

Arten fehlender Werte - Konsequenzen und Möglichkeiten der Behandlung

O. Schoffer, mit Beispielen von C. Schneider

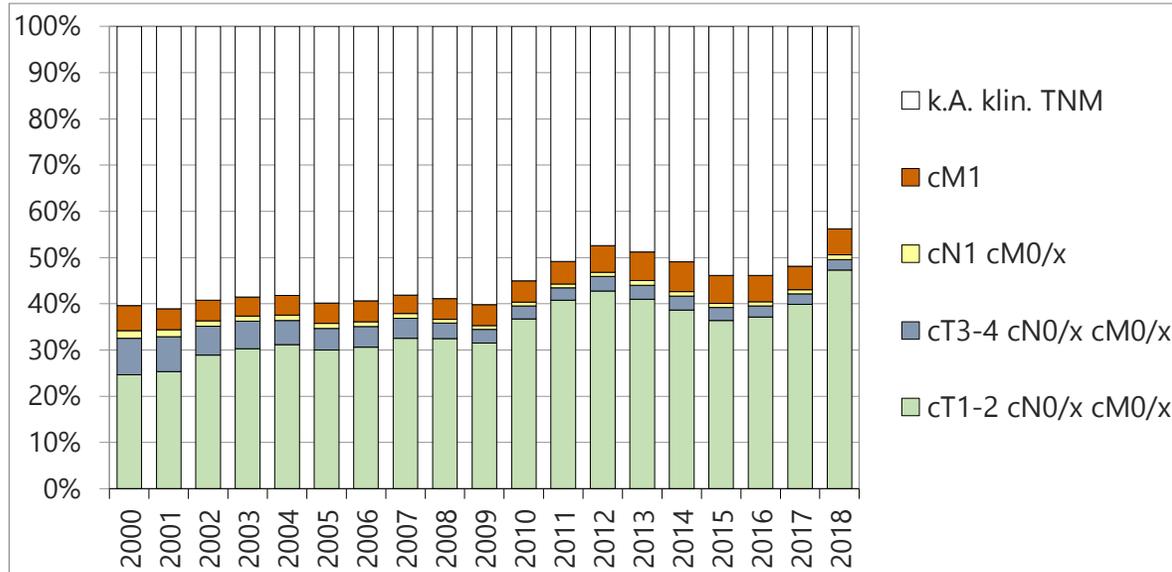
Vorbereitung zur 9. Bundesweiten Onkologischen Qualitätskonferenz (2022)

Workshop 2 - Auswertung und Methodik (22.04.2021)



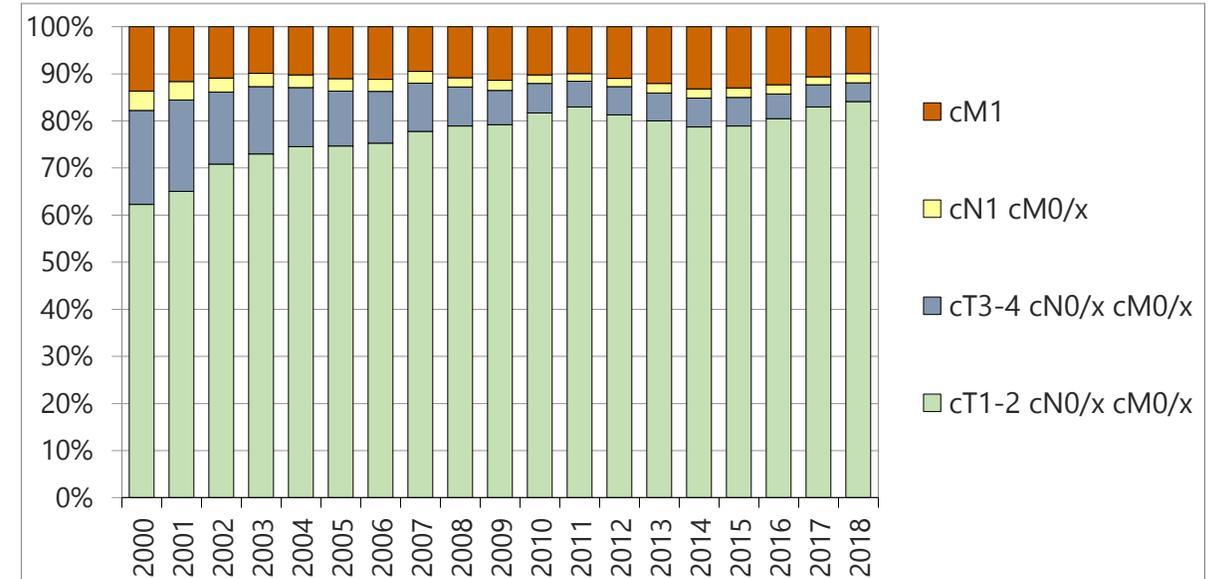
- Beispiele
- Warum sind fehlende Werte überhaupt ein Problem?
- Wie kann man mit fehlenden Werten umgehen?
- Was passiert, wenn ich fehlende Werte in der Analyse ignoriere?
- Wie erkenne ich, in welcher Situation ich bin?
- *Gegebenenfalls Exkurs: (multiple) Imputation*

Diagnosejahre 2000-2018



Einteilung der Prostatakarzinome auf Basis des klin. TNM, nach Diagnosejahr, n=368.269

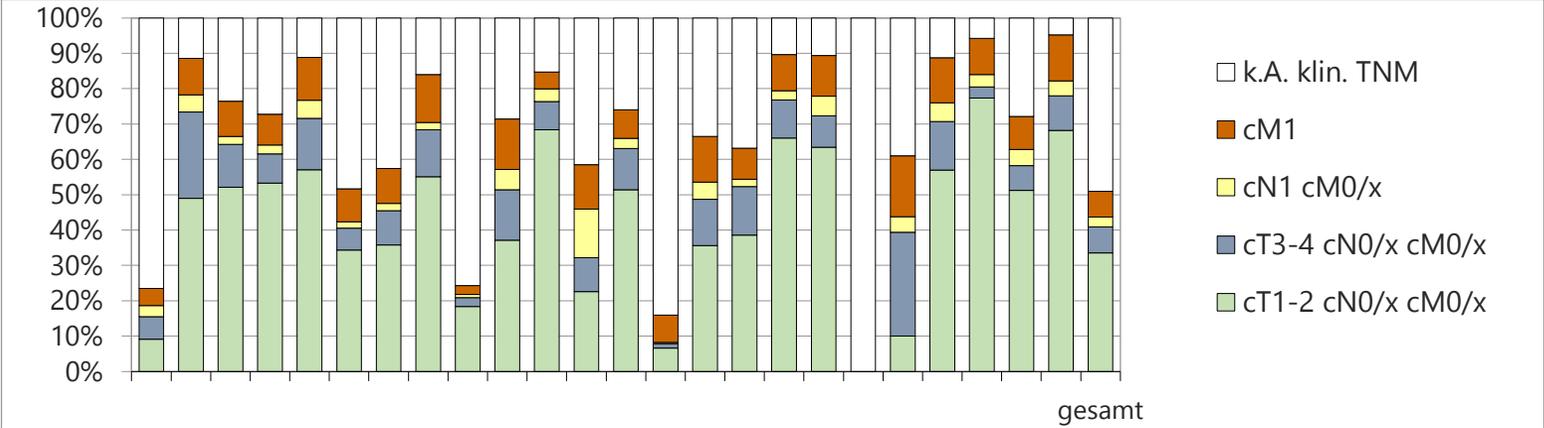
➔ fehlende Angabe, n=200.973 (54,6 %)



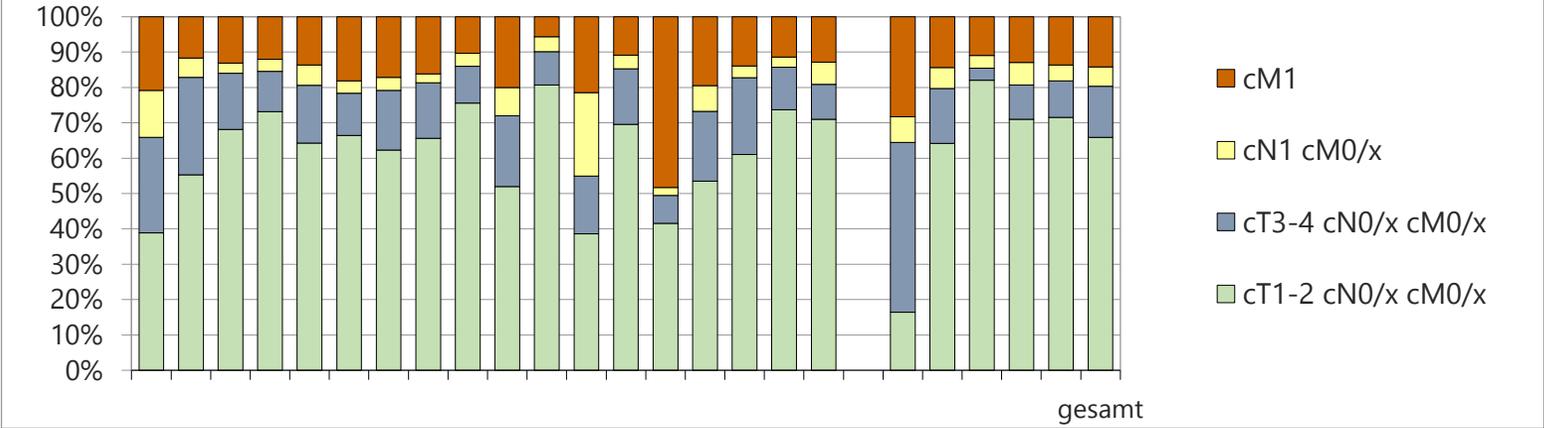
Einteilung der Prostatakarzinome auf Basis des klin. TNM, nach Diagnosejahr, ohne fehlende Angabe, n=167.296

