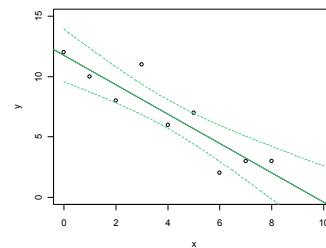




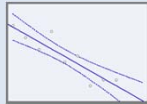
STATISTISCHE AUSWERTUNG ÖKOLOGISCHER DATENSÄTZE

REGRESSION

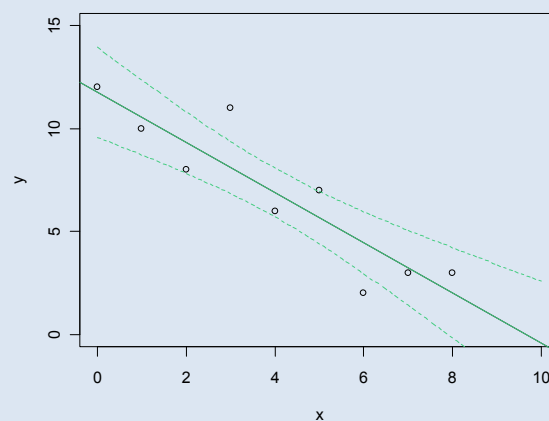


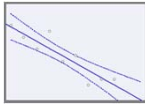
Studiengänge US, TE, WT, 3. & 7. Semester

Prof. Dr. Michael Rudner



Einfache lineare Regression





Lineare Regression

Beziehung zwischen einer abhängigen Variablen (AV)
und einer oder mehreren unabhängigen Variablen (UV).

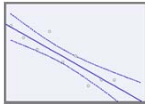
- Eine abhängige Variable
- Eine unabhängige (erklärende) Variable:
einfache (univariate) Regression
- Mehrere unabhängige Variablen:
multiple Regression
- Kausal
- Quantitativ je...desto...

$$Y = f(X)$$

$$Y = f(X_1, X_2, \dots, X_j, \dots, X_J)$$

Regression

3



Lineare Regression

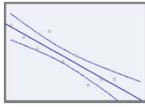
Voraussetzungen:

- Alle Variablen sollten kontinuierlich sein.
- Ggf. Umformung von kategorialen Variablen
in Dummy-Variablen (0,1)
- Wenn AV nominal (z.B. Anwesenheit: ja, nein)
dann logistische Regression
- Lineare Regression:
linearer Zusammenhang zwischen AV & UV

$$y = a + b \cdot x \text{ oder } y = b_0 + b_1 \cdot x$$

Regression

4



Lineare Regression

Vorteile:

- sehr flexibel
- Prognose für abhängige Variable

Typische Fragestellungen:

Ursachenanalysen:

Welche Faktoren beeinflussen AV?

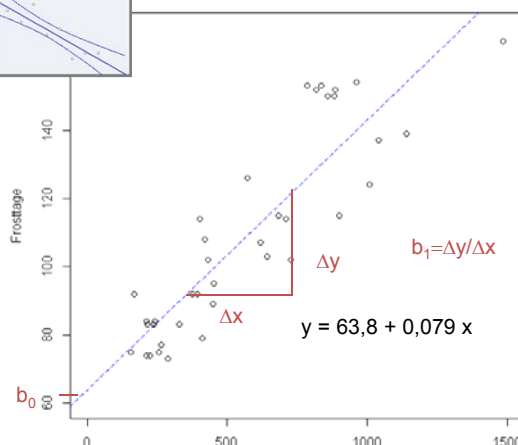
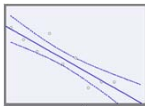
Wie verändert sich AV (abhängige Variable) entlang eines Gradienten der UV (unabhängigen Variablen) ?

Wirkungsprognosen

Wenn sich A um Wert X ändert, welchen Einfluss hat dies auf B?

Regression

5



$$\hat{Y} = b_0 + b_1 X$$

mit

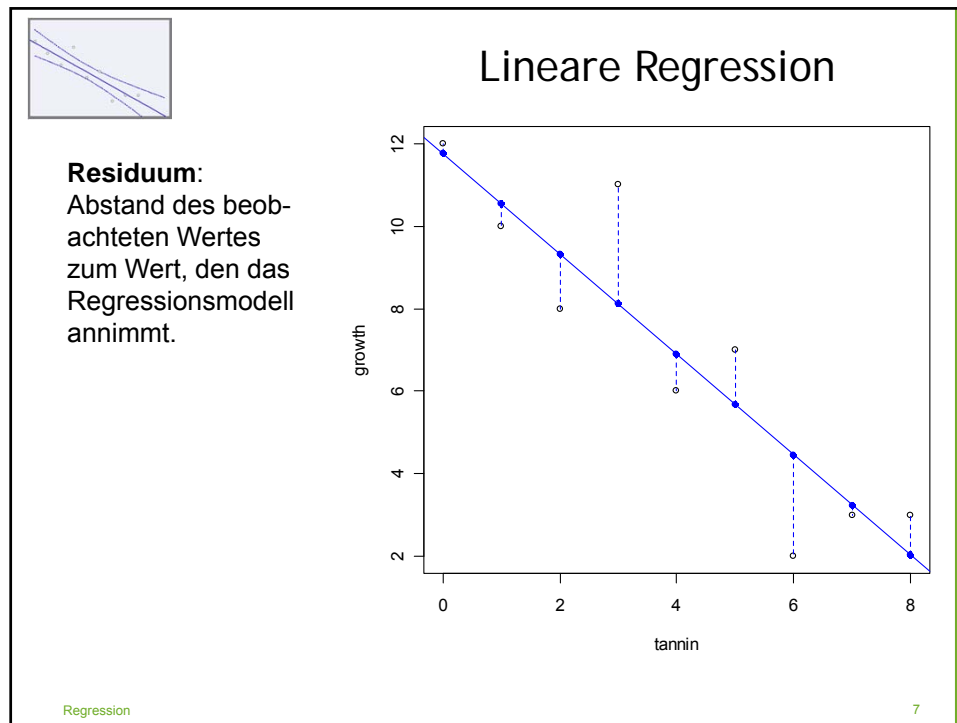
- \hat{Y} = Schätzung der abhängigen Variablen Y
- b_0 = konstantes Glied (Achsenabschnitt, *Intercept*)
- b_1 = Regressionskoeffizient (*Coefficient*)
- X = unabhängige Variable

