

# Using routinely collected data for research purposes: Challenges and mitigation strategies

Sabine Hoffmann

26.05.2026

- 1 The role of auxiliary assumptions in the analysis of routinely collected data
- 2 Challenges in the analysis of routinely collected data
  - Representativeness
  - Data quality
  - Time point alignment
  - Interventions and tests are not random
  - Multiplicity of possible analysis strategies
- 3 Mitigation strategies
- 4 Conclusion

# The role of auxiliary assumptions in the analysis of routinely collected data

# Auxiliary assumptions in research

# Auxiliary assumptions in research

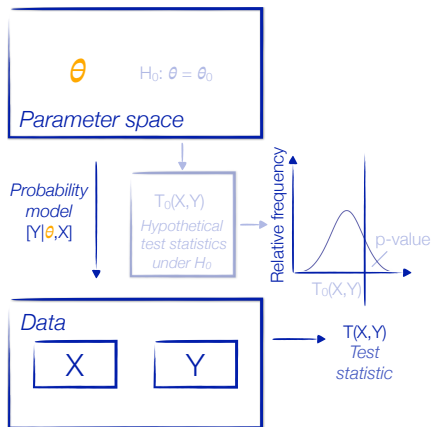
## Underdetermination of scientific theory by evidence:

It is not possible to unambiguously falsify a theory, because theories are always tested in a bundle with various auxiliary assumptions.

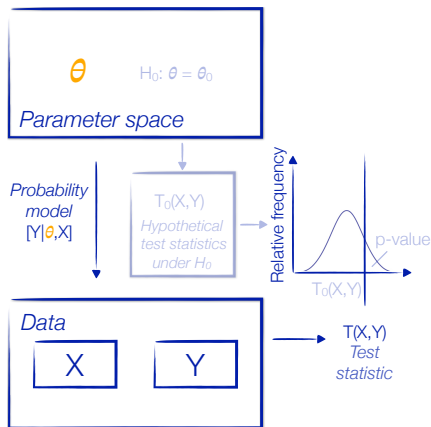
When a prediction fails, we never know whether we should blame the theory or one of the auxiliary assumptions.

(Duhem-Quine problem)

# Auxiliary assumptions in statistical modelling in medicine

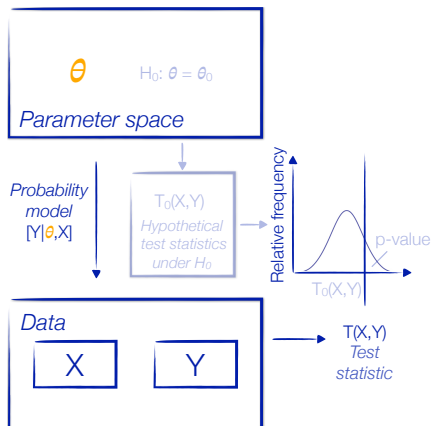


# Auxiliary assumptions in statistical modelling in medicine



⇒ When results deviate from our expectations under  $H_0$ , should we reject  $H_0$  or any of the auxiliary assumptions:

# Auxiliary assumptions in statistical modelling in medicine



⇒ When results deviate from our expectations under  $H_0$ , should we reject  $H_0$  or any of the auxiliary assumptions:

- *"The patients took the assigned treatment correctly"*
- *"All factors influencing treatment switching and the outcome were measured"*

## Auxiliary assumptions play an even bigger role in the analysis of routinely collected data

- When data are prospectively collected for research purposes:
  - The research question informs the study design
  - In experimental studies, it even determines the data generating mechanism

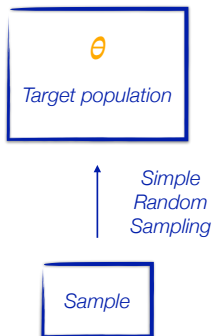
# Auxiliary assumptions play an even bigger role in the analysis of routinely collected data

- When data are prospectively collected for research purposes:
  - The research question informs the study design
  - In experimental studies, it even determines the data generating mechanism
- In routinely collected data, there is
  - Little knowledge of how the data were generated
  - Little control over measurement procedures

# Auxiliary assumptions play an even bigger role in the analysis of routinely collected data

- When data are prospectively collected for research purposes:
  - The research question informs the study design
  - In experimental studies, it even determines the data generating mechanism
- In routinely collected data, there is
  - Little knowledge of how the data were generated
  - Little control over measurement procedures
- Consequences:
  - It is in general not possible to control who receives which treatments
  - It is unknown why certain patients received certain treatments
  - Sometimes one does not even know when and for how long patients received a certain treatment

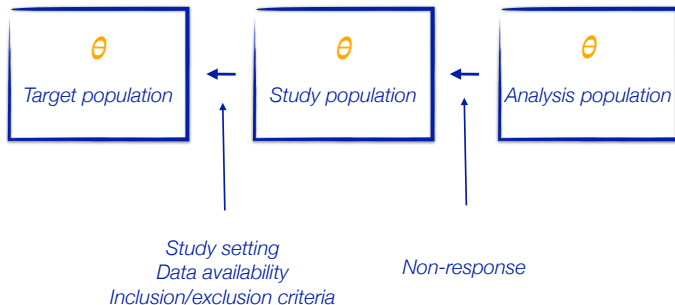
# Auxiliary assumptions play an even bigger role in the analysis of routinely collected data



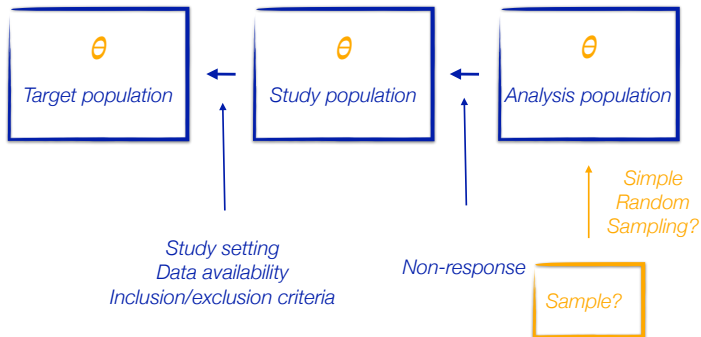
# Auxiliary assumptions play an even bigger role in the analysis of routinely collected data



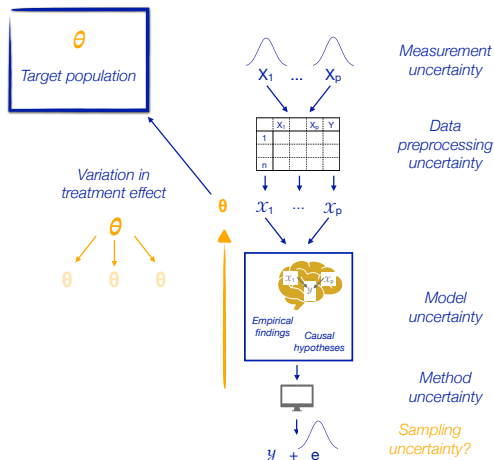
# Auxiliary assumptions play an even bigger role in the analysis of routinely collected data



# Auxiliary assumptions play an even bigger role in the analysis of routinely collected data

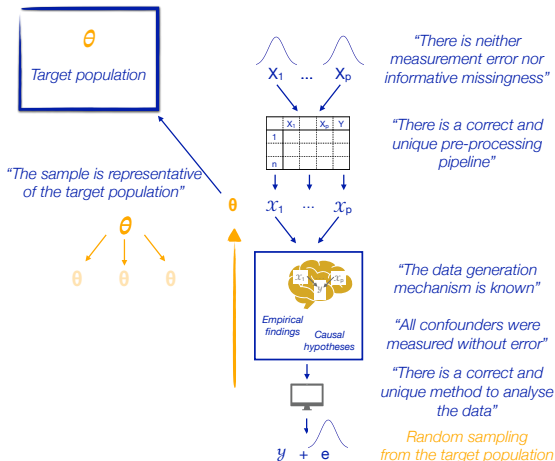


## Existing and nonexistant sources of uncertainty

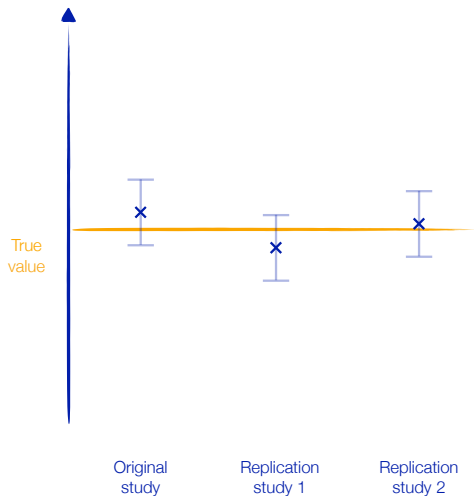


Hoffmann, S., F. Schönbrodt, R. Elsas, R. Wilson, U. Strasser, Boulesteix, A. L. (2021). The multiplicity of analysis strategies jeopardizes replicability: lessons learned across disciplines. Royal Society Open Science 8 201925

