

## Einfaches Lineares Regressionsmodell.

Grob:  $y$ -Werte (mithilfe von  $x$ -Werten) bestmöglich im Mittel beschreiben  
Interessierende Regs. Koeffizienten schätzen, wie sich der durchschnittliche Wert von  $y$  bei  $\Delta x$  verändert.  
Unter

Wollen Zusammenhang zwischen Zielgröße  $y$  und Einflussgrößen  $x$  modellieren. Durchschnittliche Wert von  $y$  bei  $\Delta x$  verändert.

Eine problematische Annahme: Wahres störiges Modell (einzigen Struktur, also linear mit bestimmten Einflussvariablen) ist bereits bekannt.  
Dies impliziert, dass wir datengenerierenden Prozess kennen! Selbst mit Fachwissen über den jeweiligen Bereich ist das niemals korrekt.

Reality: All models are wrong, but some are useful (or at least more useful than other models)

Ziel: Zusammenhang beschreiben, Effekt ( $\beta$ -Koeff.) schätzen und interpretieren, Prognostizieren / Prädiktion von neuen zukünftigen Daten

Zusammenhänge erkennen, Trends finden, Supervised Learning (Zusammenhänge zwischen Variablen aus Daten lernen)

Modell in Vektor-Schreibweise:  $i = 1, \dots, n$  Beobachtungen

Jede einzelne Beobachtung modelliert als ZV

$$y_i = \beta_0 + \beta_1 x_i + \varepsilon_i$$

unbekannte Parameter

Zufallsfehler/ Error-Term  
 $\varepsilon_i$ : bayesianisch: Alles, was wir nicht kennen  
also  $\varepsilon_i$  aus Sichtweise repräsentieren: Zufällig (in Realität liegen oft "convenient" samples aus verstecktem Beobachtd. Population)

Richtung! Nur Zielgröße  $y_{ij}$  sei Zufallsvariable. Indirekt steckt Zufall also nur in den Error-Terminen  $\varepsilon_i$ .

Willkürliche Annahme

$x_{ij}$  ist demnach **feste** bekannte Einflussvariable, wird als deterministisch betrachtet [Messfehler damit auch ausgeschlossen.]

In Realität kommen  $x$ -Werten jedoch meistens aus einer (Zufalls-) Stichprobe, ergo erinnert auch  $x_{ij}$  eine Zufallskomponente.

Ausschließen bzw. weitgehend kontrollieren lässt sich Zufälligkeit von  $x$ -Werten durch kontrollierte Experimente.  $x_{ij}$  sind dann experimentell bis zu einem gewissen Grad festgelegt.

Randomisierung (zufällige Zuweisung von  $x$ -Werten innerhalb des Experiments, im Modellingsamt aber nicht zufällig!) ist hinsichtlich möglich.

$\hookrightarrow$  Kein Confounding bzw. Störfaktoren haben keine systematische Beziehung zu Einflussgröße, Störfaktoren "mischen sich weg".

Vorteil der Annahme fester  $x$ -Werte: Berechnung von OLS-Schätzern  $\hat{\beta}$  wird mathematisch einfacher

$\beta$ -Verteilung hängt auf  $X$  (konditionell auf  $X$ ) an, ohne Verteilung von  $X$  kann es zu müssen

Konditionierung auf  $X \Rightarrow$  Zufälligkeit von  $\beta$  beruht nur auf Zufälligkeit von  $\varepsilon$

Nehme man Regressions(e)n) als zufällig an, müsste man ihre Verteilung mitmodellieren,  $\beta$ -Verteilung wäre eine unbekannte. OLS-Schätzer wäre dann erwartungstreu und (asymptotisch) normalverteilt. Brucht als Effekt gegen  $x$ , sondern die strukturellen Parameter des generierten Modells für  $(X, Y)$

In vielen statistischen Modellen interessiert man sich für gesamte Unsicherheit. Gemeinsame Zufallsverteilung von  $(Y, X)$  ist wesentlich rechtfärtiger.

(Polynome, Reciprocal-1/k, log(x)...)

