

Einfache Regression

Lernziele

- Lineares Regressionsmodell
- Anpassen des linearen Regressionsmodells, OLS
- Eigenschaften der Schätzer für das Modell
- Testen und Konfidenzintervalle
- Überprüfung der Modellannahmen
- Prognose mittels linearer Regression
- Multiples Regressionsmodell

Das lineare Regressionsmodell

Das einfache lineare Regressionsmodell lautet

$$y_t = \alpha + \beta x_t + \varepsilon_t \quad t = 1, \dots, n$$

y_t ... abhängige Variable

α ... Interzept, Konstante

x_t ... unabhängige Variable

β ... Steigung

ε_t ... Fehler, Residuum

Modellvoraussetzung

Annahmen für das einfache lineare Regressionsmodell:

- Die Beziehung zwischen y und x ist linear.
- Die x_t sind deterministische (nicht-stochastische) Variablen. Ihre Werte sind genau bekannt.
- Die Residuen ε_t sind

- unkorreliert,

$$\text{Cov}(\varepsilon_{t_i}, \varepsilon_{t_j}) = 0 \quad \text{für } i \neq j$$

- haben Erwartungswert 0 und konstante Varianz,

$$E(\varepsilon_t) = 0, \quad V(\varepsilon_t) = \sigma_\varepsilon^2, \quad \text{für alle } t$$

- und sind normalverteilt,

$$\varepsilon_t \sim N(0, \sigma_\varepsilon^2)$$

Der bedingte Erwartungswert

Der bedingte Erwartungswert $E(y_t|x_t)$ von y_t gegen x_t ist

$$E(y_t|x_t) = \alpha + \beta x_t$$

da $E(\varepsilon_t) = 0$ und x_t bekannt ist.

Das Mittel von y_t wird als Funktion von x_t angesehen.

- $E(y_t)$ bezeichnet das univariate (**unbedingte**) Mittel von y_t , i.e., der Mittelwert aller y_t über die Zeit wenn die Information x_t nicht verwendet/bekannt ist.
- $E(y_t|x_t)$ ist eine **bedingte Erwartung**. Die Information x_t wird zur Erklärung/Modellierung von y_t verwendet.

Die bedingte Varianz

- Die **unbedingte Varianz** ergibt sich als

$$V(y_t) = E([y_t - E(y_t)]^2)$$

- Die **bedingte Varianz** $V(y_t|x_t)$ ist analog

$$V(y_t|x_t) = E\left((y_t - E(y_t|x_t))^2|x_t\right) = \sigma_\varepsilon^2$$

In unserem Modell ist die bedingte Varianz von y_t konstant:

$$y_t - E(y_t|x_t) = (\alpha + \beta x_t + \varepsilon_t) - (\alpha + \beta x_t) = \varepsilon_t^2$$

Und somit

$$V(y_t|x_t) = E(\varepsilon_t^2|x_t) = E(\varepsilon_t^2) = \sigma_\varepsilon^2$$

Varianzreduktion

Zur Berechnung von $V(y_t|x_t)$ wird die Information x_t verwendet.
Daher gilt immer (sofern das Modell einen Erklärungswert hat)

$$V(y_t) > V(y_t|x_t) = \sigma_\varepsilon^2$$

Die Regression reduziert die Varianz in den Daten.
Sie erklärt somit einen Teil der Variation in y_t .

Anpassen eines Regressionsmodells

Die Parameter des Regressionsmodells werden so gewählt, daß die **Fehlerquadratsumme der Residuen minimal** wird.

$$\min_{\alpha, \beta} \sum_{t=1}^n \varepsilon_t^2 = \min_{\alpha, \beta} \sum_{t=1}^n (y_t - \hat{y}_t)^2 = \min_{\alpha, \beta} \sum_{t=1}^n (y_t - (\alpha + \beta x_t))^2$$

Das Minimierungsproblem wird mit der Methode der kleinsten Quadrate (*ordinary least squares*, **OLS**) gelöst.

Man berechnet dazu die ersten partiellen Ableitungen und setzt diese gleich Null.

$$\hat{\alpha} = \frac{1}{n} \sum y_t - \hat{\beta} \frac{1}{n} \sum x_t \quad \hat{\beta} = \frac{n \sum x_t y_t - \sum x_t \sum y_t}{n \sum x_t^2 - (\sum x_t)^2}$$

OLS-Schätzer

Man kann diese Schätzer auch anders darstellen:

$$\hat{\alpha} = \bar{y} - \hat{\beta} \bar{x} \quad \hat{\beta} = \frac{\text{Cov}(x, y)}{V(x)} = \text{Corr}(x, y) \frac{s_y}{s_x}$$

Der Schätzer für β wird Null, falls x_t und y_t unkorreliert sind. Dann wird die Variable x_t nicht in die Regressionsgleichung aufgenommen.

