

Zusammenfassung EXPD / R-Studio

Kategorielle Merkmale

- Keine Zahlenwerte im engeren Sinn
- Endliche viele verschiedene Ausprägungen

Numerische Merkmale

- Nehmen numerische Werte an

	Nominale Merkmale	Ordinale Merkmale	Diskrete Merkmale	Stetige Merkmale
Eigenschaft	Ausprägungen lassen sich nicht in natürlicher Weise anordnen	Ausprägungen lassen sich in natürlicher Weise anordnen	Wertebereich ist abzählbar , aber möglicherweise unendlich gross	Merkmale, die prinzipiell jeden beliebigen Wert in einem (evtl. unbeschränkten) Bereich annehmen können. Zwischen zwei beliebigen Werten existiert immer noch ein Wert dazwischen
Beispiele	Farbe von Gegenständen, Studiengang, Nationalität	Beförderungsklasse (Economy, Economy Plus, etc), höchster Abschluss, Zustimmungswerte bei einer Kundenumfrage	Anzahl Passagiere in einem Flugzeug, Punktzahl in der Klausur	Umsatzzahlen eines KMUs, Durchmesser einer Schraube
Bemerkung	Können in einem Datensatz als Zahlen codiert sein (Gruppe 1, 2, 3, ...). Man kann mit diesen Zahlen aber nicht rechnen	Sind häufig als Zahlen codiert (5 = stimme voll zu, ..., 1 = stimme überhaupt nicht zu). Es ist aber nur begrenzt sinnvoll, mit diesen Zahlen zu rechnen.		Praktisch können fast alle stetigen Grössen nur diskret gemessen werden, da die Genauigkeit der Messungen limitiert ist.

Histogramm – Klassenzahl

Tipp: Mit verschiedenen Breiten experimentieren und das Histogramm wählen, das die Form der Verteilung am besten wiedergibt. Die Stichprobengrösse insgesamt und pro Balken darf aber für aussagekräftiges Diagramm nicht zu klein sein.

Verschiedene **Faustregeln** zur Wahl der Klassenanzahl k

- $k = \lfloor \sqrt{n} \rfloor$
- $k \approx \sqrt{n}$
- $k = \lfloor 10 \cdot \log_{10}(n) \rfloor$
- $k \approx 1 + \log_2(n)$ (Default in R)

mit n = Anzahl Beobachtungen

- Lagemasse** beschreiben, um welchen “mittleren” Wert die Daten verteilt sind (nur für unimodal sinnvoll)
- Streuungsmaße** geben an wie “breit” die Verteilung ist, d.h. wie stark die Werte streuen.

Mittelwerte

Arithmetisches Mittel – Durchschnitt \bar{x}

$$m = (x_1 + x_2 + \dots + x_n) / n$$

Geometrisches Mittel - Mittel der Wachstumsraten (Zinssatz/Zinseszins)

$$g = (x_1 \cdot x_2 \cdot \dots \cdot x_n)^{1/n}$$

Harmonisches Mittel - Mittel von Quotienten (Ratios, Verhältnisse)

$$h = n / (1/x_1 + 1/x_2 + \dots + 1/x_n)$$

Median \tilde{x}

Quantile

Quantile teilen eine Stichprobe in einem bestimmten Verhältnis. α Quantil teilt die Stichprobe im Verhältnis $\alpha : (1 - \alpha)$. Beispiel: $Q_{90\%}$ teilt die Stichprobe in 90% : 10%

Streuemasse

Varianz (s_x^2)

Mittlere quadratische Abweichung der Beobachtungen vom arithmetischen Mittel (\bar{x})

In R \rightarrow var(A) oder sd(A)^2

Standardabweichung (s_x)

Quadratwurzel der Varianz. Gleiche Einheit wie Beobachtung (besser interpretierbar)

In R \rightarrow sd(A) oder sqrt(var(A))

MAD (median absolute deviation)

$$MAD_x = 1.4826 \cdot \text{median}(|x_1 - \tilde{x}|, |x_2 - \tilde{x}|, \dots, |x_n - \tilde{x}|)$$

Robustes Gegenstück zur Standardabweichung. Mittelwertbildung durch Mediane ersetzen. Konstante stellt sicher, dass MAD und s_x , im Fall wenn die Daten aus einer Normalverteilung stammen, das gleiche messen. In R \rightarrow mad(A)

IQR (inter quartile range)

Abstand zwischen oberem und unterem Quartil: $IQR = Q_3 - Q_1$ In R \rightarrow IQR(A)

Spannweite

Differenz zwischen Maxima und Minima. Ist jedoch anfällig auf Ausreisser. Nimmt mit Stichprobenumfang zu \Rightarrow als Streuungsmass ungeeignet In R \rightarrow max(A) – min(A)

Wann IQR, wann Standardabweichung?

- Bei rechts/linksschiefen Verteilungen liefert die Standardabweichung einen Wert, der für den Grossteil der Daten nicht repräsentativ (zu hoch ist). Die IQR liefert meist eine bessere Vorstellung der Streuung.
- Für symmetrische Verteilungen ohne Ausreisser ist die Standardabweichung ein gut geeignetes Mass.
- Die Standardabweichung ist nicht robust! Eine einzige, falsche Beobachtung kann sie grob zu verfälschen.

Boxplot

Weniger geeignet für bi- /multimodale Verteilungen: Diese wird in einem Boxplot nicht erfasst.

