

Mathe III

Regression

Ralf Herbrich

1. Konzept der Linearen Regression
2. Maximum *Likelihood* Lineare Regression
 - Methode der Kleinsten Quadrate
3. Verteilungen der Maximum *Likelihood* Schätzer
4. Hypothesentests für Regressionsparameter
 - *t*-Test für Lineare Regression
5. Analyse der Residuen
6. Exkurs: Multiple Regression

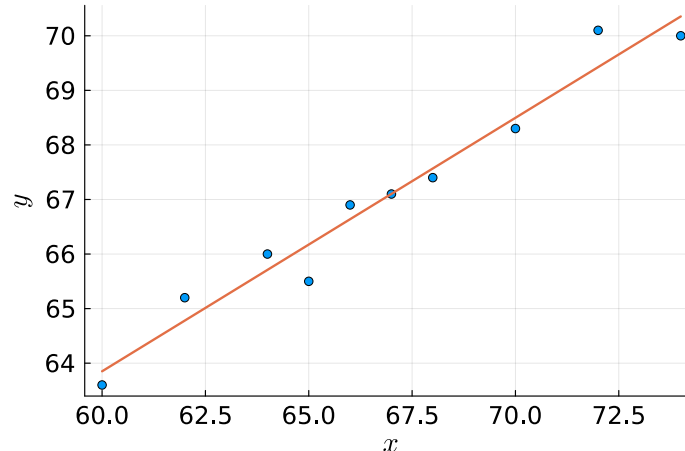
1. **Konzept der Linearen Regression**
2. Maximum *Likelihood* Lineare Regression
 - Methode der Kleinsten Quadrate
3. Verteilungen der Maximum *Likelihood* Schätzer
4. Hypothesentests für Regressionsparameter
 - *t*-Test für Lineare Regression
5. Analyse der Residuen
6. Exkurs: Multiple Regression

Motivation: Lineare Regression

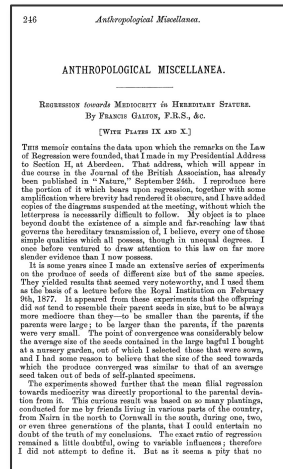
- **Beispiel (Vererbung von Größen).** Im Jahre 1886 untersuchte Sir Francis Galton, ob es einen genetischen Zusammenhang zwischen den Größen von Eltern und Kindern gibt. Dazu erstellte er einen Datensatz aus der Körpergröße von 930 erwachsenen Kindern und deren 205 Eltern. Karl Pearson wiederholte das Experiment mit 10 Vätern und deren erwachsenen Söhnen. Unterstützen die Daten die Hypothese, dass Kinder **im Mittel** kleiner sind als ihre Eltern?

- **Ansatz:**

- Wir plotten die 10 Datenpunkte (x = Größe des Vaters in Inch; y = Größe des Sohnes in Inch).
- Wir versuchen eine Gerade durch die Datenpunkte zu legen.
- Wir testen, ob die Steigung der Gerade kleiner als 1 ist.



Sir Francis Galton
(1822 – 1911)



Konzepte der Linearen Regression

- **Abhängige (*response*) Variable:** Y_i
 - **Beispiel (Vererbung von Größen).** Y_i = Größe des i -ten Kindes
- **Unabhängige (*input*) Variable:** x_i
 - Keine Zufälligkeit, da diese Variablen im Experiment ausgewählt/gesetzt werden
 - **Beispiel (Vererbung von Größen).** x_i = Größe des i -ten Elternteils
- **Model:** Zusammenhang zwischen abhängigen und unabhängigen Variablen
 - Die abhängige Variable ist die Zufallsvariable.
 - Lineares Model mit einem Fehler $Z_i \sim P$ der einer bekannten Verteilung P folgt

$$Y_i = \beta_1 \cdot x_i + \beta_0 + Z_i$$

Steigung (*slope*) → ← Schnitt mit der y -Achse (*intercept*)

- **Bemerkungen (Lineare Regression)**

- Wir nehmen an, dass $E[Z_i] = 0$ (ansonsten könnten wir $Z_i' = Z_i - E[Z_i]$ benutzen und $E[Z_i]$ zu β_0 addieren).
- Dieses Modell wird **lineare Regression** genannt, weil die Form der Funktion $E[Y_i] = \beta_1 \cdot x_i + \beta_0$ linear ist.