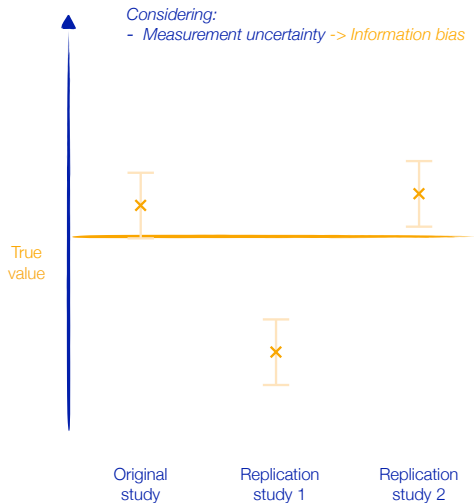
## Existing and non-existent sources of uncertainty



# If we ignore uncertainty, it leads to bias



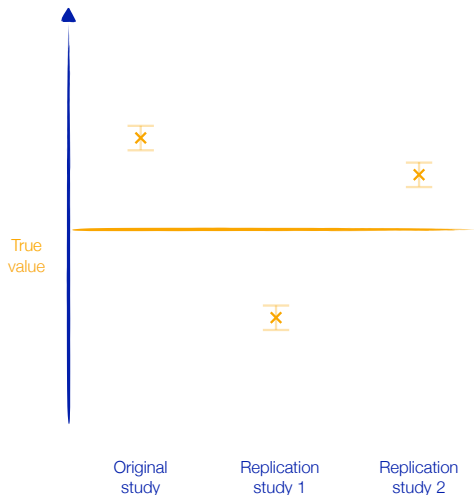
# If we ignore uncertainty, it leads to bias



# If we ignore uncertainty, it leads to bias



## Big data paradoxes (Meng (2018))



“The more the data, the surer we fool ourselves”

# Challenges in the analysis of routinely collected data

## *“The sample is representative of the target population”*

- Due to disparities in the access and use of health-care services, routinely collected data may under- or mis-represent certain subgroups [1, 2, 3, 4, 5], including
  - ethnic minorities
  - patients without medical coverage
  - low-income and rural populations

## *“The sample is representative of the target population”*

- Due to disparities in the access and use of health-care services, routinely collected data may under- or mis-represent certain subgroups [1, 2, 3, 4, 5], including
  - ethnic minorities
  - patients without medical coverage
  - low-income and rural populations
- If patients from under-represented groups are present in the data, there is a risk that they may be mis-represented, because they are
  - more likely to visit multiple institutions [6, 7, 8]
  - they receive fewer diagnostics tests and interventions [2]
  - they are more likely to be lost during data linkage

## *“The sample is representative of the target population”*

- Due to disparities in the access and use of health-care services, routinely collected data may under- or mis-represent certain subgroups [1, 2, 3, 4, 5], including
  - ethnic minorities
  - patients without medical coverage
  - low-income and rural populations
- If patients from under-represented groups are present in the data, there is a risk that they may be mis-represented, because they are
  - more likely to visit multiple institutions [6, 7, 8]
  - they receive fewer diagnostics tests and interventions [2]
  - they are more likely to be lost during data linkage
- There may be substantial variation in patient characteristics and in disease occurrence across sites in multisite studies

*“There is no measurement error or informative missingness”*

- Routinely collected data are not recorded for a specific research purpose at hand

## *“There is no measurement error or informative missingness”*

- Routinely collected data are not recorded for a specific research purpose at hand
- ⇒ Documentation practices may vary between different clinical settings, as a function of incentives and of the overall workload of the personnel collecting the data [9, 10]

## *“There is no measurement error or informative missingness”*

- Routinely collected data are not recorded for a specific research purpose at hand
- ⇒ Documentation practices may vary between different clinical settings, as a function of incentives and of the overall workload of the personnel collecting the data [9, 10]
- **Examples:**
  - On an emergency call, vital parameters may not be measured if they are irrelevant to the clinical question at hand
  - During a busy period, nurses may not find the time to record clinical events or changes in medication, or they may only find time at the end of their shift
  - Prescription orders may not be filled or consumed by the patient
  - Temporal changes in the recording of data may produce systematic differences over time

## *“There is no measurement error or informative missingness”*

- Routinely collected data are not recorded for a specific research purpose at hand
- ⇒ Documentation practices may vary between different clinical settings, as a function of incentives and of the overall workload of the personnel collecting the data [9, 10]
- **Examples:**
    - On an emergency call, vital parameters may not be measured if they are irrelevant to the clinical question at hand
    - During a busy period, nurses may not find the time to record clinical events or changes in medication, or they may only find time at the end of their shift
    - Prescription orders may not be filled or consumed by the patient
    - Temporal changes in the recording of data may produce systematic differences over time
- ⇒ Parameter estimates can be biased in any direction

*“It is possible to find an alignment between eligibility, treatment assignment and start of follow-up”*

- In routinely collected data, the first entry for a patient typically has no clear clinical meaning and it is often unknown which diagnoses or treatments the patient received before this time point

*“It is possible to find an alignment between eligibility, treatment assignment and start of follow-up”*

- In routinely collected data, the first entry for a patient typically has no clear clinical meaning and it is often unknown which diagnoses or treatments the patient received before this time point
- In a representative sample of reports of comparative non-randomised studies that assessed the effectiveness and/or safety of drug treatments, Yaacoub et al. (2024) [11] found that in 72% of studies eligibility, treatment assignment, and start of follow-up were not aligned.

*“It is possible to find an alignment between eligibility, treatment assignment and start of follow-up”*

- In routinely collected data, the first entry for a patient typically has no clear clinical meaning and it is often unknown which diagnoses or treatments the patient received before this time point
- In a representative sample of reports of comparative non-randomised studies that assessed the effectiveness and/or safety of drug treatments, Yaacoub et al. (2024) [11] found that in 72% of studies eligibility, treatment assignment, and start of follow-up were not aligned.
- In a data audit on the quality of observational study data in an international HIV research network, treatment regimens and associated dates and the timings of laboratory measurements were especially prone to error with error rates of up to 56% and 42% [12]

*“It is possible to find an alignment between eligibility, treatment assignment and start of follow-up”*

Intensive Care Med

<https://doi.org/10.1007/s00134-025-07805-4>

**ORIGINAL**

# Management of high-risk acute pulmonary embolism: an emulated target trial analysis



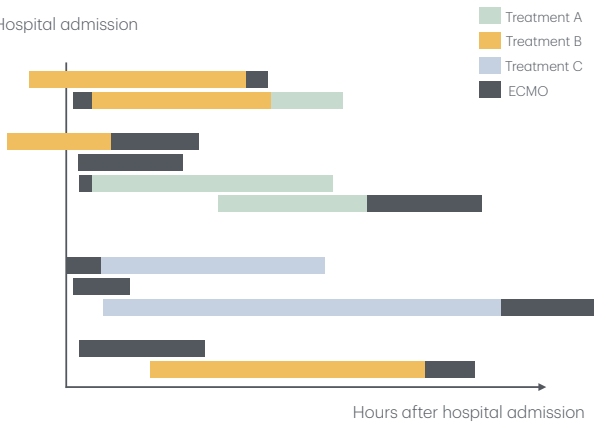
Andrea Stadlbauer<sup>1</sup>, Tom Verbelen<sup>2</sup>, Leonhard Binzenhöfer<sup>3</sup>, Tomaz Goslar<sup>4</sup>, Alexander Supady<sup>5</sup>, Peter M. Spieth<sup>6</sup>, Marko Noc<sup>4</sup>, Andreas Verstraete<sup>2</sup>, Sabine Hoffmann<sup>7</sup>, Michael Schomaker<sup>41</sup>, Julia Höpler<sup>7</sup>, Marie Kraft<sup>7</sup>, Esther Tautz<sup>5</sup>, Daniel Hoyer<sup>8</sup>, Jörn Tongers<sup>8</sup>, Franz Haertel<sup>9</sup>, Aschraf El-Essawi<sup>10</sup>, Mostafa Salem<sup>11</sup>, Rafael Henrique Rangel<sup>11</sup>, Carsten Hullermann<sup>12</sup>, Marvin Kriz<sup>13</sup>, Benedikt Schrage<sup>13</sup>, Jorge Moisés<sup>14</sup>, Manel Sabate<sup>14</sup>, Federico Pappalardo<sup>15</sup>, Lisa Crusius<sup>16</sup>, Norman Mangner<sup>16</sup>, Christoph Adler<sup>17</sup>, Tobias Tichelbäcker<sup>17</sup>, Carsten Skurk<sup>18</sup>, Christian Jung<sup>19</sup>, Sebastian Kufner<sup>20</sup>, Tobias Graf<sup>21</sup>, Clemens Scherer<sup>3</sup>, Laura Villegas Sierra<sup>3</sup>, Hannah Billig<sup>22</sup>, Nicolas Majunke<sup>23</sup>, Walter S. Speidl<sup>24</sup>, Robert Zilberszac<sup>24</sup>, Luis Chiscano-Camón<sup>25</sup>, Aitor Uribarri<sup>26</sup>, Jordi Riera<sup>25</sup>, Roberto Roncon-Albuquerque Jr<sup>27</sup>, Elizabete Terauda<sup>28</sup>, Andrejs Erglis<sup>28</sup>, Guido Tavazzi<sup>29</sup>, Uwe Zeymer<sup>30</sup>, Maike Knorr<sup>31</sup>, Juliane Kilo<sup>32</sup>, Sven Möbius-Winkler<sup>9</sup>, Robert H. G. Schwinger<sup>33</sup>, Derk Frank<sup>11</sup>, Oliver Borst<sup>34</sup>, Helene Häberle<sup>35</sup>, Frederic De Roeck<sup>36</sup>, Christian Vrints<sup>36</sup>, Christof Schmid<sup>1</sup>, Georg Nickenig<sup>22</sup>, Christian Hagl<sup>37</sup>, Steffen Massberg<sup>3</sup>, Andreas Schäfer<sup>38</sup>, Dirk Westermann<sup>39</sup>, Sebastian Zimmer<sup>22</sup>, Alain Combes<sup>40</sup>, Daniele Camboni<sup>1\*</sup>, Holger Thiele<sup>23</sup> and Enzo Lüsebrink<sup>22\*</sup> for the High-risk P. E. Investigator Group