Bivariate Darstellungen

Verschiedene Variablentypen benötigen verschiedene Handhabung. Es gibt 3 Fälle:

1. Kategoriell vs. Kategoriell
2. Metrisch vs. Kategoriell
3. Metrisch vs. Metrisch

Zwei kategorielle Variablen

Kreuztabelle → `table(dat$Haarfarbe, dat$Augenfarbe)`

Anteil am Gesamten → `prop.table(table(dat$Haarfarbe, dat$Augenfarbe))`

Anteil pro Zeile → `prop.table(table(dat$Haarfarbe, dat$Augenfarbe), margin = 1)`

Anteil pro Spalte → `prop.table(table(dat$Haarfarbe, dat$Augenfarbe), margin = 2)`

Gruppierte Balkendiagramme

Gestapelten Balkendiagramm → `barplot(table(dat$Haarfarbe, dat$Augenfarbe))`

Gruppiertes Balkendiagramm → `barplot(table(dat$Haarfarbe, dat$Augenfarbe), beside = TRUE)`

Mosaikplot

`mosaicplot(table(dat$Augenfarbe, dat$Haarfarbe))` → variablen können auch vertauscht werden = anderer Plot!

Kategorielle und metrische Variable

Kennzahlentabelle

`m <- tapply(kdata$einkauf, kdata$zivilstand, "mean")`

`s <- tapply(kdata$einkauf, kdata$zivilstand, "sd")`

`cbind(Mittelwert = m, Standardabweichung = s)`

Boxplot (nur für unimodale Verteilung geeignet)

`boxplot(einkauf ~ zivilstand, data = kdata)`

Stripchart (vor allem bei bimodaler Verteilung sehr geeignet)

`stripchart(einkauf ~ zivilstand, data=kdata, vertical=TRUE, method="stack")`

Zwei metrische Variablen

Streudiagramm (plot oder scatter.smooth mit laufenden Mitteln → Glätter)

`scatter.smooth(x=dat$alter, y=dat$einkauf, pch=16, col = rgb(0,0,0, alpha = 0.1), lpars = list(col = "red", lwd = 2))`

Aus einem Streudiagramm kann man Form, Richtung und Stärke des Zusammenhangs erkennen.

Form: streuen die Punkte um eine Gerade (linear), um eine Kurve oder existieren diverse Punktwolken?

Richtung: Je grösser die Werte einer Variablen,

...desto grösser sind die Werte der anderen Variable (positiver Zusammenhang)

...desto kleiner sind die Werte der anderen Variable (negativer Zusammenhang)

Stärke: Wie breit ist die Punktwolke? Wenn die Punktwolke breit ist, dann ist der Zusammenhang schwach. Wenn die Punktwolke schmal (Extremfall eine Gerade!) ist, dann ist der Zusammenhang stark.

Kovarianz

Der Zusammenhang zweier Variablen wird durch die Kovarianz beschrieben. Wir betrachten die Lage der Punkte bez. des Schwerpunktes (\bar{x} , \bar{y}).

- Die Kovarianz informiert nur über die Richtung des Zusammenhangs, nicht dessen Stärke.
- Der absolute Zahlenwert ist schwer zu interpretieren, da die Grösse der Kovarianz vom Massstab x/y abhängt.
- Die Kovarianz ist sehr anfällig für Ausreisser.
- Die Kovarianz von zwei identischen Variablen entspricht der Varianz.

Berechnung in R → `cov(dat$umsatz, dat$werbung)`

Es stehen 25 verschiedene Symbole zur Verfügung:

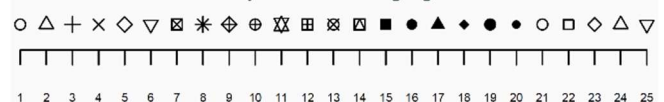
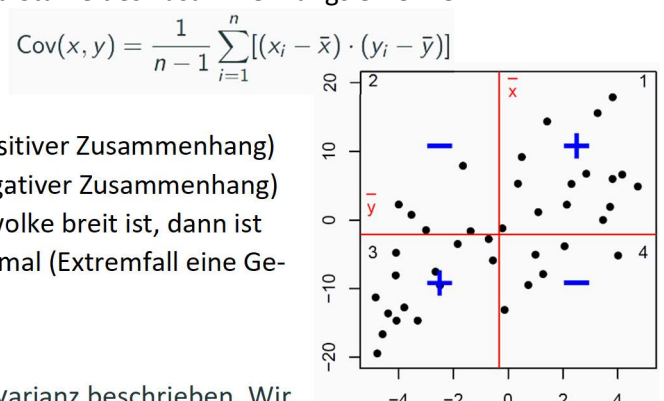


Abbildung 1: Symbole für Plot's, pch = 8



Pearson-Korrelation (einfach oder Produkt-Momenten Korrelation)

Misst die Stärke des **linearen Zusammenhangs**, d.h. wie eng die Punkte um eine Gerade liegen bzw. wie stark sie streuen, und entspricht der standardisierten Kovarianz.

Berechnung in R → `cor(variable1, variable2)`

Eigenschaften

- $r_{xy} > 0$ → positive Korrelation (Je mehr, desto mehr)
- $r_{xy} < 0$ → negative Korrelation (Je mehr, desto weniger)
- $r_{xy} = 1$ → Punkte liegen exakt auf einer Gerade positiver Steigung
- $r_{xy} = -1$ → Punkte liegen exakt auf einer Gerade mit negativer Steigung
- r_{xy} nahe bei ± 1 → Punkte streuen eng um eine Gerade
- $r_{xy} \approx 0$ → kein linearer Zusammenhang, es kann aber ein anderer Zusammenhang bestehen
- r_{xy} ist **nicht die Steigung der Geraden**

Wichtig: Die **Form** des Zusammenhangs wird von der Korrelation **nicht geprüft**, fließt aber als Annahme in die Berechnung ein! Korrelation nie ohne Blick auf ein Streudiagramm beurteilen!!!

