



# Multiple Lineare Regression (Teil 3)

Peter von Rohr

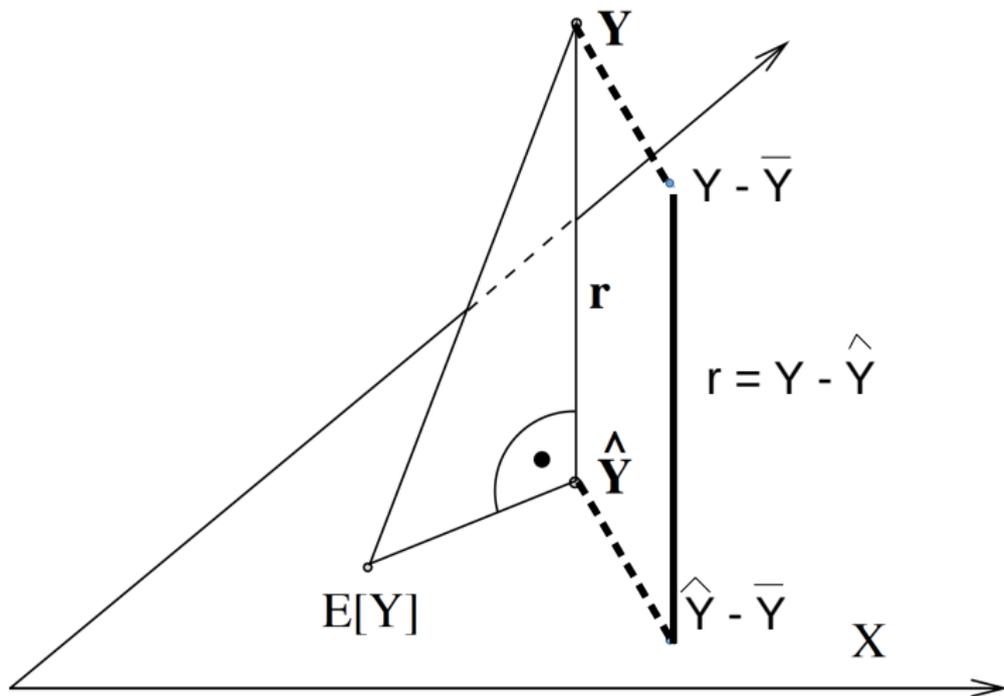
# Globaler Test eines Modells

- Beim t-Test hatten wir jede einzelne erklärende Variable getestet.
- Test, ob überhaupt eine der erklärenden Variablen einen Einfluss auf die Zielgrösse hat
- Zerlegung der Länge der totalen quadrierten Abweichungen der Beobachtungswerte  $\mathbf{y}$  um deren Mittel  $\bar{\mathbf{y}}$  in

$$\|\mathbf{y} - \bar{\mathbf{y}}\|^2 = \|\hat{\mathbf{y}} - \bar{\mathbf{y}}\|^2 + \|\mathbf{y} - \hat{\mathbf{y}}\|^2$$

wobei:  $\|\hat{\mathbf{y}} - \bar{\mathbf{y}}\|^2$  der Länge der quadrierten Abweichungen der gefitteten Werte ( $\hat{\mathbf{y}} = \mathbf{X}\hat{\beta}$ ) um das globale Mittel ( $\bar{\mathbf{y}} = \mathbf{1} * 1/n \sum_{i=1}^n y_i$ ) und  $\mathbf{r} = \mathbf{y} - \hat{\mathbf{y}}$  den Residuen entspricht

# Geometrische Begründung



# Zerlegung als Varianzanalyse (ANOVA)

- ANOVA Tabelle sieht wie folgt aus

	sums of squares	degrees of freedom	mean square
regression	$\ \hat{\mathbf{y}} - \bar{\mathbf{y}}\ ^2$	$p - 1$	$\ \hat{\mathbf{y}} - \bar{\mathbf{y}}\ ^2 / (p - 1)$
error	$\ \mathbf{y} - \hat{\mathbf{y}}\ ^2$	$n - p$	$\ \mathbf{y} - \hat{\mathbf{y}}\ ^2 / (n - p)$
total	$\ \mathbf{y} - \bar{\mathbf{y}}\ ^2$	$n - 1$	

- Relevante Teststatistik lautet

$$F = \frac{\|\hat{\mathbf{y}} - \bar{\mathbf{y}}\|^2 / (p - 1)}{\|\mathbf{y} - \hat{\mathbf{y}}\|^2 / (n - p)} \sim F_{p-1, n-p}$$

unter der globalen Nullhypothese  $H_0 : \beta_j = 0$  für alle  $j$

# Bestimmtheitsmass des Modells

- Nützliche Grösse für die Qualität eines Modells ist das Bestimmtheitsmass (coefficient of determination)

$$R^2 = \frac{\|\hat{\mathbf{y}} - \bar{\mathbf{y}}\|^2}{\|\mathbf{y} - \bar{\mathbf{y}}\|^2}$$

diese sagt aus, wieviel der totalen Variation von  $\mathbf{y}$  um  $\bar{\mathbf{y}}$  durch die Regression erklärt wird.

