

Einfache lineare Regression

<p>Modell: $Y_i = \alpha + \beta x_i + E_i$ mit E_i unabhängig $\sim N(0, \sigma^2)$, normalverteilt mit Erwartungswert = 0 und konstanter Varianz</p> <p>Schätzung: Kleinste-Quadrate-Kriterium</p> <ul style="list-style-type: none"> fit $\leftarrow \text{lm}(Y \sim X)$ coef(fit) <p>Genauigkeit und Vertrauensintervalle für Koeffizienten:</p> <ul style="list-style-type: none"> summary(fit) confint(fit, level=0.95) zusätzlich parm = 1 oder 2 oder : coef(lm1)[2] + qt(0.995, n - 2) * summary(lm1)\$coef[2,2] * c(-1,1) <p>Genauigkeit und Vertrauensintervalle für $E(y_i)$ (deckt nur Unsicherheit der Schätzung ab):</p> <ul style="list-style-type: none"> predict(fit, se.fit=TRUE, ...) predict(fit, newdata=data.frame(x=c(1, 1/10)), interval="confidence", level=0.95) -> Bsp. Windmühlen <p>Prognoseintervall für $E(y_i)$: (+ σ^2, Varianz E_i -> somit grösseres Intervall als Vertrauensintervall, deckt Variabilität aus Fehler E_i ab) In welchem Bereich liegt eine zukünftige Beobachtung zum Niveau 95%?</p> <ul style="list-style-type: none"> predict(fit, newdata=xyz, interval="prediction", level=0.95) <p># Rücktransformation $\exp(b_0 + b_1 \cdot \text{IN} + b_2 \cdot \text{IC} + b_3 \cdot \text{P} + E_i) \rightarrow \exp(b_0) \cdot N^{b_1} \cdot C^{b_2} \cdot \exp(b_3 \cdot \text{P}) + \exp(E_i)$ (kann Bias enthalten) $\exp(h)$ oder mit Korrektur: $\exp(h + \text{summary(lm1)}\\$sigma^2/2)$</p> <p>Statistische Aussagen sind nur vertrauenswürdig, wenn Modellannahmen erfüllt sind</p> <p>Hauptziele:</p> <ol style="list-style-type: none"> Eine Gerade in eine Punktwolke legen, um die Beziehung zwischen einer Zielgrösse Y und einer erklärenden Variablen x zu beschreiben. Bestimmen, ob die Daten mit einem Modell mit vorgegebenen Koeffizienten verträglich sind (statistischer Hypothesentest). Angaben, wie genau der Achsenabschnitt und die Steigung aus den Daten bestimmt werden können (Vertrauensintervall) 	<p>Allgemein</p> <p>Residuen (R_i): «Näherungswerte», Differenzen zwischen den Beobachtungen und den angepassten Werten</p> <p>Hut: Angepasste Werte (geschätzt)</p> <p>α = Achsenabschnitt</p> <p>β = Steigung</p> <p>E_i = Zufallsfehler</p> <p>Prüfen der Modelleignung</p> <p>-> Linearität vorausgesetzt</p> <p>Annahmen Fehler E_i:</p> <ol style="list-style-type: none"> Erwartungswert = 0 Alle Fehler haben die gleiche Varianz σ^2 Sind normalverteilt <p>Sind unabhängig</p> <p>R-Code</p> <p>Daten einlesen:</p> <pre>MPI <- read.table(paste("Statistisches Modellieren/Arbeitsblätter/MPIZH.dat", sep = ""), header = T)</pre> <p>Gerade (Fitter) konstruieren</p> <pre>uhr.lm <- lm(P ~ A, data=Uhr) coef(uhr.lm) summary(uhr.lm)</pre> <p>Schätzen der Geraden</p> <pre>paste("y_i =", round(lm1\$coef[1],3), "+", round(lm1\$coef[2],3), "x_i")</pre> <p>1. J. ältere Uhr: $1 \cdot x_i \rightarrow$ Preiserhöhung</p> <p>In Plot einfügen</p> <pre>plot(Uhren\$Alter, Uhren\$Preis, xlab = "Alter", ylab = "Preis") abline(uhr.lm, col="red")</pre>
<p>Kleinste Quadrate-Methode R^2: (Bestimmung Güte)</p> <p>$\frac{\text{Quadratsumme Fit}}{\text{Quadratsumme Y}}$</p> <p>-> liegt zwischen 0 - 1 (je näher an 1, desto besser der lineare Zusammenhang)</p> <p>Misst den Anteil der durch die Regression erklärten Streuung der Y-Werte. - Identisch zur Korrelations² zwischen Zielvariablen y_i und \hat{y}_i</p> <p>-> kein Mass für Eignung des Regressionsmodells</p> <p>-> misst die Stärke des linearen Zusammenhangs</p> <p>-> Modellannahme muss erfüllt sein (β_0 muss vorhanden sein)</p> <p>Technik / Naturwissenschaften: Werte > 0.9 üblich</p> <p>Geistes-/ Sozialwissenschaften: Werte um 0.6 i.O.</p>	<p>R-Code (allgemein)</p> <p>Mit Matrix x Werte generieren:</p> <pre>x.sim <- c(0,3,2,4,8) E.sim <- c(matrix(rnorm(10*100, mean=0, sd=2), ncol=100)) y.sim <- 4 + 2*x.sim + E.sim Koeff <- matrix(0, ncol=2, nrow=100)</pre> <p>for (i in 1:100) { Koeff[i,] <- coef(lm(y.sim[,i] ~ x.sim)) hist(Koeff[, 1]) # $1 = \alpha, 2 = \beta$</p> <p>Beobachtungen entfernen</p> <pre>Forbes <- lm(y ~ x, Forbes[-12,]) ODER: Forbes <- lm(y ~ x, Forbes, subset = -12)</pre> <p>R^2</p> <pre>lm(...): 2. unterste Zeile «Multiple R-squared» oder summary(lm1)\$r.squared</pre>

Diagnose Instrumente:

<p>-> Sicherstellung, dass in den Daten keine für die Theorie gefährlichen Abweichungen zu den Voraussetzungen</p> <p>-> Residuen Realisierungen von Zufallsgrösser -> Voraussetzungen können nie exakt erfüllt sein</p> <p>Tukey-Anscombe-Diagramm (Residuen gegen angepasste Werte)</p> <p>-> Rauschen um die Gerade</p> <p>Linearitätsannahme $E = 0$ eingehalten? (konstant innerhalb der stochastischen Fluktuation?)</p> <p>Robuste Methode (ohne Einfluss Ausreisser)</p> <p>Residuen sollten gleich streuen (konstante Varianz)</p> <p>2 Sägezähne: 2 Gruppen Formen beschrieben</p>	<p>Teststatistik</p> <p>R-Code</p> <pre>tc <- coef(lm1)[2] / summary(lm1)\$coef[2,2]</pre> $T = \frac{\hat{\beta} - \beta^0}{\text{se}(\hat{\beta})}$ <p>Nullhypothese:</p> <p>P-Wert: < 0.05 -> Verwerfen / ablehnen, da signifikant</p> <pre>qt(0.975, df=11) # Grenze auf 5 % Niveau mit 11 Frei.graden</pre> <p><u>Estimate – Nullhypothese</u></p> <p><u>Std. Error</u></p> <p>Grün < orange -> Nullhypothese kann nicht abgelehnt werden</p>
<p>Glätter zeigt eine Abweichung von horizontalen Geraden: (Mitte):</p> <p>Beurteilung Abweichung anhand Bootstrap-Simulation (mehr Sicherheit, da mehr Simulationen vorhanden, Zufall ausschliessen):</p> <ol style="list-style-type: none"> Zufällige Beobachtungen erzeugen (Zufallszahlen generieren) Regressionsrechnung durchführen, Glätter in das Tukey-Anscombe-Diagramm einzeichnen Schritte repetieren <p>(Kann auch auf Streuungsdiagramm und normalen QQ-Plot angewandt werden)</p> <p>Resultat:</p> <p>In Beispiel wurden 19 Kurven erzeugt, da in Anlehnung auf das Testen auf dem 5 % Niveau (1/20 -> 19 + 1 = 20)</p> <p>Nicht kritisch (siehe graue Kurven), Beobachtung oben in Ecke extrem -> Glätter muss innerhalb der stochastischen Fluktuation liegen (Graue Linien / Punkte)</p>	<p>R-Code</p> <pre>SD.lm <- lm(SZ ~ Liefvolumen, data = SD) par(mfrow=c(2,4)) plot(SD.lm) # Ohne which: alle 4 Diag. source("Statistisches Modellieren/Arbeitsblätter/RFn_Plot-ImSim.R") # Laden Funktion Plot.ImSim plot.ImSim(x.fit, SEED=567) # Bootstrap-Simulation 3 Diagramme plot.ImSim(x.fit, SEED=567, rob=TRUE) # Verdacht auf Langschw. -> SD robust schätzen: rob = TRUE</pre> <p>plot(Fit, which = 4)</p> <p>plot(lmFit, which=5, add.smooth=FALSE)</p> <p>1: Erwartungswert / 2: Normalverteilung der Fehler / 3: Varianz / 4: Distanz gegen Beobachtungsnummer / 5: Residuen gegen Hebel</p>
<p>Streuungs-Diagramm (scale-location-plot)</p> <p>-> Varianz konstant? Sägezähne vorhanden?</p> <p>Varianz ist anfälliger auf Ausreisser als Mittelwert</p> <p>Glätter auf $\sqrt{ R_i }$ anwenden</p>	<p>QQ-Plot / Normalverteilungs-Diagramm</p> <p>Sind Fehler normalverteilt? Langschwänzigkeit / Ausreisser?</p> <p>Falls Ja: Punkte streuen um eine Gerade</p> <p>-> Für die Y_i sinnlos, da verschiedene Erwartungswerte</p> <p>R-Code:</p> <pre>qqnorm() / qqline() -> Linie in Punkten</pre>
<p>Histogramm -> $N(0, \sigma^2)$</p> <p>Sind Fehler normalverteilt?</p> <p>Beurteilung, von Auge schwierig:</p> <ul style="list-style-type: none"> Bei wenig Daten (< einige 100) Form schlecht erkennbar Sensitiv auf Balkeneinteilung Nicht-lineare Strukturen werden miteinander verglichen <p>Vergleichbar mit Normalverteilungsdichte in R: hist(..., freq=F)</p> <p>Skalierte Residuen – Verteilung der Zufallsfehler überprüfen, wenn Residuen bereits benutzt</p> <p>Fehler normalverteilt: Residuen von einer R^2-Schätzung auch normalverteilt</p> <p>! Nicht gleiche Varianz: Diese hängt von erklärenden Variabel x_i ab</p> <p>Skaliert: $\frac{R_i}{\sqrt{1 - \left(\frac{1}{n} + \frac{(x_i - \bar{x})^2}{SS_X}\right)}} \sim N(0, \sigma^2)$</p> <p>Standardisiert: $\frac{R_i}{\hat{\sigma} \sqrt{1 - \left(\frac{1}{n} + \frac{(x_i - \bar{x})^2}{SS_X}\right)}}$</p>	<p>Vorgehen:</p> <ul style="list-style-type: none"> Grundsätzlich skalierte Residuen verwenden bei Varianz = 1 -> standardisierte Residuen anwenden Tukey-Anscombe-Diagramm: Rohresiduen verwenden (Unkorrelation) <p>Je weiter Beobachtung vom Schwerpunkt entfernt, desto kleiner die Varianz und desto näher geht die Regressionsgerade am Beobachtungspunkt vorbei</p>

Behandlung von Unzulänglichkeiten

<p>First Aid Transformationen Häufig: Abhängigkeit der Streuung von \hat{y} (nimmt im Diagramm zu, Normal-Plot zeigt eine rechtsschiefe Verteilung) -> Logarithmieren der Zielgröße hilft</p> <p><u>Logarithmus-Transformation</u> Für Konzentrationen und Beträge / Mengen summary(data) -> Daten anschauen, Differenz Min/ Max < als Faktor 2 -> Log lohnt sich nicht</p> <p><u>Wurzeltransformation</u> Für Zähldaten</p> <p><u>Arcus-Sinus-Wurzel-Transformation / Logit-Transformation</u> $\tilde{y} = \arcsin(\sqrt{y})$ resp. $\tilde{y} = \log\left(\frac{y + 0.005}{1.01 - y}\right)$ Für Anteile / Prozentzahlen</p> <p>Transformation der Zielvariablen ändert Form der Verteilung der Fehler -> Potenzgesetz für die ursprünglichen Größen, somit ist der Fehler proportional Falls $\beta = 1$ ist die Zielvariable proportional zu x bis auf einen multiplikativen zufälligen Fehler.</p> <p>Rücktransformation ergibt nicht immer dasselbe Ergebnis wie zuvor</p>	<p>Erwartungswert ist nicht konstant 0 Systematische Abweichungen im Erwartungswert können oft durch: - Transformation der erklärenden Variablen x - Oder durch Hinzufügen eines zusätzlichen Terms x^2 (quadr. Regression) zum Verschwinden gebracht werden.</p>
<p><u>Ausreisser - Fehler sind nicht normalverteilt</u> Auf Richtigkeit der Daten prüfen! -> Transformationen können helfen</p>	<p>Ausreisser - Fehler sind nicht normalverteilt Auf Richtigkeit der Daten prüfen! -> Transformationen können helfen</p>
<p><u>Ausreisser entfernen:</u> D.synt.lm2 <- lm(y ~ x1 + x2, data=D.synt, subset=-c(7,17,27))</p>	<p>Ausreisser entfernen: D.synt.lm2 <- lm(y ~ x1 + x2, data=D.synt, subset=-c(7,17,27))</p>
<p><u>Langschwänzigkeit - Fehler sind nicht normalverteilt</u> Extremste Beobachtungen weglassen, bis Langschwänzigkeit verschwindet (! optimistisch)</p> <p>Kleinste-Quadrate-Methode nicht optimal -> nur robuste Methoden geeignet (geringere Gewichtung von Ausreißern)</p>	<p>Langschwänzigkeit - Fehler sind nicht normalverteilt Extremste Beobachtungen weglassen, bis Langschwänzigkeit verschwindet (! optimistisch)</p>
<p><u>Unabhängigkeit zufällige Fehler</u> Bei Beobachtungen mit zeitlicher Reihenfolge, können Autokorrelationen vorhanden sein. -> Residuen in dieser Reihenfolge auftragen (keine Strukturen ersichtlich?) -> oder R_t gegen R_{t-1} in Streudiagramm auftragen - Punkte sollten frei streuen und keine Korrelation zeigen Wenn Autokorrelationen vorliegen: - P-Werte der üblichen Tests häufig grob falsch Vertrauensintervalle üblicherweise zu kurz</p>	<p>Unabhängigkeit zufällige Fehler Bei Beobachtungen mit zeitlicher Reihenfolge, können Autokorrelationen vorhanden sein. -> Residuen in dieser Reihenfolge auftragen (keine Strukturen ersichtlich?) -> oder R_t gegen R_{t-1} in Streudiagramm auftragen - Punkte sollten frei streuen und keine Korrelation zeigen Wenn Autokorrelationen vorliegen: - P-Werte der üblichen Tests häufig grob falsch Vertrauensintervalle üblicherweise zu kurz</p>

