle("Häufigkeitenverteilung Stunden Freizeit") + xlab("Ausprägung") + ylab("Anzahl Personen") +  
  theme_minimal()
```

Häufigkeitenverteilung Stunden Freizeit

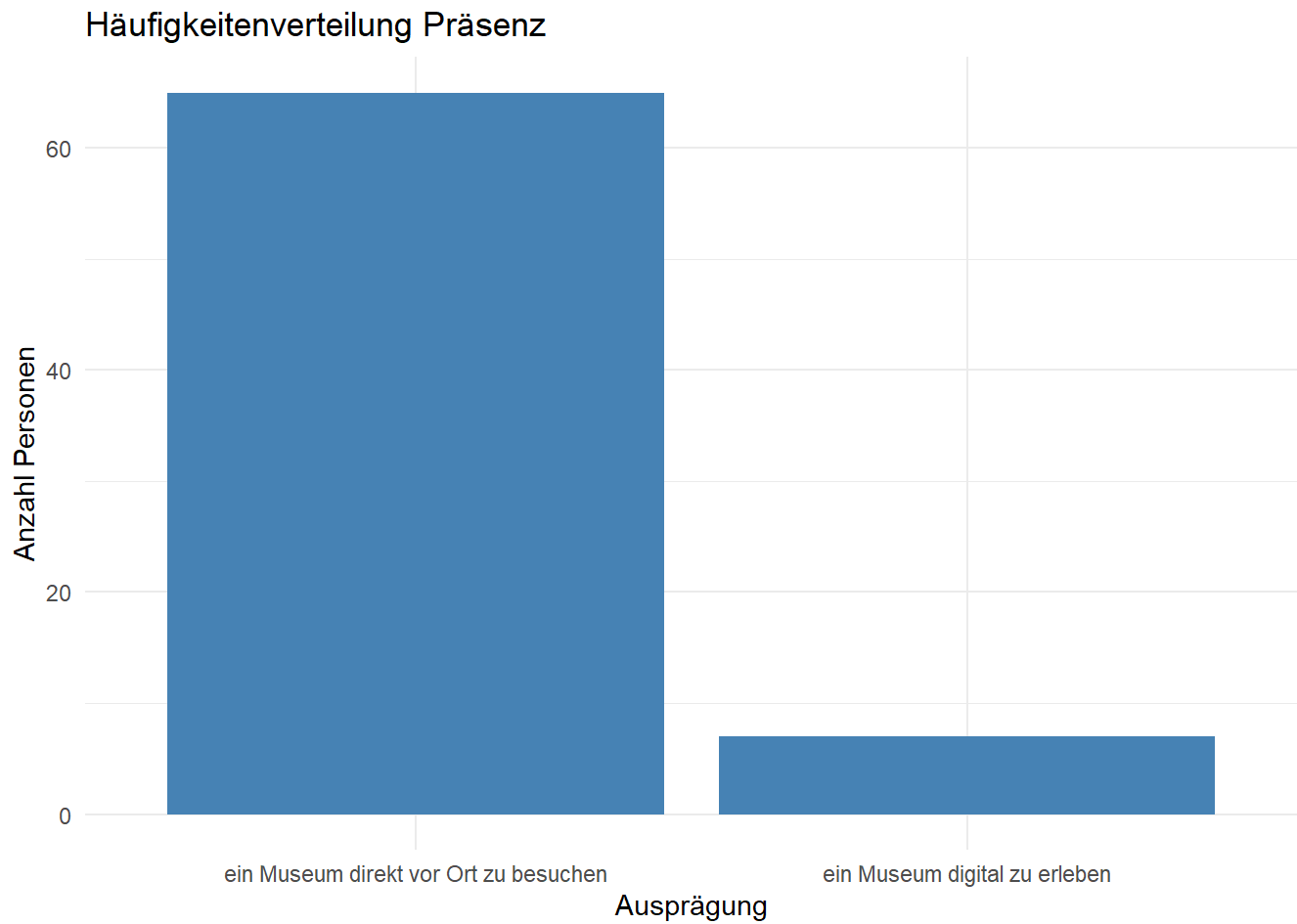


```
datanew %>% count(Besonders_wichtig) %>%  
  ggplot(aes(x=n, y=reorder(Besonders_wichtig, -n)))+  
  geom_col(fill="steelblue")+  
  ggtitle("Häufigkeitenverteilung Prioritäten") + xlab("Anzahl Personen") + ylab("Ausprägung") +  
  theme_minimal()
```

Häufigkeitenverteilug Prioritäten

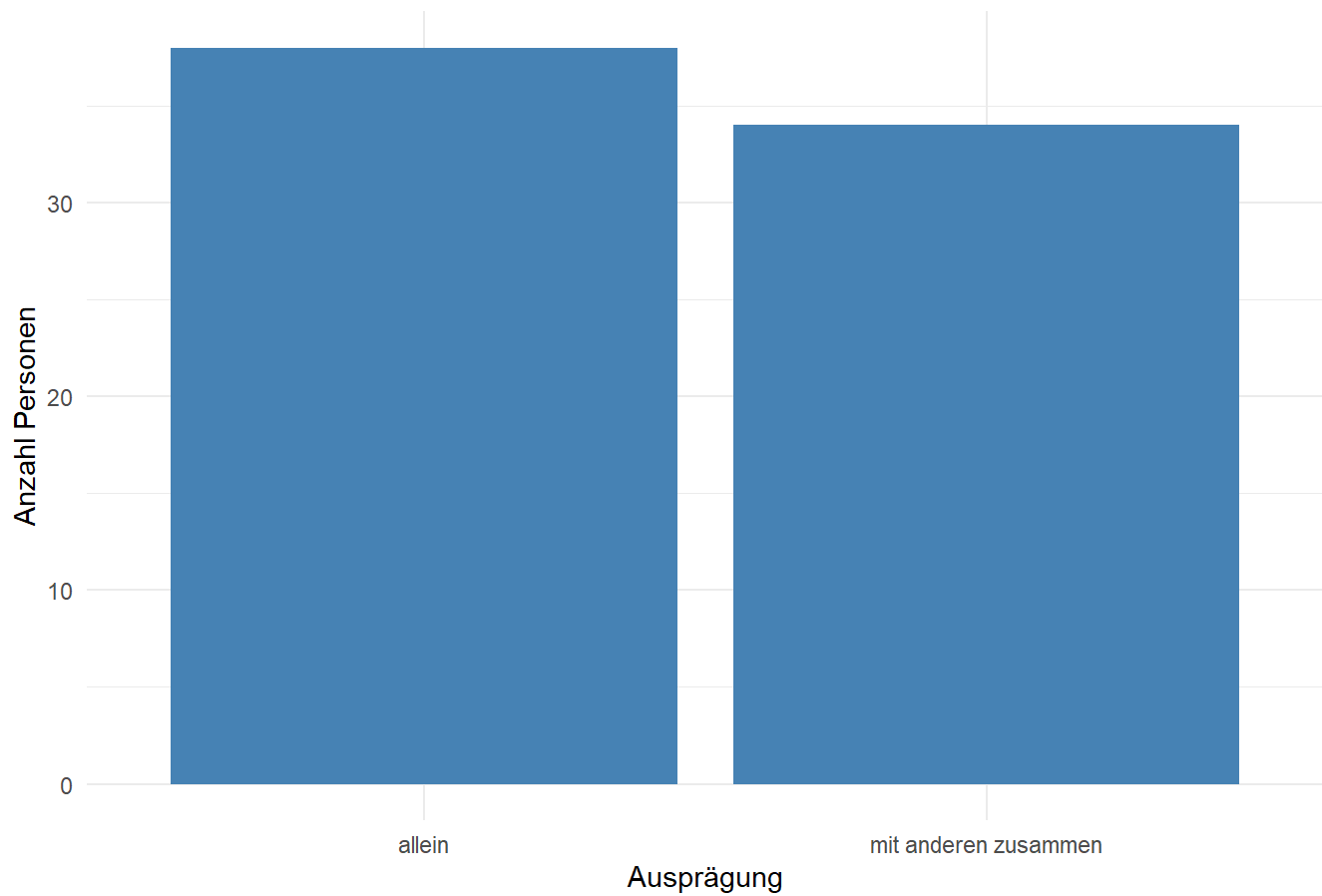


```
datanew %>% count(Praesenz) %>%  
  ggplot(aes(x=reorder(Praesenz, -n), y=n))+  
  geom_col(fill="steelblue")+  
  ggtitle("Häufigkeitenverteilung Präsenz") + xlab("Ausprägung") + ylab("Anzahl Perso  
nen") +theme_minimal()
```



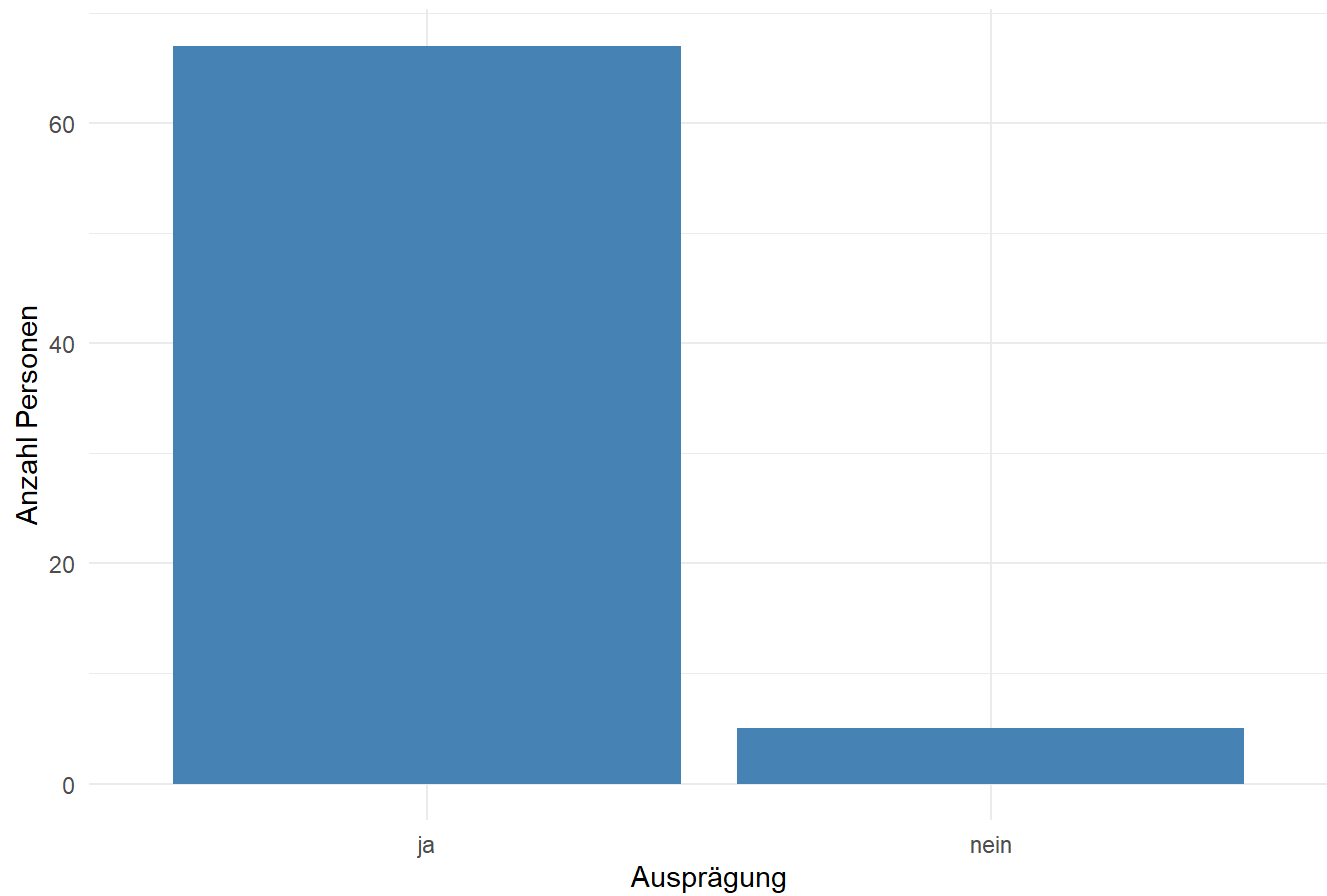
```
datanew %>% count(Gesellschaft) %>%  
  ggplot(aes(x=reorder(Gesellschaft, -n), y=n))+  
  geom_col(fill="steelblue")+  
  ggtitle("Häufigkeitenverteilung Gesellschaft bei Museumsbesuchen") +  
  xlab("Ausprägung") + ylab("Anzahl Personen") +  
  theme_minimal()
```

Häufigkeitenverteilung Gesellschaft bei Museumsbesuchen

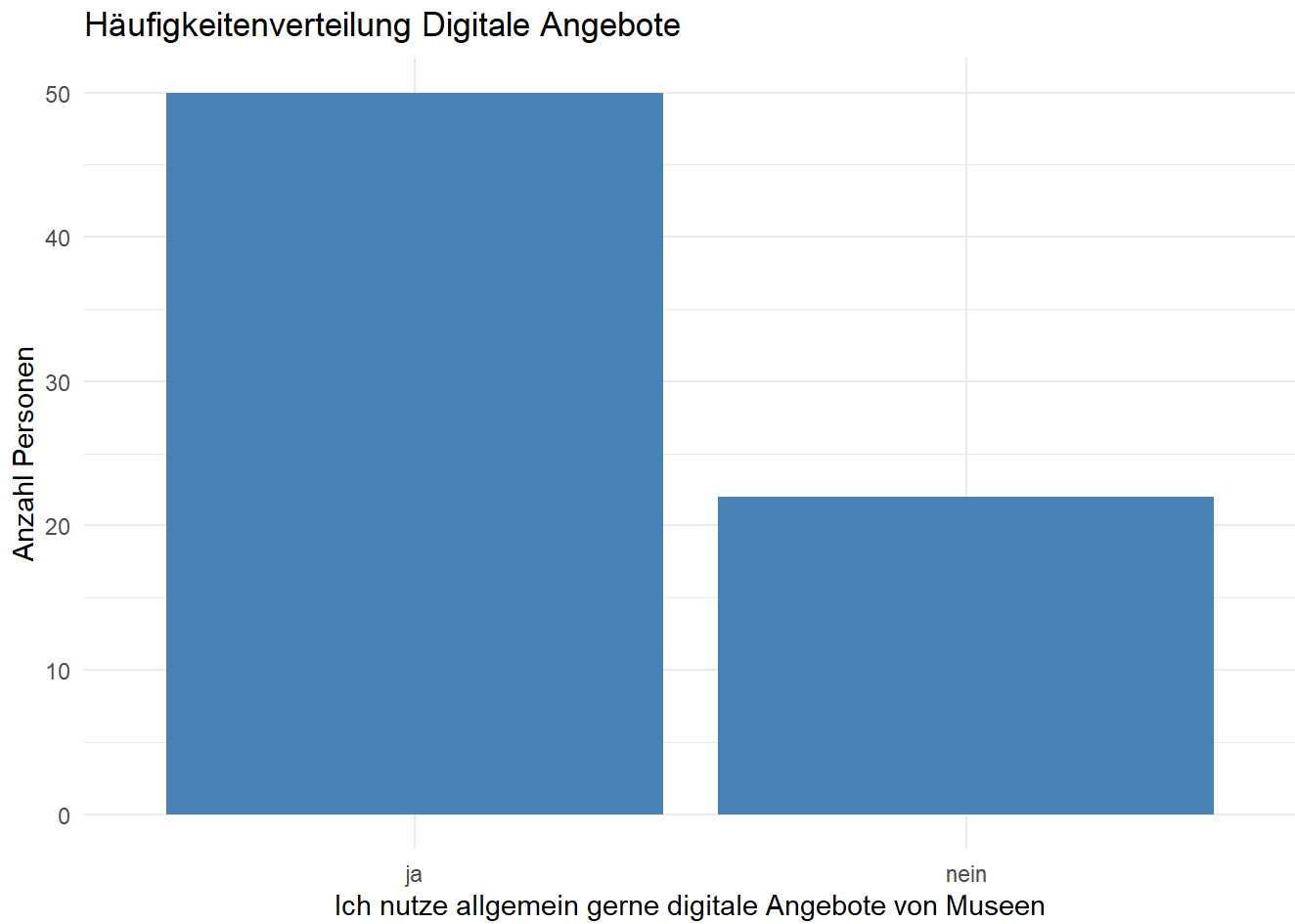


```
datanew %>% count(Empfehlung) %>%  
  ggplot(aes(x=reorder(Empfehlung, -n), y=n))+  
  geom_col(fill="steelblue")+  
  ggtitle("Häufigkeitenverteilung Weiterempfehlung des Museums")+  
  xlab("Ausprägung") +  
  ylab("Anzahl Personen") + theme_minimal()
```

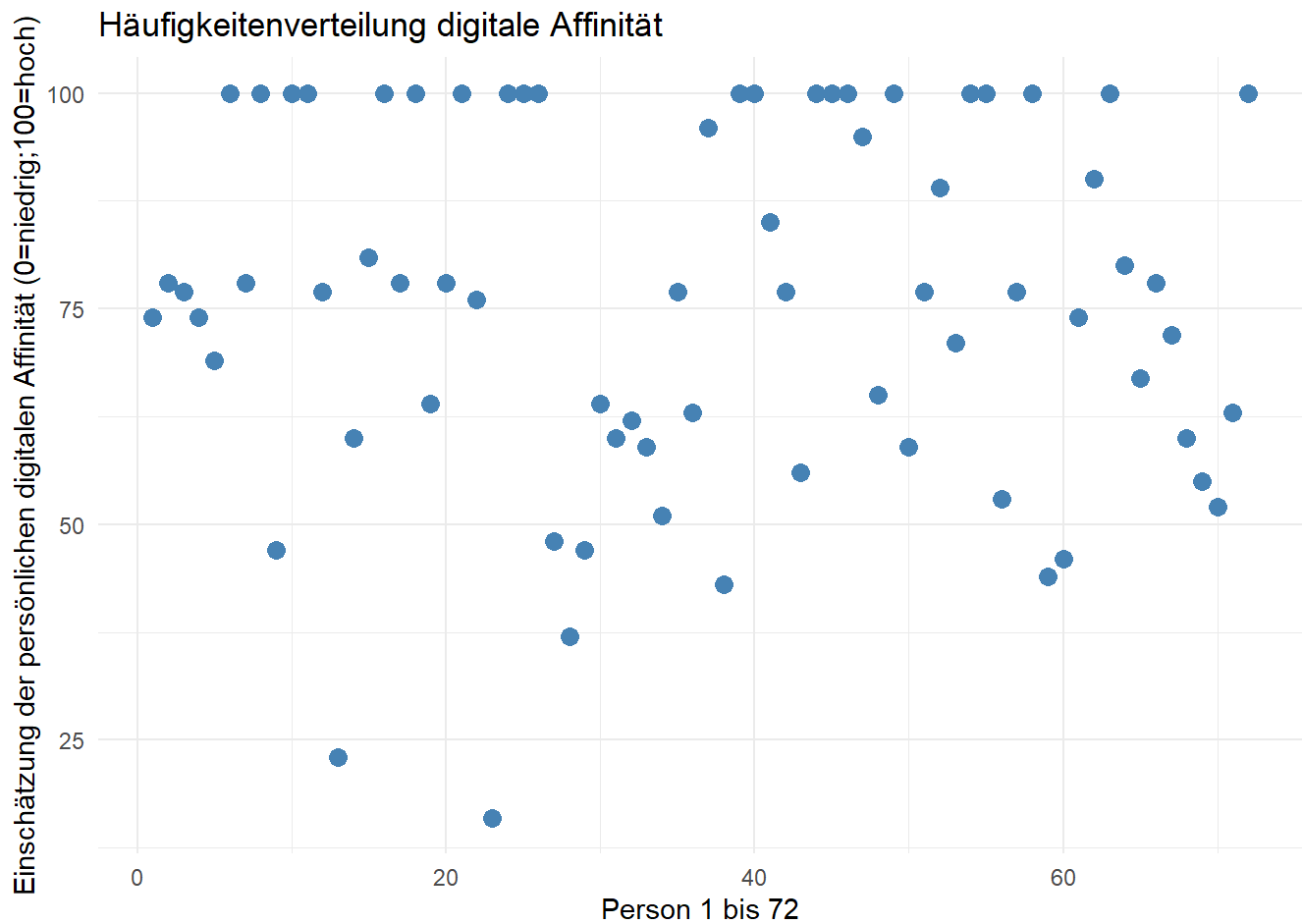
Häufigkeitenverteilung Weiterempfehlung des Museums



