

# Regressionskalibrierung

Statistische Herausforderungen im Umgang mit fehlenden bzw. fehlerbehafteten Daten

Le Minh-Anh

Ludwig Maximilians-Universität  
Institut für Statistik  
Bachelor-Seminar  
Betreuer: Prof. Dr. Thomas Augustin

5. Dezember, 2014

# Gliederung

- 1 Regressionskalibrierung-Algorithmus
  - Problemstellung
  - Algorithmus
- 2 Parameterschätzung
  - Validierungsdaten
  - Wiederholungsdaten
- 3 Bootstrapping
  - Resampling Vectors
  - Resampling Residuals
  - Bootstrap- Algorithmus
- 4 Fazit
- 5 Anhang

# Gliederung

- 1 Regressionskalibrierung-Algorithmus
  - Problemstellung
  - Algorithmus
- 2 Parameterschätzung
  - Validierungsdaten
  - Wiederholungsdaten
- 3 Bootstrapping
  - Resampling Vectors
  - Resampling Residuals
  - Bootstrap- Algorithmus
- 4 Fazit
- 5 Anhang



# Problemstellung

## fehlerbehaftete Daten

### Szenario:

Man interessiert sich für den Einfluss von  $\mathbf{X}$ ,  $\mathbf{Z}$  auf  $\mathbf{Y}$ .

Problem:  $\mathbf{X}_i$  wird fehlerhaft gemessen  $\rightarrow \mathbf{X}_i^* = \mathbf{X}_i + \mathbf{U}_i$

Eine naive Regression von  $\mathbf{Y}$  auf  $(\mathbf{X}^*, \mathbf{Z})$  führt zu verzerrte Inferenzen.

$\Rightarrow$  Anwendung von Messfehlerkorrektur-Verfahren

Für  $\mathbf{U}_i$

- [klassischer] Fehler
- nicht- differentieller Fehler

ist die *Regressionskalibrierung* anwendbar.

# Gliederung

- 1 Regressionskalibrierung-Algorithmus
  - Problemstellung
  - Algorithmus
- 2 Parameterschätzung
  - Validierungsdaten
  - Wiederholungsdaten
- 3 Bootstrapping
  - Resampling Vectors
  - Resampling Residuals
  - Bootstrap- Algorithmus
- 4 Fazit
- 5 Anhang

# Regressionskalibrierung-Algorithmus

Von Interesse ist  $E[Y|X, Z] = m_y(X, Z, \beta)$ , aber wahre  $X$  nicht vorhanden

- Schritt 1:  $E[X|X^*, Z] = m_x(X^*, Z, \gamma)$
- Schritt 2:  $E[Y|Z, X^*] \approx m_y(m_x(X^*, Z, \hat{\gamma}), Z, \beta_{RK})$   
mit wahre  $\beta \approx \beta_{RK}$
- Schritt 3: Schätze Standardabweichung von  $\hat{\beta}_{RK}$  durch Bootstrapping  
(oder andere Methoden)

# Daten

## Schritt 1 & 2

Möglichkeiten für Parameterschätzung in Schritt 1 ist abhängig von der vorliegenden Datenstruktur.

- interne Validierungsdaten
- Wiederholungsdaten
- Instrumentaldaten

# Gliederung

- 1 Regressionskalibrierung-Algorithmus
  - Problemstellung
  - Algorithmus
- 2 Parameterschätzung
  - Validierungsdaten
  - Wiederholungsdaten
- 3 Bootstrapping
  - Resampling Vectors
  - Resampling Residuals
  - Bootstrap- Algorithmus
- 4 Fazit
- 5 Anhang

# Validierungsdaten

## Schritt 1 & 2

Für einen Teil der Daten liegen wahre X-Werte vor

$i$	$Y_i$	$X_i$	$X_i^*$	$Z_i$
1	$Y_1$	$X_1$	$X_1^*$	$Z_1$
⋮	⋮	⋮	⋮	⋮
k	$Y_k$	$X_k$	$X_k^*$	$Z_k$
⋮	⋮	⋮	⋮	⋮
k+1	$Y_{k+1}$	NA	$X_{k+1}^*$	$Z_{k+1}$
⋮	⋮	⋮	⋮	⋮
n	$Y_n$	NA	$X_n^*$	$Z_n$

z.B. Untersuchung auf Prostatakrebs durch Stanzbiopsie ist teuer und schmerzhaft. Alternativ kann man eine einfache Blutprobe entnehmen, die aber weniger genau ist.

# Validierungsdaten

## Schritt 1 & 2

Für einen Teil der Daten liegen wahre X-Werte vor

$i$	$Y_i$	$X_i$	$X_i^*$	$Z_i$	$\hat{X}_i$
1	$Y_1$	$X_1$	$X_1^*$	$Z_1$	NA
$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$
k	$Y_k$	$X_k$	$X_k^*$	$Z_k$	NA
$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$
k+1	$Y_{k+1}$	NA	$X_{k+1}^*$	$Z_{k+1}$	$\hat{X}_{k+1}$
$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$
n	$Y_n$	NA	$X_n^*$	$Z_n$	$\hat{X}_n$

- Schritt 1: Rechne für  $i \in \{1, \dots, k\}$  eine Regression  $X \sim X^* + Z$

$$\rightarrow E[X|X^*, Z] = m_X(X^*, Z, \gamma)$$

Berechne  $\hat{X}_i = (\underline{1}, \underline{X}_i^*, \underline{Z}_i) \hat{\gamma}$  für  $i \in \{k+1, \dots, n\}$

# Validierungsdaten

## Schritt 1 & 2

Für einen Teil der Daten liegen wahre X-Werte vor

$i$	$Y_i$	$X_i$	$X_i^*$	$Z_i$	$\hat{X}_i$	$X_{reg,i}$	$V_i$
1	$Y_1$	$X_1$	$X_1^*$	$Z_1$	NA	$X_1$	1
$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$
k	$Y_k$	$X_k$	$X_k^*$	$Z_k$	NA	$X_k$	1
$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$
k+1	$Y_{k+1}$	NA	$X_{k+1}^*$	$Z_{k+1}$	$\hat{X}_{k+1}$	$\hat{X}_{k+1}$	0
$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$
n	$Y_n$	NA	$X_n^*$	$Z_n$	$\hat{X}_n$	$\hat{X}_n$	0

- Schritt 1: Rechne für  $i \in \{1, \dots, k\}$  eine Regression  $X \sim X^* + Z$

$$\rightarrow E[X|X^*, Z] = m_X(X^*, Z, \gamma)$$

Berechne  $\hat{X}_i = (\underline{1}, \underline{X}_i^*, \underline{Z}_i) \hat{\gamma}$  für  $i \in \{k+1, \dots, n\}$