Die Spearman-Korrelation (Rang-Korrelation)

Misst die Stärke des **monotonen Zusammenhangs**, d.h. wie nahe die Punkte um eine Kurve liegen, die von einer beliebigen, monotonen Funktion definiert ist. Funktionsweise: Man bestimmt separat die Ränge der x- und y-Werte und berechnet auf diesen Rängen die Pearson-Korrelation. Die Rangkorrelation ist der Pearson-Korrelation vorzuziehen, wenn der **Zusammenhang nicht linear**, sondern **monoton** ist, es Ausreisser geben könnte und die Werte (x_i, y_i) nicht glockenförmig verteilt sind.

Berechnung in R → `cor(variable1, variable2, method = "spearman")`

Eigenschaften:

- $r_{xy} = 1$ → Punkte liegen exakt auf einer monotonen Kurve mit positiver Steigung
- $r_{xy} = -1$! Punkte liegen exakt auf einer monotonen Kurve mit negativer Steigung
- r_{xy} nahe bei ± 1 ! Punkte streuen eng um eine **monotone Kurve**
- *restliche Eigenschaften sind gleich wie bei der Pearson-Korrelation*

Vorsicht bei kleiner Stichprobengröße! Bei 2 Punkten ist die Korrelation immer ± 1 .

Logische Operatoren

Reihenfolge der Auswertung: ! kommt vor & kommt vor |, mittels Klammern () kann die Reihenfolge angepasst werden. (Verknüpfung mit & kann als Multiplikation mit | als Addition gesehen werden)

Transformation

Bei einer Transformation werden die Werte einer Variablen mittels einer eindeutigen Zuordnung zu neuen Werten umgerechnet: $x \mapsto f(x)$

Gründe für Datentransformation:

- Zusammenfassen von Beobachtungen in Klassen
- Umrechnen von Einheiten
- Informativere Darstellung/Daten standardisieren
- Änderung unvorteilhaften Form einer Verteilung

Nominale sowie ordinale Daten

Es gibt Transformationen bei beiden jeweils mit und ohne Informationsverlust.

Lineare Transformation

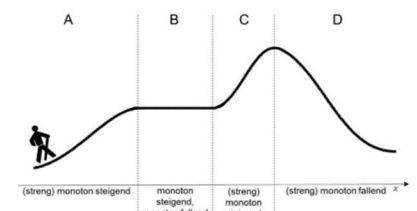
Änderung der Masseinheit (inch in cm): Transformation bei der keine Verschiebung involviert ist: $x \mapsto f(x) = a * x$

Korrektur von Messwerten (Waage verstellt): Jeder Wert wird um denselben Betrag korrigiert $x \mapsto f(x) = x + b$

Verschiebung der Skala und des Mess-Nullpunktes ($^{\circ}\text{F}$ in $^{\circ}\text{C}$): $x \mapsto f(x) = a * x + b$

Streng monotone Transformation

Bei Transformationen mit streng monoton **steigenden** Funktionen bleibt die Reihenfolge der Werte bestehen. Umsetzung durch Logarithmus oder Wurzelfunktion!



Auswirkungen von Transformationen

Je nach Transformation können sich die Kennzahlen der Variablen ändern (Mittelwert, Standardabweichung) oder kann sich die Verteilung der Variablen ändern.

Im Fall der linearen Transformation lassen sich die Auswirkungen mathematisch einfach beschreiben. Bei **streng** monoton steigenden Transformationen ist dies etwas schwieriger.

Umrechnung der linearen Transformationen

Die lineare Transformation ändert die Form der Verteilung nicht, nur die Achsenbeschriftung ändert sich.

- Mittelwert: $\bar{y} = a * \bar{x} + b$.
- Standardabweichung: $sd_y = |a| * sd_x$

Lage- und Streumasse ändern sich bei einer linearen Transformation ($x \mapsto f(x) = a * x + b$) wie folgt:

- $Lage_y = a * Lage_x + b$
- $Streuung_y = |a| * Streuung_x$

Wobei Lage: arithm. Mittel, Median, Quantile (nur für $a \geq 0$), Modus

Streuung: Standardabweichung, MAD, IQR.

Achtung: Varianz ist ein quadratisches Streumass daher $Var_y = a^2 * Var_x$

Standardisierung (Spezialfall einer linearen Transformation)

Die Standardisierung ist nützlich, um Daten zu vergleichen, wenn man sich für die Verteilungsform interessiert, Lage

und Streuung aber nicht berücksichtigen möchte: $x \mapsto z = \frac{1}{sd_x} * x - \frac{\bar{x}}{sd_x}$

Die standardisierten Variablen haben Mittelwert 0 und Standardabweichung 1

$$z = \frac{1}{sd_x} * x - \frac{\bar{x}}{sd_x} = 0 \quad sd_z = \frac{1}{sd_x} * sd_x = 1$$

Im Gegensatz zur linearen Transformation wird durch eine streng monotone Transformation die Verteilung der Variable geändert.

Was tut man mit mehr als 2 Variablen?

