

When likelihood goes wrong

Nancy Reid
University of Toronto

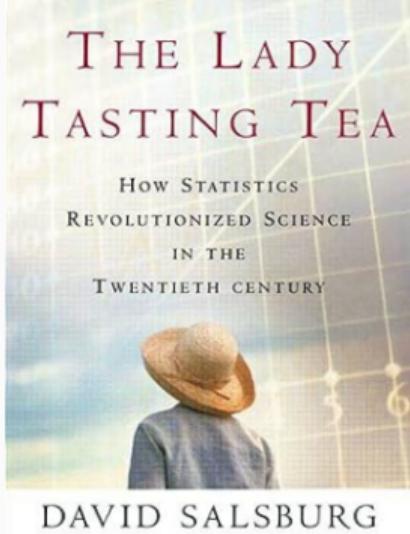
October 2 2024

Pfizer Colloquium

UConn | UNIVERSITY OF CONNECTICUT







"by the end of the twentieth century, almost all of science had shifted to using statistical models"

Salsburg, 2001

Outline

1. Models and science: a haphazard selection
2. Models and parameters
3. What's special about likelihood?
4. Some approaches to misspecification
5. Example
6. Conclusion

Models and science: a haphazard selection

Clim. Past, 20, 1387–1399, 2024
<https://doi.org/10.5194/cp-20-1387-2024>
© Author(s) 2024. This work is distributed under
the Creative Commons Attribution 4.0 License.



Climate
of the Past
Open Access
EGU

600 years of wine must quality and April to August temperatures in western Europe 1420–2019

Christian Pfister¹, Stefan Brönnimann², Andres Altweig³, Rudolf Brázdič⁴, Laurent Litzenburger⁵, Daniele Lorusso⁶,
and Thomas Pliemon⁷

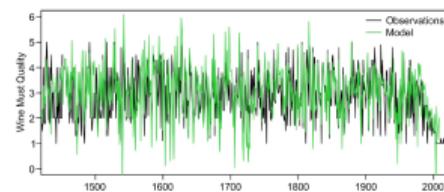


Figure 4. Observed series of wine quality (average; black) from 1420 to 2019 and series obtained with a statistical model calibrated in 1781–1800 (green). The model is explained in Sect. 3.

Clim. Past, 20, 1387–1399, 2024
<https://doi.org/10.5194/cp-20-1387-2024>
© Author(s) 2024. This work is distributed under
the Creative Commons Attribution 4.0 License.



Climate
of the Past
Open Access
EGU

600 years of wine must quality and April to August temperatures in western Europe 1420–2019

Christian Pfister¹, Stefan Brönnimann², Andres Altweig³, Rudolf Brázdil⁴, Laurent Litzenburger⁵, Daniele Lorusso⁶,
and Thomas Pliemon⁷

Scientific question: Can historical records of wine quality be used
as temperature proxies?

observational data

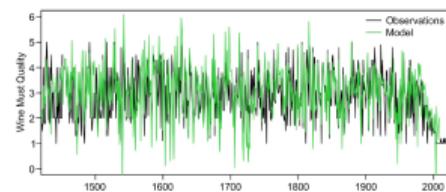


Figure 4. Observed series of wine quality (average; black) from 1420 to 2019 and series obtained with a statistical model calibrated in 1781–1800 (green). The model is explained in Sect. 3.

Clim. Past, 20, 1387–1399, 2024
<https://doi.org/10.5194/cp-20-1387-2024>
© Author(s) 2024. This work is distributed under
the Creative Commons Attribution 4.0 License.



Climate
of the Past
Open Access
EGU

600 years of wine must quality and April to August temperatures in western Europe 1420–2019

Christian Pfister¹, Stefan Brönnimann², Andres Altweig³, Rudolf Brázdič⁴, Laurent Litzenburger⁵, Daniele Lorusso⁶,
and Thomas Pliemon⁷

Scientific question: Can historical records of wine quality be used
as temperature proxies?

observational data

Statistical model: “we used a statistical [linear regression] model for wine quality
based on local temperature and precipitation”

yes, if used carefully

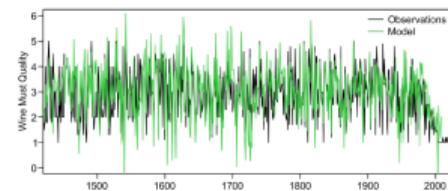


Figure 4. Observed series of wine quality (average; black) from 1420 to 2019 and series obtained with a statistical model calibrated in 1781–1800 (green). The model is explained in Sect. 3.

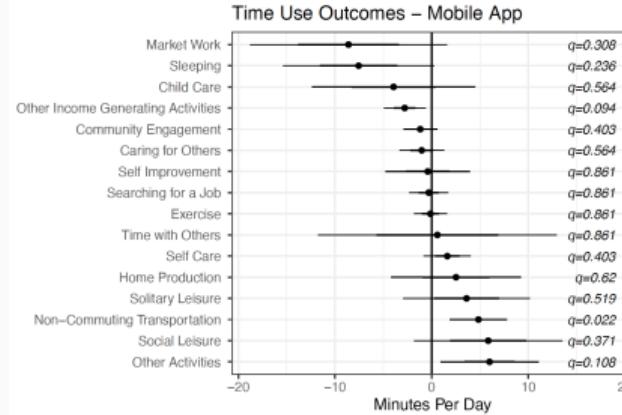
NBER WORKING PAPER SERIES

THE EMPLOYMENT EFFECTS OF A GUARANTEED INCOME: EXPERIMENTAL EVIDENCE FROM TWO U.S. STATES

Eva Vivalt
Elizabeth Rhodes
Alexander W. Bartik
David E. Broockman
Sarah Miller

Working Paper 32719
<http://www.nber.org/papers/w32719>

Figure 5: Time Use Results: Mobile App



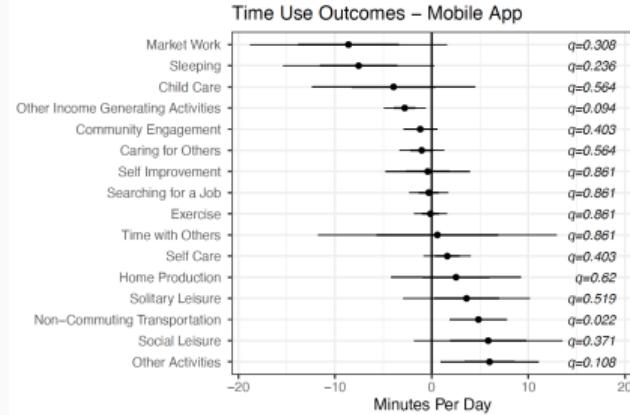
NBER WORKING PAPER SERIES

THE EMPLOYMENT EFFECTS OF A GUARANTEED INCOME: EXPERIMENTAL EVIDENCE FROM TWO U.S. STATES

Eva Vivalt
Elizabeth Rhodes
Alexander W. Bartik
David E. Broockman
Sarah Miller

Working Paper 32719
<http://www.nber.org/papers/w32719>

Figure 5: Time Use Results: Mobile App



Scientific questions: Does guaranteed income supplement affect labor market measures?

randomized controlled trial

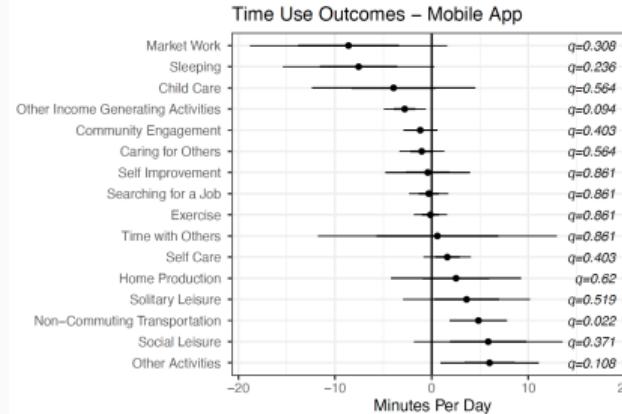
NBER WORKING PAPER SERIES

THE EMPLOYMENT EFFECTS OF A GUARANTEED INCOME: EXPERIMENTAL EVIDENCE FROM TWO U.S. STATES

Eva Vivalt
Elizabeth Rhodes
Alexander W. Bartik
David E. Broockman
Sarah Miller

Working Paper 32719
<http://www.nber.org/papers/w32719>

Figure 5: Time Use Results: Mobile App



Scientific questions: Does guaranteed income supplement affect labor market measures?

randomized controlled trial

Statistical model: $Y_i = \alpha + \beta Treated_i + \gamma^T X_i + \epsilon_i$

“support for both sides of this debate”

JAMA Oncology | Original Investigation

Bilateral Mastectomy and Breast Cancer Mortality

Vasily Giannakeas, PhD, MPH; David W. Lim, MDCM, MEd, PhD; Steven A. Narod, MD

IMPORTANCE The benefit of bilateral mastectomy for women with unilateral breast cancer in terms of deaths from breast cancer has not been shown.