Regression

b_0 Ursprung der Regressionsgerade
 \hat{Y} für $X=0$

b_1 Steigung Regressionsgerade, d.h. um wie viel ändert sich Y, wenn X sich um 1 Einheit ändert

Wichtig zur Wirkungsabschätzung



Lineare Regression

- » Residuen (e_k): Abweichungen von Regressionsgerade
- » Methode der kleinsten Abstandsquadrate (*Ordinary least squares*)
- » Ziel: Summe der quadrierten Residuen soll minimal sein, danach b_0 und b_1 wählen
- » Quadrieren, damit sich nicht negative und positive Abweichungen eliminieren

$$e_k = y_k - \hat{y}_k \quad (k=1, 2, \dots, K) \quad (4)$$

mit

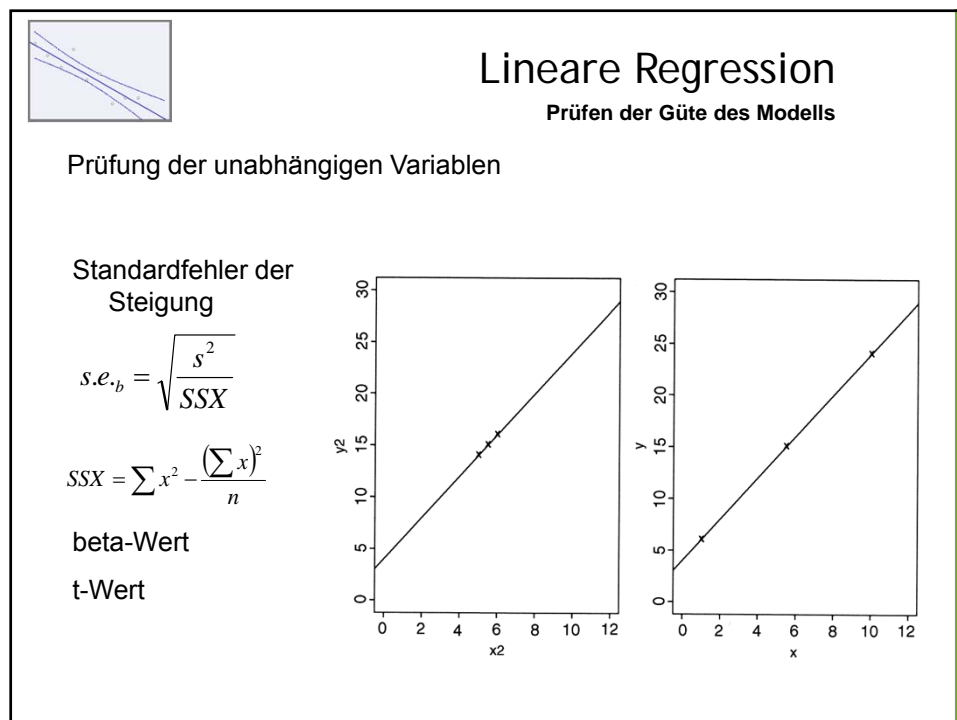
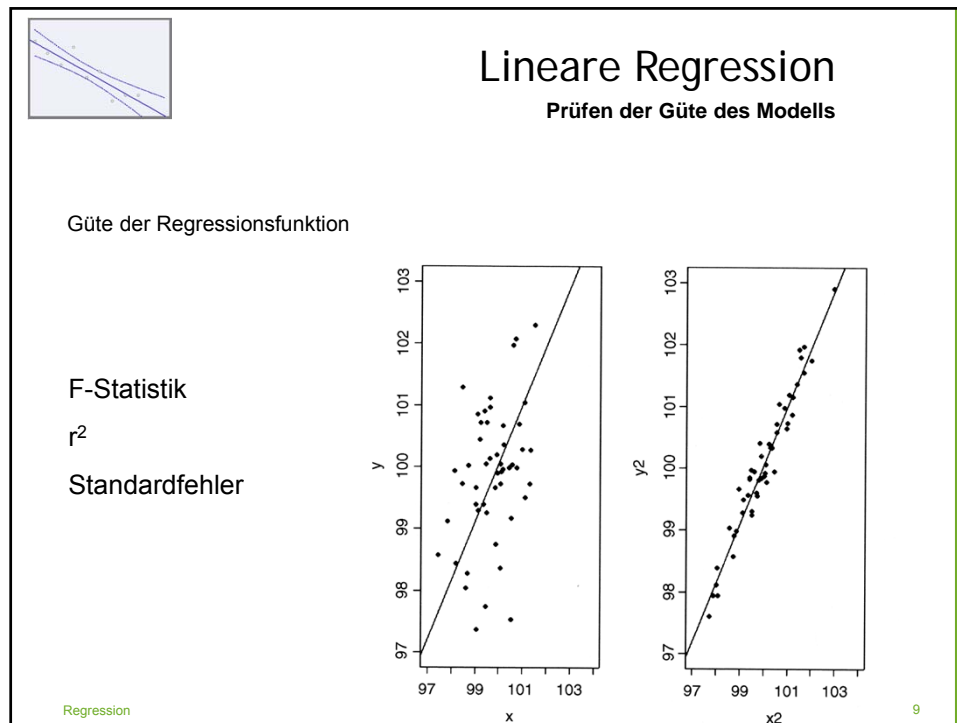
- y_k = Beobachtungswert der abhängigen Variablen Y für x_k
- \hat{y}_k = ermittelter Schätzwert von Y für x_k
- e_k = Abweichung des Schätzwertes von Beobachtungswert
- K = Zahl der Beobachtungen