Mehrere kategorielle Variablen

➔ Mosaikplot `mosaicplot(~k1 + k2 + k3 + ...)`

1 quantitative und mehrere kategorielle Variable

➔ Boxplots `boxplot(m1~k1 + k2 + k3 + ...)`

➔ Faktor/Design-Plot `plot.design`

2 quantitative und mehrere kategorielle Variable

Die 2 quantitativen Variablen werden mittels eines Streudiagramms aufgezeichnet und die kategorieller Variable zusätzlich durch Farben, Symbolform und Symbolgrösse visualisiert. `plot(m1, m2, col[k1], pch[k2])`

Reihenfolge in abnehmender Wirksamkeit ist zu beachten:

- | | | |
|-----------|--|--------------------------|
| 1. Grösse | 3. Orientierung (einer Linie, eines Rechtecks) | 5. Intensität |
| 2. Text | 4. Form (Stern, Kreis, Rechteck) | 6. Farbton/Farbsättigung |

Mehr als 2 quantitative Variablen

3D-Plots

Co-plot

Streudiagramm-Matrix

Korrelationsmatrix

Stichproben - Schlussfolgerungen

- Die Standardabweichung nimmt mit wachsendem n (Anzahl der Stichproben) ab.
- $\sqrt{n} * s_x$ ist ungefähr konstant, d.h. Standardabweichung ist umgekehrt proportional zur Wurzel aus n.

Möchte man die Standardabweichung der Mittelwerte halbieren, muss man die Stichprobengrösse vervierfachen! Erkenntnis kann im Boxplot verwendet werden, um bei Stichproben ein Intervall anzugeben, in dem der wahre Median «ziemlich sicher» liegt. **Faustregel:** Wenn sich die Kerben von zwei Boxen nicht überschneiden, so besteht ein signifikanter Unterschied zwischen den Gruppen.

Hauptkomponentenanalyse (PCA)

Idee: Finde diejenigen Achsen, welche die Daten am besten erklären und verwende nur ein paar davon (2-3).

Das Koordinatensystem wird so gedreht, dass die Varianz in Richtung der neuen x-Achse am grössten ist.

→ Drehung der Achsen, so dass nur minimal Information verloren geht.

→ Dies ist mathematisch äquivalent die Richtung zu finden, in der die grösste Streuung (Varianz) vorliegt.

Die PCA wird verwendet, um multidimensionale Daten durch die ersten zwei Hauptkomponenten auf 2 Dimensionen zu reduzieren (Visualisierung), so dass der grösste Teil der Information (Streuung) der Daten erhalten bleibt.

Schritte der PCA:

- Verschieben des Koordinatensystems in den Schwerpunkt der Daten.
- Rotation des ursprünglichen Koordinatensystems zum Koordinatensystem der Hauptkomponenten, sodass die Varianz entlang der ersten Hauptkomponente am grössten ist.
- Der grösste Teil der restlichen Varianz soll entlang der zweiten Hauptkomponenten liegen (Achse steht senkrecht zur ersten Hauptkomponente).

Die Rotationsmatrix/Ladungsmatrix rechnet die ursprünglichen Werte in neuen Werte im "gedrehtem" Koordinatensystem um. Der Absolutbetrag der "Ladungen" gibt an wie wichtig die ursprüngliche Variable für die neue Position ist. Ist die Ladung der Variable 1 mit 0.899 grösser als die von Variable 2 mit 0.44, damit hat die Variable 1 grösseres Gewicht als Variable 2. → `res <- prcomp(X, scale=F)` \n `res$x` (neue Koordinaten)
`res$rotation` gibt die neuen Wert im Koordinatensystem nach der Rotationsmatrix an.

Skalierung

```
dat$umsatzS <- (dat$umsatz - mean(dat$umsatz))/sd(dat$umsatz)
dat$werbungS <- (dat$werbung - mean(dat$werbung))/sd(dat$werbung)
Kovarianz der skalierten Daten = var(dat[,c("umsatzS", "werbungS")])
pca <- prcomp(dat[,1:2], scale = TRUE) # Identisch zu: #
prcomp(dat[,c("umsatzS", "werbungS")], scale = FALSE)
pca$rotation
```

Auswirkung der Skalierung:

- Wenn Variablen nicht skaliert, hat die Variable mit grösster Streuung grössten Einfluss auf die neue Koordinate.
- Wenn die Variablen skaliert sind, dann haben alle Variablen das gleiche Gewicht.

Faustregel: Skaliere, wenn

die Variablen verschiedene Einheiten haben explizit gewollt ist, dass alle Variablen gleiches Gewicht haben

Nicht skalieren, wenn

- alle Variablen die gleiche Einheit haben und vergleichbar sind
- Wenn man nicht genau weiss, was zu tun ist, ist es meistens besser zu standardisieren.

Dimensionsreduktion

Situation: Wir haben nun die Rotation in ein neues Koordinatensystem. Die Dimension ist hierbei immer noch gleich gross wie im gesamten Datensatz. Ziel: Dimensionsreduktion. Verwenden wir nun nur die ersten paar Hauptkomponenten und hoffen, dass die Daten dadurch möglichst gut beschrieben werden Wie gut ist unsere Näherung? Mass dazu ist die erklärte Varianz.