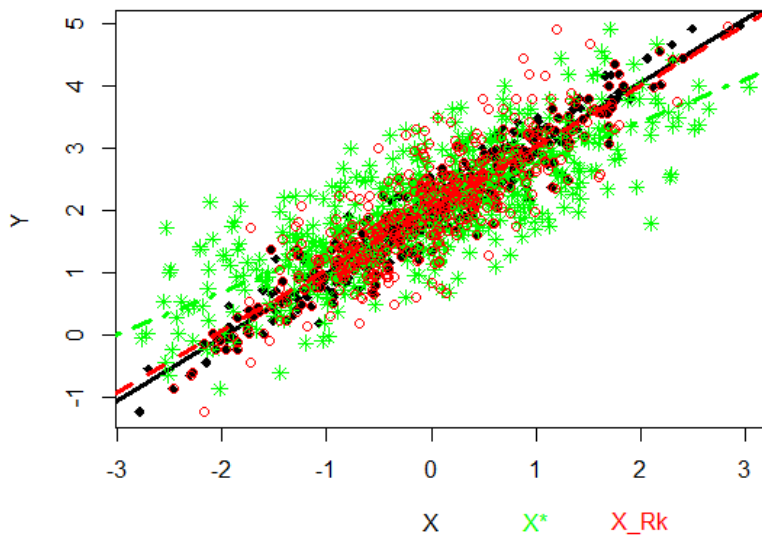
# Validierungsdaten

## Schritt 1 & 2

Für einen Teil der Daten liegen wahre X-Werte vor

$i$	$Y_i$	$X_i$	$X_i^*$	$Z_i$	$\hat{X}_i$	$X_{reg,i}$	$V_i$
1	$Y_1$	$X_1$	$X_1^*$	$Z_1$	NA	$X_1$	1
$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$
k	$Y_k$	$X_k$	$X_k^*$	$Z_k$	NA	$X_k$	1
$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$
k+1	$Y_{k+1}$	NA	$X_{k+1}^*$	$Z_{k+1}$	$\hat{X}_{k+1}$	$\hat{X}_{k+1}$	0
$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$
n	$Y_n$	NA	$X_n^*$	$Z_n$	$\hat{X}_n$	$\hat{X}_n$	0

- Schritt 2: Rechne für alle  $i$   $Y \sim X_{reg} + Z + V$   
 $\rightarrow E[Y|X_{reg}, Z, V] = m_y(X_{reg}, Z, V, \beta_{RK})$



# Gliederung

- 1 Regressionskalibrierung-Algorithmus
  - Problemstellung
  - Algorithmus
- 2 Parameterschätzung
  - Validierungsdaten
  - Wiederholungsdaten
- 3 Bootstrapping
  - Resampling Vectors
  - Resampling Residuals
  - Bootstrap- Algorithmus
- 4 Fazit
- 5 Anhang

# Wiederholungsdaten

## Schritt 1 & 2

Es liegen keine wahren Werte vor, aber Messwiederholungen (hier:  $k_i = 4$ )

Weitere Annahmen:

eine Einflussgröße  $X$  und eine Zielgröße  $Y$ ,  $E(\mathbf{Y}|\mathbf{X}) = \beta_0 + \beta_1 \mathbf{X}$

$i$	$Y_i$	$X_{i1}^*$	$X_{i2}^*$	$X_{i3}^*$	$X_{i4}^*$
1	$Y_1$	$X_{11}^*$	$X_{12}^*$	$X_{13}^*$	$X_{14}^*$
$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$
$k$	$Y_k$	$X_{k1}^*$	$X_{k2}^*$	$X_{k3}^*$	$X_{k4}^*$
$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$
$n$	$Y_n$	$X_{n1}^*$	$X_{n2}^*$	$X_{n3}^*$	$X_{n4}^*$

z.B. Wiederholte Messung von Angst anhand von Herzfrequenz, wahrer Angstwert nicht beobachtbar.

# Wiederholungsdaten

## Schritt 1 & 2

Es liegen keine wahren Werte vor, aber Messwiederholungen (hier:  $k_i = 4$ )

Weitere Annahmen:

eine Einflussgröße  $\mathbf{X}$  und eine Zielgröße  $Y$ ,  $E(\mathbf{Y}|\mathbf{X}) = \beta_0 + \beta_1 \mathbf{X}$

$i$	$Y_i$	$X_{i1}^*$	$X_{i2}^*$	$X_{i3}^*$	$X_{i4}^*$	$\bar{X}_i^*$
1	$Y_1$	$X_{11}^*$	$X_{12}^*$	$X_{13}^*$	$X_{14}^*$	$\bar{X}_{1\cdot}^*$
$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$
$k$	$Y_k$	$X_{k1}^*$	$X_{k2}^*$	$X_{k3}^*$	$X_{k4}^*$	$\bar{X}_k^*$
$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$
$n$	$Y_n$	$X_{n1}^*$	$X_{n2}^*$	$X_{n3}^*$	$X_{n4}^*$	$\bar{X}_n^*$

- Schritt 0:  $\bar{X}_i^* = \frac{X_{i1}^* + X_{i2}^* + X_{i3}^* + X_{i4}^*}{4}$

# Wiederholungsdaten

## Schritt 1 & 2

Es liegen keine wahren Werte vor, aber Messwiederholungen (hier:  $k_i = 4$ )

Weitere Annahmen:

eine Einflussgröße  $X$  und eine Zielgröße  $Y$ ,  $E(\mathbf{Y}|\mathbf{X}) = \beta_0 + \beta_1 \mathbf{X}$

$i$	$Y_i$	$X_{i1}^*$	$X_{i2}^*$	$X_{i3}^*$	$X_{i4}^*$	$\bar{X}_{i\cdot}^*$	$\hat{X}_i$
1	$Y_1$	$X_{11}^*$	$X_{12}^*$	$X_{13}^*$	$X_{14}^*$	$\bar{X}_{1\cdot}^*$	$\hat{X}_1$
$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$
$k$	$Y_k$	$X_{k1}^*$	$X_{k2}^*$	$X_{k3}^*$	$X_{k4}^*$	$\bar{X}_{k\cdot}^*$	$\hat{X}_k$
$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$
$n$	$Y_n$	$X_{n1}^*$	$X_{n2}^*$	$X_{n3}^*$	$X_{n4}^*$	$\bar{X}_{n\cdot}^*$	$\hat{X}_n$

- Schritt 1: Berechne

$$E[\widehat{X}_i | \bar{X}_i^*] \approx \frac{4\hat{\sigma}_x^2}{4\hat{\sigma}_x^2 + \hat{\sigma}_u^2} \bar{X}_i^* + \hat{\mu}_{x^*} \left(1 - \frac{4\hat{\sigma}_x^2}{4\hat{\sigma}_x^2 + \hat{\sigma}_u^2}\right) = \hat{X}_i$$

# Wiederholungsdaten

## Schritt 1 & 2

Es liegen keine wahren Werte vor, aber Messwiederholungen (hier:  $k_i = 4$ )

Weitere Annahmen:

eine Einflussgröße  $X$  und eine Zielgröße  $Y$ ,  $E(\mathbf{Y}|\mathbf{X}) = \beta_0 + \beta_1 \mathbf{X}$