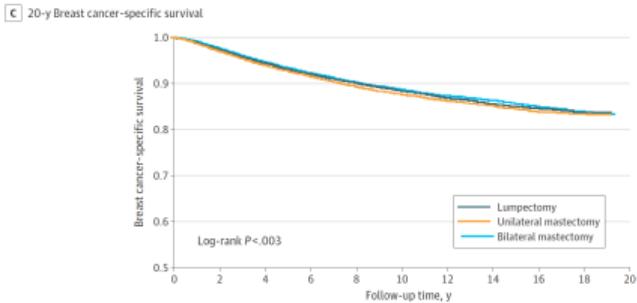
OBJECTIVES To estimate the 20-year cumulative risk of breast cancer mortality among women with stage 0 to stage III unilateral breast cancer according to the type of initial surgery performed.

DESIGN, SETTINGS, AND PARTICIPANTS This cohort study used the Surveillance, Epidemiology, and End Results (SEER) Program registry database to identify women with unilateral breast cancer (invasive and ductal carcinoma in situ) who were diagnosed from 2000 to 2019. Three closely matched cohorts of equal size were generated using 1:1:1 matching according to surgical approach. The cohorts were followed up for 20 years for contralateral breast cancer and for breast cancer mortality. The analysis compared the 20-year cumulative risk of breast cancer mortality for women treated with lumpectomy vs unilateral mastectomy vs bilateral mastectomy. Data were analyzed from October 2023 to February 2024.

EXPOSURES Type of breast surgery performed (lumpectomy, unilateral mastectomy, or bilateral mastectomy).

Editorial

Supplemental content



JAMA Oncology | Original Investigation

Bilateral Mastectomy and Breast Cancer Mortality

Vasily Giannakeas, PhD, MPH; David W. Lim, MDCM, MEd, PhD; Steven A. Narod, MD

IMPORTANCE The benefit of bilateral mastectomy for women with unilateral breast cancer in terms of deaths from breast cancer has not been shown.

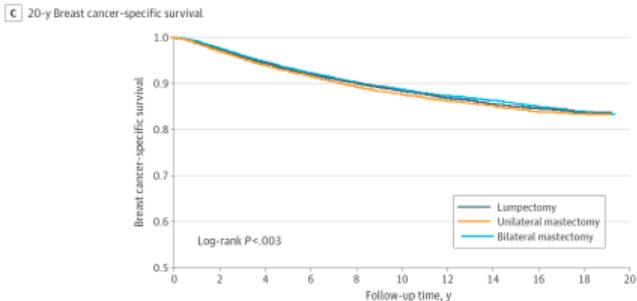
OBJECTIVES To estimate the 20-year cumulative risk of breast cancer mortality among women with stage 0 to stage III unilateral breast cancer according to the type of initial surgery performed.

DESIGN, SETTINGS, AND PARTICIPANTS This cohort study used the Surveillance, Epidemiology, and End Results (SEER) Program registry database to identify women with unilateral breast cancer (invasive and ductal carcinoma in situ) who were diagnosed from 2000 to 2019. Three closely matched cohorts of equal size were generated using 1:1 matching according to surgical approach. The cohorts were followed up for 20 years for contralateral breast cancer and for breast cancer mortality. The analysis compared the 20-year cumulative risk of breast cancer mortality for women treated with lumpectomy vs unilateral mastectomy vs bilateral mastectomy. Data were analyzed from October 2023 to February 2024.

EXPOSURES Type of breast surgery performed (lumpectomy, unilateral mastectomy, or bilateral mastectomy).

Editorial

Supplemental content



Scientific question: Does bilateral mastectomy for unilateral breast cancer improve 20-year survival?

matched case-cohort study

JAMA Oncology | Original Investigation

Bilateral Mastectomy and Breast Cancer Mortality

Vasily Giannakeas, PhD, MPH; David W. Lim, MDCM, MEd, PhD; Steven A. Narod, MD

IMPORTANCE The benefit of bilateral mastectomy for women with unilateral breast cancer in terms of deaths from breast cancer has not been shown.

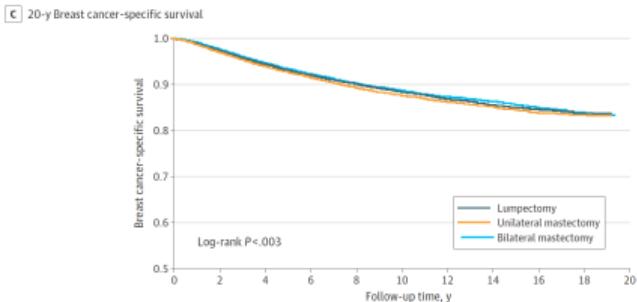
OBJECTIVES To estimate the 20-year cumulative risk of breast cancer mortality among women with stage 0 to stage III unilateral breast cancer according to the type of initial surgery performed.

DESIGN, SETTINGS, AND PARTICIPANTS This cohort study used the Surveillance, Epidemiology, and End Results (SEER) Program registry database to identify women with unilateral breast cancer (invasive and ductal carcinoma in situ) who were diagnosed from 2000 to 2019. Three closely matched cohorts of equal size were generated using 1:1 matching according to surgical approach. The cohorts were followed up for 20 years for contralateral breast cancer and for breast cancer mortality. The analysis compared the 20-year cumulative risk of breast cancer mortality for women treated with lumpectomy vs unilateral mastectomy vs bilateral mastectomy. Data were analyzed from October 2023 to February 2024.

EXPOSURES Type of breast surgery performed (lumpectomy, unilateral mastectomy, or bilateral mastectomy).

Editorial

Supplemental content



Scientific question: Does bilateral mastectomy for unilateral breast cancer improve 20-year survival?
matched case-cohort study

Statistical model: “We used the Kaplan-Meier method to estimate survival”

“preemptive surgery did not appear to reduce the risk of death”

Article

<https://doi.org/10.1038/s41550-023-01983-1>

Variability of extragalactic X-ray jets on kiloparsec scales

Received: 17 May 2022

Accepted: 27 April 2023

Published online: 29 May 2023

Eileen T. Meyer^①, Aamil Shaik¹, Yanbo Tang², Nancy Reid³, Karthik Reddy^{1,4}, Peter Breiding⁵, Markos Georganopoulos¹, Marco Chiaberge^{6,7}, Eric Perlman⁷, Devon Clautice⁷, William Sparks^{8,9}, Nat DeNigris^{1,10} & Max Trevor^{1,11}

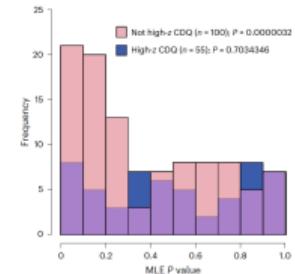


Fig. 3 | Histogram of the single-region P values from the directional test, not adjusted for multiple comparison. In pink, the subset of sources that

Yanbo Tang

Article

<https://doi.org/10.1038/s41550-023-01983-1>

Variability of extragalactic X-ray jets on kiloparsec scales

Received: 17 May 2022

Accepted: 27 April 2023

Published online: 29 May 2023

Eileen T. Meyer^①, Aamil Shaik¹, Yanbo Tang², Nancy Reid³, Karthik Reddy^{1,4}, Peter Breiding⁵, Markos Georganopoulos¹, Marco Chiaberge^{6,7}, Eric Perlman⁷, Devon Clautice⁷, William Sparks^{8,9}, Nat DeNigris^{1,10} & Max Trevor^{1,11}

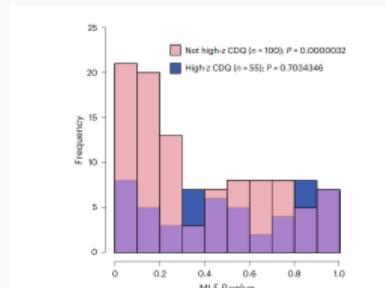


Fig. 3 | Histogram of the single-region P values from the directional test, not adjusted for multiple comparison. In pink, the subset of sources that

Yanbo Tang

Scientific question: Are observations of X-ray jets consistent with current theory?

observational data

Article

<https://doi.org/10.1038/s41550-023-01983-1>

Variability of extragalactic X-ray jets on kiloparsec scales

Received: 17 May 2022

Accepted: 27 April 2023

Published online: 29 May 2023

Eileen T. Meyer^①, Aamil Shaik¹, Yanbo Tang², Nancy Reid³, Karthik Reddy^{1,4}, Peter Breiding⁵, Markos Georganopoulos¹, Marco Chiaberge^{6,7}, Eric Perlman⁷, Devon Clautice⁷, William Sparks^{8,9}, Nat DeNigris^{1,10} & Max Trevor^{1,11}

