

Datenqualität bei Primärdatenerhebungen – Konzept und Implementation

Workshop Datenqualität TMF, 03.05.2018

Carsten Oliver Schmidt

Universitätsmedizin Greifswald

ICM-SHIP-KEF

Funktionsbereich Qualität in der Gesundheitsforschung

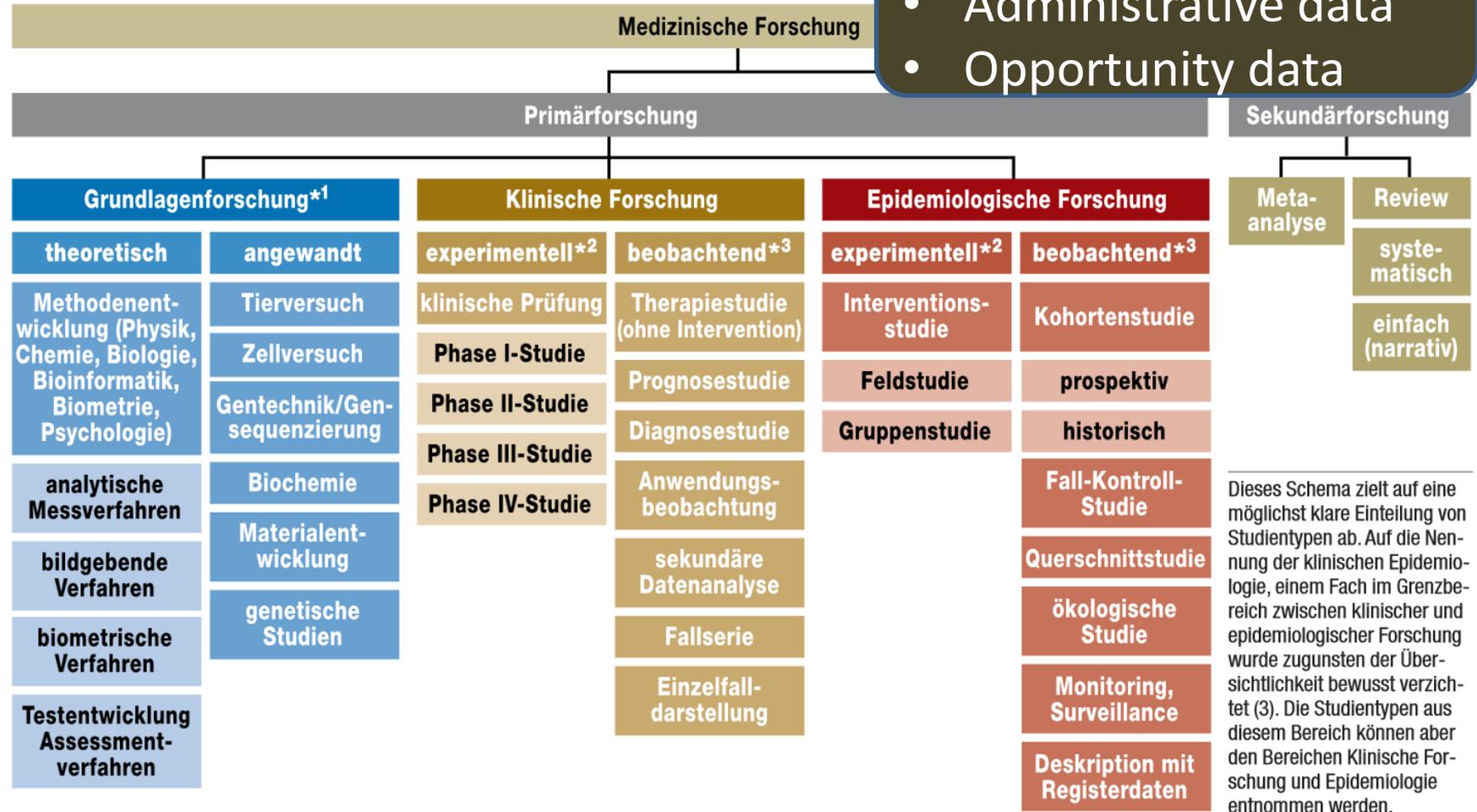
1. Anwendungsszenario und Datenqualität
2. Anwendungsbeispiel Study of Health in Pomerania
3. Konzept zur Bewertung von Datenqualität
4. Vom Konzept zur Statistik
5. Implementation in einer Kohortenstudie

- 1. Anwendungsszenario und Datenqualität**
2. Anwendungsbeispiel Study of Health in Pomerania
3. Konzept zur Bewertung von Datenqualität
4. Vom Konzept zur Statistik
5. Implementation in einer Kohortenstudie

Anwendungsszenarien für Datenqualitätsbewertungen

- Designed data
- Administrative data
- Opportunity data

GRAFIK 1



Einteilung verschiedener Studientypen

*1 häufig synonym verwendet: Experimentelle Forschung; *2 analoger Begriff: interventionell; *3 analoger Begriff: nicht interventionell/nicht experimentell

1. Anwendungsszenario und Datenqualität
- 2. Anwendungsbeispiel Study of Health in Pomerania**
3. Konzept zur Bewertung von Datenqualität
4. Vom Konzept zur Statistik
5. Implementation in einer Kohortenstudie

Erwachsene
20-79 Jahre



- Prävalenz bevölkerungsrelevanter Erkrankungen und Risikofaktoren
- Zusammenhänge zwischen Risikofaktoren, (sub-) klinischen Auffälligkeiten und Folgen analysieren

Wiederholte Messungen in SHIP

SHIP-0
Baseline

n= 4308

1997-2001

SHIP-1
5y Follow-up

n= 3300

2002-2006

SHIP-2
10y Follow-up

n= 2333

2008-2012

SHIP-3
15y Follow-up

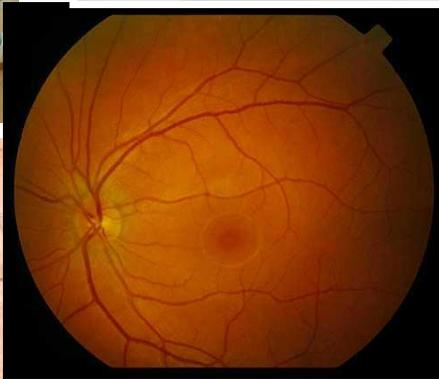
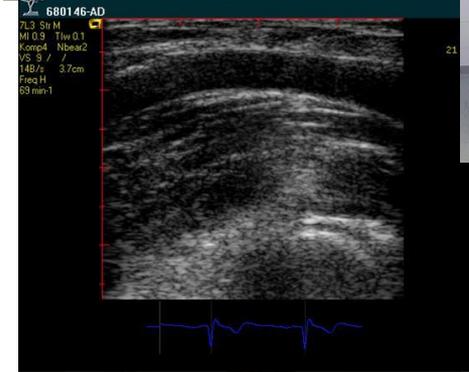
n= 1718

2014-2019

SHIP-TREND-0
Baseline

n= 4422

SHIP-TREND-1
5y Follow-up



Elemente des Qualitätsmanagements

Tabelle 1 Beispielhafte Elemente der Qualitätssicherung in Kohortenstudien

a. Maßnahmen vor der Datenerhebung

Studienhandbücher für Untersuchungen, Datenmanagement, QS

Checklisten zur Untersuchungsdurchführung

Training / (Re-) Zertifizierung von Untersuchern, Befundern und ggfs. Trainern bei multi-zentrischen Projekten

Gerätekalibrierung, Gerätevergleiche, Gerätewartung

Phantommessungen bei bildgebenden Verfahren

Gespräche/Qualitätszirkel zu qualitätsrelevanten Aspekten der Untersuchung

Pilotierung der Studie/ Prätest einzelner Untersuchungsmodulare

b. Maßnahmen während und nach der Datenerhebung

Kontrolle der Teilnahmevoraussetzungen (Identität und Einverständnisse)

Standardisierte Dateneingaben (z.B. webbasiert mit automatischer Plausibilitätskontrolle)

Mehrfachbefundungen (v.a. bei bildgebenden Verfahren)

Mehrfacheingaben (Standard bei Eingabe von Papiervorlagen, z.B. Fragebögen)

On-Site Monitoring (Räumlichkeiten sowie Untersuchungsprozess)

Standardisierte und zentralisierte Erfassung von Auffälligkeiten (z.B. datenbankbasiert)

Qualitätsberichte

Monitoring der laufenden Datenerhebung (Fokus auf Rohdaten)

Datenmanagement, syntaxbasierte Datenbereinigungen

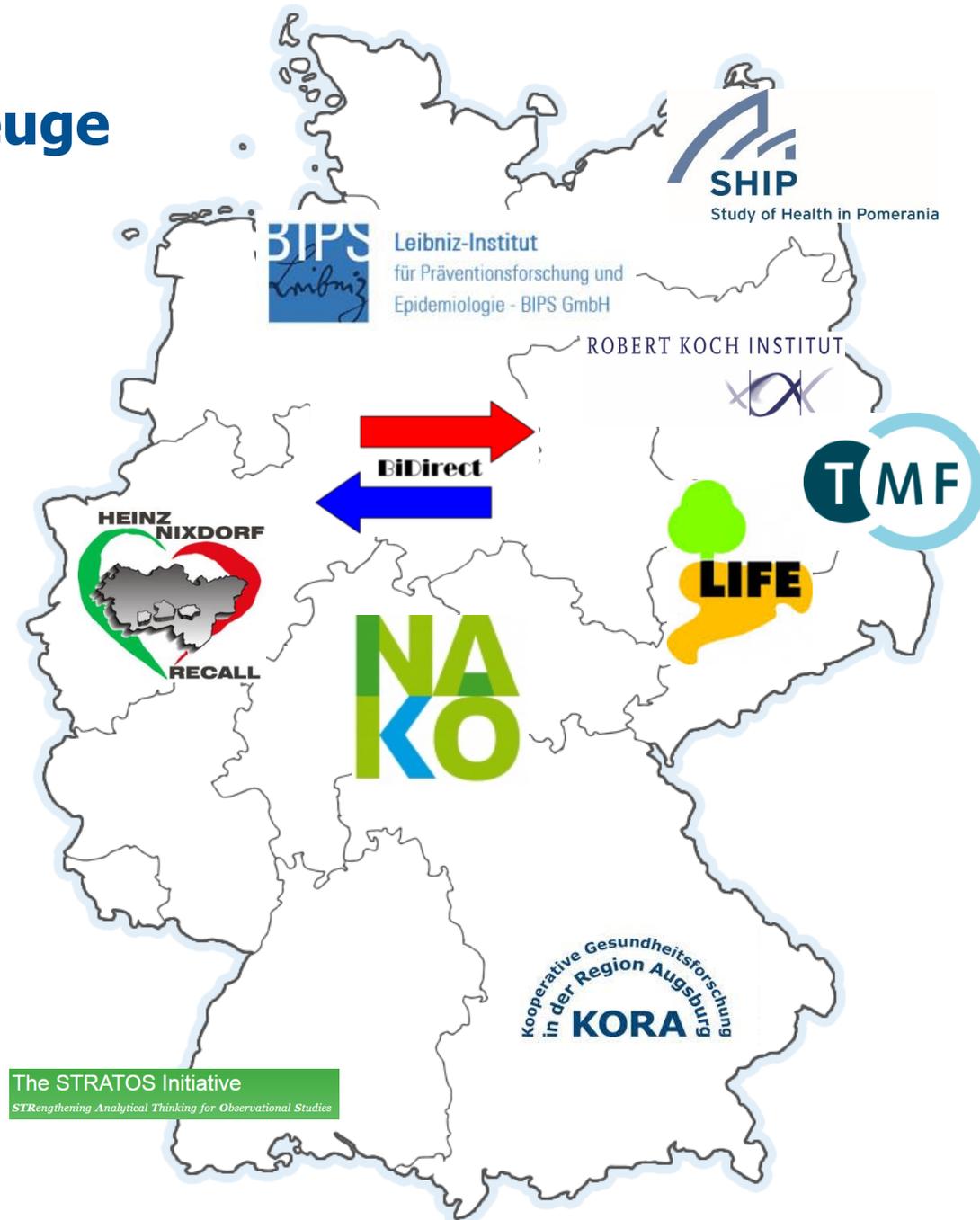
Externes Monitoring (Advisory Board)

Quelle: Schmidt 2014 Qualitätssicherung in Kohortenstudien. Aus Leitlinie zum Adaptiven Datenmanagement in Kohortenstudien und Registern. TMF e.V. |

1. Anwendungsszenario und Datenqualität
2. Anwendungsbeispiel Study of Health in Pomerania
- 3. Konzept zur Bewertung von Datenqualität**
4. Vom Konzept zur Statistik
5. Implementation in einer Kohortenstudie

Standards und Werkzeuge zur Beurteilung der Datenqualität in komplexen epidemiologischen Studien

Schmidt, Carsten Oliver; Prof. Dr. Dr. Bamberg, Fabian; Prof. Dr. Berger, Klaus; Prof. Dr. Hoffmann, Wolfgang; Prof. Dr. Jöckel, Karl-Heinz; Prof. Dr. Kurth, Bärbel-Maria; Prof. Dr. Löffler, Markus; Prof. Dr. Meisinger, Christa; Prof. Dr. Pigeot, Iris; Prof. Dr. Stausberg, Jürgen; Prof. Dr.

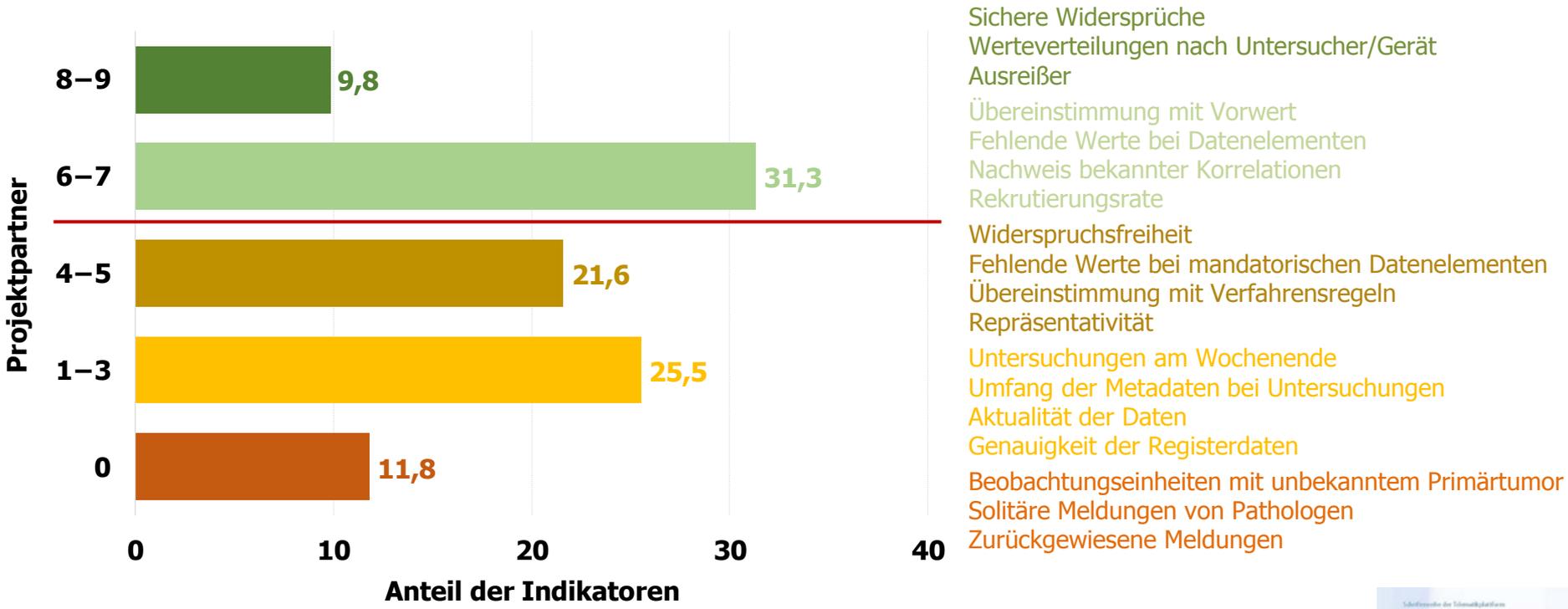


The STRATOS Initiative
STRengthening Analytical Thinking for Observational Studies

Kennziffern zur Qualität – TMF Leitlinie, Auszug

Qualitätsindikator: Integrität	ID neu
Wertevertelung	
<p style="text-align: right;">Bevorzugung bestimmter Endziffern</p>	
<p style="text-align: right;">Wertevertelung der durch Untersucher erfassten Parameter</p>	
<p style="text-align: right;">Wertevertelung der durch Geräte erfassten Parameter</p>	
<p style="text-align: right;">Wertevertelung von Befunden</p>	
<p style="text-align: right;">Wertevertelung von Parametern zwischen Zentren</p>	
Fehlende Einträge	
<p style="text-align: right;">Fehlende Module</p>	TMF-1012
<p style="text-align: right;">Fehlende Werte bei Datenelementen</p>	TMF-1013
<p style="text-align: right;">Fehlende Werte bei mandatorischen Datenelementen</p>	TMF-1014
<p style="text-align: right;">Fehlende Werte bei optionalen Datenelementen</p>	TMF-1015
<p style="text-align: right;">Datenelemente mit Wert unbekannt o. ä.</p>	TMF-1016
Ausreißer bei stetigen Datenelementen	TMF-1018
Werte, die die Messbarkeitsgrenzen von Verfahren unter- oder überschreiten	TMF-1019
Unerlaubte Werte	
<p style="text-align: right;">Unerlaubte Werte bei qualitativen Datenelementen¹</p>	TMF-1021
<p style="text-align: right;">Unerlaubte Werte bei qualitativen Datenelementen zur Kodierung von Missings</p>	TMF-1022
<p style="text-align: right;">Unerlaubte Werte bei quantitativen Datenelementen</p>	TMF-1024

Berücksichtigung inhaltliche Bereiche aus TMF Leitlinie



- Viele **relevante inhaltliche Bereiche** werden getroffen
- **Abstraktionsniveau** der Indikatoren sehr heterogen
- **Definitionen** aus Ausführungen teilweise unklar
- **Registerbezug** sehr eng
- Eignung vorgeschlagener **Berechnungen** („Raten“) teilweise fraglich
- Interpretation der **Ebenen** Integrität, Organisation und Richtigkeit bei Primärdatenerhebung unklar
- **Relevanz** der Indikatoren nicht benannt

Anforderungen

- Indikatoren mit vergleichbarer Komplexität
- Möglichst wenig Überlappung zwischen Indikatoren
- Hierarchische Strukturierung
 - Um verschiedene Komplexitäten abzubilden
- Bestehende TMF Indikatoren sollen gemappt werden
- Definition zur Vereinfachung der Anwendbarkeit

Heterogenes Abstraktionsniveau bestehender Indikatoren



Konzeptüberlegung End Digit Preference

Application examples

Laboratory parameters may follow a normal distribution. A particular device may generate a large number of measurement values equalling "3".

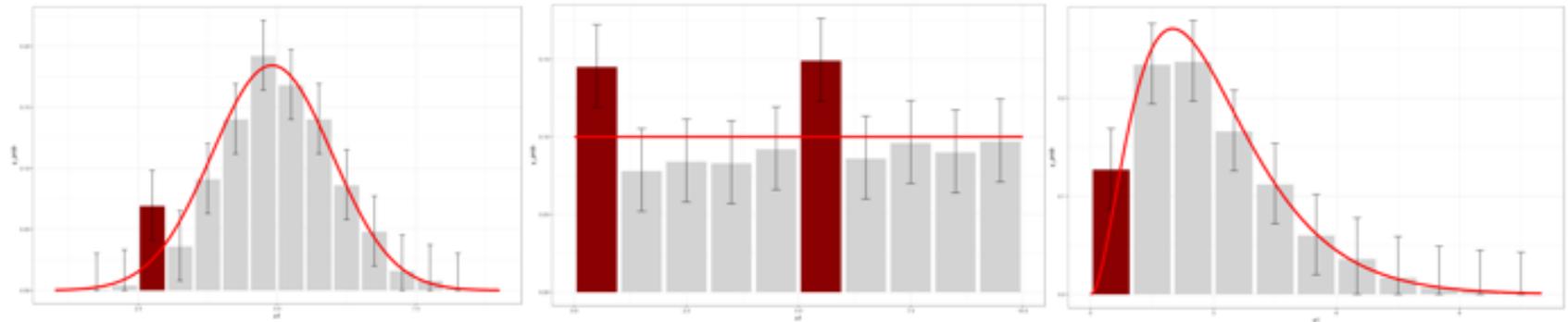


Figure 1: Normal distribution; Figure 2: Uniform distribution; Figure 3: Gamma distribution.

End digit preference: Body weight measurements may be rounded to "0" or "5" instead of using last digits provided by digital balance (Figure 2).

Heterogenes Abstraktionsniveau bestehender Indikatoren



Group level

Indicator domain

Indicator

Unexpected probability distribution

Subtype

End-digit-preference

Group level

Indicator domain Unexpected value distributions

Indicator Unexpected probability distribution

Subtype End-digit-preference

Begriffe zur Beschreibung von Datenqualität

Completeness	Accessibility	Trueness
Data completeness	Timeliness	Correctness
Model completeness	Currency	Accuracy
Data volume	Volatility	Validity
Uniqueness	Availability	Reliability
App propr. Amount of data	Granularity	Reputation
Comprehensiveness	Resolution	Objectivity
Naturalness	Consistency	Plausibility
Rate of enrollment	Concordance	Precision
Utility	Conformance	Agreement
Contextualization	Integrity	Verifiability
Usefulness	Conciseness	Comparability
Relevance	Spatial stability	Standardization
Appropriateness	Predictive value	Generalizability
Informativeness	Coherence	Redundancy
Maintainability	Traceability	Believability
Responsiveness	Interpretability	Credibility
Usability	Complexity	Flexibility
Security	Cohesiveness	Portability

- Group level** ...provides an ontological framework.
- Indicator domain** ...provides a descriptive classification of methods to approach data issues.
- Indicator** ...is the level at which data quality indicators are defined.
- Subtype** ...classifies different application scenarios of an indicator which do not merit the definition of own indicators.

Completeness

The degree to which data values are present in a data collection.

Accuracy

The closeness of agreement between data values and the reference values.

Consistency

The degree to which data values are free of contradictions or convention breaks.

....

Unit Missingness

The degree to which measurements from an entire data collection are missing.

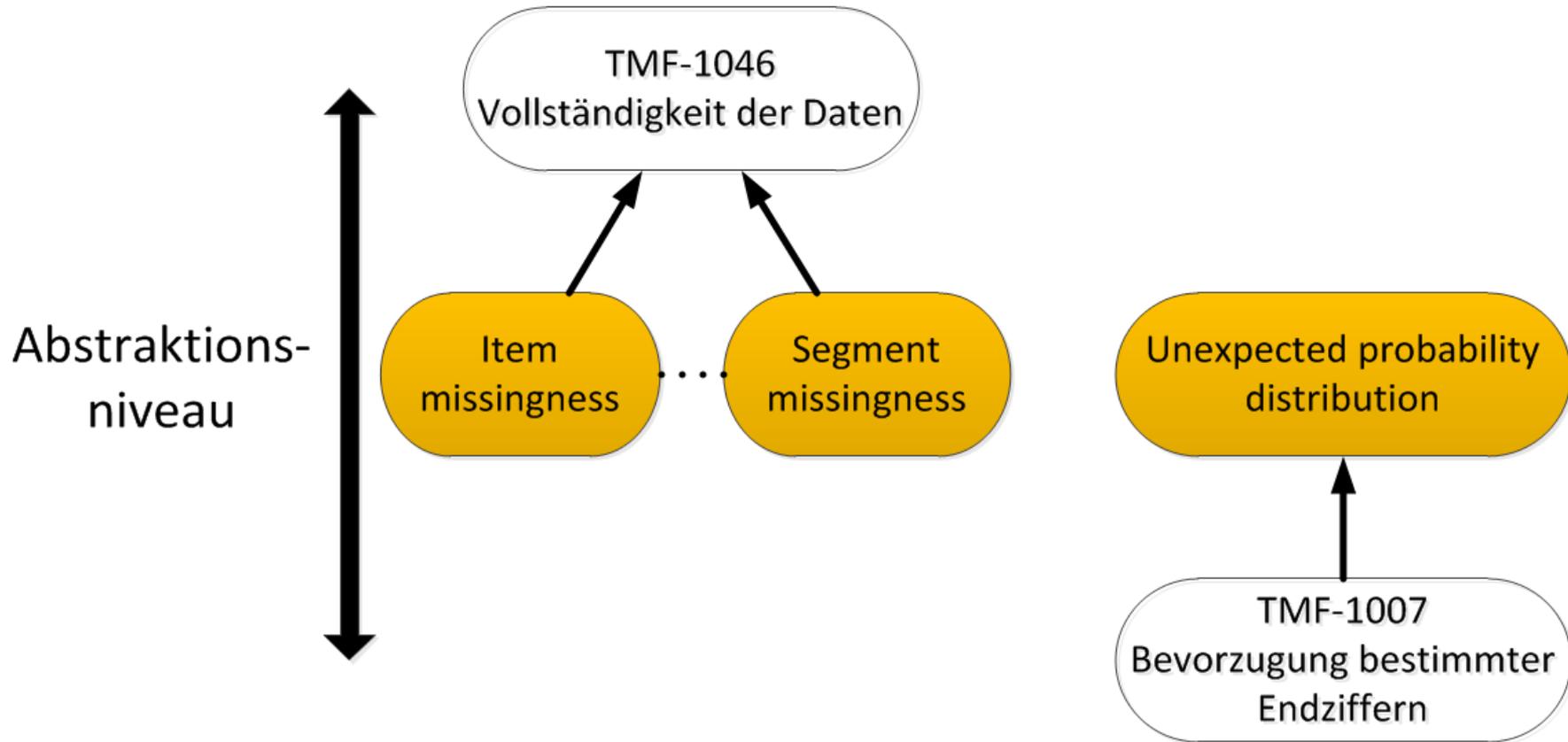
Segment missingness

The degree to which measurements from an entire segment (e.g. examinations or case report) of a data collection are missing.

Item missingness

The degree to which measurements are partially missing in segments of a data collection.

Heterogenes Abstraktionsniveau bestehender Indikatoren



1. Anwendungsszenario und Datenqualität
2. Anwendungsbeispiel Study of Health in Pomerania
3. Konzept zur Bewertung von Datenqualität
- 4. Vom Konzept zur Statistik**
5. Implementation in einer Kohortenstudie



Completeness

Boolean, absolute, relative Häufigkeiten

Consistency

Boolean, absolute, relative Häufigkeiten

Accuracy

Vielfältige Metriken

ICC, Korrelationen

(Nicht-) parametrische Regressionen

Sensitivität, Spezifität, NPW, PPW

Stat. Tests

...

Completeness

Consistency

Accuracy

Prüffokus: Data values

Boolean, absolute, relative Häufigkeiten

Boolean, absolute, relative Häufigkeiten

Vielfältige Metriken

ICC, Korrelationen

(Nicht-) parametrische Regressionen

Sensitivität, Spezifität, NPW, PPW

Stat. Tests

...

Prüffokus:

Verteilungen, Assoziationen



Metadata

Konzeptüberlegung End Digit Preference

Application examples

Laboratory parameters may follow a normal distribution. A particular device may generate a large number of measurement values equalling "3".

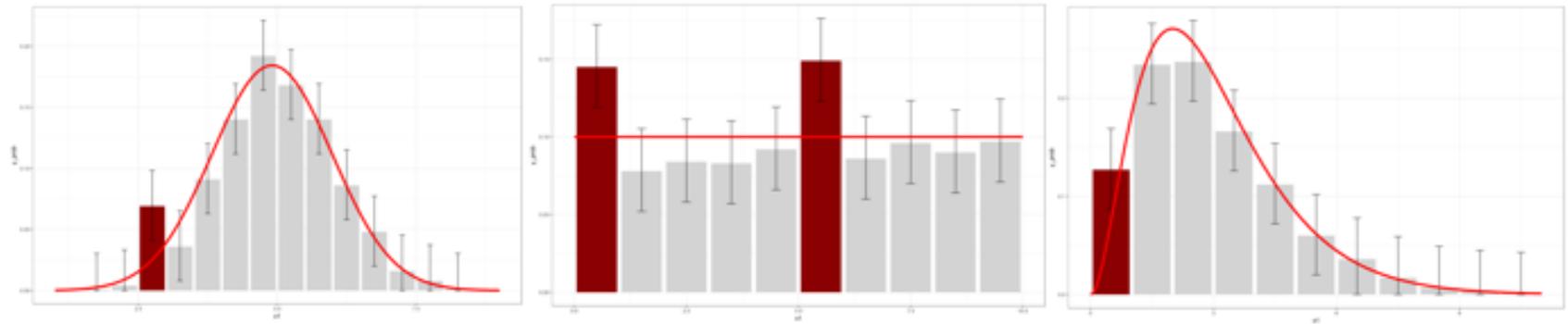


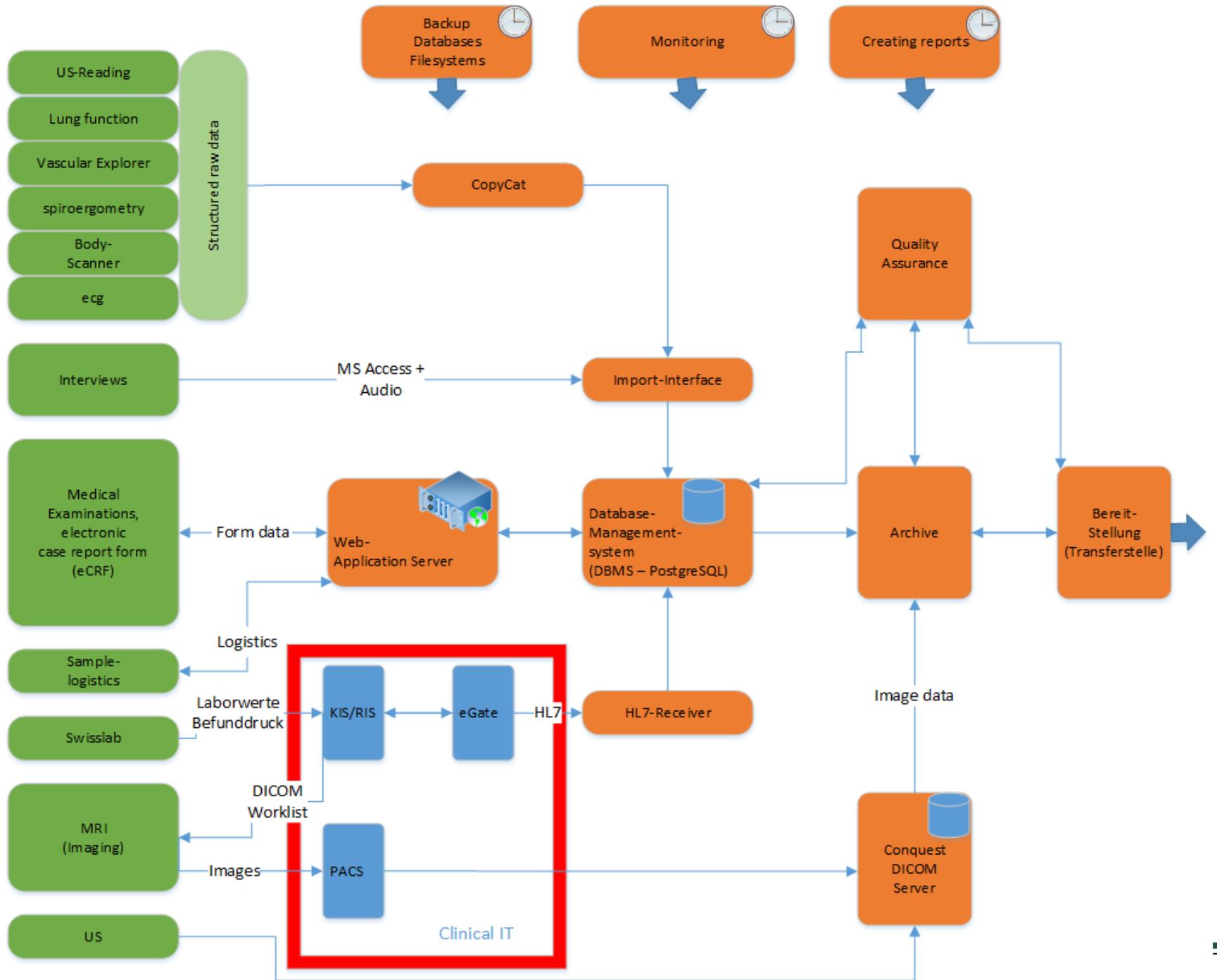
Figure 1: Normal distribution; Figure 2: Uniform distribution; Figure 3: Gamma distribution.

End digit preference: Body weight measurements may be rounded to "0" or "5" instead of using last digits provided by digital balance (Figure 2).

Anwendungen Metadaten in SHIP

	Description	Application examples
Auxiliary variables	related to study design and implementation	
Observer, Device	Identifier of the examiner, reader, etc.	Observer/Device-Effects
Auxiliary variables	related to the environmental measurement conditions	
Date-time stamps	Date and time stamps related to a visit or examination, pre-analytic processing time	Compliance with procedural rules Time trends
Environmental conditions	e.g. temperature, humidity, luminance	Compliance with procedural rules
Metadata attributes	related to completeness	
Missing codes	List of reasons for missing measurements	Item missingness
Jump codes	Conditionally missing measurements	To compute item missingness
Metadata attributes	related to data consistency / precision	
Value list	Variables with predefined categories	Inadmissible measurements
Validity limits	Upper and/or lower validity limits	Inadmissible measurements
Distribution	Expected prob. distribution	Unexpected prob. distribution
Metadata attributes	related to selection of quality statistics	
Data type	e.g. categorical, count, continuous, string	to select appropriate statistics
Metadata attributes	related to interpretation /standardized reporting	
Measurement unit	Continuous variables, e.g. mg/l	Implausible measurements

1. Anwendungsszenario und Datenqualität
2. Anwendungsbeispiel Study of Health in Pomerania
3. Konzept zur Bewertung von Datenqualität
4. Vom Konzept zur Statistik
- 5. Implementation in einer Kohortenstudie**



Framework: Initiale Datenanalyse

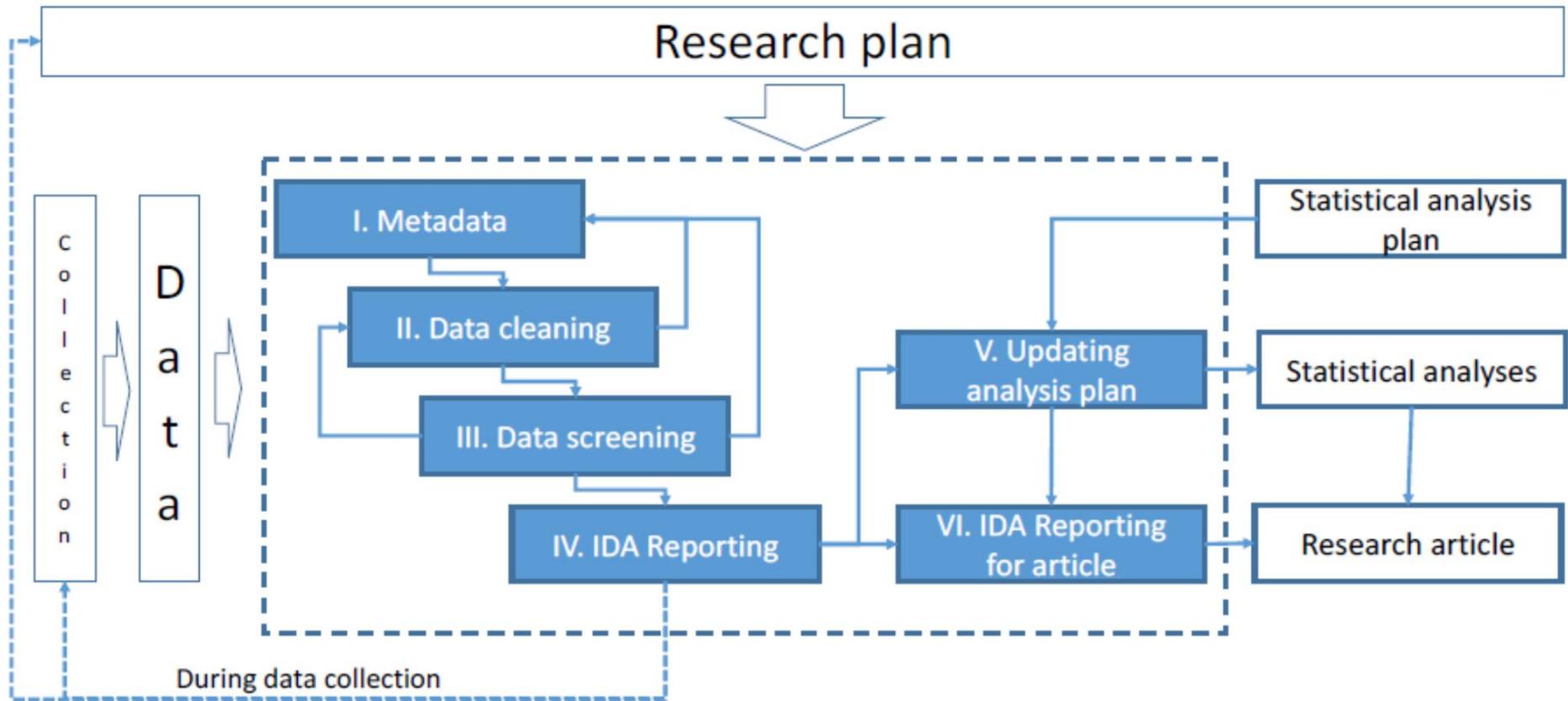
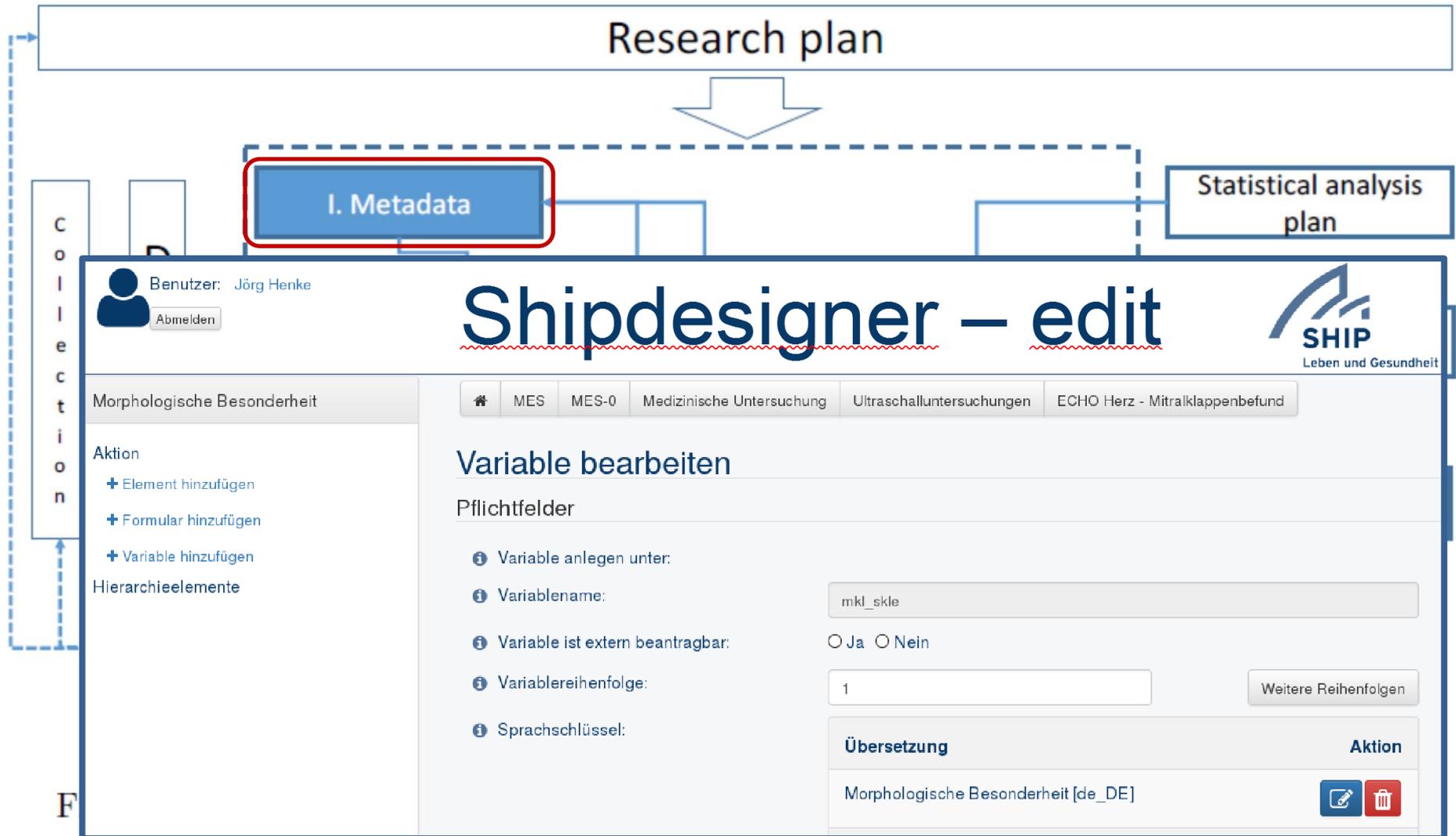


Figure 1: The main connections between the IDA steps and external components

Framework: Initiale Datenanalyse



Framework: Initiale Datenanalyse

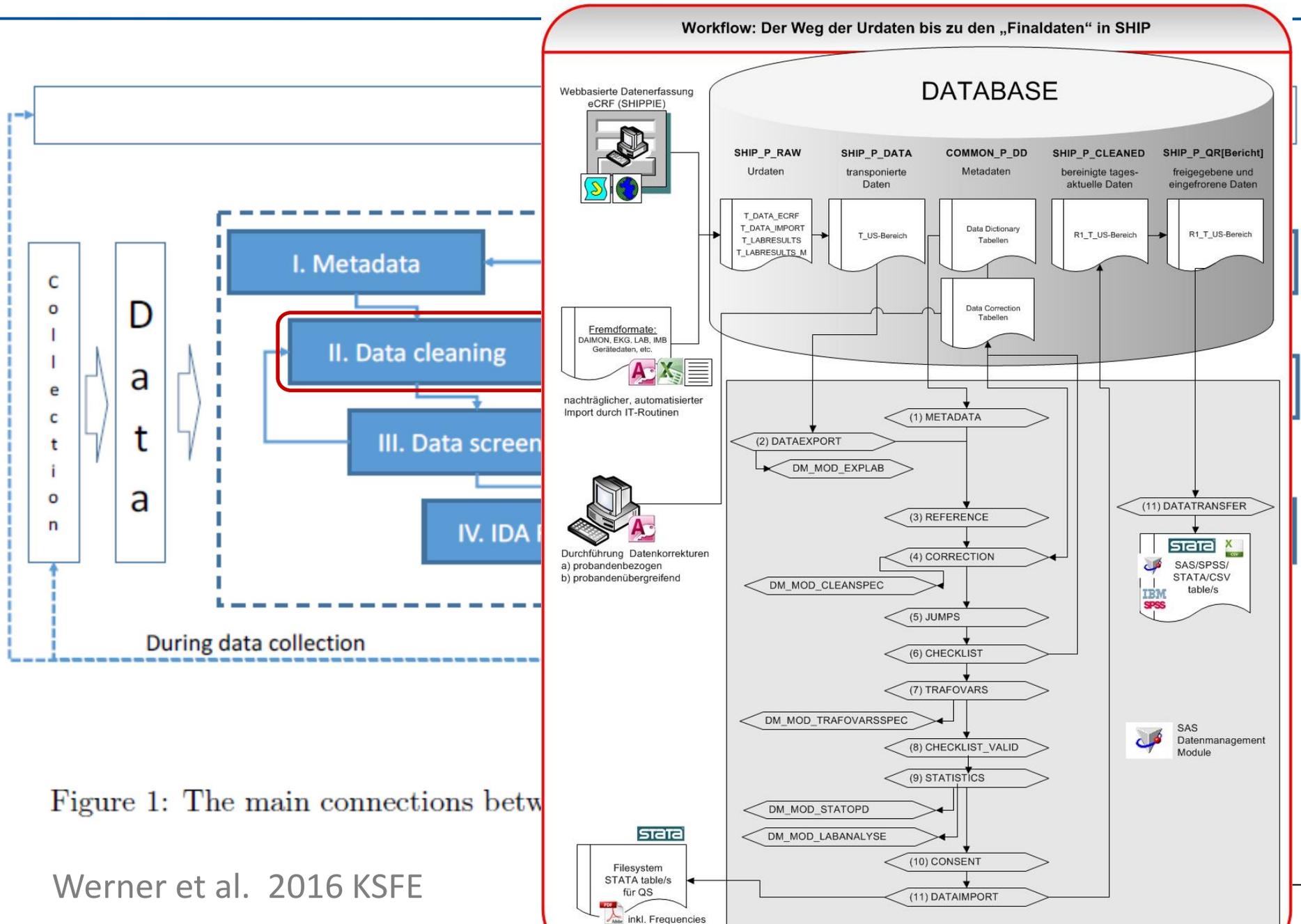


Figure 1: The main connections between

Modul – CHECKLIST / Prüffallkorrektur

- Schritt 2: Durchführung der Datenkorrektur

start_form data_correction_obs_id

100030144387, ECHOREADINGDATA_S, er_mkr2_deceit2, alter Wert: Missing, USNI
100030149755, ECHOREADINGDATA_S, er_lv_lv, alter Wert: Missing, USNR: -, DAT:
100030149755, ECHOREADINGDATA_S, er_lv_rv, alter Wert: Missing, USNR: -, DAT: -
100030149755, ECHOREADINGDATA_S, er_lv_pws, alter Wert: Missing, USNR: -, DAT:
100030149755, ECHOREADINGDATA_S, er_sv, alter Wert: Missing, USNR: -, DAT: -

auswählen

Proband 100030144387 **Korrekturstatus** 2 - Korrektur durchgeführt

Untersuchungsdatum 23.06.2016 13:02:00 **weitere Handlung** 4 - allg. Korrektur

Untersuchernummer 150 **neuer Wert** Format falls Datumsvariable:
TT.MM.JJJJ HH:MM:00
wenn Datum/Uhrzeit fehlt:
jew. Missing aus Missingliste

Oberbereich ECHOREADINGDATA_S **Grund keine Korrektur / Unsicherheit**

Label MK (Ruhe 2sek): Dezelerationszeit [ms]

Variablen-/Formularname er_mkr2_deceit **Fehlertyp** 3 - Fehlwert, Software-/Gerätefehler

Korrekturtyp 5 - einzelne Variable leer **Entscheidungsgrund** 3 - Rücksprache mit Untersucher/-in

zu prüfender Wert . **Besonderheiten?**

Missingliste

99974	Missing durch Software
99977	Nicht durchführbar
99981	Reading noch nicht erfolgt
99982	Variable gehört nicht zur Studie
99983	erlaubter Sprung
99984	Kein Material
99986	Variable nicht mehr verwendet
99987	Nicht bestimmbar
99989	Daten werden geprüft

fertig Eingabe verwerfen

„Zero Tolerance“ für Auffälligkeiten, z.B. bei Consistency

Framework: Initiale Datenanalyse

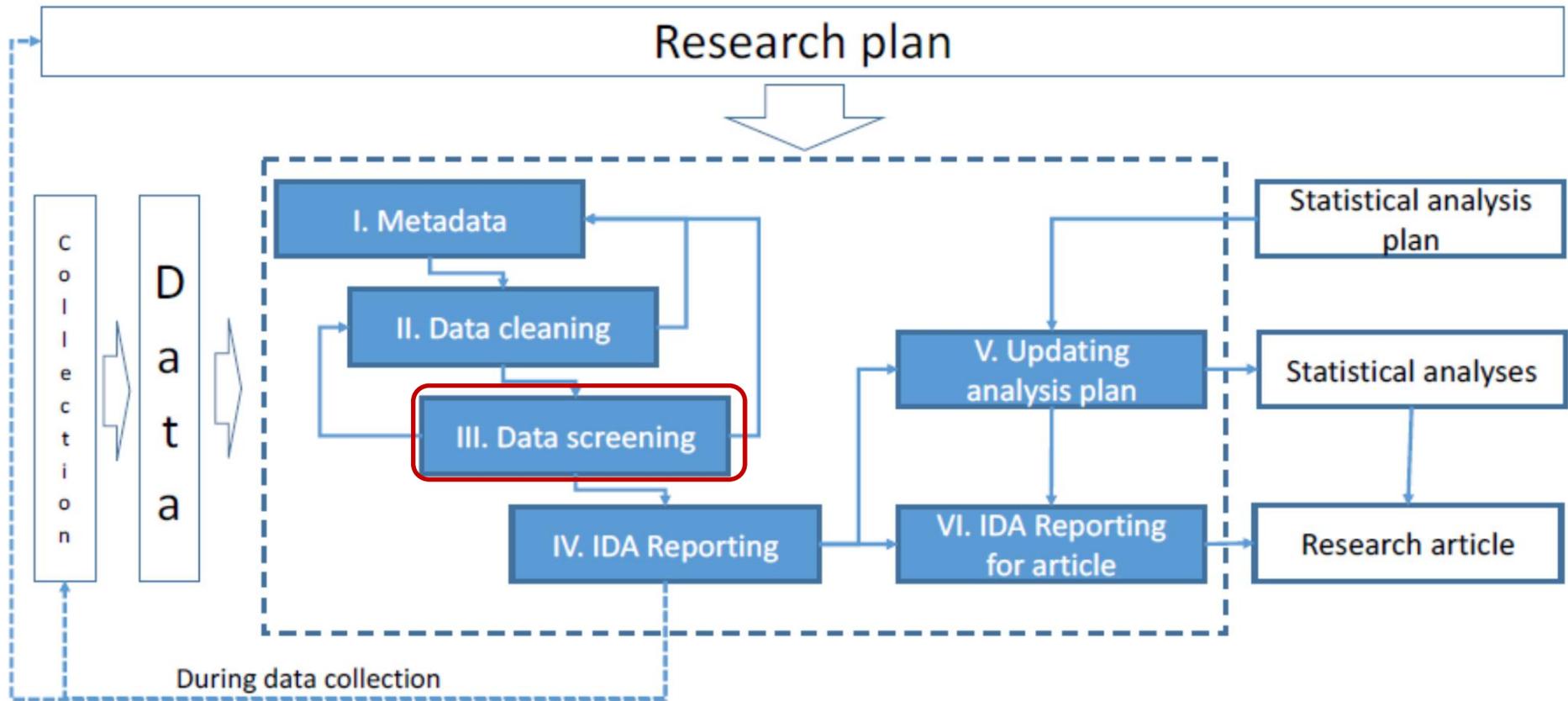
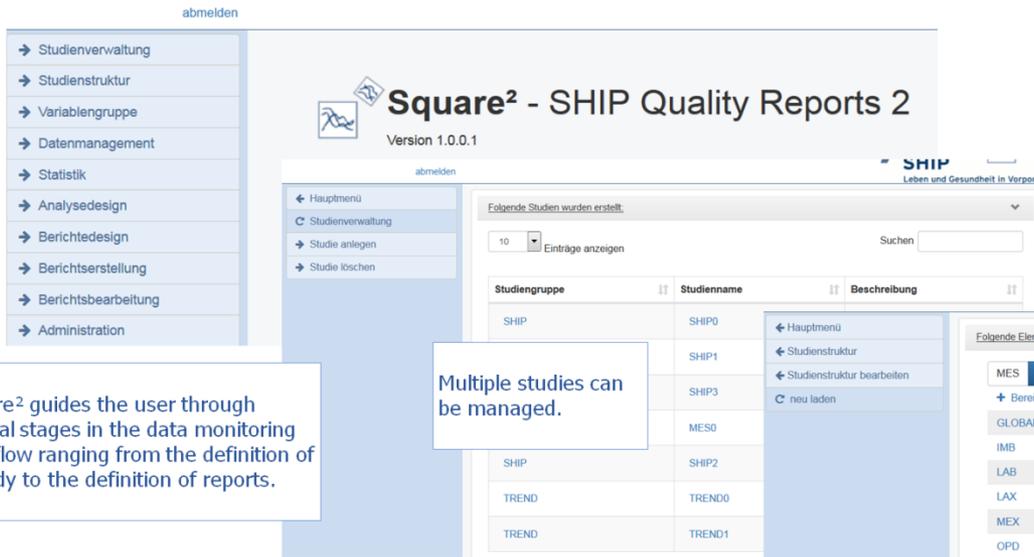


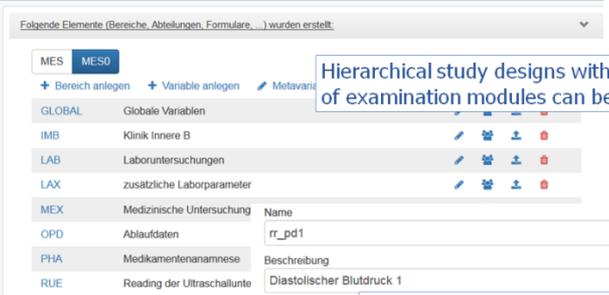
Figure 1: The main connections between the IDA steps and external components



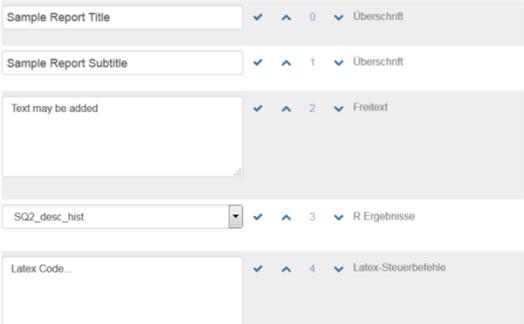
Square² guides the user through several stages in the data monitoring workflow ranging from the definition of a study to the definition of reports.

Multiple studies can be managed.

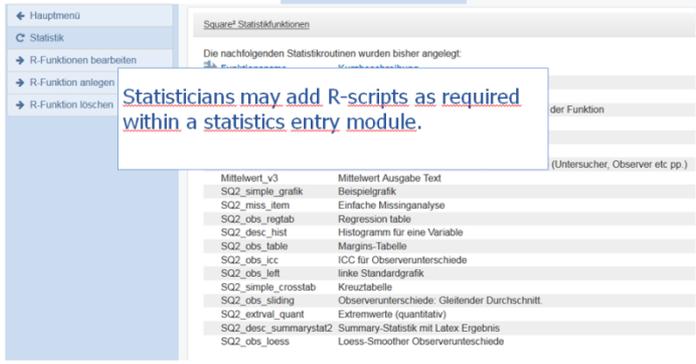
SQUARE² comprises an extensive user rights and roles concept to be usable in large scale studies with many different roles such as PIs, senior and junior quality officers, statisticians, examiners, and data base managers.



Hierarchical study designs with a wide range of examination modules can be managed.



Quality reports can be configured based on a predefined set of standard report elements such as titles, subtitles, text blocks and statistical output.



Statisticians may add R-scripts as required within a statistics entry module.

A meta-data design entry is available for single study variables. Meta variables may be entered via a GUI but may also be imported via existing data dictionaries.

report 33

report
Statistikroutinenübersicht HGP S3

C
3
C
1
C
1

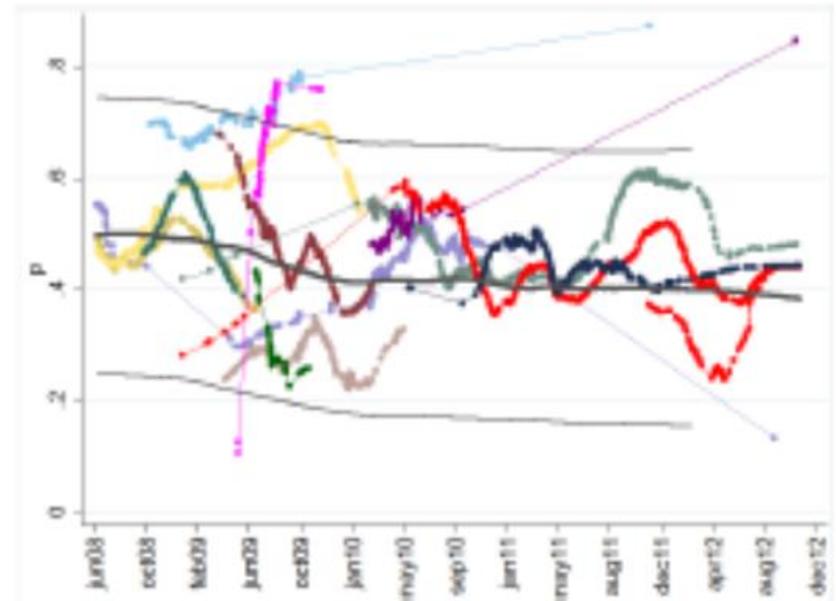
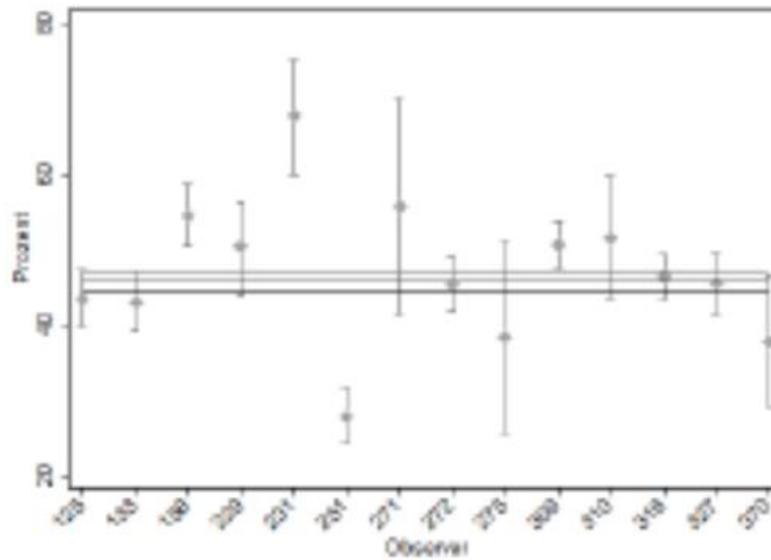


Abbildung 1.2: Mittlere Prävalenz für MRT_KOPF_DEMY adjustiert für Alter und Geschlecht nach logistischer Regression.

1. Datenqualitätsanalyse muss Anwendungskontext berücksichtigen
2. Einheitliche Konzepte zur Datenqualitätsanalyse erforderlich
3. Umfassendes Metadatenmanagement zur Implementation automatisierter Datenqualitätschecks essentiell
4. Metadatenstandards für Datenqualitätschecks erforderlich
5. Umsetzung und Folgen aus Datenqualitätschecks auf den Ebenen Completeness, Consistency, Accuracy sehr unterschiedlich
6. Umfassende infrastrukturelle Anforderungen zur Umsetzung automatisierter Datenqualitätschecks erforderlich
7. Standards für Routinen und Tools wichtig für Datenqualitätsanalyse



<http://www.medizin.uni-greifswald.de/icm/>

carsten.schmidt@uni-greifswald.de

Anwendungen Metadaten