$i$	$Y_i$	$X_{i1}^*$	$X_{i2}^*$	$X_{i3}^*$	$X_{i4}^*$	$\bar{X}_i^*$	$\hat{X}_i$
1	$Y_1$	$X_{11}^*$	$X_{12}^*$	$X_{13}^*$	$X_{14}^*$	$\bar{X}_1^*$	$\hat{X}_1$
$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$
$k$	$Y_k$	$X_{k1}^*$	$X_{k2}^*$	$X_{k3}^*$	$X_{k4}^*$	$\bar{X}_k^*$	$\hat{X}_k$
$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$
$n$	$Y_n$	$X_{n1}^*$	$X_{n2}^*$	$X_{n3}^*$	$X_{n4}^*$	$\bar{X}_n^*$	$\hat{X}_n$

- Schritt 1: Berechne

$$E[\widehat{X}_i | \bar{X}_i^*] \approx \frac{4\hat{\sigma}_x^2}{4\hat{\sigma}_x^2 + \hat{\sigma}_u^2} \bar{X}_i^* + \hat{\mu}_{x^*} \left(1 - \frac{4\hat{\sigma}_x^2}{4\hat{\sigma}_x^2 + \hat{\sigma}_u^2}\right) = \hat{X}_i$$

$$\hat{\Sigma}_{uu} = \hat{\sigma}_u^2 = \frac{\sum_{i=1}^n \sum_{j=1}^4 (\mathbf{x}_{ij}^* - \bar{\mathbf{x}}_{i\cdot}^*) (\mathbf{x}_{ij}^* - \bar{\mathbf{x}}_{i\cdot}^*)^t}{3n}$$

# Wiederholungsdaten

## Schritt 1 & 2

Es liegen keine wahren Werte vor, aber Messwiederholungen (hier:  $k_i = 4$ )

Weitere Annahmen:

eine Einflussgröße  $X$  und eine Zielgröße  $Y$ ,  $E(\mathbf{Y}|\mathbf{X}) = \beta_0 + \beta_1 \mathbf{X}$

$i$	$Y_i$	$X_{i1}^*$	$X_{i2}^*$	$X_{i3}^*$	$X_{i4}^*$	$\bar{X}_{i\cdot}^*$	$\hat{X}_i$
1	$Y_1$	$X_{11}^*$	$X_{12}^*$	$X_{13}^*$	$X_{14}^*$	$\bar{X}_{1\cdot}^*$	$\hat{X}_1$
$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$
$k$	$Y_k$	$X_{k1}^*$	$X_{k2}^*$	$X_{k3}^*$	$X_{k4}^*$	$\bar{X}_{k\cdot}^*$	$\hat{X}_k$
$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$
$n$	$Y_n$	$X_{n1}^*$	$X_{n2}^*$	$X_{n3}^*$	$X_{n4}^*$	$\bar{X}_{n\cdot}^*$	$\hat{X}_n$

- Schritt 1: Berechne

$$\hat{\mu}_X = \hat{\mu}_{X^*} = \frac{\sum_{i=1}^n \bar{X}_{i\cdot}^*}{n}$$



# Wiederholungsdaten

## Schritt 1 & 2

Es liegen keine wahren Werte vor, aber Messwiederholungen (hier:  $k_i = 4$ )  
 Weitere Annahmen:

eine Einflussgröße  $X$  und eine Zielgröße  $Y$ ,  $E(\mathbf{Y}|\mathbf{X}) = \beta_0 + \beta_1 \mathbf{X}$

$i$	$Y_i$	$X_{i1}^*$	$X_{i2}^*$	$X_{i3}^*$	$X_{i4}^*$	$\bar{X}_{i\cdot}^*$	$\hat{X}_i$
1	$Y_1$	$X_{11}^*$	$X_{12}^*$	$X_{13}^*$	$X_{14}^*$	$\bar{X}_{1\cdot}^*$	$\hat{X}_1$
$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$
$k$	$Y_k$	$X_{k1}^*$	$X_{k2}^*$	$X_{k3}^*$	$X_{k4}^*$	$\bar{X}_{k\cdot}^*$	$\hat{X}_k$
$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$
$n$	$Y_n$	$X_{n1}^*$	$X_{n2}^*$	$X_{n3}^*$	$X_{n4}^*$	$\bar{X}_{n\cdot}^*$	$\hat{X}_n$

- Schritt 1: Berechne

$$\hat{\mu}_X = \hat{\mu}_{X^*} = \frac{\sum_{i=1}^n \bar{X}_{i\cdot}^*}{n}$$

$$\hat{\Sigma}_{XX} = \hat{\sigma}_X^2 = \frac{\sum_{i=1}^n (\bar{X}_{i\cdot}^* - \hat{\mu}_{X^*})(\bar{X}_{i\cdot}^* - \hat{\mu}_{X^*})^t}{(n-1)} - \frac{\hat{\Sigma}_{UU}}{4}$$

# Wiederholungsdaten

## Schritt 1 & 2

- Wiederholungsmessungen ermöglichen die Schätzung von  $\Sigma_{uu}$
- auch ohne einen einzigen wahren Wert von  $X$  kann  $X$  geschätzt werden

$$E[\widehat{\mathbf{X}}_i | \widehat{\mathbf{X}}_i^*] \approx \underbrace{\frac{4\widehat{\sigma}_x^2}{4\widehat{\sigma}_x^2 + \widehat{\sigma}_u^2}}_{\widehat{\gamma}_1} \mathbf{X}_i^* + \underbrace{\widehat{\mu}_{x^*} \left(1 - \frac{4\widehat{\sigma}_x^2}{4\widehat{\sigma}_x^2 + \widehat{\sigma}_u^2}\right)}_{\widehat{\gamma}_0} = \widehat{\mathbf{X}}_i$$

# Approximations-Schritt

- Schritt 2:

Ersetze die nicht beobachtete Variable  $\mathbf{X}$  durch die im vorherigen Schritt durchgeführte Schätzung, d.h. ersetze im Hauptmodell  $\mathbf{X}$  durch  $m_{\mathbf{x}}(\mathbf{X}^*, \hat{\gamma})$ . Führe anschließend eine Standardanalyse durch, um die Parameterschätzer zu erhalten. Somit erhält man:

$$E[\mathbf{Y}|\overline{\mathbf{X}}^*] \approx m_{\mathbf{y}}(\underbrace{m_{\mathbf{x}}(\overline{\mathbf{X}}^*, \hat{\gamma})}_{\hat{\mathbf{x}}}, \beta_{RK})$$

# Approximations-Schritt

- Schritt 2:

$$\begin{aligned} E(Y|\bar{X}^*) &= E(\{E(Y|X, \bar{X}^*)\} |\bar{X}^*) \\ &= E(\{E(Y|X)\} |\bar{X}^*) \\ &= E(\{\beta_0 + \beta_1 X\} |\bar{X}^*) \\ &= \beta_0 + \beta_1 E(X|\bar{X}^*) \end{aligned}$$