© 2025 The Author(s)

ID	A	B	C	ECMO	Mortality
1	0	1	0	1	1
2	1	1	0	1	0
3	1	0	0	0	1
4	0	1	0	1	1
5	0	0	0	1	1
6	1	0	0	1	0
7	1	0	0	1	1
8	0	0	1	0	1
9	0	1	0	0	0
10	0	0	1	1	0
11	0	0	0	1	1
12	0	0	1	1	1
13	0	0	1	0	0
14	0	0	0	1	1
15	0	1	0	1	1

*“It is possible to find an alignment between eligibility, treatment assignment and start of follow-up”*

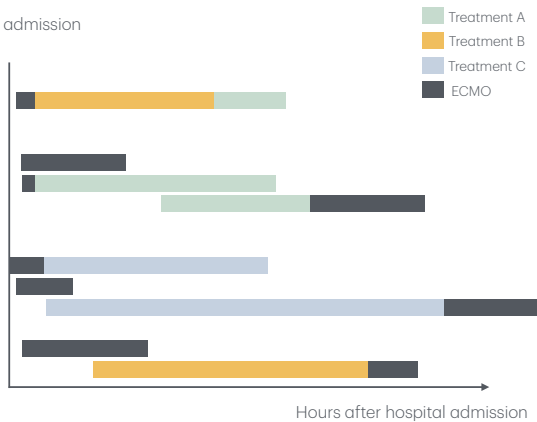
Hospital admission



ID	A	B	C	ECMO	Mortality
1	0	1	0	1	1
2	1	1	0	1	0
3	1	0	0	0	1
4	0	1	0	1	1
5	0	0	0	1	1
6	1	0	0	1	0
7	1	0	0	1	1
8	0	0	1	0	1
9	0	1	0	0	0
10	0	0	1	1	0
11	0	0	0	1	1
12	0	0	1	1	1
13	0	0	1	0	0
14	0	0	0	1	1
15	0	1	0	1	1

*“It is possible to find an alignment between eligibility, treatment assignment and start of follow-up”*

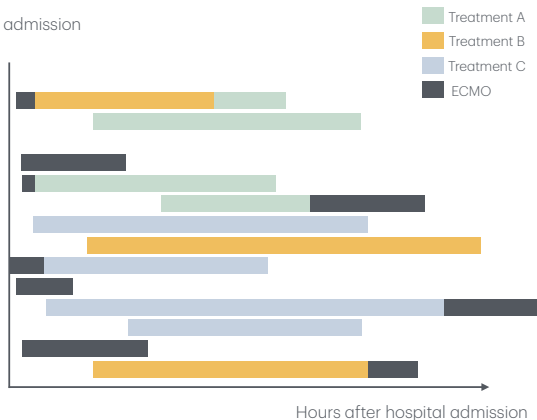
Hospital admission



ID	A	B	C	ECMO	Mortality
1	0	1	0	1	1
2	1	1	0	1	0
3	1	0	0	0	1
4	0	1	0	1	1
5	0	0	0	1	1
6	1	0	0	1	0
7	1	0	0	1	1
8	0	0	1	0	1
9	0	1	0	0	0
10	0	0	1	1	0
11	0	0	0	1	1
12	0	0	1	1	1
13	0	0	1	0	0
14	0	0	0	1	1
15	0	1	0	1	1

*“It is possible to find an alignment between eligibility, treatment assignment and start of follow-up”*

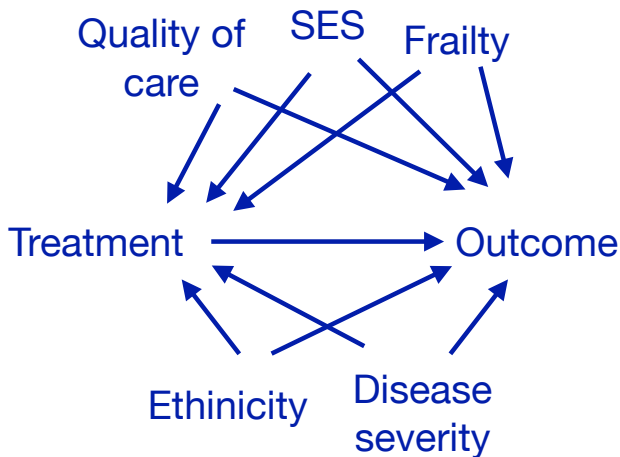
Hospital admission

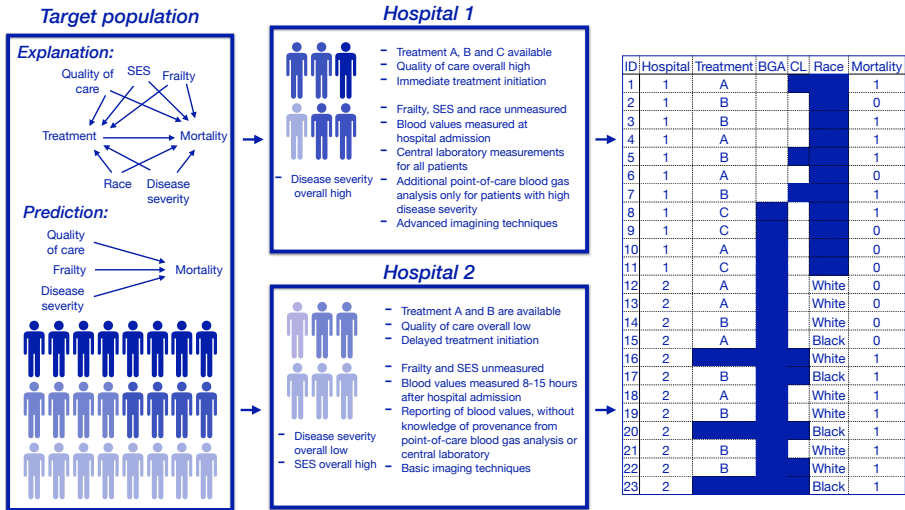


ID	A	B	C	ECMO	Mortality
1	0	1	0	1	1
2	1	1	0	1	0
3	1	0	0	0	1
4	0	1	0	1	1
5	0	0	0	1	1
6	1	0	0	1	0
7	1	0	0	1	1
8	0	0	1	0	1
9	0	1	0	0	0
10	0	0	1	1	0
11	0	0	0	1	1
12	0	0	1	1	1
13	0	0	1	0	0
14	0	0	0	1	1
15	0	1	0	1	1

*“All relevant confounders were measured without error”*

*“All relevant confounders were measured without error”*





Hoffmann, S., Morris, T., Herrmann, M., Heinze, G., Wynants, L., Van Calster, B., Bischl, B., Schmid, M., Shaw, P.A., Mathes, T., Naudet, F., Harrell, F.E. (...), Thiele, H., Lüsebrink, E. (2026). Using routinely collected data for research purposes: Challenges and mitigation strategies. *The BMJ In print*

# Relying on auxiliary assumptions leads to a multiplicity of possible analysis strategies

- “There is no measurement error or informative missingness”
- “It is possible to find an alignment between eligibility, treatment assignment and start of follow-up”
- “All relevant confounders were measured without error”

# Relying on auxiliary assumptions leads to a multiplicity of possible analysis strategies

## Underdetermination of scientific theory by evidence:

It is not possible to unambiguously falsify a theory, because theories are always tested in a bundle with various auxiliary assumptions.

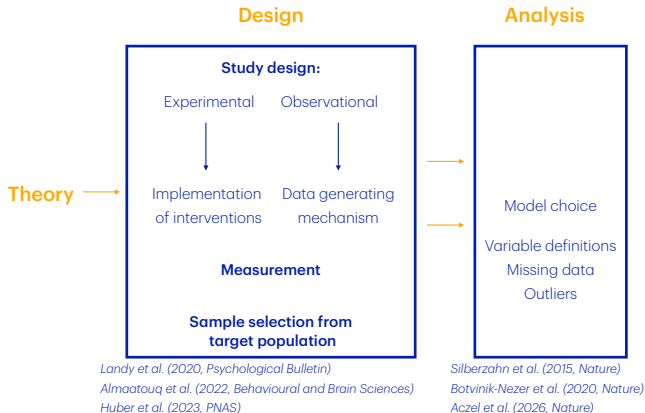
When a prediction fails, we never know whether we should blame the theory or one of the auxiliary assumptions.

(Duhem-Quine problem)

# Relying on auxiliary assumptions leads to a multiplicity of possible analysis strategies

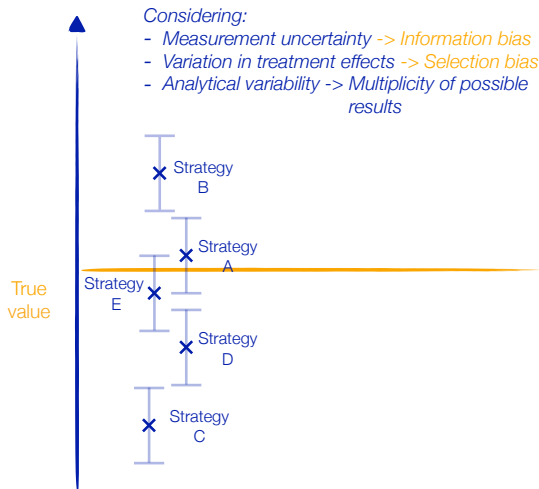
Underdetermination of scientific theory by evidence

Type 2: Underdetermination of evidence by scientific theory

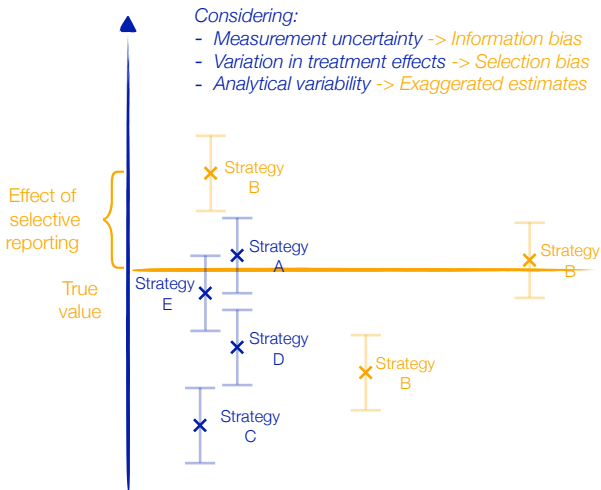




# Selective reporting among a multiplicity of possible analysis strategies leads to non-replicability and overconfidence



# Selective reporting among a multiplicity of possible analysis strategies leads to non-replicability and overconfidence



# Selective reporting among a multiplicity of possible analysis strategies leads to conflicting findings [13, 14, 15]

Surgery 165 (2019) 953–957



Surgery 165 (2019) 1199–1202



## Colon/Rectum

## Does retrieval bag use during laparoscopic appendectomy reduce postoperative infection?

Adam C. Fields, MD<sup>a,c</sup>, Pamela Lu, MD<sup>a,b</sup>, Deanna L. Palenzuela, BS<sup>a</sup>, Ronald Bleday, MD<sup>a</sup>, Joel E. Goldberg, MD, MPH<sup>a</sup>, Jennifer Irani, MD<sup>a</sup>, Jennifer S. Davids, MD<sup>a</sup>, Nelya Melnitchouk, MD, MSc<sup>a,b,c</sup>

<sup>a</sup>Department of Surgery, Brigham and Women's Hospital, Harvard Medical School, Boston, MA

<sup>b</sup>Center for Surgery and Public Health, Department of Surgery, Brigham

<sup>c</sup>Department of Surgery, University of Massachusetts Memorial Health



Presented at the Academic Surgical Congress 2019

## Utilization of a specimen retrieval bag during laparoscopic appendectomy for both uncomplicated and complicated appendicitis is not associated with a decrease in postoperative surgical site infection rates

Scott A. Turner, MD<sup>\*</sup>, Hee Soo Jung, MD, FACS, John E. Scarborough, MD, FACS




Table. Comparison of 2 Studies on the Association of a Specimen Retrieval Bag With Surgical Site Infection Rates in Laparoscopic Appendectomy

Criteria	Source	Turner et al. <sup>a</sup> 2019
Inclusion criteria	Fields et al. <sup>8</sup> 2019	Turner et al. <sup>4</sup> 2019
Analytic sample reported, No.	11 475	10 357
Primary outcome	Postoperative intra-abdominal abscess	Any SSI (superficial, deep, organ space)
Primary predictor	Use of retrieval bag	Use of retrieval bag
Covariates included, with operationalization	Age (continuous) Sex (dichotomized) BMI (continuous) Race (categorized as: White, Black, Asian, other) Diabetes (dichotomized) Hypertension (dichotomized) COPD (dichotomized) Smoker (dichotomized) Functional status (dichotomized) Steroid use (dichotomized) Weight loss (dichotomized) Preoperative sepsis (dichotomized) Wound class 3/4 (dichotomized) Complicated appendicitis (dichotomized) ASA class 3/4 (dichotomized) Operative time (continuous) White blood cell count (continuous)	Age (dichotomized at 65 y) Sex (dichotomized) Obesity (categorical: not obese, class I/II/III obesity, missing) Not included Diabetes (dichotomized) Not included Not included Not included Not included Steroid use (dichotomized) Not included Unclear if included Not included 2 indicator variables: presence of abscess and presence of perforation Not included Operative time dichotomized at 75th percentile Not included
Coefficient on primary predictor	OR (95% CI): 0.6 (0.42–0.95)	OR (95% CI): 1.15 (0.78–1.69)
	P value: .03	P value: .49

# Mitigation strategies to address challenges in the analysis of routinely collected data

# Transparency is both more important and more difficult in the analysis of routinely collected data

## ANALYSIS

 Check for updates

<sup>1</sup> CHU Rennes, Inserm, Institut de Recherche en Santé, Environnement et Travail-UMR\_S 1085, University of Rennes, Rennes, France

<sup>2</sup> Institut Universitaire de France, Paris, France

<sup>3</sup> Department of Biomedical Informatics, Harvard Medical School, Boston, Massachusetts, USA

<sup>4</sup> Nuffield Department of Primary Care Health Sciences, University of Oxford, Oxford, UK

## Improving the transparency and reliability of observational studies through registration

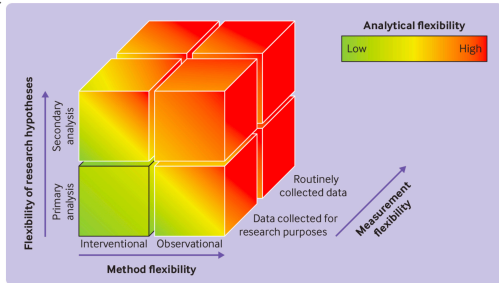
**Florian Naudet and colleagues** argue that routine registration of observational research is needed and suggest how current processes can be adapted to facilitate it

Florian Naudet,<sup>1,2</sup> Chirag J Patel,<sup>3</sup> Nicholas J DeVito,<sup>4</sup> Gérard Le Goff,<sup>5</sup> Ioana A Cristea,<sup>6</sup> Alain Brailion,<sup>7</sup> Sabine Hoffmann<sup>8-9</sup>

From the use of booster doses against covid-19<sup>1</sup> to

This already happens for certain studies, such as

an international post-authorization safety study



BMJ: first published as 10.1136/bmj-2023-076123

# Mitigation strategies



**Hoffmann, S.,** Morris, T., Herrmann, M., Heinze, G., Wynants, L., Van Calster, B., Bischl, B., Schmid, M., Shaw, P.A., Mathes, T., Naudet, F., Harrell, F.E. (...), Thiele, H., Lüsebrink, E. (2026). Using routinely collected data for research purposes: Challenges and mitigation strategies. *The BMJ In print*

# Mitigation strategies

Diagnosis	
Representativeness	Comparison of sample characteristics with target population and across centres
Time point alignment	Target trial emulation
Data quality	Discussion with people familiar with data collection
	Extensive data quality checks
Interventions not random	Interviews with physicians who are responsible for treatment decision
Multiplicity of analysis strategies	Comparison of results when two analysts independently analyse data including pre-processing steps

**Hoffmann, S., Morris, T., Herrmann, M., Heinze, G., Wynants, L., Van Calster, B., Bischl, B., Schmid, M., Shaw, P.A., Mathes, T., Naudet, F., Harrell, F.E. (...), Thiele, H., Lüsebrink, E. (2026).** Using routinely collected data for research purposes: Challenges and mitigation strategies. *The BMJ In print*

# Mitigation strategies

	Diagnosis	Design	
Representativeness	Comparison of sample characteristics with target population and across centres	Use of validated definitions for sample selection, time points, exposures, confounders, outcomes and or validation using prospectively collected and/or external data	Randomisation
Time point alignment	Target trial emulation		
Data quality	Discussion with people familiar with data collection	Standardisation and training	
	Extensive data quality checks		
Interventions not random	Interviews with physicians who are responsible for treatment decision	Combined analysis of RCD with prospective data from RCTs	
Multiplicity of analysis strategies	Comparison of results when two analysts independently analyse data including pre-processing steps	Pre-registration of detailed protocol including statistical analysis plan	

**Hoffmann, S., Morris, T., Herrmann, M., Heinze, G., Wynants, L., Van Calster, B., Bischl, B., Schmid, M., Shaw, P.A., Mathes, T., Naudet, F., Harrell, F.E. (...), Thiele, H., Lüsebrink, E. (2026).** Using routinely collected data for research purposes: Challenges and mitigation strategies. *The BMJ In print*

# Mitigation strategies

	Diagnosis	Design	Analysis	
Representativeness	Comparison of sample characteristics with target population and across centres	Use of validated definitions for sample selection, time points, exposures, confounders, outcomes and or validation using prospectively collected and/or external data	Matching and weighting techniques	
Time point alignment	Target trial emulation		Cloning, censoring and weighting	
Data quality	Discussion with people familiar with data collection		Standardisation and training	Correction for measurement error and missing values
	Extensive data quality checks			Quantitative bias analysis and falsification techniques (negative controls or outcomes)
Interventions not random	Interviews with physicians who are responsible for treatment decision		Combined analysis of RCD with prospective data from RCTs	Accounting for sources of uncertainty leading to multiplicity of analysis strategies
Multiplicity of analysis strategies	Comparison of results when two analysts independently analyse data including pre-processing steps	Pre-registration of detailed protocol including statistical analysis plan		

Randomisation

Hoffmann, S., Morris, T., Herrmann, M., Heinze, G., Wynants, L., Van Calster, B., Bischl, B., Schmid, M., Shaw, P.A., Mathes, T., Naudet, F., Harrell, F.E. (...), Thiele, H., Lüsebrink, E. (2026). Using routinely collected data for research purposes: Challenges and mitigation strategies. *The BMJ In print*

# Steps to produce reliable results on routinely collected data

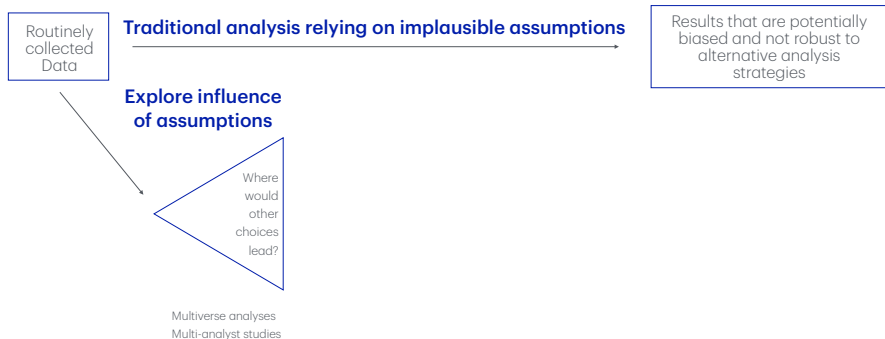
Routinely  
collected  
Data

**Traditional analysis relying on implausible assumptions**

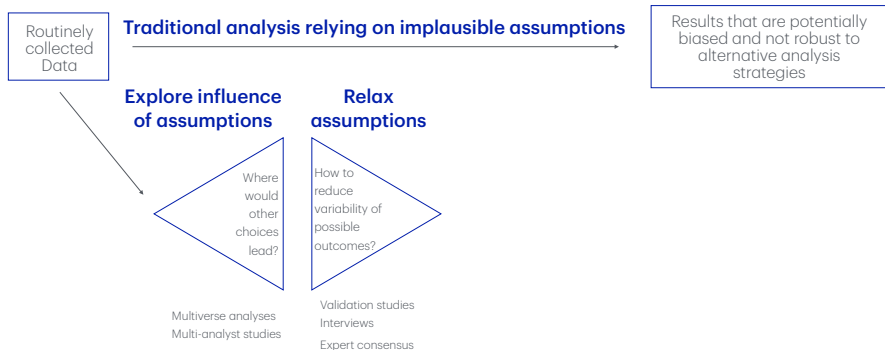


Results that are potentially  
biased and not robust to  
alternative analysis  
strategies

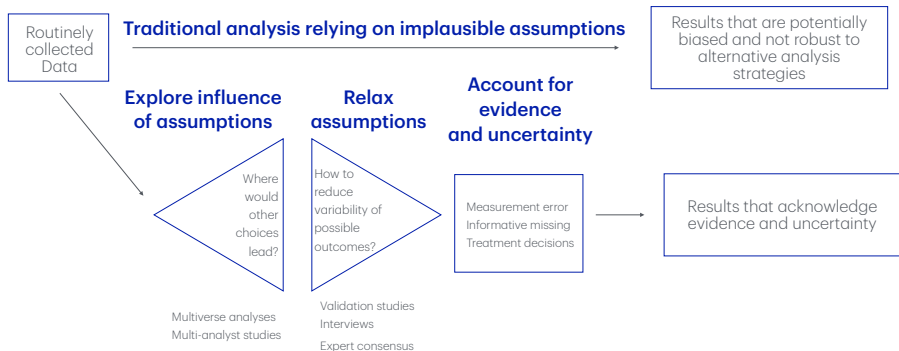
# Steps to produce reliable results on routinely collected data



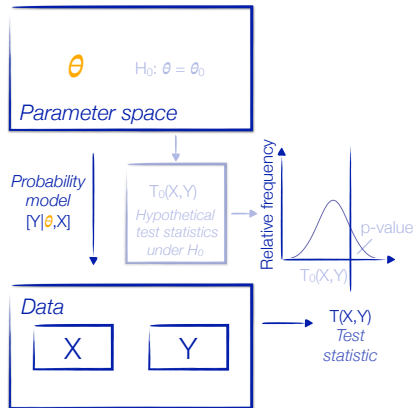
# Steps to produce reliable results on routinely collected data



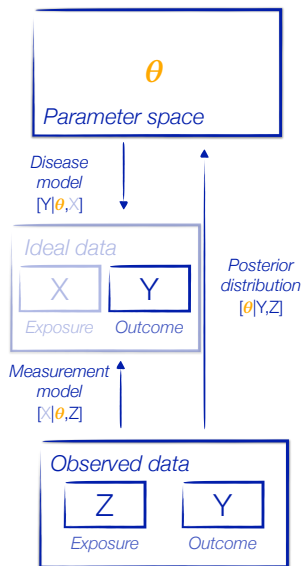
# Steps to produce reliable results on routinely collected data



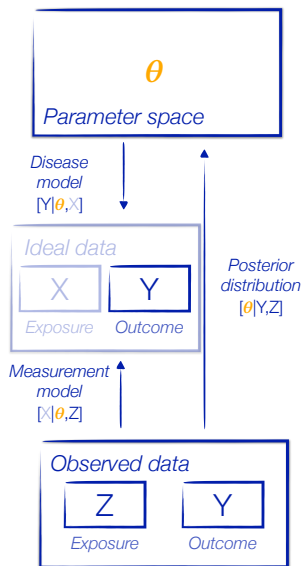
## Accounting for evidence and uncertainty



## Accounting for evidence and uncertainty

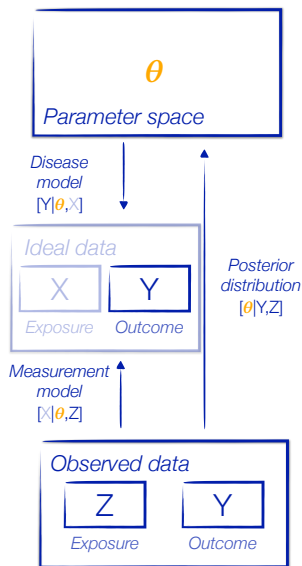


## Accounting for evidence and uncertainty



Suggestions on how to relax auxiliary assumptions:

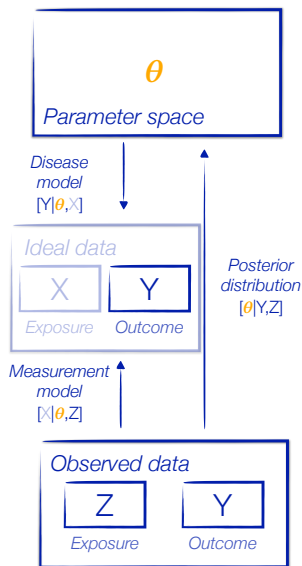
## Accounting for evidence and uncertainty



Suggestions on how to relax auxiliary assumptions:

- Think (very hard) about the data generating mechanism

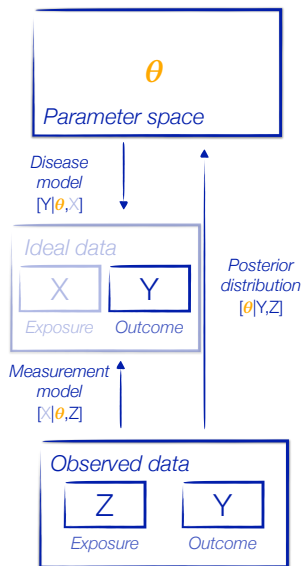
# Accounting for evidence and uncertainty



Suggestions on how to relax auxiliary assumptions:

- Think (very hard) about the data generating mechanism
- Conduct validation studies and interviews with medical personnel to relax auxiliary assumptions

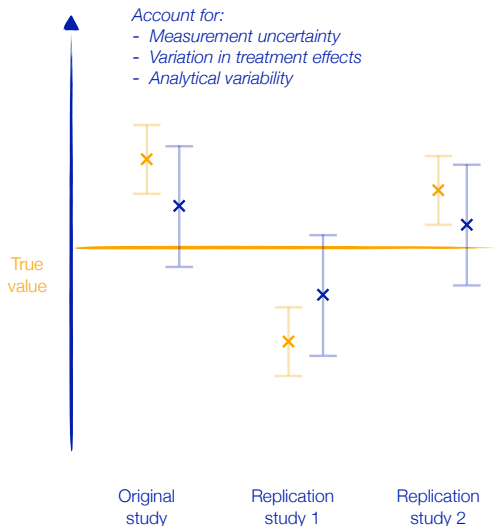
# Accounting for evidence and uncertainty



Suggestions on how to relax auxiliary assumptions:

- Think (very hard) about the data generating mechanism
- Conduct validation studies and interviews with medical personnel to relax auxiliary assumptions
- Account for all sources of evidence and all sources of uncertainty

# Relax auxiliary assumptions to reduce overconfidence and improve replicability



# Conclusion

# Conclusion

- Auxiliary assumptions play a very important role in the analysis of routinely collected data because researchers usually lack knowledge of the data generating mechanism and they have no control over measurement procedures

# Conclusion

- Auxiliary assumptions play a very important role in the analysis of routinely collected data because researchers usually lack knowledge of the data generating mechanism and they have no control over measurement procedures
- To produce reliable and credible results, it is important to verify and to relax auxiliary assumptions and to account for all evidence and uncertainty in the analysis of routinely collected data

# Conclusion

- Auxiliary assumptions play a very important role in the analysis of routinely collected data because researchers usually lack knowledge of the data generating mechanism and they have no control over measurement procedures
- To produce reliable and credible results, it is important to verify and to relax auxiliary assumptions and to account for all evidence and uncertainty in the analysis of routinely collected data
- Qualitatively discussing biases arising from the the violation of auxiliary assumptions makes scientific uncertainty essentially inaccessible to readers who have limited knowledge and understanding of methodological weaknesses and challenges

Thank you for your attention!



Benjamin A Goldstein, Nrupen A Bhavsar, Matthew Phelan, and Michael J Pencina.

Controlling for informed presence bias due to the number of health encounters in an electronic health record.

*American journal of epidemiology*, 184(11):847–855, 2016.



Milena A Gianfrancesco and Neal D Goldstein.

A narrative review on the validity of electronic health record-based research in epidemiology.

*BMC medical research methodology*, 21(1):234, 2021.



Ariana Mihan and Harriette GC Van Spall.

Interventions to enhance digital health equity in cardiovascular care.

*Nature Medicine*, 30(3):628–630, 2024.



Ariana Mihan, Ambarish Pandey, and Harriette GC Van Spall.

Mitigating the risk of artificial intelligence bias in cardiovascular care.

*The Lancet Digital Health*, 6(10):e749–e754, 2024.



Ban Al-Sahab, Alan Leviton, Tobias Loddenkemper, Nigel Paneth, and Bo Zhang.

Biases in electronic health records data for generating real-world evidence: An overview.

*Journal of Healthcare Informatics Research*, 8(1):121–139, 2024.



Nicholas C Arpey, Anne H Gaglioti, and Marcy E Rosenbaum.

How socioeconomic status affects patient perceptions of health care: a qualitative study.

*Journal of primary care & community health*, 8(3):169–175, 2017.



Judy H Ng, Faye Ye, Lauren M Ward, Samuel C, ÁúChris, À Haffer, and Sarah Hudson Scholle.

Data on race, ethnicity, and language largely incomplete for managed care plan members.

*Health Affairs*, 36(3):548–552, 2017.



Milena A Gianfrancesco, Suzanne Tamang, Jinoos Yazdany, and Gabriela Schmajuk.

Potential biases in machine learning algorithms using electronic health record data.

*JAMA internal medicine*, 178(11):1544–1547, 2018.



Kaiwen Ni, Hongling Chu, Lin Zeng, Nan Li, and Yiming Zhao.  
Barriers and facilitators to data quality of electronic health records used for clinical research in china: a qualitative study.

*BMJ open*, 9(7):e029314, 2019.



Robert A Verheij, Vasa Curcin, Brendan C Delaney, and Mark M McGilchrist.

Possible sources of bias in primary care electronic health record data use and reuse.

*Journal of medical Internet research*, 20(5):e185, 2018.



Sally Yaacoub, Raphael Porcher, Anna Pellat, Hillary Bonnet, Viet-Thi Tran, Philippe Ravaud, and Isabelle Boutron.

Characteristics of non-randomised studies of drug treatments: cross sectional study.

*BMJ medicine*, 3(1):e000932, 2024.



Stephany N Duda, Bryan E Shepherd, Cynthia S Gadd, Daniel R Masys, and Catherine C McGowan.

Measuring the quality of observational study data in an international hiv research network.

*PLoS one*, 7(4):e33908, 2012.



Adam C Fields, Pamela Lu, Deanna L Palenzuela, Ronald Bleday, Joel E Goldberg, Jennifer Irani, Jennifer S Davids, and Nelya Melnitchouk.

Does retrieval bag use during laparoscopic appendectomy reduce postoperative infection?

*Surgery*, 165(5):953–957, 2019.



Scott A Turner, Hee Soo Jung, and John E Scarborough.

Utilization of a specimen retrieval bag during laparoscopic appendectomy for both uncomplicated and complicated appendicitis is not associated with a decrease in postoperative surgical site infection rates.

*Surgery*, 165(6):1199–1202, 2019.



Christopher P Childers and Melinda Maggard-Gibbons.

Same data, opposite results?: A call to improve surgical database research.

*JAMA surgery*, 156(3):219–220, 2021.



Yeganeh Khazaei, Helmut Küchenhoff, Sabine Hoffmann, Diella Sylighi, and Raphael Rehms.

Using a bayesian hierarchical approach to study the association between non-pharmaceutical interventions and the spread of covid-19 in germany.

*Scientific Reports*, 13(1):18900, 2023.



Raphael Rehms, Nicole Ellenbach, Eva Rehfues, Jacob Burns, Ulrich Mansmann, and **Hoffmann, Sabine**.

**A Bayesian hierarchical approach to account for evidence and uncertainty in the modeling of infectious diseases: An application to COVID-19.**

*Biometrical Journal*, 66(1):2200341, 2024.

# Criteria to assure the quality of studies on routinely collected data

	<b>Minimal requirements</b>	<b>Acceptable</b>	<b>Ideal approach</b>
<b>Eligibility</b>	Detailed comparison of patient characteristics with target population and across centres		
<b>Data quality</b>	Extensive data quality checks and reporting of data pre-processing and missing patterns		
<b>Time point alignment</b>	Report exact timing of all measurements and of treatment trajectories		
<b>Interventions and tests not random</b>	Directed acyclic graph including unmeasured and mismeasured confounders		
<b>Multiplicity of analysis</b>	Pre-registration of		

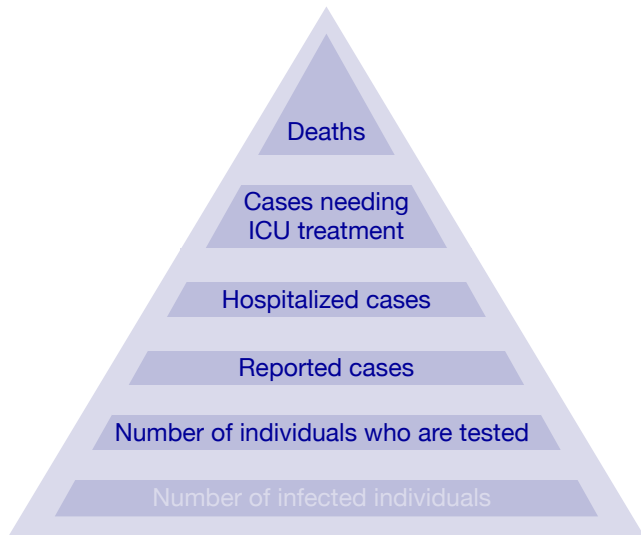
# Criteria to assure the quality of studies on routinely collected data

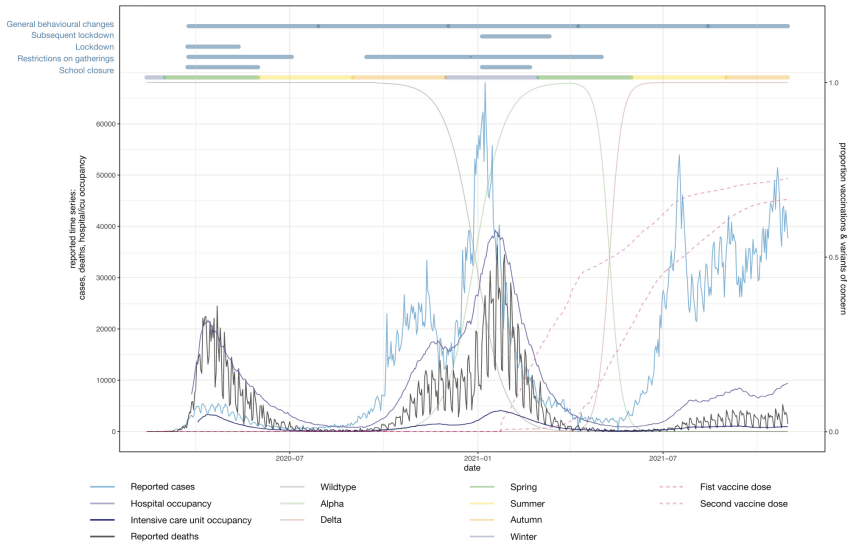
	<b>Minimal requirements</b>	<b>Acceptable</b>	<b>Ideal approach</b>
<b>Eligibility</b>	Detailed comparison of patient characteristics with target population and across centres	Inverse probability weighting or multilevel regression modelling with post-stratification	
<b>Data quality</b>	Extensive data quality checks and reporting of data pre-processing and missing patterns	Audits, standardisation and training in data collection, validation data to quantify errors	
<b>Time point alignment</b>	Report exact timing of all measurements and of treatment trajectories	Target trial emulation	
<b>Interventions and tests not random</b>	Directed acyclic graph including unmeasured and mismeasured confounders	Quantitative bias analysis and falsification endpoints of negative controls	
<b>Multiplicity of analysis</b>	Pre-registration of	Multi-analyst studies or extensive multiverse	

# Criteria to assure the quality of studies on routinely collected data

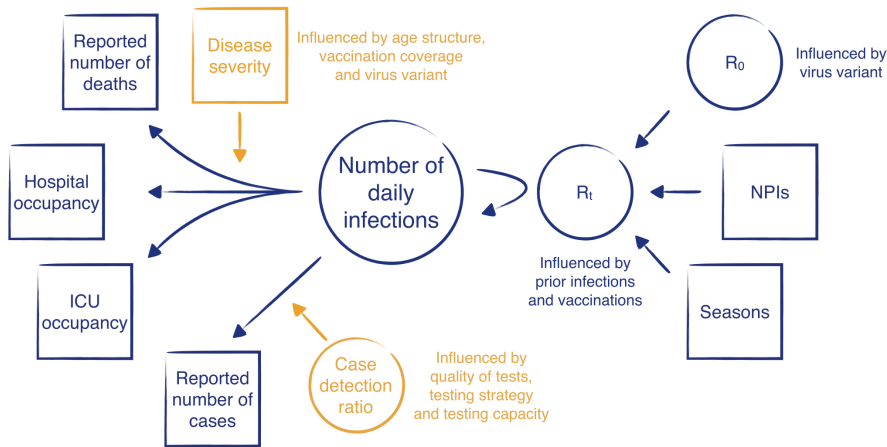
	<b>Minimal requirements</b>	<b>Acceptable</b>	<b>Ideal approach</b>
<b>Eligibility</b>	Detailed comparison of patient characteristics with target population and across centres	Inverse probability weighting or multilevel regression modelling with post-stratification	Combination of data with more representative data sources
<b>Data quality</b>	Extensive data quality checks and reporting of data pre-processing and missing patterns	Audits, standardisation and training in data collection, validation data to quantify errors	Account for complex measurement error and informative missing data
<b>Time point alignment</b>	Report exact timing of all measurements and of treatment trajectories	Target trial emulation	Validate accuracy of time stamps, report reasons for treatment switches
<b>Interventions and tests not random</b>	Directed acyclic graph including unmeasured and mismeasured confounders	Quantitative bias analysis and falsification endpoints of negative controls	Combined analysis with RCT, high quality documentation of treatment decisions
<b>Multiplicity of analysis</b>	Pre-registration of	Multi-analyst studies or extensive multiverse	Uncertainty intervals that account for

# Sources of information in the modelling of COVID-19





[16, 17]



[16, 17]

# A hierarchical model of COVID-19 propagation

- The disease model:

$$C_{t,m} = \sum_{u < t} I_{u,m} (F_{\xi^c}(t - u + 1) - F_{\xi^c}(t - u))$$

# A hierarchical model of COVID-19 propagation

- The disease model:

$$C_{t,m} = \sum_{u < t} I_{u,m} (F_{\xi^C}(t - u + 1) - F_{\xi^C}(t - u))$$

- The death model:

$$D_{t,m} \sim$$

Negative Binomial  $(\pi_m^D \sum_{u < t} C_{u,m} (F_{\xi^D}(t - u + 1) - F_{\xi^D}(t - u)), \phi_d)$

# A hierarchical model of COVID-19 propagation

- The disease model:

$$C_{t,m} = \sum_{u < t} I_{u,m} (F_{\xi^C}(t - u + 1) - F_{\xi^C}(t - u))$$

- The death model:

$$D_{t,m} \sim$$

Negative Binomial  $(\pi_m^D \sum_{u < t} C_{u,m} (F_{\xi^D}(t - u + 1) - F_{\xi^D}(t - u)), \phi_d)$

- The reporting model:

$$C_{t,m}^R \sim$$

Negative Binomial  $(\rho_{t,m} \sum_{u < t} C_{u,m} (F_{\xi^R}(t - u + 1) - F_{\xi^R}(t - u)), \phi_c)$

# A hierarchical model of COVID-19 propagation

- The disease model:

$$C_{t,m} = \sum_{u < t} I_{u,m} (F_{\xi^C}(t - u + 1) - F_{\xi^C}(t - u))$$

- The death model:

$$D_{t,m} \sim$$

$$\text{Negative Binomial} \left( \pi_m^D \sum_{u < t} C_{u,m} (F_{\xi_m^D}(t - u + 1) - F_{\xi_m^D}(t - u)), \phi_d \right)$$

- The reporting model:

$$C_{t,m}^R \sim$$

$$\text{Negative Binomial} \left( \rho_{t,m} \sum_{u < t} C_{u,m} (F_{\xi_m^R}(t - u + 1) - F_{\xi_m^R}(t - u)), \phi_c \right)$$

# A hierarchical model of COVID-19 propagation

- The disease model:

$$C_{t,m} = \sum_{u < t} I_{u,m} (F_{\xi^C}(t-u+1) - F_{\xi^C}(t-u))$$

- The death model:

$$D_{t,m} \sim$$

$$\text{Negative Binomial} \left( \pi_m^D \sum_{u < t} C_{u,m} (F_{\xi_m^D}(t-u+1) - F_{\xi_m^D}(t-u)), \phi_d \right)$$

- The reporting model:

$$C_{t,m}^R \sim$$

$$\text{Negative Binomial} \left( \rho_{t,m} \sum_{u < t} C_{u,m} (F_{\xi_m^R}(t-u+1) - F_{\xi_m^R}(t-u)), \phi_c \right)$$

- The hospitalization model:

$$H_{t,m} \sim$$

$$\text{Negative Binomial} \left( \pi_m^H \sum_{u < t} C_{u,m} (F_{\xi^H}(t-u+1) - F_{\xi^H}(t-u)), \phi_h \right)$$

# A hierarchical model of COVID-19 propagation

- The death model:

$$D_{t,m} \sim$$

$$\text{Negative Binomial} \left( \pi_m^D \sum_{u < t} C_{u,m} (F_{\xi_m^D}(t - u + 1) - F_{\xi_m^D}(t - u)), \phi_d \right)$$

- The reporting model:

$$C_{t,m}^R \sim$$

$$\text{Negative Binomial} \left( \rho_{t,m} \sum_{u < t} C_{u,m} (F_{\xi_m^R}(t - u + 1) - F_{\xi_m^R}(t - u)), \phi_c \right)$$

- The hospitalization model:

$$H_{t,m} \sim$$

$$\text{Negative Binomial} \left( \pi_m^H \sum_{u < t} C_{u,m} (F_{\xi_m^H}(t - u + 1) - F_{\xi_m^H}(t - u)), \phi_h \right)$$

- The renewal model:

$$I_{t,m} \sim$$

$$\text{Negative Binomial} \left( R_{t,m} \sum_{u < t} I_{u,m} (F_{\gamma}(t - u + 1) - F_{\gamma}(t - u)), \phi_i \right)$$

# A hierarchical model of COVID-19 propagation

- The death model:

$$D_{t,m} \sim$$

$$\text{Negative Binomial} \left( \pi_m^D \sum_{u < t} C_{u,m} (F_{\xi_m^D}(t-u+1) - F_{\xi_m^D}(t-u)), \phi_d \right)$$

- The renewal model:

$$I_{t,m} \sim$$

$$\text{Negative Binomial} \left( R_{t,m} \sum_{u < t} I_{u,m} (F_\gamma(t-u+1) - F_\gamma(t-u)), \phi_i \right)$$

$$R_{t,m} = R_m^0 \cdot \exp \left( - \sum_{k=1}^K \alpha_{k,m} \ln_{k,t,m} \right)$$

$$R_m^0 \sim \mathcal{N}(R^0, \sigma_R)$$

$$\alpha_{k,m} \sim \mathcal{N}(\alpha_k, \sigma_{\alpha_k})$$

# A hierarchical model of COVID-19 propagation

- The death model:

$$D_{t,m} \sim$$

$$\text{Negative Binomial} \left( \pi_m^D \sum_{u < t} C_{u,m} (F_{\xi_m^D}(t - u + 1) - F_{\xi_m^D}(t - u)), \phi_d \right)$$

- The renewal model:

$$I_{t,m} \sim \text{Negative Binomial}(\tau_m, \phi_i) \text{ for } t = 1$$

$$I_{t,m} \sim$$

$$\text{Negative Binomial} \left( R_{t,m} \sum_{u < t} I_{u,m} (F_\gamma(t - u + 1) - F_\gamma(t - u)), \phi_i \right)$$

$$R_{t,m} = R_m^0 \cdot \exp \left( - \sum_{k=1}^K \alpha_{k,m} \ln_{k,t,m} \right)$$

# A hierarchical model of COVID-19 propagation

- The death model:

$$D_{t,m} \sim$$

$$\text{Negative Binomial} \left( \pi_m^D \sum_{u < t} C_{u,m} (F_{\xi_m^D}(t-u+1) - F_{\xi_m^D}(t-u)), \phi_d \right)$$

- The renewal model:

$$I_{t,m} \sim$$

$$\text{Negative Binomial} \left( R_{t,m} \sum_{u < t} I_{u,m} (F_\gamma(t-u+1) - F_\gamma(t-u)), \phi_i \right)$$

$$R_{t,m} = R_{t,m}^0 \cdot \exp \left( - \sum_{k=1}^K \alpha_{k,m} \ln_{k,t,m} \right)$$

$$R_{t,m}^0 = R_m^0 \cdot (1 - p_{t,m}^\alpha - p_{t,m}^\delta) + (1 + \beta^\alpha) \cdot R_m^0 \cdot p_{t,m}^\alpha \\ + (1 + \beta^\delta) \cdot R_m^0 \cdot p_{t,m}^\delta$$

# A hierarchical model of COVID-19 propagation

- The death model:

$$D_{t,m} \sim$$

$$\text{Negative Binomial} \left( \pi_{m,t}^D \sum_{u < t} C_{u,m} (F_{\xi_m^D}(t-u+1) - F_{\xi_m^D}(t-u)), \phi_d \right)$$

- The renewal model:

$$I_{t,m} \sim$$

$$\text{Negative Binomial} \left( R_{t,m} \sum_{u < t} I_{u,m} (F_{\gamma}(t-u+1) - F_{\gamma}(t-u)), \phi_i \right)$$

$$R_{t,m} = R_{t,m}^0 \cdot \exp \left( - \sum_{k=1}^K \alpha_{k,m} \ln_{k,t,m} \right) \cdot (1 - c_{t,m}^1 - c_{t,m}^2 \cdot (1 - c_{t,m}^1))$$

$$c_{t,m}^1 = \frac{\sum_{u < t} I_{u,m}}{N_m} \cdot (1 - \beta^{\text{reinf}}) \text{ and}$$

$$c_{t,m}^2 = \frac{\sum_{u < t} \text{Vacc}_{u,m}^1 \cdot \beta^{v1} + \text{Vacc}_{u,m}^2 \cdot \beta^{v2}}{N_m}$$

# A hierarchical model of COVID-19 propagation

- The death model:

$$D_{t,m} \sim$$

$$\text{Negative Binomial} \left( \pi_{m,t}^D \sum_{u < t} C_{u,m} (F_{\xi_m^D}(t-u+1) - F_{\xi_m^D}(t-u)), \phi_d \right)$$

- The renewal model:

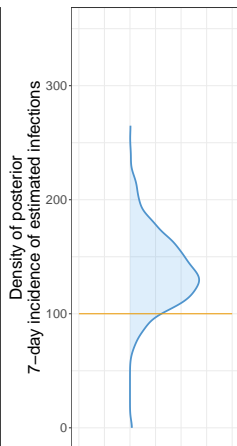
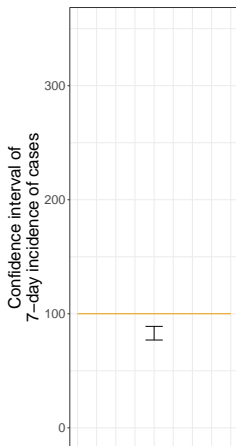
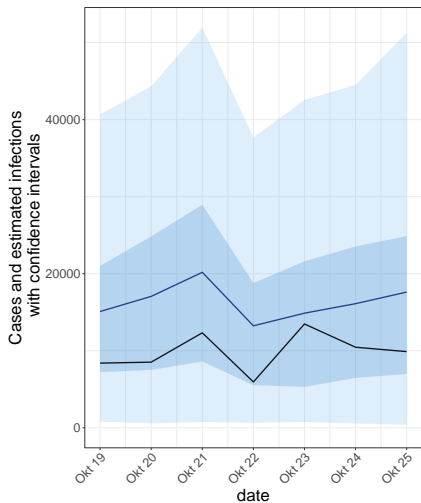
$$I_{t,m} \sim$$

$$\text{Negative Binomial} \left( R_{t,m} \sum_{u < t} I_{u,m} (F_{\gamma}(t-u+1) - F_{\gamma}(t-u)), \phi_i \right)$$

$$R_{t,m} = R_{t,m}^0 \cdot \exp \left( - \sum_{k=1}^K \alpha_{k,m} \ln_{k,t,m} \right) \cdot (1 - c_{t,m}^1 - c_{t,m}^2 \cdot (1 - c_{t,m}^1))$$

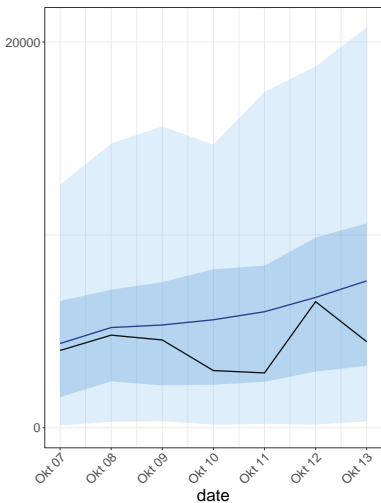
$$c_{t,m}^1 = \frac{\sum_{u < t} I_{u,m}}{N_m} \cdot (1 - \beta^{\text{reinf}}) \text{ and}$$

$$c_{t,m}^2 = \frac{\sum_{u < t} \text{Vacc}_{u,m}^1 \cdot \beta^{v1} + \text{Vacc}_{u,m}^2 \cdot \beta^{v2}}{N_m}$$

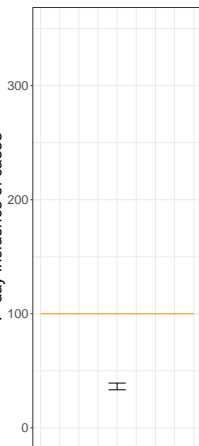


[16, 17]

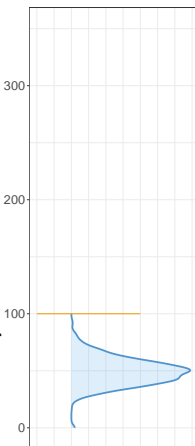
Cases and estimated infections with confidence intervals



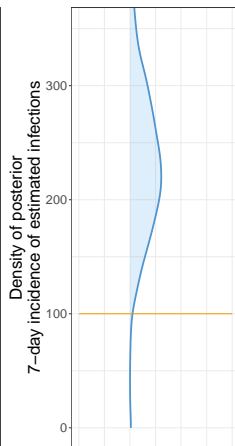
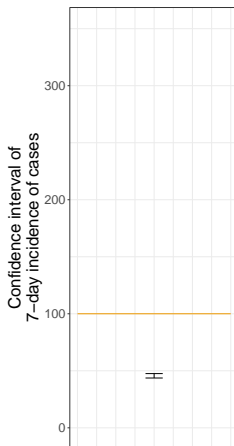
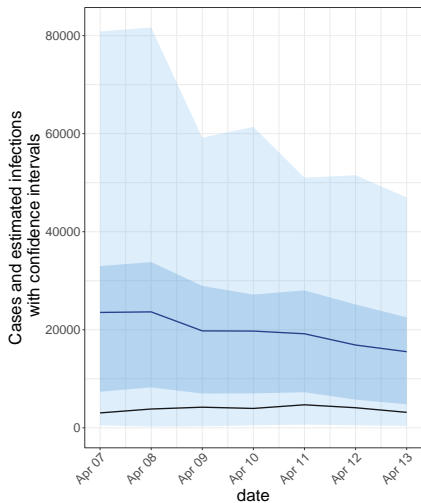
Confidence interval of 7-day incidence of cases



Density of posterior 7-day incidence of estimated infections



[16, 17]



[16, 17]