Mathe III

Unit 12a –
Regression

1. Konzept der Linearen Regression
2. **Maximum *Likelihood* Lineare Regression**
 - Methode der Kleinsten Quadrate
3. Verteilungen der Maximum *Likelihood* Schätzer
4. Hypothesentests für Regressionsparameter
 - *t*-Test für Lineare Regression
5. Analyse der Residuen
6. Exkurs: Multiple Regression

Lineare Regression mit normalverteiltem Fehler

- **Annahme:** Die Fehler sind unabhängig normalverteilt

$$Z_i \sim \mathcal{N}(0, \sigma^2)$$

- **Wiederholung (Likelihood).** Ist $(\Omega, \mathcal{F}, \{P_\theta \mid \theta \in \Theta\})$ ein parametrisches Modell und sei $y \in \Omega$ eine Stichprobe, so heißt die Funktion $\mathcal{L}: \Theta \times \Omega \rightarrow [0, +\infty)$ mit $\mathcal{L}(\theta, y) = p_\theta(y)$ die zugehörige *Likelihood*, wobei p_θ die (Zähl)dichte von P_θ ist.
- **Frage:** Was sind die Parameter der Regression und wie ist die *Likelihood* definiert?

1. Die Parameter sind

$$\theta = \begin{bmatrix} \beta_0 \\ \beta_1 \\ \sigma^2 \end{bmatrix}$$

← Schnitt mit der y -Achse (*intercept*)
← Steigung (*slope*)
← Fehlervarianz

2. Da die Fehler unabhängig verteilt sind, gilt, dass

$$\begin{aligned} \log \mathcal{L}((\beta_0, \beta_1, \sigma^2), \{y_1, \dots, y_n\}) &= \sum_{i=1}^n \log(\mathcal{N}(y_i; \beta_1 \cdot x_i + \beta_0, \sigma^2)) \\ &= \sum_{i=1}^n -\frac{1}{2} \cdot \left[\log(2\pi\sigma^2) + \frac{(y_i - (\beta_1 \cdot x_i + \beta_0))^2}{\sigma^2} \right] \end{aligned}$$

Mathe III

Unit 12a –
Regression

Maximum *Likelihood* Lineare Regression

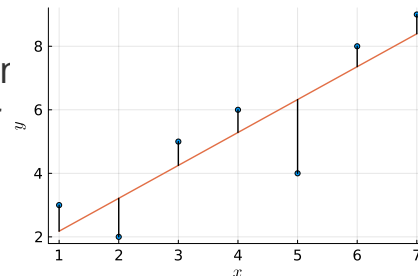
- **Satz (Maximum *Likelihood* Lineare Regression).** Gegeben ein lineares Regressionsmodell $Y = \beta_1 \cdot x + \beta_0 + Z$ mit $Z \sim \mathcal{N}(0, \sigma^2)$ und ein Datensatz von n Beobachtungen (x_i, y_i) . Dann entspricht der Maximum *Likelihood* Schätzer $\beta_{0,\text{MLE}}$ und $\beta_{1,\text{MLE}}$ von β_0 und β_1 dem Minimierer der kleinsten Quadrate

$$(\beta_{0,\text{MLE}}, \beta_{1,\text{MLE}}) = \arg \min_{\beta_0 \in \mathbb{R}, \beta_1 \in \mathbb{R}} \sum_{i=1}^n (y_i - (\beta_1 \cdot x_i + \beta_0))^2$$

- **Beweis:** Wenn wir alle Ausdrücke, die nicht von β_0 und β_1 abhängen, ignorieren, ergibt sich

$$\log \mathcal{L}((\beta_0, \beta_1, \sigma^2), \{y_1, \dots, y_n\}) = \sum_{i=1}^n -\frac{1}{2} \cdot \left[\log(2\pi\sigma^2) + \frac{(y_i - (\beta_1 \cdot x_i + \beta_0))^2}{\sigma^2} \right]$$

$$\Rightarrow -\frac{1}{2\sigma^2} \cdot \sum_{i=1}^n (y_i - (\beta_1 \cdot x_i + \beta_0))^2$$



Mathe III

Unit 12a –
Regression

Da der Vorfaktor strikt negativ ist, ist der Minimierer von die Maximum *Likelihood* Lösung.