```
datanew %>% count(`Digitale Angebote`) %>%  
  ggplot(aes(x=reorder(`Digitale Angebote`, -n), y=n))+  
  geom_col(fill="steelblue")+  
  ggtitle("Häufigkeitenverteilung Digitale Angebote") +  
  xlab('Ich nutze allgemein gerne digitale Angebote von Museen') +  
  ylab("Anzahl Personen") + theme_minimal()
```

```
datanew %>% ggplot(aes(y=datanew$digitale_Affinitaet, x=c(1:72)))+  
  geom_point(col="steelblue",size=3)+  
  ggtitle("Häufigkeitenverteilung digitale Affinität") +  
  xlab('Person 1 bis 72') +  
  ylab("Einschätzung der persönlichen digitalen Affinität (0=niedrig;100=hoch)") +  
  theme_minimal()
```

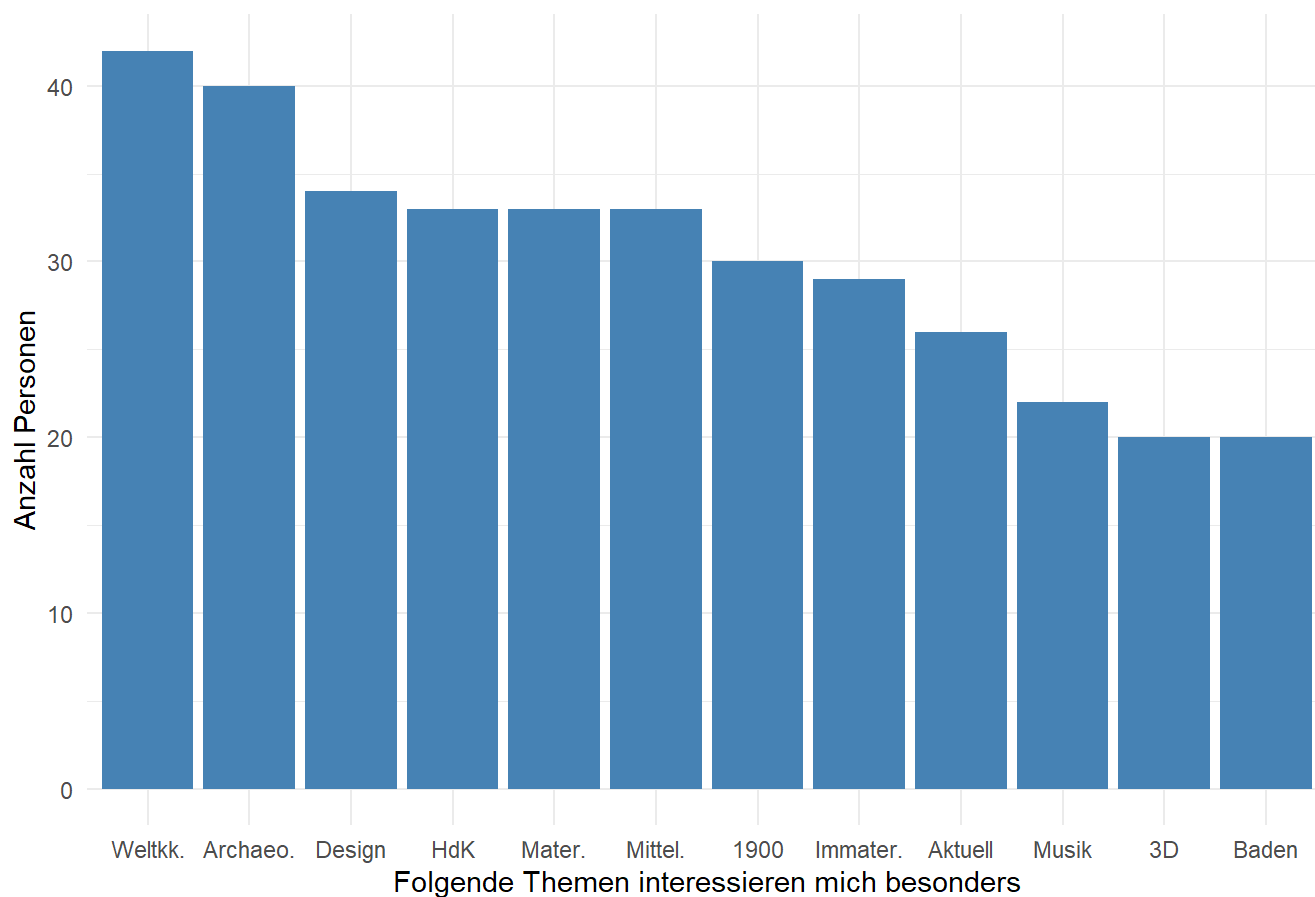


```

values <- c(length(which(datanew$Interesse_1900 == 1)),
            length(which(datanew$Interesse_Mittelalter == 1)),
            length(which(datanew$Interesse_Weltkultur == 1)),
            length(which(datanew$Interesse_3D == 1)),
            length(which(datanew$Interesse_Archaeologie == 1)),
            length(which(datanew$Interesse_Design == 1)),
            length(which(datanew$Interesse_Baden == 1)),
            length(which(datanew$Interesse_Aktuell == 1)),
            length(which(datanew$Interesse_Musik == 1)),
            length(which(datanew$Interesse_Materiell == 1)),
            length(which(datanew$Interesse_Immateriell == 1)),
            length(which(datanew$Interesse_HdK == 1)))
names <- c("1900", "Mittel.", "Weltkk.", "3D", "Archaeo.", "Design",
           "Baden", "Aktuell", "Musik", "Mater.", "Immater.", "HdK")
dataplot <- data.frame(values, names)
ggplot(dataplot, aes(x= reorder(names, -values), values)) +
  geom_col(fill="steelblue") +
  ggtitle("Häufigkeitenverteilung Interesse") +
  xlab('Folgende Themen interessieren mich besonders') +
  ylab("Anzahl Personen") +
  theme_minimal()

```

Häufigkeitenverteilung Interesse

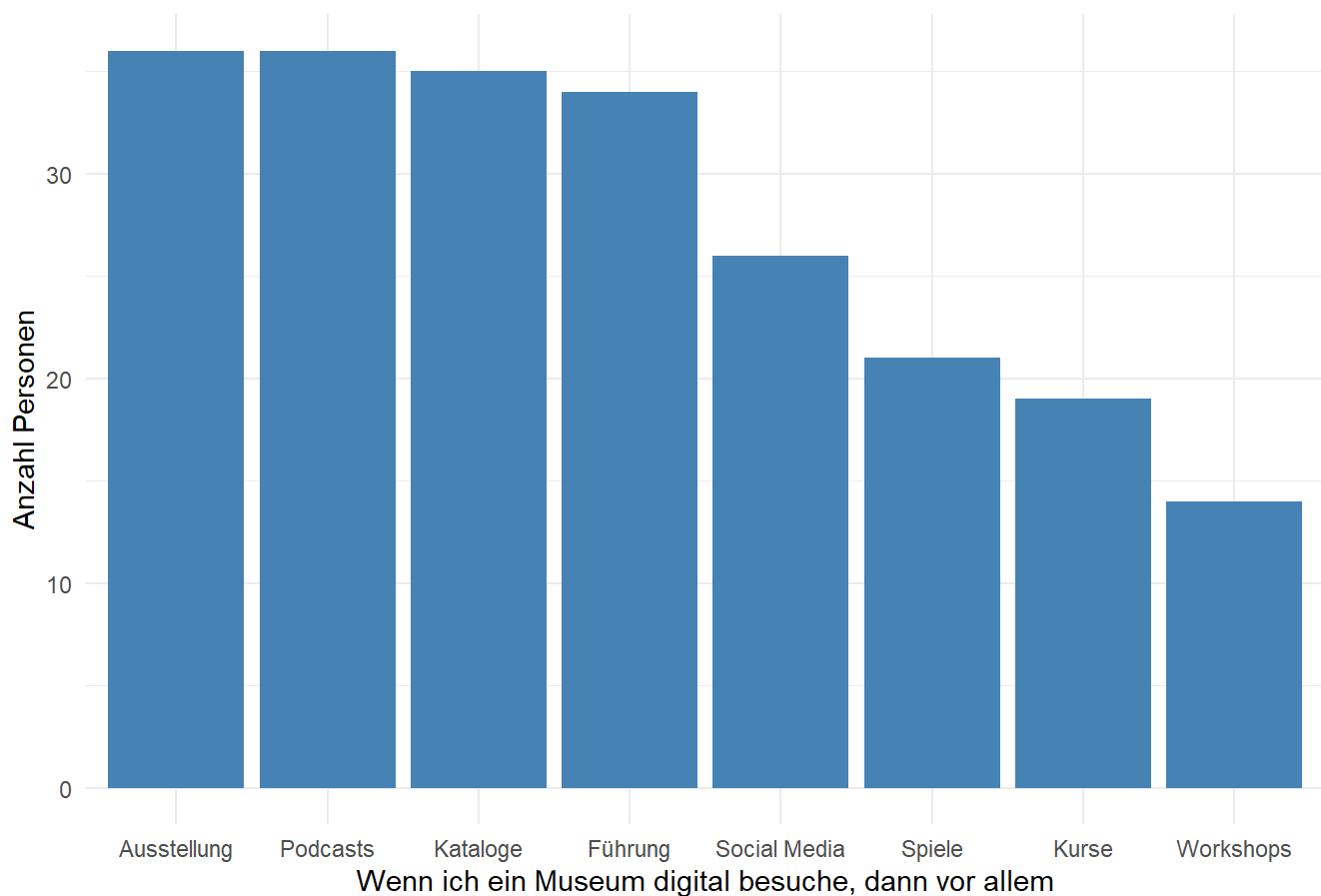


```

values <- c(length(which(datanew$Digital_Fuehrung == 1)),
            length(which(datanew$Digital_Workshops == 1)),
            length(which(datanew$Digital_Ausstellungen == 1)),
            length(which(datanew$Digital_Podcasts == 1)),
            length(which(datanew$Digital_Kataloge == 1)),
            length(which(datanew$Digital_Kurse == 1)),
            length(which(datanew$Digital_Spiele == 1)),
            length(which(datanew$Digital_SocialMedia == 1)))
names <- c("Führung", "Workshops", "Ausstellung", "Podcasts",
           "Kataloge", "Kurse", "Spiele", "Social Media")
dataplot <- data.frame(values, names)
ggplot(dataplot, aes(x= reorder(names, -values), values)) +
  geom_col(fill="steelblue") +
  ggtitle("Häufigkeitenverteilung Digital") +
  xlab('Wenn ich ein Museum digital besuche, dann vor allem') +
  ylab("Anzahl Personen") +
  theme_minimal()

```

Häufigkeitenverteilung Digital

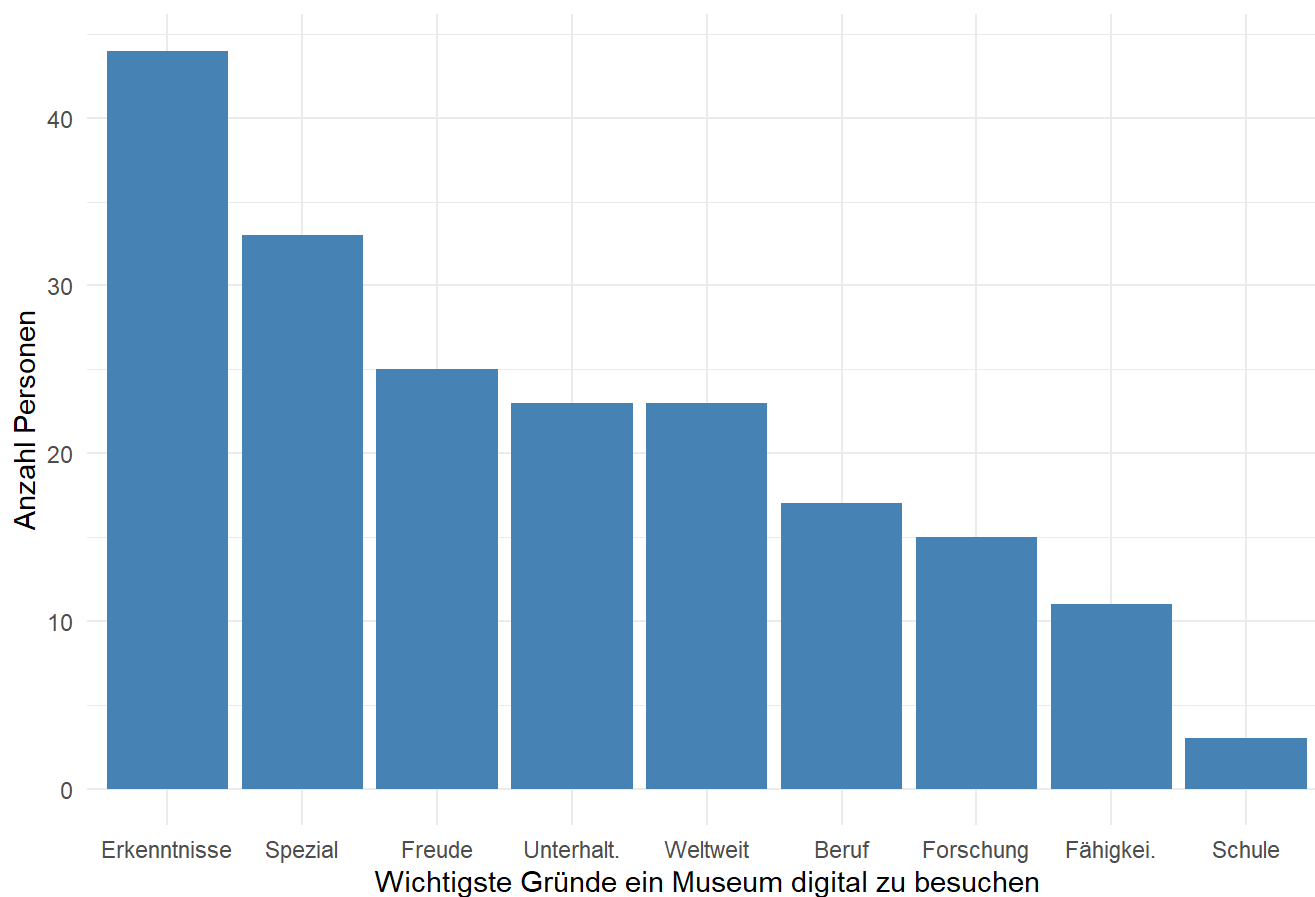


```

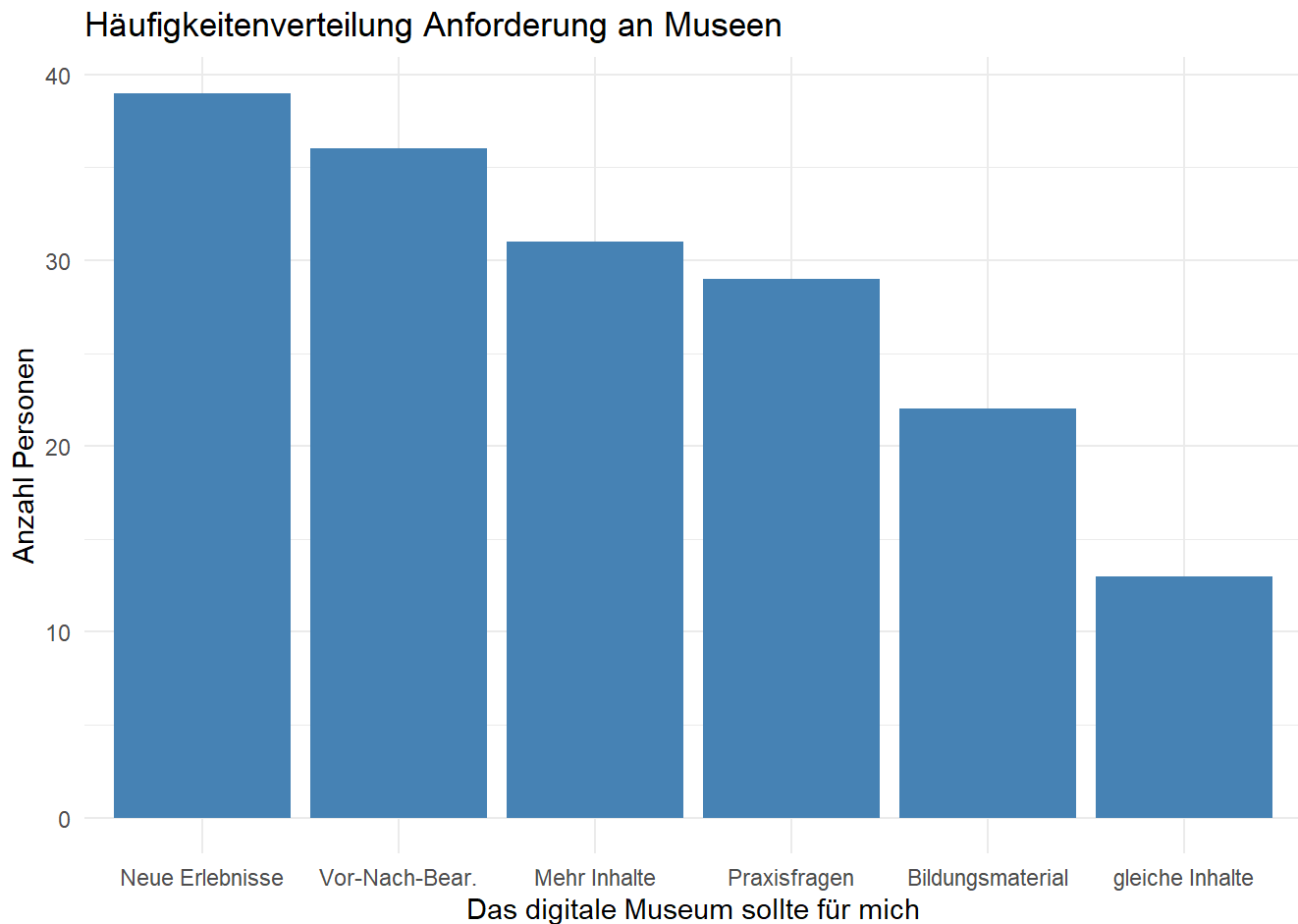
values <- c(length(which(datanew$Grunddigital_Freude == 1)),
            length(which(datanew$Grunddigital_Beruf == 1)),
            length(which(datanew$Grunddigital_Erkenntnisse == 1)),
            length(which(datanew$Grunddigital_Faehigkeiten == 1)),
            length(which(datanew$Grunddigital_Schule == 1)),
            length(which(datanew$Grunddigital_Spezial == 1)),
            length(which(datanew$Grunddigital_Forschung == 1)),
            length(which(datanew$Grunddigital_Unterhaltung == 1)),
            length(which(datanew$Grunddigital_Weltweit == 1)))
names <- c("Freude", "Beruf", "Erkenntnisse", "Fähigkei.", "Schule",
           "Spezial", "Forschung", "Unterhalt.", "Weltweit")
dataplot <- data.frame(values, names)
ggplot(dataplot, aes(x= reorder(names, -values), values)) +
  geom_col(fill="steelblue") +
  ggtitle("Häufigkeitenverteilung Grund für digital") +
  xlab('Wichtigste Gründe ein Museum digital zu besuchen') +
  ylab("Anzahl Personen") +
  theme_minimal()

```

Häufigkeitenverteilung Grund für digital

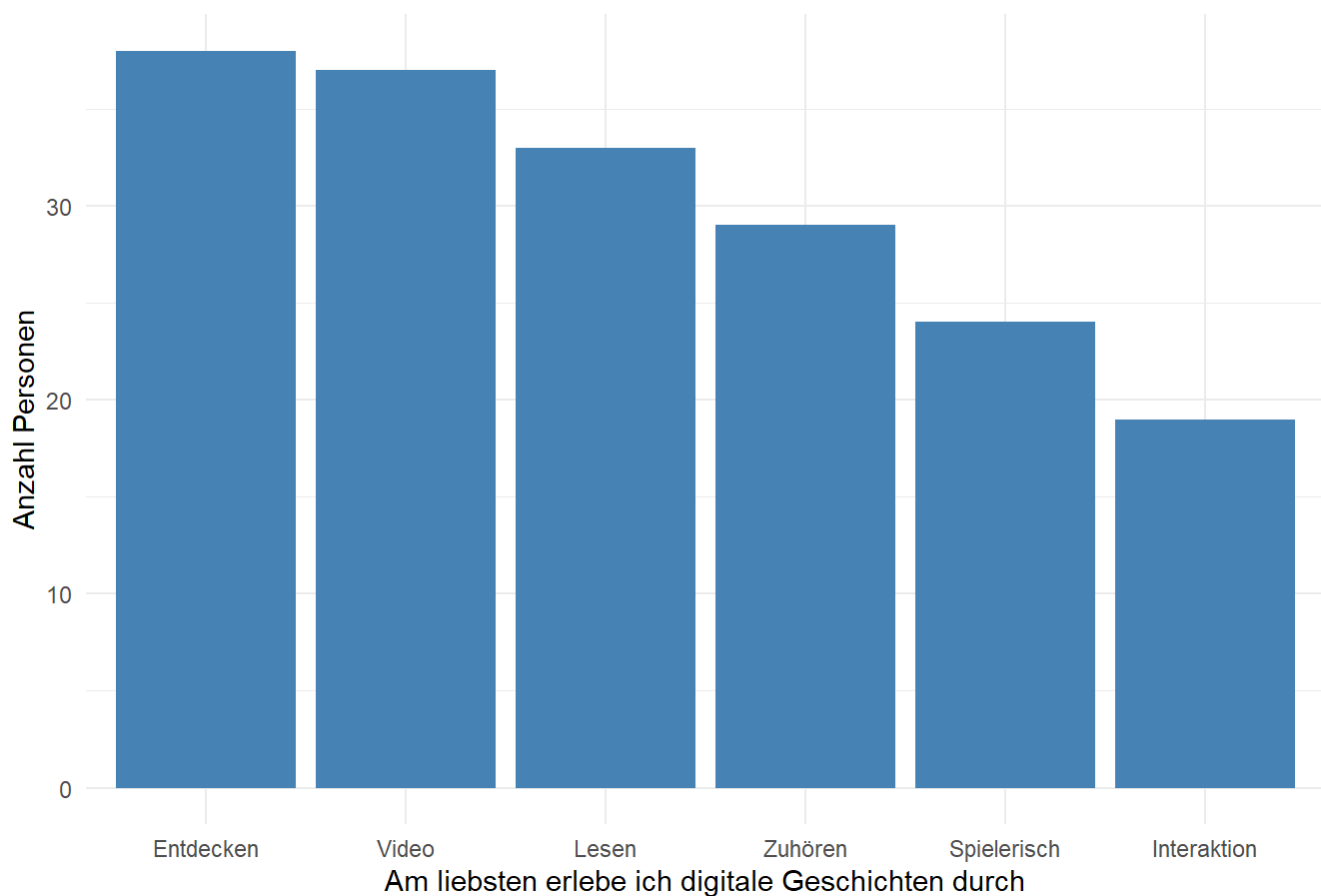