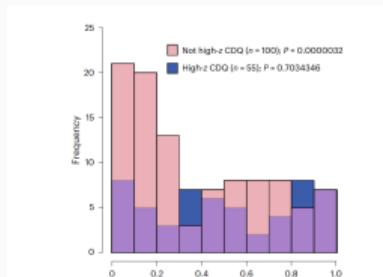


Fig. 3 | Histogram of the single-region P values from the directional test, not adjusted for multiple comparison. In pink, the subset of sources that

Yanbo Tang

Scientific question: Are observations of X-ray jets consistent with current theory?

observational data

Statistical model: compare background and source counts using Poisson distribution:

$$x_i \sim Po(a_i \beta_i), \quad y_i \sim Po(b_i \beta_i + b_i f_i \mu_i), \quad H_0 : \mu_i \equiv 0$$

Article

<https://doi.org/10.1038/s41550-023-01983-1>

Variability of extragalactic X-ray jets on kiloparsec scales

Received: 17 May 2022

Accepted: 27 April 2023

Published online: 29 May 2023

Eileen T. Meyer¹ , Aamil Shaik¹, Yanbo Tang², Nancy Reid³, Karthik Reddy^{1,4}, Peter Breiding⁵, Markos Georganopoulos¹, Marco Chiaberge^{6,7}, Eric Perlman⁷, Devon Clautice⁷, William Sparks^{8,9}, Nat DeNigris^{1,10} & Max Trevor^{1,11}

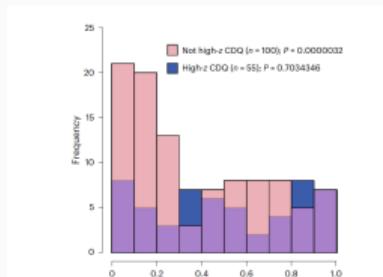


Fig. 3 | Histogram of the single-region P values from the directional test, not adjusted for multiple comparison. In pink, the subset of sources that

Yanbo Tang

Scientific question: Are observations of X-ray jets consistent with current theory?

observational data

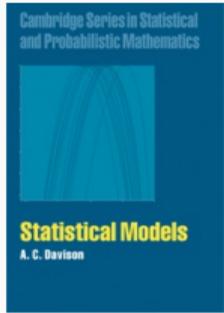
Statistical model: compare background and source counts using Poisson distribution:

$$x_i \sim Po(a_i \beta_i), \quad y_i \sim Po(b_i \beta_i + b_i f_i \mu_i), \quad H_0 : \mu_i \equiv 0$$

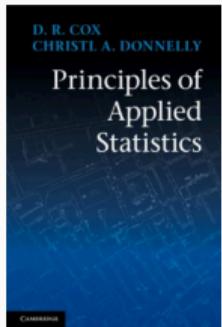
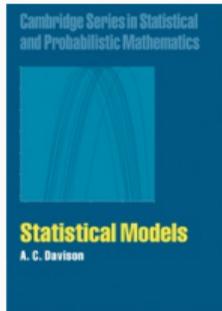
Models and parameters

Why these models?

- standard in the literature of that field income
- standard in the publications of that lab breast cancer
- follow some prescription:
 - binary response — use logistic regression
 - time to event — use PH model
 - time series — use ARMA wine
 - repeated measures — use random effects
 - ...
- motivated by theory: economic, physical, ... X-ray jets



- the key feature of a statistical model is that variability is represented using probability distributions
- the art of modelling lies in finding a balance that enables the questions at hand to be answered or new ones posed



- the key feature of a statistical model is that variability is represented using probability distributions
- the art of modelling lies in finding a balance that enables the questions at hand to be answered or new ones posed
- probability models as an aid to the interpretation of data
- perturbations of no intrinsic interest distort an otherwise exact measurement
- substantial natural variability in the phenomenon under study

Statistical Science
1990, Vol. 5, No. 2, 163–168

Model Specification: The Views of Fisher and Neyman, and Later Developments

E. L. Lehmann

empirical, or predictive models, contrasted with explanatory models

Role of Models in Statistical Analysis

D. R. Cox

Abstract. A number of distinct roles are identified for probability models used in the analysis of data. Examples are outlined. Some general issues arising in the formulation of such models are discussed.

indirect models

Statistical Science
1990, Vol. 5, No. 2, 163–168

Model Specification: The Views of Fisher and Neyman, and Later Developments

E. L. Lehmann

empirical, or predictive models, contrasted with explanatory models

indirect models

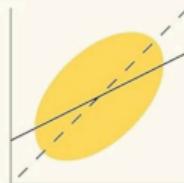
Role of Models in Statistical Analysis

D. R. Cox

Abstract. A number of distinct roles are identified for probability models used in the analysis of data. Examples are outlined. Some general issues arising in the formulation of such models are discussed.

Statistical Models

Theory and Practice
REVISED EDITION



David A. Freedman

The emphasis throughout is on the connection
– or lack of connection –
between the models and the real phenomena.

The role of parameters

- probability models very likely be parameterized
- thus defining a class of models
- parameters may be finite- or infinite-dimensional

$$\{f(y; \theta); \theta \in \Theta\}$$

parametric vs nonparametric

The role of parameters

- probability models very likely be parameterized
- thus defining a class of models $\{f(y; \theta); \theta \in \Theta\}$
- parameters may be finite- or infinite-dimensional parametric vs nonparametric
- ideally one or more parameters represent key aspects of the model for the application at hand
- other parameters complete the specification
- the meaning of various parameters varies with the application

The role of parameters

- probability models very likely be parameterized
- thus defining a class of models $\{f(y; \theta); \theta \in \Theta\}$
- parameters may be finite- or infinite-dimensional parametric vs nonparametric
- ideally one or more parameters represent key aspects of the model for the application at hand
- other parameters complete the specification
- the meaning of various parameters varies with the application
- this sounds simpler than it is

The role of parameters

- probability models very likely be parameterized
 - thus defining a class of models
 - parameters may be finite- or infinite-dimensional
- $\{f(y; \theta); \theta \in \Theta\}$
parametric vs nonparametric
- ideally one or more parameters represent key aspects of the model
 - other parameters complete the specification
 - the meaning of various parameters varies with the application
- for the application at hand

The Annals of Statistics
2002, Vol. 30, No. 5, 1225–1310

WHAT IS A STATISTICAL MODEL?¹

BY PETER McCULLAGH

University of Chicago

- this sounds simpler than it is

The role of parameters

- probability models very likely be parameterized
 - thus defining a class of models
 - parameters may be finite- or infinite-dimensional
 - ideally one or more parameters represent key aspects of the model
 - other parameters complete the specification
 - the meaning of various parameters varies with the application
- $\{f(y; \theta); \theta \in \Theta\}$
parametric vs nonparametric
for the application at hand

The Annals of Statistics
2002, Vol. 30, No. 5, 1225–1310

WHAT IS A STATISTICAL MODEL?¹

BY PETER McCULLAGH

University of Chicago

- this sounds simpler than it is

What's special about likelihood?

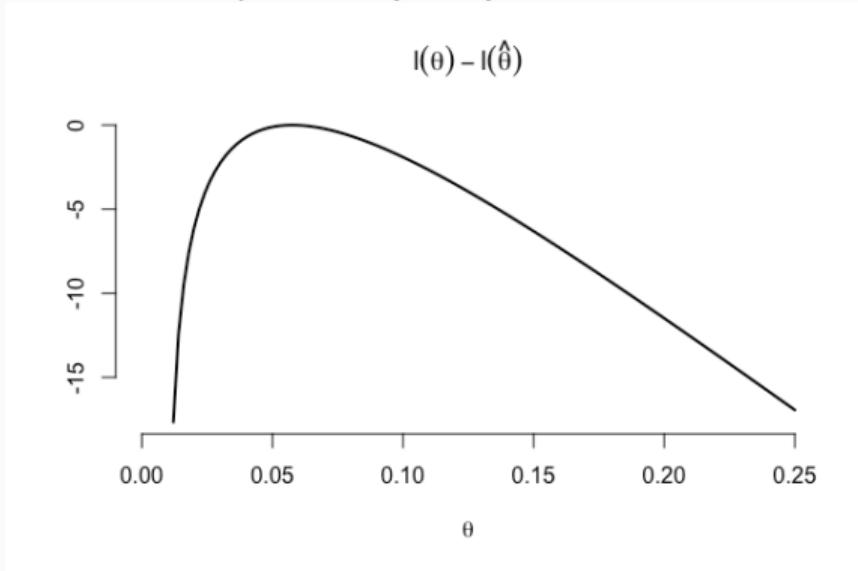
The likelihood function

- it puts the emphasis on the model: $L(\theta; y) \propto f(y; \theta)$ inverse problem
- provides a convenient way to compare parameter values e.g. $L(\theta)/L(\hat{\theta})$

