

---

# Einführung R

**WS 2023-24**

**DI Emil Marinov**

# Übersicht

1.	Warum R?	3 – 4
2.	Basics - Grundlegende Funktionen und Datentypen	5 – 14
3.	Datentransformation mit dplyr und tidyr	15 – 23
4.	Visualisierung mit ggplot2	24 – 33

# Warum R?



# Was ist R?

- Software für statistische Analysen und Visualisierungen
- Programmiersprache
- open source, freie Software, lizenziert unter GNU General Public Licence
- läuft unter Windows, MacOS und UNIX
- <https://www.r-project.org/>
- erweiterbar durch zusätzliche Bibliotheken (packages):  
> 18.000 frei verfügbar auf CRAN (Comprehensive R Archive Network: <https://cran.r-project.org/>)

# Warum R?



- reproduzierbare Ergebnisse
- leicht automatisierbar und integrierbar mit anderen Technologien
- einfacher Umgang mit großen Datensätzen aus verschiedenen Quellen (Excel ist beschränkt auf 1.048.576 Zeilen und 16.384 Spalten)
- Erstellung von automatisierten Berichten und interaktiven Apps
- kostenlos
- weit verbreitet und dadurch zahlreiche Handbücher, Kurse, Internet-Foren vorhanden
- großer Funktionsumfang im Bereich der statistischen Auswertungen und Visualisierungen
- einfach erweiterbar durch R Packages

(siehe auch <https://www.northeastern.edu/graduate/blog/r-vs-excel/>)

# Basics

# Basics

Definition von Variablen:

```
x <- 5  
y <- "apple"
```

Operatoren:

+	plus
-	minus
*	mal
/	geteilt
^	hoch
==	ist gleich
!=	ist nicht gleich

<=	kleiner gleich
>=	größer gleich
&	und (elementweise)
	oder (elementweise)
%in%	ist in
is.na(x)	Test auf fehlenden Wert

# Basics

## Mathematische Funktionen

<b>max(x)</b>	größtes Element
<b>min(x)</b>	kleinstes Element
<b>sum(x)</b>	Summe aller Elemente
<b>round(x, n)</b>	Runden auf n Dezimalstellen
<b>signif(x, n)</b>	Runden auf n signifikante Stellen
<b>mean(x)</b>	arithmetisches Mittel
<b>median(x)</b>	Median
<b>var(x)</b>	Varianz
<b>sd(x)</b>	Standardabweichung



# Basics

## Datentypen und Umwandlungen

<b>Logical</b>	<b>as.logical</b>	logische Werte (TRUE, FALSE)
<b>Integer</b>	<b>as.integer</b>	ganze Zahlen
<b>Numeric</b>	<b>as.numeric</b>	reelle Zahlen
<b>Character</b>	<b>as.character</b>	Zeichenkette (String, Text)
<b>Factor</b>	<b>as.factor</b>	kategoriale Variable mit vorgegebenen Ausprägungen
<b>Vector</b>	<b>as.vector</b>	Vektor mit einheitlichen Datentyp
<b>List</b>	<b>as.list</b>	Liste mit Daten (unterschiedlicher Datentyp möglich)
<b>Matrix</b>	<b>as.matrix</b>	Tabelle mit einheitlichem Datentyp
<b>Data frame</b>	<b>as.dataframe</b>	Tabelle, deren Spalten unterschiedliche Datentypen haben können

# Vektoren

## Vektoren erzeugen

<code>x &lt;- c(1, 3, 7, 5, 9)</code>	1 3 7 5 9	Vektor durch Aufzählen definieren
<code>1:5</code>	1 2 3 4 5	ganzzahlige Zahlenfolge
<code>seq(2,3, by=0.5)</code>	2 2.5 3	Zahlenfolge erzeugen
<code>rep(1:2, times=3)</code>	1 2 1 2 1 2	Wiederholung von Zahlen