# Approximations-Schritt

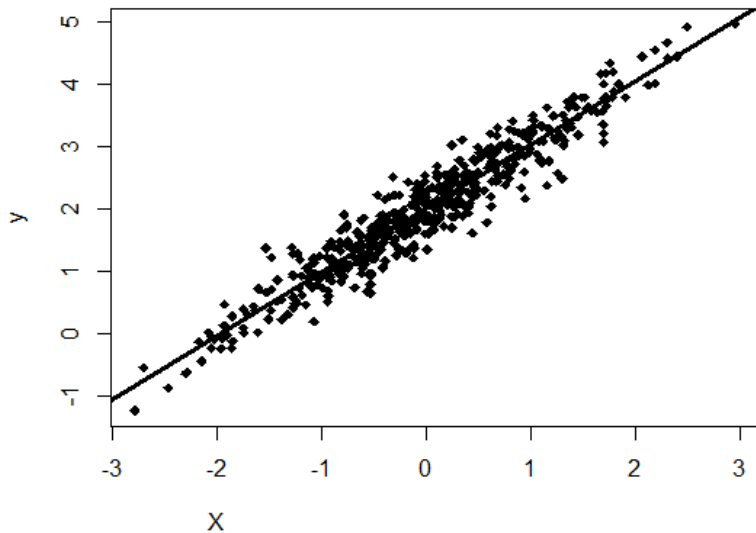
- Schritt 2:

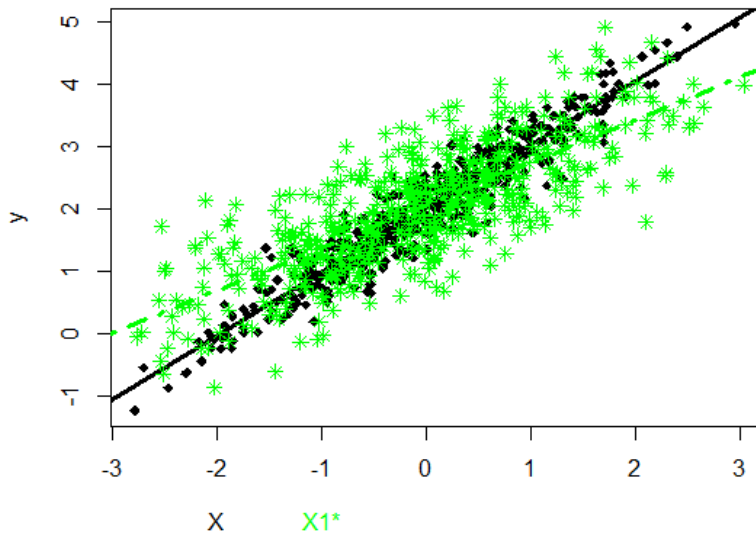
$$\begin{aligned}
 E(Y|\bar{X}^*) &= E(\{E(Y|X, \bar{X}^*)\}|\bar{X}^*) \\
 &= E(\{E(Y|X)\}|\bar{X}^*) \\
 &= E(\{\beta_0 + \beta_1 X\}|\bar{X}^*) \\
 &= \beta_0 + \beta_1 E(X|\bar{X}^*) \\
 &\approx \beta_{RK_0} + \beta_{RK_1} \left[ \underbrace{\left(\frac{4\hat{\sigma}_x^2}{4\hat{\sigma}_x + \hat{\sigma}_u}\right)}_{\hat{\gamma}_1} \bar{X}^* + \hat{\mu}_{X^*} \underbrace{\left(1 - \frac{4\hat{\sigma}_x^2}{4\hat{\sigma}_x + \hat{\sigma}_u}\right)}_{\hat{\gamma}_0} \right] \\
 &\approx \underbrace{\beta_{RK_0} \beta_{RK_1} \hat{\gamma}_0}_{\beta_{naiv_0}} + \underbrace{\beta_{RK_1} \hat{\gamma}_1}_{\beta_{naiv_1}} \bar{X}^*
 \end{aligned}$$

## Approximations-Schritt

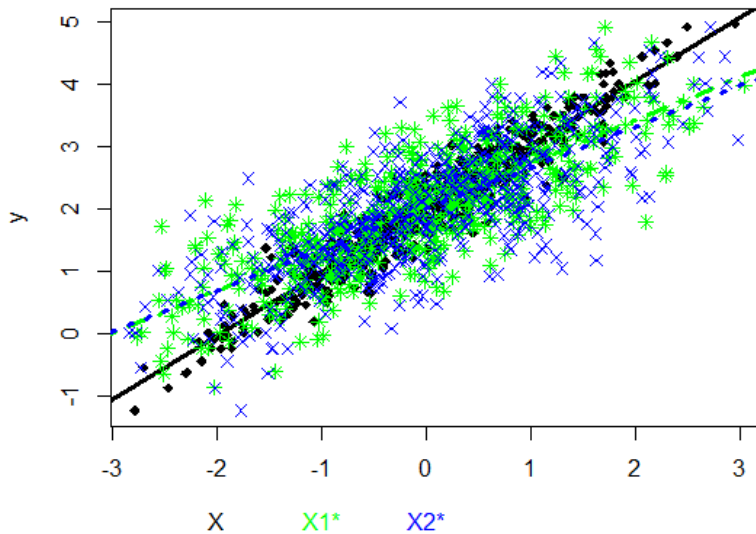
Die Schätzer für  $\beta_0$  &  $\beta_1$  können extrahiert werden, vorausgesetzt die Schätzung von  $\mathbf{X}$  in Schritt 1 der RK ist gültig.

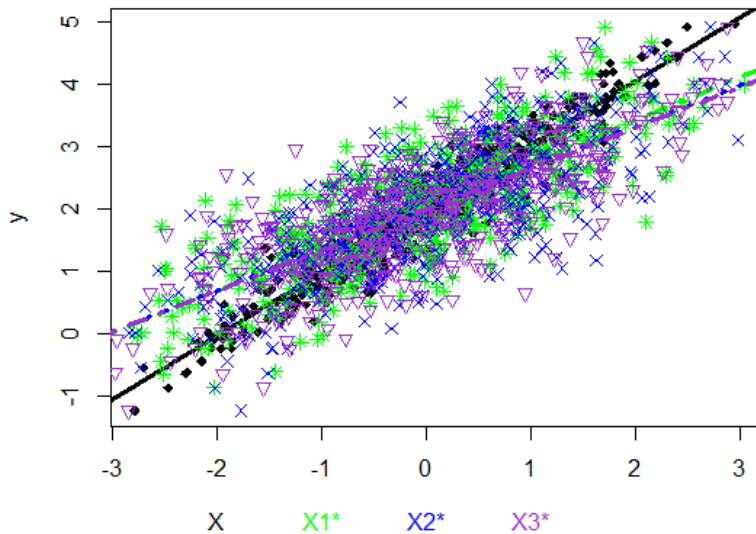
$$\hat{\beta}_1 \approx \hat{\beta}_{RK_1} = \frac{\hat{\beta}_{naiv_1}}{\hat{\gamma}_1}, \quad \hat{\beta}_0 \approx \hat{\beta}_{RK_0} = \frac{\hat{\beta}_{naiv_0}}{\hat{\gamma}_0 \hat{\beta}_1},$$