1. Konzept der Linearen Regression
2. Maximum *Likelihood* Lineare Regression
 - **Methode der Kleinsten Quadrate**
3. Verteilungen der Maximum *Likelihood* Schätzer
4. Hypothesentests für Regressionsparameter
 - *t*-Test für Lineare Regression
5. Analyse der Residuen
6. Exkurs: Multiple Regression

Methode der Kleinsten Quadrate

- **Satz (Methode der Kleinsten Quadrate).** Gegeben ein Datensatz (x, y) von n Beobachtungen (x_i, y_i) , sind die Minimierer $\widehat{\beta}_0$ und $\widehat{\beta}_1$ der Summe der quadratischen Abstände $\sum_{i=1}^n (y_i - (\beta_1 \cdot x_i + \beta_0))^2$ explizit darstellbar durch

$$\widehat{\beta}_1 = \frac{\sum_{i=1}^n x_i \cdot y_i - n \cdot \bar{x} \cdot \bar{y}}{\sum_{i=1}^n x_i^2 - n \cdot (\bar{x})^2}$$

$$\widehat{\beta}_0 = \bar{y} - \widehat{\beta}_1 \cdot \bar{x}$$

- **Bemerkung (Methode der Kleinsten Quadrate).** In \mathbb{R} gibt es die Funktion `lm` um diese Berechnung direkt auszuführen

```
x = c(1, 2, 3, 4, 5, 6, 7)
y = c(3, 2, 5, 6, 4, 8, 9)
fit = lm(y ~ x)
fit
```

```
Coefficients:
(Intercept)      1.143
x              1.036
```

Mathe III

Unit 12a –
Regression

Beweis: Methode der Kleinsten Quadrate

- **Beweis:** Wir betrachten die erste Ableitung von $\sum_{i=1}^n (y_i - (\beta_1 \cdot x_i + \beta_0))^2$ nach β_0 und β_1

$$\frac{d \sum_{i=1}^n (y_i - (\beta_1 \cdot x_i + \beta_0))^2}{d\beta_0} = -2 \sum_{i=1}^n (y_i - (\beta_1 \cdot x_i + \beta_0))$$

$$\frac{d \sum_{i=1}^n (y_i - (\beta_1 \cdot x_i + \beta_0))^2}{d\beta_1} = -2 \sum_{i=1}^n x_i \cdot (y_i - (\beta_1 \cdot x_i + \beta_0))$$

Wenn wir die erste Ableitung nach β_0 Null setzen, ergibt sich

$$0 = \sum_{i=1}^n y_i - \widehat{\beta}_1 \cdot \sum_{i=1}^n x_i - n \cdot \widehat{\beta}_0 \Leftrightarrow \widehat{\beta}_0 = \bar{y} - \widehat{\beta}_1 \cdot \bar{x}$$

Wenn wir die erste Ableitung nach β_1 zu Null setzen, ergibt sich

$$0 = \sum_{i=1}^n x_i \cdot y_i - \widehat{\beta}_1 \cdot \sum_{i=1}^n x_i^2 - \widehat{\beta}_0 \cdot n \cdot \bar{x} \Leftrightarrow \widehat{\beta}_1 \cdot \left(\sum_{i=1}^n x_i^2 - n(\bar{x})^2 \right) + n \cdot \bar{y} \cdot \bar{x} = \sum_{i=1}^n x_i \cdot y_i$$

$$\Leftrightarrow \widehat{\beta}_1 = \left(\sum_{i=1}^n x_i \cdot y_i - n \cdot \bar{y} \cdot \bar{x} \right) / \left(\sum_{i=1}^n x_i^2 - n(\bar{x})^2 \right)$$



Adrien-Marie Legendre
(1752 - 1833)



Carl Friedrich Gauss
(1777 - 1855)