```
values <- c(length(which(datanew$Anforderungen_Praxis == 1)),
            length(which(datanew$Anforderungen_Ausstellungen == 1)),
            length(which(datanew$Anforderungen_gleicheInhalte == 1)),
            length(which(datanew$Anforderungen_Bildungsmaterial == 1)),
            length(which(datanew$Anforderungen_MehrInhalte == 1)),
            length(which(datanew$Anforderungen_NeueErlebnisse == 1)))
names <- c("Praxisfragen", "Vor-Nach-Bear.", "gleiche Inhalte",
           "Bildungsmaterial", "Mehr Inhalte", "Neue Erlebnisse")
dataplot <- data.frame(values, names)
ggplot(dataplot, aes(x= reorder(names, -values), values)) +
  geom_col(fill="steelblue") +
  ggtitle("Häufigkeitenverteilung Anforderung an Museen") +
  xlab('Das digitale Museum sollte für mich') +
  ylab("Anzahl Personen") +
  theme_minimal()
```



```
values <- c(length(which(datanew$Erlebnis_Video == 1)),
            length(which(datanew$Erlebnis_Lesen == 1)),
            length(which(datanew$Erlebnis_Zuhoeren == 1)),
            length(which(datanew$Erlebnis_Entdecken == 1)),
            length(which(datanew$Erlebnis_Interaktion == 1)),
            length(which(datanew$Erlebnis_Spielerisch == 1)))
names <- c("Video", "Lesen", "Zuhören", "Entdecken", "Interaktion", "Spielerisch")
dataplot <- data.frame(values, names)
ggplot(dataplot, aes(x= reorder(names, -values), values)) +
  geom_col(fill="steelblue") +
  ggtitle("Häufigkeitenverteilung Erlebnisse") +
  xlab('Am liebsten erlebe ich digitale Geschichten durch') +
  ylab("Anzahl Personen") +
  theme_minimal()
```

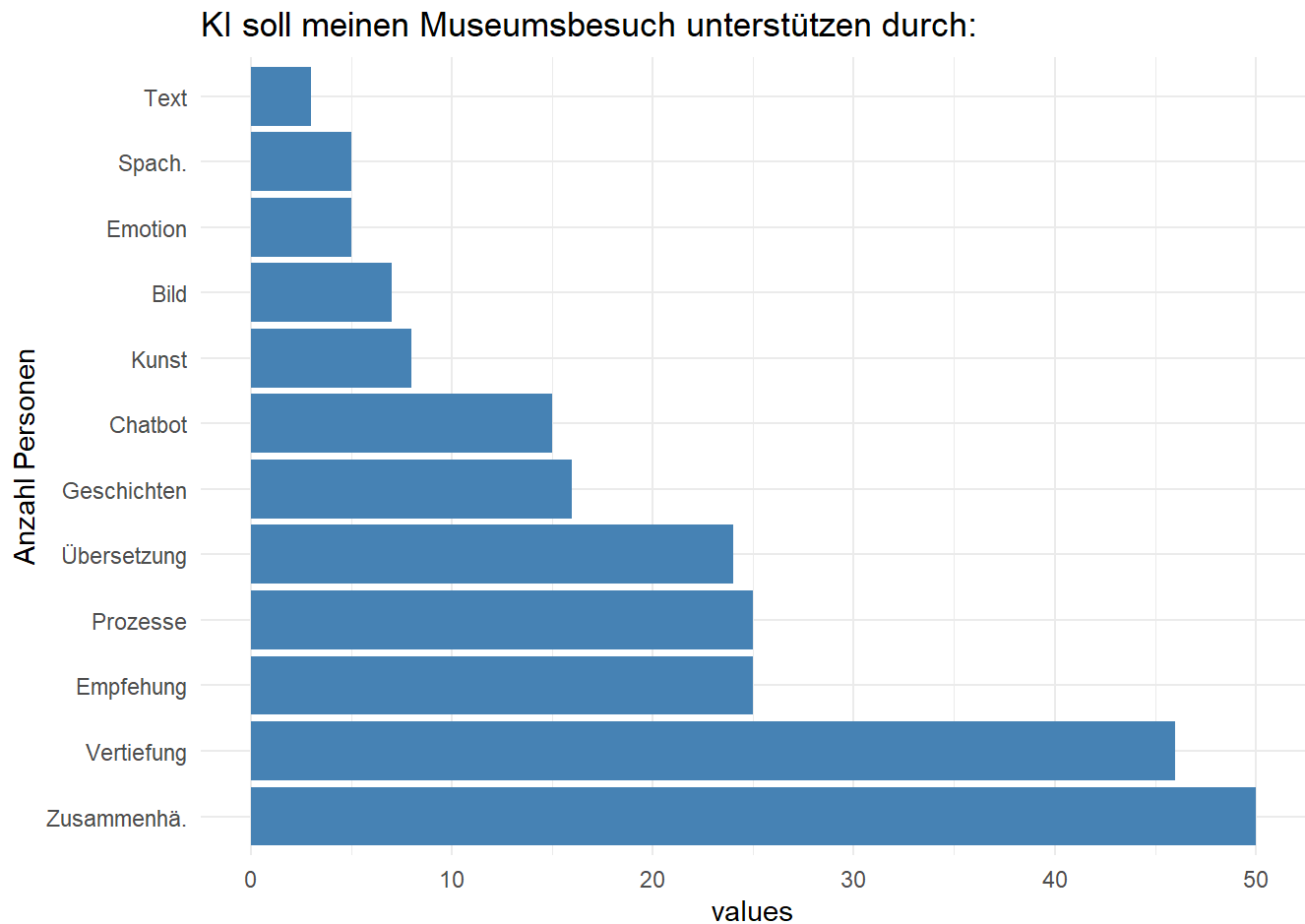
Häufigkeitenverteilung Erlebnisse



```

values <- c(length(which(datanew$KI_Uebersetzung == 1)),
            length(which(datanew$KI_Bildererkennung == 1)),
            length(which(datanew$KI_indv_Empfehlung == 1)),
            length(which(datanew$KI_Texterstellung == 1)),
            length(which(datanew$KI_Vertiefungsinfos == 1)),
            length(which(datanew$KI_Spracherkennung == 1)),
            length(which(datanew$KI_Chatbot == 1)),
            length(which(datanew$KI_Zusammenhaenge_sichtbar_machen == 1)),
            length(which(datanew$KI_Emotionserkennung == 1)),
            length(which(datanew$KI_Geschichten_generieren == 1)),
            length(which(datanew$KI_generierte_Kunst == 1)),
            length(which(datanew$KI_neue_kreative_Prozesse == 1)))
names <- c("Übersetzung", "Bild", "Empfehlung", "Text", "Vertiefung", "Spach.",
           "Chatbot", "Zusammenhä.", "Emotion", "Geschichten", "Kunst", "Prozesse")
dataplot <- data.frame(values, names)
ggplot(dataplot, aes(x= values, y=reorder(names, -values)))+
  geom_col(fill="steelblue")+
  ggtitle("KI soll meinen Museumsbesuch unterstützen durch: ") +
  ylab("Anzahl Personen")+
  theme_minimal()

```



4. Erstellen relevanter Teildatensätze

Für das Modell und die anschließende Interpretation wird der Datensatz in relevante Teildatensätze zerlegt. Die Auswahl der Merkmale erfolgte auch in Absprache mit dem Badischen Landesmuseum nach dessen Prioritäten.

Die Variable "Bildungsabschluss" haben wir in Datensatz 2,3,4 rausgenommen, da hauptsächlich homogene Ausprägungen vorhanden waren.


```

#Dataset 1. "Der Demographische"
dataset1 <- subset(datanew,select=c(Geschlecht,Alter,Bildungsabschluss,Freizeit_Std,A
nzahl_Besuche,digitale_Affinitaet,Praesenz, Gesellschaft))

#Dataset 2. "Interesse & digitales Erlebnis"
dataset2 <- subset(datanew,select=c(Geschlecht,Alter,Interesse_Archaeologie, Interess
e_Weltkultur, Interesse_Mittelalter, Interesse_Design, Interesse_Baden, Interesse_190
0, Interesse_Aktuell, Interesse_Musik, Interesse_Materiell, Interesse_Immateriell, In
teresse_HdK, Interesse_3D,Erlebnis_Video,Erlebnis_Lesen,Erlebnis_Zuhoeren,Erlebnis_En
tdecken,Erlebnis_Interaktion,Erlebnis_Spielerisch))

#Dataset 3. "Gründe für digitale Besuche & Nutzung "
dataset3 <- subset(datanew,select=c(Geschlecht,Alter,Digital_Fuehrung, Digital_Worksh
ops, Digital_Ausstellungen, Digital_Podcasts, Digital_Kataloge, Digital_Kurse, Digita
l_Spiele, Digital_SocialMedia, Grunddigital_Freude, Grunddigital_Unterhaltung, Grundd
igital_Erkenntnisse, Grunddigital_Faehigkeiten, Grunddigital_Schule, Grunddigital_Spe
zial, Grunddigital_Forschung, Grunddigital_Beruf, Grunddigital_Weltweit))

#Dataset 4. "Interesse, digitales Erlebnis & Grunddigital"
dataset4 <- subset(datanew,select=c(Geschlecht,Alter,Interesse_Archaeologie, Interess
e_Weltkultur, Interesse_Mittelalter, Interesse_Design, Interesse_Baden, Interesse_190
0, Interesse_Aktuell, Interesse_Musik, Interesse_Materiell, Interesse_Immateriell, In
teresse_HdK, Interesse_3D,Erlebnis_Video,Erlebnis_Lesen,Erlebnis_Zuhoeren,Erlebnis_En
tdecken,Erlebnis_Interaktion,Erlebnis_Spielerisch, Grunddigital_Freude, Grunddigital_
Unterhaltung, Grunddigital_Erkenntnisse, Grunddigital_Faehigkeiten, Grunddigital_Schu
le, Grunddigital_Spezial, Grunddigital_Forschung, Grunddigital_Beruf, Grunddigital_We
ltweit))

#Ressourcen freigeben
rm(rawdata,headers,i,names,row.has.na,values,dataplot)

```

Formatierung der Werte und Auswahl des Datensatzes

Faktorielle Werte in numerische Werte umwandeln und die Werte auf eine Skala von 0 bis 1 skalieren.

```

df <- dataset2
to_int <- c("Geschlecht", "Alter")
df[to_int] <- lapply(datanew[to_int], as.integer)
rm(to_int)
X <- data.frame(lapply(df, function(x) {
  if (is.factor(x)) {
    x <- as.integer(x)
  }
  x
}))
X <- as.matrix(scale_dataset(df, type = "minmax"))

```

5. Clustering

Hyperparameter festlegen & SOM erstellen

Verschiedene "Einstellungen" für die Optimierung des Modells. Hierzu gehören zum Beispiel die Dimensionen des Grids, die Anordnung (hier hexagonal) und die Nachbarschaftsfunktion (hier Gauß'sche Nachbarschaftsfunktion).

```
# Set seed to ensure reproducibility
set.seed(222)

# Hyperparameter Grid
xdim <- 4
ydim <- 4

# Raster/ Matrix erstellen
# Hexagonales Grid & Gauß'sche Nachbarschaftsfunktion

som_grid <- somgrid(xdim = xdim, ydim = ydim, topo = "hexagonal", neighbourhood.fct =
"gaussian")
som_model <- som(X, grid = som_grid, alpha = c(0.05, 0.01), radius = 2.5)
som_numbers <- xdim*ydim
```

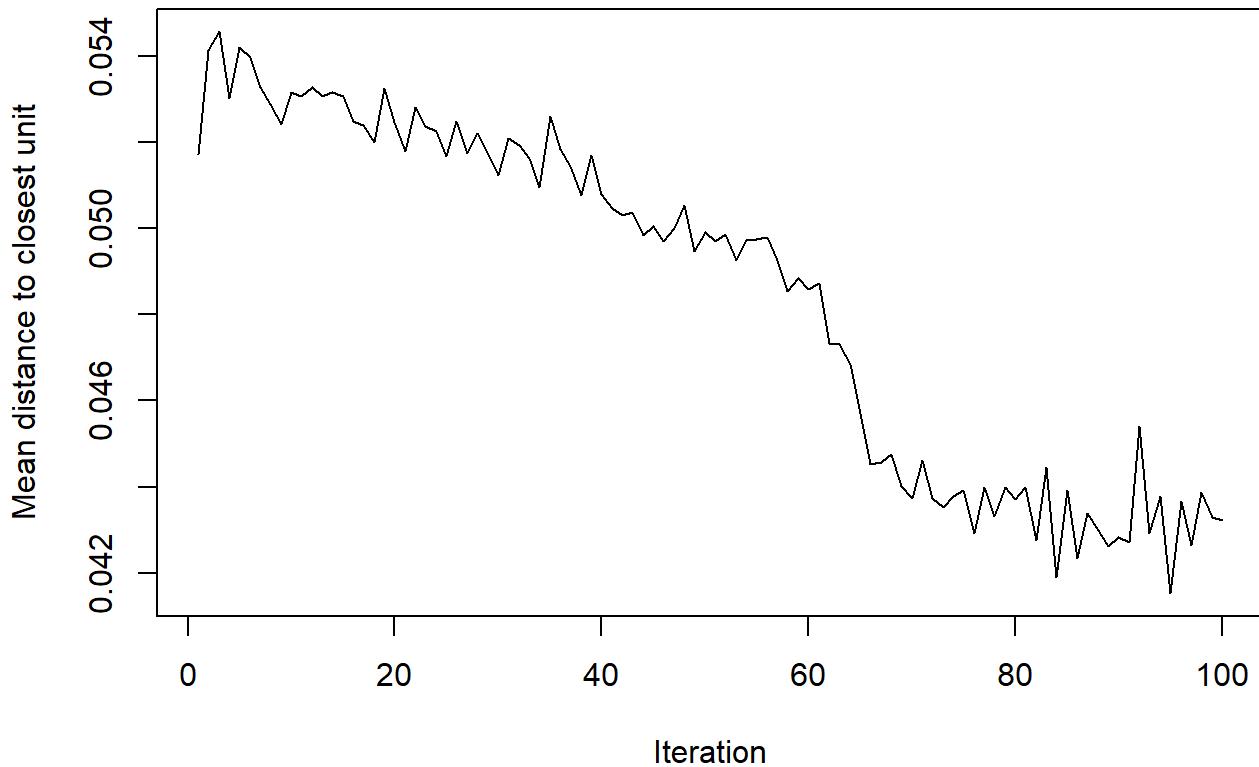
Qualitätskontrolle

Trainingsverlauf des SOM-Modells

Darstellung der durchschnittlichen (euklidischen) Distanz zwischen Input und Best-Matching-Unit über den Verlauf der Durchläufe (Iterationen)

```
plot(som_model, type = "changes")
```

Training progress



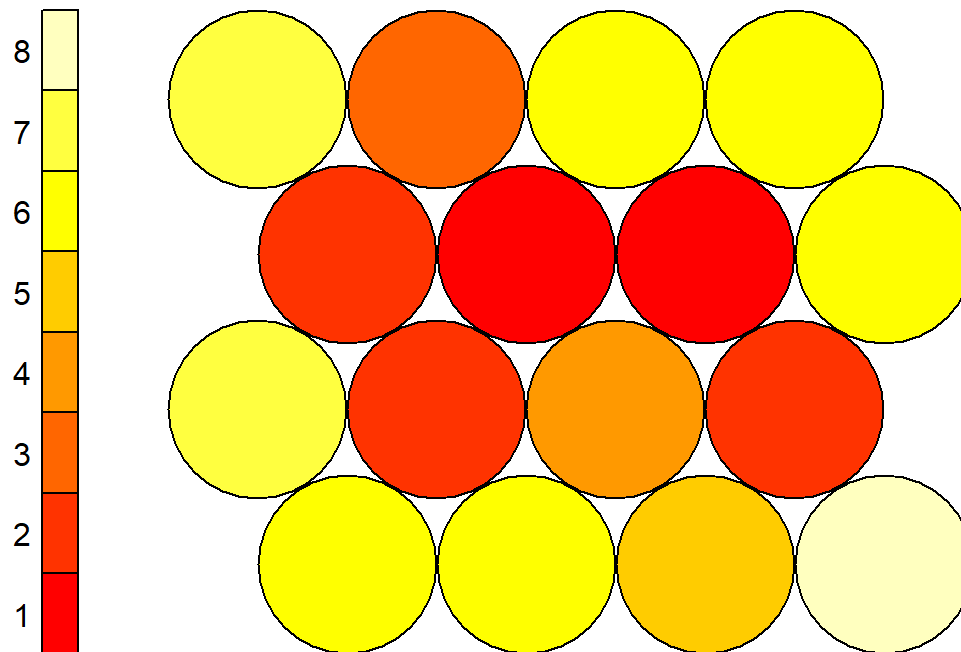
Man sieht deutlich wie die durchschnittliche Distanz mit steigender Anzahl Durchläufe abnimmt. Das Model wird bis zu einem bestimmten Punkt zunehmend genauer.

Anzahl der Beobachtungen per Grid-Unit

Darstellung der einzelnen Grid-Units in ihrer hexagonalen Anordnung. In diesem Plot ist die Anzahl an Personen dargestellt, welche der jeweiligen Grid-Unit zugeordnet wurden.

