

$$1. \eta = X\beta + \varepsilon \quad (\text{Strukturannahme - ohne geht garnichts})$$

(2) impliziert

- 2. $E(\varepsilon_i) = 0$
- 3. $\text{Var}(\varepsilon_i) = \sigma^2 \text{ const.}$

Zufalls-Stichprobe

- 4. ε_i prozeßweise unabhängig
- 5. $\varepsilon_i \sim N(0, \sigma^2)$

a.) Welche Annahmen 1.) - 5.) sind für die in den Zeilen angegebenen Schätzer, Eigenschaften bzw. Größen notwendig? **Multiples Regressionsmodell:** $\eta_{\text{gewohnt}, i} = \beta_0 + \beta_1 x_{1,i} + \beta_2 x_{2,i} + \dots + \beta_k x_{k,i} + \varepsilon_i$

	„wahrscheinlichste“ Gerade.	1.)	2.)	3.)	4.)	5.)
MLE weniger fehleranfällig, weil wir Annahmen (2) bis(5) gründlich vorher prüfen müssen	KQ-Schätzer	2.) Welche Gerade passt am besten zu Daten?	x	(2) bis(5) $\sqrt{n} \hat{\beta} \xrightarrow{P} \beta_{MLE}$		
	ML-Schätzer	2.) Welche Gerade passt am besten zu Daten und Annahmen (2) bis (5)?	x	x	x	x
	$E(\hat{\beta}) = \beta$		x	x		
	$\hat{\beta} \sim N(\beta, \sigma^2 (X'X)^{-1})$		x	x	x	(x)*
	Konfidenzintervalle (von Regressionskoeffizienten, Prädiktionen, ...)		x	x	x	(x)*
	Prädiktionsintervalle (von Prädiktionen, ...)		x	x	x	x

*Asymmetrisch ist jeder MLE-Schätzer neutraleswertet?? großer n kann falsche Verteilungssumme (ε_i) kompensieren

*Diese Eigenschaften gelten auch für ausreichend großen Stichprobenumfang n approximativ. Für kleine Stichprobenumfänge ist die Normalverteilungsannahme jedoch zwingend notwendig!

Hinweis: Besonders anzumerken ist, dass für den KQ-Schätzer nichts weiter benötigt wird als die Modellgleichung. Nimmt man zusätzlich noch an, dass der Erwartungswert der Residuen 0 ist, sind die Schätzer für β sogar erwartungstreu.

Um Konfidenz-Intervalle bilden zu können sind alle 5 Annahmen notwendig (oder 1.)-4.) und eine große Stichprobe). Für Prädiktions-Intervalle sind stets alle 5 Annahmen notwendig (insbesondere genügt hier auch eine ausreichend große Stichprobe NICHT).

b.) Unverzerrte Trendabschätzung (1-2 retten) vs. Unverzerrte Trend-Extrapolation (3-5 retten)

Es wird angenommen, dass sich die Zielgröße (Körpergewicht) als eine Linearkombination aus den Einflussgrößen (Alter, Körpergröße) beschreiben lässt - also z.B. konkret dass sich das Gewicht von Grundschulkindern direkt nach der Formel in der Angabe berechnen lässt (und nicht z.B. Gewicht = Alter Größe oder sonst einer anderen Formel).

Diese Annahme scheint als Laie beurteilt nicht vollkommen unplausibel und kann daher zunächst so angenommen werden.

Hinweise: Diese Annahme ist vor allem von Wissenschaftlern aus dem jeweiligen Fachgebiet zu beurteilen. Es gibt jedoch Methoden der Modelldiagnostik (Residuen-vs.-Kovariablen-Plot), welche noch später im Kurs behandelt werden, mit der man diese Annahme im Nachhinein beurteilen kann. Auch gibt es die Möglichkeit Einflussvariablen nicht-linear ins Modell aufzunehmen (z.B. via Splines, welche ebenso später im Kurs behandelt werden).

- Annahme 2.) $E(\varepsilon_i) = 0$ für alle verfügbare Annahmen

Die Residuen sind zufällige Fehler. Man nimmt an, dass diese Fehler im Erwartungswert 0 sind. Anders ausgedrückt: Die Modellgleichung $Y = X\beta$ spiegelt den Zusammenhang wieder. Es kommt keine systematische Störgröße ϵ mehr hinzu, von der man erwarten kann, dass sie im Mittel nicht 0 ist.

ist $X'X$ nicht mehr (stabil) invertierbar.

Wobei KQ-Schätzmethode generell nur nicht auf parametrisches Modell angewendet. Nein (aber es wäre ideal, wenn es so ist). Sie dürfen nur nicht perfekt oder näherungsweise perfekt linear abh. sein. $[\text{corr}(\cdot) = 1]$, dann kann $X'X$ nicht mehr (stabil) invertierbar.

4.) Müssen alle auf genannten Korrelationen x_1, \dots, x_k unterliegen? Nur nicht auf parametrisches Modell angewendet.