Nach der Schätzung lautet unser Modell:

$$y_t = \hat{\alpha} + \hat{\beta} x_t + \hat{\varepsilon}_t$$

$\hat{\varepsilon}_t$ sind die geschätzten Fehler.

Verteilung der Schätzer

Unter Annahme (Nullhypothese), daß unser Modell das *wahre Modell* ist, sind die OLS-Schätzer $\hat{\alpha}$ und $\hat{\beta}$ normalverteilt mit Erwartungswert α bzw. β :

$$\hat{\alpha} \sim N(\alpha, \sigma_{\hat{\alpha}}^2) \quad \hat{\beta} \sim N(\beta, \sigma_{\hat{\beta}}^2)$$

mit

$$\sigma_{\hat{\alpha}}^2 = \sigma_{\varepsilon}^2 \frac{\sum x_t^2}{n \sum (x_t - \bar{x})^2} = \sigma_{\varepsilon}^2 \left(\frac{1}{n} + \frac{\bar{x}^2}{\sum (x_t - \bar{x})^2} \right)$$
$$\sigma_{\hat{\beta}}^2 = \sigma_{\varepsilon}^2 \frac{1}{\sum (x_t - \bar{x})^2}$$

Verteilung der Schätzer / (2)

σ_{ε}^2 wird durch die **Stichprobenvarianz des Fehlers** s^2 ersetzt:

$$\hat{\sigma}_{\varepsilon}^2 = s^2 = \frac{1}{n-2} \sum \hat{\varepsilon}_t^2, \quad \text{mit } \hat{\varepsilon}_t = y_t - (\hat{\alpha} + \hat{\beta} x_t)$$

Damit wird erhalten wir die Schätzer $\hat{\sigma}_{\hat{\alpha}}^2$ und $\hat{\sigma}_{\hat{\beta}}^2$ für $\sigma_{\hat{\alpha}}^2$ bzw. $\sigma_{\hat{\beta}}^2$.

Eigenschaften der OLS-Schätzer

Die OLS-Lösung liefert den **BLUE**, den besten linearen unverzerrten Schätzer (*best linear unbiased estimator*):

- **Best**: der Schätzer hat die kleinste Varianz.
- **Linear**: das Modell ist linear.
- **Unbiased** oder **unverzerrt**: Erwartungswert des Schätzers ist der Wert des Parameters ist. (Ansonsten wäre er verzerrt.)
- **Estimator**: Schätzer

***t*-Test der geschätzten Parameter**

Wir können obige Verteilung benutzen um die Hypothese

$$H_0: \beta = b \quad H_1: \beta \neq b$$

zu testen. Die Prüfgröße (*Teststatistik*)

$$\frac{\hat{\beta} - b}{\hat{\sigma}_{\beta}} \sim t_{n-2}$$

ist dabei *t*-verteilt mit $n - 2$ Freiheitsgraden.

(Analog für α .)

Konfidenzintervalle

Das $(1 - \alpha)$ -Konfidenzintervall für β ergibt sich aus

$$[\hat{\beta} - t_{n-2, 1-\alpha/2} \hat{\sigma}_{\beta}, \hat{\beta} + t_{n-2, 1-\alpha/2} \hat{\sigma}_{\beta}]$$

Das Bestimmtheitsmaß

Wir wollen messen, wie gut sich das Modell an die Daten anpasst und ob unsere Schätzung die Annahmen des Regressionsmodells erfüllt.

Das **Bestimmtheitsmaß** R^2 ist definiert durch als der Quotient aus der Summe der Abweichungsquadrate im geschätzten Modell und der Summe der Abweichungsquadrate der beobachteten Daten.

$$R^2 = \frac{\sum(\hat{y}_t - \bar{y})^2}{\sum(y_t - \bar{y})^2}$$

R^2 gibt den Anteil der durch die Regression erklärten Varianz an.