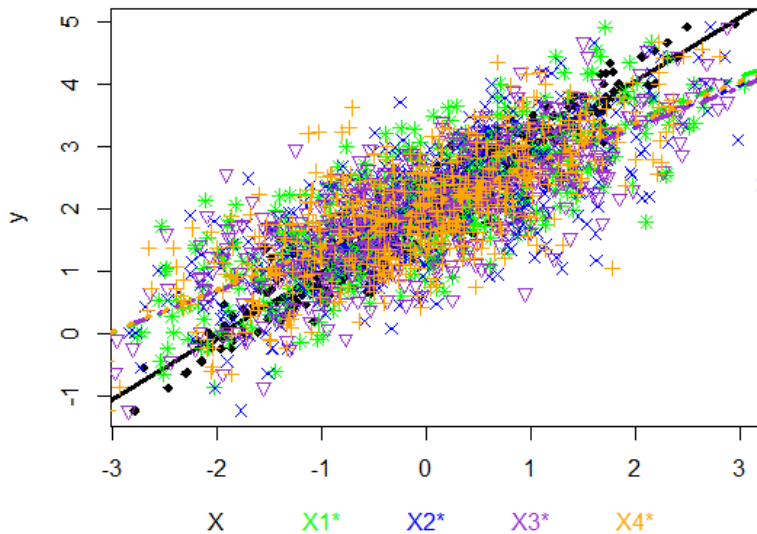


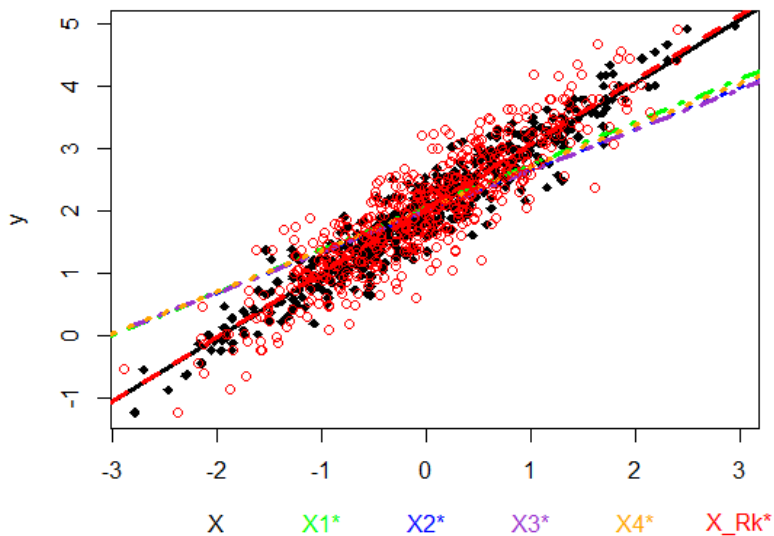












# Bootstrapping

- Schritt 3: Bootstrapping

Für die resultierende  $\hat{\beta}_{RK}$  soll im dritten Schritt der Standardfehler mit Bootstrapping (oder andere Methoden) geschätzt werden.

- Warum?

Da für die Schätzung von  $\mathbf{Y}$  wiederum eine Schätzung von  $\mathbf{X}$  eingesetzt wurde.

⇒ resultierende p-Werte und Standardabweichungen in statistischen Programmen nur approximativ als “ersten Eindruck” anzusehen.

# Bootstrapping

## Schritt 3

### Bootstrapping

- parametrische Verfahren  
Ziehen aus einer angenommenen Verteilung
- nonparametrische Verfahren  
Ziehen mit zurücklegen aus den vorliegenden Daten.  
(Stichprobenumfang  $\equiv$  Umfang vorliegende Daten)
  - resampling Vectors
  - resampling Residuals

Vorteil von Bootstrapping ist, dass bekannte statistische Verfahren auf Bootstrapp-Stichproben angewendet werden kann.

z.B. Regressionskalibrierung.

# Bootstrapping bei Messfehlerkorrektur

## Schritt 3

- Daten können in verschiedenen Formen vorliegen.  
z.B. Mischung aus interne Validierungsdaten, Instrumentaldaten, Wiederholungsdaten (2,3,4,...Messwiederholungen)
- Unterschiedliche Datenstrukturen → Unterschiedliche Informationen
- Für Bootstrapping bei Messfehlerkorrektur gilt:
  - Teildatensätze bilden (gruppieren nach den verschiedenen Strukturen)  
aus Teildatensätze Bootstrap-Stichproben ziehen  
→ ermöglicht Ziehungen aus homogener Umgebung  
→ geringere Varianz

# Gliederung

- 1 Regressionskalibrierung-Algorithmus
  - Problemstellung
  - Algorithmus
- 2 Parameterschätzung
  - Validierungsdaten
  - Wiederholungsdaten
- 3 Bootstrapping
  - Resampling Vectors
  - Resampling Residuals
  - Bootstrap- Algorithmus
- 4 Fazit
- 5 Anhang



# Resampling Vectors

Validierungsdaten:

$i$	$Y_i$	$X_i$	$X_i^*$	$Z_i$
1	$Y_1$	$X_1$	$X_1^*$	$Z_1$
$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$
$\frac{k}{k+1}$	$Y_k$	$X_k$	$X_k^*$	$Z_k$
$\frac{k+1}{k+1}$	$Y_{k+1}$	NA	$X_{k+1}^*$	$Z_{k+1}$
$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$
n	$Y_n$	NA	$X_n^*$	$Z_n$

Vektorweise ziehen mit zurücklegen aus

$\{(Y_i, X_i, X_i^*, Z_i)\}_{i=1}^k$  bzw.  $\{(Y_i, X_i^*, Z_i)\}_{i=k+1}^n$

- Vorteil: Kaum Annahmen müssen getroffen werden.  
besondere Beziehungen müssen nicht explizit berücksichtigt werden.  
z.B. wenn  $\epsilon_i$  von  $Z_i$  abhängt
- Nachteil: Bootstrapsstichprobe enthält nicht die gleiche Variablenmenge.  
z.B. High-Leverage Point mehrmals oder gar nicht enthalten.