# Vertrauensintervall der Schätzung

- Basierend auf der Teststatistik des t-Tests

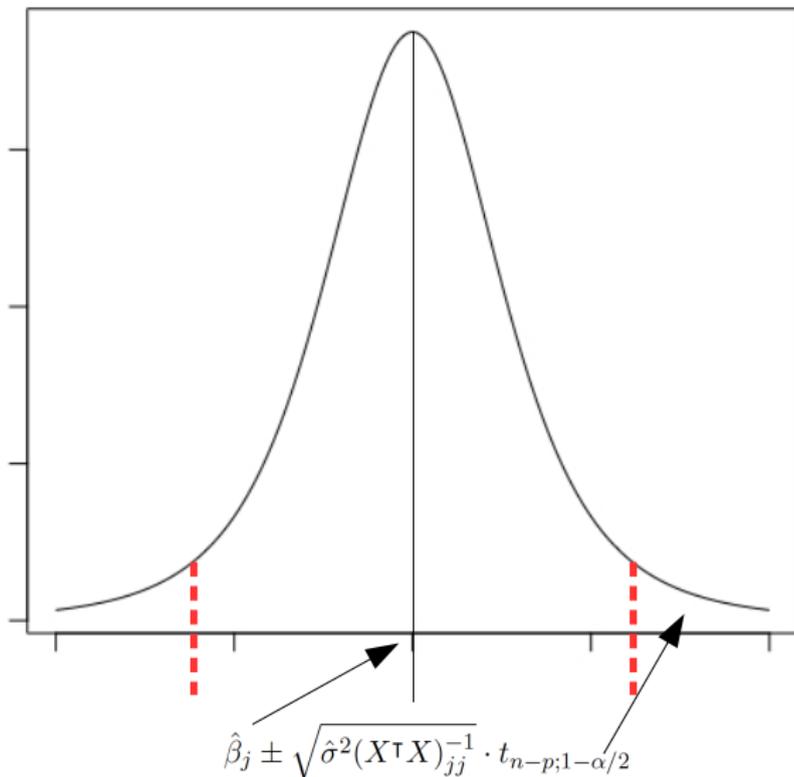
$$T_j = \frac{\hat{\beta}_j}{\sqrt{\hat{\sigma}^2(\mathbf{X}^T \mathbf{X})_{jj}^{-1}}} \sim t_{n-p}$$

- Vertrauensintervall für den unbekannt Parameter  $\beta_j$  als

$$\hat{\beta}_j \pm \sqrt{\hat{\sigma}^2(\mathbf{X}^T \mathbf{X})_{jj}^{-1}} * t_{n-p, 1-\alpha/2}$$

→ somit beinhaltet das Intervall zwischen den angegebenen Grenzen den wahren Wert mit Wahrscheinlichkeit  $1 - \alpha$ , wobei  $t_{n-p, 1-\alpha/2}$  das  $1 - \alpha/2$  Quantil der Verteilung  $t_{n-p}$  darstellt

# Vertrauensintervall im Bild



# R Output

```
Call:
lm(formula = LOGRUT ~ ., data = asphalt1)
```

1

```
Residuals:
    Min       1Q   Median       3Q      Max
-0.48348 -0.14374 -0.01198  0.15523  0.39652
```

2

Coefficients:

	Estimate	Std. Error	t value	Pr(> t )
(Intercept)	-5.781239	2.459179	-2.351	0.027280 *
LOGVISC	-0.513325	0.073056	-7.027	2.90e-07 ***
ASPH	1.146898	0.265572	4.319	0.000235 ***
BASE	0.232809	0.326528	0.713	0.482731
RUN	-0.618893	0.294384	-2.102	0.046199 *
FINES	0.004343	0.007881	0.551	0.586700
VOIDS	0.316648	0.110329	2.870	0.008433 **

3

---

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

```
Residual standard error: 0.2604 on 24 degrees of freedom
Multiple R-Squared: 0.9722,    Adjusted R-squared: 0.9653
F-statistic: 140.1 on 6 and 24 DF,  p-value: < 2.2e-16
```