Beispiele: Einteilung der Prostatakarzinome auf Basis der klinischen T-, N-, M-Kategorie

Diagnosejahre 2010-2018 - nach Registern (Reihenfolge zufällig und auf den Folien unterschiedlich)



Einteilung der Prostatakarzinome auf Basis klin. TNM, n=207.322



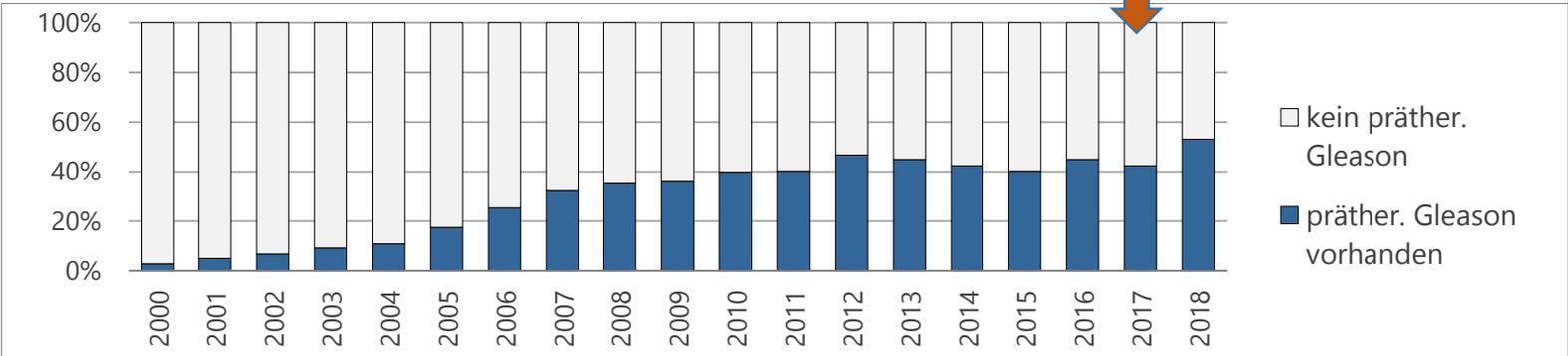
Einteilung der Prostatakarzinome auf Basis klin. TNM, ohne fehlende Angabe, n=105.616



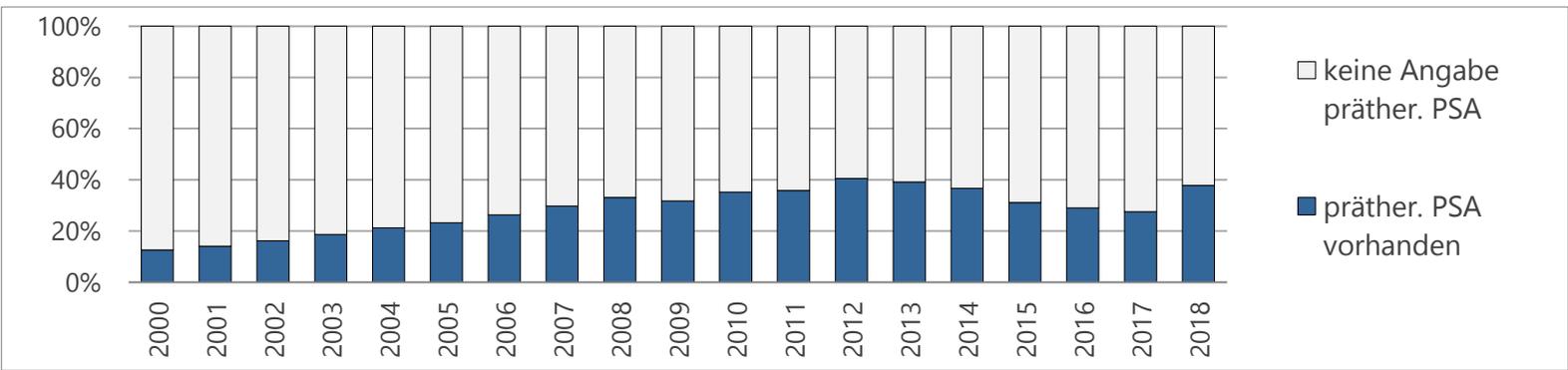
Beispiele: Prostatakarzinom - wichtige Angaben: Gleason und PSA-Wert

Diagnosejahre 2000-2018

Modul Prostatakarzinom 08/2017



Prostatakarzinome, Angabe eines prätherapeutischen Gleason, nach Diagnosejahr, n=368.269

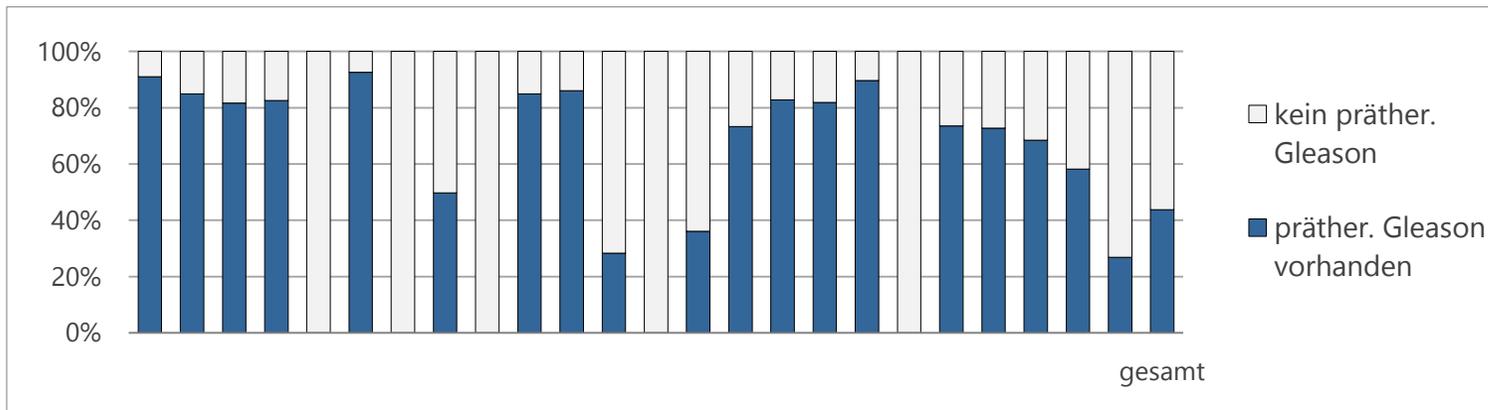


Prostatakarzinome, Angabe eines prätherapeutischen PSA-Wertes, nach Diagnosejahr, n=368.269

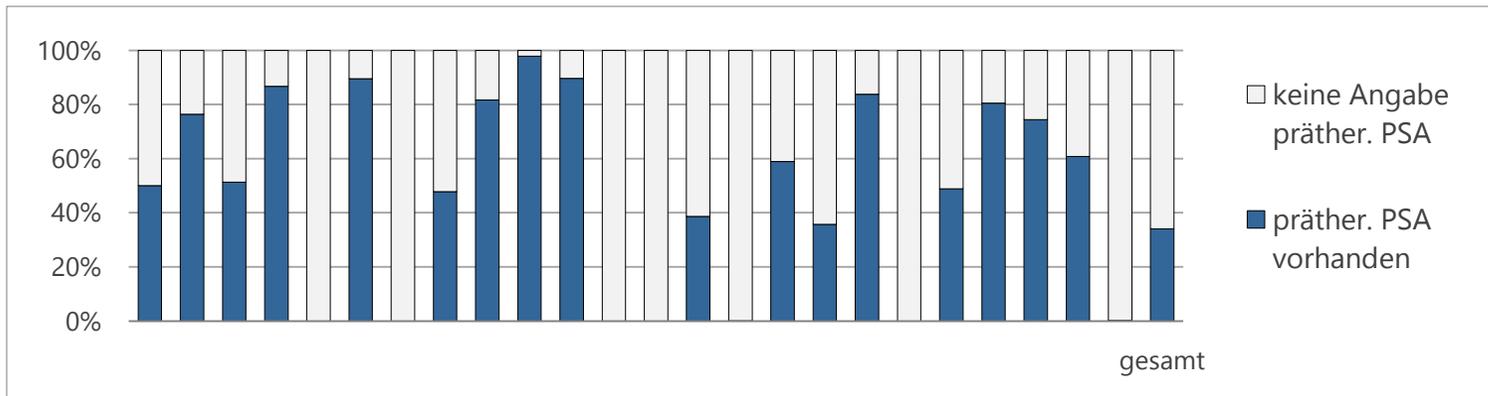


Beispiele: Prostatakarzinom - wichtige Angaben: Gleason und PSA-Wert

Diagnosejahre 2010-2018 - nach Registern (Reihenfolge zufällig und auf den Folien unterschiedlich)



Prostatakarzinome, Angabe eines prätherapeutischen Gleason, n=207.322



Prostatakarzinome, Angabe eines prätherapeutischen PSA-Wertes, n=207.322

Beispiele: Problem der Fallreduzierung durch fehlende Angaben

- lokal begrenztes PCa: cT1-2 cN0/x cM0/x
- lokal fortgeschrittenes PCa: cT3-4 cN0/x cM0/x
- regionär metastasiertes PCa: cN1 cM0/x
- metastasiertes PCA: cM1



Risikoeinteilung nach D'Amico auf Basis TNM, PSA, Gleason



niedriges/mittleres/hohes Risiko



unterschiedliche Therapieempfehlungen

Diagnosejahre 2010-2018:

207.322 Prostatakarzinome

101.754 Prostatakarzinome

**82.774 (81,3%)
lokal begrenzte Prostatakarzinome**

**45.715 Prostatakarzinome mit
Risikoeinschätzung**

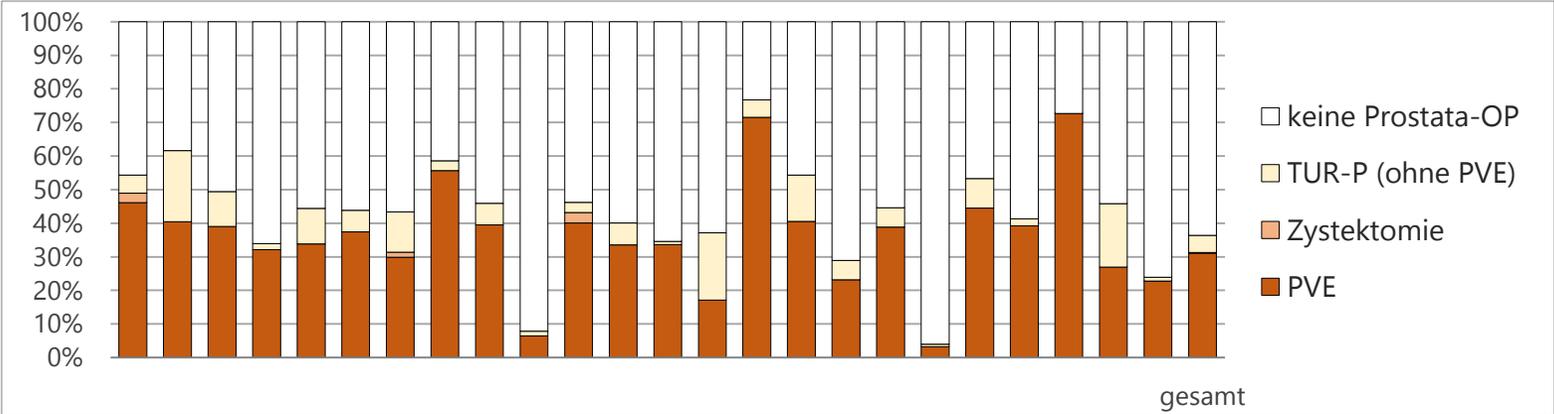
**Fehlendes klin. TNM:
105.568 (50,9%)**

**ohne PSA oder ohne Gleason:
37.059 (44,8 %)**

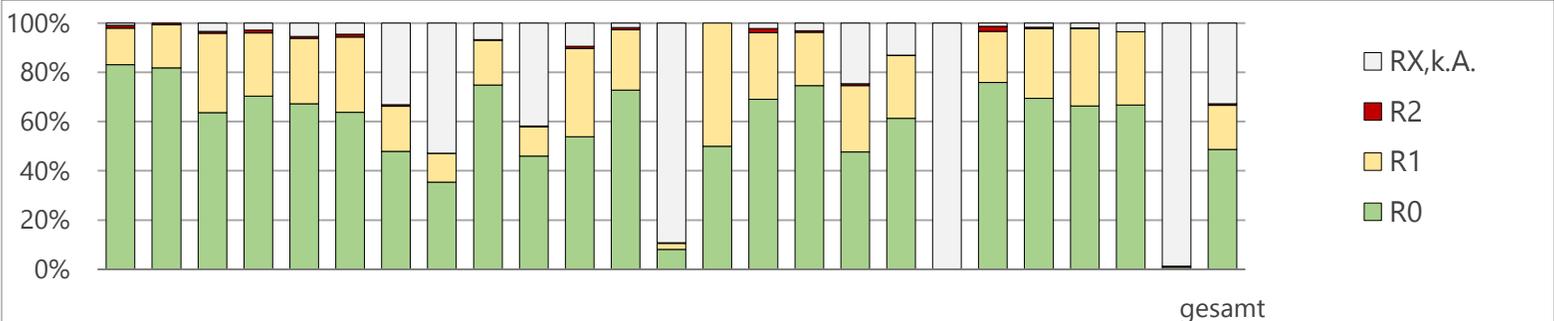
↔ geschätzte Anzahl bei vollständiger Meldung: 168.650



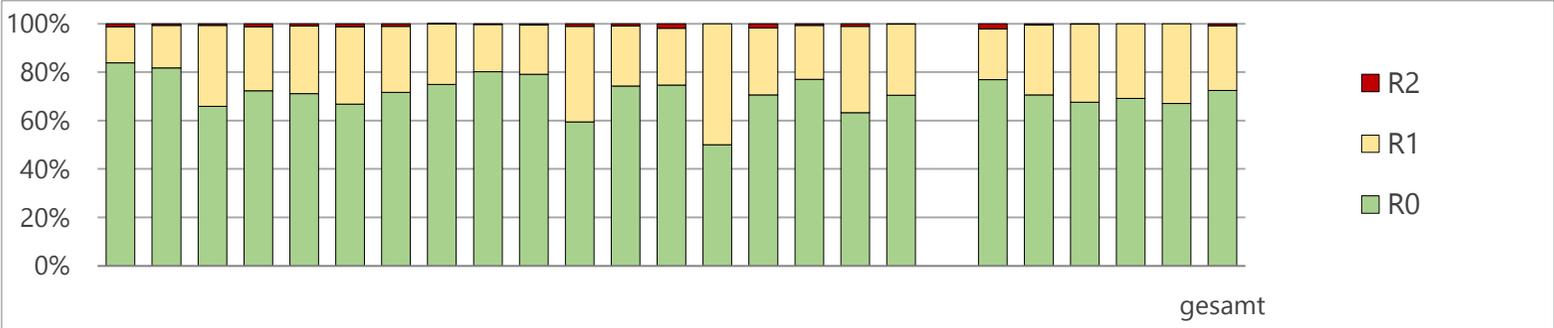
Beispiele: Operative Therapie, Diagnosejahre 2010-2018 – nach Registern



Operative Therapie der Prostatakarzinome, 2010-2018, n=207.322



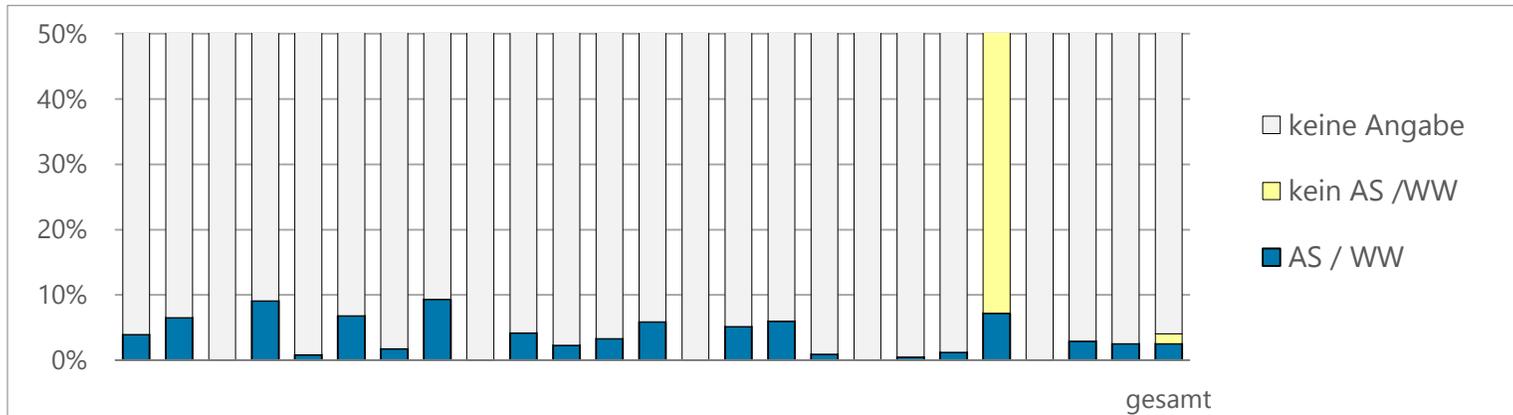
R-Klassifikation bei Prostatakarzinomen mit PVE/Zystektomie, n=64.848



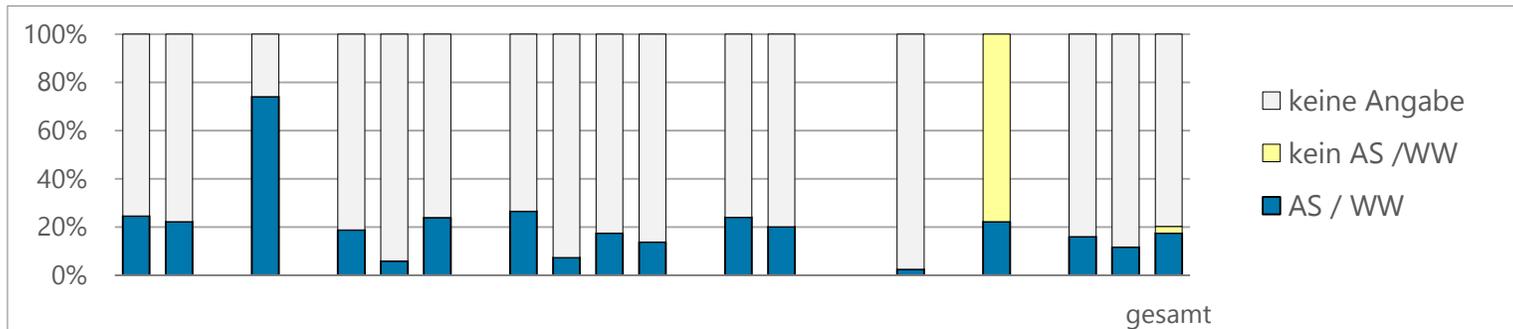
R-Klassifikation bei Prostatakarzinomen mit PVE/Zystektomie, Ausschluss von fehlenden Werten, n=43.585