Multiple lineare Regression

<p>$Y_i = \beta_0 + (\beta_1 x_i)^{(1)} + (\beta_2 x_i)^{(2)} + \dots + (\beta_m x_i)^{(m)} + E_i$</p> <p>Beispiel: $\log(ersch) = \beta_0 + \beta_1 \log(dist) + \beta_2 \log(lad) + E_i$</p> <p>Zufällige Fehler E_i: Unabhängig Koeffizienten β_j (unb. Parameter): normalverteilt Schätzung der Koeffizienten β_j: Kl. Quadrate Meth.</p> <p>Erklärende Variablen xi: - Keinen bestimmten Datentyp - Keine bestimmte Verteilung (keine Zufallsvar.) - Keine Voraussetzung über ihre Abhängigkeit untereinander</p>	<p>Bestimmtheitsmass (Multiple R-Squared) - Analog einf. lin. Regression - Misst Anteil der durch Regressfunktion erklärten Streuung an der Streuung der Y-Werte - Entspricht dem Quadrat der Korrelation zwischen den beobachteten und den angepassten Werten (linearer Zusammenhang)</p> $R^2 = \frac{\sum_{i=1}^n (\hat{y}_i - \bar{y})^2}{\sum_{i=1}^n (y_i - \bar{y})^2}$ <p>Nach kleinsten Quadrate Methode geschätzte Koeffizienten: - Minimieren Quadratsumme der Residuen - Maximieren die Korrelation zwischen den angepassten Werten und den Beobachtungen der Zielgröße (maximale Wert = multiple Korrelation)</p>
<p>R-Code: Vertrauensintervall: confint(SP.lm2, level=0.95) confint(MPIZH.lm1, parm=2, level=0.95) von Hand: $2.97 + c(-1, 1) * 0.36 * qt(0.99, 43)$ aus Summary(): Estimate / Std. Error / F-Statistics (DF)</p> <p>Level: $0.98 \rightarrow 0.02/2 + 0.98 = 0.99$</p>	<p>Prognoseintervall: x.MPI <- data.frame(HZ=4.5, KPI=100.5) -> Reihenfolge egal predict(MPI.lm, newdata=x.MPI, interval="prediction", level=0.95)</p>

!! Multiple Regression ist nicht gleich der Summe der einfachen Regression

```

> SprngS2[1:4,]
  Stelle ladung dist ersch lLadung lDist lErsch
1 St1 2.18 188 0.32 0.7793249 5.236442 -1.1394343
2 St2 3.33 183 0.53 1.2029723 5.209486 -0.6348783
3 St1 3.33 177 0.50 1.2029723 5.176150 -0.6931472
4 St1 3.33 53 7.61 1.2029723 3.970292 1.9035990

> Spr2.lm3 <- lm(lErsch ~ lDist + lLadung + Stelle, data=SprngS2)
> summary(Spr2.lm3)
...
Coefficients:
(Intercept)  5.78051  0.64967  8.898  3.25e-11
lDist       -1.33779  0.14073  -9.506  4.97e-12
lLadung      0.69179  0.29666  2.332  0.0246 *
StelleSt2   -0.37832  0.17257  2.192  0.0340 *
StelleSt3   0.04996*  0.14657  0.341  0.7349
StelleSt4   0.25511  0.17216  1.482  0.1459
          
```

Handwritten notes:
- Achsenabschnitt von Stelle 1
- Verschiebung des Achsenabschnittes
- Differenz
- Stellen: Unterschiede zwischen den Achsenabschnitten

$\log(ersch) = \beta_0 + \beta_1 \log(dist) + \beta_2 \log(ladung) + \beta_{3,1} St1 + \beta_{3,2} St2 + \beta_{3,3} St3 + \beta_{3,4} St4 + E_i$

$\hat{\sigma} \rightarrow$ sollte möglichst klein sein

Vielfalt der Modellierungsmöglichkeiten

Polynomiale Regression (Spezialfall von mult. lin. Regression)
Bsp.: Polynom 2. Ordnung wird zur Beschreibung eines Zusammenhangs verwendet: $Y_i = \beta_0 + \beta_1 x_i + \beta_2 x_i^2 + E_i$
-> Erklärende Variablen $x_i^{(1)}$ & $x_i^{(2)}$ einfügen
-> Lin. Reg. modell: $Y_i = \beta_0 + \beta_1 x_i^{(1)} + \beta_2 x_i^{(2)} + E_i$

«linear» im Begriff der multiplen linearen Regression bezieht sich darauf, dass die Koeffizienten linear in der Formel vorkommen!

Nichtlineare Funktionen und lineare Regression
Oft müssen Zielvariablen und / oder die erklärenden Variablen transformiert werden. -> «linearisierbaren Gleichungen» (Nach First Aid Transformationen)

$$y = \frac{1}{a + b * \exp(-x)} \leftrightarrow \frac{1}{y} = a + b * \exp(-x)$$

$$y = \frac{a * x}{b + x} \leftrightarrow \frac{1}{y} = \frac{1}{a} + \frac{b}{a} * \frac{1}{x}$$

$$y = a * x^b \leftrightarrow \ln(y) = \ln(a) + b * \ln(x)$$

$$y = a * \exp(b * g(x)) \leftrightarrow \ln(y) = \ln(a) + b * g(x)$$

Vergleich von Regressionsmodellen mit F-Test
Einfluss von Koeffizienten: «Kein Einfluss»: alle Koeffizienten (Stellen) = 0 (Nullhypothese $b_{jq} = 0$)

$$T = \frac{\frac{SS_E^* - SS_E}{q}}{\frac{SS_E}{n - p}}$$

SS_E^* = Quadratsumme des Fehlers im red. Modell
 SS_E = Quadratsumme des Fehlers aus dem alternativen Modell
q, n - p = Freiheitsgrade

R-Code:
anova(Spr2.lm3, Spr2.lm2) 2 Modelle gegeneinander
drop1(Spr.lm3, test="F") -> ähnlich Summary, jedoch auch für Faktorvariablen, Test Variable wird weggelassen im gleichen Regressionsmodell (eignet sich nicht um signifikante Einflüsse auf Zielvariablen festzustellen, sobald Faktor Variablen vorkommen)

F-Test bei einfacher Regression:
Entspricht derselben Größe wie t-Test zum Steigungskoeffizienten

Binär erklärende Variablen
Eine erklärende Variable kann binär sein, also auf die Werte 0 und 1 beschränkt sein.
-> Regressionsmodell beschreibt 2 unabhängige Stichproben (ungepaarter Zweiprobe-t-Test)

$$Y_i = \beta_0 + E_i \quad \text{für } x = 0$$

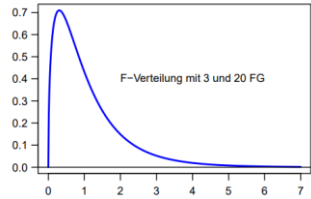
$$Y_i = \beta_0 + \beta_1 + E_i \quad \text{für } x = 1$$

Beispiel: Allfälliger Unterschied der Lage / sind zwei Geraden gleich?

$$Y_i = \beta_0 + \beta_1 x_i^{(1)} + \beta_2 x_i^{(2)} + E_i$$

$x_i^{(1)} = \log(\text{Distanz})$
 $x_i^{(2)} = \text{bin. Variable (0 oder 1) für Messstelle}$

Modell umformulieren:
 $Y_i = \alpha + \beta x_i + \Delta \alpha g_i + \Delta \beta g_i + E_i$
-> $g_i = 0$ falls Gruppe A / $g_i = 1$ falls Gruppe B
Fallweise aufgeschrieben:
 $g_i = 0 : Y_i = \alpha + \beta x_i + E_i$
 $g_i = 1 : Y_i = (\alpha + \Delta \alpha) + (\beta + \Delta \beta) x_i + E_i$
-> Steigung stimmt überein, wenn $\Delta \beta = 0$
-> Geraden stimmen überein, wenn $\Delta \alpha = 0$ und $\Delta \beta = 0$



Verteilung:
F-Verteilung (rechtsschief, unimodal)
Nur für positive Werte definiert, je schiefer verteilt je kleiner ist q

Globaler F-Test: In letzter Zeile R-Outputs (summary)
->immer die gleiche Antwort wie t-test

```

> anova(Spr2.lm3, Spr2.lm2)
Analysis of Variance Table

Model 1: lErsch ~ lDist + lLadung + Stelle
Model 2: lErsch ~ lDist + lLadung

Res.Df  RSS    Df Sum of Sq    F    Pr(>F)
1      42 4.7964    -0.78205  2.2827  0.09287
2      45 5.5785    -0.78205  2.2827  0.09287

Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

'RSS' = Summe der quadrierten Fehler (SSE oder SSR)
'Sum of Sq' = SST - SSR

> drop1(Spr.lm3, test="F")
Single term deletions

Model:
lErsch ~ lDist + lLadung + Stelle
<none>             4.7964   -98.560
lDist              10.321  15.1170  -45.458  90.3722  4.973e-12 ***
lLadung             1      0.621   5.4175  -94.716  5.4379   0.02457 *
Stelle              3      0.782   5.5785  -97.310  2.2827   0.09287 .

Beachten Sie, dass F_{1, n-p} * (t_{n-p})^2 gleiche P-Werte
    
```

