

Additives Modell:

$$E(Y|X_1, \dots, X_p) = b_0 + f_1(x_1) + f_2(x_2) + \dots + f_p(x_p)$$

$$\text{bzw. } Y = b_0 + \dots + \varepsilon$$

Lineares Modell:

$$E(Y|X_1, \dots, X_p) = b_0 + b_1 x_1 + \dots + b_p x_p$$

Vereinf. lineares Modell:

$$E(Y|X_1, \dots, X_p) = f(b_0 + b_1 x_1 + \dots + b_p x_p)$$

Wdh. Schätzer:

$$(Y_1, X_1), \dots, (Y_n, X_n) \text{ iid} \Rightarrow \hat{b}_1 \text{ (wird geschätzt)}$$

$$T := g(X_1, Y_1, X_2, Y_2, \dots, X_n, Y_n) : (\Omega, \mathcal{F}, P) \rightarrow \mathbb{R}$$

$T$  wird in Abhängigkeit der Stichprobe berechnet

$\Rightarrow$  variiert je nach Stichprobe

$$E_{b_1}(T) = b_1$$

Univariate lineare Regression:

$$Y = b_0 + b_1 X + \varepsilon, \quad E(\varepsilon) = 0, \quad X, \varepsilon \text{ unabh.}$$

$$\text{KQ-Schätzer: } \hat{b}_1 = \frac{r_{xy}}{s_x} \cdot s_y$$

$$\hat{b}_0 = \bar{y} - \hat{b}_1 \bar{x}$$

$$E(\hat{b}_0) = b_0$$

$$E(\hat{b}_1) = b_1$$

$$\text{Var}(\hat{b}_1) = \frac{\sigma^2}{\sum (X_i - \bar{X})^2}$$

Multivariate lineare Regression:

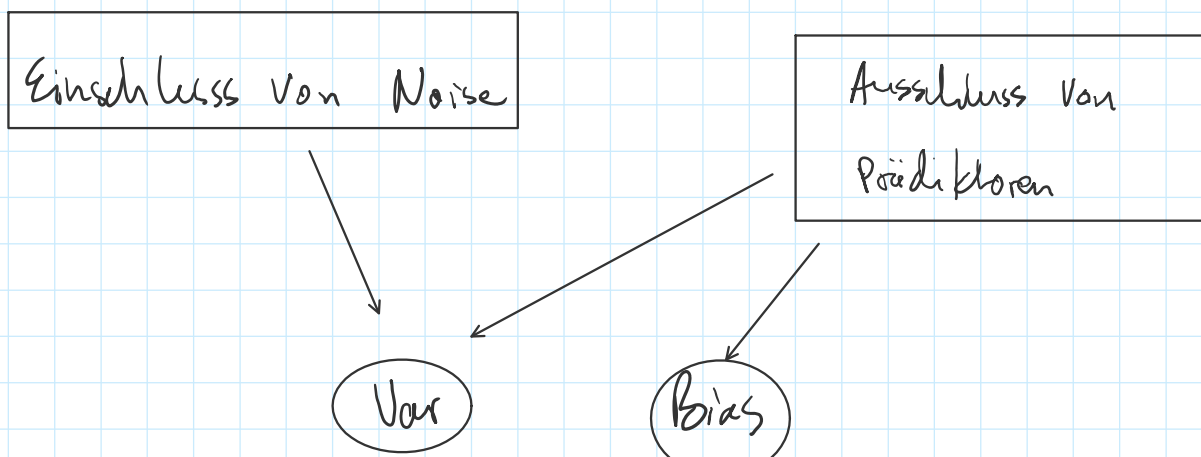
$$Y = b_0 + b_1 X_1 + b_2 X_2 + \varepsilon$$

$\varepsilon$ -freie  $\checkmark$

$$\text{Var}(\hat{b}_1) = \frac{\sigma^2}{(1 - R^2) \cdot \sum_{i=1}^n (X_i - \bar{X})^2} \quad \text{mit } R^2 \text{ Bestimmtheitsmaß}$$

$$\text{aus } X_1 = \tilde{b}_0 + \tilde{b}_2 X_2 + \varepsilon$$

Fehlerquellen:



Var

Bias

bsp.: Modellierung einer Noise-Var. im lin. Regr. Modell

• datengen. Modell:  $Y = b_1 X_1 + \epsilon = b_1 X_1 + 0 \cdot X_2 + \epsilon$

• Analysemodell:  $Y = b_1 X_1 + b_2 X_2 + \epsilon$

• Bias:  $E(\hat{b}_1) = b_1$

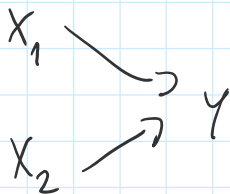
• Varianz:  $Var(\hat{b}_1) = \frac{\sigma^2}{(1-R^2) \sum_{i=1}^n (X_i - \bar{X})^2} \geq \frac{\sigma^2}{\sum_{i=1}^n (X_i - \bar{X})^2}$

bsp.: Fehlende Einflussvariable

• datengen. Modell  $Y = b_1 X_1 + b_2 X_2 + \epsilon = b_1 X_1 + \tilde{\epsilon}$

• Analysemodell  $Y = b_1 X_1 + \epsilon$

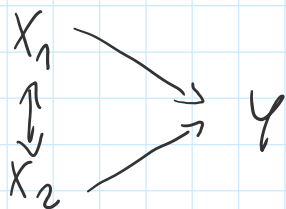
a) unabh. Prädiktoren



Bias:  $E(\hat{b}_1) = b_1$

Var:  $Var(\hat{b}_1) = \frac{\tilde{\sigma}^2}{\sum_{i=1}^n (X_i - \bar{X})^2} \geq \frac{\sigma^2}{\sum_{i=1}^n (X_i - \bar{X})^2}$  mit  $\tilde{\sigma}^2 = Var(\tilde{\epsilon})$

b) abh. Prädiktoren (Bsp.:  $X_2 = a \cdot X_1$ )



datengen. Modell:  $y = b_1 x_1 + b_2 x_2 + \varepsilon = (b_1 + b_2 \cdot a) x_1 + \varepsilon$

Bias:  $E(\hat{b}_1) = b_1 + b_2 \cdot a \Rightarrow R^2 \rightarrow 1$

Var: 
$$\text{Var}(\hat{b}_1) = \frac{\sigma^2}{\sum_{i=1}^n (x_i - \bar{x})^2} < \frac{\sigma^2}{\underbrace{(1-R^2)}_{( < 0 )} \sum_{i=1}^n (x_i - \bar{x})^2}$$

### Nichtparametrische Tests:

Vergleich: Gauss-Test (bekannter parametrischer Test)

$X_1, \dots, X_n$  iid.  $\Leftrightarrow X_i \sim N(\mu, \sigma^2)$ ,  $\sigma^2$  bekannt

$H_0: \{ \mu = 0 \}$  vs  $H_1: \{ \mu \neq 0 \}$

$$T = \frac{\frac{1}{n} \sum_{i=1}^n X_i}{\sigma} \cdot \sqrt{n} \underset{H_0}{\sim} N(0, 1)$$

Jetzt:

Ein-Stichproben-Test:

Rangtests

$X_1, \dots, X_n$  iid. ZV,  $X_i$  stetig,  $X_{\text{med}} = \text{Median von } X_i$

$H_0: \{X_{\text{med}} = S_0\}$  vs  $H_1: \{X_{\text{med}} \neq S_0\}$ , Sign. niveau  $\alpha$

Unter  $H_0$  gilt:  $P(X_i < S_0) = P(X_i \geq S_0) = 0.5$

$Y_i := \begin{cases} 0, & X_i < S_0 \\ 1, & X_i \geq S_0 \end{cases} \Rightarrow \text{unter } H_0 \text{ ist } Y_i \sim B(1, 0.5)$

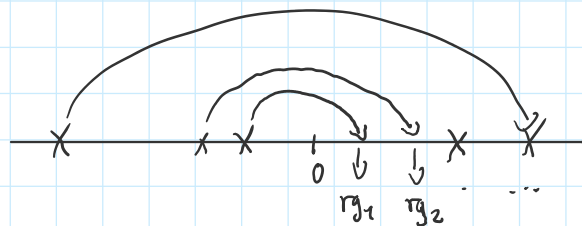
$$T := \sum_{i=1}^n Y_i \underset{H_0}{\sim} B(n, 0.5)$$