## auf Elemente von Vektoren zugreifen

<code>x[4]</code>	5	4. Element eines Vektors
<code>x[-4]</code>	1 3 7 9	alle Element außer dem 4. Element
<code>x[2:4]</code>	3 7 5	2. – 4. Element eines Vektors
<code>x[c(1, 5)]</code>	1 9	1. und 5. Element
<code>x[x &lt; 5]</code>	1 3 7	alle Elemente bis zum ersten Datenpunkt 5
<code>x[x %in% c(1, 5)]</code>	1 5	alle Elemente, die in (1, 5) enthalten sind

# Vektoren

```
x <- c(1, 3, 7, 5, 3, 4, 5, 3)
```

## Funktionen für Vektoren

<b>sort(x)</b>	1 3 3 3 4 5 5 7	Werte sortieren
<b>table(x)</b>	1 3 4 5 7 1 3 1 2 1	absolute Häufigkeiten der Werte
<b>rev(x)</b>	3 5 4 3 5 7 3 1	Vektor umdrehen
<b>unique(x)</b>	1 3 7 5 4	verschiedene Werte eines Vektors

# Listen

## Liste definieren

```
l <- list(x=1:5, y=c("a", "b"))
```

## auf Elemente einer Liste zugreifen

<code>l[[2]]</code>	"a" "b"	zweites Element einer List
<code>l[1]</code>	1 2 3 4 5	neue Liste mit erstem Element der Liste l
<code>l\$x</code>	1 2 3 4 5	auf Element mit Namen zugreifen
<code>l["y"]</code>	"a" "b"	neue Liste mit dem Element "y"

## Schleife über Listenelemente

```
l <- lapply(list_name, function(x){....})
```

# Matrizen

## Matrix definieren

```
x <- 1:9  
m <- matrix(x, nrow = 3, ncol = 3)  
m
```

	[,1]	[,2]	[,3]
[1,]	1	4	7
[2,]	2	5	8
[3,]	3	6	9

## auf Elemente einer Matrix zugreifen

m[2, ]	2 5 8	Auswahl einer Zeile
m[ , 1]	1 2 3	Auswahl einer Spalte
m[2, 3]	8	Auswahl eines Elements

## Matrix-operationen

t (m)	Matrix transponieren
m %*% n	Matrix multiplizieren

# Dataframes

## Matrix definieren

```
df <- data.frame(x = 1:3,  
                 y = c("a", "b", "c"))
```

```
> df  
  x y  
1 1 a  
2 2 b  
3 3 c
```

## auf Elemente eines Dataframes zugreifen

<code>df\$x</code>	1 2 3	Spalte mit Namen auswählen
<code>df[[2]]</code>	"a" "b" "c"	Spalte mit Index auswählen
<code>df[1,2]</code>	"a"	eine Element mit Zeilen und Spaltenindex auswählen

## Funktionen für Dataframes

<code>head(df)</code>	ersten 6 Zeilen eines Dataframes
<code>nrow(df)</code>	Zeilenanzahl

# Kontrollstrukturen und Funktionen

## FOR-Schleife

```
for(i in 1:n){  
    ....  
}
```

## WHILE-Schleife

```
while(condition){  
    ....  
}
```

## IF-Abfrage

```
if(condition){  
    ....  
} else {  
    ....  
}
```

## Funktion definiern

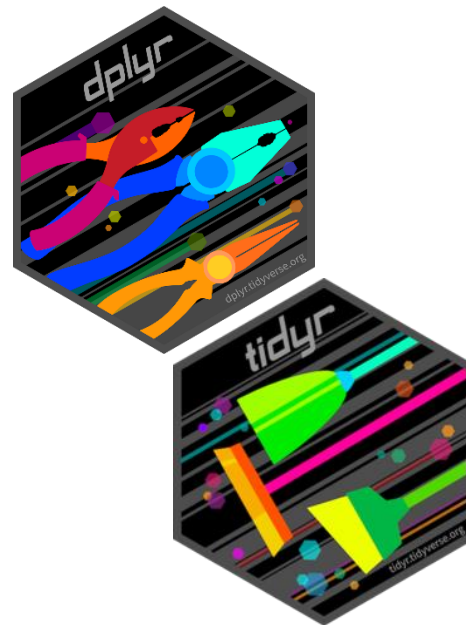
```
func_name <- function(var){  
    ....  
    return(new_var)  
}
```