Mathe III

Unit 12a -
Regression

Motivation: Lineare Regression

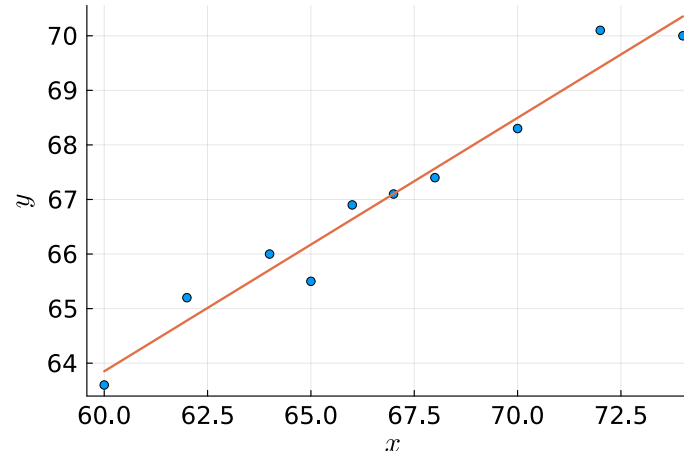
- **Beispiel (Vererbung von Größen).** Im Jahre 1886 untersuchte Sir Francis Galton, ob es einen genetischen Zusammenhang zwischen den Größen von Eltern und Kindern gibt. Dazu erstellte er einen Datensatz aus der Körpergröße von 930 erwachsenen Kindern und deren 205 Eltern. Karl Pearson wiederholte das Experiment mit 10 Vätern und deren erwachsenen Söhnen. Unterstützen die Daten die Hypothese, dass Kinder **im Mittel** kleiner sind als ihre Eltern?

- **Maximum Likelihood Regression:**

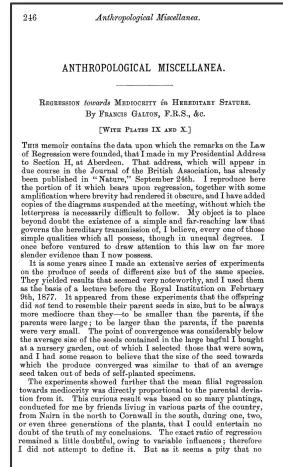
```
lm(formula = y ~ x)
```

Coefficients:
(Intercept) 35.9768 x 0.4646

Initiale Bestätigung für $\hat{\beta}_1 < 1$.
Aber: Wie sicher können wir uns aufgrund der Stichprobe sein?



Sir Francis Galton
(1822 – 1911)



1. Konzept der Linearen Regression
2. Maximum *Likelihood* Lineare Regression
 - Methode der Kleinsten Quadrate
3. **Verteilungen der Maximum *Likelihood* Schätzer**
4. Hypothesentests für Regressionsparameter
 - *t*-Test für Lineare Regression
5. Analyse der Residuen
6. Exkurs: Multiple Regression

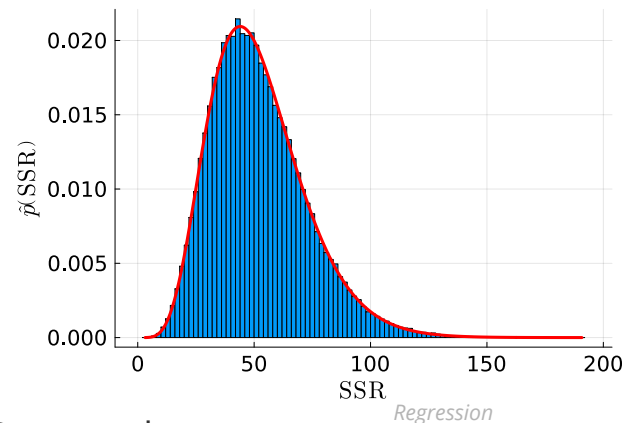
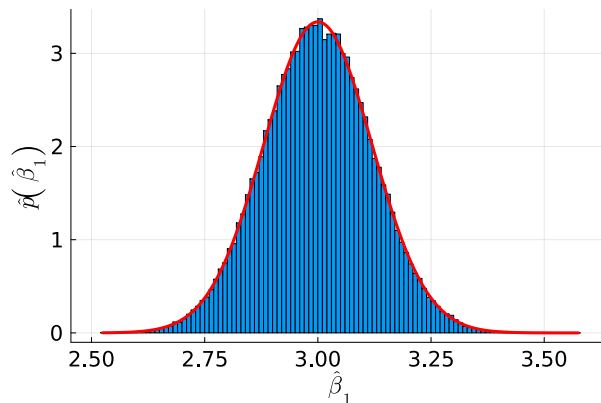
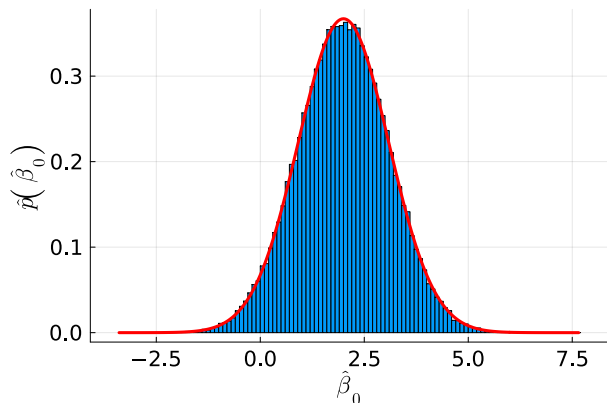
Motivation: Verteilung von Regressionsparametern