# Gliederung

- 1 Regressionskalibrierung-Algorithmus
  - Problemstellung
  - Algorithmus
- 2 Parameterschätzung
  - Validierungsdaten
  - Wiederholungsdaten
- 3 Bootstrapping
  - Resampling Vectors
  - Resampling Residuals
  - Bootstrap- Algorithmus
- 4 Fazit
- 5 Anhang

# Resampling Residuals- Validierungsdaten

## Validierungsdaten

$i$	$Y_i$	$X_i$	$X_i^*$	$Z_i$
1	$Y_1$	$X_1$	$X_1^*$	$Z_1$
$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$
$k$	$Y_k$	$X_k$	$X_k^*$	$Z_k$
$k+1$	$Y_{k+1}$	NA	$X_{k+1}^*$	$Z_{k+1}$
$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$
$n$	$Y_n$	NA	$X_n^*$	$Z_n$

- Annahme: Zwei Regressionsmodelle liegen vor
  - $Y_i \sim (Z_i, X_i), \epsilon_i \sim \text{iid mit } \Sigma_i \approx \Sigma$
  - $X_j^* \sim (Z_j, X_j), \epsilon_j \sim \text{iid mit } \Sigma_j \approx \Sigma'$

# Resampling Residuals- Validierungsdaten

## Validierungsdaten

i	$Y_i$	$X_i$	$X_i^*$	$Z_i$	$Y_i^{(1)}$	$\dots$	$Y_i^{(M)}$
1	$Y_1$	$X_1$	$X_1^*$	$Z_1$	$Y_1^{(1)}$	$\dots$	$Y_1^{(M)}$
$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$
k	$Y_k$	$X_k$	$X_k^*$	$Z_k$	$Y_k^{(1)}$	$\dots$	$Y_k^{(M)}$
$\frac{k+1}{n}$	$Y_{k+1}$	NA	$X_{k+1}^*$	$Z_{k+1}$	$\vdots$	$\vdots$	$\vdots$
$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$
n	$Y_n$	NA	$X_n^*$	$Z_n$	$\vdots$	$\vdots$	$\vdots$

- Bootstrapping  $Y_i^{(m)}$ ,  $M = \#$  Bootstrapschprobe und  $m \in \{1, \dots, M\}$ 
  - 1.  $\epsilon_i = Y_i - m_y(Z_i, X_i, \hat{B})$  für  $i \in \{1, \dots, k\}$
  - 2.  $B = \{(\epsilon_i - \bar{\epsilon})\}_i^k$
  - 3. k mal Ziehen mit zurücklegen aus B  $\rightarrow \{\epsilon_i^{(m)}\}_i^k$
  - 4.  $Y_i^{(m)} = m_y(Z_i, X_i, \hat{B}) + \epsilon_i^{(m)}$  für  $i \in \{1, \dots, k\}$

## Resampling Residuals- Validierungsdaten

## Validierungsdaten

$i$	$Y_i$	$X_i$	$X_i^*$	$Z_i$	$Y_i^{(1)}$	$\dots$	$Y_i^{(M)}$	$X_i^{*(1)}$	$\dots$	$X_i^{*(M)}$
1	$Y_1$	$X_1$	$X_1^*$	$Z_1$	$Y_1^{(1)}$	$\dots$	$Y_1^{(M)}$	$X_1^{*(1)}$	$\dots$	$X_1^{*(M)}$
$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$
k	$Y_k$	$X_k$	$X_k^*$	$Z_k$	$Y_k^{(1)}$	$\dots$	$Y_k^{(M)}$	$X_k^{*(1)}$	$\dots$	$X_k^{*(M)}$
$\frac{k+1}{-}$	$Y_{k+1}$	NA	$X_{k+1}^*$	$Z_{k+1}$	$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$
$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$
n	$Y_n$	NA	$X_n^*$	$Z_n$	$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$

- Bootstrapping  $Y_i^{(m)}$ ,  $M = \#$  Bootstrapschprobe und  $m \in \{1, \dots, M\}$ 
  - 1.  $\epsilon_i = Y_i - m_y(Z_i, X_i, \hat{B})$  für  $i \in \{1, \dots, k\}$
  - 2.  $B = \{(\epsilon_i - \bar{\epsilon})\}_i^k$
  - 3. k mal Ziehen mit zurücklegen aus B  $\rightarrow \{\epsilon_i^{(m)}\}_i^k$
  - 4.  $Y_i^{(m)} = m_y(Z_i, X_i, \hat{B}) + \epsilon_i^{(m)}$  für  $i \in \{1, \dots, k\}$
- Bootstrapping  $X_i^{*(m)}$ , analog

# Gliederung

## 1 Regressionskalibrierung-Algorithmus

- Problemstellung
- Algorithmus

## 2 Parameterschätzung

- Validierungsdaten
- Wiederholungsdaten

## 3 Bootstrapping

- Resampling Vectors
- Resampling Residuals
- Bootstrap- Algorithmus

## 4 Fazit

## 5 Anhang

## Bootstrap- Algorithmus

- Schritt 1: M Bootstrapstichproben ziehen
- Schritt 2: Schritt 1 & 2 des Regressionskalibrierungs-Algorithmus anwenden

Man erhält somit nach M Durchläufen die Parameter

$$\hat{\beta}_{RK_k}^{(1)}, \dots, \hat{\beta}_{RK_k}^{(M)}$$

- Schritt 3: Aus den vorliegenden  $\hat{\beta}_{RK_k}$ s kann nun die Standardabweichung  $\hat{\sigma}_{\beta_{RK_k}}$  geschätzt werden

$$\hat{\sigma}_{RK_k}^2 = \widehat{\text{var}}(\hat{\beta}_{RK_k}) = \frac{1}{M-1} \sum_{m=1}^M (\hat{\beta}_{RK_k}^{(m)} - \overline{\hat{\beta}_{RK_k}})(\hat{\beta}_{RK_k}^{(m)} - \overline{\hat{\beta}_{RK_k}})^t.$$

# Nachteile-Regressionskalibrierung

## Nachteile

- nur approximatives Verfahren
- die Berechnung der Regression von  $\mathbf{X}$  auf  $(\mathbf{X}^*, \mathbf{Z})$  stellt eine Herausforderung dar, da  $\mathbf{X}$  nicht beobachtbar ist
- Schätzer sind nicht unbedingt konsistent (abhängig vom Modell)



# Vorteile-Regressionskalibrierung

## Vorteile I

- auf viele Modelle anwendbar (GLM)
- effektive Methode im Umgang mit fehlerhaften gemeinsam Einflussgrößen
- einfache Berechnung
- Regression  $Y \sim X+Z$  möglich, obwohl wahres  $X$  nicht beobachtet
- anschließende Standardanalysen noch möglich
- Reduzierung der Bias
- keine extra Implementierung in statistische Programme nötig

Vielen Dank für eure Aufmerksamkeit.