Die totale Varianz berechnet sich für p Variablen als

Durch die Rotation wird die totale Varianz nicht verändert.

```
pca <- prcomp(dat[,1:2], scale = FALSE)
```

```
sum(diag(cov(dat[,1:2])))
```

```
sum(diag(cov(pca$x)))
```

ergeben beide das Gleiche!

$$\text{Var}_{\text{total}} = \sum_{j=1}^p \text{Var}(X_j)$$

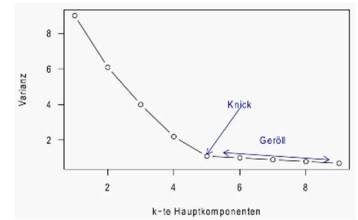
$$p_k = \left(\sum_{j=1}^k \text{Var}(Z_j) \right) / \text{Var}_{\text{total}}$$

Als Gütekriterium der Approximation wählen wir daher den Anteil der Varianz der durch die ersten k-Hauptkomponenten erklärt wird:

1. Kriterium: k sollte so gewählt werden, dass ca. 80% der totalen Varianz durch die berücksichtigten Hauptkomponenten erklärt wird.

2. Kriterium Varianzen in abnehmender Ordnung grafisch betrachten mit screeplot. Die Position des Knicks im Plot ergibt einen weiteren Hinweis, wie viele PCs benötigt werden. Die PCs nach dem Knick tragen nicht viel zur Gesamtinformation bei.

```
pca <- prcomp(...)
screeplot(pca)
```



Interpretation der PCA

Im Score-Plot werden die ersten beiden Hauptkomponenten, die die meiste Varianz erklären visualisiert.

```
load("data/body.Rdata")
pca <- prcomp(body, scale=TRUE)
library(ggfortify)
autoplot(pca)
```

Der Biplot entspricht dem Score-Plot mit standardisierten Daten und zeigt zusätzlich noch rote Pfeile (Projektionen der ursprünglichen Merkmale in der Ebene der beiden Hauptkomponenten):

Pfeillänge = Varianz der Variable (Bei Skalierung haben alle Pfeile die gleiche Länge)

Bei korrekter Projektion zeigt der Winkel zwischen den Pfeilen die Korrelation. Je kleiner der Winkel, desto größer die Korrelation zwischen den beiden Variablen.

```
pca <- prcomp(body, scale=TRUE)
biplot(pca)
```

Paket ggfortify gibt es Funktion autoplot, welche sich gut zum Darstellen/Interpretieren einer PCA eignet (ggplot).

```
library(ggfortify)
autoplot(pca, loadings=TRUE, loadings.label=TRUE, label = TRUE, label.hjust = -0.3)
```

Boxplot mit Kerben/Notch (boxplot(dat, notch = T))

Die Kerbe überdeckt mit einer Wahrscheinlichkeit von 95% den Median der Population oder Grundgesamtheit (nicht der Stichprobe) der im Allgemeinen aber nicht bekannt ist. Wenn sich zwei Kerben nicht überschneiden, dann ist die Evidenz gross, dass die Stichproben von unterschiedlichen Populationen stammen.

Mosaicplot – zwei kategoriale Variablen

Merkmale:

- Fläche ist proportional zur Anzahl Beobachtung für die Merkmalskombination in der Stichprobe.
- Breite der Säulen ist proportional zur relativen Häufigkeit der ersten Variable, die Höhe proportional zur zweiten Variable.
- Die Anzahl Beobachtungen ist nicht ablesbar.
- Vertauschen der Variablen führt zu einem anderen Mosaikplot

Stripchart - Kategoriell und metrisch

Bei bimodalen Verteilungen sind Stripcharts eine Alternative zu dem Boxplot.

```
par(mfrow=c(1,2))
stripchart(einkauf ~ zivilstand, data=kdata, vertical=TRUE, method="stack")
stripchart(einkauf ~ kaufkraft, data=kdata, vertical=TRUE, method="stack")
```

Manipulationen von Strings

gsub

Entfernt gewünschte Zeichen aus den Ausprägungen der einzelnen Variablen

```
censUSA$education <- gsub(pattern = " ", replacement = "", x = censUSA$education)
censUSA$occupation <- gsub(pattern = " ", replacement = "", x = censUSA$occupation)
censUSA$education <- factor(censUSA$education)
censUSA$occupation <- factor(censUSA$occupation)
```

substring

Entfernt von den Variablen/Kategorien die ersten hier bspw. fünf Buchstaben des Namens!

```
names(Erhebung) <- substring(text = names(Erhebung), first = 5)
```

nchar(nzz) → Anzahl der Zeichen im String

nzz <- gsub(pattern = "ä", replacement = "ae", x = nzz) → suchen und ersetzen durch

nzz <- gsub(pattern = "ö", replacement = "oe", x = nzz)

nzz <- gsub(pattern = "ü", replacement = "ue", x = nzz)

strsplit

Sonderzeichen z.B. ? oder . müssen zwischen [] stehen, wenn explizit nach ihnen gesucht werden soll.

nzzList <- strsplit(x = nzz, split = "[.]") → Text wird in einzelne Sätze umgewandelt

nzzVec <- unlist(nzzList)

→ wandelt ihn in Vektor um, dessen Elemente gerade die Sätze des eingelesenen Textes enthalten