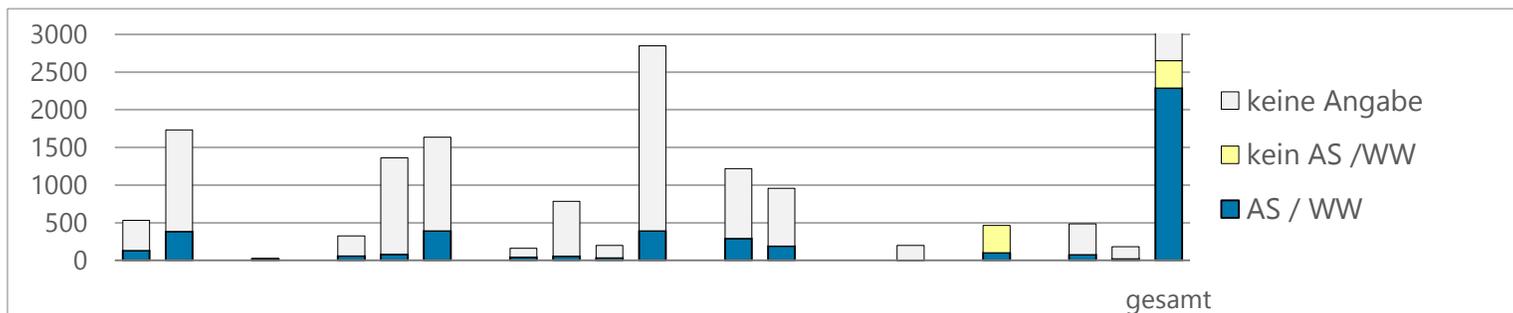
Beispiele: Überwachungsstrategie, Diagnosejahre 2010-2018 – nach Registern



Überwachungsstrategie bei Prostatakarzinomen (keine Einschränkung auf lokal begrenztes Prostatakarzinom mit niedrigem Risiko), n=207.322

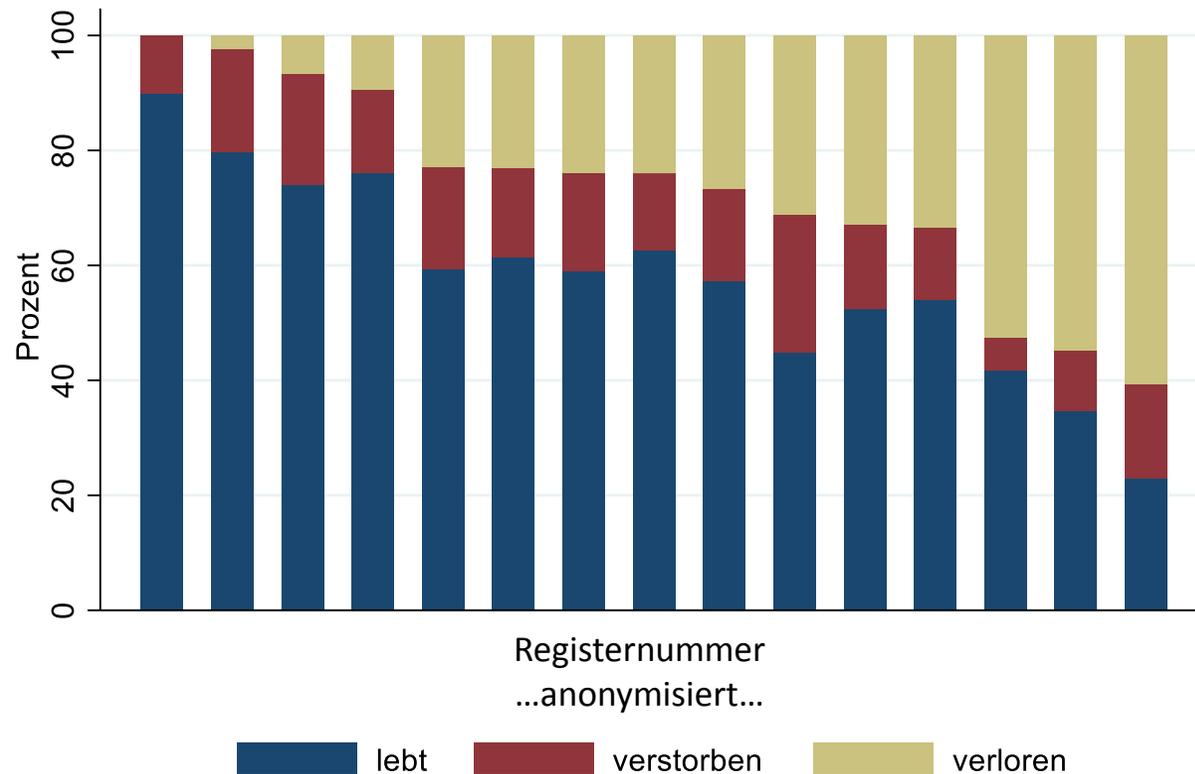


Überwachungsstrategie bei lokal begrenztes Prostatakarzinom mit niedrigem Risiko, n=13.116



Beispiele: Lost-to-follow-up

- Daten zum malignen Melanom 2000-2016
- Definition „verloren“:
 - kein Sterbedatum und
 - kein nach 01.01.2016 liegendes Lebenddatum und
 - Follow-up nach Diagnosestellung < 5 Jahre

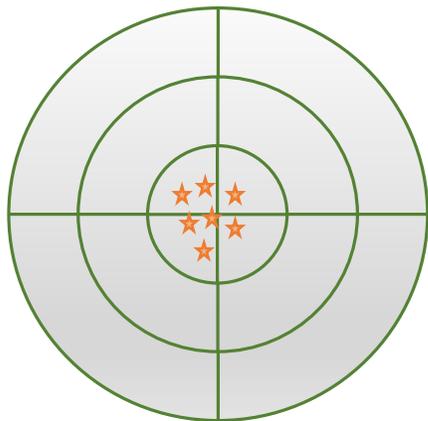


Beispiele: Lessons learned

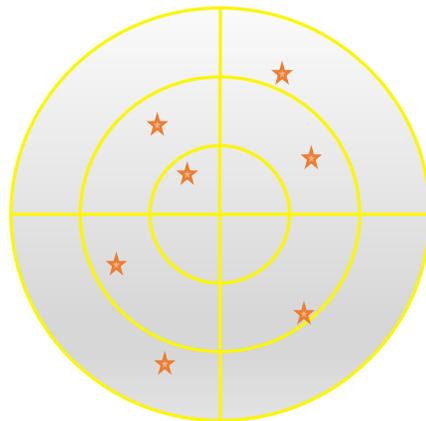
- Trends werden mit/ohne Abbildung fehlender Werte unterschiedlich sichtbar
- Ob die Verteilung von Ausprägungen ohne fehlende Werte einen Rückschluss auf die Verteilung der Gesamtheit (Verallgemeinerung) erlaubt, ist fraglich
- Je höher der Anteil fehlender Werte, desto unsicherer wird die Verallgemeinerung
- "keine Therapie" ist kein Meldeanlass, daher ist keine Unterscheidung zwischen "nicht angewendet" und "angewendet, aber nicht durchgeführt" möglich
- Zwischen Registern gibt es teilweise deutliche Unterschiede → Gründe können in Melderstruktur, gesetzlichen Grundlagen, Software etc. liegen
- Einige interessierende Merkmale außerhalb des Basisdatensatzes werden mitunter von Registern generell nicht dokumentiert

Warum sind fehlende Werte überhaupt ein Problem?

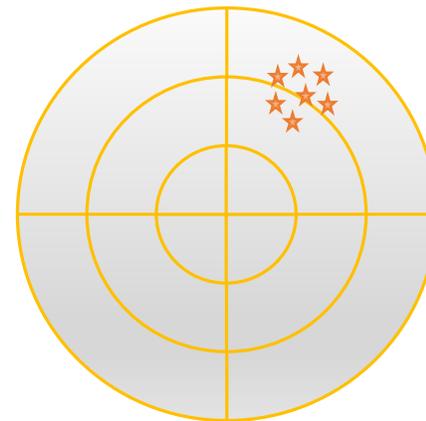
- Reduzierte Fallzahl
 - Reduzierte Schätzgenauigkeit (höhere Schätzvarianz)
 - In multivariaten Analysen können fehlende Werte für verschiedene Variablen und an verschiedenen Positionen besonders zu starker Fallzahlreduktion führen
- Generalisierbarkeit
 - Problematisch, falls fehlende Werte nicht zufällig, sondern nach einer Systematik zustande kommen: Gefahr verzerrter Schätzungen (Bias: <https://catalogofbias.org>)



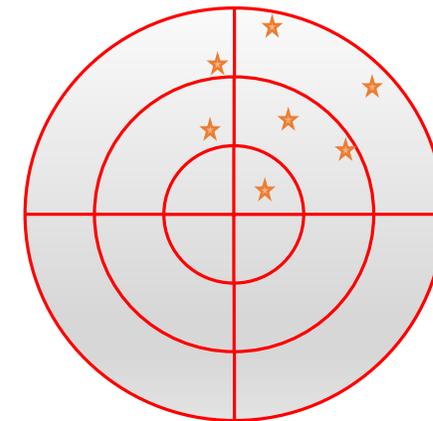
optimal



hohe
Schätzvarianz



Bias



Bias + hohe
Schätzvarianz

Wie kann man mit fehlenden Werten umgehen?

- MCAR (**missing completely at random**)
 - Fehlende Werte einer Variablen entstehen komplett zufällig
 - Struktur fehlender Werte hat nichts mit Ausprägungen der betreffenden Variablen zu tun und
 - Struktur fehlender Werte hängt von keinen anderen Variablen ab
 - Beispiel
 - Angaben zum Wohnort gehen unsystematisch aufgrund von Eingabefehlern verloren
 - Ergebnisse sind generalisierbar
 - Analyse unter Ausschluss der fehlenden Werte liefert unverzerrte Ergebnisse
 - insbesondere complete case analysis: Ausschluss aller Beobachtungen/Individuen mit mindestens einem fehlenden Wert in den analyserelevanten Variablen
 - Übertragung der Aussagen auf Grundgesamtheit ist möglich
 - Fallzahl muss trotzdem im Blick behalten werden
 - Realisationen unverzerrter Punktschätzer können „weit“ vom wahren Wert entfernt sein, insbesondere wenn aufgrund kleiner Fallzahl eine hohe Schätzvarianz gegeben ist



Wie kann man mit fehlenden Werten umgehen?