The likelihood function

- it puts the emphasis on the model: $L(\theta; y) \propto f(y; \theta)$
- provides a convenient way to compare parameter values

inverse problem
e.g. $L(\theta)/L(\hat{\theta})$

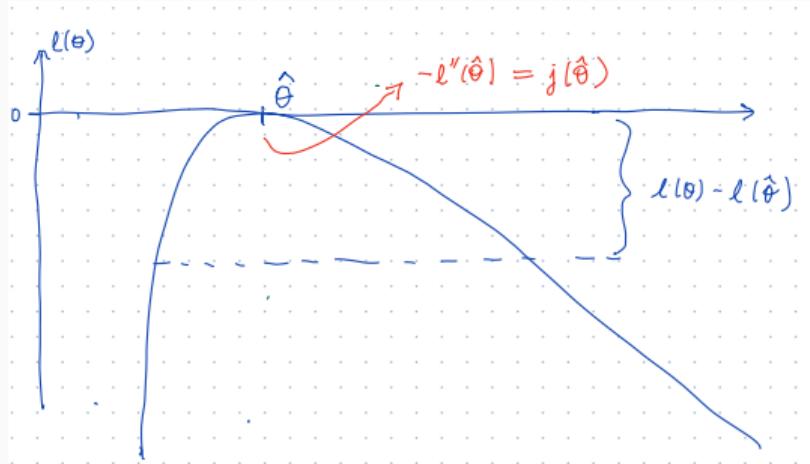


Pfizer vaccine
 $Bin(162 + 8, \theta)$ via 2 Poissons

The likelihood function

- it puts the emphasis on the model: $L(\theta; y) \propto f(y; \theta)$
- provides a convenient way to compare parameter values
- provides **reliable** summary measures

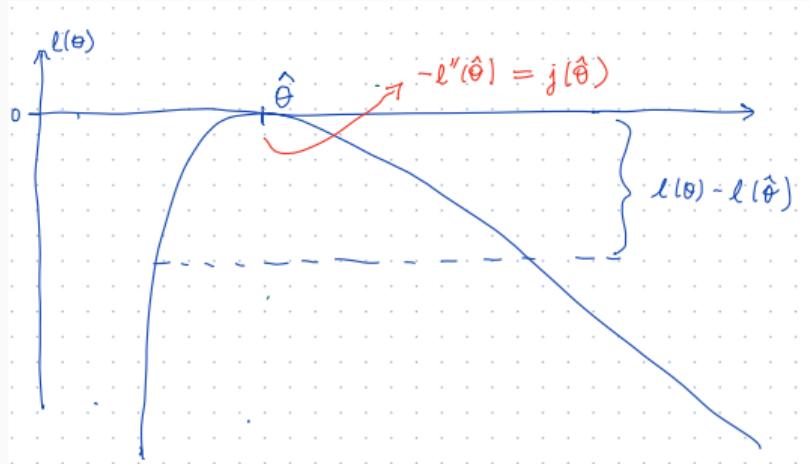
inverse problem
e.g. $L(\theta)/L(\hat{\theta})$



The likelihood function

- it puts the emphasis on the model: $L(\theta; y) \propto f(y; \theta)$
- provides a convenient way to compare parameter values
- provides **reliable** summary measures

inverse problem
e.g. $L(\theta)/L(\hat{\theta})$



- can be converted to a probability, given a prior probability for θ

Inference

- Data and model: $y = y_1, \dots, y_n$ independent; $Y \sim f(y; \theta)$, $\theta \in \mathbb{R}^p$ $y \in \mathbb{R}^n$
- Likelihood function: $L(\theta; y) \propto f(y; \theta) = \prod_{i=1}^n f(y_i; \theta)$
- Log-likelihood function: $\ell(\theta; y) = \log L(\theta; y) = \sum_{i=1}^n \log f(y_i; \theta)$

Inference

- Data and model: $y = y_1, \dots, y_n$ independent; $Y \sim f(y; \theta)$, $\theta \in \mathbb{R}^p$ $y \in \mathbb{R}^n$
- Likelihood function: $L(\theta; y) \propto f(y; \theta) = \prod_{i=1}^n f(y_i; \theta)$
- Log-likelihood function: $\ell(\theta; y) = \log L(\theta; y) = \sum_{i=1}^n \log f(y_i; \theta)$
- Maximum likelihood estimator: $\hat{\theta} = \hat{\theta}(y) = \arg \sup_{\theta} \ell(\theta; y)$ $\ell'(\hat{\theta}) = 0$
- Score function $\ell'(\theta; y) = \sum_{i=1}^n \partial \ell(\theta; y_i) / \partial \theta^T$ $i(\theta) = \text{var}\{\partial \ell(\theta; Y) / \partial \theta^T\}$ CLT
- observed Fisher information $j(\theta) = -\partial^2 \ell(\theta; y) / \partial \theta \partial \theta^T$ $E\{j(\theta)\} = i(\theta)$

Inference

- Data and model: $y = y_1, \dots, y_n$ independent; $Y \sim f(y; \theta)$, $\theta \in \mathbb{R}^p$ $y \in \mathbb{R}^n$
- Likelihood function: $L(\theta; y) \propto f(y; \theta) = \prod_{i=1}^n f(y_i; \theta)$
- Log-likelihood function: $\ell(\theta; y) = \log L(\theta; y) = \sum_{i=1}^n \log f(y_i; \theta)$
- Maximum likelihood estimator: $\hat{\theta} = \hat{\theta}(y) = \arg \sup_{\theta} \ell(\theta; y)$ $\ell'(\hat{\theta}) = 0$
- Score function $\ell'(\theta; y) = \sum_{i=1}^n \partial \ell(\theta; y_i) / \partial \theta^T$ $i(\theta) = \text{var}\{\partial \ell(\theta; Y) / \partial \theta^T\}$ CLT
- observed Fisher information $j(\theta) = -\partial^2 \ell(\theta; y) / \partial \theta \partial \theta^T$ $E\{j(\theta)\} = i(\theta)$
- in “large samples”

$$\ell'(\theta) \sim N_p\{\mathbf{0}, i(\theta)\}, \quad \hat{\theta} \sim N_p\{\theta, j^{-1}(\hat{\theta})\}, \quad 2\{\ell(\hat{\theta}) - \ell(\theta)\} \sim \chi_p^2$$

Too much notation!

CLT:

$$\frac{1}{\sqrt{n}} \ell'(\theta) \xrightarrow{d} N\{\mathbf{o}, i_1(\theta)\}$$

large-sample approx:

$$\ell'(\theta) \sim N_p\{\mathbf{o}, i(\theta)\}, \quad \hat{\theta} \sim N_p\{\theta, \mathbf{j}^{-1}(\hat{\theta})\}, \quad 2\{\ell(\hat{\theta}) - \ell(\theta)\} \sim \chi_p^2$$

Too much notation!

CLT:

$$\frac{1}{\sqrt{n}} \ell'(\theta) \xrightarrow{d} N\{\mathbf{0}, i_1(\theta)\}$$

large-sample approx:

$$\ell'(\theta) \sim N_p\{\mathbf{0}, i(\theta)\}, \quad \hat{\theta} \sim N_p\{\theta, j^{-1}(\hat{\theta})\}, \quad 2\{\ell(\hat{\theta}) - \ell(\theta)\} \sim \chi^2_p$$

Coefficients:

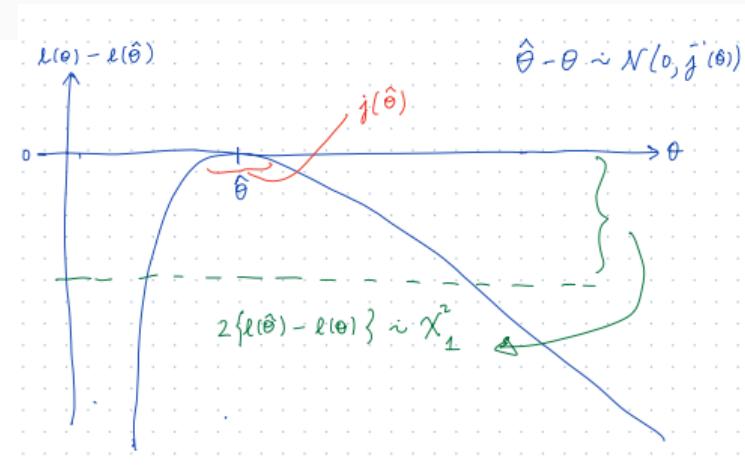
	Estimate	Std. Error	z value	Pr(> z)
(Intercept)	-3.079	0.987	-3.12	0.0018 **
aged1	-0.292	0.754	-0.39	0.6988
stage1	1.373	0.784	1.75	0.0799 .
grade1	0.872	0.816	1.07	0.2850
xray1	1.801	0.810	2.22	0.0263 *
acid1	1.684	0.791	2.13	0.0334 *

Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

(Dispersion parameter for binomial family taken to be 1)

Null deviance: 40.710 on 22 degrees of freedom

Residual deviance: 18.069 on 17 degrees of freedom



Aside: A bit too simple

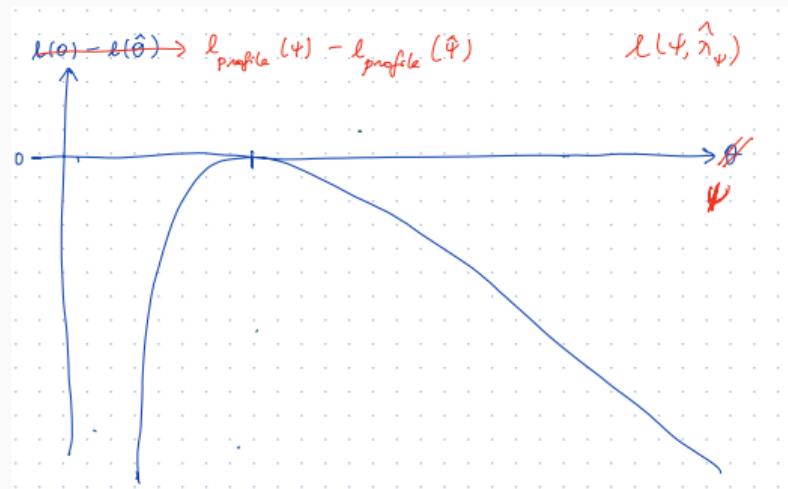
- model $f(y; \theta)$, $\theta \in \mathbb{R}^p$

Aside: A bit too simple

- model $f(y; \theta)$, $\theta \in \mathbb{R}^p$
- $\theta = (\psi, \lambda)$ parameters of interest nuisance parameters

Aside: A bit too simple

- model $f(y; \theta)$, $\theta \in \mathbb{R}^p$
- $\theta = (\psi, \lambda)$ parameters of interest nuisance parameters
- results above used profile log-likelihood function $\ell_{\text{profile}}(\psi) = \ell(\psi, \hat{\lambda}_\psi)$



What can go wrong?

- too many parameters

What can go wrong?

- too many parameters

$$p \sim n, \quad p/n \rightarrow C, \quad p/n \rightarrow \infty$$

What can go wrong?

- too many parameters $p \sim n, \quad p/n \rightarrow C, \quad p/n \rightarrow \infty$
- weird parameter space

What can go wrong?

- too many parameters
- weird parameter space

$$p \sim n, \quad p/n \rightarrow C, \quad p/n \rightarrow \infty$$

$$pf(y; \theta_1) + (1-p)f(y; \theta_2), \quad 0 \leq p \leq 1$$

What can go wrong?

- too many parameters $p \sim n, \quad p/n \rightarrow C, \quad p/n \rightarrow \infty$
- weird parameter space $pf(y; \theta_1) + (1 - p)f(y; \theta_2), \quad 0 \leq p \leq 1$
- computational intractability

What can go wrong?

- too many parameters $p \sim n, \quad p/n \rightarrow C, \quad p/n \rightarrow \infty$
- weird parameter space $pf(y; \theta_1) + (1 - p)f(y; \theta_2), \quad 0 \leq p \leq 1$
- computational intractability $L(\theta, \tau; y) = \int_{\mathbb{R}^k} f(y | z; \theta) f(z; \tau) dz$

What can go wrong?

- too many parameters $p \sim n, \quad p/n \rightarrow C, \quad p/n \rightarrow \infty$
- weird parameter space $pf(y; \theta_1) + (1 - p)f(y; \theta_2), \quad 0 \leq p \leq 1$
- computational intractability $L(\theta, \tau; y) = \int_{\mathbb{R}^k} f(y | z; \theta) f(z; \tau) dz$
- model is misspecified

What can go wrong?

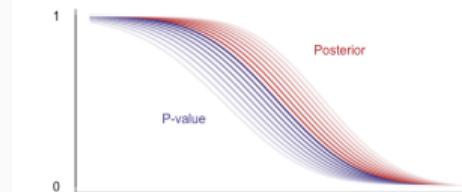
- too many parameters $p \sim n, \quad p/n \rightarrow C, \quad p/n \rightarrow \infty$
- weird parameter space $pf(y; \theta_1) + (1 - p)f(y; \theta_2), \quad 0 \leq p \leq 1$
- computational intractability $L(\theta, \tau; y) = \int_{\mathbb{R}^k} f(y | z; \theta) f(z; \tau) dz$
- model is misspecified true $Y \sim m(y), \quad \nexists \theta_* \text{ w. } f(\cdot; \theta_*) = m(\cdot)$

What can go wrong?

- too many parameters $p \sim n, \quad p/n \rightarrow C, \quad p/n \rightarrow \infty$
- weird parameter space $pf(y; \theta_1) + (1 - p)f(y; \theta_2), \quad 0 \leq p \leq 1$
- computational intractability $L(\theta, \tau; y) = \int_{\mathbb{R}^k} f(y | z; \theta) f(z; \tau) dz$
- model is misspecified true $Y \sim m(y), \quad \nexists \theta_* \text{ w. } f(\cdot; \theta_*) = m(\cdot)$
- likelihood is not a probability

What can go wrong?

- too many parameters $p \sim n, \quad p/n \rightarrow C, \quad p/n \rightarrow \infty$
- weird parameter space $pf(y; \theta_1) + (1 - p)f(y; \theta_2), \quad 0 \leq p \leq 1$
- computational intractability $L(\theta, \tau; y) = \int_{\mathbb{R}^k} f(y | z; \theta) f(z; \tau) dz$
- model is misspecified true $Y \sim m(y), \quad \nexists \theta_* \text{ w. } f(\cdot; \theta_*) = m(\cdot)$
- likelihood is not a probability



What can go wrong?

- too many parameters $p \sim n, \quad p/n \rightarrow C, \quad p/n \rightarrow \infty$
- weird parameter space $pf(y; \theta_1) + (1 - p)f(y; \theta_2), \quad 0 \leq p \leq 1$
- computational intractability $L(\theta, \tau; y) = \int_{\mathbb{R}^k} f(y | z; \theta) f(z; \tau) dz$
- model is misspecified $\text{true } Y \sim m(y), \quad \nexists \theta_* \text{ w. } f(\cdot; \theta_*) = m(\cdot)$
- likelihood is not a probability

Some approaches to misspecification

- true model $m(\mathbf{y})$ fitted model $f(\mathbf{y}; \theta)$ $\mathbf{y} = (y_1, \dots, y_n)$
- maximum likelihood estimator $\hat{\theta}$ $\ell(\theta; \mathbf{y}) \equiv \log f(\mathbf{y}; \theta)$
 $\hat{\theta} \equiv \arg \sup_{\theta} \ell(\theta; \mathbf{y})$

- true model $m(\mathbf{y})$ fitted model $f(\mathbf{y}; \theta)$ $\mathbf{y} = (y_1, \dots, y_n)$
- maximum likelihood estimator $\hat{\theta}$ $\ell(\theta; \mathbf{y}) \equiv \log f(\mathbf{y}; \theta)$
 $\hat{\theta} \equiv \arg \sup_{\theta} \ell(\theta; \mathbf{y})$
- $\hat{\theta}$ converges to the “closest true value” KL-divergence

$$\theta_m^0 = \arg \min_{\theta} \int m(\mathbf{y}) \log \left\{ \frac{m(\mathbf{y})}{f(\mathbf{y}; \theta)} \right\} d\mathbf{y}$$

- true model $m(\mathbf{y})$ fitted model $f(\mathbf{y}; \theta)$ $\mathbf{y} = (y_1, \dots, y_n)$
- maximum likelihood estimator $\hat{\theta}$ $\ell(\theta; \mathbf{y}) \equiv \log f(\mathbf{y}; \theta)$
 $\hat{\theta} \equiv \arg \sup_{\theta} \ell(\theta; \mathbf{y})$
- $\hat{\theta}$ converges to the “closest true value” KL-divergence

$$\theta_m^0 = \arg \min_{\theta} \int m(\mathbf{y}) \log \left\{ \frac{m(\mathbf{y})}{f(\mathbf{y}; \theta)} \right\} d\mathbf{y}$$

- $\hat{\theta}$ has asymptotic normal distribution, but is not fully efficient “sandwich variance”

$$\text{a.var.} (\hat{\theta}) = G^{-1}(\theta_m^0), \quad G(\theta) = J(\theta)I^{-1}(\theta)J(\theta)$$