Anzahl Freiheitsgrade
*Hand: (4.7964 - 5.5785) * 3 = (4.7964 / 4.7)*
= 2.2857
Nullhypothese nicht verworfen
Wert aus T-Test

Modell und Schätzung in Matrix-Schreibweise

$$Y_i = \beta_0 + (\beta_1 x_i)^{(1)} + (\beta_2 x_i)^{(2)} + \dots + (\beta_m x_i)^{(m)} + E_i$$

Als Vektor schreiben:

mit

$$Y = \begin{bmatrix} Y_1 \\ Y_2 \\ \vdots \\ Y_n \end{bmatrix}, E = \begin{bmatrix} E_1 \\ E_2 \\ \vdots \\ E_n \end{bmatrix}, \beta = \begin{bmatrix} \beta_0 \\ \beta_1 \\ \vdots \\ \beta_m \end{bmatrix} \text{ und } X = \begin{bmatrix} 1 & x_1^{(1)} & x_1^{(2)} & \dots & x_1^{(m)} \\ 1 & x_2^{(1)} & x_2^{(2)} & \dots & x_2^{(m)} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & x_n^{(1)} & x_n^{(2)} & \dots & x_n^{(m)} \end{bmatrix}$$

$$Y = X\beta + E$$

Residuen: $R = Y - X\beta$
 Summe der Quadrate: $R^T R$
 Erwartungswert: $X\beta$
 Varianz: $\sigma^2 I$
 KQ: $\hat{\beta} = (X^T X)^{-1} X^T Y$
 → Voraussetzung: $(X^T X)$ invertierbar oder nicht singulär 3

→ Das Prinzip der Maximalen Likelihood besteht darin, die Parameter so zu wählen, dass die Wahrscheinlichkeit (oder Dichte), die beobachtete Zielvariablenwerte zu erhalten, maximiert wird.
 → ML & KQ führen zu gleichem Schätzer für β

Zufallsvektoren – Relevant? Skript S. 64 - 70

Erwartungswert

$$\mathbb{E}(Y) := [\mathbb{E}(Y^{(1)}), \mathbb{E}(Y^{(2)}), \dots, \mathbb{E}(Y^{(p)})]^T$$

Varianz-Kovarianz-Matrix / Kovarianz-Matrix

$$\text{cov}(Y) = \text{var}(Y) := \begin{bmatrix} \text{var}(Y^{(1)}) & \text{cov}(Y^{(1)}, Y^{(2)}) & \dots & \text{cov}(Y^{(1)}, Y^{(p)}) \\ \text{cov}(Y^{(2)}, Y^{(1)}) & \text{var}(Y^{(2)}) & \dots & \text{cov}(Y^{(2)}, Y^{(p)}) \\ \vdots & \vdots & \ddots & \vdots \\ \text{cov}(Y^{(p)}, Y^{(1)}) & \text{cov}(Y^{(p)}, Y^{(2)}) & \dots & \text{var}(Y^{(p)}) \end{bmatrix}$$

Zusammengefasst:

$$\text{var}(a + BY) = B \cdot \text{var}(Y) \cdot B^T = B \Sigma B^T$$

 Das gilt auch für eindimensionale Z; die Matrix B besteht dann aus einer einzigen Zeile, die wir als transponierten Vektor b^T aufschreiben. So wird

$$\text{var}(a + b^T Y) = b^T \text{var}(Y) b = b^T \Sigma b$$

 ein Resultat, das oft nützlich ist.

Für die Zufallsvariable Y gilt
 • $\text{var}(Y) = \mathbb{E}((Y - \mu)^2)$
 • $\text{cov}(Y^{(i)}, Y^{(k)}) = \mathbb{E}((Y^{(i)} - \mu^{(i)})(Y^{(k)} - \mu^{(k)}))$
 Für den Zufallsvektor Y gilt analog
 • $\text{var}(Y) = \mathbb{E}((Y - \mu)(Y - \mu)^T)$
 Für die Zufallsvariable Y gilt bei linearen Transformationen
 • $\mathbb{E}(a + bY) = a + b\mathbb{E}(Y)$
 • $\text{var}(a + bY) = b^2 \text{var}(Y)$
 Für den Zufallsvektor Y gilt analog
 • $\mathbb{E}(a + BY) = a + B\mathbb{E}(Y)$
 • $\text{var}(a + BY) = B \cdot \text{var}(Y) \cdot B^T$ *kurz* $= B \Sigma Y B^T$
 • und falls B nur ein Zeilenvektor ist: $\text{var}(a + b^T Y) = b^T \text{var}(Y) b = b^T \Sigma b$

Sensitivität und Robustheit

Theoretische Verteilung der Residuen

Da die Nebendiagonalelemente in $\text{var}(R) = (I - H) \sigma^2$ im Allgemeinen nicht null sind, sind die Residuen im Gegensatz zu den Fehlern E_i korreliert,

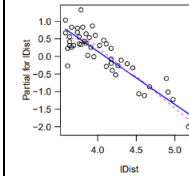
$$\text{cov}(R_i, R_k) = -\sigma^2 H_{ik}$$

Diese Korrelation zwischen den Residuen beeinflusst jedoch die Strukturen in der grafischen Residuen-Analyse kaum. Auch nicht jene, wo wir den zeitlichen oder räumlichen Korrelationsstrukturen (stochastische Abhängigkeit, siehe Abschnitt 3.4) nachgingen.

Residuen-Analyse bei der multiplen Regression

1. Grafiken analysieren. Allenfalls zeitliche Abhängigkeiten klären.
2. Fazit – Welches ist die relevanteste Unstimmigkeit und wie kann sie bereinigt werden?

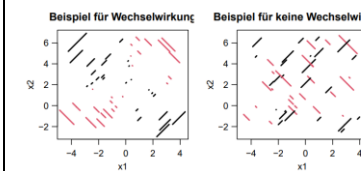
Überprüfung Abh. Streuung:



Beziehung zwischen Residuen und erkl. Variablen

Überprüfung Wechselwirkung:

Voraussetzung: Effekte 2 erkl. Variablen addieren sich



Gewichtete lineare Regression

Annahme, dass Varianzen der Residuen nicht konstant

$\sigma_i^2 = \text{var}(E_i) \rightarrow$ Lösungsansatz: $\sigma_i^2 = \frac{\sigma^2}{w_i}$
 → Mittelwert streut weniger je grösser die Stichprobe -> $w_i =$ Anzahl Stichproben
 → Generell: Für Beurteilung von Verteilung und Streuung der Fehler skalierte Residuen verwenden!!

Weitere Überprüfungsmethoden Residuenanalyse

```

Schauen, ob Streuung um Winkelhalbierende:
par(mfrow=c(3,3))
plot(fitted(asp.lm), log(asp$RUT));
abline(a=0, b=1, lty=2)
    
```

```

fittes vs resid / Erkl. Variablen vs resid
plot(fitted(asp.lm), resid(asp.lm), col=asp$RUN+2)
plot(Ox$TS, resid(Ox.lm1)); abline(h=0, lty=3)
plot(Ox$Day, resid(Ox.lm1), type="h") Zeit Analyse
    
```

```

Residuen vs erkl. Var. (um horizontale Null-Linie):
plot(log(asp[, 'VISC']), resid(asp.lm), col=asp$RUN+2)
    
```

R-Code:

In multiplen Regression -> Plot pro erkl. Variable erstellen:
`scatter.smooth(dat$ILadung, resid(fit), lpar = list(col = 2)); abline(h = 0))`

Partierer-Residuen-Plot:

```

termplot(Spr2.lm2, partial.resid=TRUE,
smooth=panel.smooth, ylim="free",
col.res="black", col.term="blue",
col.smth="magenta")
    
```