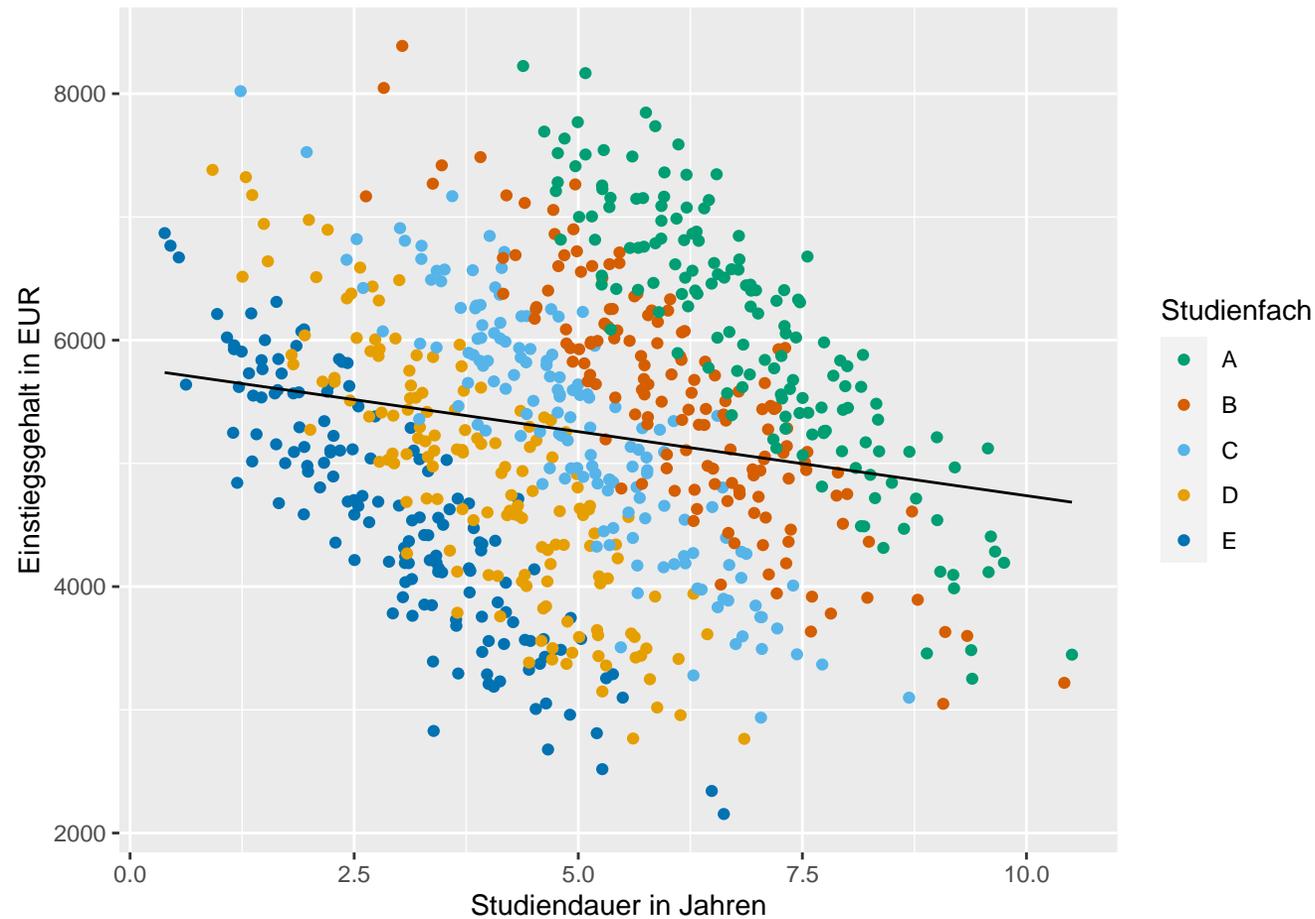
- MAR (**missing at random**)
 - Fehlende Werte einer Variablen entstehen zufällig innerhalb einer bekannten Struktur
 - Struktur fehlender Werte hat nichts mit Ausprägungen der betreffenden Variablen zu tun und
 - Struktur fehlender Werte hängt von anderen Variablen *im Datensatz* (beobachtete Variablen) ab
 - Beispiel:
 - Angaben zu Therapie und Vitalstatus werden in Zentren häufiger dokumentiert als in anderen Einrichtungen und Zentrumsstatus ist bekannt
 - Ergebnisse sind generalisierbar, wenn Einflussgrößen auf die Struktur fehlender Werte adäquat berücksichtigt wurden
 - (multiple) Imputation [**→ Exkurs**]
 - Stratifikation
 - Ggf. auch Adjustierung

Wie kann man mit fehlenden Werten umgehen?

- **NMAR (not missing at random)**
 - Fehlende Werte einer Variablen entstehen systematisch, aber nicht durch eine bekannte Struktur erklärbar
 - Struktur fehlender Werte hat mit Ausprägungen der betreffenden Variablen zu tun oder/und
 - Struktur fehlender Werte hängt von anderen Variablen *außerhalb des Datensatzes* (unbeobachtete Variablen) ab
 - Beispiele
 - UICC-Stadium wird für invasive Tumoren häufiger dokumentiert als für in situ
 - Dokumentation des Allgemeinzustands je nach (unbeobachtetem) Bildungsstatus
 - Keine Generalisierbarkeit \Rightarrow Gefahr von Bias
 - Ggf. Analyse mit fehlenden Werten als eigene Kategorie
 - Benennung als Limitation in Ergebnisdarstellung und –interpretation (\rightarrow „Diskussion“ in wiss. Beiträgen)
 - Sensitivitätsanalysen

Was passiert, wenn ich fehlende Werte in der Analyse ignoriere?

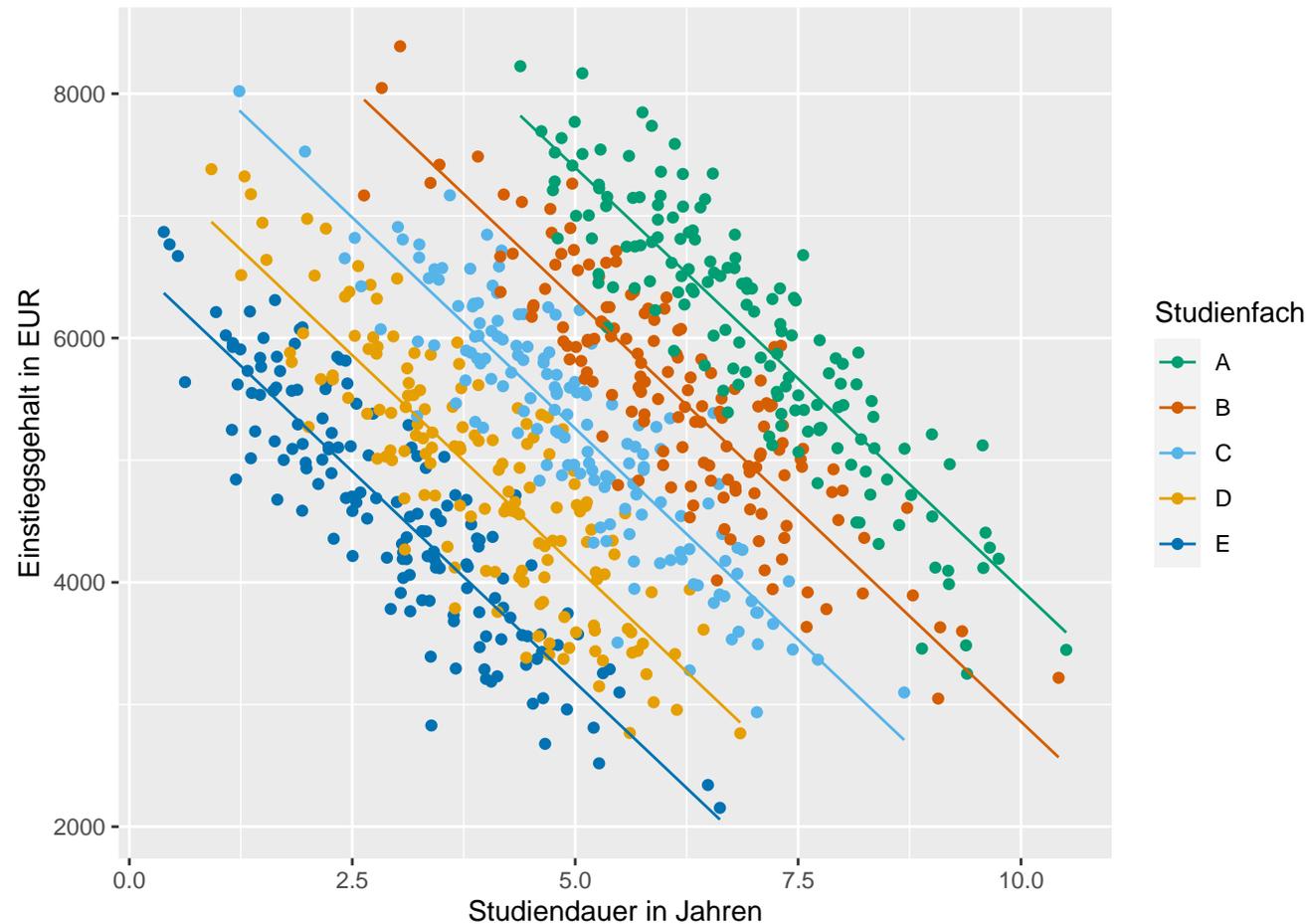
- Simulation: Variante des Simpson-Paradoxons



- Werden relevante Informationen nicht genutzt, entstehen irreführende Schätzungen

Was passiert, wenn ich fehlende Werte in der Analyse ignoriere?