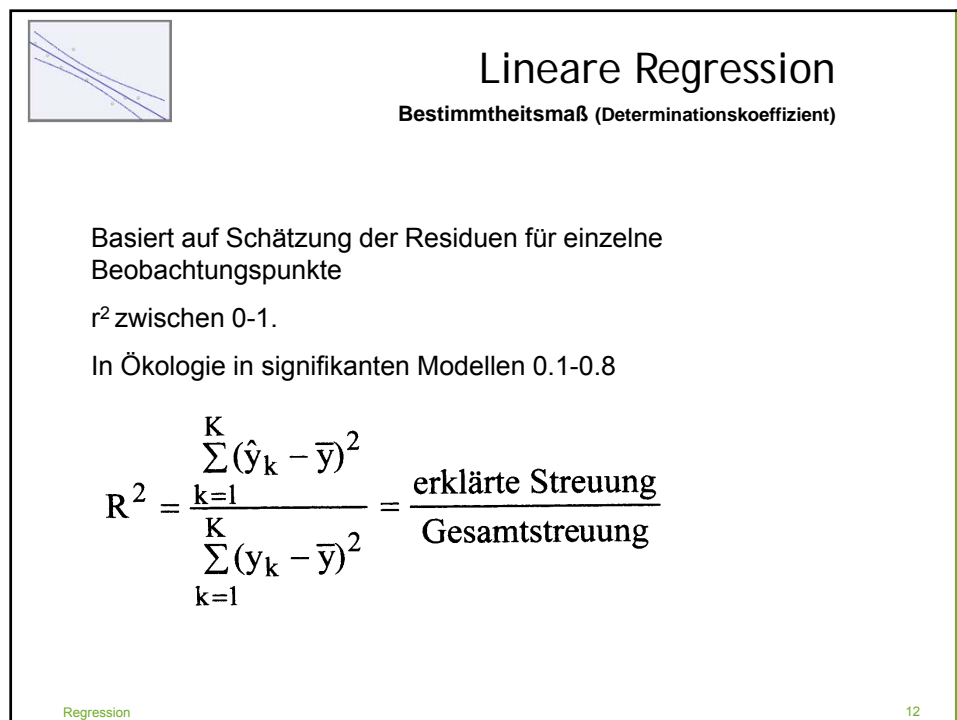
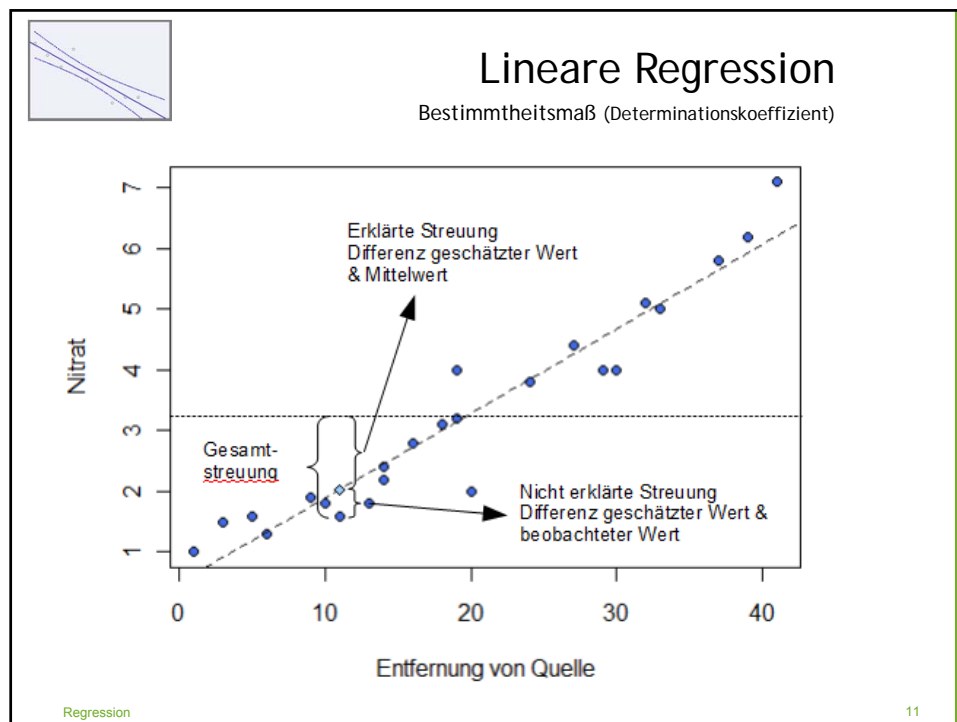
Durch Umformung von (4) und unter Einbeziehung von (2) läßt sich folgende Funktion bilden:

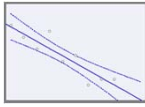
$$Y = \hat{Y} + e$$

$$= b_0 + b_1 X + e \quad (5)$$

8







Lineare Regression

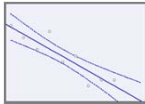
Bestimmtheitsmaß (Determinationskoeffizient)

- » r^2 = Abschätzung des Effekts, wie viel der Streuung wird durch UV erklärt
- » r^2 = Quadrat des Korrelationskoeffizienten r zwischen beobachteten und geschätzten Y-Werten
- » r^2 wird von Anzahl der UVs beeinflusst, daher bei vielen UVs häufig korrigiertes r^2 benutzt, dies nimmt mit Anzahl der UVs ab, da Erklärungsanteil jeder UV zufällig bedingt sein kann

Maß wie gut geschätzte Werte beobachtete Werte abbilden

Regression

13



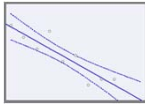
Lineare Regression

F-Statistik

- » Besitzt **Modell** Vorhersagequalität für Grundgesamtheit?
- » Prüfung gegen Nullhypothese, dass UV keinen Einfluss auf AV hat
- » Wenn Nullhypothese stimmt: F-Wert = 0
Daher ist F-Statistik ein Test, ob F-Wert unterschiedlich von 0

Regression

14



Lineare Regression

F-Statistik

$$F_{\text{emp}} = \frac{R^2 / J}{(1 - R^2) / (K - J - 1)}$$

- » J = Anzahl UV
- » K = Stichprobengröße
- » Vergleich beobachteter mit theoretischem F-Wert
- » Vertrauenswahrscheinlichkeit (0.95 oder 0.99)
- » Irrtumswahrscheinlichkeit (5%, 1%)
- » $\alpha = 1 - \text{Vertrauenswahrscheinlichkeit}$
- » oder: zugehörigen p-Wert ermitteln

Regression

15

Lineare Regression - Output

```
> tann<-read.table("Crawley\\tannin.txt",header=T)
> x<-tann$tannin
> y<-tann$growth
> tann.lr <- lm(y~x)
> summary(tann.lr)
```

Call:

lm(formula = y ~ x)

Residuals:

Min	1Q	Median	3Q	Max
-2.4556	-0.8889	-0.2389	0.9778	2.8944

Coefficients:

	Estimate	Std. Error	t value	Pr(> t)
(Intercept)	11.7556	1.0408	11.295	9.54e-06 ***
x	-1.2167	0.2186	-5.565	0.000846 ***