# Einlesen und Schreiben von Dateien

Dateiformat	Funktion	Beschreibung
csv	<b>read.csv</b> ("file.csv") <b>write.csv</b> (df, "file.csv")	CSV Datei
xlsx	<b>library</b> (openxlsx) <b>read.xlsx</b> ("file.xlsx", sheet = 1) <b>write.xlsx</b> (df, "file.xlsx")	Excel Datei
Rdata	<b>load</b> ("file.RData") <b>save</b> (df, "file.Rdata")	R spezifische Datei
fst	<b>library</b> (fst) <b>read_fst</b> ("file.fst") <b>write_fst</b> (df, "file.fst")	R spezifischer "fast storage" Dateityp (nur für Tabellen als data frame möglich)



# Datentransformationen mit dplyr und tidyr





# dplyr

- Package laden: `library(dplyr)`
- **dplyr** stellt einfache Funktionen für die wichtigsten Aufgaben bei der Datenmanipulation zur Verfügung.
- Voraussetzung: “**tidy data**”
  - jede Variable ist in einer eigenen Spalte
  - jede Beobachtung (Fall) ist in einer eigenen Zeile
- dplyr Funktionen verwenden „**pipes**“:  
`x %>% func(y)` entspricht `func(x, y)`
- Pipes können verschachtelt sein. Das Ergebnis von einem Schritt dient als Input für den nächsten Schritt. `x %>% func1(y) %>% func2(z)`



# dplyr

## Zeilen manipulieren

<code>df %&gt;% <b>filter</b>(Var1 &gt; 1)</code>	Zeilen basierend auf Variablenwerten auswählen
<code>df %&gt;% <b>arrange</b>(Var1)</code> <code>df %&gt;% <b>arrange</b>(<b>desc</b>(Var1))</code>	Tabelle sortieren (aufsteigend oder absteigend nach Variablenwerten)
<code>df %&gt;% <b>distinct</b>()</code>	doppelte Zeilen entfernen
<code>df %&gt;% <b>slice</b>(10:15)</code>	Zeilen mittels Index auswählen

## Spalten manipulieren

<code>df %&gt;% <b>select</b>(Var1, Var2)</code>	Spalten auswählen
<code>df %&gt;% <b>mutate</b>(Var3 = Var1 + Var2)</code>	neue Spalten berechnen
<code>df %&gt;% <b>rename</b>(Var_new = Var_old)</code>	Spalte umbenennen
<code>df %&gt;% <b>pull</b>(Var1)</code>	Werte einer Spalte als Vektor extrahieren



# dplyr

## Fälle zusammenfassen

<code>df %&gt;% <b>summarize</b>(avg = mean(Var1))</code>	Tabelle aggregieren
<code>df %&gt;% <b>count</b>(Var1)</code>	Häufigkeiten einer Variable bestimmen
<code>df %&gt;% <b>summarize_all</b>(mean)</code>	Funktion auf alle Spalten anwenden

## Fälle gruppieren

<code>df %&gt;% <b>group_by</b>(Var1) %&gt;% <b>summarize</b>(avg = mean(Var2))</code>	Zeilen gruppieren und Werte nach Gruppen getrennt auswerten
<code>grouped_df %&gt;% <b>ungroup</b>()</code>	Dataframe degruppieren



# dplyr

## Tabellen kombinieren



<code>df1 %&gt;% left_join(df2, by = "Key1")</code>	<b>left join</b> von Dataframes
<code>df1 %&gt;% right_join(df2, by = "Key1")</code>	<b>right join</b> von Dataframes



# dplyr

## Tabellen kombinieren

`inner_join()` 

`full_join()` 

<pre>df1 %&gt;% inner_join(df2, by = "Key1")</pre>	<b>inner join</b> von Dataframes
<pre>df1 %&gt;% full_join(df2, by = "Key1")</pre>	<b>outer join</b> von Dataframes
<pre>df_list = list(name1 = df1, name2 = df2, name3 = df3) bind_rows(df_list, .id = "Name")</pre>	Liste von Dataframes zusammenfügen und ID-Spalte mit Listennamen erzeugen



