# Datenanalyse



.consulting .solutions .partnership

.msg
systems

ML / AI

Informatik

Statistik /
Mathematik

Datenanalyse

IT-Beratung

Wissenschaft

Fachwissen

# US Census Income Data

Explorative Datenanalyse

adult.data

```
1   39, State-gov, 77516, Bachelors, 13, Never-married, Adm-clerical, Not-in-family, White, Male, 2
2   50, Self-emp-not-inc, 83311, Bachelors, 13, Married-civ-spouse, Exec-managerial, Husband, White
3   38, Private, 215646, HS-grad, 9, Divorced, Handlers-cleaners, Not-in-family, White, Male, 0, 0,
4   53, Private, 234721, 11th, 7, Married-civ-spouse, Handlers-cleaners, Husband, Black, Male, 0, 0
5   28, Private, 338409, Bachelors, 13, Married-civ-spouse, Prof-specialty, Wife, Black, Female, 0,
6   37, Private, 284582, Masters, 14, Married-civ-spouse, Exec-managerial, Wife, White, Female, 0,
7   49, Private, 160187, 9th, 5, Married-spouse-absent, Other-service, Not-in-family, Black, Female
8   52, Self-emp-not-inc, 209642, HS-grad, 9, Married-civ-spouse, Exec-managerial, Husband, White,
9   31, Private, 45781, Masters, 14, Never-married, Prof-specialty, Not-in-family, White, Female, 1
10  42, Private, 159449, Bachelors, 13, Married-civ-spouse, Exec-managerial, Husband, White, Male,
11  37, Private, 280464, Some-college, 10, Married-civ-spouse, Exec-managerial, Husband, Black, Mal
12  30, State-gov, 141297, Bachelors, 13, Married-civ-spouse, Prof-specialty, Husband, Asian-Pac-Is
```

- Erste 12 Instanzen mit
- 15 Variablen

# Daten

| | id | employerKind | fnlwgt | degree | yearsOfEd | maritalStatus | occupation | relationshipRole |
|---|---|---|---|---|---|---|---|---|
| 1 | 39 | State-gov | 77516 | Bachelors | 13 | Never-married | Adm-clerical | Not-in-family |
| 2 | 50 | Self-emp-not-inc | 83311 | Bachelors | 13 | Married-civ-spouse | Exec-managerial | Husband |
| 3 | 38 | Private | 215646 | HS-grad | 9 | Divorced | Handlers-cleaners | Not-in-family |
| 4 | 53 | Private | 234721 | 11th | 7 | Married-civ-spouse | Handlers-cleaners | Husband |
| 5 | 28 | Private | 338409 | Bachelors | 13 | Married-civ-spouse | Prof-specialty | Wife |
| 6 | 37 | Private | 284582 | Masters | 14 | Married-civ-spouse | Exec-managerial | Wife |
| 7 | 49 | Private | 160187 | 9th | 5 | Married-spouse-absent | Other-service | Not-in-family |
| 8 | 52 | Self-emp-not-inc | 209642 | HS-grad | 9 | Married-civ-spouse | Exec-managerial | Husband |
| 9 | 31 | Private | 45781 | Masters | 14 | Never-married | Prof-specialty | Not-in-family |
| 10 | 42 | Private | 159449 | Bachelors | 13 | Married-civ-spouse | Exec-managerial | Husband |
| 11 | 37 | Private | 280464 | Some-college | 10 | Married-civ-spouse | Exec-managerial | Husband |
| 12 | 30 | State-gov | 141297 | Bachelors | 13 | Married-civ-spouse | Prof-specialty | Husband |

*CensusData.R* — *rawData* — 32,561 observations of 15 variables

- Erste 12 Instanzen als Data-Frame mit
- 8 von 15 Variablen

| incomeGroup | academicLvl | incomeMoreThan50K | capitalDeviation | yearsOfEdStdUnits | workingHoursAWeekStdUnits |
|---|---|---|---|---|---|
| <=50K | Bachelor | TRUE | 2.7562412 | 1.13472134 | -0.03542890 |
| <=50K | Bachelor | TRUE | -0.3071748 | 1.13472134 | -2.22211900 |
| <=50K | Highschool | TRUE | -0.3071748 | -0.42005317 | -0.03542890 |
| <=50K | 1 | TRUE | -0.3071748 | -1.19744043 | -0.03542890 |
| <=50K | Bachelor | TRUE | -0.3071748 | 1.13472134 | -0.03542890 |
| <=50K | Master | TRUE | -0.3071748 | 1.52341497 | -0.03542890 |
| <=50K | 1 | TRUE | -0.3071748 | -1.97482769 | -1.97915343 |
| >50K | Highschool | FALSE | -0.3071748 | -0.42005317 | 0.36951371 |
| >50K | Master | FALSE | 3.5009221 | 1.52341497 | 0.77445632 |
| >50K | Bachelor | FALSE | 3.1020895 | 1.13472134 | -0.03542890 |
| >50K | College | FALSE | -0.3071748 | -0.03135955 | 3.20411198 |
| >50K | Bachelor | FALSE | -0.3071748 | 1.13472134 | -0.03542890 |

32,561 observations of 20 variables

5 Sekundärvariablen

Arbeitsstunden je Woche

**Kapitalmehrung**

Häufigkeit

logarithmierte, individuelle Gewinne

Kapitalminderung

Häufigkeit

logarithmierte, individuelle Verluste

Kapitaländerung

Jahre der Bildung
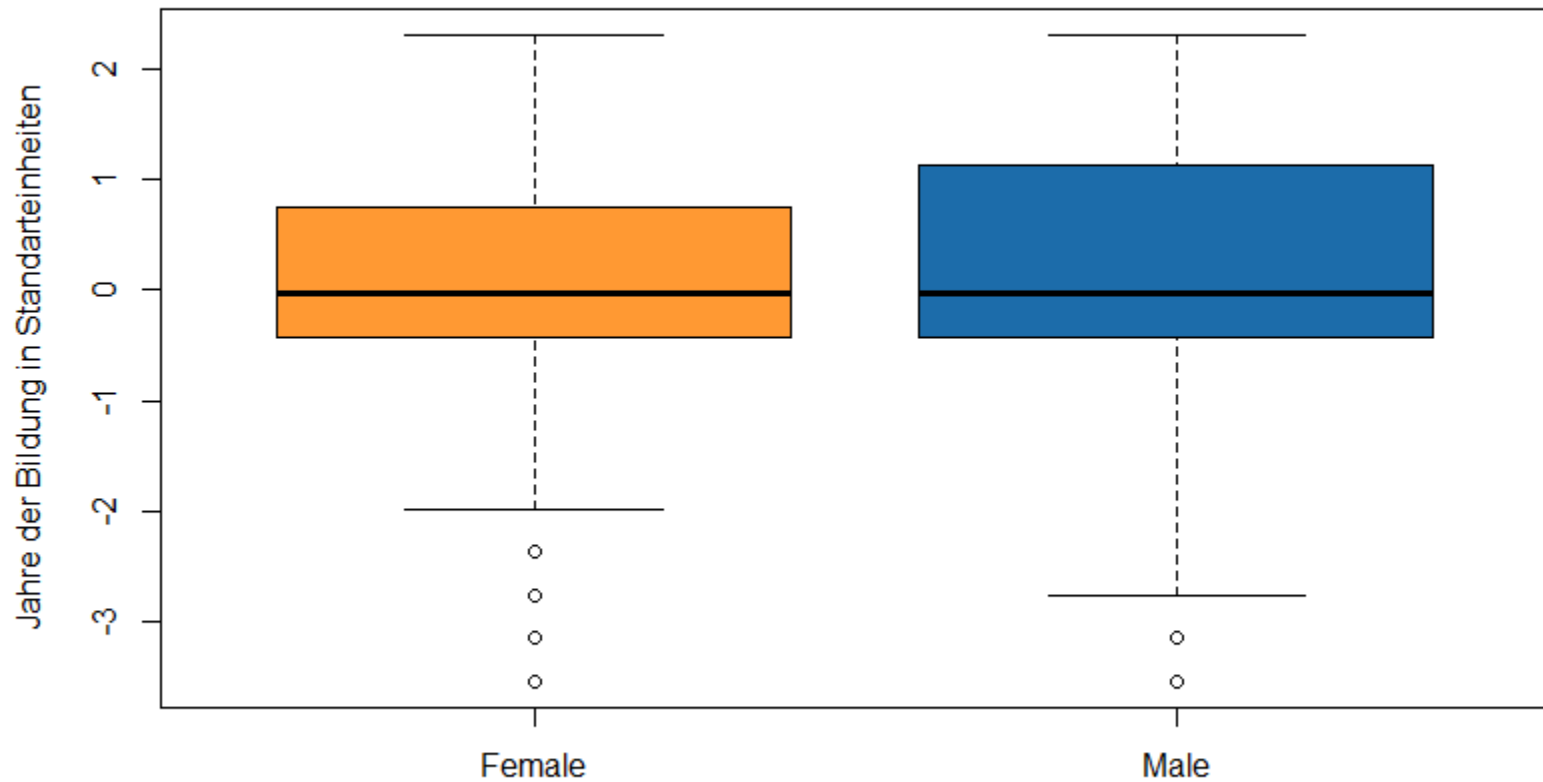
Jahre der Bildung nach Geschlecht

Arbeitsstunden nach Bildungsgrad

Arbeitsstunden nach Anstellungsart

# US Census Income Data

Klassifikation

- ■ Trainingsdaten
- ■ Testdaten

Klassifikation: Sensitivität & Spezifität

|  | tatsächlich nicht hohes Einkommen | tatsächlich hohes Einkommen |
|---|---|---|
| geschätzt nicht hohes Einkommen | 13% | 5% |
| geschätzt hohes Einkommen | 10% | 72% |

# Klassifizierung

Länge der Entscheidungsvektoren

# Anwendungsbeispiele

Stanford University

Prof. Sebastian Thrun mit autonom fahrenden VW Touareg „Stanley"

github
SOCIAL CODING

**LearningR**

Training Dataset for R Beginners

Last updated 3 minutes ago

**AddHealth-Data-Analysis**

The analysis of biases and influencer in attendance of religious services

Last updated a month ago

github.com/danielschulz/LearningR

UNIVERSITY of CALIFORNIA | IRVINE

UCI Machine Learning Repository
archive.ics.uci.edu/ml

kaggle.com

data.gov

Wöchentliche Todesstatistiken

# US Census Income Data

Übersicht: Code-Sektionen Schritt für Schritt

```
1
2   # SETUP WORKSPACE
3
4   library(e1071)
5   set.seed(4711)
6
7   # clean
8   rm(list = ls()[!(ls() %in% PERSISTENT_CONSTANTS)])
9
```

Workspace einrichten

```
19
20  # INIT DATA
21
22  # load data
23  dataLocation = "..\\..\\..\\input\\data\\adult.data"
24  rawData = read.csv2(dataLocation, header=FALSE, encoding="ANSI", sep=",", strip.white=TRUE,
25                      na.strings=c("", " ", "?", " ?"))
26
27  dataColumnHeaders = c("id", "employerKind", "fnlwgt", "degree", "yearsOfEd", "maritalStatus",
28                        "occupation", "relationshipRole", "ethnicity", "sex", "capitalGain",
29                        "capitalLoss", "workingHoursAWeek", "homeland", "incomeGroup")
30  names(rawData) = dataColumnHeaders
31
32  rm(list = c("dataColumnHeaders", "dataLocation"))
33
```

Daten laden, Headernamen zuweisen

```
35
36  # FORMAT DATA
37
38  # format data types
39  rawData$id = as.numeric(rawData$id)
40  rawData$employerKind = as.factor(rawData$employerKind)
41  rawData$degree = as.factor(rawData$degree)
42
43  # assign secondary variable academic level
44  rawData$academicLvl = "none"
45  rawData$academicLvl = as.factor(rawData$academicLvl)
46
47  rawData$academicLvl = ifelse ("Doctorate" == rawData$degree || "Prof-school" == rawData$degree,
48                                "PhD", rawData$academicLvl)
49  rawData$academicLvl = ifelse ("Masters" == rawData$degree, "Master", rawData$academicLvl)
50  rawData$academicLvl = ifelse ("Bachelors" == rawData$degree, "Bachelor", rawData$academicLvl)
51  rawData$academicLvl = ifelse ("Some-college" == rawData$degree, "College", rawData$academicLvl)
52  rawData$academicLvl = ifelse ("HS-grad" == rawData$degree, "Highschool", rawData$academicLvl)
53
```

- Daten-Typen zuweisen
- Sekundärvariablen einfügen

```
54
55  # assign secondary variable income to be more than 50000 USD / yr
56  rawData$incomeMoreThan50K = FALSE
57  rawData$incomeMoreThan50K = as.logical(rawData$incomeMoreThan50K)
58  rawData$incomeMoreThan50K = ifelse ("<=50K" == rawData$incomeGroup, TRUE, rawData$incomeMoreThan50K)
59
60
61  # assign secondary variable capital deviation / difference in standard units
62  rawData$capitalDeviation = rawData$capitalGain - rawData$capitalLoss
63  rawData$capitalDeviation = scale(log(rawData$capitalDeviation + 1))
64
65  # assign secondary variable working hours / wk in standard units
66  rawData$yearsOfEdStdUnits = scale(rawData$yearsOfEd)
67  rawData$workingHoursAWeekStdUnits = scale(rawData$workingHoursAWeek)
68
69  # format data types
70  rawData$maritalStatus = as.factor(rawData$maritalStatus)
71  rawData$occupation = as.factor(rawData$occupation)
72  rawData$relationshipRole = as.factor(rawData$relationshipRole)
73  rawData$ethnicity = as.factor(rawData$ethnicity)
74  rawData$sex = as.factor(rawData$sex)
75  rawData$homeland = as.factor(rawData$homeland)
76  rawData$capitalDeviation = as.numeric(rawData$capitalDeviation)
77  rawData$workingHoursAWeekStdUnits = as.numeric(rawData$workingHoursAWeekStdUnits)
78  rawData$yearsOfEdStdUnits = as.numeric(rawData$yearsOfEdStdUnits)
79
```

- Daten-Typen zuweisen
- Sekundärvariablen einfügen

# US Census Income Data

```
81
82  # DROP COLUMNS
83  dropColumns = c("id", "fnlwgt", "yearsOfEd", "workingHoursAWeek", "capitalGain",
84                  "capitalLoss", "incomeGroup")
85  rawData = rawData[,!(names(rawData) %in% dropColumns)]
86
87  # remove dropping column from workspace value list
88  rm(list = c("dropColumns"))
89
```

Nicht benötigte Spalten entfernen

```
 91
 92   # SAMPLE TRAINING AND TEST DATA
 93   rawData$clazz = sample(1:5, dim(rawData)[1], replace=TRUE)
 94   rawData$clazz = as.factor(rawData$clazz)
 95
 96   data = rawData
 97   data = na.omit(data) # drop missing value instances
 98
 99   train = subset(data, 1 == data$clazz)
100   test = subset(data, 1 != data$clazz)
101
102   dropColumns = c("clazz")
103   train = train[,!(names(train) %in% dropColumns)]
104   test = test[,!(names(test) %in% dropColumns)]
105   data = data[,!(names(data) %in% dropColumns)]
106
107   # remove dropping column from workspace value list
108   rm(list = c("dropColumns", "rawData"))
109
```

Trainings- und Testdaten erzeugen

```
111
112  # TRAIN CLASSIFICATION MODEL SUPPORT VECTOR MACHINES AND EVALUATE ACCURANCY
113  svm = svm(train$incomeMoreThan50K ~ ., train, type="C-classification", probability=TRUE,
114          gamma=0.0001, cost=100000)
115  pr = predict(svm, test, probability=TRUE)
116  # plot(formula=train$capitalDeviation ~ train$workingHoursAWeekStdUnits, data=train)
117  # plot(formula=test$capitalDeviation ~ test$workingHoursAWeekStdUnits, data=test)
118
119  table = table(classifications = pr, test$incomeMoreThan50K)
120  table
121
122  # chisquare = chisq.test(table)
123  # chisquare
124  # summary(chisquare)
125
126
127  sumInTable = 0
128
129 ▾ for (i in c(1:4)) {
130    sumInTable = sumInTable + table[i]
131  }
132 ▾ for (i in c(1:4)) {
133    table[i] = table[i] / sumInTable
134  }
135
136  # prediction accurancy is one the main diagonal table[1] + table[4] or for table t: t_11 + t_22
137  table
138
139  rm(list = c("i", "sumInTable", "chisquare"))
140  |
```

SVM-Classifizierung trainieren und testen

# Resumée

# Datenanalyse

- Google´s Chef-Ökonom Hal Varian
  - „The next sexy job"
  - „The ability to take data – to be able to understand it, to process it, to extract value from it, to communicate it – that´s going to be a hugely important skill."
  - New York Times, 2009

- „Hot new gig in tech" – Fortune

# Vielen Dank für Ihre Aufmerksamkeit

.consulting .solutions .partnership

.msg systems