4

# R Output Bedeutung

- 1 Funktionsaufruf mit welchem das Resultatobjekt erzeugt wurde. Wichtig, falls Resultate als R-objekt (.rda) gespeichert werden
- 2 Verteilung der Residuen aufgrund der Quantile
- 3 Schätzwert und Schätzfehler für die Parameter  $\beta_j$  zu jeder erklärenden Variablen. Werte der t-Teststatistik
- 4 Schätzung der Rest-Standardabweichung  $\sigma$ . Zusätzliche Modellinformationen, wie F-Teststatistik,  $R^2$  und das um Anzahl erklärende Variablen korrigierte  $\bar{R}^2$ , wobei

$$\bar{R}^2 = R^2 - (1 - R^2) \frac{p - 1}{n - p}$$

# Überprüfung der Modellannahmen anhand Analyse der Residuen

- Residuen  $r_i = y_i - \hat{y}_i$  als Approximation der unbekanntem Fehler  $\epsilon_i$  bei der Überprüfung der Modellannahmen verwenden
- **Tukey-Anscombe** Plot: zeigt Residuen  $r_i$  versus gefittete Werte  $\hat{y}_i$ . Dieser sollte keine erkennbaren Muster aufweisen

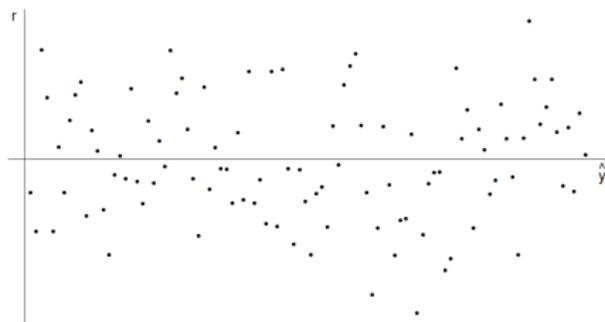


Figure 1.4: Ideal Tukey-Anscombe plot: no violations of model assumptions.

# Probleme bei Modellannahmen

Folgende Plots deuten auf Probleme hin

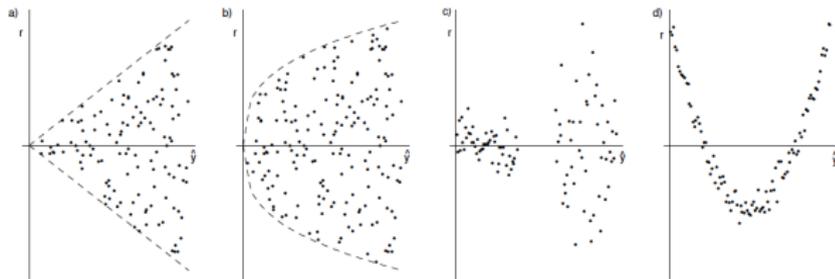
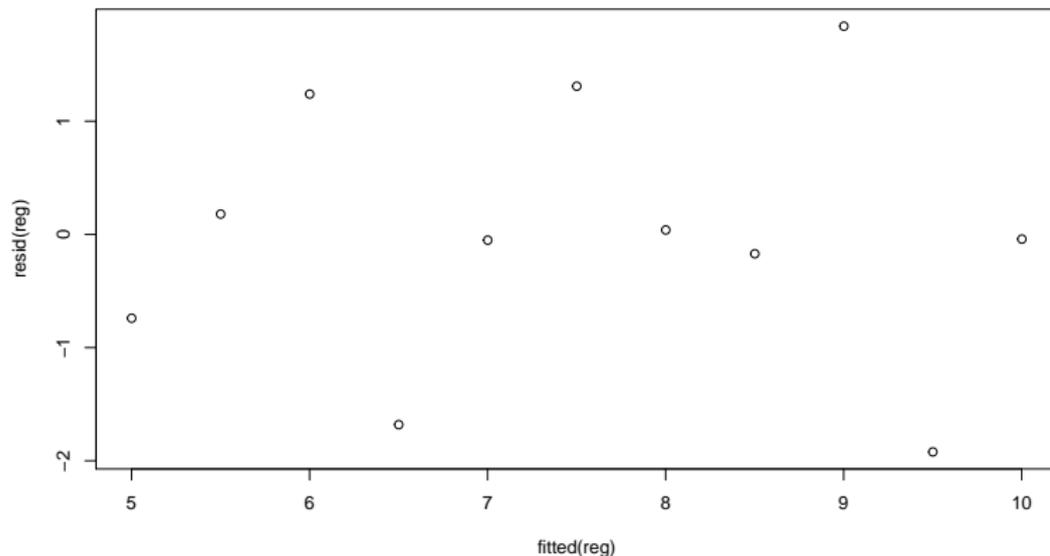


Figure 1.5: a) linear increase of standard deviation, b) nonlinear increase of standard deviation, c) 2 groups with different variances, d) missing quadratic term in the model.