Entfernt überflüssige Leerzeichen am Anfang einiger Sätze und fügt dem Satzenden einen Punkt hinzu.

```
vSel <- which(substr(nzzVec,1,1) == " ")
```

```
nzzVec[vSel] <- substr(nzzVec[vSel], start = 2, stop = nchar(nzzVec)[vSel])
```

```
nzzVec <- paste(nzzVec, ".", sep = "")
```

```
nzzVec
```

grepl

In wie vielen Sätzen kommt der String AHV mindestens einmal vor?

```
sum(grepl(pattern = "AHV", x = nzzVec))
```

plot Argumente

type

"p" Punkte

"l" verbundene Linien

"o" verbundene Linien und Punkte

"b" nicht verbundene Linien mit Punkten

"c" nicht verbundene Linien ohne Punkte

"s" Linien als Treppe

"h" vertikale Linien (ähnlich zum Barplot)

"n" keine Linien und Punkte (leere Grafik)

log

log = "y"

log = "x"

log = "xy"

legend()

```
legend(x= 4.5, # x-Koordinate, oder "topright"
```

```
      y= 35, # y-Koordinate
```

```
      title = "Fahrzeugantrieb",
```

```
      legend= c("4wd", "Front", "Heck"), # Text
```

```
      col= c("black", "#DF536B", "#61D04F"), # Farbe
```

```
      pch= 17, # Symbole lty u. lwd auch möglich
```

```
      horiz = T # horizontale Auflistung
```

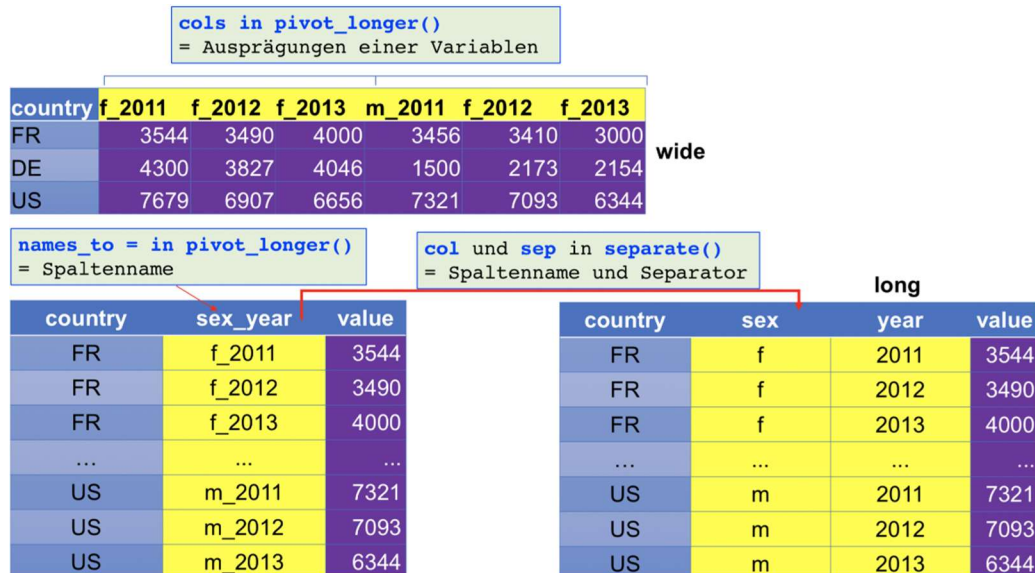
Reshape

Daten vom wide ins long Format transformieren, sodass mit diesen auch «gearbeitet» werden kann. Dabei werden die Funktionen `pivot_longer` sowie `separate` aus dem Paket `tidyr` verwendet.

```
library(tidyr)
```

```
longTBC <- pivot_longer(data = TBCases,
  cols = 2:7, #f_2011 – m_2013
  names_to = "sex_year") #Name der neuen Spalte
```

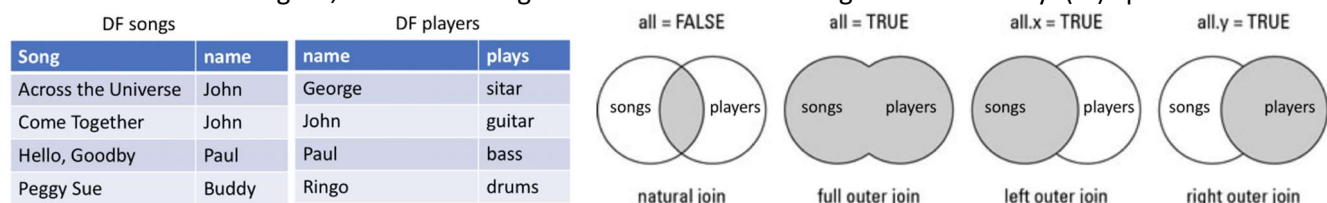
```
longTBC <- separate(data = longTBC,
  col = "sex_year", #welche Spalte
  sep = "_", #was trennt die Elemente
  into = c("sex", "year")) #neue Spaltennamen
```



Zusammenführen von Datensätzen → `merge()`

Damit Datensätze zusammengeführt werden können, sollten beide mindestens eine Spalte mit ähnlichem Inhalt (nicht zwingen mit dem gleichen Namen) aufweisen.

x, y: x und y Data Frame
by, by.x und by.y: Name(n) der Spalte(n) des x, y Data Frames, welche für die Vereinigung verwendet werden. Wenn nichts angegeben, wird automatisch anhand gemeinsamer Spaltennamen vereinigt
all: logical, Art der Zusammenführung (siehe Grafik unten)
all.x: logical, TRUE = left outer join (siehe Grafik unten)
all.y: logical, TRUE = right outer join (siehe Grafik unten)
sort: logical, TRUE = vereinigter Datensatz wird aufsteigend nach den `by=()` Spalten sortiert