$I = \text{var}_m(\ell')$, $J = \text{E}_m(-\ell'')$

- true model $m(\mathbf{y})$ fitted model $f(\mathbf{y}; \theta)$ $\mathbf{y} = (y_1, \dots, y_n)$
- maximum likelihood estimator $\hat{\theta}$ $\ell(\theta; \mathbf{y}) \equiv \log f(\mathbf{y}; \theta)$
 $\hat{\theta} \equiv \arg \sup_{\theta} \ell(\theta; \mathbf{y})$
- $\hat{\theta}$ converges to the “closest true value” KL-divergence

$$\theta_m^0 = \arg \min_{\theta} \int m(\mathbf{y}) \log \left\{ \frac{m(\mathbf{y})}{f(\mathbf{y}; \theta)} \right\} d\mathbf{y}$$

- $\hat{\theta}$ has asymptotic normal distribution, but is not fully efficient “sandwich variance”

$$\text{a.var.} (\hat{\theta}) = G^{-1}(\theta_m^0), \quad G(\theta) = J(\theta)I^{-1}(\theta)J(\theta)$$
$$I = \text{var}_m(\ell'), \quad J = \text{E}_m(-\ell'')$$

- change the inference goal, proceed more or less as usual

- true model $m(\mathbf{y})$ fitted model $f(\mathbf{y}; \theta)$ $\mathbf{y} = (y_1, \dots, y_n)$
- maximum likelihood estimator $\hat{\theta}$ $\ell(\theta; \mathbf{y}) \equiv \log f(\mathbf{y}; \theta)$
 $\hat{\theta} \equiv \arg \sup_{\theta} \ell(\theta; \mathbf{y})$
- $\hat{\theta}$ converges to the “closest true value” KL-divergence

$$\theta_m^0 = \arg \min_{\theta} \int m(\mathbf{y}) \log \left\{ \frac{m(\mathbf{y})}{f(\mathbf{y}; \theta)} \right\} d\mathbf{y}$$

- $\hat{\theta}$ has asymptotic normal distribution, but is not fully efficient “sandwich variance”

$$\text{a.var.} (\hat{\theta}) = G^{-1}(\theta_m^0), \quad G(\theta) = J(\theta)I^{-1}(\theta)J(\theta)$$

$$I = \text{var}_m(\ell'), J = E_m(-\ell'')$$

- change the inference goal, proceed more or less as usual
“we used robust standard errors”

2. More flexible inference functions

Composite likelihood

- **true model** $m(\mathbf{y}_i) = f(\mathbf{y}_i; \theta), \mathbf{y}_i \in \mathbb{R}^d$ fitted model $\prod_{A \in \mathcal{A}} f(y_{iA}; \theta)$ subsets A

2. More flexible inference functions

Composite likelihood

- **true model** $m(\mathbf{y}_i) = f(\mathbf{y}_i; \theta)$, $\mathbf{y}_i \in \mathbb{R}^d$ fitted model $\prod_{A \in \mathcal{A}} f(y_{iA}; \theta)$ subsets A
- Example: pairwise likelihood $\mathbf{y} = (\mathbf{y}_1, \dots, \mathbf{y}_n)$

$$L_{pair}(\theta; \mathbf{y}) = \prod_{i=1}^n \prod_{s \neq t} f_2(y_{is}, y_{it}; \theta)$$

2. More flexible inference functions

Composite likelihood

- **true model** $m(\mathbf{y}_i) = f(\mathbf{y}_i; \theta)$, $\mathbf{y}_i \in \mathbb{R}^d$ fitted model $\prod_{A \in \mathcal{A}} f(y_{iA}; \theta)$ subsets A
- Example: pairwise likelihood $\mathbf{y} = (\mathbf{y}_1, \dots, \mathbf{y}_n)$

$$L_{pair}(\theta; \mathbf{y}) = \prod_{i=1}^n \prod_{s \neq t} f_2(y_{is}, y_{it}; \theta)$$

- Example AR(1) likelihood $\mathbf{y} = (y_1, \dots, y_n)$

$$L_{cond}(\theta; \mathbf{y}) = \prod_{i=1}^n f(y_i \mid y_{i-1}; \theta)$$

interpretation of θ

2. More flexible inference functions

Composite likelihood

- **true model** $m(\mathbf{y}_i) = f(\mathbf{y}_i; \theta)$, $\mathbf{y}_i \in \mathbb{R}^d$ **fitted model** $\prod_{A \in \mathcal{A}} f(y_{iA}; \theta)$ subsets A
- Example: pairwise likelihood $\mathbf{y} = (\mathbf{y}_1, \dots, \mathbf{y}_n)$

$$L_{pair}(\theta; \mathbf{y}) = \prod_{i=1}^n \prod_{s \neq t} f_2(y_{is}, y_{it}; \theta)$$

- Example AR(1) likelihood $\mathbf{y} = (y_1, \dots, y_n)$

$$L_{cond}(\theta; \mathbf{y}) = \prod_{i=1}^n f(y_i | y_{i-1}; \theta)$$

- Example pseudo-likelihood in spatial models interpretation of θ
condition on near neighbours; Besag 74

... More flexible inference functions

Quasi-likelihood and **generalized estimating equations**

$$g\{E(y_i | \mathbf{x}_i)\} = g(\mu_i) = \mathbf{x}_i^T \boldsymbol{\beta}, \quad \text{var}(y_i | \mathbf{x}_i) = \sigma^2 V(\mu_i)$$

- estimating equation for $\boldsymbol{\beta}$

full distribution unspecified

$$\sum_{i=1}^n \frac{\partial \mu_i(\boldsymbol{\beta})}{\partial \boldsymbol{\beta}} \frac{(y_i - \mu_i)}{V(\mu_i)} = \mathbf{o}$$

column vector

... More flexible inference functions

Quasi-likelihood and **generalized estimating equations**

$$g\{E(y_i | \mathbf{x}_i)\} = g(\mu_i) = \mathbf{x}_i^T \boldsymbol{\beta}, \quad \text{var}(y_i | \mathbf{x}_i) = \sigma^2 V(\mu_i)$$

- estimating equation for $\boldsymbol{\beta}$

full distribution unspecified

$$\sum_{i=1}^n \frac{\partial \mu_i(\boldsymbol{\beta})}{\partial \boldsymbol{\beta}} \frac{(y_i - \mu_i)}{V(\mu_i)} = \mathbf{o}$$

column vector

Quadratic inference functions

Qu, Lindsay, Li 2000; Hector 2023

- replace $V^{-1}(\mu_i)$ above with an expansion in basis functions
- apply generalized method of moments

3. More flexible models

- identify one or more parameters of interest here β
- use a highly flexible specification form for other aspects of the model

3. More flexible models

- identify one or more parameters of interest here β
- use a highly flexible specification form for other aspects of the model
- Example: proportional hazards regression instantaneous failure rate

$$h(t; x, \beta) = h_0(t) \exp(x^T \beta)$$

3. More flexible models

- identify one or more parameters of interest here β
- use a highly flexible specification form for other aspects of the model
- Example: proportional hazards regression instantaneous failure rate

$$h(t; \mathbf{x}, \boldsymbol{\beta}) = h_0(t) \exp(\mathbf{x}^T \boldsymbol{\beta})$$

- Example: empirical likelihood $T(F)$ to be specified; e.g. $E_F(Y_i)$

$$\max_F L(F; \mathbf{y}), \text{ subject to } T(F) = \theta$$

$$L(F; \mathbf{y}) = \prod_{i=1}^n F(y_i)$$

3. More flexible models

- identify one or more parameters of interest here β
- use a highly flexible specification form for other aspects of the model
- Example: proportional hazards regression instantaneous failure rate

$$h(t; \mathbf{x}, \boldsymbol{\beta}) = h_0(t) \exp(\mathbf{x}^T \boldsymbol{\beta})$$

- Example: empirical likelihood $T(F)$ to be specified; e.g. $E_F(Y_i)$

$$\max_F L(F; \mathbf{y}), \text{ subject to } T(F) = \theta$$

$$L(F; \mathbf{y}) = \prod_{i=1}^n F(y_i)$$

- Example: semi-parametric regression

$$E(y | T, \mathbf{x}) = \psi T + \omega(\mathbf{x})$$

3. More flexible models

- identify one or more parameters of interest here β
- use a highly flexible specification form for other aspects of the model
- Example: proportional hazards regression instantaneous failure rate

$$h(t; x, \beta) = h_0(t) \exp(x^T \beta)$$

- Example: empirical likelihood $T(F)$ to be specified; e.g. $E_F(Y_i)$

$$\max_F L(F; \mathbf{y}), \text{ subject to } T(F) = \theta$$

$$L(F; \mathbf{y}) = \prod_{i=1}^n F(y_i)$$

- Example: semi-parametric regression

$$E(y | T, x) = \psi T + \omega(x)$$

- does parameter of interest have a stable interpretation model assumption

- Possible model $E(y | T, x) = \psi T + \omega(x)$ binary treatment T
- Define an estimand of interest limit of E -estimator, Robins et al 92