# tidyr

- Package laden:

```
library(tidyr)
```

- Tabellen umformatieren:

```
gather(df, key, value)
```

“wide”-Format → “long”-Format

```
spread(df, key, value)
```

“long”-Format → “wide”-Format

table4a

country	1999	2000
A	0.7K	2K
B	37K	80K
C	212K	213K



country	year	cases
A	1999	0.7K
B	1999	37K
C	1999	212K
A	2000	2K
B	2000	80K
C	2000	213K

key value

```
gather(table4a, `1999`, `2000`,  
       key = "year", value = "cases")
```

table2

country	year	type	count
A	1999	cases	0.7K
A	1999	pop	19M
A	2000	cases	2K
A	2000	pop	20M
B	1999	cases	37K
B	1999	pop	172M
B	2000	cases	80K
B	2000	pop	174M
C	1999	cases	212K
C	1999	pop	1T
C	2000	cases	213K
C	2000	pop	1T

key value



country	year	cases	pop
A	1999	0.7K	19M
A	2000	2K	20M
B	1999	37K	172M
B	2000	80K	174M
C	1999	212K	1T
C	2000	213K	1T

```
spread(table2, type, count)
```

# Datenvisualisierung mit **ggplot2**





# ggplot2

- Package laden: `library(ggplot2)`
- **ggplot2** baut auf einer Grammatik von Grafikelementen auf.
- Jede Grafik wird schrittweise aus mehreren Schichten von Komponenten gebaut.



```
ggplot(data, aes(...)) + geom_point(...)  
                        + facet_wrap(...)  
                        + .....
```

# ggplot2

- Daten und „Aesthetics“: `ggplot(data, aes(...)) + ...`
- Optionen für aes():

<code>x = Var1</code>	Variable zu x-Achse zuweisen
<code>y = Var2</code>	Variable zu y-Achse zuweisen
<code>color = Var3</code>	Variable zu Linienfarbe zuweisen
<code>fill = Var4</code>	Variable zu Füllfarbe zuweisen
<code>size = Var5</code>	Variable zu Größe zuweisen
<code>shape = Var6</code>	Variable zu Form zuweisen

# ggplot2

## Geometrien

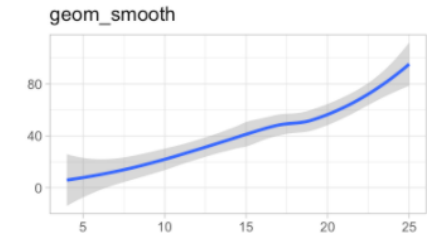
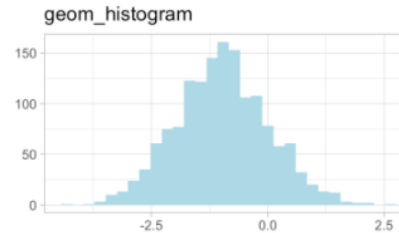
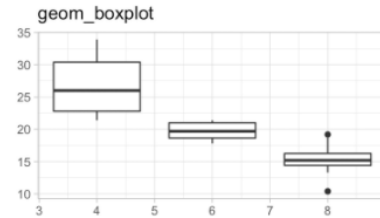
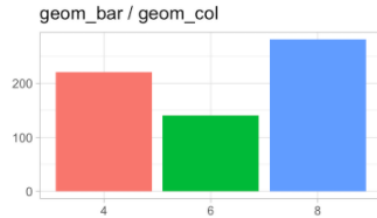
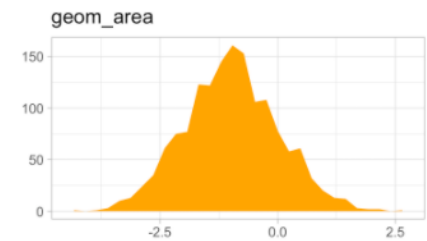
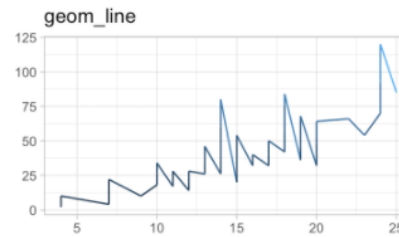
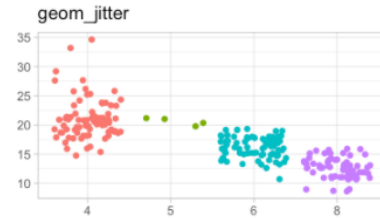
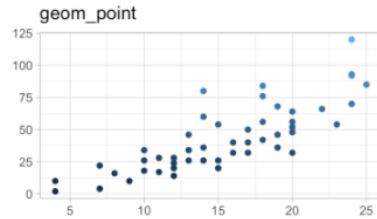
```
ggplot(data, aes(...)) + geom_...
```