- **Beispiel (Verteilung von Regressionsparametern).** Wir simulieren die Verteilung von $\widehat{\beta}_0$ und $\widehat{\beta}_1$ indem wir 100,000-mal $n = 15$ Datenpunkte $x_i = i$ für $\beta_0 = 2$, $\beta_1 = 3$ und $Z_i \sim \mathcal{N}(0,4)$ simulieren. Dann ergibt sich

$$\widehat{\beta}_0 = \bar{y} - \widehat{\beta}_1 \cdot \bar{x}$$

$$\widehat{\beta}_1 = \frac{\sum_{i=1}^n x_i \cdot y_i - n \cdot \bar{x} \cdot \bar{y}}{\sum_{i=1}^n x_i^2 - n \cdot (\bar{x})^2}$$

$$SSR = \sum_{i=1}^n (y_i - (\widehat{\beta}_1 \cdot x_i + \widehat{\beta}_0))^2$$



- Es erscheint so, dass $\widehat{\beta}_0$ und $\widehat{\beta}_1$ auch normalverteilt sind und die Summe der quadratischen Residuen χ^2 -verteilt sind!
- Aber mit welchen Parametern?

Verteilung von $\widehat{\beta}_1$ (Erwartungswert)

- **Satz (Verteilung von $\widehat{\beta}_1$).** Im linearen Regressionsmodell $Y = \beta_1 \cdot x + \beta_0 + Z$ mit unabhängig normalverteiltem Fehler $Z \sim \mathcal{N}(0, \sigma^2)$ gilt für den MLE von β_1 für ein Stichprobe der Größe n

$$\widehat{\beta}_1 \sim \mathcal{N}\left(\beta_1, \frac{\sigma^2}{\sum_{i=1}^n x_i^2 - n(\bar{x})^2}\right)$$

- **Beweis für $E[\widehat{\beta}_1]$:** Wir benutzen die Linearität des Erwartungswertes

$$\begin{aligned} E[\widehat{\beta}_1] &= E_{Y_1, \dots, Y_n} \left[\frac{\sum_{i=1}^n x_i \cdot Y_i - \cancel{n \cdot \bar{Y} \cdot \bar{x}}}{\sum_{i=1}^n x_i^2 - n(\bar{x})^2} \right] \quad \text{---} = n \cdot \frac{1}{n} \sum_{i=1}^n Y_i \\ &= E_{Y_1, \dots, Y_n} \left[\frac{\sum_{i=1}^n (x_i - \bar{x}) \cdot Y_i}{\sum_{i=1}^n x_i^2 - n(\bar{x})^2} \right] \quad \text{--- da } E[Y_i] = E[\beta_1 \cdot x_i + \beta_0 + Z] \text{ und } E[Z] = 0 \\ &= \frac{\sum_{i=1}^n (x_i - \bar{x}) \cdot (\beta_1 x_i + \beta_0)}{\sum_{i=1}^n x_i^2 - n(\bar{x})^2} \\ &= \frac{\beta_1 \cdot \sum_{i=1}^n x_i \cdot (x_i - \bar{x}) + \beta_0 \cdot \sum_{i=1}^n (x_i - \bar{x})}{\sum_{i=1}^n x_i^2 - n(\bar{x})^2} \quad \text{---} = 0, \text{ da } \sum_{i=1}^n x_i = n \cdot \bar{x} \text{ und } \sum_{i=1}^n \bar{x} = n \cdot \bar{x} \\ &= \beta_1 \cdot \frac{\sum_{i=1}^n x_i^2 - n(\bar{x})^2}{\sum_{i=1}^n x_i^2 - n(\bar{x})^2} \end{aligned}$$

Mathe III

Unit 12a –
Regression

Verteilung von $\widehat{\beta}_1$ (Varianz)

- **Satz (Verteilung von $\widehat{\beta}_1$).** Im linearen Regressionsmodell $Y = \beta_1 \cdot x + \beta_0 + Z$ mit unabhängig normalverteiltem Fehler $Z \sim \mathcal{N}(0, \sigma^2)$ gilt für den MLE von β_1

$$\widehat{\beta}_1 \sim \mathcal{N}\left(\beta_1, \frac{\sigma^2}{\sum_{i=1}^n x_i^2 - n(\bar{x})^2}\right)$$

- **Beweis für $V[\widehat{\beta}_1]$:** Wir benutzen die Unabhängigkeit der Fehler

$$V[\widehat{\beta}_1] = V_{Y_1, \dots, Y_n} \left[\frac{\sum_{i=1}^n (x_i - \bar{x}) \cdot Y_i}{\sum_{i=1}^n x_i^2 - n(\bar{x})^2} \right]$$

Wegen
Varianzzerlegungssatz

$$= \frac{\sum_{i=1}^n (x_i - \bar{x})^2 \cdot V[Y_i]}{(\sum_{i=1}^n x_i^2 - n(\bar{x})^2)^2}$$

wegen Unabhängigkeit der
Fehler

$$= \frac{\sum_{i=1}^n x_i^2 - n(\bar{x})^2}{(\sum_{i=1}^n x_i^2 - n(\bar{x})^2)^2} \cdot \sigma^2$$

da $V[Y_i] = V[\beta_1 \cdot x_i + \beta_0 + Z] = V[Z]$

$$= \frac{\sigma^2}{\sum_{i=1}^n x_i^2 - n(\bar{x})^2}$$

Mathe III

Unit 12a –
Regression

Verteilung von $\widehat{\beta}_0$ (Erwartungswert)

- **Satz (Verteilung von $\widehat{\beta}_0$).** Im linearen Regressionsmodell $Y = \beta_1 \cdot x + \beta_0 + Z$ mit unabhängig normalverteiltem Fehler $Z \sim \mathcal{N}(0, \sigma^2)$ gilt für den MLE von β_0 für ein Stichprobe der Größe n

$$\widehat{\beta}_0 \sim \mathcal{N}\left(\beta_0, \frac{\sigma^2 \cdot \sum_{i=1}^n x_i^2}{n \cdot (\sum_{i=1}^n x_i^2 - n(\bar{x})^2)}\right)$$

- **Beweis für $E[\widehat{\beta}_0]$:** Wir benutzen die Linearität des Erwartungswertes und den vorherigen Satz

$$\begin{aligned} E[\widehat{\beta}_0] &= E_{Y_1, \dots, Y_n}[\bar{Y} - \widehat{\beta}_1 \bar{x}] = \beta_1 \\ &= \frac{1}{n} \sum_{i=1}^n E[Y_i] - \beta_1 \cdot \bar{x} \\ &= \frac{1}{n} \sum_{i=1}^n (\beta_1 \cdot x_i + \beta_0) - \beta_1 \cdot \bar{x} \\ &= \beta_0 \end{aligned}$$

$= \beta_1 \cdot x_i + \beta_0$, da $E[Y_i] = E[\beta_1 \cdot x_i + \beta_0 + Z]$ und $E[Z] = 0$
 $= \beta_1 \cdot \bar{x}$

Mathe III

Unit 12a –
Regression

Verteilung von $\widehat{\beta}_0$

- **Satz (Verteilung von $\widehat{\beta}_0$).** Im linearen Regressionsmodell $Y = \beta_1 \cdot x + \beta_0 + Z$ mit unabhängig normalverteiltem Fehler $Z \sim \mathcal{N}(0, \sigma^2)$ gilt für den MLE von β_0 für ein Stichprobe der Größe n

$$\widehat{\beta}_0 \sim \mathcal{N}\left(\beta_0, \frac{\sigma^2 \cdot \sum_{i=1}^n x_i^2}{n \cdot (\sum_{i=1}^n x_i^2 - n(\bar{x})^2)}\right)$$

- **Beweis für $V[\widehat{\beta}_0]$:** Wir benutzen die Unabhängigkeit der Fehler und den vorherigen Satz

$$\begin{aligned}
 V[\widehat{\beta}_0] &= V_{Y_1, \dots, Y_n}[\bar{Y} - \widehat{\beta}_1 \cdot \bar{x}] = \frac{1}{n} \sum_{i=1}^n Y_i = \frac{\sigma^2}{\sum_{i=1}^n x_i^2 - n(\bar{x})^2} \\
 &= V_{Y_1, \dots, Y_n} \left[\frac{1}{n} \sum_{i=1}^n Y_i \right] + (\bar{x})^2 \cdot V_{Y_1, \dots, Y_n}[\widehat{\beta}_1] + 2 \cdot \text{Cov} \left[\frac{1}{n} \sum_{i=1}^n Y_i, \widehat{\beta}_1 \cdot \bar{x} \right] \\
 &= \frac{\sigma^2}{n} + (\bar{x})^2 \cdot \frac{\sigma^2}{\sum_{i=1}^n x_i^2 - n(\bar{x})^2} \\
 &= \frac{\sigma^2 \cdot (\sum_{i=1}^n x_i^2 - n(\bar{x})^2) + \sigma^2 \cdot n(\bar{x})^2}{n \cdot (\sum_{i=1}^n x_i^2 - n(\bar{x})^2)} = \frac{\sigma^2 \cdot \sum_{i=1}^n x_i^2}{n \cdot (\sum_{i=1}^n x_i^2 - n(\bar{x})^2)}
 \end{aligned}$$

wegen Unabhängigkeit der Fehler

Mathe III

Unit 12a - Regression

Verteilung von SSR

- **Satz (Verteilung von SSR).** Im linearen Regressionsmodell $Y = \beta_1 \cdot x + \beta_0 + Z$ mit unabhängig normalverteiltem Fehler $Z \sim \mathcal{N}(0, \sigma^2)$ gilt für die Verteilung der Summe der quadratischen Fehler bei einer Stichprobe der Größe n

$$\frac{SSR}{\sigma^2} \sim \chi^2(n - 2)$$

- **Bemerkungen (Verteilung von Regressionsparametern)**

- Intuitiv ist die Anzahl der Freiheitsgrade $n - 2$ weil SSR , $\widehat{\beta}_0$, und $\widehat{\beta}_1$ unabhängig sind und die beiden Schätzer $\widehat{\beta}_0$ und $\widehat{\beta}_1$ jeweils einen Freiheitsgrad „entfernen“.
- Die Verteilung der Schätzer erlaubt ein Konfidenzintervall für $\widehat{\beta}_0$ und $\widehat{\beta}_1$ anzugeben, wenn wir die Varianz der Fehler als bekannt annehmen.
- Wenn wir ein zweiseitiges Konfidenzintervall annehmen, dann gilt zum Konfidenzniveau $1 - \alpha$, dass

$$\beta_1 \in \left[\widehat{\beta}_1 + z_{\frac{\alpha}{2}} \cdot \frac{\sigma}{\sqrt{\sum_{i=1}^n x_i^2 - n(\bar{x})^2}}, \widehat{\beta}_1 + z_{1-\frac{\alpha}{2}} \cdot \frac{\sigma}{\sqrt{\sum_{i=1}^n x_i^2 - n(\bar{x})^2}} \right]$$

- Der Begriff „Regression“ geht auf Sir Francis Galton zurück!

Mathe III

Unit 12a –
Regression

Viel Spaß bis zur nächsten Vorlesung!