```
plot(som_model, type = "counts")
```

Counts plot



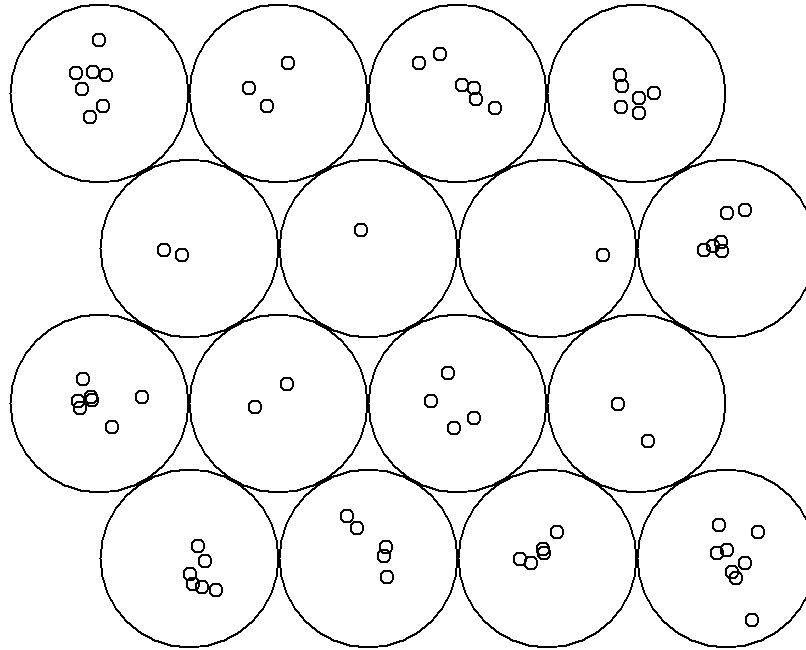
Man sieht deutlich, dass es Grid-Units gibt, zu denen keine Person zugeordnet wurden. Dies ist nicht ungewöhnlich. Bei den restlichen Grid-Units schwankt die Anzahl Personen zwischen 5 (rot) und 10 (sandgelb).

Anzahl der Beobachtungen per Grid-Unit

Darstellung der Zuordnung der einzelnen Beobachtungen (Personen) zu der Best-Matching-Unit (BMU). In der Konsole wird zusätzlich ausgegeben welche Person zu welcher Grid-Unit zugeordnet wurde.

```
plot(som_model, type = "mapping")
```

Mapping plot



```
som_model$unit.classif
```

```
## [1] 7 2 8 1 16 4 16 16 15 13 13 12 3 15 14 13 13 13 8 9 12 4 3 4 3
## [26] 1 4 2 14 2 14 2 5 15 9 16 5 16 1 2 6 3 6 12 5 4 1 4 5 5
## [51] 15 5 12 13 11 7 13 7 5 15 15 1 4 1 10 12 3 7 16 12 4 2
```

```
table(som_model$unit.classif) # The total number of assigned input samples per each g
rid unit
```

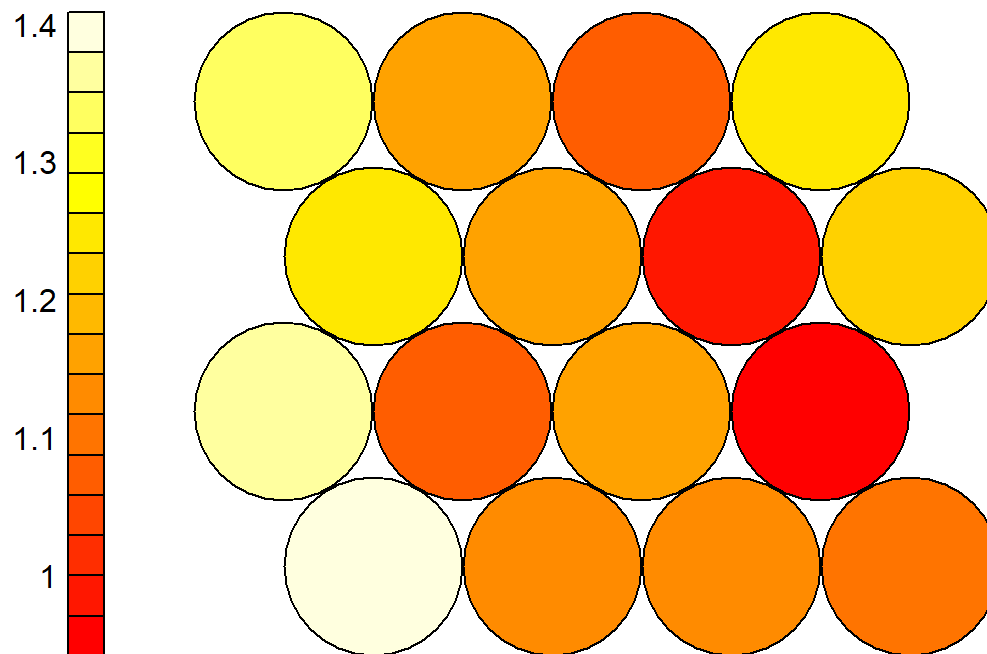
```
##
## 1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 16
## 6 6 5 8 7 2 4 2 2 1 1 6 7 3 6 6
```

Distanzen zwischen den Grids

Darstellung der (aufsummierten) Distanzen zwischen den einzelnen Grid-Units und ihren direkten Nachbarn.
(U-matrix plot)

```
plot(som_model, type = "dist.neighbours", main = "SOM neighbour distances")
```

SOM neighbour distances

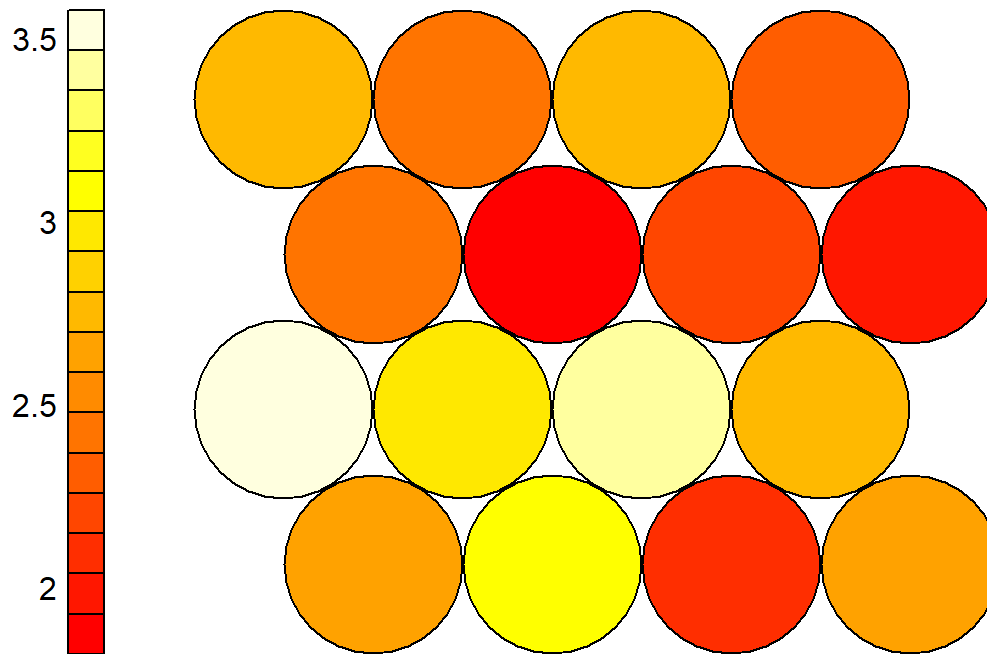


Streuung der zugeordneten Beobachtungen zu dem Codebook-Vektor

Dieser Plot zeigt wie weit die Beobachtungen (Personen) von den Grid-Units (genauer: den Codebook-Vektoren der Grid-Units), zu welchen sie zugeordnet wurden, entfernt sind. Dies ist in einfachen Worten die Streuung der Personen um den "Mittelpunkt" der Grid-Unit und zeigt wie gut die Personen von der Grid-Unit repräsentiert werden. Je kleiner die Distanz, desto genauer die Zuordnung.

