Universität Stuttgart Institut für Sozialwissenschaften Vorlesung: Statistische Modellbildung II

Dozent: Thomas Krause, M.A. Wintersemester 2018/19

Anna Götze Matrikelnummer: 3421518 Datum: 10.11.2018

Studiengang: EPSF M.A., 1. Semester

Übungsaufgabe 2

1. Wozu werden Standardisierungen durchgeführt und wie wird dabei vorgegangen? Erläutern Sie zudem exemplarisch wozu b* benutzt wird und wie man diesen interpretiert!

Standardisierungen von Regressionskoeffizienten werden durchgeführt, um die Stärke der Koeffizienten innerhalb eines Regressionsmodells untereinander zu vergleichen, unabhängig von der Messeinheit/Messskala. Der Vergleich zwischen Modellen sowie eine direkte Interpretation oder eine Prognose für die abhängige Variable sind mit standardisierten Koeffizienten jedoch nicht möglich. Dabei wird folgendermaßen vorgegangen: der unstandardisierte Regressionskoeffizient (b) einer unabhängigen Variable (x) wird mit der Standardabweichung der Variable (s_x) multipliziert und durch die Standardabweichung der abhängigen Variable (s_y) geteilt (siehe (1)).

$$b^* = b * \frac{s_x}{s_y} \tag{1}$$

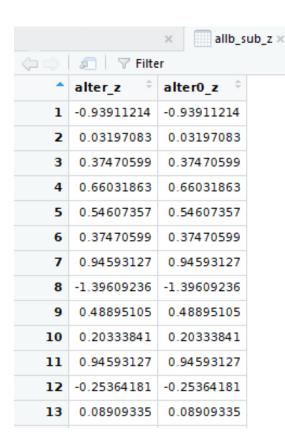
Beispiel: in einem Regressionsmodell mit den (metrischen) Variablen Alter und Einkommen kann nur mit unstandardisierten partiellen Regressionskoeffizienten nicht bestimmt werden, welche Variable einen stärkeren Einfluss auf die abhängige Variable hat, da die unabhängigen Variablen in verschiedenen Einheiten gemessen werden und je nach Skalierung auch unterschiedlich hoch ausfallen können. Bei standardisierten Koeffizienten besteht dieses Problem nicht und es kann eine Aussage darüber getroffen werden, ob das Alter oder das Einkommen einen stärkeren Einfluss auf die abhängige Variable hat.

2. Führen Sie eine z-Standardisierung für die Originalaltersvariable (alter_z) und die auf Null gesetzte Altersvariable (alter_0z) sowie für "unsere" Bildungsvariable (0 bis 4). [Daten: ALLBUS 2014]

```
#Datensatz einlesen
allbus <- read sav("C:/Users/Anna/Documents/Uni S/Statistische-Modellbildung-II-WS1819-mas-
ter/data/allbus2014.sav")
#Subdatensatz erstellen
allb_sub <- select(allbus, V84, V86, V81, V420)
  #Relevante Variablen raussuchen und umbenennen
allb_sub <- rename(allb sub, alter = V84,</pre>
          bildung = V8\overline{6},
geschl = V81,
        einkommen = V420)
  #Niedrigstes Alter auf O setzen, Bildung aufbereiten mit 6 und 7 als missings, Geschlecht
dummvcodieren
allb sub <- mutate(allb sub, alter0 = alter - 18,
            bildung rec = ifelse(bildung == 6 | bildung == 7, NA, bildung - 1),
            geschl rec = ifelse(geschl == 2, 0, 1))
# Subdatensatz erstellen für z-standardisierte Variablen und Variablen auswählen
allb sub z <- allb sub
  select(einkommen, alter, alter0, geschl rec, bildung rec)
# Variablen z-standardisieren
allb_sub_z <- mutate(allb_sub_z, alter_z = scale(alter),</pre>
       alter0_z = scale(alter0),
       bildung_z = scale(bildung_rec),
       einkommen z = scale(einkommen))
```

2.a) Vergleichen Sie die Zahlenwerte, Mean und die Standardabweichung von alter_z und alter_0z und erklären Sie Ihre "Beobachtung".

```
allb sub z \leftarrow select(allb sub z, alter z, alter0 z)
descr(allb sub z)
## Basic descriptive statistics
                               n NA.prc <mark>mean sd</mark> se
                                                        md trimmed
                     label
             type
                                                                               range skew
                                           0 1 0.02 0.03
                                                             -0.01 4.17 (-1.8-2.37) 0.06
                                  0.09
  alter_z numeric alter_z 3468
 alter0_z numeric alter0_z 3468
                                   0.09
                                           0 1 0.02 0.03
                                                             -0.01 4.17 (-1.8-2.37) 0.06
```



Betrachtet man den Datensatz in R (siehe Tabelle links), sieht man, dass die Zahlenwerte von alter_z und alter0_z identisch sind. Ebenso verhält es sich mit dem arithmetischen Mittel und der Standardabweichung (gelb hervorgehoben im Output oben).