Wechselwirkungs-Diagramm

```

Package sfsmisc
p.res.2(x=SprengS2$lDist, y=SprengS2$lLadung,
z=resid(Spr2.lm2), xlab="log(Distanz)",
ylab="log(Ladung)", scol=c(1,1), size=1.2, slwd=2,
main="")
    
```

R-Code:

Einzeichnen in Plot (rote Line, skaliert):

```

G <- data.frame(SPEED=GASKETS$SPEED, RESID=resid(G.lm1))
scatter.smooth(G$SPEED, sqrt(abs(G$RESID)), span=1)
    
```

Gewichtung ermitteln:

```

(G.varRes <- aggregate(RESID ~ SPEED, data=G, FUN=var))
# Gegenüberstellung der beiden Variablen
    
```

$G.\text{varRes}[2]/c(100,150,200)$ # (100,150,200) = Variable Speed, Resultat noch nicht konstant, somit weiter ausprobieren:

```

G.varRes[2]/(c(100,150,200)^2) # +/- konstant, jedoch sehr klein
G.varRes[2]/(c(100,150,200)/100)^2 # auch +/- konstant, wenn
pro 100 Einheiten Produktionsgeschwindigkeit?
-> verwende als Gewicht 1/(SPEED/100)^2
    
```

Fitter gewichten:

```

GASKETS$w <- 1/((GASKETS$SPEED/100)^2)
> G.lm2 <- lm(DEFECTS~SPEED, weights=w, data=GASKETS)
    
```

Residual standard error: 4.219 -> bei Gewicht = 1

Einflussreiche Beobachtungen
Cook's Distance (Ausreisser-Analyse)

$$\tilde{R}_i = \frac{R_i}{\hat{\sigma} \sqrt{1 - H_{ii}}}$$

Gibt es eine zu einflussreiche Beobachtung?

Wertebereich für H_{ii} $0 \leq H_{ii} \leq 1$ -> Funktioniert wie Hebelarm, misst wie untypisch die Beobachtung auf die erklärenden Variablen ist

- Beobachtungen mit einem Hebelarm $> 2 \frac{p}{n}$ haben zu grosse Hebelwirkung
- Beobachtungen mit einem Hebelarm unter 0.2 sind unbedenklich
- Beobachtungen mit einem Hebelarm über 0.5 sollten vermieden werden

«böartige Beobachtung»
 «An Grenze zu einflussreicher Beobachtung»
 «Hebelpunkte»

Diagramm
 Standardisierte Residuen gegen H_{ii}

11: grosser Ausreisser in y-Richtung
 34: zu grosser Hebelarm ($h_i > 0.2$) (siehe Leverage)
 80: Grosser Hebel, jedoch kleines standardisiertes Residuum daher ungefährlich

Robuste Anpassungsmethoden
 Untersuchung Robustheit:

- Einflussfunktion
- Bruchpunkt

➔ Robuster Schätzer hat beschränkte Sensitivität und Bruchpunkt, der möglichst nahe beim maximal möglichen Wert von 1/2 liegt

Regressions-M-Schätzer
 Knick in der Funktion legt fest, ab wo extreme Beobachtungen an Einfluss verlieren

Schätzer für lineares Regressionsmodell:

- Beobachtungen mit grossen Residuen müssen ignoriert werden

R-Code:
Robuste M-Anpassung:
 Mfit <- rlm(y ~ x, method = "M", data=AQ)

Regressions-MM-Schätzer: (Package: robustbase)
 lmrob(y ~ x1 + x2, data=D.syn, setting="KS2014")
 besonders effiziente Vertrauensintervalle: setting = "KS2014"
 par(mfrow=c(2,3))
 plot(lmrob)

Robuster F-Test
 Beruht auf robusten Devianz
 Summen der Quadrate im F-Test durch Devianzen ersetzen:

lm1 <- lmrob(FoHF ~ ., data=FoHF2, setting="KS2014")
 anova(lm1, FoHF ~ LSE + GM + EM + SS, test="Deviance")

Variablenselektion und Modellbildung

Welche Variablen sollen wie ins Regressionsmodell?
 Kriterien basierte Variablenselektion ist vorzuziehen.

Warum nicht alle Variablen?

- Einfachheit
- Reduktion Schätzvariabilität (unnötige Variablen verschlechtern die Genauigkeit der Schätzung)
- Bessere Interpretierbarkeit (Multikollinearität führt nicht zu eindeutigen Lösungen)
- Vorhersage – weniger erkl. Variablen, weniger Aufwand beim Vor- und Aufbereiten

Modellwahlkriterien

- Modellgenauigkeit / Modellkomplexität
- ➔ Sind konträr – Zusätzliche erkl. Variablen führen stets zu höherer Modellgenauigkeit

Lineare Regression mit KQ-Anpassung:

- Modellgenauigkeit mit der Summe der Residuenquadrate
- Modellkomplexität mit der Anzahl Koeffizienten

Bestimmtheitsmass R^2 korrigieren (adjusted R^2)
 R^2 wird grösser je mehr Variablen hinzugefügt werden -> untauglich (Komplexität nicht berücksichtigt)
 Korrigiertes Bestimmtheitsmass:

$$R^2_{adj} := 1 - \frac{SS_E / (n - p)}{SS_Y / (n - 1)}$$

Mallow's Cp-Statistik:
 Minimiert in gewisser Weise den Vorhersagefehler

$$C_p := \frac{SS_E}{\hat{\sigma}_p^2} + 2p - n = (n - p) \left(\frac{\hat{\sigma}_p^2}{\hat{\sigma}_p^2} - 1 \right) + p$$

Informations-Kriterium von Akaike (AIC)
 Gutes verallgem. Kriterium (auch Zeitreihen)

$$AIC = -2 (\text{maximierte Log-Likelihood}) + 2 \cdot (\text{Anzahl geschätzter Parameter})$$

$$= n \log \left(\frac{1}{n} SS_E \right) + 2p + \text{Konstante}$$

wobei SS_E : die Summe der quadrierten Residuen ist.
 (Achtung: p ist hier die Anzahl der geschätzten Parametern inklusive $\hat{\sigma}$.)

Variablenselektion mit P-Werten
 Ist ein bestimmter Term im Modell nötig / nützlich / überflüssig?

Problem des multiplen Testens:

- ➔ Bei Faktorvariablen den t-Test durch den F-Test ersetzen (prüft, ob der ganze Block von der entspr. Variabel weggelassen werden kann)

Vorwärts-Selektion (auf P-Wert gestütztes Kriterium)

1. Modell wählen: $Y_i = \beta_0 + E_i$
2. Jeweils Variabel, welche den kleinsten P-Wert besitzt ins Modell aufnehmen.
3. Stopp: Sobald keine Verbesserung mehr möglich

Rückwärts-Selektion

1. Mit dem «vollen» Modell starten
2. Jeweils eine Variabel aus dem Modell entfernen (Beginnen mit Variabel, welche am unwichtigsten)
3. Stopp: Sobald keine Verbesserung mehr möglich

Schrittweise Selektion:

- Kombination aus Vorwärts- und Rückwärtselimination, lokales optimales Modell
- Stopp: Sobald keine Verbesserung mehr möglich

Vollständige Modellselektion (all-subset selection)
 Global optimierte Kriterien, berechnet für alle möglichen Modellvarianten
 2^m mögliche lineare Modell-Gleichungen, bei grossem m ist Rechenaufwand zu gross

All-Subset-Verfahren mit C_p -Kriterium:
 Wenn das Modell alle notwendigen Variablen enthält, Kriteriumwert = $E(C_p) = p$
 Fehlen notwendige Variablen: $E(C_p) > p$

Allgemein

- ➔ Verfahren führen nicht alle zum selben Modell
- ➔ P-Werte keine Aussagen über Modellgenauigkeit / Komplexität
- ➔ Schranke von 5 % im Stoppkriterium ist willkürlich

Anmerkungen:

- Empfehlung: All-Subset Selection
- Schrittweise Selektion: Grosses Modell
- Vorwärts Selektion: Datensätze mit vielen erklärenden Variablen und wenig Beobachtungen

Faktorvariablen/Polynomen/Wechselwirkungstermen:

- **Faktorvariablen:** Faktorstufen dürfen nicht einzeln entfernt werden (nur als ganze Variable)
- **Wechselw.terme:** immer mit Haupteffekt
- **Polynomen:** alle Terme bis zur maximalen Ordnung behalten, niedrigen Terme nicht entf.
- ➔ R-Funktion regsubsets(...) berücksichtigt dies nicht

R-Code:

Vorwärts-Selektion:
 AKW.lm <- lm(Y ~ 1, data = ...)
 add1(AKW.lm, scope = ~ lgG + D + BW + sqrtN + KG, test="F")
 AKW.lm1 <- update(AKW.lm0, . ~ . + KG + lgG)
 # Kleinster p-Wert schrittweise hinzufügen

Rückwärts-Selektion:
 AKW.lm0 <- lm(Y ~ ., data = ...)
 drop1(AKW.lm0, test="F")
 AKW.lm1 <- update(AKW.lm0, . ~ . - WZ - BW) # Grösster P-Wert schrittweise entfernen

Selektion basierend auf AIC (step(...)):

1. Vorwärtsselektion:
 M0 <- lm(Y ~ 1, data=DF)
 step(M0, scope=list(lower=~ 1, upper=~ x1 + xn), direction="forward") *hintere Teil von lm(~ ...)*

2. Rückwärts-Elimination:
 MF <- lm(Y ~ x1 + xn, data=DF) oder lm(Y ~ ., ...)
 step(MF, direction="backward")

3. Schrittweise Selektion:
 Mm <- lm(Y ~ x2, data=DF) *kann auch ~ 1 sein*
 X <- step(Mm, scope=list(lower=~ 1, upper=~ x1 + xn, direction="both") # scope: immer kleinstes & grösstes Modell angeben
 X\$anova # Zusammenfassung über Entfernung / Hinzufügen

All-Subset-Verfahren mit C_p -Kriterium:
 library(leaps) library(car)
 AKW.Cp <- regsubsets(lk ~ lg + BW + wN + KG, nbest=6, nvmax=10, data=AKW) # nbest immer angeben
 summary(AKW.Cp) # Übersicht
 h <- subsets(AKW.Cp, statistic="cp", legend="interactive", min.size=4, main="Mallow Cp", cex.subsets=0.7, las=1) # statistic = "cp", "bic", "adjr2", usw.
 abline(a=1, b=1, lty=2) # Intercept wird in Subset size nicht mitgezählt -> + eine Einheit
 plot(AKW\$Cp, scale="Cp") # Plot mit schwarzen Blöcken

Bemerkungen:

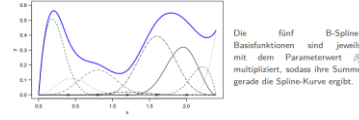
- Argument **nbest**: Anzahl Modell mit p Koeffizienten, die berücksichtigt werden. (default = 1)
- Argument **nvmax**: Modelle mit maximal p = nvmax werden berücksichtigt; Anzahl Modell mit p Koeffizienten, die berücksichtigt werden. (default = 8)
- Die Funktion **subsets()** zählt den Achsenabschnitt nicht mit, deshalb ist abline(a=1, b=1) die Winkelhalbierende.

➔ Die in Frage kommenden Modelle, müssen um die Winkelhalbierende $C_p = p$ streuen

➔ Potenzielle Modelle: minimaler C_p -Wert und / oder minimales p (kleinste Anzahl Parameter)

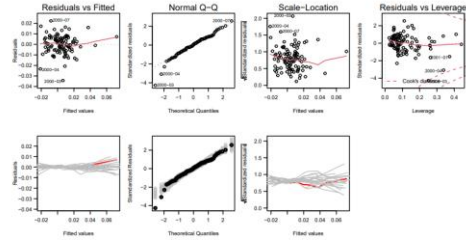
<p>Genügt nicht immer, weil:</p> <ul style="list-style-type: none"> - Auswahl der Variablen zufällig (mehrere Modelle in Betracht ziehen) - Beste Modell muss nicht unbedingt alle Voraussetzungen erfüllen (Residuenanalyse) - Aufgrund übergrosser Anzahl Möglichkeiten: Nicht alle potentiell nützlichen Transformationen können berücksichtigt werden 	<p>Bayes Informationskriterium (BIC, Alternative zu AIC) Bestraft grössere Modelle härter</p> <ul style="list-style-type: none"> - BIC: welche Prädiktoren leisten einen Beitrag, schlankes gut interpretierbares Modell - AIC: Modell wird für die Vorhersage von zukünftigen Werten eingesetzt, Prädiktoren sind weniger zentral $BIC = -n \log \left(\frac{1}{n} SS_E \right) + \log(n) p^\circ + \text{Konstante}$ <p style="text-align: center;"><code>step(..., k=log(nrow(dataset)))</code></p>
<p>Welches Modell ist das richtige? Werte der Modellwahlkriterien sind mit Unsicherheit behaftet – daher kann das beste Modell nur zufälligerweise das Richtige sein.</p> <p>➔ Mehrere Modelle in Betracht ziehen</p>	
<p>Alternative Ansätze für Variablenlektion:</p> <p>PRESS Summe der quadrierten Differenzen zwischen beobachteten und vorhergesagten Zielwerten.</p> <ul style="list-style-type: none"> - «leave-one-out» Ansatz (spez. Kreuzvalidierung) - Mass für Vorhersagegüte - Grosser Rechenaufwand bei grossen Modellen <p>Kein Overfitting kann negative Werte annehmen</p> $R_{pred}^2 := 1 - \frac{PRESS}{SS_Y}$ <p>LASSO Variablen selektieren (nur einflussreiche)</p> <ul style="list-style-type: none"> - Shrinkage Schätzung/ Regularisierungsmethode (Schrumpft Koeffizienten in Richtung 0) - Stellt sicher, dass kein Overfitting vorliegt (keine unsicheren Vorhersagen, da zu gut ans Modell angepasst) 	<p>R-Code:</p> <pre>library(lars) h.x <- model.matrix(IG ~ D + WZ + BZ, data=AKW) AKW.lasso <- lars(x=h.x, y=AKW\$IK) plot(AKW.lasso) # Bedeut. Koef.: 1. der auf 0 (Plot) Optimal: 5 – 10 fache Kreuzvalidierung (1/5 der Daten wird weggelassen), auch für hochdimensionale Daten geeignet: AKW.lasso.cv <- cv.lars(x=h.x, y=AKW\$IK, K = 10) min <- which.min(mf.lassocv\$cv) # Plot siehe S. 6</pre> <p>Adaptive LASSO: (mit Gewichtung)</p> <pre>library(lasso2) t.r <- l1ce(K ~ ., data=t.d, bound=seq(0.05, 1, 0.05)) plot(t.r) summary(t.r[5])</pre> <p>Mit folgenden Befehlen auf gleichnamigen Paketen ebenfalls möglich: lars, glmnet, sealasso</p>
<p>Kollinearität Hohe Korrelationen zwischen erkl. Variablen zugelassen, führen zu Problemen bei Interpretation & Modellierung.</p> <p>Erklärende Variable lässt sich bei Kollinearität annähernd als Linearkombination der anderen darstellen. Gilt die Beziehung exakt, gibt es keine eindeutige Lösung bei der Kleinste-Quadrate Schätzung.</p>	<p>Auswirkungen der Kollinearität</p> <ul style="list-style-type: none"> - Hohe Kollinearität führt zu grossen Standardfehlern bei geschätzten Koeffizienten ➔ Mit $\sqrt{VIF_j}$ aufgeblasen (Optimal Faktor 1) - Viele / alle Variablen gem. t-Test nicht signifikant - Gewisse Richtungen: sehr kleine oder grosse Prognosefehler -> Progn.intervalle bestimmen - Interpretation Effekte der einzelnen Variablen auf Zielvariable nicht möglich
<p>Konditionszahl $K = \frac{\lambda_{max}}{\lambda_{min}}$ (auch für Entdeckung von Multikollinearitäten geeignet)</p> <p>Maximaler und minimaler Eigenwert von $X^T X$ Zahl 100 – 1'000: Moderate bis starke Multikollinearität Zahl > 1'000: Multikollinearität schwerwiegend</p>	<p>Kollinearität muss bereinigt werden</p> <p>Was tun gegen Kollinearität? Wenn immer möglich, soll man Beobachtungen so durchführen, dass das Problem vermieden wird. Ansonsten:</p> <ul style="list-style-type: none"> - Variablen linear transformieren; d.h. z.B. stark korrelierte Variablen ersetzt man z.B. durch ihre Summe und ihre Differenz - Weitere Möglichkeiten: - $0.5 * (\text{Jet}\\$x1 + \text{Jet}\\$x2)$; $\text{Jet}\\$dif <- \text{Jet}\\$x1 - \text{Jet}\\$x2$ - Relation: $\text{seatpos}\\$rSeated <- \text{seatpos}\\$Seated / \text{seatpos}\\Ht; $\text{seatpos}\\$rArm <- \text{seatpos}\\$Arm / \text{seatpos}\\$Ht$ - Variable mit dem höchsten VIF aus dem Modell entfernen (= «Aputation») -> Korrelationen treten fast zwingend auf, wenn Vergleich zur Anzahl Beob. viele erkl. Var. vorhanden sind. - Setze so genannte Shrinkage-Schätzer ein wie z.B. - Hauptkomponentenregression (principal component regression), ridge regression oder LASSO, elastic net. Solche Schätzer haben sich vorteilhaft bei Prognosemodellen erwiesen. Allerdings leidet Interpretierbarkeit stark.
<p>Variance inflation factor (VIF) Bestimmtheitsmass zeigt:</p> <ul style="list-style-type: none"> - Wie stark eine solche Beziehung ist - Ist also ein sinnvolles Mass für Kollinearität - Gibt an, welche Variable das Problem verursacht $VIF_j = \frac{1}{(1-R_j^2)}$ <p>Nur für num. e. rkl. Variablen definiert Faustregel: Falls > 5-10, Probleme mit Kollinearität</p> <p>R-Code: library(car) round(vif(AKW.lm), 2)</p>	