Das Bestimmtheitsmaß / (2)

$$\begin{array}{ccccc} (y_t - \bar{y}) & = & (y_t - \hat{y}_t) & + & (\hat{y}_t - \bar{y}) \\ \uparrow & & \uparrow & & \uparrow \\ \text{totale} & & \text{nicht-erklärte} & & \text{erklärte} \\ \text{Abweichung} & & \text{Abweichung} & & \text{Abweichung} \end{array}$$

Quadrieren und Summieren ergibt

$$\begin{array}{ccccc} \sum (y_t - \bar{y})^2 & = & \sum (y_t - \hat{y}_t)^2 & + & \sum (\hat{y}_t - \bar{y})^2 \\ \uparrow & & \uparrow & & \uparrow \\ \text{SST} & & \text{SSE} & & \text{SSR} \end{array}$$

SST ... *total sum of squares*

SSE ... *sum of squared errors*

SSR ... *sum of squares from regression*

(Der Term $2 \sum (y_t - \hat{y}_t)(\hat{y}_t - \bar{y})$ ist Null.)

F-test

Das Bestimmtheitsmaß R^2 kann ebenfalls getestet werden.

Der Test $H_0: R^2 = 0$ ist dabei äquivalent dem Test $H_0: \beta = 0$.

Die Prüfgröße lautet

$$\frac{R^2}{1 - R^2} (n - 2) \sim F_{1, n-2}$$

und ist F -verteilt mit 1 und $n - 2$ Freiheitsgraden.

Der Test heißt kurz **F-Test**.

Überprüfung der Modellannahmen

- **Unkorreliertheit der Residuen.** Aus der Autokorrelation könnte man die Residuen prognostizieren. Damit wäre unser Modell nicht vollständig.
- **Konstanz der Varianz der Residuen.** Wenn die Varianzen in verschiedenen Teilperioden unterschiedlich sind, sind die geschätzten Varianzen von $\hat{\alpha}$ und $\hat{\beta}$ nicht korrekt.
- **Abweichungen von der Normalverteilung** kann mittels Histogrammen und Test auf Normalverteilung entdeckt werden.

Prognose mittels linearer Regression

Zur Prognose müssen die Werte der unabhängigen Variable $x_{t+\tau}$, $\tau = 1, \dots, k$ bekannt sein.

- **Prognosefunktion**

$$\hat{y}_{t+\tau} = \hat{\alpha} + \hat{\beta} x_{t+\tau}$$

- **Standardfehler der Prognose**

$$\text{se}(\hat{y}_{t+\tau}) = \hat{\sigma}_\varepsilon \sqrt{1 + \frac{1}{t} + \frac{(x_{t+\tau} - \bar{x})^2}{\sum_{i=1}^t (x_i - \bar{x})^2}}$$

Der Standardfehler der Prognose ist minimal, wenn $x_{t+\tau} = \bar{x}$ ist.

- **Prognoseintervall**

$$[\hat{y}_{t+\tau} - t_{t-2} \cdot \text{se}(\hat{y}_{t+\tau}), \hat{y}_{t+\tau} + t_{t-2} \cdot \text{se}(\hat{y}_{t+\tau})]$$

Multiples Regressionsmodell

Wir erweitern das Modell auf mehrere erklärende Variable

$$y_t = \alpha + \beta_1 x_{1t} + \beta_2 x_{2t} + \dots + \beta_m x_{mt} + \varepsilon_t$$

Die Koeffizienten $\alpha, \beta_1, \beta_2, \dots, \beta_m$ werden wieder mittels OLS geschätzt.

Es gibt eine Reihe von Verfahren um die Anzahl m der notwendigen Variablen zu bestimmen (so wenig wie möglich!).

Das multiple Bestimmtheitsmaß

Das **multiple Bestimmtheitsmaß** ist analog zum einfachen Bestimmtheitsmaß definiert.

Liegt eine Regression mit m unabhängigen Variablen, x_{jt} , mit $j = 1, \dots, m$ vor, so ist der Test auf $R^2 = 0$ äquivalent dem Test, dass alle $\beta_j = 0$ sind (bzw. dass kein x_j einen Beitrag zur Varianzreduktion in y zu leisten im Stande ist).

F-Test:

$$H_0: \beta_1 = \beta_2 = \dots = \beta_m = 0 \quad (R^2 = 0)$$

$$H_1: \text{zumindest ein } \beta_j \neq 0 \quad (R^2 \neq 0)$$

Die Prüfgröße

$$\frac{R^2}{1 - R^2} \frac{n - (m + 1)}{m} \sim F_{m, n - (m + 1)}$$

ist F -verteilt mit m und $n - (m + 1)$ Freiheitsgraden.