Daran kann man erkennen, dass es für die Standardisierung keine Rolle spielt, in welcher Einheit oder auf welcher Skala eine Variable, hier das Alter, gemessen wurde: man erhält das gleiche Ergebnis für das Alter, das mit den "wahren" Alterswerten der Beobachtungen beschrieben wird und für das Alter0, wo von allen "wahren" Alterswerten der Wert 18 abgezogen wurde.

2.b) Führen Sie eine Regression von Einkommen auf Alter_0 und Bildung (Modell 1) und eine Regression von Einkommen auf alter_0z und bildung_z (Modell 2) durch und vergleichen Sie die b-Koeffizienten.

Tabelle 1: Einflüsse auf Einkommen

	Modell 1 (unstandardisiert)		Modell 2 (standardisiert)	
(Intercept)	7.17	***	0	
	(0.28)		(0.02)	
Alter0	0.04	***	0.14	***
	(0.01)		(0.02)	
Bildung	1.20	***	0.29	***
	(0.07)		(0.02)	
R^2	0.08		0.08	
Adj. R²	0.08		0.08	
n	3039		3039	
RMSE	4.74		0.96	

^{***} p<0.001, ** p<0.01, *p<0.05

Während die b-Koeffizienten in Modell 1 eine Interpretation ermöglichen, sind die b*-Koeffizienten in Modell 2 nur untereinander vergleichbar. Die standardisierten Werte kommen sich "näher", da sie immer zwischen -1 und 1 liegen, während die unstandardisierten Werte auch außerhalb liegen, was man am unstandardisierten Wert für Bildung, 1,2 erkennt.

2.c) Wie erklären Sie die Werte b und b* in Modell 2? TIPP: Verwenden Sie bei Modell 2 das ztransformierte Einkommen als abhängige Variable.

Die b*-Werte in Modell 2 sagen aus, dass der Einfluss von Bildung auf Einkommen etwas höher ist als der von Alter auf Einkommen, bzw. dass die Erklärkraft von Bildung bezüglich der Varianz von Einkommen etwas höher ist als die von Alter. Das sieht man daran, dass der Wert 0,29 (Bildung) höher ist als der Wert 0,14 (Alter).

(Keine Werte von b in Modell 2 vorhanden, da sie in Modell 1 sind?)

3. Erstellen Sie ein multivariates Regressionsmodell mit Y=Einkommen. Versuchen Sie dabei den R²-Wert so groß wie nur irgendwie möglich zu bekommen. Jeder schmutzige Trick der Sozialforschung ist erlaubt (und in diesem Fall erwünscht). Fügen Sie die entsprechenden Teile des SPSS-Outputs in Ihre Abgabe ein.

Einzige Einschränkung: Keine Regression von Y auf Y.

```
\#Strategien für großes R² (die normalerweise berücksichtigt/vermieden werden sollten): \#1) Viele UVs
```

^{#2)} UVs, die kausal nachgeordnet sind

^{#3)} UVs mit inhaltlicher Nähe zu AV

^{#4)} Beobachtungen, die linearer Beziehung widersprechen, ausschließen

^{#5)} X-Variablen mit höchster Varianz

^{#6)} AV mit kleinster Residualvarianz

^{#1), 2), 3)} und 5):

[#] Einkommen (Nettoeinkommen v417). Einflüsse: Geschlecht v81, Alter (metrisch gemessen) v84, Schulabschluss v86, Berufsstatus ISCED 2011 v102, Arbeitsstunden pro Woche v118, Geburt in D v377, Erhebungsgebiet westost v7

[#] Nachgeordnete UVs/inhaltliche Nähe (vermutet): Zahl der Bücher im Haushalt v524, Wohnfläche in m² v594, Sprechanlage v843