Additive Modelle – nicht lineare Funktion

<p>Was wenn Funktion h nicht linear in den Parametern β ist?</p> <p>Anpassung durch:</p> <ul style="list-style-type: none"> - Nichtlineare Regression, dies muss jedoch von Fachwissen motiviert werden (also nicht relevant) - Nichtparametrische Schätzung des funktionalen Zusammenhangs (Glätter) <p>Modelle der Form: $Y_i = h(x_i, \beta) + \epsilon_i$</p>	<p>R-Code: Splines</p> <pre>Regression-Splines: library(splines) lm(y ~ bs(x, df=...), data=dat) plot(NOx ~ Equiv, data=exhaust, las=1, xlab="Äquivalenzverhältnis", ylab="Stickoxidkonz.", mgp=c(2,0,9,0), cex=0.7) h.range <- range(exhaust\$Equiv) h.knots4 <- seq(h.range[1], h.range[2], length=6)[-c(1,6)] # 4 Knoten, da 1 und 6 weggelassen werden bs(x=exhaust\$Equiv, knots=h.knots4, degree=3) # generiert Matrix mit kubischen Spline-Basen</pre> <p>Plot kann auch analog Polynomiale Regression erstellt werden!</p>
<p>Spline-Interpolation</p> <ul style="list-style-type: none"> - Stückweise Polynome der Ordnung k - Verbindungspunkte der Stücke = Knoten - Kubische Spline (k=3) ausreichend - $g(x) = \beta_1 * \beta_1^{(2)}(x) + \beta_2 * \beta_2^{(2)}(x) + \dots + \beta_q * \beta_q^{(2)}(x)$ - Sinnvoller Daten zu glätten als zu interpolieren - Glattheit der Kurve wird über die Wahl der Anzahl Stützpunkte q bestimmt - q klein: Kurve sehr glatt, q gross: eher Interpolation <p>Die B-Spline-Basisfunktionen (dünne Kurven) für eine Spline-Kurve (blaue fettere Kurve) mit 5 inneren Stützstellen (x auf der Nulllinie).</p>  <p>Die fünf B-Spline-Basisfunktionen sind jeweils mit dem Parameterwert β_j multipliziert, sodass ihre Summe gerade die Spline-Kurve ergibt.</p>	<p>Polynomiale Regression</p> <pre>MU.p3 <- lm(accel ~ poly(times,3), data=mcycle) hx <- data.frame(times=seq(min(mcycle\$times), max(mcycle\$times), length=200)) plot(accel ~ times, data=mcycle, las=1, main="Regression") lines(hx\$times, predict(MU.p3, newdata=hx), col="red") legend(40, -100, paste("Ordnung =", c(3,6,12)), lty=rep(1,3), col=c("red", "blue", "green"))</pre> <p>Smoothing Splines: smooth.spline(x,y, cv=TRUE) smooth.spline(data\$x, data\$y, df=7) # ohne Angabe df: optimale Anzahl FG smooth.spline(x=jitter(exhaust\$Equiv), y=exhaust\$NOx, cv=TRUE)</p> <p>Kubische Splines: interpSpline(dat\$x, dat\$y)</p> <p>Thin plate regression splines: gam(y ~ s(x), data=dat) library(mgcv) gam(Zielvariable ~ s(x1) + s(x2) + ... + s(xm), ...) # Identisch mit lowess() bei nur einer erkl. Variable plot.gam() par(mfrow=c(2,2)) # für partielle Residuenplots plot(P.gam, se=TRUE) plot(P.gam, se=TRUE, residuals=TRUE) ODER: library(gam) LB.gam2 <- gam(lBurntime ~ lo(I(Nitrogen) + lo(I(Chlorine) + lo(I(Potassium), data=LB) # Variablen vorher logarithmieren</p>
<p>Smoothing Splines</p> <p>Wie soll q (Glattheit der Spline-Kurve) gewählt werden?</p> <ul style="list-style-type: none"> - Zielkonflikt: Modellgenauigkeit / Modellglattheit wird durch λ geregelt: - λ gross: Gerade / glatteste Funktion - $\lambda = 0$: unterglättete Smoothing-Spline-Kurve - Optimale Glattheit: optimale Wahl λ - Kreuzvalidierungsverfahren zur Wahl λ 	
<p>Thin plate regression splines</p> <ul style="list-style-type: none"> - Knotenfreien Basisfunktionen - Mehrere erkl. Variablen zugleich glätten - optimal <p>Lokale Regression (LOWESS / LOESS) Jede Funktion h in einer kleinen Umgebung von einem vorgegebenen Punkt linear approximieren. (in Fenster aufteilen)</p>	
<p>Umsetzung: Gewichtete KQ-Methode Zu bestimmende Punkte:</p> <ul style="list-style-type: none"> - Wahl der Fensterbreite b b klein: Approximationsfehler sehr klein, grosse Varianz der Vorhersage b gross: kleine Varianz, lineare Approximation kann ungenügend sein -> vorgegebene Anzahl d Punkte im Fenster: d = 2/3 * Anzahl Punkte n - Wahl der Gewichtsfunktion K(u) So wählen, dass z.B. Berechnung vereinfacht wird oder Umsetzung mit LOWESS / LOESS: Tukeys Gewichtsfunkt. - Kontrollieren den Einfluss von schlechten Beobachtungen - c klein: desto schneller verlieren extreme Beobachtungen an Einfluss -> Verlust von Effizienz, grosse Standardab. - Cleveland: Empfehlung c = 4.05 - Behandlung von Ausreissern (in y-Richtung) 	<p>Lokale Regression (LOWESS / LOESS): loess(y ~ x, data=dat, degree=1, span=...) # mit family="symmetric" für robuste Anpassung, "gaussian" für Gauss / 1 = lokale linear, 2 = lokal quadratisch / Fensterbreite 20 % = span 0.2 -> je höher desto glätter</p> <pre>plot(exhaust\$Equiv, exhaust\$NOx, mgp=c(2,0,9,0), cex=0.7, xlab="Äquivalenzverhältnis", ylab="Stickoxidkonzentration") exhaust.loess <- loess(NOx ~ Equiv, data=exhaust, span=0.3, degree=2, family="symmetric") # Beobachtungen d = span xnew <- seq(0.5, 1.3, length=100) lines(xnew, predict(exhaust.loess, xnew), col="magenta") lowess(x1, Zielvariable) lo(..., span=0.5, degree=1) # LOWESS Glätter sti(...)</pre>