+ <code>geom_point()</code>	Punktdiagramm
+ <code>geom_line()</code>	Liniendiagramm
+ <code>geom_bar()</code>	Säulendiagramm aus Rohdaten
+ <code>geom_col()</code>	Säulendiagramm aus Häufigkeitswerten
+ <code>geom_histogram()</code>	Histogramm
+ <code>geom_boxplot()</code>	Boxplot
+ <code>geom_density()</code>	Verteilungsdichte
+ <code>geom_vline()</code>	vertikale Linie
+ <code>geom_hline()</code>	horizontale Linie

# ggplot2



## Geometrien

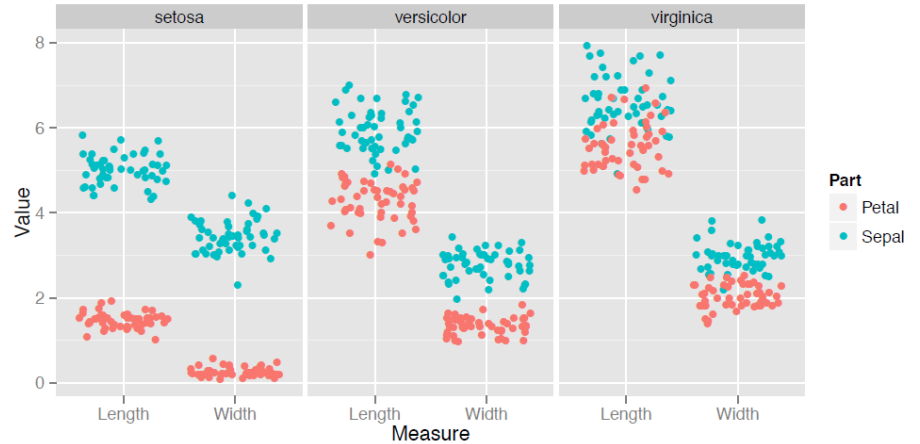


# ggplot2



## Facetten

```
ggplot(data, aes(...)) +  
  geom_... +  
  facet_...
```



+ `facet_grid(x ~ y)`

Facetten in Zeilen und Spalten

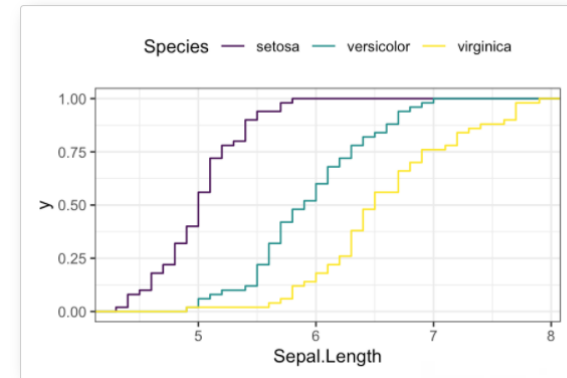
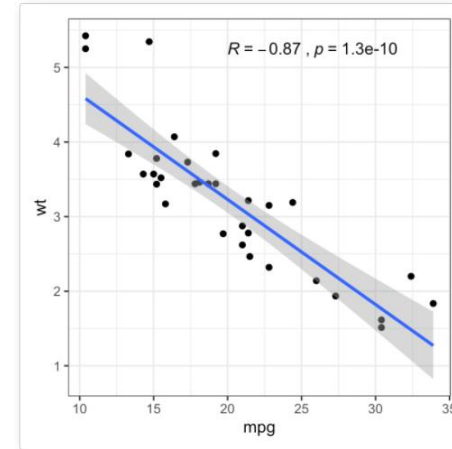
+ `facet_wrap(. ~ x)`

Facetten in einem rechteckigen Layout

# ggplot2

## Statistiken

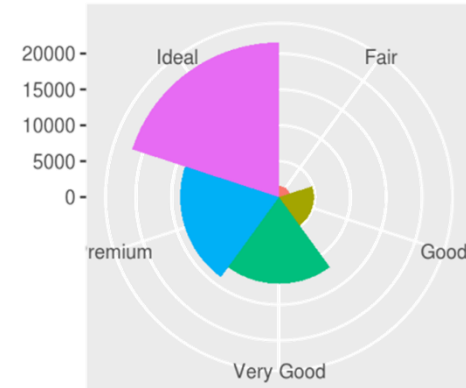
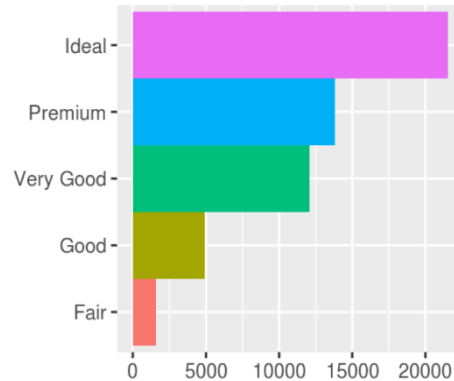
+ <code>geom_smooth()</code>	Kurven glätten
+ <code>stat_ecdf()</code>	empirische Verteilungsfunktion
+ <code>stat_summary()</code>	Zusammenfassung von y-Werten für verschiedene x-Werte



# ggplot2

## Koordinaten

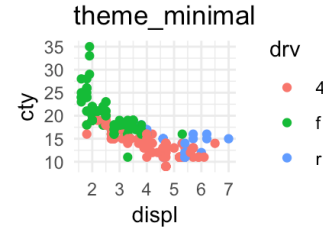
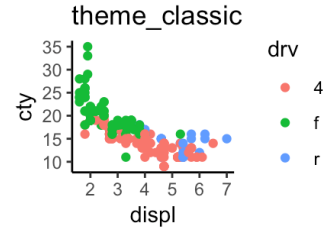
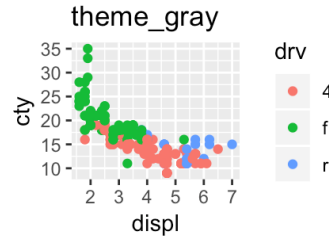
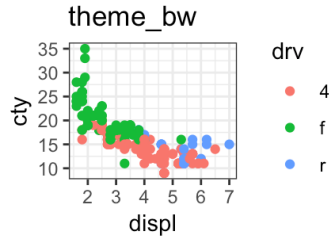
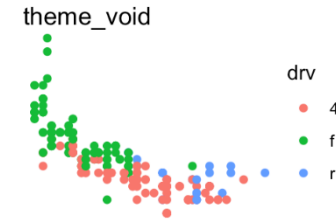
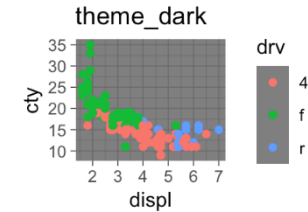
+ <code>coord_flip()</code>	Achsen vertauschen
+ <code>coord_polar()</code>	Polarkoordinaten



# ggplot2

## Themen

+ <code>theme_gray()</code>	grauer Hintergrund (default)
+ <code>theme_bw()</code>	schwarz-weiß
+ <code>theme_minimal()</code>	minimal
+ <code>theme_void()</code>	leer





# ggplot2

## Skalen

+ <code>scale_color_manual(values = c(...))</code>	Linienfarben eines Plots verändern
+ <code>scale_fill_manual(values = c(...))</code>	Füllfarben eines Plots verändern
+ <code>scale_x_log10()</code>	logarithmische Skala auf der x-Achse