# Tukey-Anscombe Plot in R

```
data(anscombe)
reg <- lm(y1 ~ x1, data = anscombe)
plot(fitted(reg), resid(reg))
```

# Tukey-Anscombe Plot - Das Resultat



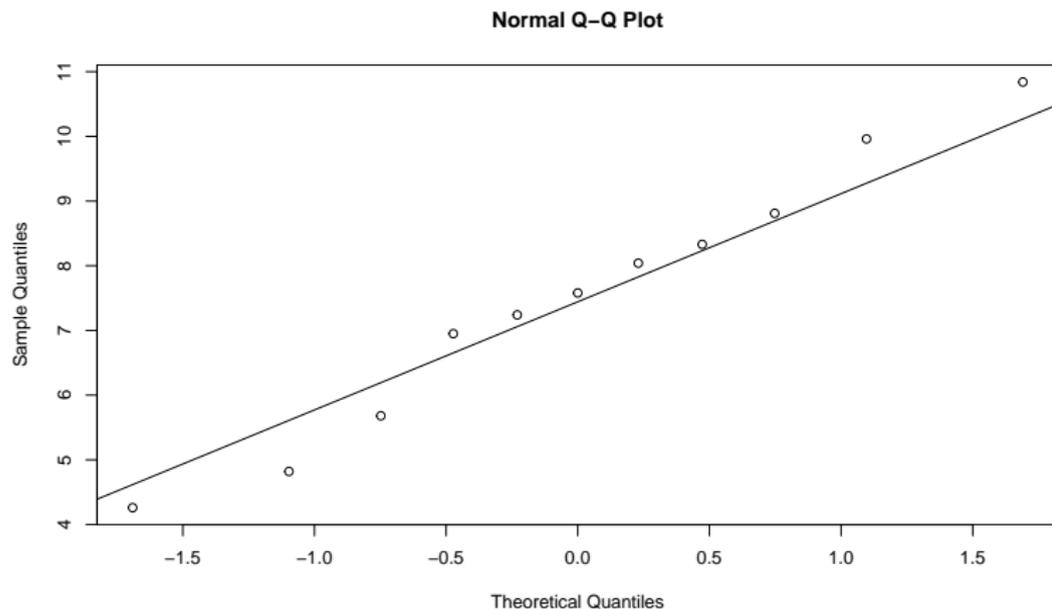
# QQ (quantile-quantile) Plot

- Überprüfung der Verteilung der Zufallsvariablen (Zielgrösse und Residuen)
- Empirische Verteilung der Residuen ( $y$ -Achse) wird gegen theoretische Quantile der Normalverteilung ( $x$ -Achse) aufgezeichnet
- Falls Normalverteilung zutrifft, dann liegen alle Punkte auf einer Linie

# In R:

```
qqnorm(anscombe$y1)
```

```
qqline(anscombe$y1)
```



# Probleme mit Verteilung

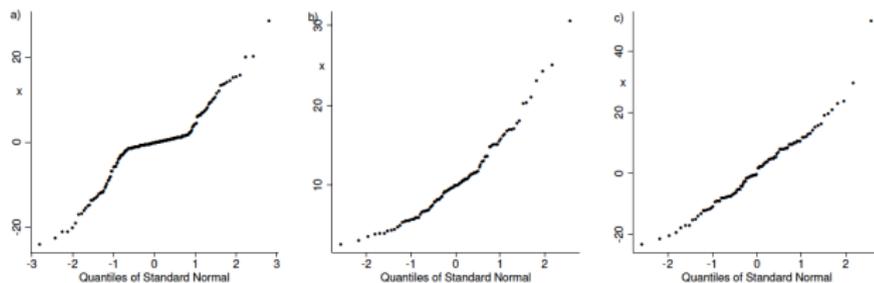


Figure 1.7: QQ-plots for a) long-tailed distribution, b) skewed distribution, c) dataset with outlier.

# Quellen

Tukey-Anscombe Plots und QQ-Plots stammen aus dem Skript:

*Computational Statistics*

*Peter Bühlmann and Martin Mächler*

*Seminar für Statistik ETH Zürich*

*Version of January 31, 2014*