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

Residual standard error: 1.693 on 7 degrees of freedom

Multiple R-Squared: 0.8157, Adjusted R-squared: 0.7893

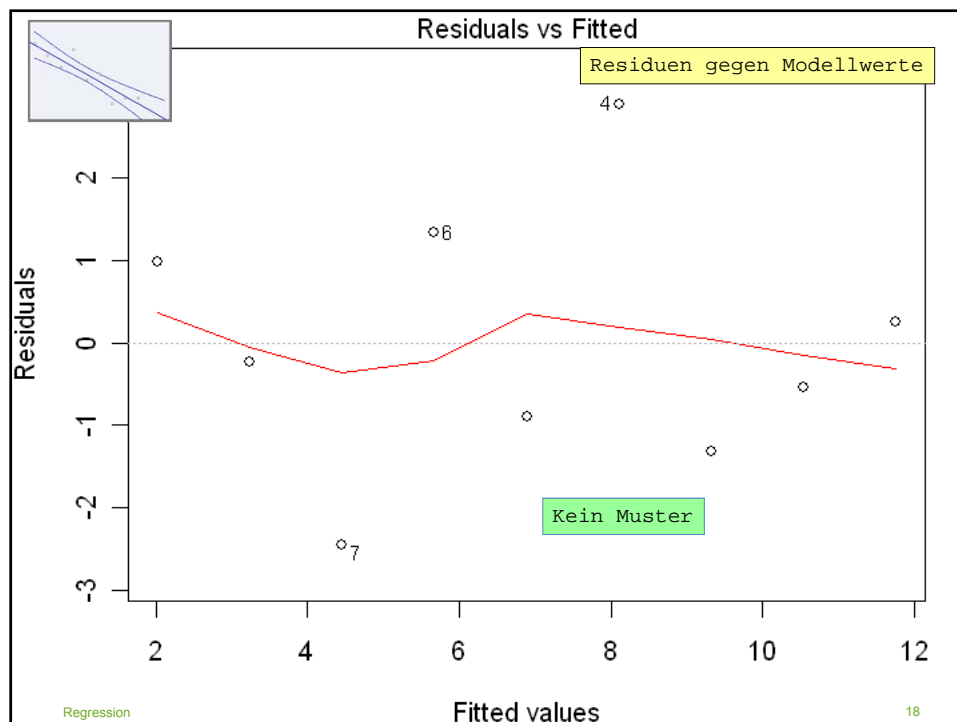
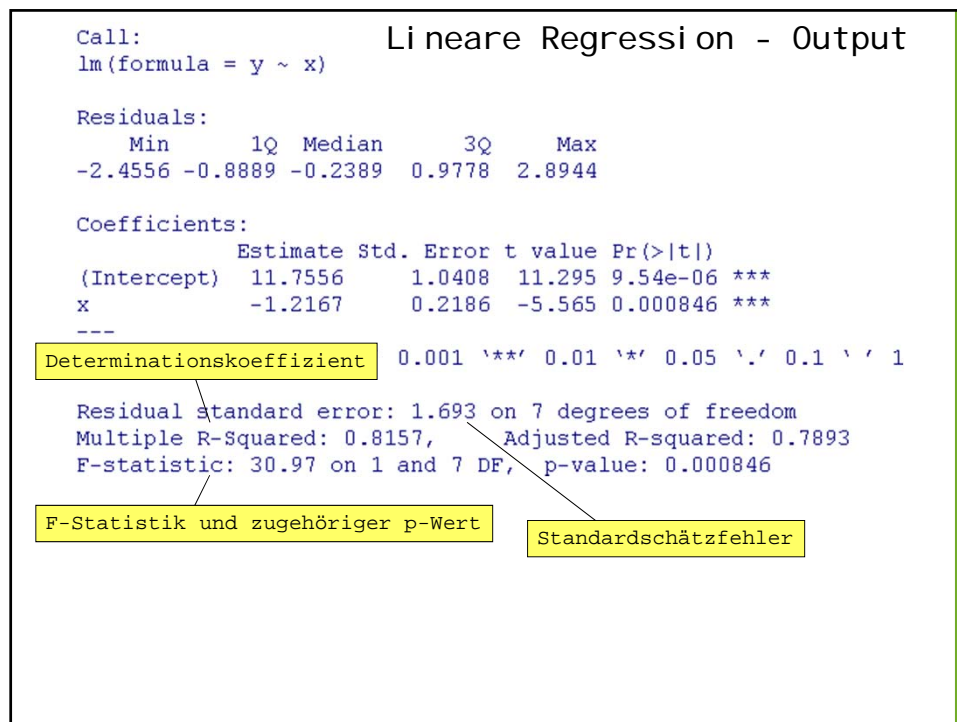
F-statistic: 30.97 on 1 and 7 DF, p-value: 0.000846