Zusammenfassung STMO

Plots mit klassischen Methoden

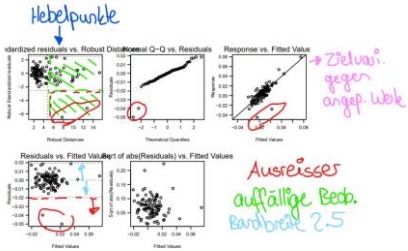


Abweichungen nicht mehr normal:

Varianz > als * Faktor 4

Standardabweichung > als * Faktor 2

Plots mit robusten Anpassungsmethoden



Vorgehen Modellentwicklung

- (i) Problem (d.h. Auftrag, Ziel, Zweck) verstehen; gibt es schon Modellansätze?
- (ii) Daten beschaffen, kennenlernen und aufbereiten
 - Codieren von fehlenden Werten abklären (Achtung: Manchmal sind diese mit '-99' oder mit '9999' codiert.)
Umgang mit fehlenden Werten in der Analyse festlegen.
 - Bedeutung der Zahl 0 in den verschiedenen Variablen vereinheitlichen.
 - Datenqualität beim Zusammenführen mehrerer Datensätze hinterfragen. Die Gesamtqualität ist nie höher als das schwächste Glied.
 - Daten gemäss first-aid Transformationen behandeln, ausser es gibt triftige Gründe dagegen (z.B. bestehendes Modell).
- (iii) Erste Anpassung; vorzugsweise mit robusten Methoden
- (iv) Residuen-Analyse;
 - Tragen die Daten dazu bei, das Problem zu lösen?
ev. zurück nach (ii) oder (i)
- (v) Variablenlektion, allenfalls Kollinearitäten behandeln
- (vi) Modelleignung klären
 - Residuen-Analyse mit selektierten Modellen
 - Modelle mit Fachwissen abgleichen, falls mit dem Modell Zusammenhänge beschrieben und erklärt werden sollen
 - 'out-of-sample' Validierung in Betracht ziehen (d.h. Validierung mit noch nicht verwendeten Daten), vor allem wenn das Modell zur Prognose eingesetzt werden soll.

R-Code Plots (Allgemein)

Plot mit abline und eingezeichnetem Vertrauensintervall
windows(8,4) [separates Fenster öffnen](#)
identify(x, y) [Punkt in Grafik auswählen \(Gibt Nr zurück\)](#)

plot(Forbes\$x, Forbes\$y)
points(Forbes\$x[12], Forbes\$y[12], col="red", pch=16) [die 12. Beobachtung blau einfärben](#)
abline(ForbesR.lm, col="blue", lty=1) [Kleinste Quadrate Schätzer als Linie einfügen](#)

abline(v=325.81, lty=4, col="red") [Linie bei vorgegebenem Punkt eintragen](#)

```
x0 <- data.frame(x= seq(min(Forbes$x), max(Forbes$x), length=50))
Data Frame
ForbesR.cia <- predict(ForbesR.lm, newdata=x0, Vertrauensintervall
interval="confidence", level=0.99)
lines(x0$x, ForbesR.cia[, "upr"], col="red") Vertrauensintervall in
Plot einzeichnen
lines(x0$x, ForbesR.cia[, "lwr"], col="red")
```

```
Plot mit Bestimmtheitsmass (R2) eingezeichnet
plot(fitted(MPI.lm), MPI$MPI)
abline(a=0, b=1, lty=2, col="red")
h <- predict(MPI.lm, newdata=x.MPI, interval="prediction", level=0.95)
points(h[, "fit"], h[, "fit"], col="red", pch=16) -> wenn keine Punkte in
der Nähe -> Voraussetzungen (gegeben) nicht erfüllt
```

```
Streudiagramm für mehrere Variablen (verwirrender Plot aus EXPD,
Korrelationen zwischen Variablen)
pairs(cat[, 3:1])
```

```
scatter3d-Plot
library(car)
scatter3d(y ~ x1 + x2, data = syn, axis.scales=FALSE)
```

```
Zufallszahlen generieren:
for(i in 1:6) {x <- rnorm(10); qqnorm(x); qqline(x, col="grey", lty=2)}
normalverteilt
for(i in 1:3) {x <- rt(1000, df=3); qqnorm(x); qqline(x, col="grey", lty=2)}
t-verteilt, Freiheitsgrade
x <- rchisq(20, df=7) chiquadrat-verteilt
qqnorm(x, main="n=20, df=7"); qqline(x, col="grey", lty=2)
```

```
Differenz zwischen Prognose- / Vertrauensintervall «upr» & «lwr»
x <- predict(lm(y ~ x2, catheter), interval="prediction") Prognose- /
Vertrauensintervall bilden
round(cbind(pi.unten=x[, "lwr"], fit=x[, "fit"], pi.oben=x[, "upr"],
pi.laenge=x[, "upr"]-x[, "lwr"]), 1) Runden und in Tabelle Differenz
berechnen
```

-> übersteigt Differenz das erwartete Intervall: Modell muss verbessert werden (hinzufügen von weiteren erkl. Variablen)

```
Intervall mit Rücktransformation
predict(FourCm.lm1, newdata=FourCm.new, interval="confidence")
exp(predict(FourCm.lm1, newdata=FourCm.new,
interval="confidence"), c("lwr", "upr"))
Pi in der Regel geeigneter -> da stochastische Variabilität des Fehlers Ei
berücksichtigt wird und nicht nur die Ungenauigkeit, welche sich aus
der Schätzung des Modells ergibt
```

```
Sortieren / Filtern etc
farm1 <- farm[farm$region==121,] nur die mit Region = 121
farm1 <- farm1[,-1] Region (1. Spalte) entfernen
farm1$industry <- as.factor(farm1$industry) in Faktorvariable
umwandeln
is.factor(farm1$industry) Prüfen ob Faktorvariable
str(farm1) summary(farm1)
dummy.coeff(farm1.fit1) Alle geschätzten und festgelegten
Koeffizientwerte
table(farm$region) Sortieren nach Region
```

```
Plots AB 4
plot(MPIZH$KPI, resid(MPIZH.lm1)) KPI gegen Residuen
plot(MPIZH$HZ, resid(MPIZH.lm1)) HZ gegen Residuen
termplot(MPIZH.lm1, partial.resid=T, smooth=panel.smooth)

plot(resid(MPIZH.lm1), type="h")
x.n <- nrow(MPIZH)
scatter.smooth(resid(MPIZH.lm1)[1:(x.n-1)], resid(MPIZH.lm1)[2:x.n])

plot(MPIZH$KPI, resid(MPIZH.lm1)); abline(v=100, col="blue", lty=3)
Linie bei 100 - (Residuen / Index-Plot)
plot(MPIZH$KPI, type="l"); abline(h=100, col="blue", lty=3)
Linie bei 100 - Index-Plot
```

```
load(paste(DPfad, "GASKETS.Rdata", sep="")) Rdata
Dokument einlesen
```

Optimalen Bestimmung des tuning Parameters β

```
par(mar=c(4,4,2,3), mgp=c(2.5, 0.8, 0))
set.seed(4567)
mfm.lasso.cv <- cv.lars(x=h.x, y=CEDHEC$FoHF, K = 10)
(h.wMin <- which.min(mfm.lasso.cv$cv)) # = 98
h.sel <- which(mfm.lasso.cv$cv <= (mfm.lasso.cv$cv[h.wMin] +
mfm.lasso.cv$cv.error[h.wMin]))[1]

h.sel
abline(v=mfm.lasso.cv$index[c(h.wMin, h.sel)],
col=c("magenta", "blue"))
abline(h=mfm.lasso.cv$cv[h.wMin] + mfm.lasso.cv$cv.error[h.wMin],
col="blue")

c(mfm.lasso.cv$index[h.sel], mfm.lasso.cv$cv[h.sel])
```

```
Koeffizientenschätzung an Stelle optimalen b: coef(AKW.lasso,
s=AKW.lasso.cv$index[h.sel], mode="fraction")
```