```
plot(som_model, type = "quality")
```

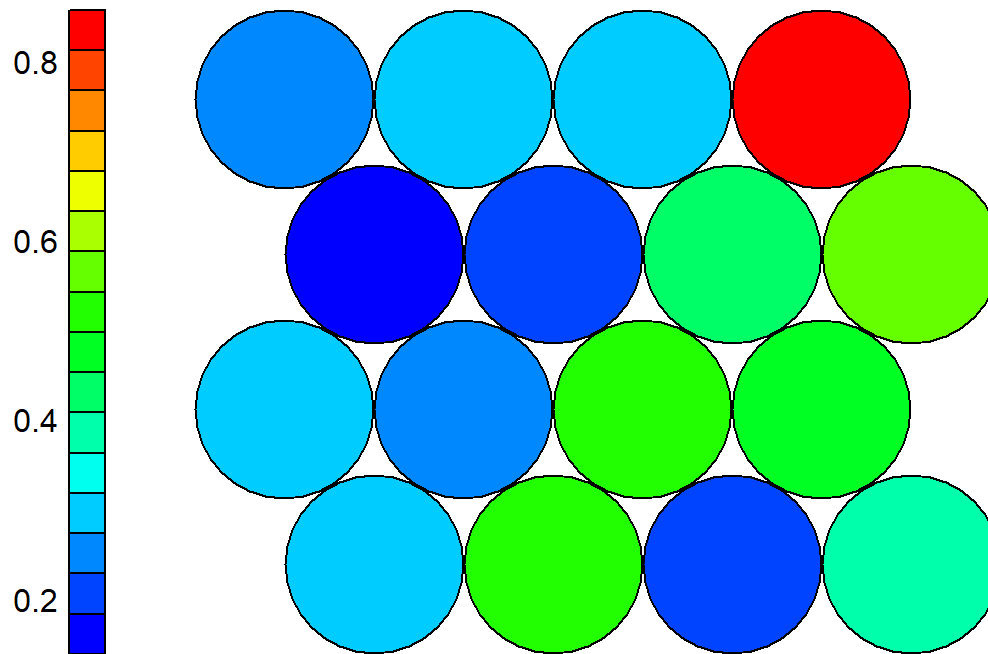
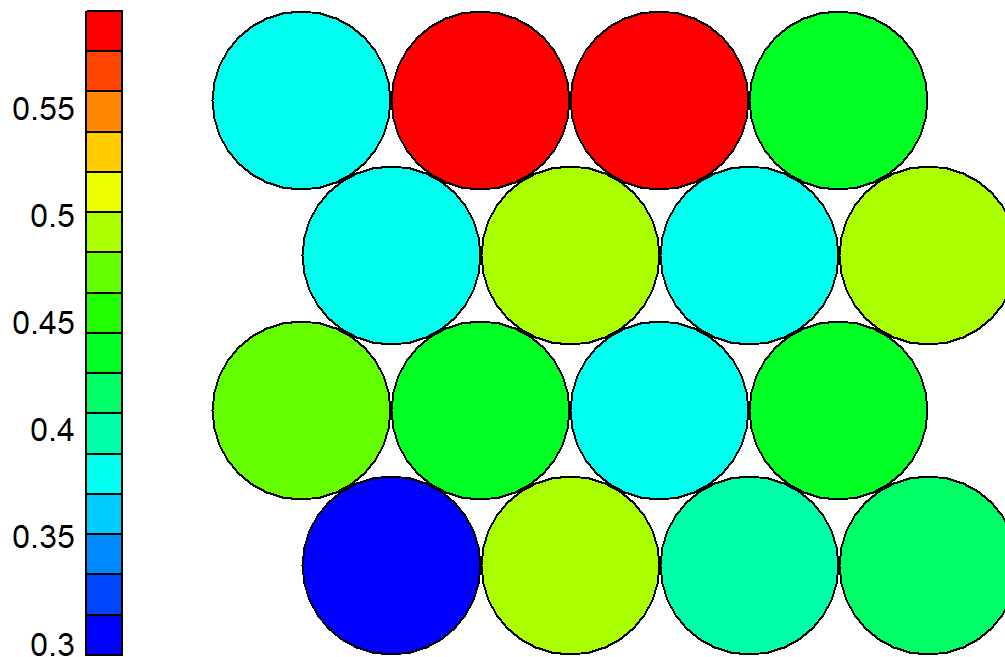
Quality plot



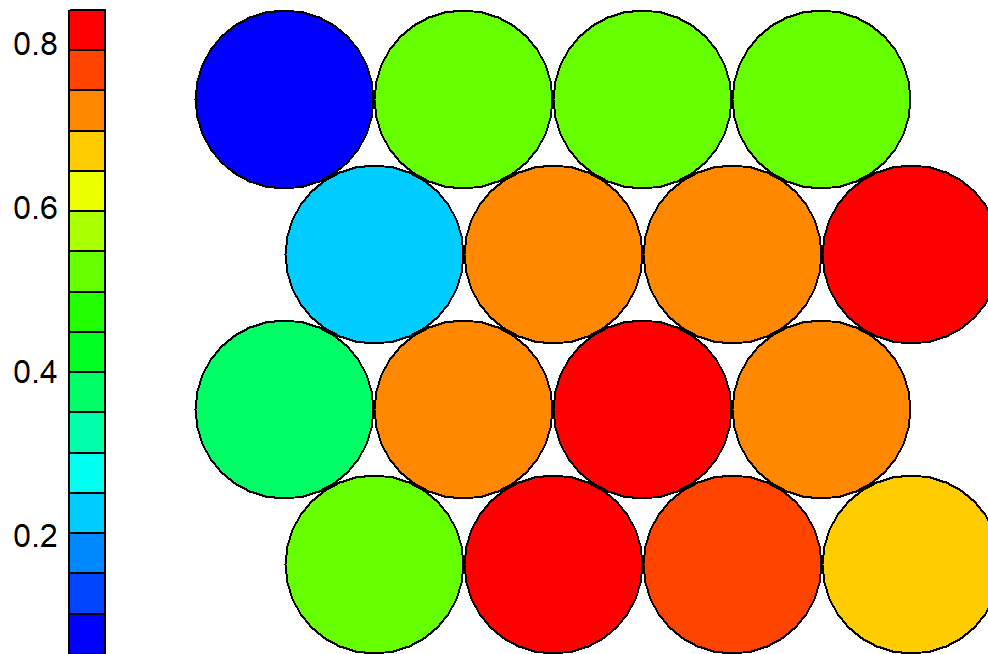
Höhe der Ausprägungen einzelner Variablen pro Grit-Unit

Für jedes Merkmal gibt es einen einzelnen Plot. Die Plots zeigen die Verteilung der Ausprägungen des Merkmals zwischen den Grid-Units. Beispiel Geschlecht: Die oberen zwei Reihen an Grid-Units haben eine niedrige Ausprägung. Zu diesen Grid-Units wurden (vor allem) Frauen zugeordnet. In der dritten Zeile lässt sich keine eindeutige Aussage treffen, während in der letzten Zeile die Ausprägungen hoch sind. Dies lässt auf Männer schließen. (Für die Bedeutung der Ausprägungen siehe Präsentation (Kapitel "Visualisierung"))

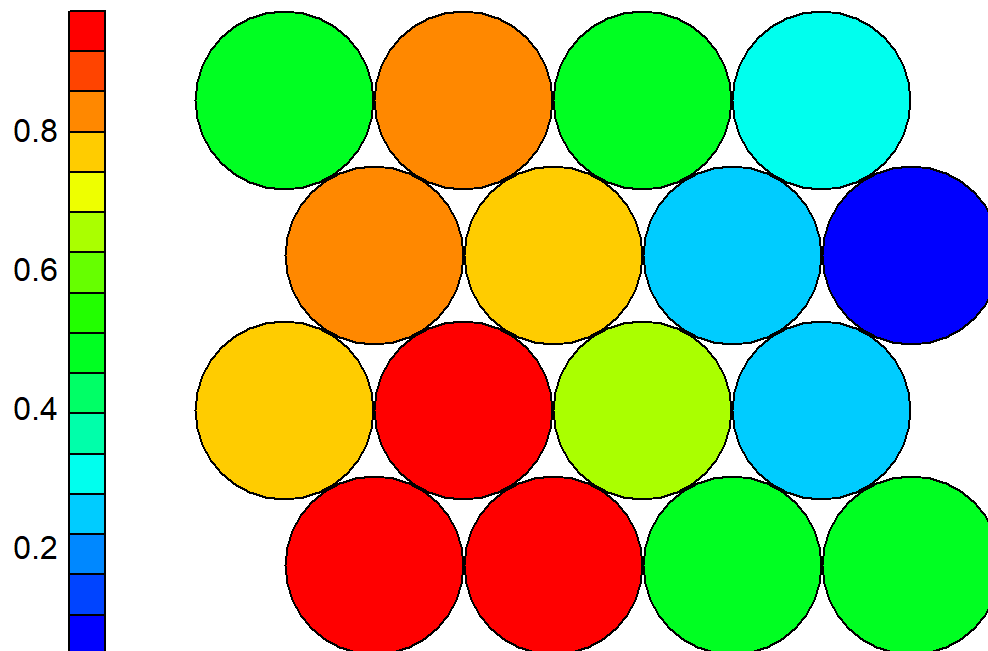
```
coolBlueHotRed <- function(n, alpha = 1) { rainbow(n, end = 4 / 6, alpha = alpha)[n:
1] }
xcodes <- getCodes(som_model)
xnames <- colnames(X)
for (i in seq_along(data.frame(X))) {
  plot(som_model, type = "property", property = xcodes[, i],
      main = paste0("Heatmap: ", xnames[i]), palette.name = coolBlueHotRed)
}
```

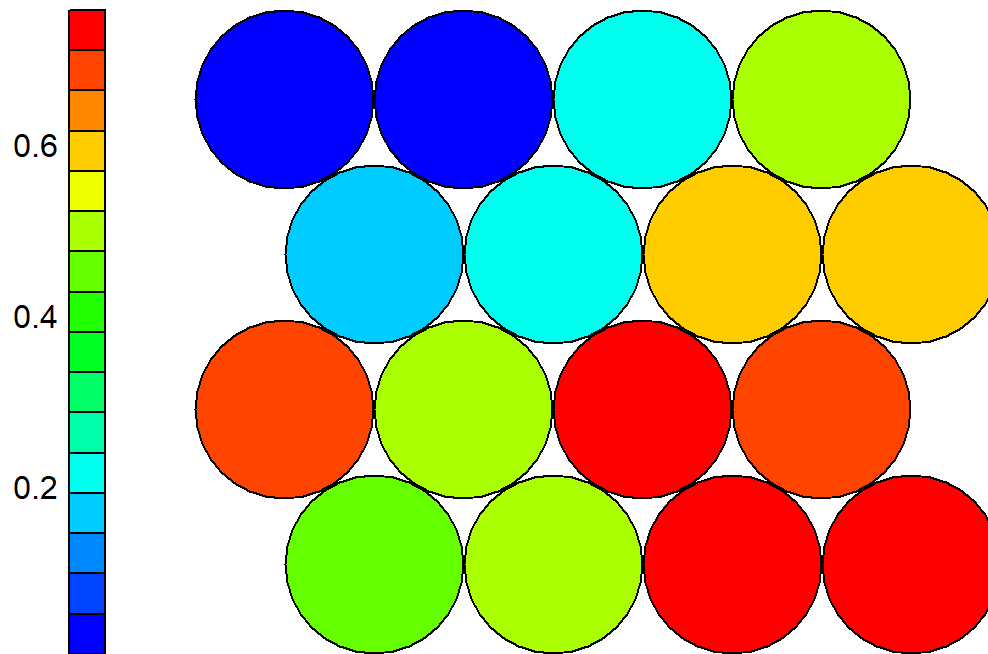
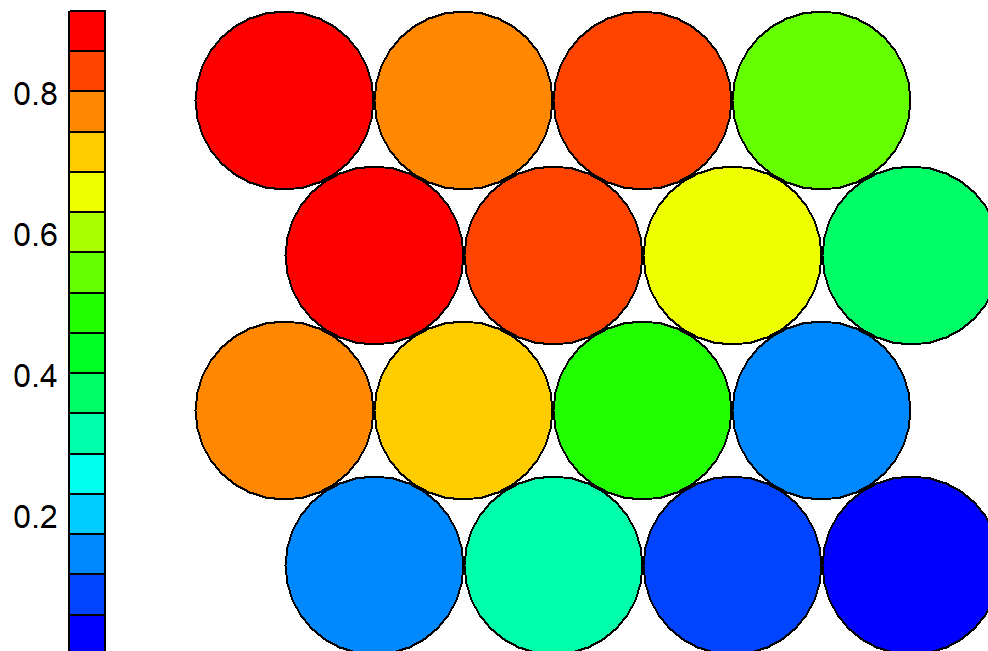
Heatmap: Geschlecht**Heatmap: Alter**

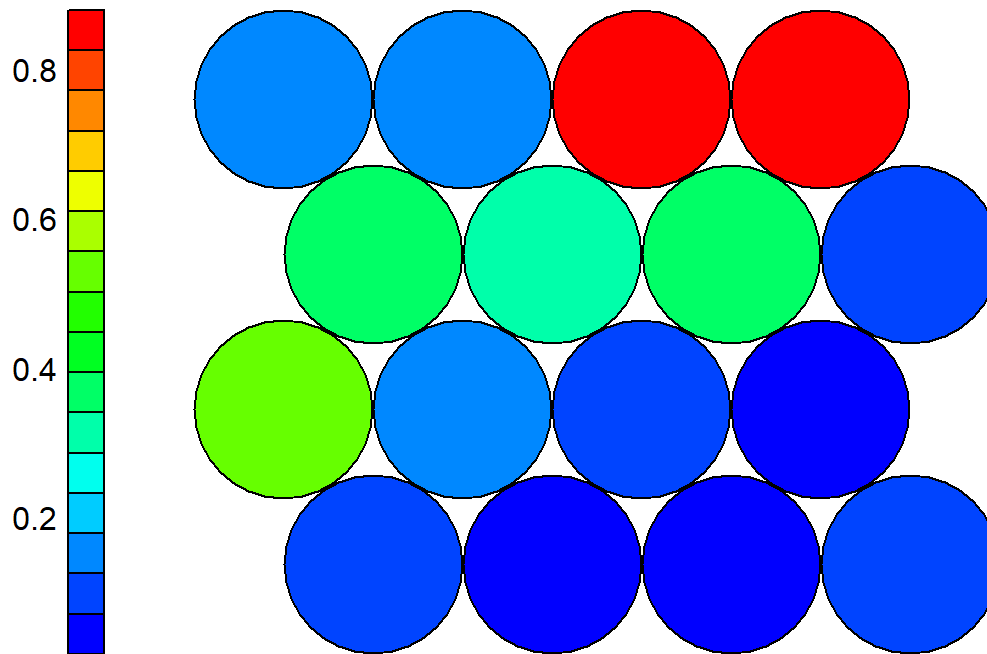
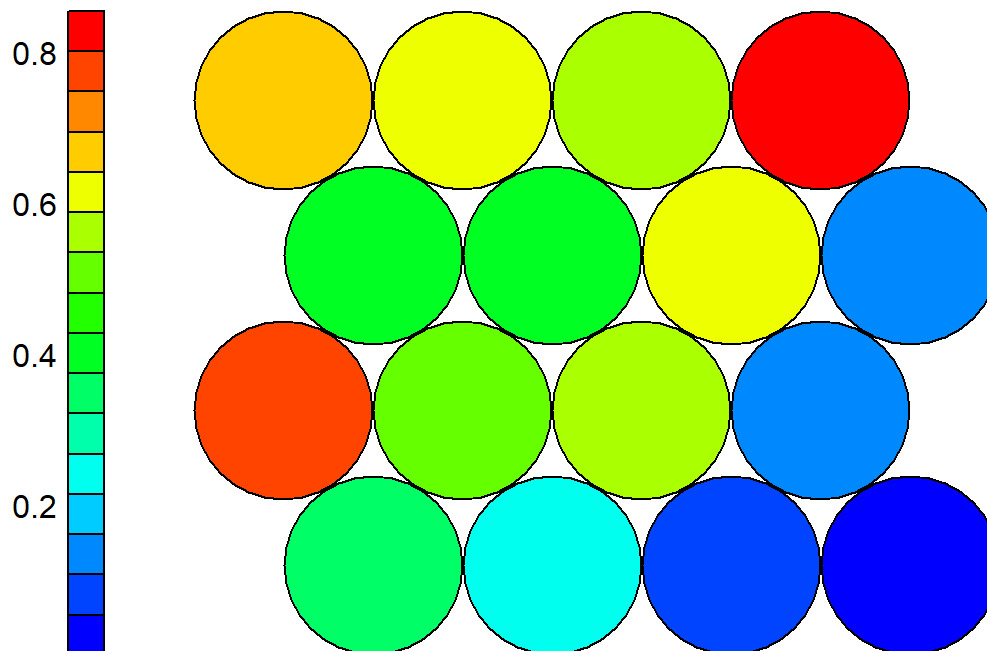
Heatmap: Interesse_Archaeologie

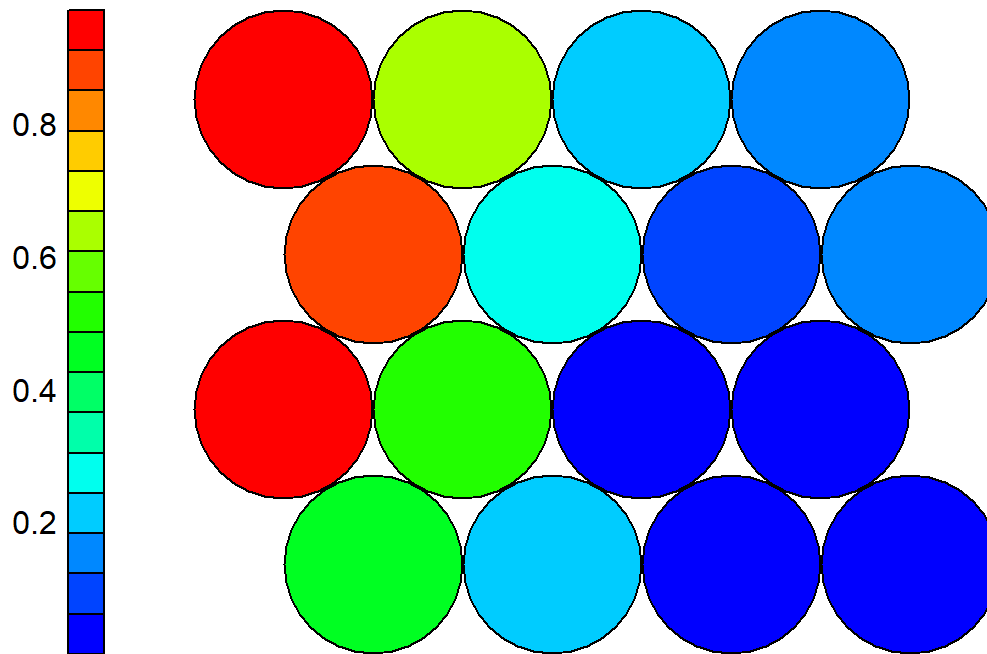
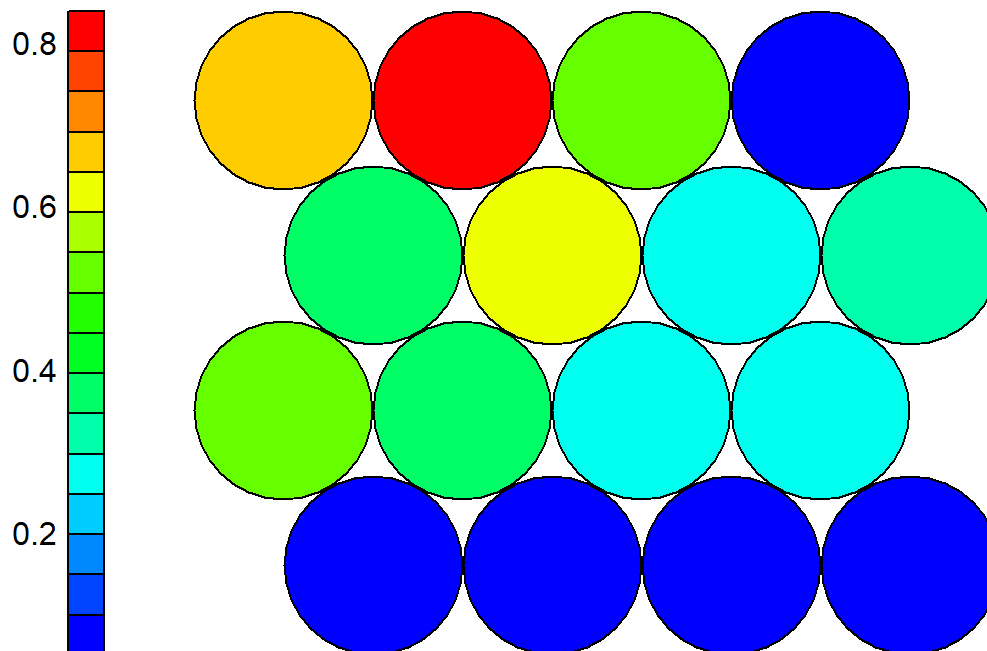


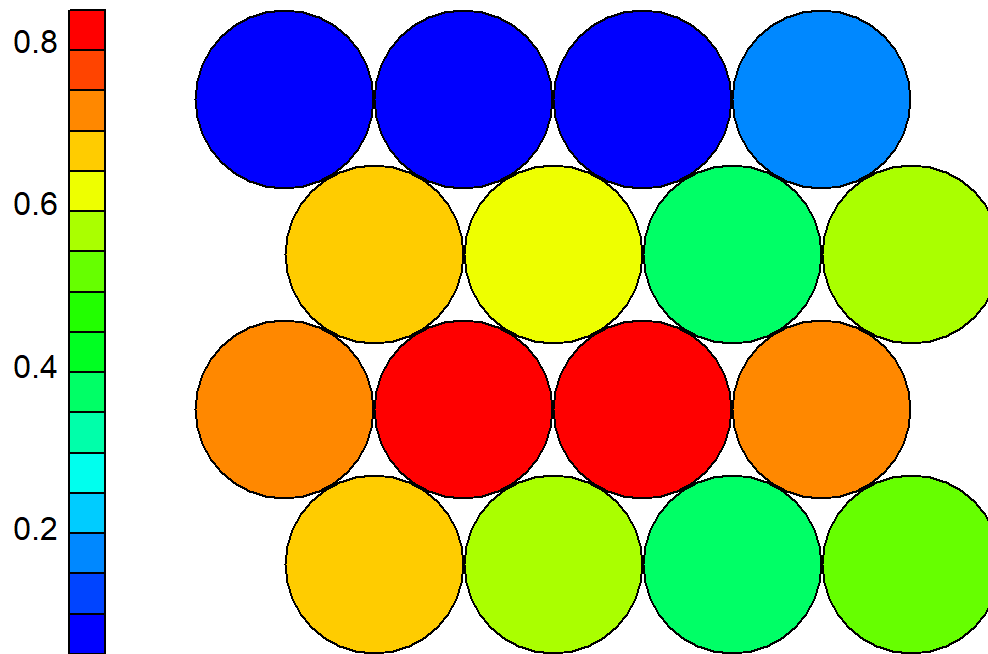
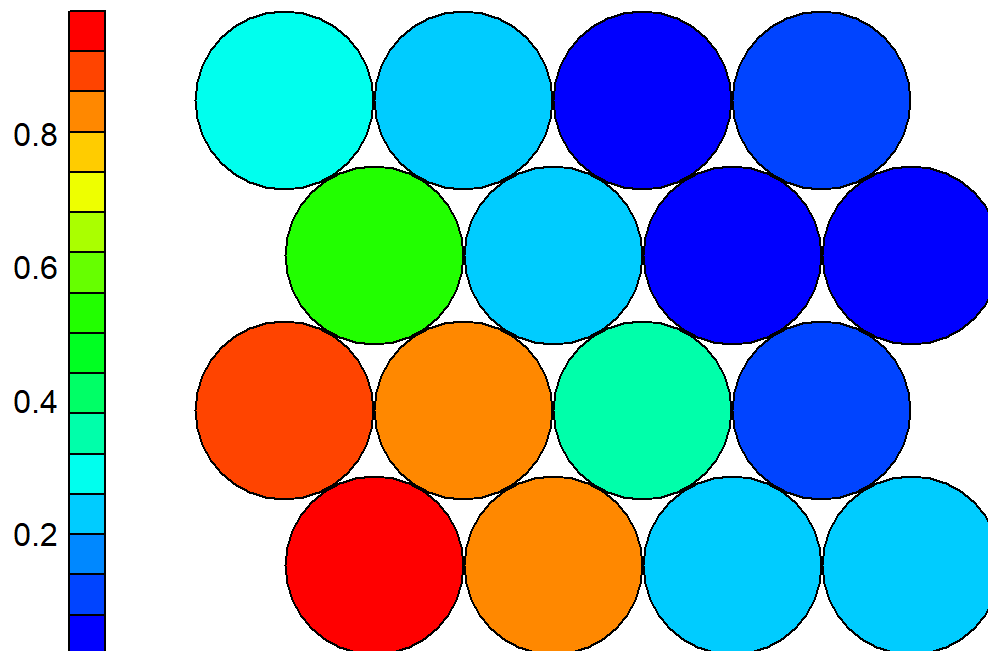
Heatmap: Interesse_Weltkultur

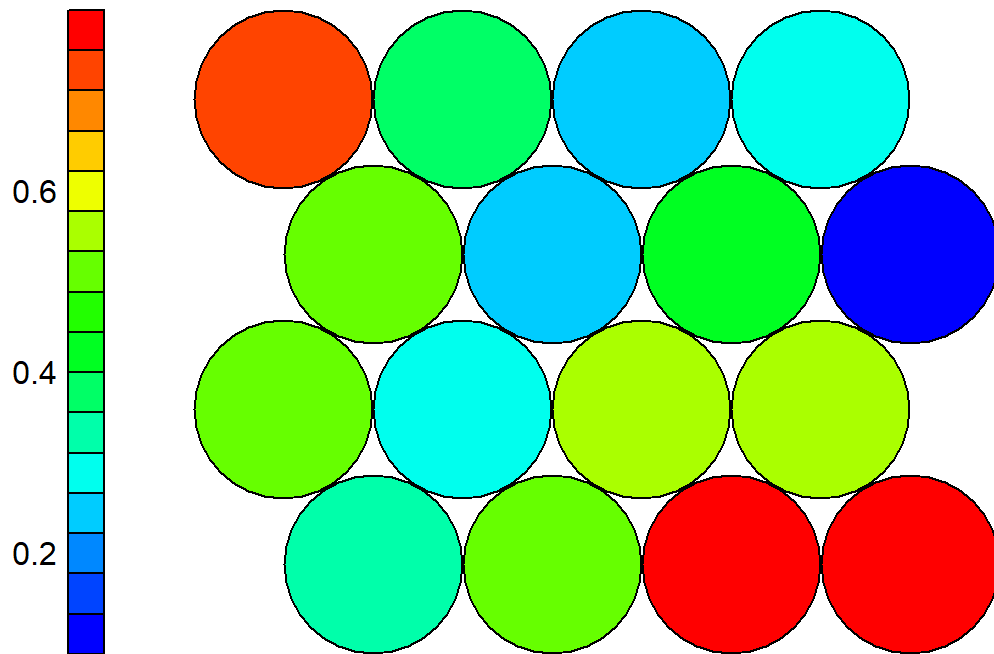
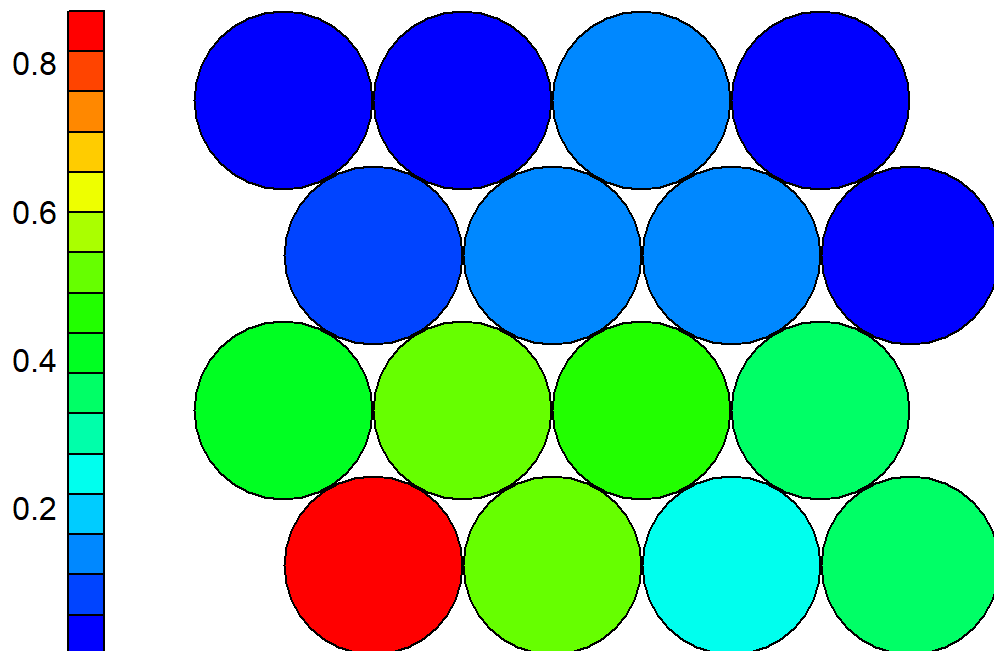


Heatmap: Interesse_Mittelalter**Heatmap: Interesse_Design**

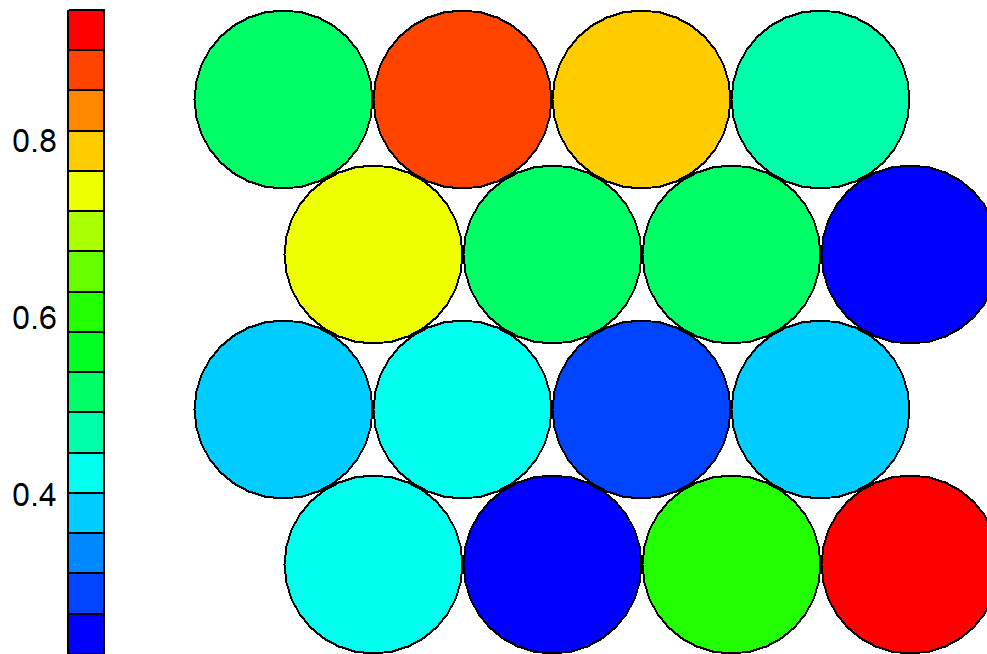
Heatmap: Interesse_Baden**Heatmap: Interesse_1900**

Heatmap: Interesse_Aktuell**Heatmap: Interesse_Musik**

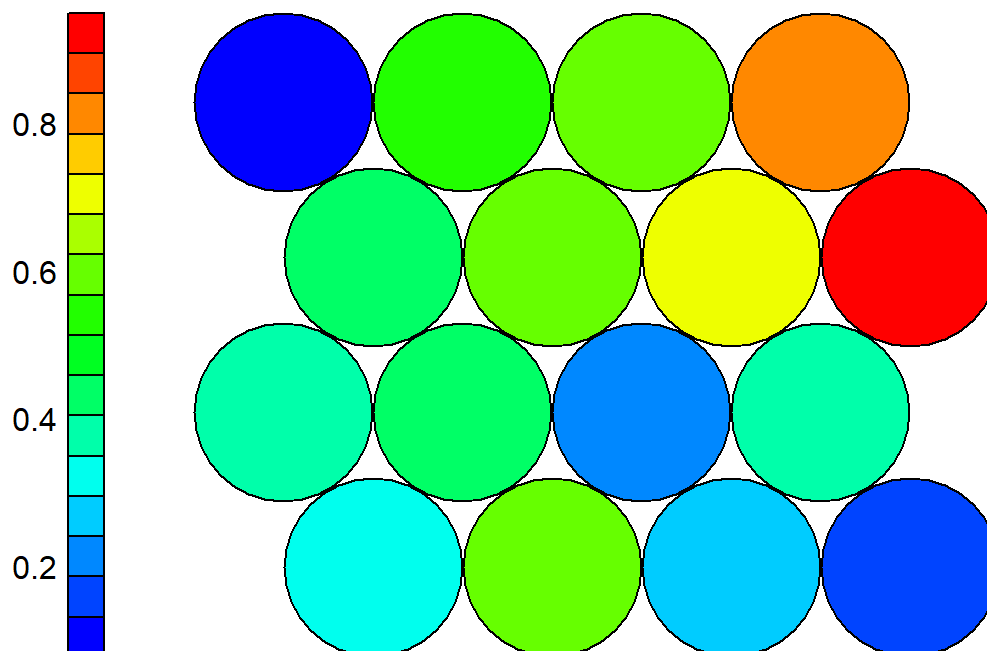
Heatmap: Interesse_Materiell**Heatmap: Interesse_Immateriell**

Heatmap: Interesse_HdK**Heatmap: Interesse_3D**

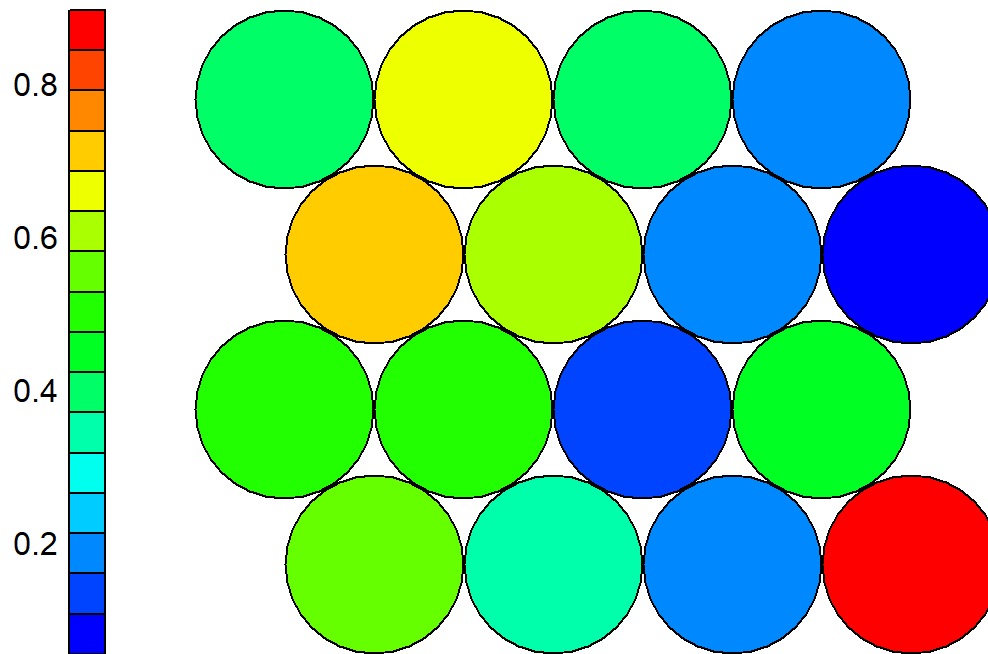
Heatmap: Erlebnis_Video



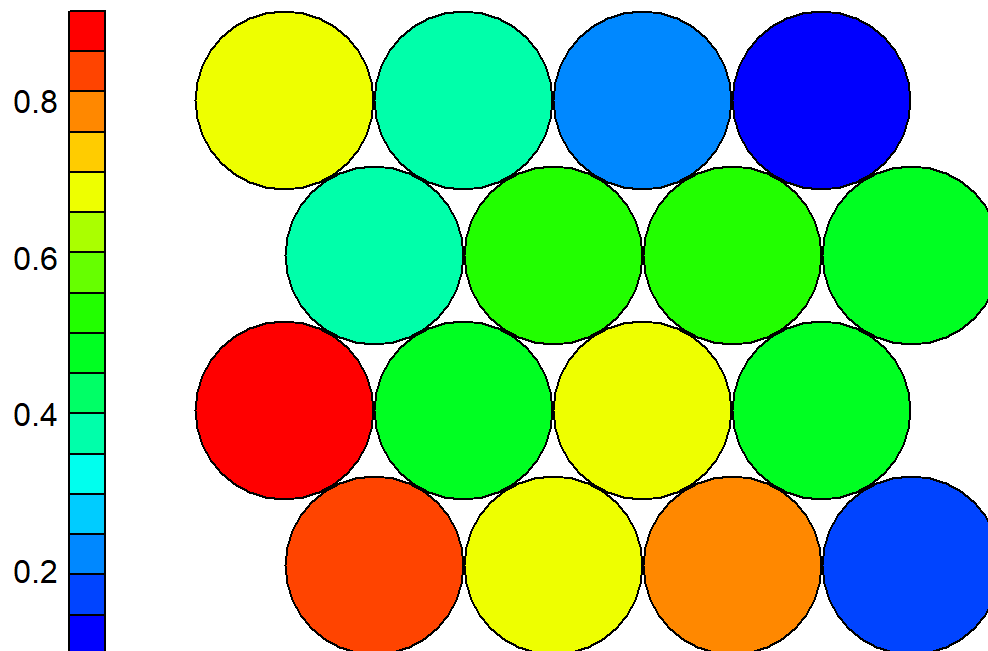
Heatmap: Erlebnis_Lesen



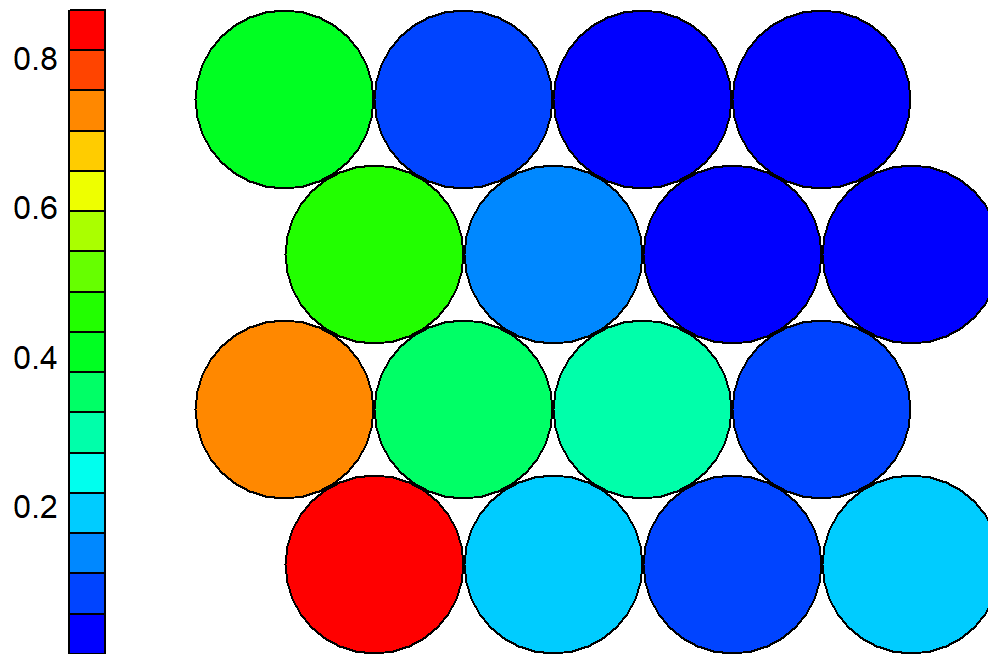
Heatmap: Erlebnis_Zuhoeren



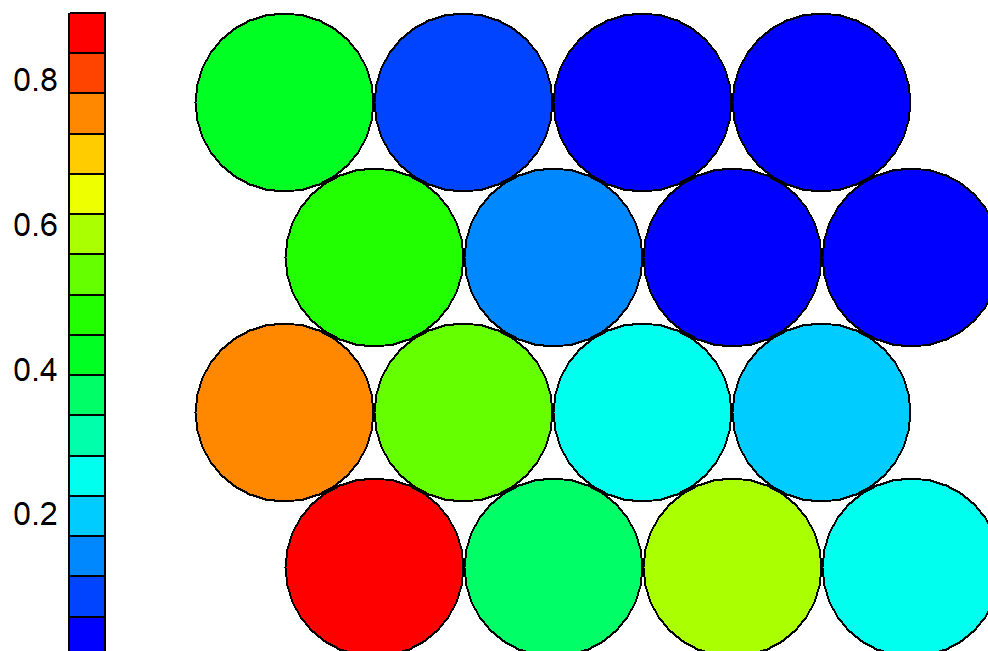
Heatmap: Erlebnis_Entdecken



Heatmap: Erlebnis_Interaktion



Heatmap: Erlebnis_Spielerisch

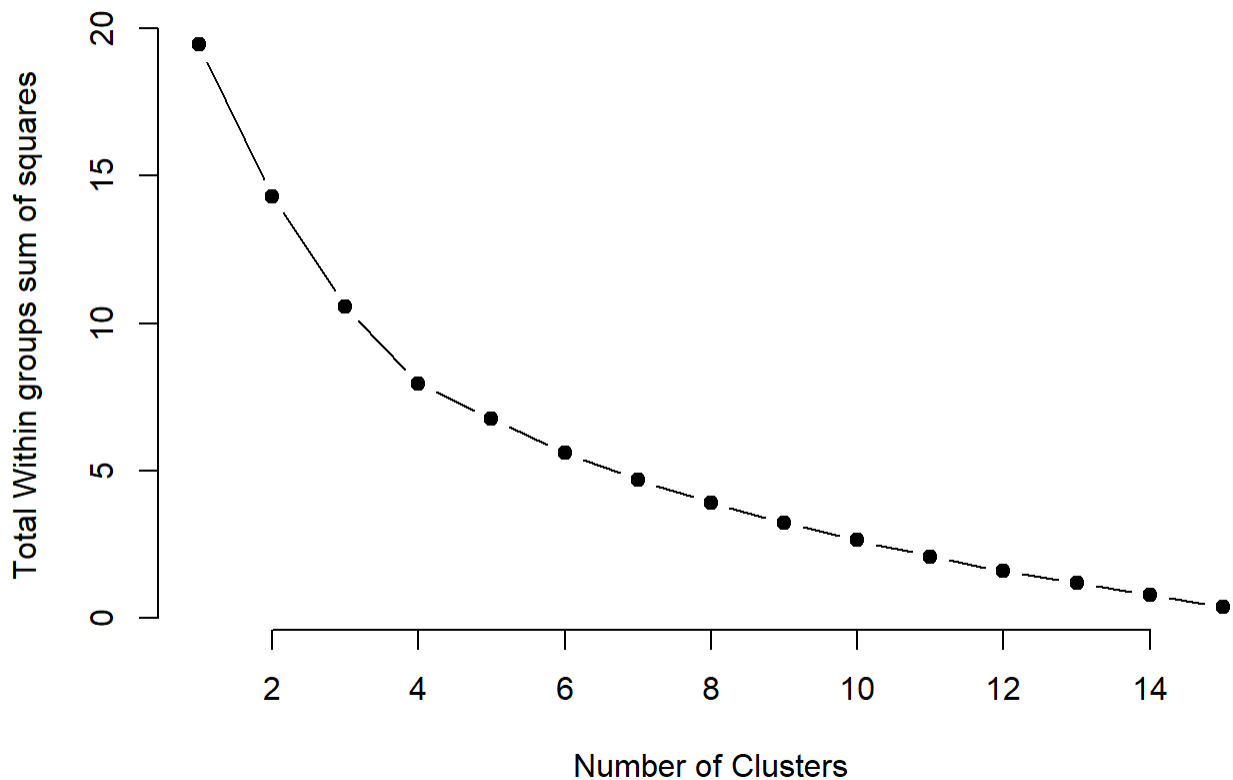


Optimale Clusteranzahl bestimmen

Die Optimale Anzahl an Cluster kann mit verschiedenen Methoden bestimmt werden. Drei von ihnen werden im Folgenden angewendet, aber nicht weiter in die Ergebnisse implementiert. Bei den Methoden handelt es sich um: - Elbow Curve - Silhouette Score - Gap Statistic

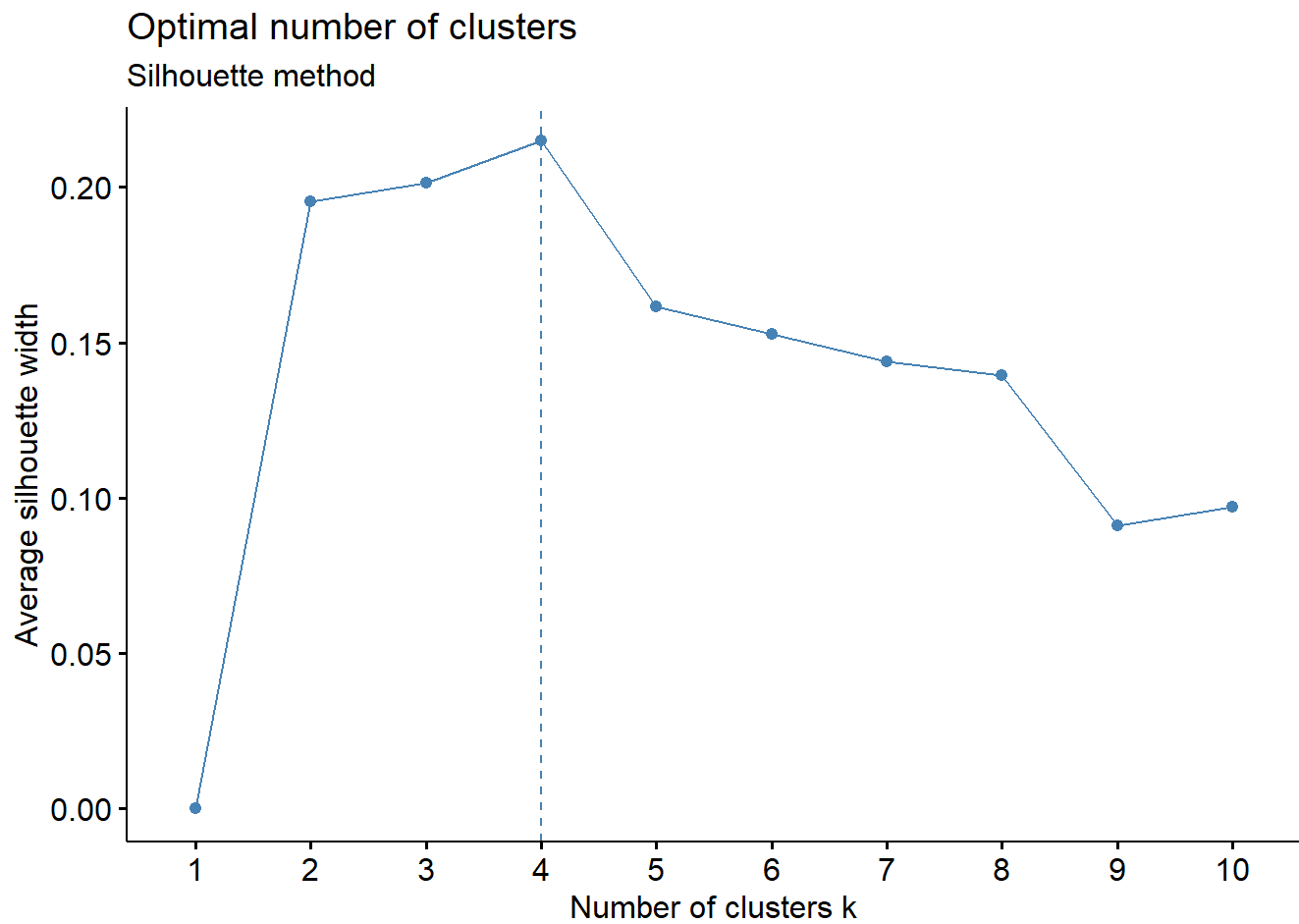
```
# Elbow curve
# See where the curve flattens. After this point there's no more explanation of most
of the variance in data.
set.seed(20)
clusterdata <- getCodes(som_model)