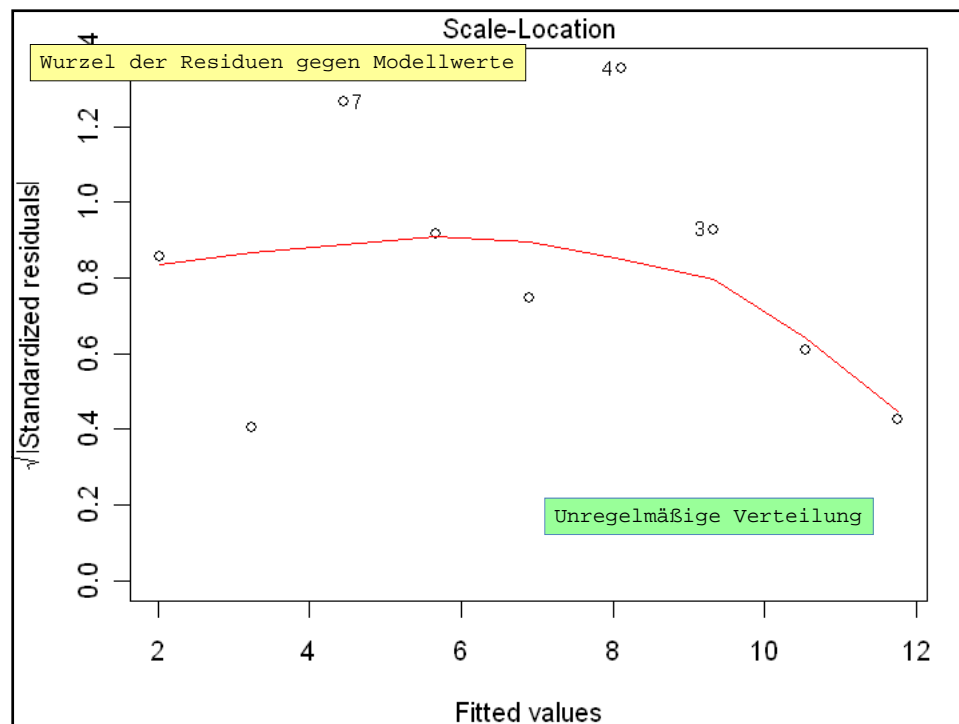
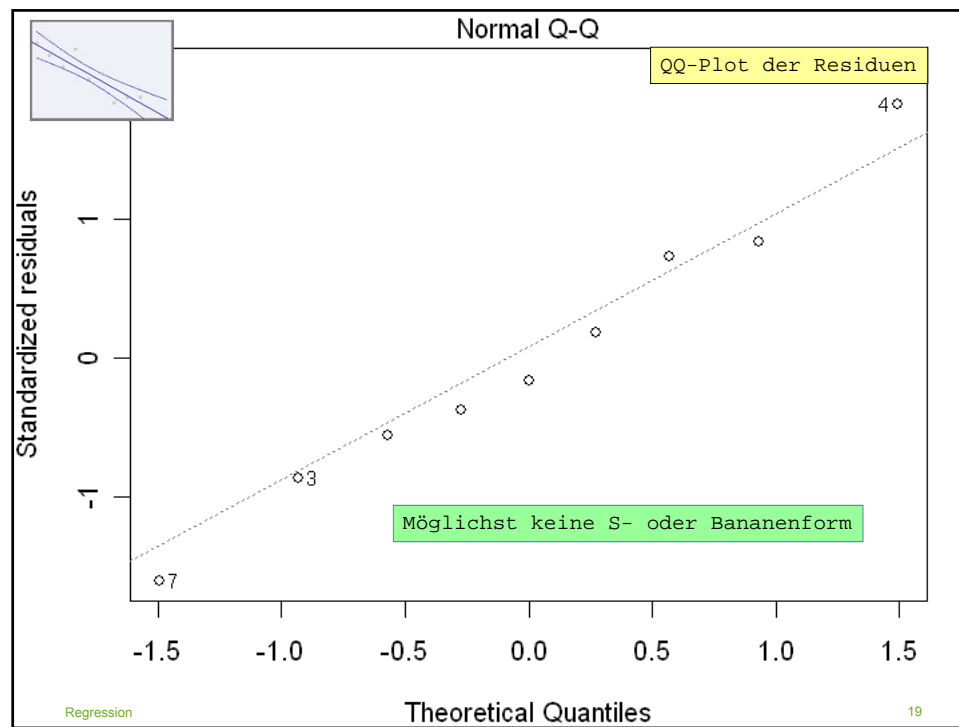
Verteilung der Residuen

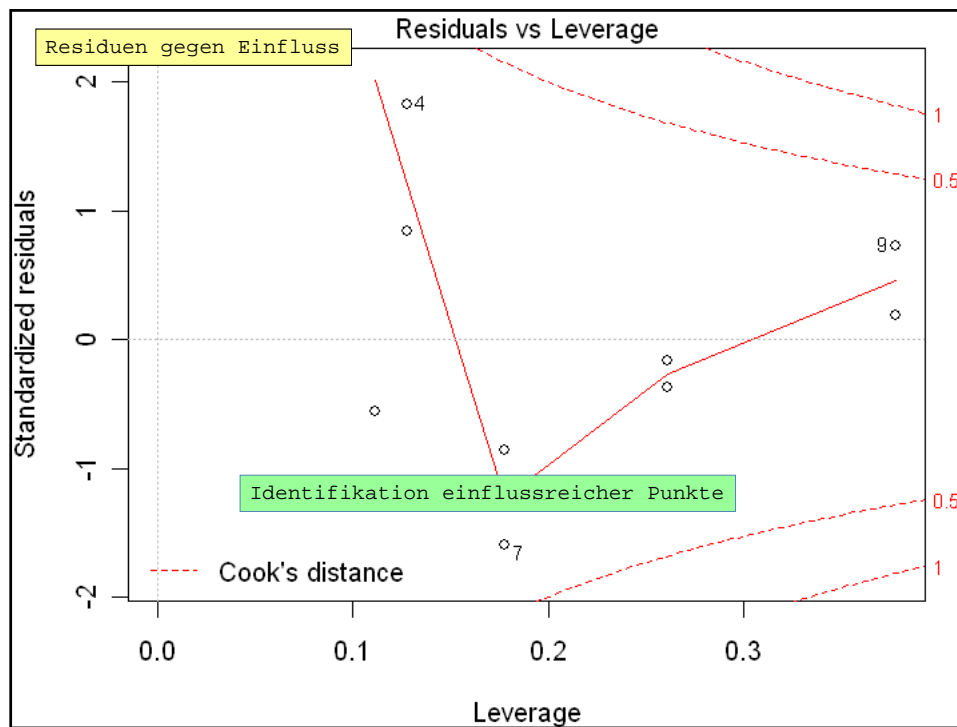
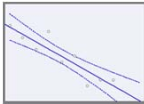
Regressionskoeffizienten

Fehler der Regressionskoeffizienten

p-Wert





Lineare Regression

Standard-Schätzfehler

$$s.e._{\hat{y}} = \sqrt{s^2 \cdot \left[\frac{1}{n} + \frac{(x - \bar{x})^2}{SSX} \right]} = s \cdot \sqrt{\frac{1}{n} + \frac{(x - \bar{x})^2}{SSX}}$$