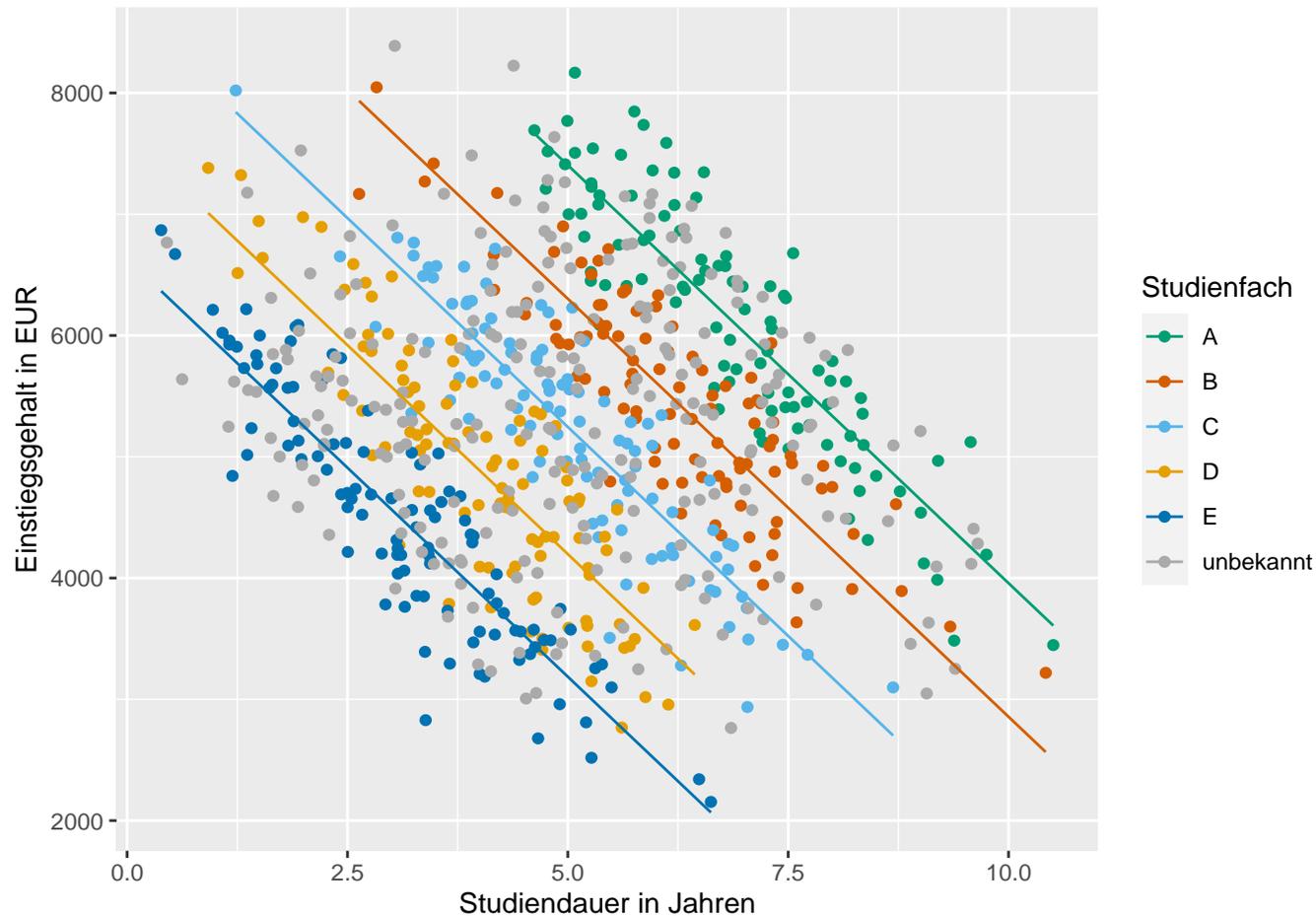
- Simulation: Variante des Simpson-Paradoxons



- Bei Einbeziehung aller relevanten Informationen ist eine valide Schätzung möglich

Was passiert, wenn ich fehlende Werte in der Analyse ignoriere?

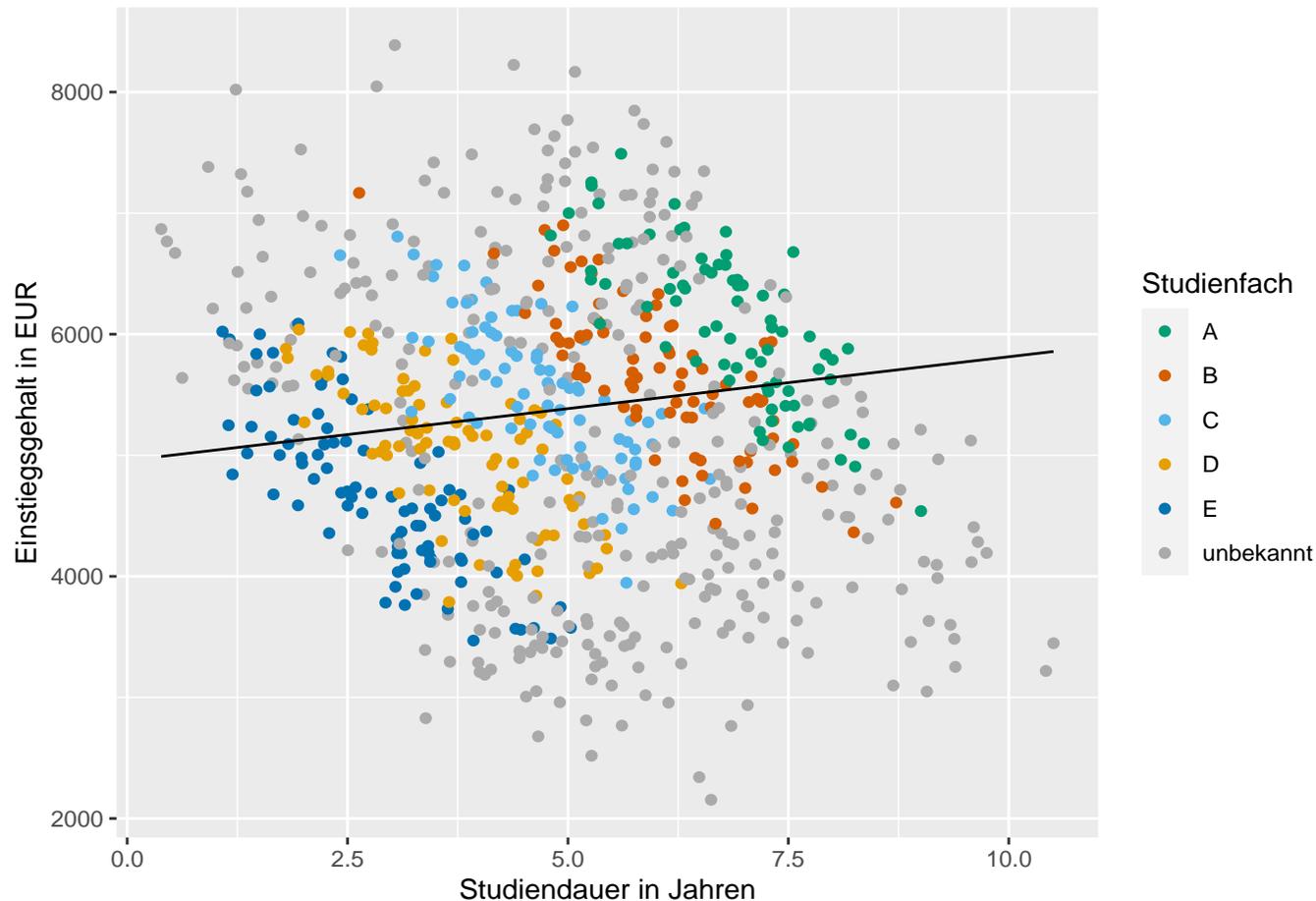
- Simulation: Variante des Simpson-Paradoxons
 - Fehlende Werte entstehen zufällig (MCAR)



- Die Schätzung ist auch bei fehlenden Werten valide möglich, sofern
 - die fehlenden Werte keiner Systematik unterliegen und
 - noch hinreichend viele Werte zur Schätzung zur Verfügung stehen

Was passiert, wenn ich fehlende Werte in der Analyse ignoriere?