Eine der am häufigsten verwendeten Zusammenführungen ist der sog. left outer join, durch den der x DF (im Beispiel songs) mit Beobachtungen aus y DF (players) über gemeinsame Variable(n) (`by = c("...")`) erweitert wird. Fehlende Beobachtungen in y DF (im Beispiel players) werden zu NA-Werten.

natural join

```
newdat <- merge(x = songs,
  y = players,
  by = "name",
  all = FALSE)
```

	name	song	plays
1	John	Across the Universe	guitar
2	John	Come Together	guitar
3	Paul	Hello, Goodbye	bass

left outer join

```
newdat <- merge(x = songs,
  y = players,
  by = "name",
  all.x = TRUE)
```

	name	song	plays
1	Buddy	Peggy Sue	<NA>
2	John	Across the Universe	guitar
3	John	Come Together	guitar
4	Paul	Hello, Goodbye	bass

ggplot2

asthetics:

- color
- shape
- fill
- linetype
- stroke
- size

Pipe-Operator %>%

Insbesondere im Data Science Bereich trifft man in R-Code immer wieder auf den sogenannten Pipe-Operator %>%. Der %>%-Operator wurde 2014 mit dem Paket magrittr eingeführt und erfreute sich rasch grosser Beliebtheit.

Benutzt wird die Pipe meistens über das dplyr-Paket bzw. die Paketsammlung tidyverse. Eine Sammlung von R-Paketen, die für Data Science entwickelt wurden. Sie helfen beim managen und «streamlinen» des Daten-Workflows.

Seit neustem (R-Version 4.1.0) wurde der Pipe-Operator auch in Base R integriert. |>

head(mtcars, n = 2) = mtcars |> head(n = 2)

head(summary(as.factor(rownames(mtcars))), n = 3) as.factor() |>

summary() |>

mtcars |> head(n = 2)

rownames() |>

Alternative Standardisierung

- Varianzen bzw. Korrelationen sind sensitiv bezüglich Ausreissern → Ergebnis der PCA kann von Ausreissern stark verfälscht werden!
- In der Praxis Standardisierung meistens mit Mittelwert und Standardabweichung.
- Häufig wäre es besser / sicherer die Standardisierung robust durchzuführen (z.B. Median als Lagemass und MAD als Skalenmass) → Ausreisser in den einzelnen Variablen werden besser sichtbar.
 - Die entsprechende Standardisierung muss für prcomp manuell durchgeführt werden und die Hauptkomponenten-Analyse anschliessend mit den manuell standardisierten Daten durchgeführt werden prcomp(dat_robust, scale=FALSE).
 - Ein zweiter Ansatz ist es, die Hauptkomponentenanalyse direkt mit robusten Schätzern durchzuführen. In R gibt es Pakete, die das für uns übernehmen → PcaHubert aus Paket rrcov.

Die robuste PCA wird häufig auch verwendet, um Ausreisser in hochdimensionalen Daten zu detektieren.

Es gibt zwei Arten von Ausreissern.

- Ausreisser im PCA-Raum
 - Grosse Distanz zum Ursprung der PCA
 - Score-Distanz
- Orthogonal Ausreisser (Punkte im Raum der berücksichtigen Hauptkomponenten nicht richtig erfasst werden). Rekonstruktions-Fehler (PCA minimiert die quadrierte Summe dieser Distanzen). Punkte mit grossem Rekonstruktions-Fehler sind Ausreisser!

$$SD_i = \sqrt{\sum_{j=1}^k \frac{y_{ij}^2}{\lambda_j}} \quad (\text{optionale Cutoff-Line } \sqrt{\chi_{k,0.975}^2})$$

$$OD_i = ||x_i - \hat{\mu} - P \cdot y_i^T||$$

Zusammenfassung Hauptkomponentenanalyse

- Drehung (orthogonale Transformation) des p-dimensionalen Raums der zentrierten Originalvariablen in das Hauptkomponentensystem.
- Hauptkomponenten werden so konstruiert, dass sie untereinander unkorreliert sind und in absteigender Weise einen möglichst grossen Teil der Totalvarianz abdecken.
- Die ersten Hauptkomponenten decken ein Grossteil der Varianz ab.
- Visualisierung → Strukturen entdecken → Ausreisser detektieren
- Die wichtigsten Hauptkomponenten können für weiterführende Analysen benutzt werden (Machine Learning).
- Nutzen bzgl. einer Dimensionsreduktion ist besonders hoch, wenn die Originalvariablen stark korreliert sind.

Nachteile:

- Ergebnisse abhängig von der Skalierung.
- Anzahl der auszuwählenden bedeutenden Hauptkomponenten ist nicht eindeutig.
- Hauptkomponenten schwierig zu interpretieren.