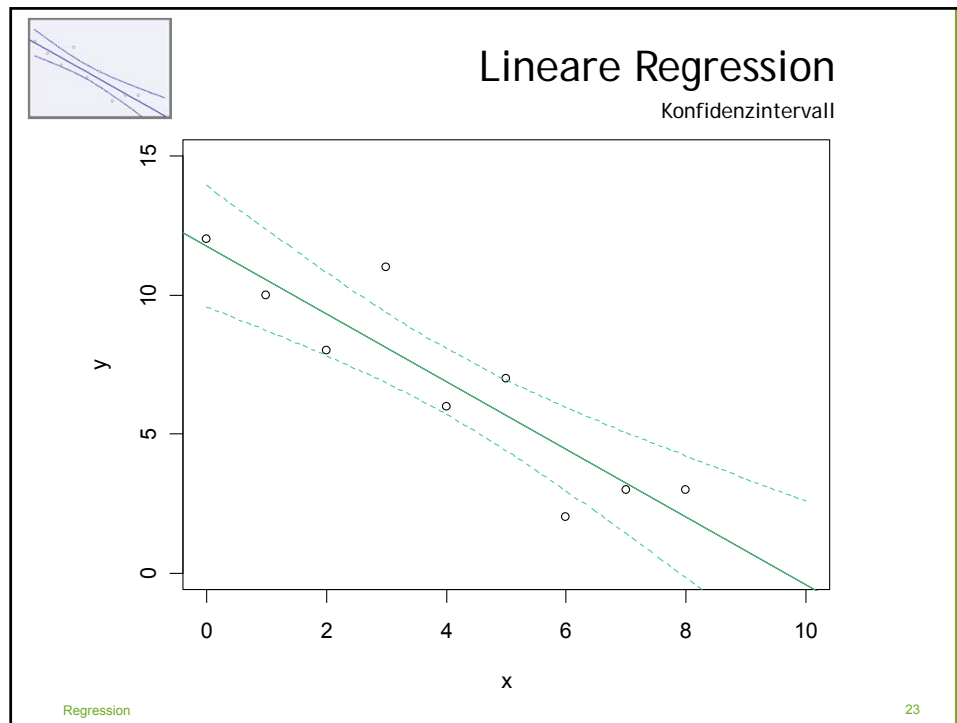
Für Konfidenzintervall der Regressionsgeraden:

$$\Delta_{krit} = \hat{y}_j \pm t_{(1-\alpha/2)} \cdot s.e._{\hat{y}}$$

$$\Delta_{krit} = \hat{y}_j \pm t_{(1-\alpha/2)} \cdot s \cdot \sqrt{\frac{1}{n} + \frac{(x - \bar{x})^2}{SSX}}$$

Regression

22



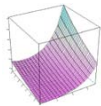
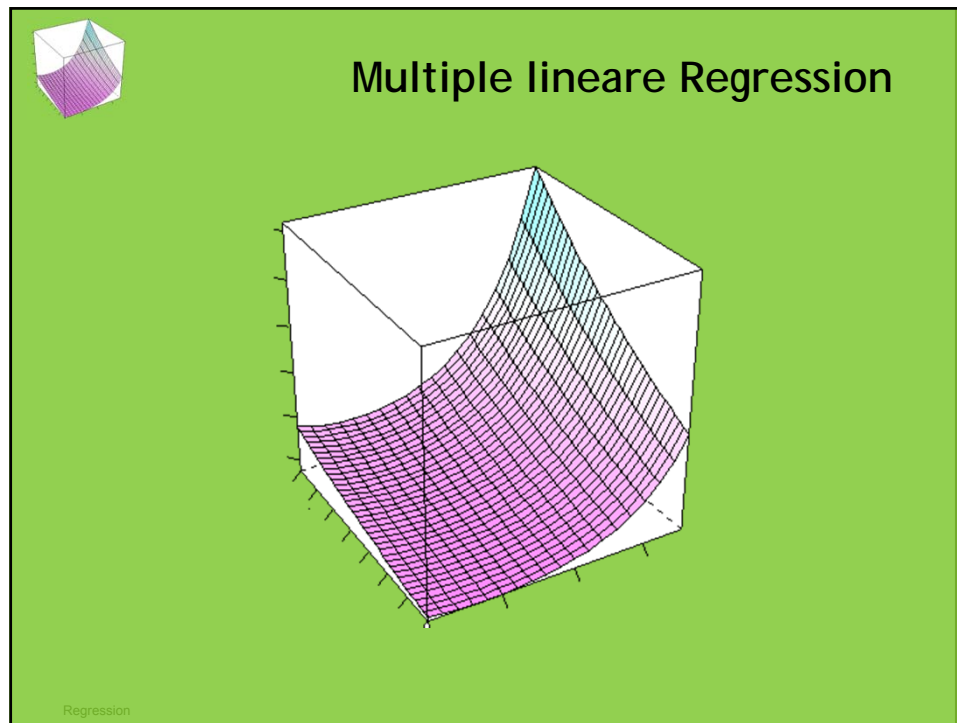
Lineare Regression

Zusammenfassung

- Regressions-Parameter werden aus Daten geschätzt
- **Vielzahl möglicher Modelle für einen Datensatz**
- **Modellauswahl ist ein wesentlicher Schritt**
- Diagnostische Plots um zu prüfen, ob das gewählte Modell passt
- Hauptprobleme:
 - **nicht-konstante Varianz** und
 - **nicht-normalverteilte Fehler**
- Mögliche Lösungen für gekrümmte Zusammenhänge zwischen y und x:
 - Quadrierter Term
 - Transformation
- Vergessen Sie nicht, die Modellwerte rückzutransformieren, bevor Sie die Regressionskurve zeichnen!