- Possible model $E(y | T, x) = \psi T + \omega(x)$ binary treatment T

- Define an estimand of interest limit of E -estimator, Robins et al 92

$$\frac{E[\pi(x)\{1 - \pi(x)\}\{E(y | T = 1, x) - E(y | T = 0, x)\}]}{E[\pi(x)\{1 - \pi(x)\}]}$$

propensity score $\pi(x) = \text{pr}(T = 1 | x)$

- Possible model $E(y | T, x) = \psi T + \omega(x)$ binary treatment T

- Define an estimand of interest limit of E -estimator, Robins et al 92

$$\frac{E[\pi(x)\{1 - \pi(x)\}\{E(y | T = 1, x) - E(y | T = 0, x)\}]}{E[\pi(x)\{1 - \pi(x)\}]}$$

propensity score $\pi(x) = \text{pr}(T = 1 | x)$

- reduces to ψ under this model
- is a meaningful quantity when the model is incorrect e.g. interaction between T and x

- Possible model $E(y | T, x) = \psi T + \omega(x)$ binary treatment T

- Define an estimand of interest limit of E -estimator, Robins et al 92

$$\frac{E[\pi(x)\{1 - \pi(x)\}\{E(y | T = 1, x) - E(y | T = 0, x)\}]}{E[\pi(x)\{1 - \pi(x)\}]}$$

propensity score $\pi(x) = \text{pr}(T = 1 | x)$

- reduces to ψ under this model
- is a meaningful quantity when the model is incorrect e.g. interaction between T and x
- linear model \rightarrow generalized linear model $g\{E(y | T, x)\} = \psi T + \omega(x)$

$$\frac{E(\pi(x)\{1 - \pi(x)\}[g\{E(y | T = 1, x)\} - g\{E(y | T = 0, x)\}])}{E[\pi(x)\{1 - \pi(x)\}]}$$

- Possible model $E(y | T, x) = \psi T + \omega(x)$ binary treatment T

- Define an estimand of interest limit of E -estimator, Robins et al 92

$$\frac{E[\pi(x)\{1 - \pi(x)\}\{E(y | T = 1, x) - E(y | T = 0, x)\}]}{E[\pi(x)\{1 - \pi(x)\}]}$$

propensity score $\pi(x) = \text{pr}(T = 1 | x)$

- reduces to ψ under this model
- is a meaningful quantity when the model is incorrect e.g. interaction between T and x
- linear model \rightarrow generalized linear model $g\{E(y | T, x)\} = \psi T + \omega(x)$

$$\frac{E(\pi(x)\{1 - \pi(x)\}[g\{E(y | T = 1, x)\} - g\{E(y | T = 0, x)\}])}{E[\pi(x)\{1 - \pi(x)\}]}$$

Example

- survival times for n matched pairs (y_{1i}, y_{2i})
- random assignment of pair members to treatment/control
- nuisance parameters describing the pairs $\lambda_1, \dots, \lambda_n$
- parameter of interest is the treatment effect

- survival times for n matched pairs (y_{1i}, y_{2i})
- random assignment of pair members to treatment/control
- nuisance parameters describing the pairs $\lambda_1, \dots, \lambda_n$
- parameter of interest is the treatment effect
- model y_{1i} exponential with rate λ_i/ψ
 y_{2i} exponential with rate $\lambda_i\psi$
- ψ common parameter of interest; λ_i pair-specific nuisance parameters

- survival times for n matched pairs (y_{1i}, y_{2i})
- random assignment of pair members to treatment/control
- nuisance parameters describing the pairs $\lambda_1, \dots, \lambda_n$
- parameter of interest is the treatment effect
- model y_{1i} exponential with rate λ_i/ψ
 y_{2i} exponential with rate $\lambda_i\psi$
- ψ common parameter of interest; λ_i pair-specific nuisance parameters
- consider modelling nuisance parameters λ_i as random effects

- model y_{1i} exponential with rate λ_i/ψ
 y_{2i} exponential with rate $\lambda_i \psi$
- random effects model: $\lambda_i \sim \text{Gamma}(\alpha, \beta)$ shape, rate
- integrated likelihood

$$L(\psi, \alpha, \beta; \mathbf{y}) = \int f(\mathbf{y}; \psi, \boldsymbol{\lambda}) g(\boldsymbol{\lambda}; \alpha, \beta) d\boldsymbol{\lambda}$$

- model y_{1i} exponential with rate λ_i/ψ
 y_{2i} exponential with rate $\lambda_i \psi$
- random effects model: $\lambda_i \sim \text{Gamma}(\alpha, \beta)$ shape, rate
- integrated likelihood

$$L(\psi, \alpha, \beta; \mathbf{y}) = \int f(\mathbf{y}; \psi, \boldsymbol{\lambda}) g(\boldsymbol{\lambda}; \alpha, \beta) d\boldsymbol{\lambda}$$

- in the integrated model, ψ is orthogonal to (α, β) $\hat{\psi}$ asy. ind't of $(\hat{\alpha}, \hat{\beta})$

- model y_{1i} exponential with rate λ_i/ψ
 y_{2i} exponential with rate $\lambda_i \psi$
- random effects model: $\lambda_i \sim \text{Gamma}(\alpha, \beta)$ shape, rate
- integrated likelihood

$$L(\psi, \alpha, \beta; \mathbf{y}) = \int f(\mathbf{y}; \psi, \boldsymbol{\lambda}) g(\boldsymbol{\lambda}; \alpha, \beta) d\boldsymbol{\lambda}$$

- in the integrated model, ψ is orthogonal to (α, β) $\hat{\psi}$ asy. ind't of $(\hat{\alpha}, \hat{\beta})$
- even better: this is the case for any random effects distribution not just Gamma

- model y_{1i} exponential with rate λ_i/ψ
 y_{2i} exponential with rate $\lambda_i \psi$
- random effects model: $\lambda_i \sim \text{Gamma}(\alpha, \beta)$ shape, rate
- integrated likelihood

$$L(\psi, \alpha, \beta; \mathbf{y}) = \int f(\mathbf{y}; \psi, \boldsymbol{\lambda}) g(\boldsymbol{\lambda}; \alpha, \beta) d\boldsymbol{\lambda}$$

- in the integrated model, ψ is orthogonal to (α, β) $\hat{\psi}$ asy. ind't of $(\hat{\alpha}, \hat{\beta})$
- even better: this is the case for any random effects distribution not just Gamma
- conclude MLE $\hat{\psi} \xrightarrow{P} \psi$ even if random effects model is misspecified could be inefficient

- model y_{1i} exponential with rate λ_i/ψ
 y_{2i} exponential with rate $\lambda_i \psi$
- random effects model: $\lambda_i \sim \text{Gamma}(\alpha, \beta)$ shape, rate
- integrated likelihood

$$L(\psi, \alpha, \beta; \mathbf{y}) = \int f(\mathbf{y}; \psi, \boldsymbol{\lambda}) g(\boldsymbol{\lambda}; \alpha, \beta) d\boldsymbol{\lambda}$$

- in the integrated model, ψ is orthogonal to (α, β) $\hat{\psi}$ asy. ind't of $(\hat{\alpha}, \hat{\beta})$
- even better: this is the case for any random effects distribution not just Gamma
- conclude MLE $\hat{\psi} \xrightarrow{P} \psi$ even if random effects model is misspecified could be inefficient

1. for pair (y_{1i}, y_{2i})

$$L(\psi, \alpha, \beta; y_{1i}, y_{2i}) = \frac{\Gamma(\alpha + 2)}{\Gamma(\alpha)} \frac{\beta^\alpha}{(y_{1i}/\psi + \psi y_{2i} + \beta)^{\alpha+2}}$$

2. orthog

$$\frac{\partial^2}{\partial \psi \partial \alpha} \log L(\psi, \alpha, \beta; y_{1i}, y_{2i}) = \frac{y_{2i} - y_{1i}/\psi^2}{y_{2i}\psi + y_{1i}/\psi + \beta}$$

- 3.

$$E_m \left\{ \frac{y_{2i} - y_{1i}/\psi^2}{y_{2i}\psi + y_{1i}/\psi + \beta} \right\} = 0$$

4. for Gamma random effects, but also for any random effects distribution

interpretation of α, β when not Gamma

- true model $m(\mathbf{y})$ with parameter ψ and true value ψ_*
- fitted model $f(\mathbf{y}; \psi, \lambda)$ same parameter of interest, (many) nuisance parameters
 interpretation of ψ is stable

- true model $m(\mathbf{y})$ with parameter ψ and true value ψ_*
- fitted model $f(\mathbf{y}; \psi, \lambda)$ same parameter of interest, (many) nuisance parameters
 - interpretation of ψ is stable
- maximum likelihood estimates $(\hat{\psi}, \hat{\lambda}) \xrightarrow{P} (\psi_m^o, \lambda_m^o)$
 - KL divergence
- assume no value of $\lambda \in \Lambda$ gives back $m(\cdot)$
 - misspecified

- true model $m(\mathbf{y})$ with parameter ψ and true value ψ_*
- fitted model $f(\mathbf{y}; \psi, \lambda)$ same parameter of interest, (many) nuisance parameters
 - interpretation of ψ is stable
- maximum likelihood estimates $(\hat{\psi}, \hat{\lambda}) \xrightarrow{P} (\psi_m^o, \lambda_m^o)$
 - KL divergence
- assume no value of $\lambda \in \Lambda$ gives back $m(\cdot)$
 - misspecified
- Does $\hat{\psi} \xrightarrow{P} \psi_*$?