# Anhang

- *RK mit Instrumentaldaten*
- *Vergleich Valid Wdh1 Wdh4*
- *Vergleich RK Valid mit/ohne Dummyvariable*
- *Überprüfen der Schätzung in Schritt 1 der RK*

# Instrumentaldaten

## Schritt 1 & 2

### Instrumentaldaten $\mathbf{T}$

- $\mathbf{T}$  ist abhängig von  $\mathbf{X}$
- $\mathbf{T}$  ist unkorreliert mit Fehler  $\mathbf{U} = \mathbf{X}^* - \mathbf{X}$
- $\mathbf{T}$  unkorreliert mit  $\epsilon = \mathbf{Y} - E[\mathbf{Y}|\mathbf{Z}, \mathbf{X}]$

Außerdem soll gelten  $\mathbf{T}$  ist *unverzerrt* für  $\mathbf{X}$  d.h. eine Regression von  $\mathbf{T} \sim \mathbf{Z} + \mathbf{X}^*$  entspricht einer Regression von  $\mathbf{X} \sim \mathbf{Z} + \mathbf{X}^*$

$$E[\mathbf{T}|\mathbf{X}^*, \mathbf{Z}] = E[\mathbf{X}|\mathbf{X}^*, \mathbf{Z}]$$

# Instrumentaldaten

## Schritt 1 & 2

Beispiel:

- $Y$ : Brustkrebs
- $X$ : langfristige durchschnittliche Aufnahme von Nährstoffen
- $T$ : durchschnittliche Aufnahme von Nährstoffen, extrahiert aus einem Ernährungs-Tagebuch über vier Wochen (professionell dokumentiert)
- $X^*$ : durchschnittliche Aufnahme von Nährstoffen, extrahiert aus einem Fragebogen über Ernährung

# Instrumentaldaten

## Schritt 1 & 2

Für einen Teil der Daten liegen T-Werte vor,  
Wobei T unverzerrt für X ist.

i	$Y_i$	$T_i$	$X_i^*$	$Z_i$
1	$Y_1$	$T_1$	$X_1^*$	$Z_1$
⋮	⋮	⋮	⋮	⋮
k	$Y_k$	$T_k$	$X_k^*$	$Z_k$
$\dots$	$Y_{k+1}$	NA	$X_{k+1}^*$	$Z_{k+1}$
⋮	⋮	⋮	⋮	⋮
n	$Y_n$	NA	$X_n^*$	$Z_n$

- Schritt 1: Rechne für  $i \in \{1, \dots, k\}$  eine Regression  $T \sim X^* + Z$   
 $\rightarrow E[T|X^*, Z] = m_T(X^*, Z, \gamma) = E[X|X^*, Z]$   
 Berechne  $\hat{X}_i = (\underline{1}, \underline{X}_i^*, \underline{Z}_i) \hat{\gamma}$  für  $i \in \{k+1, \dots, n\}$

# Instrumentaldaten

## Schritt 1 & 2

Für einen Teil der Daten liegen T-Werte vor,  
Wobei T unverzerrt für X ist.

i	$Y_i$	$T_i$	$X_i^*$	$Z_i$	$\hat{X}_i$
1	$Y_1$	$T_1$	$X_1^*$	$Z_1$	NA
⋮	⋮	⋮	⋮	⋮	⋮
k	$Y_k$	$T_k$	$X_k^*$	$Z_k$	NA
k+1	$Y_{k+1}$	NA	$X_{k+1}^*$	$Z_{k+1}$	$\hat{X}_{k+1}$
⋮	⋮	⋮	⋮	⋮	⋮
n	$Y_n$	NA	$X_n^*$	$Z_n$	$\hat{X}_n$

- Schritt 1: Rechne für  $i \in \{1, \dots, k\}$  eine Regression  $T \sim X^* + Z$

$$\rightarrow E[T|X^*, Z] = m_T(X^*, Z, \gamma) = E[X|X^*, Z]$$

Berechne  $\hat{X}_i = (\underline{1}, \underline{X}_i^*, \underline{Z}_i) \hat{\gamma}$  für  $i \in \{k+1, \dots, n\}$

# Instrumentaldaten

## Schritt 1 & 2

Für einen Teil der Daten liegen T-Werte vor,  
Wobei T unverzerrt für X ist.

i	$Y_i$	$T_i$	$X_i^*$	$Z_i$	$\hat{X}_i$	$X_{reg_i}$	$VT_i$
1	$Y_1$	$T_1$	$X_1^*$	$Z_1$	NA	$T_1$	1
⋮	⋮	⋮	⋮	⋮	⋮	⋮	⋮
k	$Y_k$	$T_k$	$X_k^*$	$Z_k$	NA	$T_k$	1
⋮	⋮	⋮	⋮	⋮	⋮	⋮	⋮
k+1	$Y_{k+1}$	NA	$X_{k+1}^*$	$Z_{k+1}$	$\hat{X}_{k+1}$	$\hat{X}_{k+1}$	0
⋮	⋮	⋮	⋮	⋮	⋮	⋮	⋮
n	$Y_n$	NA	$X_n^*$	$Z_n$	$\hat{X}_n$	$\hat{X}_n$	0

- Schritt 1: Rechne für  $i \in \{1, \dots, k\}$  eine Regression  $T \sim X^* + Z$

$$\rightarrow E[T|X^*, Z] = m_T(X^*, Z, \gamma) = E[X|X^*, Z]$$

Berechne  $\hat{X}_i = (\underline{1}, \underline{X}_i^*, \underline{Z}_i) \hat{\gamma}$  für  $i \in \{k+1, \dots, n\}$



# Instrumentaldaten

## Schritt 1 & 2

Für einen Teil der Daten liegen T-Werte vor,  
Wobei T unverzerrt für X ist.

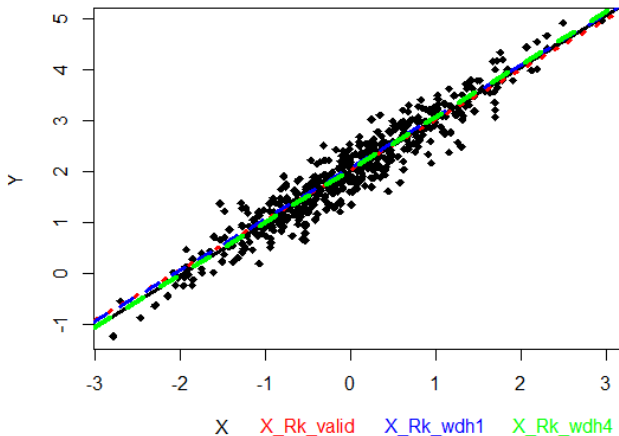
$i$	$Y_i$	$T_i$	$X_i^*$	$Z_i$	$\hat{X}_i$	$X_{reg_i}$	$VT_i$
1	$Y_1$	$T_1$	$X_1^*$	$Z_1$	NA	$T_1$	1
⋮	⋮	⋮	⋮	⋮	⋮	⋮	⋮
k	$Y_k$	$T_k$	$X_k^*$	$Z_k$	NA	$T_k$	1
⋮	⋮	⋮	⋮	⋮	⋮	⋮	⋮
k+1	$Y_{k+1}$	NA	$X_{k+1}^*$	$Z_{k+1}$	$\hat{X}_{k+1}$	$\hat{X}_{k+1}$	0
⋮	⋮	⋮	⋮	⋮	⋮	⋮	⋮
n	$Y_n$	NA	$X_n^*$	$Z_n$	$\hat{X}_n$	$\hat{X}_n$	0

- Schritt 2: Rechne für alle  $i$   $Y \sim X_{reg} + Z + VT$   
 $\rightarrow E[Y|X_{reg}, Z, VT] = m_Y(X_{reg}, Z, VT, \beta_{RK})$