kmax <- 15L
wss <- sapply(1L:kmax, function(k) { kmeans(clusterdata, k, nstart = 20L)$tot.withins
s })
plot(1L:kmax, wss, type = "b", pch = 19L, frame = FALSE,
     xlab = "Number of Clusters", ylab = "Total Within groups sum of squares")
```

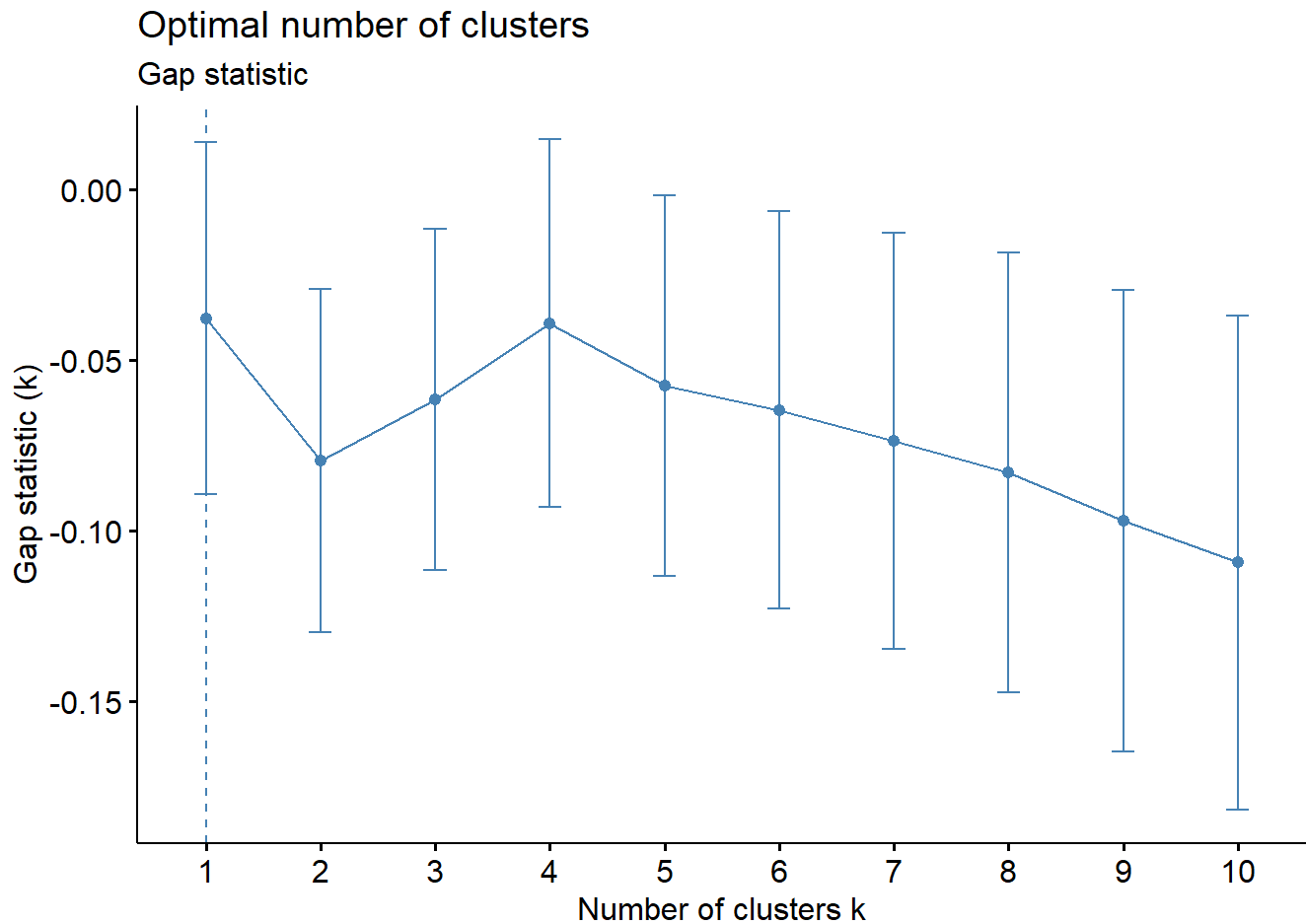