Wilcoxon - Vorzeichen-Rang-Test:

$X_1, \dots, X_n$  iid. ZV, stetig, symmetrisch,  $X_{\text{med}} = \text{Median der } X_i$

$H_0: \{X_{\text{med}} = S_0\}$  vs  $H_1: \{X_{\text{med}} \neq S_0\}$

$$D_i := X_i - S_0$$



Bilde Ränge der  $|D_i|$ ,  $rg(|D_i|)$

$$W^+ := \sum_{i: D_i > 0} rg(|D_i|) \quad W^- := \sum_{i: D_i < 0} rg(|D_i|)$$

Es gilt:  $W^+ + W^- = \frac{n(n+1)}{2}$ . Unter  $H_0$   $E(W^+) = \frac{n(n+1)}{4}$ .

$$T = \min(W^+, W^-)$$

$$T = \min(w^+, w^-)$$

Zweistichproben test:

Wilcoxon - Rangsummentest:

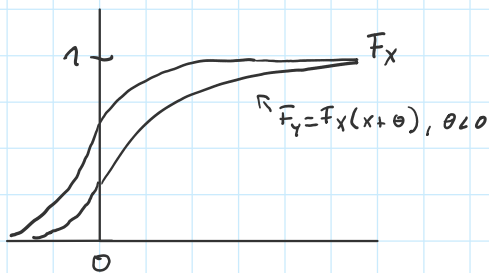
$$X_1, \dots, X_n \text{ iid.}$$

$$Y_1, \dots, Y_m \text{ "}$$

$$X_1, \dots, X_n, Y_1, \dots, Y_m \text{ unabh.}$$

$X_i$  stetig mit Verteil.  $F_X$ ,  $Y_i$  stetig mit Verteil.  $F_Y$

Es gelte das Lokations-Shift-Modell:  $F_X(x) = F_Y(x + \theta) \quad \forall x$



Frage: "Sind die Verteilungen gleich?"

$$H_0: \{ \theta = 0 \} \text{ vs. } H_1: \{ \theta \neq 0 \}$$

Bilde Ränge  $rg(X_1), rg(X_2), \dots, rg(Y_1), \dots, rg(Y_m)$

$$T = \sum_{i=1}^n rg(X_i)$$

$$T = \sum_{i=1}^n \eta(X_i)$$

Permutationstests:

Verb. Stichproben:  $(X_i, Y_i), i=1, \dots, n$  iid. ZV

$D_i := X_i - Y_i$  symmetrisch um  $\theta$

Fragestellung: "Gibt es systematische Tendenzen?"

$$H_0: \{ \theta = 0 \} \text{ vs } H_1: \{ \theta \neq 0 \}$$

Unter  $H_0$  gilt:  $X_i - Y_i \sim Y_i - X_i$

$$T = \frac{1}{n} \sum_{i=1}^n D_i = \frac{1}{n} \sum_{i=1}^n (X_i - Y_i) \sim \frac{1}{n} \left( \sum_{i \in M} (X_i - Y_i) + \sum_{i \notin M} (Y_i - X_i) \right) = T_M$$

für  $M \subset \{1, \dots, n\}$

$$p = p_0(T \geq t) = \frac{|\{M \subset \{1, \dots, n\}, T_M \geq t\}|}{2^n}$$

2 unabh. Stichproben (wie bei Rangsummentest):

$$Z := (X_1, \dots, X_n, Y_1, \dots, Y_m) = (Z_1, \dots, Z_{n+m})$$

Unter  $H_0$  gilt:  $Z \sim Z_{\pi} = (Z_{\pi(1)}, Z_{\pi(2)}, \dots, Z_{\pi(n+m)})$

$$T := \frac{1}{n} \sum_{i=1}^n X_i \underset{H_0}{\sim} \frac{1}{n} \sum_{i=1}^n Z_{\pi(i)} = T_{\pi}$$

$$p = P_0(T \geq t) = \frac{|\{ \pi; T_\pi \geq t \}|}{\binom{n+m}{n}}$$