Linearmodell muss nur in **Parameter** getrennt! Man kann auch komplexe  $x$ -Termen mit linearen Parametern modellieren. Linearkombination der Kovariablen(n)

Weitere Annahmen des einfachen linearen Modells:  $y_i = f(x_i) + \varepsilon_i = E(y_i|x_i) + \varepsilon_i$ , also  $y_i = \beta_0 + \beta_1 x_i + \varepsilon_i$ : (1)

Bedachte:  $E(\varepsilon) = 0$ , aber  $\varepsilon \neq 0$  kann man nicht ausschließen!

Unrichtig: Identisch verteilte Beobachtungen  $(y_i, x_i)$  i.i.d., also auch Fehlerterme  $\varepsilon_i$  sind i.i.d.

$E(\varepsilon_i^2) = \text{Var}(\varepsilon_i^2) = \sigma^2$

$E(\varepsilon_i) = 0$  (2)

$\text{Var}(\varepsilon_i) = \sigma^2$  konstant (3)  $\leftarrow$

$\hookrightarrow E(\varepsilon_i) = \frac{1}{n} \sum \varepsilon_i = 0$  per  $kQ$ -Jedogen mit  $\beta$

$E(\varepsilon_i) = 0$  Bei zufälligen Regressoren wäre es strenger:

$E(\varepsilon_i | x_i) = 0$  trivial erfüllt bei festem  $x$

$\hookrightarrow$  Existenz der Einflussrichtung

$\Rightarrow$  Varianz der Residuen = Null (oder per  $kQ$ -Design!)

auch gilt  $\bar{\varepsilon} = \bar{\eta} = \bar{\eta}$  dann folgerichtig übernommt Intercept  $\beta_0$ :

Wie wichtig diese Annahme sind (insbes. (3) bis (5)) hängt stark davon ab, was wir mit dem linearen Modell bezeichnen, und manches lässt sich auch korrigieren.

im mathematischen Modell:  $Statt y = \beta_0 + \beta_1 x + yz + u$  schätzen wir  $y = \beta_0 + \beta_1 x + \varepsilon$

Confounding Variables: Variablen, die nicht im Modell enthalten sind, mit einer oder mehreren  $x$ -Variablen zusammenhängen und zuden.

Einfluss auf Zielgröße  $y$  haben (nehmen wir sie  $z$ )

$\Rightarrow$  (implizite bei festem  $x$ ) Eigentümlichkeitsannahme:  $E(\varepsilon_i | x_i) = 0$  ist verletzt,  $x$  korreliert mit Störfaktoren und damit folgt  $\text{cov}(x_i, \varepsilon_i) \neq 0$   $\text{cov}'(x_i) \neq 0$

Teil des (den man aus  $x$  nicht perfekt herausgenommen,  $\xrightarrow{\text{Restauflösung Projektion}}$ )  $\text{Var}(y - \text{Proj}(y|X)) = 0$   $\xrightarrow{\text{Von } Z \text{ auf } X}$   $\text{Var}(y - \text{Proj}(y|X)) = 0$   $\xrightarrow{\text{Restauflösung Projektion}}$

→ Effekt der eingeschlossenen Variablen wird in den Errorterm geschrieben, der systematisch und  $(E(\varepsilon_i) + 0)$  für manche  $i$  konst.

OMITTED VARIABLE: BIAS: Bias verschwindet nur wenn  $\text{Var}(x_i) \cdot \text{Var}(y - \text{Proj}(y|X)) = 0$   $\xrightarrow{\text{Von } Z \text{ auf } X}$   $\text{Var}(y - \text{Proj}(y|X)) = 0$   $\xrightarrow{\text{Restauflösung Projektion}}$

Ein Teil des Effekts von  $z$  auf  $y$  wird durch eine/mehrere Regressoren ( $x$ ) aufgegessen,  $\beta_1$  wird nach oben verzerrt, da  $\text{Var}(\beta_1)$  wird aufgezerrt, also Ineffizienz (außer  $\text{Var}(\beta_1) = 0$ )