```
# Silhouette score

require(factoextra)
fviz_nbclust(clusterdata, kmeans, method = "silhouette") +
  labs(subtitle = "Silhouette method")
```



```
# Gap statistic  
# https://statweb.stanford.edu/~gwalther/gap  
  
require(factoextra)  
set.seed(123)  
fviz_nbclust(clusterdata, kmeans, nstart = 25L, method = "gap_stat", nboot = 500L) +  
  labs(subtitle = "Gap statistic")
```



Cluster Visualisierung

Für die Visualisierung der Ergebnisse gibt es verschiedene Möglichkeiten. Im folgenden Werden zwei Möglichkeiten angewandt und kurz vorgestellt.

1. Möglichkeit: Fan-Diagramme

Darstellung der Grid-Units in ihrer hexagonale Anordnung. Es wurde je ein Diagramm für die Anzahl Cluster 1 bis 10 erstellt. Die Farbe der Grid-Units repräsentiert das Cluster zu dem diese Grid-Unit (und damit auch die ihr zugeordneten Personen) zugeordnet wurden. Die Fan-Diagramme innerhalb der Grid-Units zeigen die Ausprägungen der Merkmale (also des Codebook-Vektors). Ist ein Fan klein ist dessen Ausprägung gering. Da die Werte auf 0 bis 1 skaliert wurden ist eine hohe Ausprägung 1 und eine geringe Ausprägung 0. (Die Bedeutung der Ausprägung entnehmen Sie bitte der Legende). Innerhalb eines Clusters kann so nach Überschneidungen und Merkmalskombinationen gesucht werden.

```
# Form clusters on grid
# Try several cluster algorithms and different numbers of clusters k

max_cluster <- 10L
clusterdata <- getCodes(som_model)

# Capture outputs within a list structure
som_clusters <- list(
  model = list(),
  kmeans = list(),
  hierarchical = list()
)

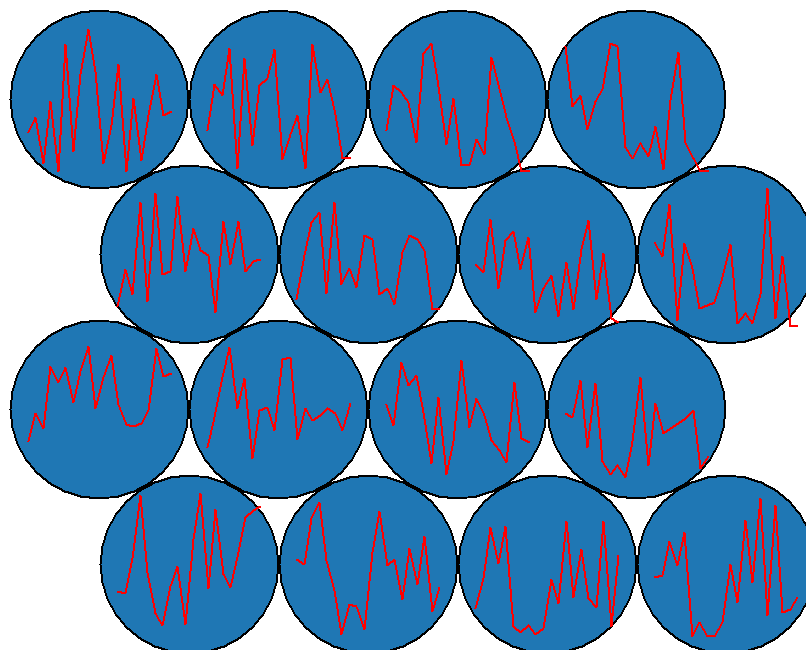
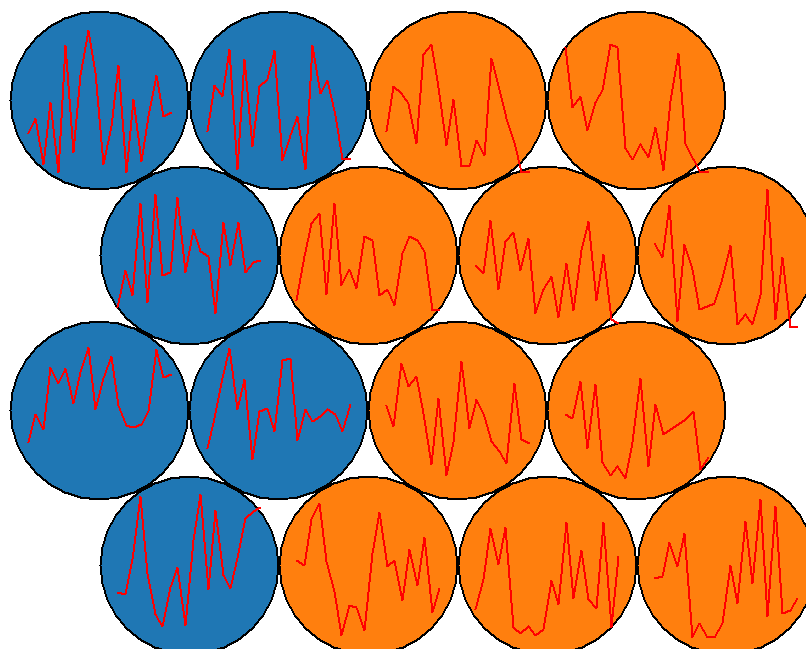
for (k in seq_len(max_cluster)) {

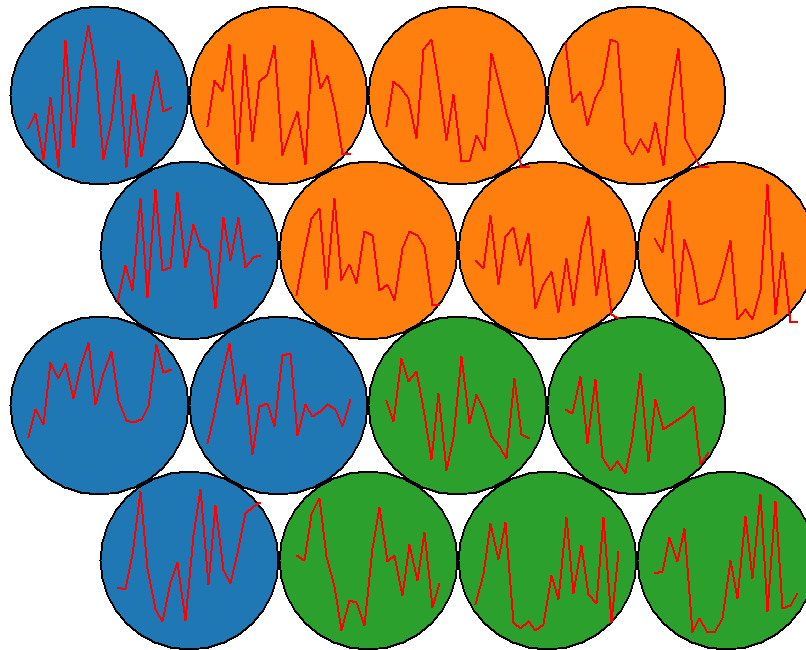
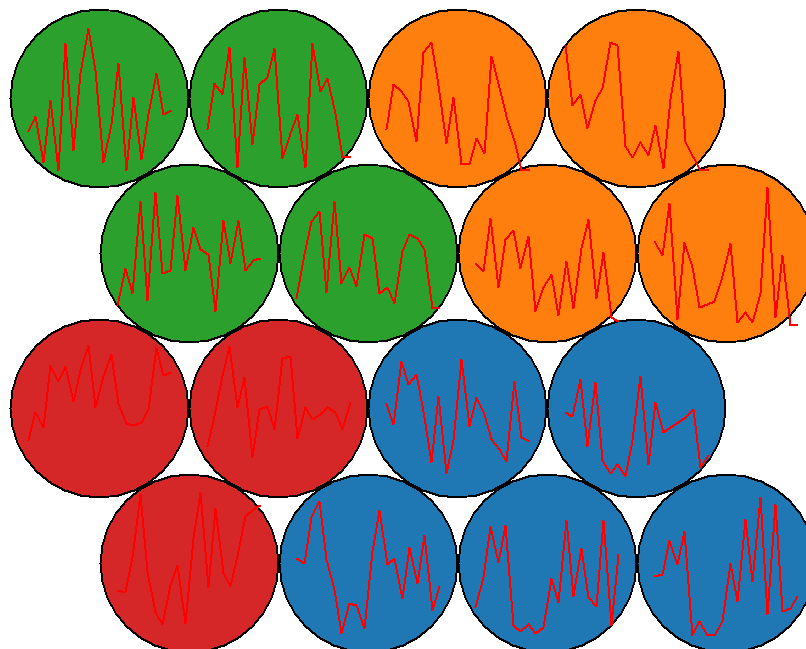
  ## k-means clustering
  som_cluster_kmeans <- kmeans(clusterdata, centers = k, iter.max = 100L, nstart = 1
0)$cluster
  som_clusters$kmeans[[toString(k)]] <- som_cluster_kmeans

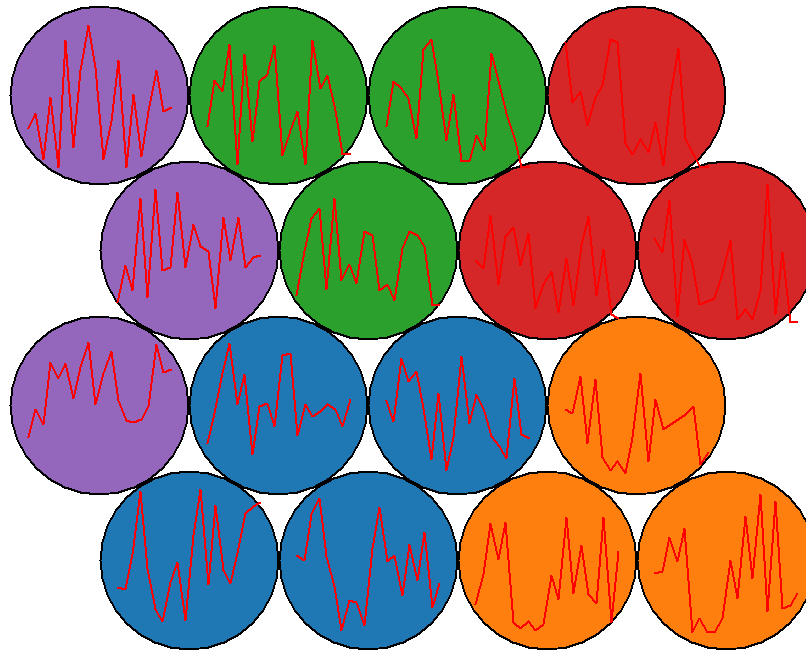
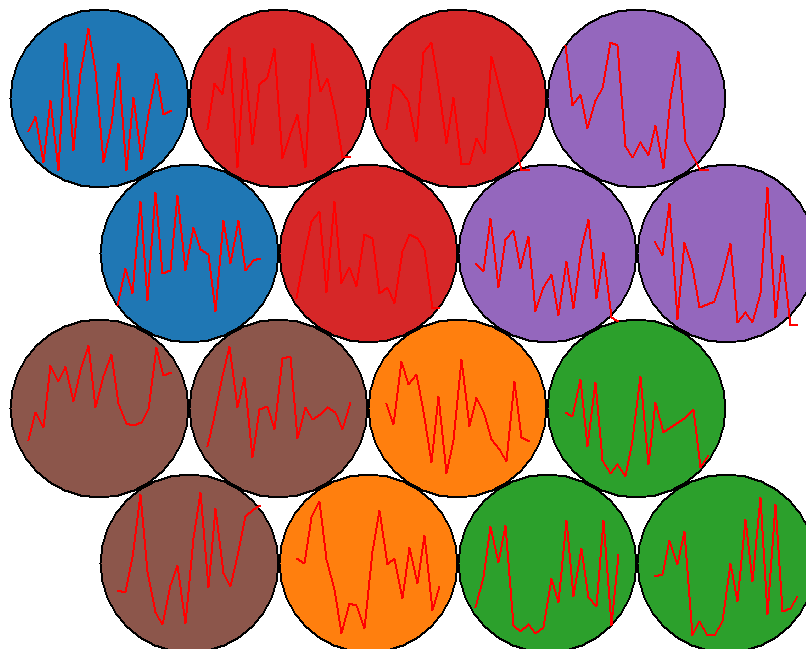
  ## Hierarchical clustering
  som_cluster_hierarchical <- cutree(hclust(dist(clusterdata)), k = k)
  som_clusters$hierarchical[[toString(k)]] <- som_cluster_hierarchical
}

# Plot clusters
rgb_palette <- c("#1f77b4", "#ff7f0e", "#2ca02c", "#d62728", "#9467bd", "#8c564b", "#
e377c2", "#33245a", "#ca4455", "#bf1123")
plotSOM <- function(clusters, title) {
  plot(som_model, type = "codes", bgscol = rgb_palette[clusters], keepMargins = F, col
= NA, main = title)
}

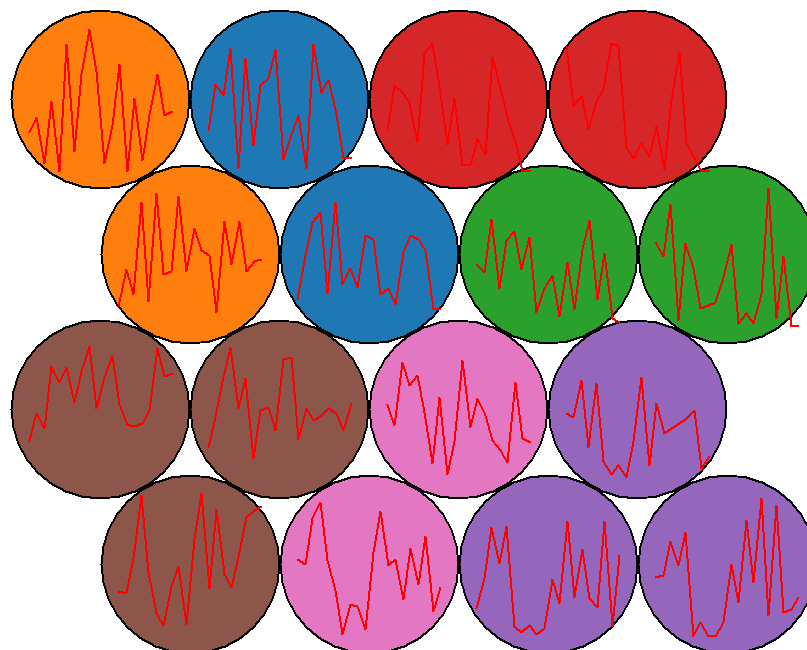
for (k in seq_len(max_cluster)) plotSOM(som_clusters$kmeans[[toString(k)]], paste0("k
means: ", k))
```

kmeans: 1**kmeans: 2**

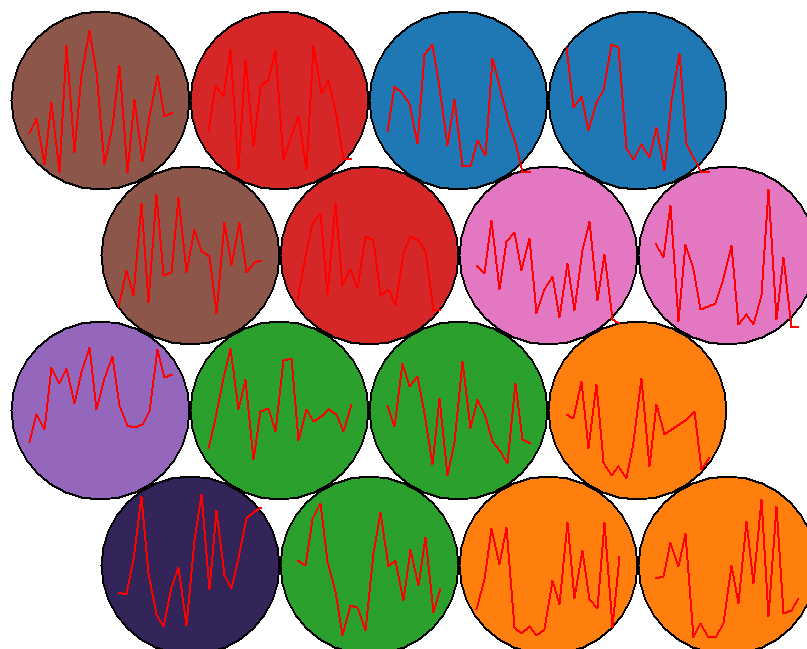
kmeans: 3**kmeans: 4**

kmeans: 5**kmeans: 6**

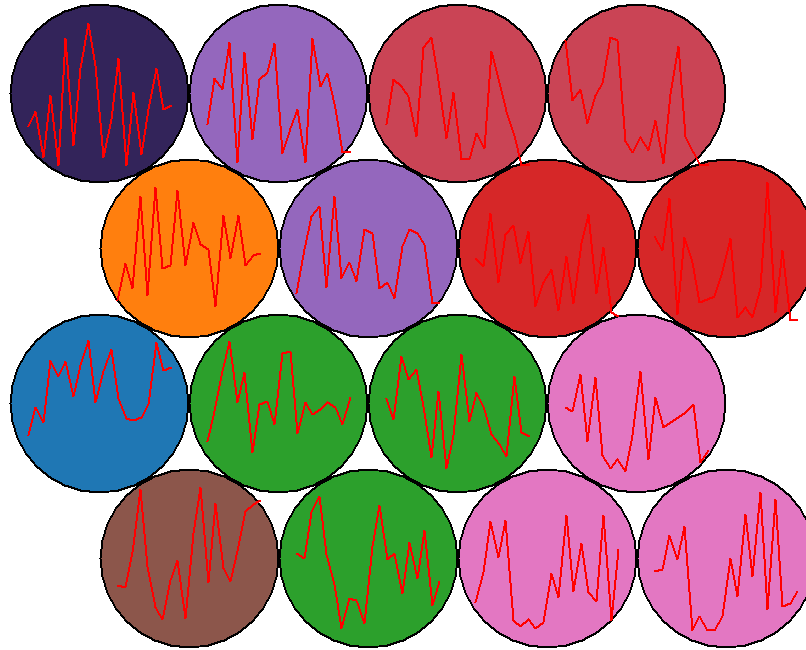
kmeans: 7



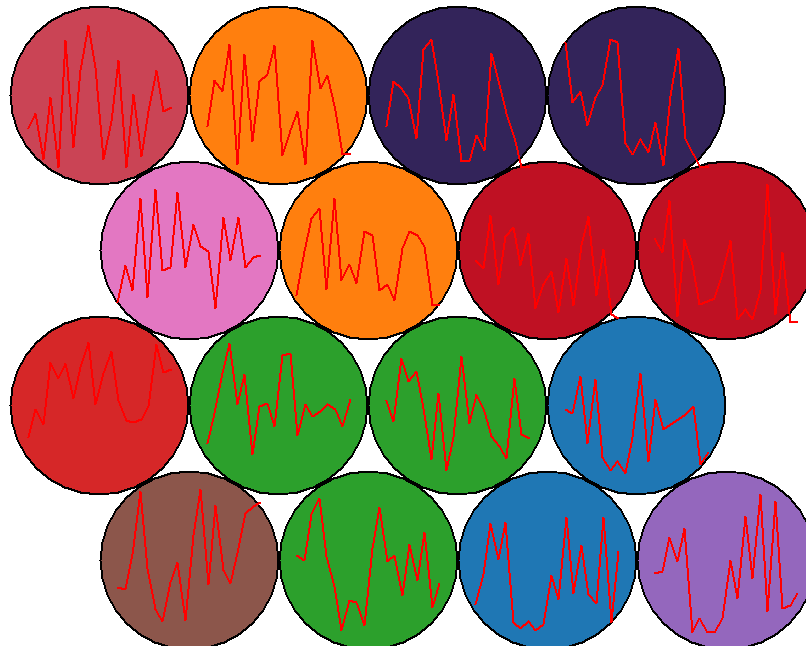
kmeans: 8



kmeans: 9

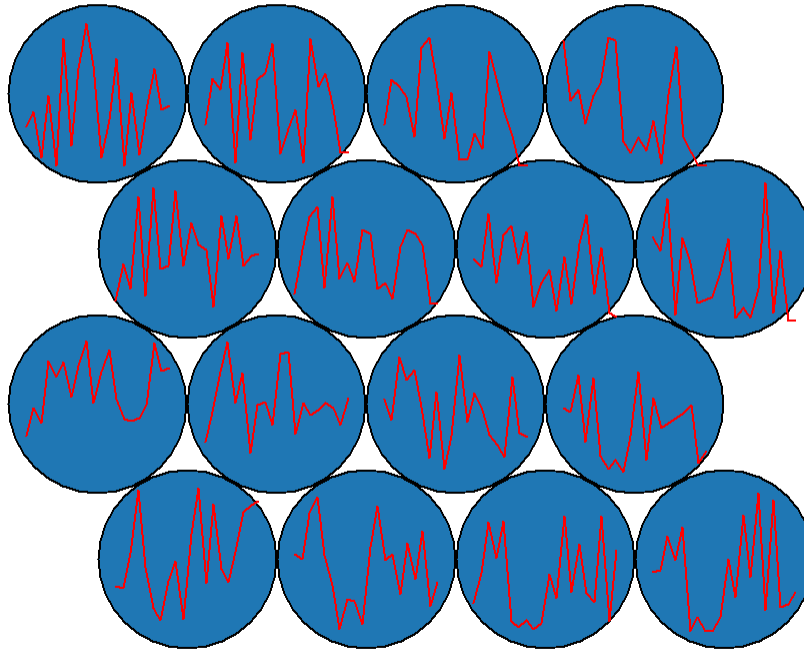


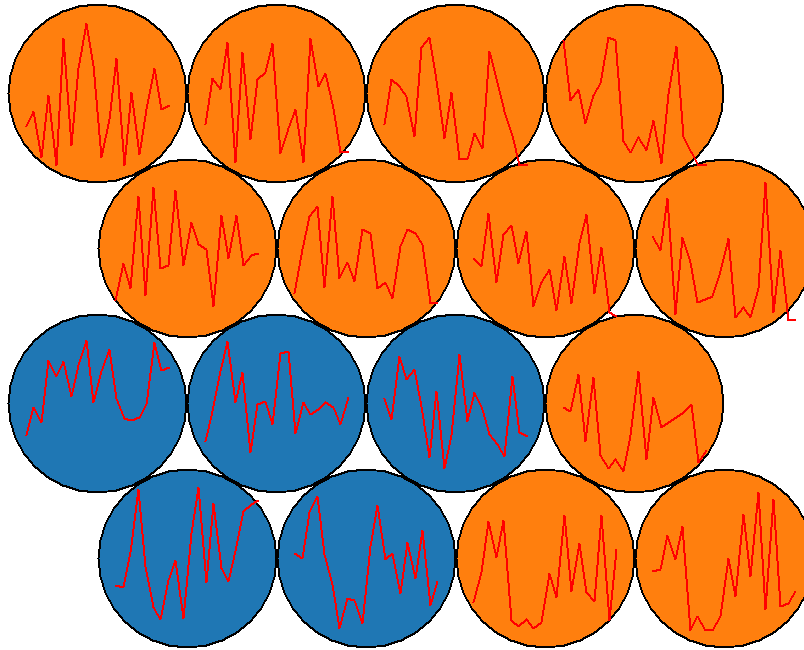
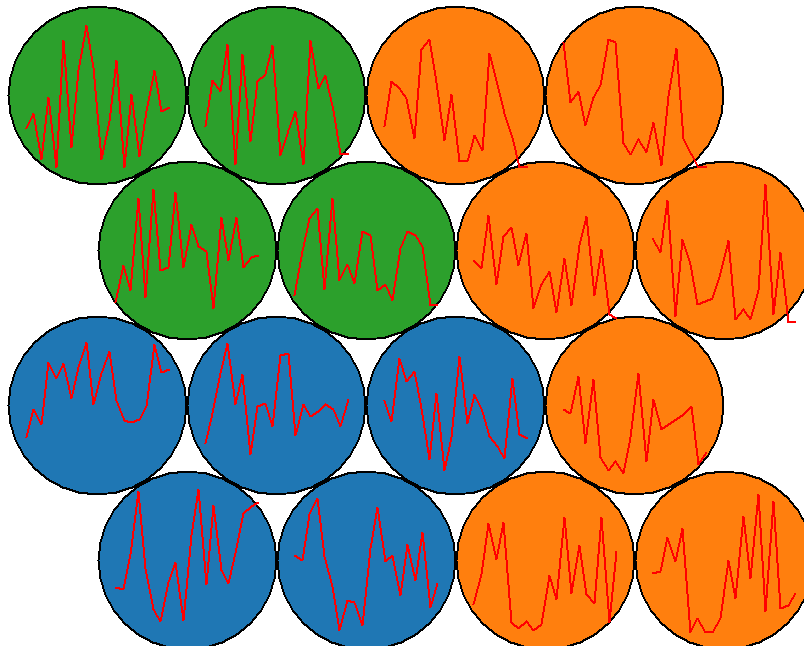
kmeans: 10

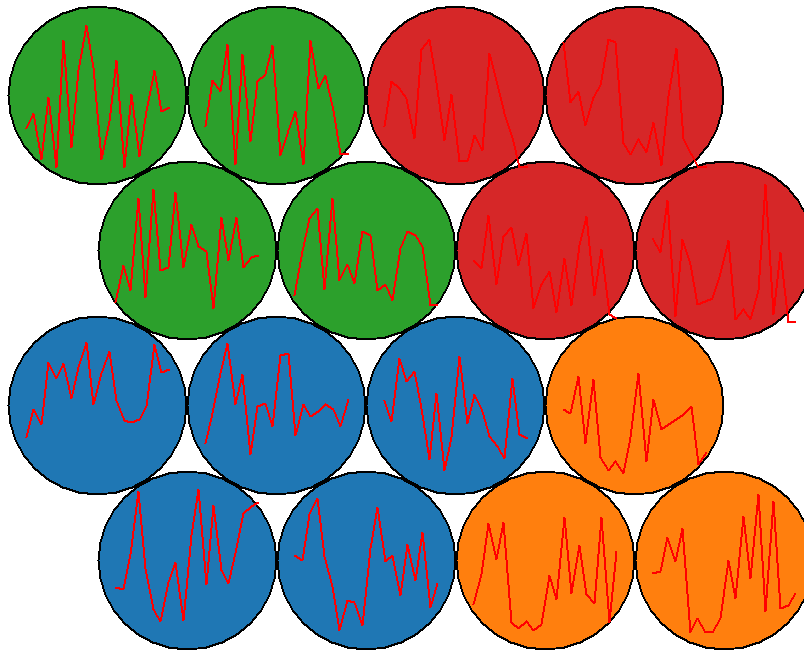
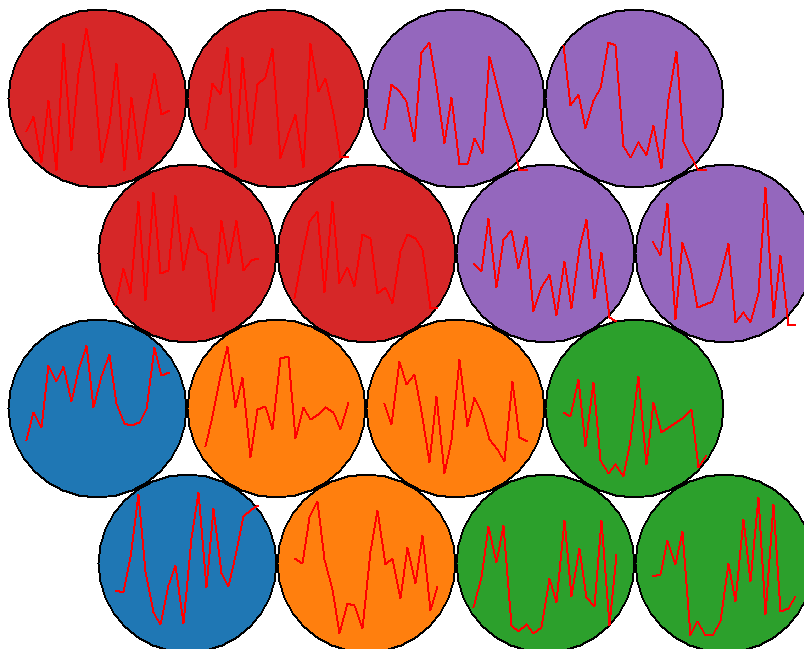


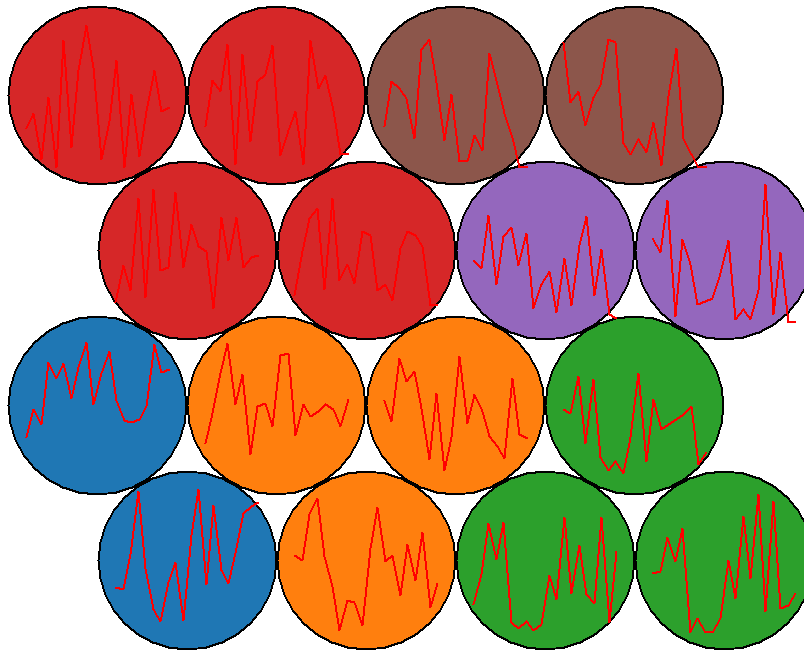
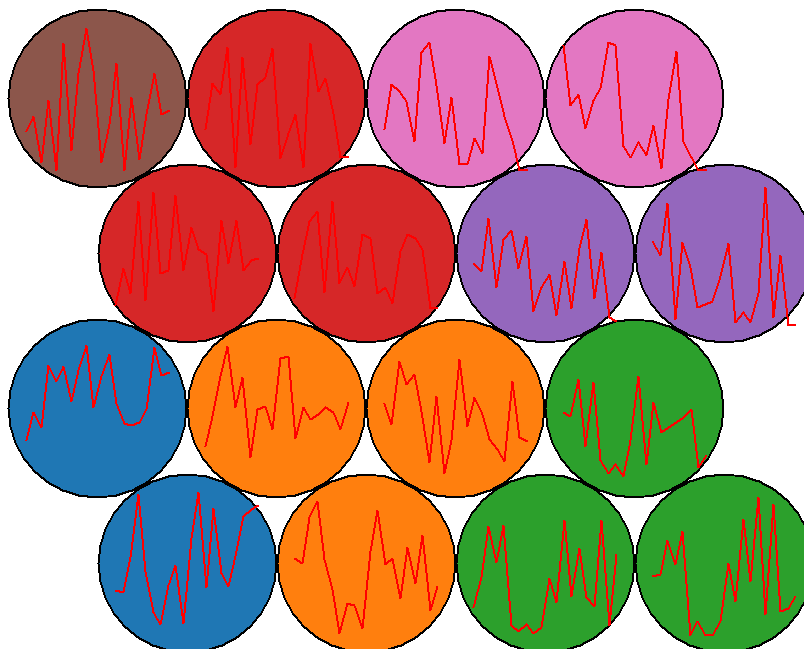
```
for (k in seq_len(max_cluster)) plotSOM(som_clusters$hierarchical[[toString(k)]], pas  
te0("hierarchical: ", k))
```

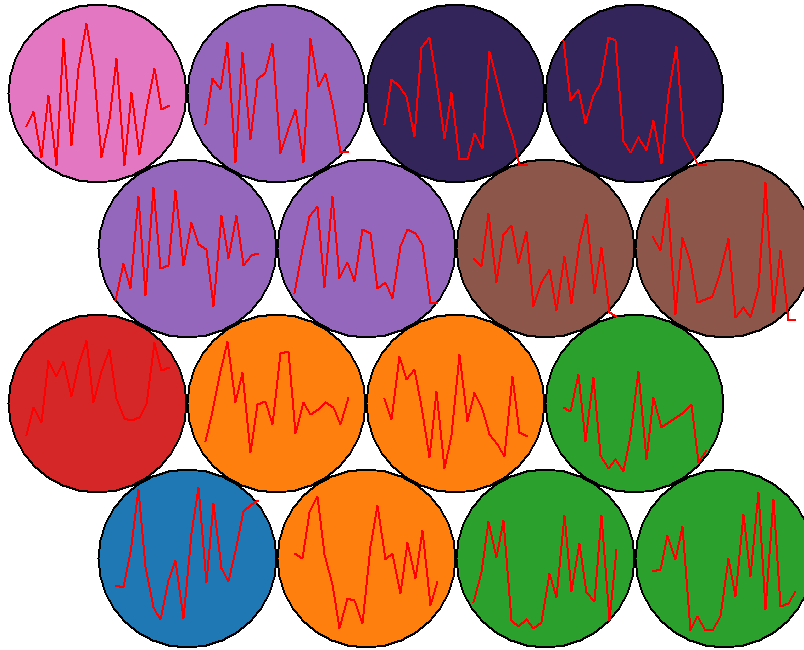
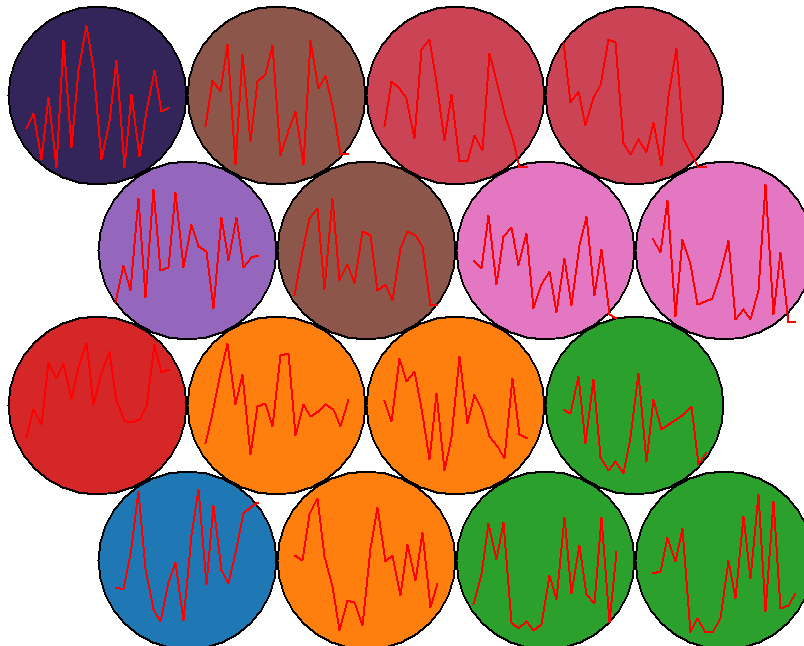
hierarchical: 1



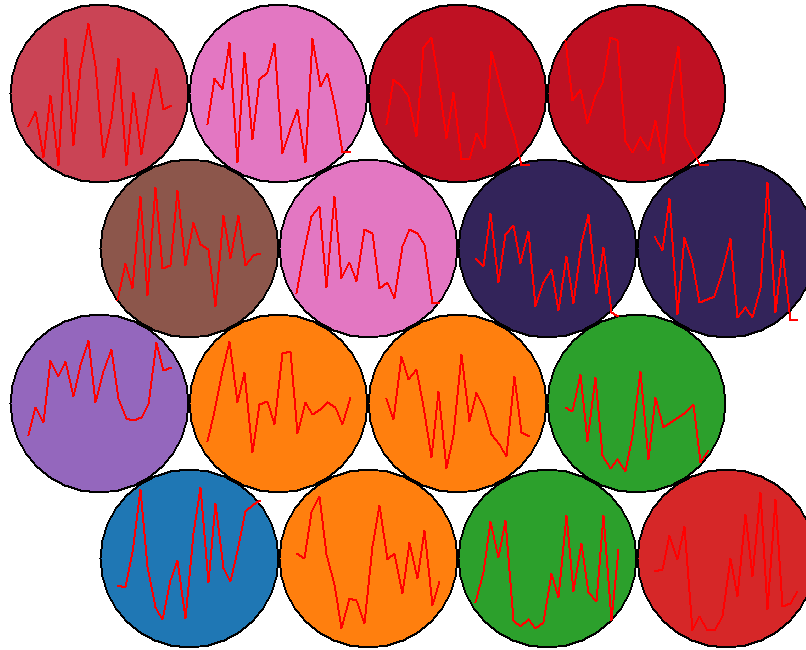
hierarchical: 2**hierarchical: 3**

hierarchical: 4**hierarchical: 5**

hierarchical: 6**hierarchical: 7**

hierarchical: 8**hierarchical: 9**

hierarchical: 10



2. Möglichkeit: Balkendiagramm

Jede Grid-Unit wird durch einen Balken repräsentiert. Die Ausprägungen der Merkmale dieser Grid-Unit sind in der Länge der Balken der einzelnen Merkmale zu entnehmen. Die Grid-Units können nach der Ausprägung aller Merkmale sortiert werden. Daher gibt es mehrere Balkendiagramme. So kann nach Gemeinsamkeiten und Trends zwischen den Grid-Units gesucht werden und auf diese Art und Weise einzelne Cluster ausgemacht werden.


```

codes <- getCodes(som_model)
plotdata <- data.frame(cluster = as.factor(seq_len(NROW(codes))), codes)
plotdata <- plotdata[order(plotdata$Alter),]
plotdata_long <- melt(plotdata, id.vars = "cluster")

for(i in unique(plotdata_long$variable)){
  level_order = plotdata_long %>%
  filter(variable == i) %>%
  group_by(cluster) %>%
  summarize(val = sum(value), .groups = "drop") %>%
  arrange(val) %>%
  pull(cluster)

  plotdata_long = mutate(plotdata_long, cluster = factor(cluster, levels = level_order))
}

p <- (plotdata_long %>%
  ggplot(aes(x = cluster, y = value, fill = variable, group = (cluster))) +
  geom_col(position = "stack", color = "black", alpha = .75) +
  coord_flip() +
  ggtitle("Grids ordered by", i) + xlab("Grids") + ylab("C-Vektoren"))
print(p)

```