Regression

24



Multiple lineare Regression

- Eine abhängige Variable
- Mehrere unabhängige Variablen

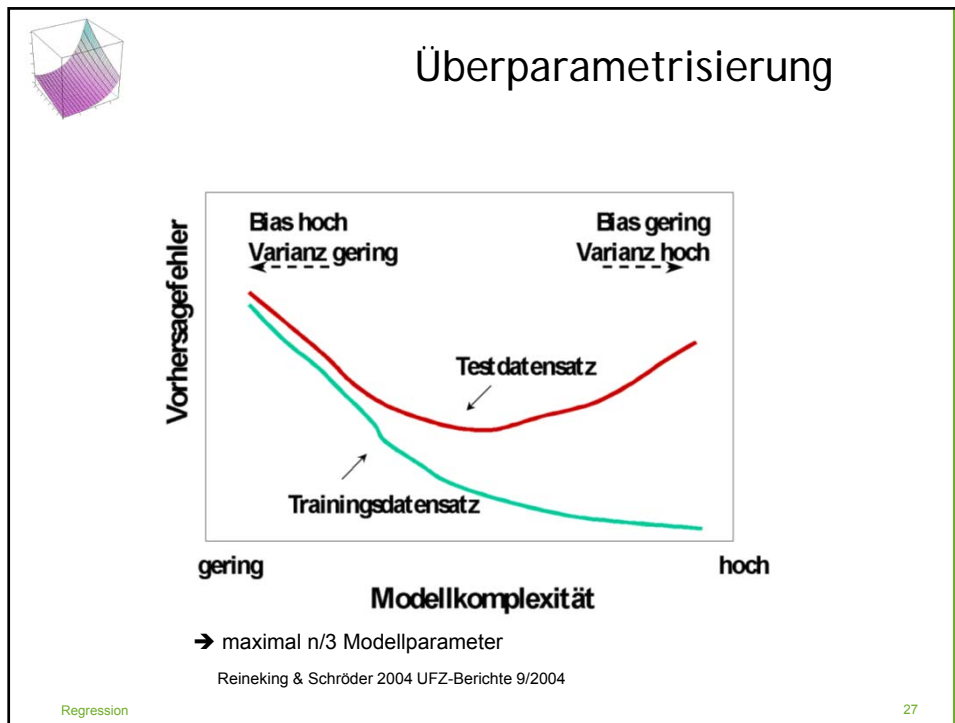
Wesentliche Fragen:

- Linearer Zusammenhang
- Interaktion der Variablen
- Korrelationen der Variablen
- Überparametrisierung

?

Regression

26



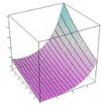
Multiple lineare Regression

Vorgehen

- Überlegen, welche Variablen sinnvollerweise ins Modell genommen werden sollen
- Prüfen der Korrelationen der Variablen untereinander
 - Problem der Multikollinearität
 - Stichwort: *Variance Inflation Factor*
- Beschränken der Anzahl der Modellparameter auf maximal $n/3$

Regression

28



Multiple lineare Regression

Korrelation der Variablen

		s	t	v	q
AV	s	1,00			
	t	-0,29	1,00		
UV	v	0,68	0,06	1,00	
	q	0,05	0,61	0,41	1,00

Starke Korrelation

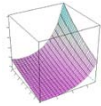
Man sollte stark korrelierte UV nicht gemeinsam in ein Modell aufnehmen.

Herausnehmen der Variable „Entfernung zur Quelle“

AV: Sauerstoffgehalt s;
UV: Wassertemperatur t, Fließgeschwindigkeit v, Abstand von der Quelle q

Regression

29

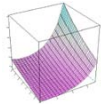


Modellbildung

- Ziel: Kleinstmögliches Modell (*Ockham's razor*)
- Volles Modell aufstellen
- Schrittweise rückwärts unbedeutende Variablen entfernen.
- Hinweis aus t-Statistik
- Minimales zufriedenstellendes Modell weiter prüfen
- ggf. Transformationen der Werte
- ggf. entfernen zu einflussstarker Werte
- Modellvergleich mit F-Test

Regression

30



Responseoberfläche

- Erzeugen einer Oberfläche der Modellwerte in Abhängigkeit von zwei erklärenden Variablen
- weitere Prädiktoren müssen berücksichtigt werden, da sie Einfluss auf den Achsenabschnitt haben
- also zwei Variablen zur Darstellung auswählen, und für die restlichen Variablen Referenzwerte auswählen, die anzugeben sind.
- ggf. mehrere Grafiken mit unterschiedlichen Referenzwerten erzeugen