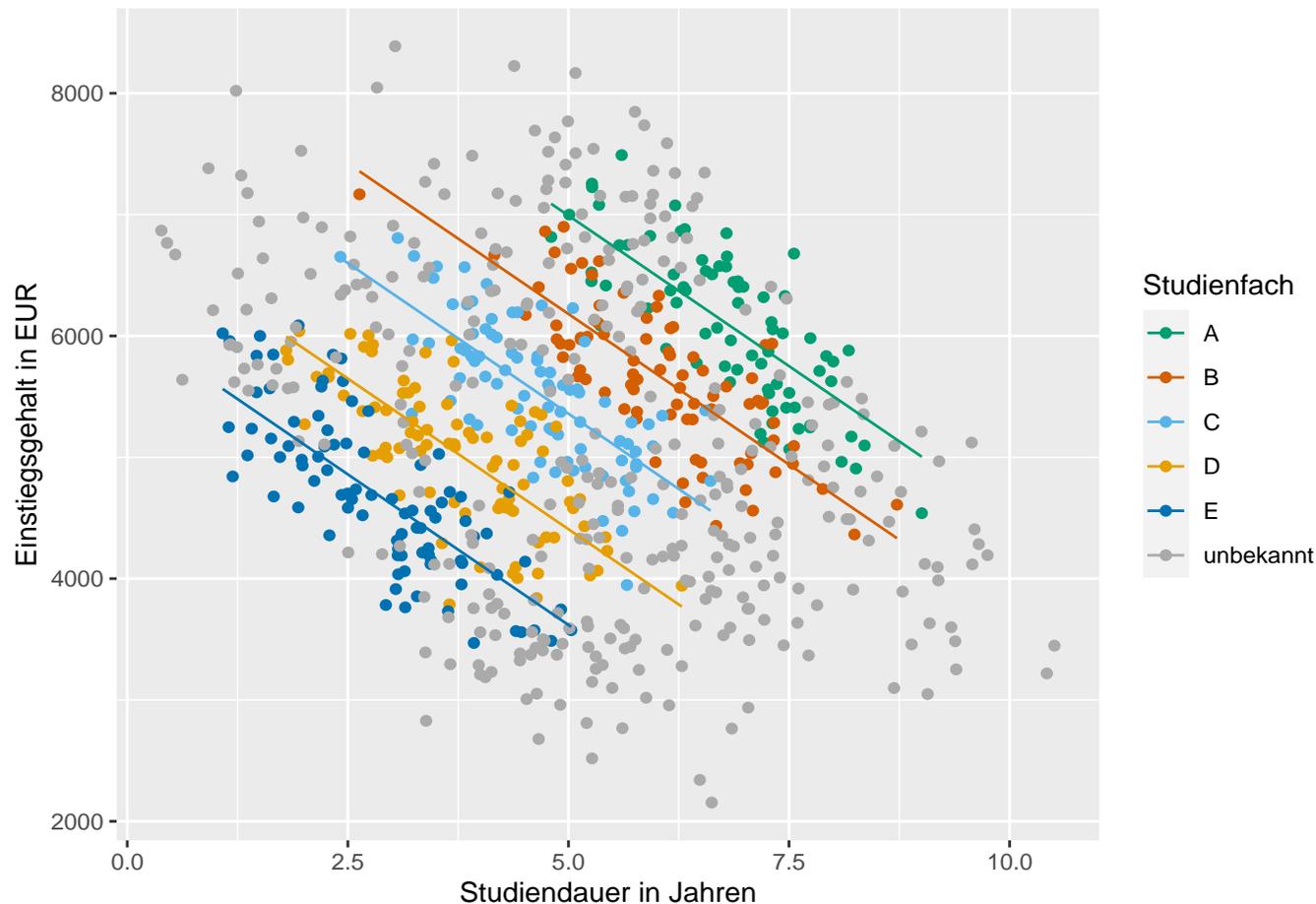
- Simulation: Variante des Simpson-Paradoxons
 - Fehlende Werte entstehen systematisch (MAR oder NMAR)



- Die Schätzung kann verzerrt (irreführend) sein, wenn
 - die fehlenden Werte einer Systematik unterliegen
 - selbst wenn noch hinreichend viele Werte zur Schätzung zur Verfügung stehen

Was passiert, wenn ich fehlende Werte in der Analyse ignoriere?

- Simulation: Variante des Simpson-Paradoxons
 - Fehlende Werte entstehen systematisch (MAR oder NMAR)



- Mitunter kann die Situation verbessert werden, wenn
 - die Systematik der fehlenden Werte erkannt und
 - diese Information in die Schätzung einbezogen wird,
 - jedoch kann auch dann noch Verzerrung verbleiben, falls die Systematik nicht vollständig aufgeklärt werden konnte

Was passiert, wenn ich fehlende Werte in der Analyse ignoriere?

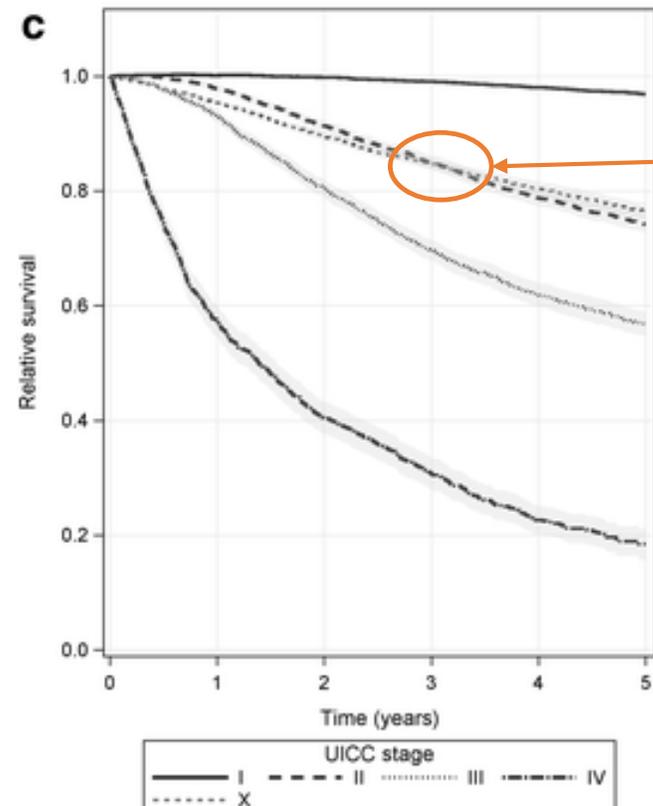
- Kann man in den Situationen MAR und NMAR abschätzen, welche Auswirkungen fehlende Werte auf die Analyse haben?
 - Anteil der fehlenden Werte beachten
 - ein kleiner Anteil fehlender Werte kann zwar Bias nicht verhindern, aber bestenfalls wird der Effekt auf Analyseergebnisse irrelevant klein → qualitative Aussage unbeeinflusst
 - Sensitivitätsanalysen z.B. wie in Hellmund et al. (2020)
 - <https://www.mdpi.com/2072-6694/12/9/2354> → Supplement S2

Cutoff value: maximum share of no documented therapies per registry & year (%)	50	67	75
<i>n</i> (share of included patients)	1053 (27%)	1970 (51%)	2642 (68%)
	HR (95%-CI)	HR (95%-CI)	HR (95%-CI)
Age at metastasis	1.018 (1.012; 1.024)	1.013 (1.009; 1.017)	1.011 (1.007; 1.014)
Sex: female (ref: male)	0.924 (0.791; 1.080)	0.870 (0.781; 0.970)	0.884 (0.806; 0.970)
Region: Eastern Germany (ref: Western Germany)	1.739 (1.375; 2.200)	1.535 (1.331; 1.771)	1.559 (1.387; 1.751)

Quelle: Hellmund, Schmitt, Roessler, Meier, Schoffer: "Targeted and Checkpoint Inhibitor Therapy of Metastatic Malignant Melanoma in Germany, 2000-2016", *Cancers* 2020; 12(9).

Wie erkenne ich, in welcher Situation ich bin?

- NMAR kann „per definitionem“ nicht getestet werden, trotzdem...
 - Konzeptionelle Gedanken → Kenntnis des Meldeprozesses („Wie werden Daten generiert?“)
 - Gesunder Menschenverstand in der Auswertung



Proportionalitätsannahme für UICC X
deutlich verletzt (Kurven scheiden sich)

Quelle: Schoffer, Schüle, Arand, Arnholdt, Baaske, et al.:
“Tumour stage distribution and survival of malignant melanoma in
Germany 2002–2011”, *BMC Cancer* 2016; 16(1).

Wie erkenne ich, in welcher Situation ich bin?

- Pragmatisch wird oft von MAR ausgegangen
 - MCAR-Situation ist nur selten gegeben \Rightarrow sollte nicht unbegründet vorausgesetzt werden
 - Selbst in NMAR-Situation werden Techniken aus MAR-Situation angewendet, um durch Bereinigung der MAR-Situation „Schadensbegrenzung“ zu erreichen

Wie erkenne ich, in welcher Situation ich bin?

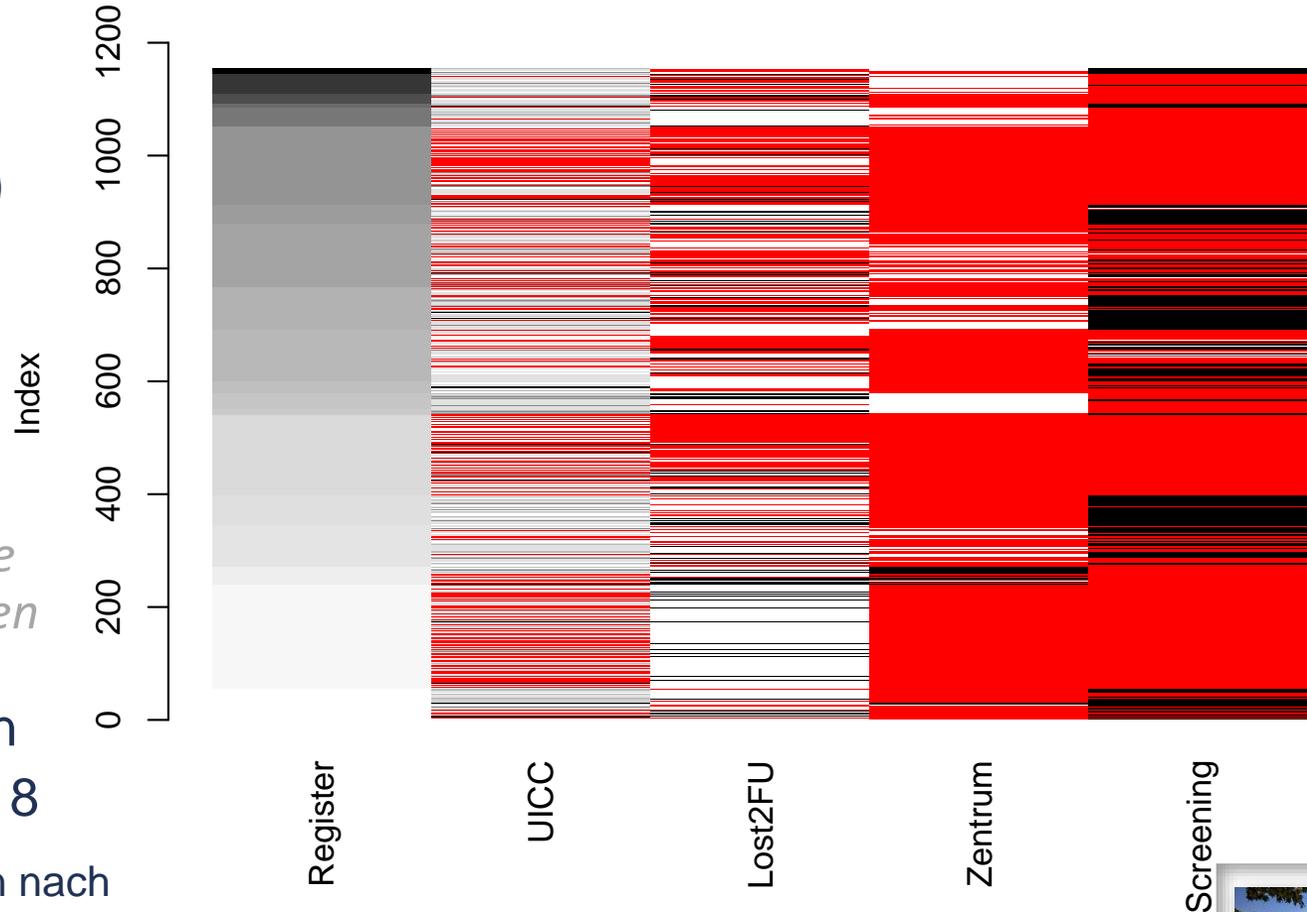
- Unterscheidung zwischen MAR und MCAR

- Analyse der Struktur fehlender Werte

- Visualisierung (Matrixplot)
 - *Hinweis 1: Eine „Zeile“ des Matrixplots steht für eine Beobachtung*
 - *Hinweis 2: fehlende Werte sind rot dargestellt, verschiedene nichtfehlende Ausprägungen in Grautönen*

Beispiel: 1%-Stichprobe der Daten zum malignen Melanom 2000-2018

⇒ Fehlende Werte für Screening clustern nach Registern



Wie erkenne ich, in welcher Situation ich bin?

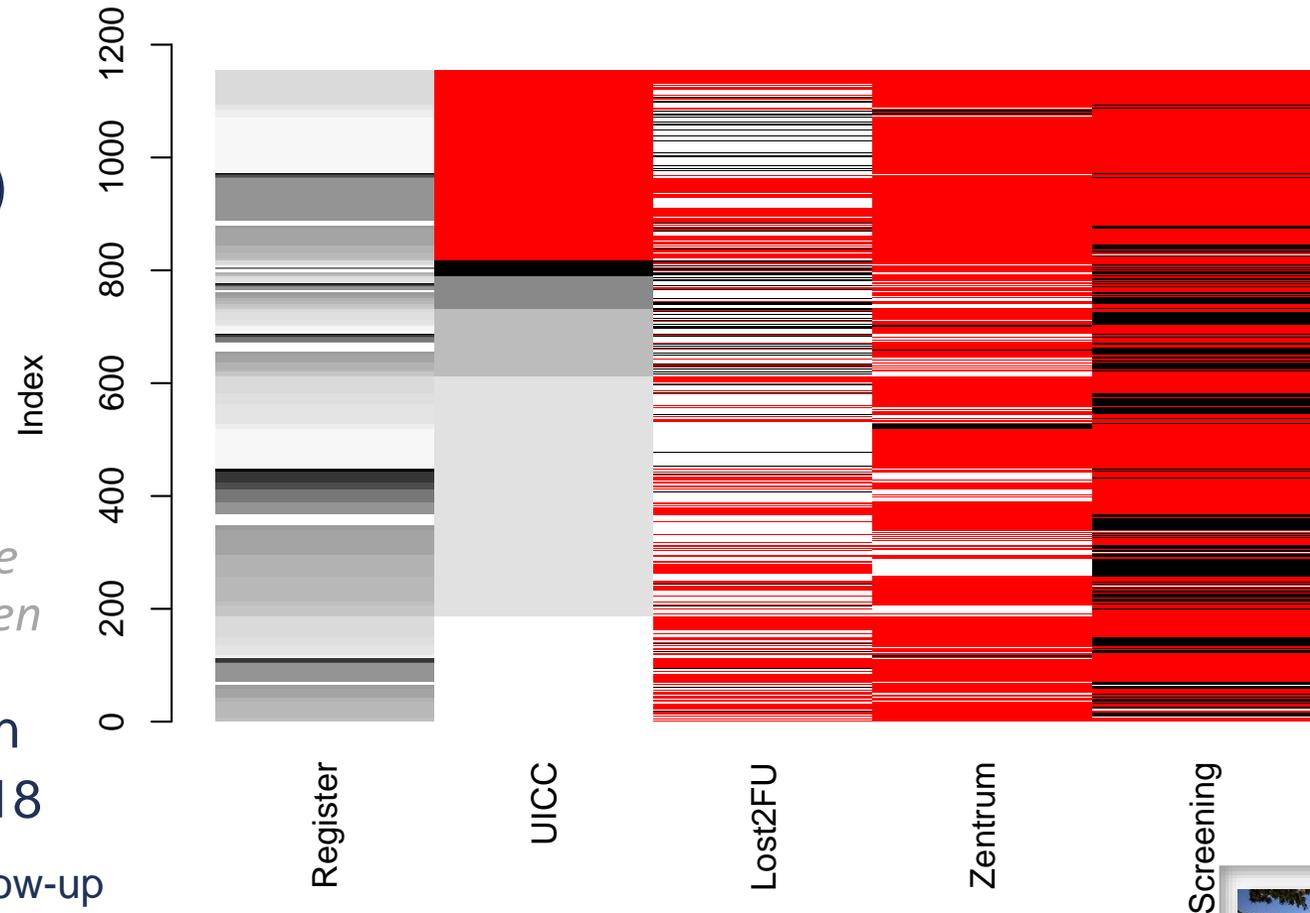
- Unterscheidung zwischen MAR und MCAR

- Analyse der Struktur fehlender Werte

- Visualisierung (Matrixplot)
 - *Hinweis 1: Eine „Zeile“ des Matrixplots steht für eine Beobachtung*
 - *Hinweis 2: fehlende Werte sind rot dargestellt, verschiedene nichtfehlende Ausprägungen in Grautönen*

Beispiel: 1%-Stichprobe der Daten zum malignen Melanom 2000-2018

⇒ Anteil fehlender Werte für Lost-to-Follow-up und Zentrumsbehandlung variiert je nach UICC-Stadium (insb. UICC X)



Wie erkenne ich, in welcher Situation ich bin?

- Unterscheidung zwischen MAR und MCAR
 - Analyse der Struktur fehlender Werte
 - Visualisierung (Matrixplot)
 - Zusammenhangsmaße (Korrelation etc.) bei Codierung „bekannt“ vs. „unbekannt“ mit anderen Variablen
 - Test nach Little

- „Naive“ Ersetzung
 - Last observation carried forward
 - Ersetzung durch Gesamtmittelwert/-median/-modalwert je nach Messniveau
- Ausnutzung der Kenntnis der Struktur fehlender Werte
 - (Gruppen-)Mittelwerte/Mediane/Modalwerte
 - Prädiktion (Punktschätzer) aus complete-case-Regressionsmodell



- Problem: Unterschätzung der Varianz
 - Ersetzte Werte haben keine oder nur minimale Varianz
 - Varianz, welche „wahre“ statt der fehlenden Werte hätten, bleibt unberücksichtigt

- Ersetzung nach (bedingten) Wahrscheinlichkeitsregeln
 - Nutzung einer Modellierung, welche zu Vorhersageverteilungen oder -wahrscheinlichkeiten statt zu Punktprognosen führt
 - Spezifisch für jede Parameterkonstellation bzw. Gruppe
 - Beispiel: Logistische Regression für dichotome Zielvariable
 - Mehrfache(!) Simulation von Zufallswerten aus den spezifischen Verteilungen
 - Datenauswertung wird für jeden Simulationsdurchlauf separat durchgeführt
 - Zusammenführung von Schätzergebnissen aus den Iterationen der Datenauswertung mittels Rubins Rule (Rubin, 1987)
 - Gesamtschätzung: Mittelwert der Einzelschätzungen
 - Gesamtvarianz: gepoolt aus Intra- und Inter-Iterations-Varianz
- 
- Ziel: Berücksichtigung der Varianz sowie Reduzierung der Schätzunsicherheit

Quelle: Rubin, 1987: *Multiple Imputation for Nonresponse in Surveys*. New York: John Wiley & Sons.

- Fehlende Werte können zu unsicheren oder sogar irreführenden Schlussfolgerungen führen

- ...wenn sie nicht komplett zufällig entstehen
- ...wenn sie einen relevanten Anteil der Beobachtungen betreffen

⇒ Ignorieren ist zumeist keine gute Idee

- Aber wir sind dem Problem nicht machtlos ausgeliefert

- Quantifizierung der Auswirkungen mittels Sensitivitätsanalysen
- Aufklärung der Abhängigkeiten für Entstehung von fehlenden Werten
- Ersetzung fehlender Werte → Ausnutzung der erkannten Abhängigkeiten



**Vielen Dank für Ihre
Aufmerksamkeit!**



Kontakt

olaf.schoffer@ukdd.de

constanze.schneider@kkrbb.de



A. Stang:

- in einem Regressionsmodell, in dem wir einen Indikator für die Missingness setzen, gibt es zwei unterschiedliche Bewertungen: sofern wir einen kausalen Effektschätzen wollen, ist es inakzeptabel; sofern nur individual Risk Prediction im Vordergrund steht, könnte es ein hilfreiches Vorgehen sein, sofern man auf Imputationen verzichtet.

A. Katalinic:

- Literaturhinweis: Imputation of missing values of tumour stage in population-based cancer registration
 - <https://bmcmedresmethodol.biomedcentral.com/articles/10.1186/1471-2288-11-129>