Exkurs zu Median (ϵ_i) $E(\epsilon_i) = 0 \rightarrow$ Per KQ-Design ist immer $MW(\epsilon_i) = 0$, da steht keine Information drin

Schauen uns $Med(\epsilon_i)$ an, da kann man schiefe Verteilung der Fehler [$E(\epsilon_i) \neq 0$] tatsächlichenken

Davon gehen wir aus.

(3) stetig linear in Likelihood drin

- Annahme 3.) $Var(\epsilon_i) = \sigma^2$: Homoskedastizität (Varianzhomogenität)

Es wird angenommen, dass die Varianz für jede Beobachtung identisch ist.

Diese Annahme könnte verletzt sein. Z.B. könnte bei größeren Kindern das Gewicht stärker streuen, als bei kleineren Kindern. Bei größeren Kindern könnte also eine höhere Varianz vorliegen.

Hinweise: Auch dies können ggf. Wissenschaftler aus dem jeweiligen Fachbereich besser beurteilen.

Prüfen kann man diese Annahme mittels später im Kurs vorgestellten Residuen-Plots (z.B. Residuen-vs.-Fitted values oder Residuen-vs.Kovariable-Plot).

Die Varianz-Annahme kann auch aufgeweicht werden (z.B. via verallgemeinerte lineare Modelle, welche später im Kurs vorge stellt werden).

(4) nötig für ML: Bildung der Likelihood $\prod_{i=1}^n \dots$ erfordert Unabh. der Beobachtungen ϵ_i zum Aus tauschen

- Annahme 4.) ϵ_i paarweise unabhängig

Es wird angenommen, dass die Beobachtungen der Zielgröße bei gegebenen Einflussgrößen unabhängig sind.

Diese Annahme ist bei der Wahl einer einigen kleinen Grundschule unwahrscheinlich. Es könnten z.B. Kinder der selben Familie in die Schule gehen und vielleicht genetisch bedingt zu Fettleibigkeit neigen. Diese Beobachtungen wären also z.B. nicht unabhängig voneinander.

Hinweise: Auch bei Verletzung dieser Annahme gibt es modellierungstechnische Möglichkeiten dies zu berücksichtigen. Z.B. mittels Gemischten Modellen (Mixed Models - z.B. für einen weiteren zufälligen Effekt für jede Familie - später in diesem Kurs) oder Korrelationen zwischen Beobachtungen (Generalized Estimation Equations (GEE) -> Vorlesung: Generalisierte Regression).

(5) MLE braucht Verteilung mit konst. Varianz, aber es muss nicht unbedingt M sein, denn jetzt Anwendung auf new data! dasselbe gilt für Prädiktion. Man sollte sich um Verteilung der Fehler (egal welche) sehr Sorgen!

Dass die übrige Reststreuung normalverteilt ist, ist durchaus denkbar. Störgrößen sind in der Natur häufig normalverteilt. Auch der Wertebereich der Zielgröße lässt nicht auf mögliche Verletzung der Normalverteilungsannahme hindeuten.

Explizit aufpassen muss man, wenn z.B. offensichtlich ist, dass die Größe nicht normalverteilt sein kann. Dies ist z.B. der Fall wenn die Zielgröße binär ist, nur ganze Zahlen annehmen kann oder häufig Werte an den Grenzen des möglichen Wertebereichs hat.

Hintergrund: Die Normalverteilung ist eine stetige, symmetrische Verteilung mit Werten von $-\infty$ bis ∞ . In den gerade genannten Fällen können die Residuen (= tatsächlicher y-Wert - geschätzter y-Wert) nicht oder nur schwer einer Normalverteilung folgen.

Überprüfen kann man die Normalverteilungsannahme z.B. mit dem später im Kurs vorgestelltem Q-Q-Plot.

c) Solange die Annahmen 1.) und 2.) erfüllt sind, können die Regressionsparameter problemlos interpretiert werden. Insbesondere für die Interpretation von Konfidenzintervallen, p-Werten etc. müssen aber entsprechende Annahmen erfüllt sein! (Es geht um die konkrete Bezifferung von Wahrscheinlichkeiten). Falls dies (teilweise) nicht der Fall ist, können diese wahrscheinlichkeitsbasierten Kennzahlen evtl. noch als Indikatoren für selbige verwendet werden, aber keinesfalls explizit mit Wahrscheinlichkeits-Interpretation!

e) KörperfgröBe und Alter sind korrelierte Einflussgrößen. Was folgt daraus für Interpretation der einzelnen Regr. Koeffizienten, c.p.?

Interpretation zunächst normal wie üblich: z.B. "Wenn Alter um ein Jahr steigt, nimmt Gewicht im Mittel um 3kg zu"

Vorsicht! Modellschätzung besagt nicht, dass Kind pro Jahr im Schnitt 3kg zunimmt. Kind wird auch wahrscheinlich insgesamt pro Jahr > 3kg an Gewicht zunehmen.

Brüder: Es geht nur um Effekt auf verschiedene elterliche Kinder GLEICHER KÖRPERGRÖSSE!! Kombination möglich! Wenn Gewichtszunahme durchschn. eines Kind halbes Jahr älter wird und 4 cm wächst => $\frac{1}{2}$ älter + 4 cm groß