- true model $m(\mathbf{y})$ with parameter ψ and true value ψ_*
- fitted model $f(\mathbf{y}; \psi, \lambda)$ same parameter of interest, (many) nuisance parameters
 - interpretation of ψ is stable
- maximum likelihood estimates $(\hat{\psi}, \hat{\lambda}) \xrightarrow{P} (\psi_m^0, \lambda_m^0)$
 - KL divergence
- assume no value of $\lambda \in \Lambda$ gives back $m(\cdot)$
 - misspecified
- Does $\hat{\psi} \xrightarrow{P} \psi_*$? need $E_m\{\partial\ell(\psi_*, \lambda_m^0)/\partial\psi\} = 0$ (1)
 - λ_m^0 unknown

- true model $m(\mathbf{y})$ with parameter ψ and true value ψ_*
- fitted model $f(\mathbf{y}; \psi, \lambda)$ same parameter of interest, (many) nuisance parameters
 - interpretation of ψ is stable
- maximum likelihood estimates $(\hat{\psi}, \hat{\lambda}) \xrightarrow{P} (\psi_m^0, \lambda_m^0)$
 - KL divergence
- assume no value of $\lambda \in \Lambda$ gives back $m(\cdot)$
 - misspecified
- Does $\hat{\psi} \xrightarrow{P} \psi_*$? need $E_m\{\partial\ell(\psi_*, \lambda_m^0)/\partial\psi\} = 0$ (1)
 - λ_m^0 unknown
- can be easier to show $E_m\{\partial\ell(\psi_*, \lambda)/\partial\psi\} = 0 \quad \forall \lambda$ (2)

- true model $m(\mathbf{y})$ with parameter ψ and true value ψ_*
- fitted model $f(\mathbf{y}; \psi, \lambda)$ same parameter of interest, (many) nuisance parameters
interpretation of ψ is stable
- maximum likelihood estimates $(\hat{\psi}, \hat{\lambda}) \xrightarrow{P} (\psi_m^0, \lambda_m^0)$
KL divergence
- assume no value of $\lambda \in \Lambda$ gives back $m(\cdot)$ misspecified
- Does $\hat{\psi} \xrightarrow{P} \psi_*$? need $E_m\{\partial\ell(\psi_*, \lambda_m^0)/\partial\psi\} = 0$ (1) λ_m^0 unknown
- can be easier to show $E_m\{\partial\ell(\psi_*, \lambda)/\partial\psi\} = 0 \quad \forall \lambda$ (2)
- Result: $(1) \equiv (2) \iff \psi_*$ is m -orthogonal to Λ :

$$\forall \lambda \quad E_m \left\{ \frac{\partial^2 \ell(\psi, \lambda)}{\partial \psi \partial \lambda} \right\} = 0$$

- true model $m(\mathbf{y})$ with parameter ψ and true value ψ_*
- fitted model $f(\mathbf{y}; \psi, \lambda)$
- Result : orthogonal parameters lead to consistent estimate m -orthogonal
- Some weaker conditions will suffice for m -orthogonality, hence consistency
- Result : if the parametrization of f has a particular symmetry, then ψ_* is m -orthogonal to nuisance, hence consistency
- Example: generalized scale model $f(y; \sigma) = (1/\sigma)f_0(y/\sigma)$

$$f(y; \psi\lambda) = \frac{1}{\psi} f_0\left(\frac{y}{\psi}; \lambda\right), \quad f(y; \lambda/\psi) = \psi f_0(\psi y; \lambda)$$

Tentative conclusions, further work

- Results above only establish consistency
 - asymptotic variance is much more difficult
- although estimating it might be okay

Tentative conclusions, further work

- Results above only establish consistency
- asymptotic variance is much more difficult although estimating it might be okay
- in the matched pairs examples, nuisance parameters treated as arbitrary constants can be eliminated by transformation to conditional or marginal distributions
- effectively assuming an arbitrary (nonparametric) mixing distribution
- less efficient when the random effects model is correct

Tentative conclusions, further work

- Results above only establish consistency
- asymptotic variance is much more difficult although estimating it might be okay
- in the matched pairs examples, nuisance parameters treated as arbitrary constants can be eliminated by transformation to conditional or marginal distributions
- effectively assuming an arbitrary (nonparametric) mixing distribution
- less efficient when the random effects model is correct
- orthogonality under assumed model $E_\theta\{-\partial^2\ell(\theta)/\partial\theta\partial\theta^T\} = \mathbf{0}$ $\theta = (\psi, \lambda)$
- **m-orthogonality** under true model $E_m\{-\partial^2\ell(\theta)/\partial\theta\partial\theta^T\} = \mathbf{0}$
- connection to Neyman orthogonality? decorrelated score
$$\partial\ell(\psi, \lambda)/\partial\psi - \mathbf{w}^T\partial\ell(\psi, \lambda)/\partial\lambda, \quad \mathbf{w} = I_{\psi\lambda}I_{\lambda\lambda}^{-1}$$
- extension to general estimating equations important in 2-debiased ML

Conclusion

What can go wrong?

- too many parameters

$$p \sim n, \quad p/n \rightarrow C, \quad p/n \rightarrow \infty$$

- weird parameter space

$$pf(y; \theta_1) + (1 - p)f(y; \theta_2), \quad 0 \leq p \leq 1$$

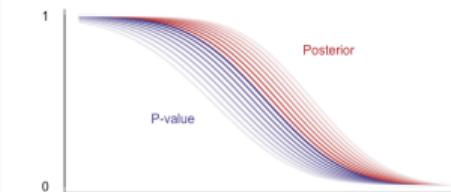
- computational intractability

$$L(\theta, \tau; y) = \int_{\mathbb{R}^k} f(y | z; \theta) f(z; \tau) dz$$

- model is misspecified

true $Y \sim m(y)$, $\nexists \theta_* \text{ w. } f(\cdot; \theta_*) = m(\cdot)$

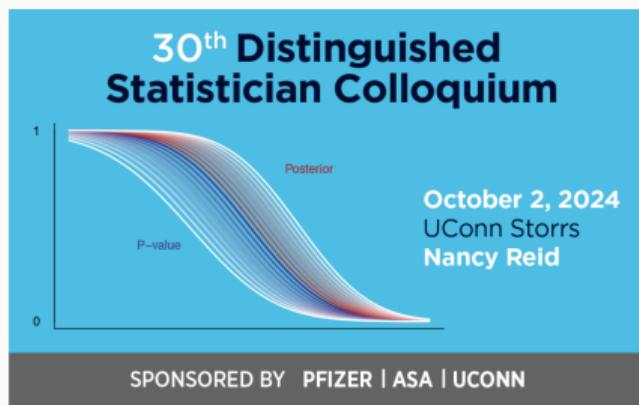
- likelihood is not a probability



- too many parameters
 - new asymptotic theory ($p \sim n$) Sur & Candès 19; Zhao et al 22
 - regularization ($p > n$) Lasso, SCAD, MCP
- weird parameter space
 - different asymptotic theory,
e.g. $\chi_D^2 \rightarrow \sum \lambda_j \chi_{1j}^2$ Battey & McCullagh 24
- computational intractability
 - composite likelihood Genton et al 15
- likelihood is not a probability
 - priors can be very influential

... and so much more!

Thank you!



References i

- Battey, H.S. & Cox, D.R. (2020). High dimensional nuisance parameters: an example from parametric survival analysis. *Information Geometry* **3** 119–148. matched pairs exp
- Battey, H.S. & McCullagh, P. (2024). An anomaly arising in the analysis of processes with more than one source of variability. *Biometrika* **111**, 677–689.
- Battey, H.S. & Reid, N. (2024). On the role of parametrization in models with a misspecified nuisance component. [arxiv PNAS](#), to appear.
- Besag, J. (1974). Spatial interaction and the statistical analysis of lattice systems (with discussion). *J. R. Statist. Soc. B* **36**, 192–236.
- Chernozhukov, V. et al. (2018). