

Prof. Dr. Rainer Stollhoff

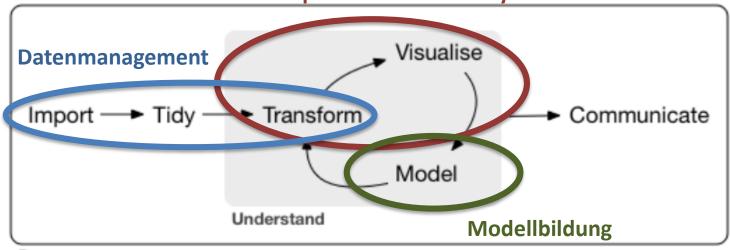
Vgl.

R for Data Science, Grolemund & Wickham, http://r4ds.had.co.nz/exploratory-data-analysis.html

### Übersicht



#### **Explorative Datenanalyse**



Program



- Beschreibung des Datensatzes
- Verteilungen der Werte einzelner Variablen (Univariate Analysen)
- Zusammenhänge zwischen Variablen (Bivariate Analysen)



#### Beschreibung des Datensatzes

- Typische Fragen:
  - Wieviele Beobachtungen sind im Datensatz?
  - Wieviele Variablen sind im Datensatz
  - Welchen Datentyp haben die Variablen und was für Wertebereiche?
  - View() zur Anzeige des gesamten Datensatzes und head() für die ersten Zeilen
  - dim() zur Anzeige der Zeilen und Spalten
  - str () für Informationen zu den enthaltenen Variablen
  - summary() für einen Überblick über die Wertebereiche der Variablen
- Verteilungen der Werte einzelner Variablen (Univariate Analysen)
- Zusammenhänge zwischen Variablen (Bivariate Analysen)

# Beschreibung des Datensatzes

```
> ## Zeigt die ersten paar Zeilen
> head(mpg)
# A tibble: 6 x 11
  manufacturer model displ year
                                   cyl trans
                                                                  hwy fl
                                                                            class
               <chr> <dbl> <int> <int> <chr>
                                                    <chr> <int> <int> <chr> <chr>
                       1.8 <u>1</u>999
                                                                   29 p
1 audi
                                      4 auto(15)
                                                                             compact
                            <u>1</u>999
 audi
                       1.8
                                      4 manual(m5) f
                                                                   29 p
                                                                            compact
                             2008
                                      4 manual(m6) f
 audi
                                                                   31 p
                                                                            compact
4 audi
                             2008
                                      4 auto(av) f
                                                                   30 p
               a4
                                                                            compact
                       2.8
                            1999
                                      6 auto(15) f
                                                                   26 p
5 audi
                                                                            compact
                       2.8
                            1999
                                      6 manual(m5) f
                                                                   26 p
6 audi
                                                                            compact
> ## Größe des Datensatzes
> dim(mpg)
[1] 234 11
> ## Struktur der Variablen
> str(mpg)
tibble [234 x 11] (S3: tbl_df/tbl/data.frame)
$ manufacturer: chr [1:234] "audi" "audi" "audi" "audi" ...
               : chr [1:234] "a4" "a4" "a4" "a4" ...
 $ model
 $ displ
               : num [1:234] 1.8 1.8 2 2 2.8 2.8 3.1 1.8 1.8 2 ...
               : int [1:234] 1999 1999 2008 2008 1999 1999 2008 1999 1999 2008 ...
 $ year
 $ cy1
               : int [1:234] 4 4 4 4 6 6 6 4 4 4 ...
               : chr [1:234] "auto(15)" "manual(m5)" "manual(m6)" "auto(av)" ...
: chr [1:234] "f" "f" "f" ...
 $ drv
               : int [1:234] 18 21 20 21 16 18 18 18 16 20 ...
 $ ctv
               : int [1:234] 29 29 31 30 26 26 27 26 25 28 ...
 $ hwv
               : chr [1:234] "p" "p" "p" "p" ...
: chr [1:234] "compact" "compact" "compact" ...
 $ f1
> ## Statistische Zusammenfassung des Wertebereichs der Variablen
> summary(mpg)
 manufacturer
                        model
                                            displ
                                                              year
                                                                             cyl
                                                                                            trans
 Length: 234
                    Length: 234
                                        Min. :1.600
                                                        Min. :1999
                                                                        Min. :4.000
                                                                                         Length: 234
 Class :character
                    Class :character
                                        1st Qu.:2.400
                                                         1st Qu.:1999
                                                                        1st Qu.:4.000
                                                                                         Class :character
                                        Median :3.300
                                                         Median:2004
 Mode :character
                    Mode :character
                                                                        Median :6.000
                                                                                         Mode :character
                                        Mean :3.472
                                                         Mean : 2004
                                                                        Mean :5.889
                                        3rd Qu.:4.600
                                                         3rd Qu.:2008
                                                                        3rd Qu.:8.000
                                              :7.000
                                                         Max. :2008
                                                                        Max. :8.000
                                        Max.
     drv
                          cty
                                          hwy
                                                                            class
 Lenath: 234
                          : 9.00
                                            :12.00
                                                     Lenath: 234
                                                                         Length: 234
 Class :character
                    1st Qu.:14.00
                                     1st Qu.:18.00
                                                     Class :character
                                                                         Class:character
 Mode :character
                    Median :17.00
                                     Median :24.00
                                                     Mode :character
                                                                         Mode :character
                                            :23.44
                           :16.86
                                     Mean
                    3rd Qu.:19.00
                                     3rd Qu.:27.00
                           :35.00
                                     Max.
                                          :44.00
>
```





- Beschreibung des Datensatzes
- Verteilungen der Werte einzelner Variablen (Univariate Analysen)
  - Typische Fragen:
    - Welche Werte nimmt die Variable an?
    - Was sind typische / untypische Werte?
    - Gibt es Häufungen?
    - arrange() zum Sortieren der Daten
    - summarise() zum Aggregieren der Daten z.B. mean, max, min
    - count () zum Berechnen von Häufigkeiten
    - geom bar () bei kategorischen Variablen
    - geom histogram(), geom freqpoly() beistetigen Variablen
- Zusammenhänge zwischen Variablen (Bivariate Analysen)

### **Univariate Analyse**

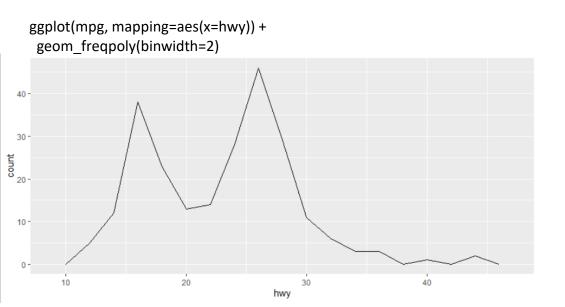


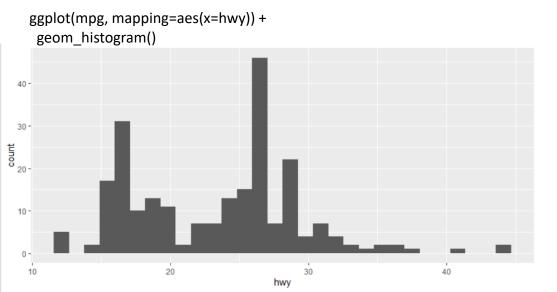
```
> ## Wer hat die niedrigste Reichweite?
> arrange(mpg, cty)
# A tibble: 234 x 11
   manufacturer model
                                     displ year
                                                  cyl trans
                                                                          cty
                                                                                hwy fl
                                                                                           class
   <chr>
                <chr>
                                     <db1> <int> <int> <chr>
                                                                  <chr> <int> <int> <chr>
                                                                                          <chr>
                                           2008
                                                                                 12 e
  dodge
                dakota pickup 4wd
                                      4.7
                                                     8 auto(15)
                                                                                           pickup
                durango 4wd
  dodge
                                           2008
                                                     8 auto(15)
                                                                                 12 e
                                                                                           suv
                ram 1500 pickup 4wd
                                                                                           pickup
  dodge
                                      4.7
                                           2008
                                                     8 auto(15)
                                                                                 12 e
  dodge
                ram 1500 pickup 4wd
                                      4.7
                                           2008
                                                     8 manual(m6) 4
                                                                                 12 e
                                                                                           pickup
                grand cherokee 4wd
                                      4.7
                                            2008
                                                     8 auto(15)
                                                                                 12 e
 5 jeep
                                                                                           suv
 6 chevrolet
                c1500 suburban 2wd
                                      5.3
                                           2008
                                                     8 auto(14)
                                                                           11
                                                                                 15 e
                                                                                           suv
  chevrolet
                k1500 tahoe 4wd
                                      5.3
                                           2008
                                                     8 auto(14)
                                                                           11
                                                                                 14 e
                                                                                           suv
 8 chevrolet
                k1500 tahoe 4wd
                                      5.7
                                           1999
                                                     8 auto(14)
                                                                           11
                                                                                 15 r
                                                                                           suv
 9 dodae
                caravan 2wd
                                      3.3 2008
                                                     6 auto(14)
                                                                           11
                                                                                 17 e
                                                                                           minivan
10 dodge
                dakota pickup 4wd
                                      5.2 <u>1</u>999
                                                     8 manual(m5) 4
                                                                           11
                                                                                 17 r
                                                                                           pickup
# ... with 224 more rows
                                                                                                          <chr>>
> ## Und wer die höchste?
> arrange(mpg, desc(cty))
                                                                                                        1 audi
# A tibble: 234 x 11
                                         cyl trans
   manufacturer model
                           displ
                                  year
                                                         drv
                                                                 cty
                                                                       hwy fl
                                                                                  class
                                                                                                        3 dodge
   <chr>
                <chr>>
                           <db1> <int> <int> <chr>
                                                         <chr> <int> <int> <chr> <chr>
                                                                                                        4 ford
                                  1999
 1 volkswagen
                new beetle 1.9
                                           4 manual(m5) f
                                                                  35
                                                                        44 d
                                                                                  subcompact
                                                                                                        5 honda
 2 volkswagen
                jetta
                             1.9
                                  1999
                                           4 manual(m5) f
                                                                  33
                                                                        44 d
                                                                                  compact
                                                                                                        6 hyundai
 3 volkswagen
                new beetle 1.9
                                  1999
                                           4 auto(14) f
                                                                        41 d
                                                                                  subcompact
                                                                                                          jeep
 4 honda
                             1.6
                                  1999
                                           4 manual(m5) f
                                                                  28
                                                                        33 r
                civic
                                                                                  subcompact
                                  2008
                                           4 manual(m5) f
                                                                        37 r
 5 toyota
                corolla
                             1.8
                                                                                  compact
 6 honda
                civic
                             1.8
                                  2008
                                           4 manual(m5) f
                                                                  26
                                                                        34 r
                                                                                  subcompact
                                  <u>1</u>999
                                                                  26
                                                                        35 r
                                                                                                       10 mercury
 7 toyota
                corolla
                             1.8
                                           4 manual(m5) f
                                                                                  compact
                                                                                                       11 nissan
 8 toyota
                corolla
                                  2008
                                           4 auto(14) f
                                                                  26
                                                                        35 r
                                                                                  compact
9 honda
                civic
                                  1999
                                           4 manual(m5) f
                                                                  25
                                                                        32 r
                                                                                  subcompact
                                                                                                       12 pontiac
10 honda
                             1.8
                                  2008
                                           4 auto(15)
                                                                  25
                                                                        36 r
                civic
                                                                                  subcompact
                                                                                                       13 subaru
# ... with 224 more rows
```

```
> ## Was ist die durchschnittliche Reichweite?
> summarise(mpg,mean(cty))
# A tibble: 1 x 1
   `mean(cty)`
         \langle db 1 \rangle
         16.9
> ## Was sind die Mittelwerte?
> summarise_if(mpg,is.numeric,funs(mean))
# A tibble: 1 x 5
  displ year cyl cty
  <db1> <db1> <db1> <db1> <db1>
1 3.47 <u>2</u>004. 5.89 16.9 23.4
> ## Wieviele Autos gibt es pro Hersteller?
> count(mpg, manufacturer)
# A tibble: 15 x 2
   manufacturer
                     n
                 <int>
                    18
 2 chevrolet
                    19
                    37
                    25
                    14
 8 land rover
 9 lincoln
                    13
                     5
                    14
                    34
14 toyota
                    27
15 volkswagen
```

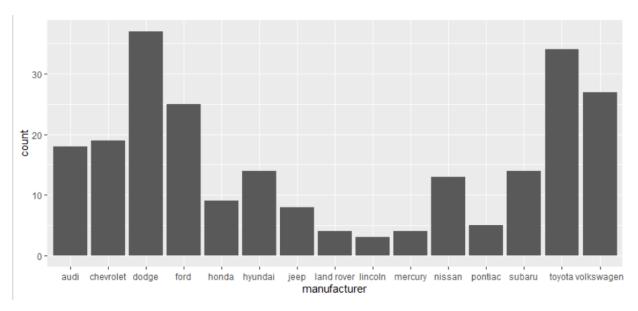
# Univariate Analyse







ggplot(mpg, mapping=aes(x=manufacturer)) +
 geom\_bar()





- Beschreibung des Datensatzes
- Verteilungen der Werte einzelner Variablen (Univariate Analysen)

#### Zusammenhänge zwischen Variablen (Bivariate Analysen)

- Typische Fragen:
  - Gibt es Zusammenhänge zwischen Variablen?
  - Wenn ja, sind diese positiv/negativ, sind diese stark oder schwach ausgeprägt?
  - Gibt es nichtlineare Zusammenhänge?
  - summarise() in Verbindung mit group by()
  - cor() und cov() zum Berechnen statistischer Zusammenhangsmaße
  - stat\_bin(x=stet,color=kat) oder geom\_boxplot(x=kat,y=stet)
     für x stetig und y kategorisch
  - geom\_count(x=kat1, y=kat2) für x und y kategorisch
  - geom\_point(x=stet1, y=stet2) für x und y stetig

### Bivariate Analyse



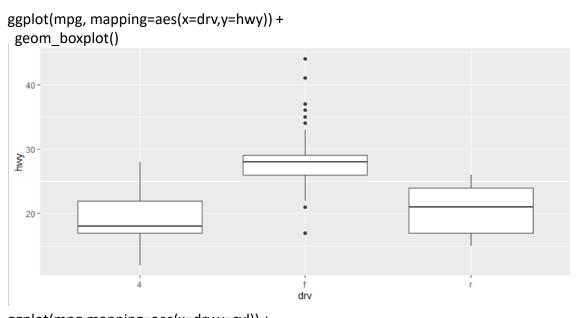
```
> ## Minimale, Maximale und durchschnittliche Reichweite je Hersteller
> summarise(
    group_by(mpg,manufacturer),
min(cty), max(cty), mean(cty))
`summarise()` ungrouping output (override with `.groups` argument)
# A tibble: 15 x 4
   manufacturer `min(cty)` `max(cty)` `mean(cty)`
   <chr>
                                   <int>
                                                <db1>
                       <int>
 1 audi
                          15
                                      21
                                                 17.6
 2 chevrolet
                          11
                                                 15
 3 dodge
                                                 13.1
 4 ford
                                                 14
 5 honda
                                                 24.4
 6 hyundai
                                      21
                                                 18.6
                                      17
 7 jeep
                                                 13.5
 8 land rover
                                      12
                                                 11.5
 9 lincoln
                                                 11.3
10 mercury
                                      14
                                                 13.2
                                      23
                                                 18.1
11 nissan
12 pontiac
                                      18
                                                 17
13 subaru
                                      21
                                                 19.3
                          11
                                      28
14 toyota
                                                 18.5
15 volkswagen
                                      35
                                                 20.9
```

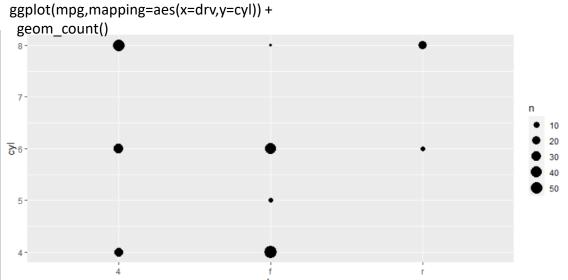
- > ## Korrelationen zwisschen den numerischen Werten
- > cor(select\_if(mpg,is.numeric))

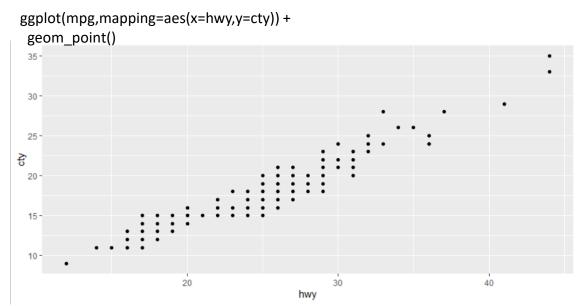
	displ	year	cyl	cty	hwy
displ	1.0000000	0.147842816	0.9302271	-0.79852397	-0.766020021
year	0.1478428	1.000000000	0.1222453	-0.03723229	0.002157643
cyl	0.9302271	0.122245347	1.0000000	-0.80577141	-0.761912354
cty	-0.7985240	-0.037232291	-0.8057714	1.00000000	0.955915914
hwy	-0.7660200	0.002157643	-0.7619124	0.95591591	1.000000000

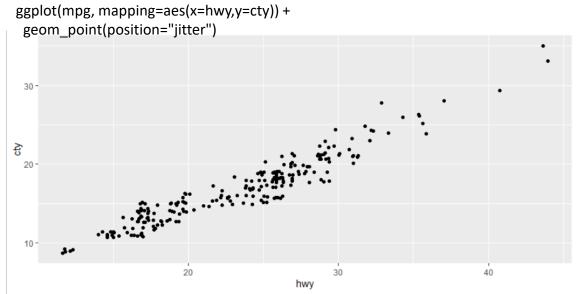
# **Bivariate Analyse**











# Multivariate Analyse

Technische
Hochschule
Wildau
Technical University
of Applied Sciences

ggplot(mpg, mapping=aes(x=hwy,y=displ,shape=drv,color=class)) +
 geom\_point(position="jitter")