coral3	gray27	gray39	gray87	gray99	lightpink1	mistyrose1	pink4	slategray1
	gray26	gray38	gray86	gray98	lightpink	mistyrose	pink3	slategray
	gray25	gray37	gray85	gray97	lightgray	mintcream	pink2	slateblue4
	gray24	gray36	gray84	gray96	lightgreen	midnightblue	pink1	slateblue3
chocolate4	gray23	gray35	gray83	gray95	lightgray	mediumvioletred	peru	slateblue2
	gray22	gray34	gray82	gray94	lightgoldenrod4	mediumturquoise		slateblue1
	gray21	gray33	gray81	gray93	lightgoldenrod3	mediumspringgreen	peachpuff4	slateblue
	gray20	gray32	gray80	gray92	lightgoldenrod2	mediumslateblue	peachpuff3	skyblue4
chartreuse4	gray19	gray31	gray79	gray91	lightgoldenrod1	mediumseagreen	peachpuff2	skyblue3
	gray18	gray30	gray78	gray90	lightgoldenrod1	mediumpurple4	peachpuff1	skyblue2
	gray17	gray29	gray77	gray89	lightcyan4	mediumpurple3	peachpuff	skyblue1
	gray16	gray28	gray76	gray88	lightcyan4	mediumpurple2	papayawhip	whitesmoke
chartreuse1	gray15	gray27	gray75	gray87	lightcyan3	mediumpurple1	skyblue	wheat3
	gray14	gray26	gray74	gray86	lightcyan2	mediumpurple	sienna4	wheat2
	gray13	gray25	gray73	gray85	lightcyan1	mediumorchid4	sienna3	wheat1
	gray12	gray24	gray72	gray84	lightcyan	mediumorchid3	sienna1	
cadetblue3	gray11	gray23	gray71	gray83	lightcoral	mediumorchid2	sienna	violetred4
	gray10	gray22	gray70	gray82	lightblue4	mediumorchid1	seashell4	violetred3
	gray9	gray21	gray69	gray81	lightblue3	mediumorchid	seashell3	violetred2
	gray8	gray20	gray68	gray80	lightblue2	mediumblue	seashell2	violetred1
cadetblue1	gray7	gray19	gray67	gray79	lightblue1	mediumaquamarine	seashell1	violet
	gray6	gray18	gray66	gray78	lightblue	maroon4	seashell	turquoise4
	gray5	gray17	gray65	gray77	lemonchiffon4	maroon3	seagreen4	turquoise3
	gray4	gray16	gray64	gray76	lemonchiffon3	maroon2	seagreen3	turquoise3
brown4	gray3	gray15	gray63	gray75	lemonchiffon2	maroon1	seagreen2	turquoise2
	gray2	gray14	gray62	gray74	lemonchiffon1	maroon	seagreen1	turquoise1
	gray1	gray13	gray61	gray73	lemonchiffon	magenta4	seagreen	tomato4
	gray0	gray12	gray60	gray72	lawngreen	magenta3	sandybrown	tomato3
blueviolet	gray	gray11	gray59	gray71	lavenderblush4	magenta2	salmon4	tomato2
	goldenrod4	gray10	gray58	gray70	lavenderblush3	magenta1	salmon3	tomato1
	goldenrod3	gray9	gray57	gray69	lavenderblush2	magenta	salmon2	tomato
	goldenrod2	gray8	gray56	gray68	lavenderblush1	linen	salmon1	
blue3	goldenrod1	gray7	gray55	gray67	lavenderblush	limegreen	salmon	thistle3
	goldenrod	gray6	gray54	gray66	lavender	lightyellow4	saddlebrown	thistle2
	gold4	gray5	gray53	gray65	khaki4	lightyellow3	royalblue4	thistle1
	gold3	gray4	gray52	gray64	khaki3	lightyellow2	royalblue3	tan4
black	gold2	gray3	gray51	gray63	khaki2	lightyellow1	royalblue2	tan3
	gold1	gray2	gray50	gray62	khaki1	lightyellow	royalblue1	tan2
	gold	gray1	gray49	gray61	khaki	lightsteelblue4	royalblue	tan1
	ghostwhite	gray0	gray48	gray60	ivory3	lightsteelblue3	rosybrown4	tan
bisque4	gainsboro	gray46	gray46	gray58	ivory2	lightsteelblue2	rosybrown3	steelblue4
	forestgreen	gray45	gray45	gray57	ivory1	lightsteelblue1	rosybrown2	steelblue3
	floralwhite	gray44	gray44	gray56	ivory	lightslategray	rosybrown1	steelblue2
	firebrick4	gray43	gray43	gray55	indianred4	lightslategray	red4	steelblue1
azure4	firebrick3	gray42	gray42	gray54	indianred3	lightslateblue	red3	steelblue
	firebrick2	gray41	gray41	gray53	indianred2	lightskyblue4	red2	springgreen4
	firebrick1	gray40	gray40	gray52	indianred1	lightskyblue3	red1	springgreen3
	firebrick	gray39	gray39	gray51	indianred	lightskyblue2	red	springgreen2
aquamarine4	dodgerblue4	gray38	gray38	gray50	holipink4	lightskyblue1	purple4	springgreen1
	dodgerblue3	gray37	gray37	gray49	holipink3	lightskyblue	purple3	springgreen
	dodgerblue2	gray36	gray36	gray48	holipink2	lightsalmon4	purple2	springgreen
	dodgerblue1	gray35	gray35	gray47	holipink1	lightsalmon3	purple1	springgreen
antiquewhite4	dodgerblue	gray34	gray34	gray46	holipink	lightsalmon2	purple	springgreen
	dimgray	gray33	gray33	gray45	honeydew4	lightsalmon1	powderblue	
	dimgray	gray32	gray32	gray44	honeydew3	lightsalmon1	plum4	slategray4
	deepskyblue4	gray31	gray31	gray43	honeydew2	lightsalmon	plum3	slategray3
antiquewhite1	deepskyblue3	gray30	gray30	gray42	honeydew1	lightsalmon	plum2	slategray2
	deepskyblue2	gray29	gray29	gray41	honeydew	lightsalmon	plum1	
	deepskyblue1	gray28	gray28	gray40	gray100	lightsalmon	plum	
	deepskyblue					lightsalmon		