allbus14 <- read_sav("C:/Users/Anna/Documents/Uni S/Statistische-Modellbildung-II-WS1819-master/data/allbus2014.sav")

```
allbus14 <- rename(allbus14, Einkommen=V417, Weiblich=V81, Alter=V84, Abschluss=V86, IS-
CED2011=V102, hArbeit=V118, GeburtD = V377, Ostdt=V7, Buecher=V524, Flaeche=V594, Sprech=V843)
allbus14 s <- select(allbus14, Einkommen, Weiblich, Alter, Abschluss, ISCED2011, hArbeit, Ge-
burtD, Ostdt, Buecher, Flaeche, Sprech, Einkommen2)
# Beobachtungen mit missings löschen
allbus14 s<-allbus14 s[!(allbus14 s$Einkommen==99997 | allbus14 s$Einkommen==99990),]
allbus14 s<-allbus14 s[!(allbus14 s$Alter==999),]
schluss==99),]
allbus14 s<-allbus14 s[!(allbus14 s$ISCED2011==94 | allbus14 s$ISCED2011==99),]
allbus14_s<-allbus14_s[!(allbus14_s$hArbeit==999.6 | allbus14_s$hArbeit==999.9),]
allbus14_s<-allbus14_s[!(allbus14_s$GeburtD==9),]
allbus14 s<-allbus14 s[!(allbus14 s$0stdt==9),]
allbus14 s<-allbus14 s[!(allbus14 s$Buecher==98 | allbus14 s$Buecher==99),]
allbus14_s<-allbus14_s[!(allbus14_s$Sprech==8),]
allbus14 s<- allbus14 s[complete.cases(allbus14 s), ]</pre>
# Variablen neu codieren (Dummies)
allbus14_s <- mutate (allbus14 s, Weiblich = ifelse(Weiblich == 2, 1, 0),
                 GeburtD = ifelse(GeburtD ==1, 1, 0),
                 Ostdt= ifelse(Ostdt == 2, 1, 0),
                 Sprech = ifelse(Sprech == 1, 1, 0))
#Regressionsmodell
mod4 <- lm(Einkommen ~ Weiblich + Alter + Abschluss + ISCED2011 + hArbeit + GeburtD + Ostdt +
Buecher + Flaeche + Sprech, data = allbus14 s)
texreq(list(mod4))
#R<sup>2</sup> ist 0,17.
\hline
& Model 1 \\
\hline
(Intercept) & $-1516.08^{***}$ \\
          & $(322.54)$
                         \\
          & $-585.75^{***}$ \\
Weiblich
          & $(99.59)$
                            \\
Alter
          & $18.31^{***}$
                            \\
          & $(3.95)$
                            \\
Abschluss & $35.49$
                             //
          & $(59.13)$
                            \\
ISCED2011 & $164.84^{***}$
                            //
                             \\
          & $(38.06)$
          & $31.64^{***}$
hArbeit
                             //
          & $(4.51)$
                             GeburtD
          & $110.84$
                             //
          & $(144.67)$
                             //
Ostdt
          & $-310.40^{**}$
                             //
           & $(100.11)$
                             //
          & $99.58^{**}$
                             //
Buecher
          & $(32.35)$
                             //
          & $3.04^{***}$
                             //
Flaeche
          & $(0.81)$
                             //
           & $79.21$
                             //
Sprech
           & $(90.94)$
                             //
\hline
         & 0.17
Adj. R$^2$ & 0.17
                             \\
Num. obs. & 1529
                             \\
RMSE
          & 1733.24
                            \\
\hline
\mathcal{L}_{1}_{scriptsize}
```

```
#4): Lineare Beziehung zwischen einzelnen Variablen anschauen
plot (Einkommen~Alter, data=allbus14 s)
plot (Einkommen~Abschluss, data=allbus14 s)
plot (Einkommen~hArbeit, data=allbus14_s) #Notiz: alle über 70h löschen, die verdienen wieder
weniger?
plot (Einkommen~Buecher, data=allbus14_s) plot (Einkommen~Flaeche, data=allbus14_s) #Notiz: alle über 300m² löschen, die verdienen weni-
ger?
plot (Einkommen~Weiblich, data=allbus14 s)
plot (Einkommen~GeburtD, data=allbus14 s)
plot (Einkommen~Ostdt, data=allbus14_s)
plot (Einkommen~Sprech, data=allbus14 s)
# Dummies sagen nicht so viel über Linearität
#Ausreißer löschen, auch von AV (#6)
\verb|allbus14_s<-allbus14_s[!(allbus14_s$Einkommen>30000),||\\
allbus14_s<-allbus14_s[!(allbus14_s$hArbeit>70),]
allbus14 s<-allbus14 s[!(allbus14 s$Flaeche>300),]
allbus14 s<- allbus14 s[complete.cases(allbus14 s), ]
#Nochmal Modell probieren
mod5 <- lm(Einkommen ~ Weiblich + Alter + Abschluss + ISCED2011 + hArbeit + GeburtD + Ostdt +
Buecher + Flaeche + Sprech, data = allbus14 s)
texreg(list(mod5))
#R2 bei 0,4!
\hline
 & Model 1 \\
\hline
(Intercept) & $-1326.56^{***} \\
            & $(175.64)$
                                & $-521.44^{***}$
Weiblich
                                \\
            & $(53.13)$
                                & $17.37^{***}$
Alter
                                & $(2.09)$
                                & $33.19$
Abschluss
                                \\
            & $(31.28)$
                                \\
ISCED2011
            & $208.99^{***}$
                                & $(20.14)$
                                \\
            & $28.99^{***}$
hArbeit
                                \\
            & $(2.55)$
                                \\
            & $178.71^{*}$
GeburtD
                                \\
            & $(76.43)$
                                \\
            & $-456.11^{***}$
Ostdt
                                11
            & $(53.15)$
                                \\
Buecher
            & $41.22^{*}$
                                \\
            & $(17.23)$
                                \\
Flaeche
            & $2.30^{***}$
                                \\
            & $(0.51)$
                                \\
            & $120.06^{*}$
                                \\
Sprech
            & $(48.18)$
                                \\
\hline
          & 0.40
Adj. R$^2$ & 0.40
                                //
Num. obs.
            & 1507
                                \\
RMSE
            & 911.20
                                11
\hline
\multicolumn{2}{1}{\scriptsize{$^{***}p<0.001$, $^{**}p<0.01$, $^*p<0.05$}}
#Schauen, ob kategorisiertes Einkommen besser ist: (alles neu weil kategorisiertes und direk-
tes Einkommen sich ausschließen wegen missings)
```

allbus14 <- read sav("C:/Users/Anna/Documents/Uni S/Statistische-Modellbildung-II-WS1819-master/data/allbus2014.sav")

allbus14 <- rename(allbus14, Weiblich=V81, Alter=V84, Abschluss=V86, ISCED2011=V102, hArbeit=V118, GeburtD = V377, Ostdt=V7, Buecher=V524, Flaeche=V594, Sprech=V843, Einkommen2=V418) allbus14_s <- select(allbus14, Weiblich, Alter, Abschluss, ISCED2011, hArbeit, GeburtD, Ostdt, Buecher, Flaeche, Sprech, Einkommen2)

```
# Beobachtungen mit missings löschen
\verb| allbus14_s<-allbus14_s|! (allbus14_s\\$Einkommen2==95 | allbus14_s\\$Einkommen2==99),] |
allbus14_s<-allbus14_s[!(allbus14_s$Alter==999),]
allbus14 s<-allbus14 s[!(allbus14 s$Abschluss==6 | allbus14 s$Abschluss==7 | allbus14 s$Ab-
schluss==99),]
allbus14_s<-allbus14_s[!(allbus14 s$ISCED2011==94 | allbus14 s$ISCED2011==99),]
allbus14 s<-allbus14 s[!(allbus14 s$hArbeit==999.6 | allbus14 s$hArbeit==999.9),]
allbus14_s<-allbus14_s[!(allbus14_s$GeburtD==9),]
allbus14 s<-allbus14 s[!(allbus14 s$0stdt==9),]</pre>
allbus14_s<-allbus14_s[!(allbus14_s$Buecher==98 | allbus14_s$Buecher==99),]
allbus14_s<-allbus14_s[!(allbus14_s$Flaeche==9998 | allbus14_s$Flaeche==9999),] allbus14_s<-allbus14_s[!(allbus14_s$prech==8),]
allbus14 s<- allbus14 s[complete.cases(allbus14 s), ]
# Variablen neu codieren (Dummies)
allbus14_s \leftarrow mutate (allbus14_s, Weiblich = ifelse(Weiblich == 2, 1, 0),
                   GeburtD = \overline{ifelse} (GeburtD ==1, 1, 0),
                   Ostdt= ifelse(Ostdt == 2, 1, 0),
                   Sprech = ifelse(Sprech == 1, 1, 0))
mod6 <- lm(Einkommen2 ~ Weiblich + Alter + Abschluss + ISCED2011 + hArbeit + GeburtD + Ostdt +
Buecher + Flaeche + Sprech, data = allbus14 s)
texreg(list(mod6))
\hline
& Model 1 \\
\hline
(Intercept) & $3.42$
                            \\
           & $(1.89)$
                            \\
           & $-2.38^{***}$ \\
Weiblich
           & $(0.58)$
                            11
           & $0.02$
                            \\
Alter
           & $(0.03)$
                            11
Abschluss & $0.77^{*}$
                            11
           & $(0.38)$
                            11
ISCED2011 & $0.67^{**}$
                            11
           & $(0.23)$
                            11
           & $0.08^{***}$
                            11
hArbeit.
           & $(0.02)$
                            11
           & $2.22^{*}$
                            11
Geburt.D
           & $(0.93)$
                            \\
           & $-3.18^{***}$ \\
Ost.dt.
           & $(0.58)$
                            \\
           & $0.00$
Buecher
                            11
           & $(0.19)$
                            11
           & $0.00$
                            \\
Flaeche
           & $(0.01)$
                            \\
           & $-0.17$
                            11
Sprech
            & $(0.54)$
                            11
\hline
R$^2$ & 0.42
Adj. R$^2$ & 0.39
                            \\
Num. obs. & 209
                            \\
RMSE
           & 3.75
                            \\
\hline
\#R^2 ist 0,42.
```