Anhang Übersicht

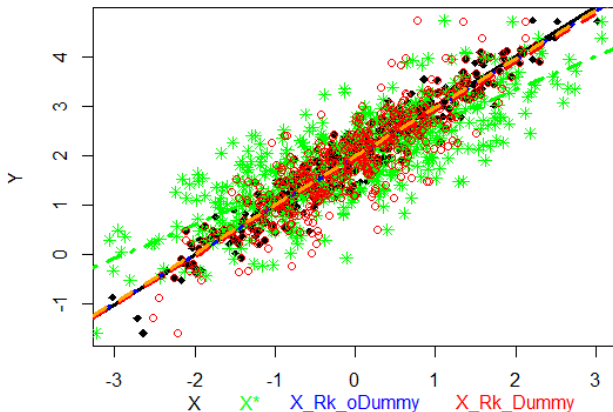
Vergleich: Valid  $\sim$  Wdh1  $\sim$  Wdh4

Kennzahl	Valid	Wdh1	Wdh4
$MSE(\beta_0)$	0.0009120709	1.149171367	1.0528207133
$MSE(\beta_1)$	0.0015110924	0.004843075	0.0019111721



# Vergleich: Validierungsdaten mit/ohne Dummy

Modell	ohne Dummy	mit Dummy
	$Y \sim \hat{X}$	$Y \sim \hat{X} + \text{Valid}(Ja/Nein)$
MSE( $\beta_0$ )	0.0010143196	0.9939534458
MSE( $\beta_1$ )	1.0529453460	0.0013049964



# Überprüfung der Schätzung in Schritt 1

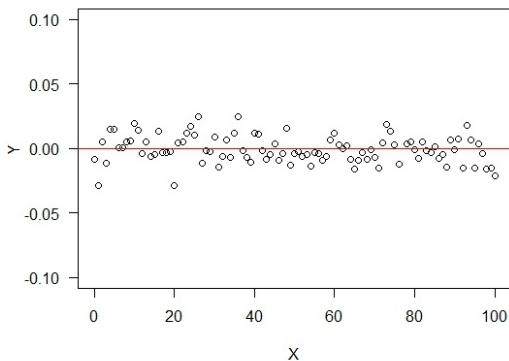
Schätzung von  $\mathbf{X}$  in Schritt 1 der RK kann überprüft werden

- durch gewöhnliche Regressionsdiagnosen bei
  - Validierungsdaten
  - Instrumentaldaten
- mithilfe von [partielle] Wiederholungsdaten
  - Wiederholungsdaten

# Überprüfung der Schätzung in Schritt 1

## Validierungsdaten & Instrumentaldaten

### Residuenplot



### Validierungsdaten

- $X \equiv$  wahre  $X_i$
- $Y$   
 $\equiv \epsilon_i = X_i - E[X_i|X_i^*, Z_i]$

### Instrumentaldaten

- $X \equiv T$
- $Y$   
 $\equiv \epsilon_i = T_i - E[T_i|X_i^*, Z_i]$

# Überprüfung der Schätzung in Schritt 1

[partielle] Wiederholungsdaten

Überprüfen der Schätzung in Schritt 1 der RK anhand von [partiellen] Wiederholungsdaten

$i$	$Y_i$	$X_{i1}^*$	$X_{i2}^*$
1	$Y_1$	$X_{11}^*$	$X_{12}^*$
$\vdots$	$\vdots$	$\vdots$	$\vdots$
$k$	$Y_k$	$X_{k1}^*$	$X_{k2}^*$
$\vdots$	$\vdots$	$\vdots$	$\vdots$
$\bar{k}+1$	$Y_{k+1}$	$X_{k+1,1}^*$	NA
$\vdots$	$\vdots$	$\vdots$	$\vdots$
$n$	$Y_n$	$X_{n1}^*$	NA

- Wie überprüft man ob das eine gute Schätzung für  $\mathbf{X}$  ist ohne die wahren  $\mathbf{X}$  beobachtet zu haben?

# Überprüfung der Schätzung in Schritt 1

[partielle] Wiederholungsdaten

Es liegen klassische Fehler vor

$$\begin{aligned} \mathbf{X}_{i1}^* &= \mathbf{X}_i + \mathbf{U}_{i1} \\ \mathbf{X}_{i2}^* &= \mathbf{X}_i + \mathbf{U}_{i2} \end{aligned}$$

und somit zunächst

$$\begin{aligned} E[\mathbf{X}_{i2}^* | \mathbf{Z}_i, \mathbf{X}_{i1}^*] &= E[\mathbf{X}_i + \mathbf{U}_{i2} | \mathbf{Z}_i, \mathbf{X}_{i1}^*] \\ &= \underbrace{E[\mathbf{X}_i | \mathbf{Z}_i, \mathbf{X}_{i1}^*]}_{\text{Schritt 1 der RK}} + \underbrace{E[\mathbf{U}_{i2} | \mathbf{Z}_i, \mathbf{X}_{i1}^*]}_{V_i}. \end{aligned}$$

Im klassischem Fehlermodell gilt für die bedingte Zufallsvariable  $V_i$   
 $E[V_i]=0$ .

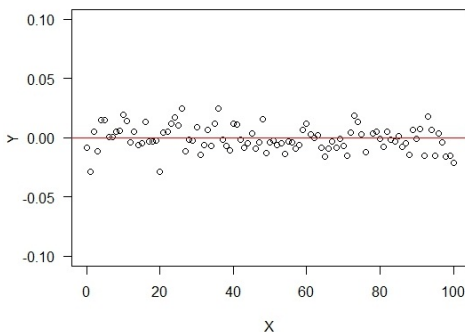
# Überprüfung der Schätzung in Schritt 1

[partielle] Wiederholungsdaten

$$\epsilon_i^* = \mathbf{X}_{i2}^* - E[\mathbf{X}_{i2}^* | \mathbf{Z}_i, \mathbf{X}_{i1}^*] =$$

$$(\mathbf{X}_i + \mathbf{U}_{i2}) - (E[\mathbf{X}_{i2} | \mathbf{Z}_i, \mathbf{X}_{i1}^*] + E[\mathbf{U}_{i2} | \mathbf{Z}_i, \mathbf{X}_{i1}^*])$$

$$(\mathbf{X}_i - E[\mathbf{X}_i | \mathbf{Z}_i, \mathbf{X}_{i1}^*]) + (\mathbf{U}_{i1} - E[\mathbf{U}_{i2} | \mathbf{Z}_i, \mathbf{X}_{i1}^*]) = \epsilon_i + \tilde{\epsilon}_i$$



[partielle] Wiederholungsdaten

- $X \equiv X_{i2}^* = X_i + U_{i2}$

- $Y \equiv \epsilon_i^* = \epsilon_i + \tilde{\epsilon}_i$

→ Aus Residuenplot

Tendenz ersichtlich ob

Schätzung in Schritt 1 der

RK eine gute Schätzung ist.

Anhang Übersicht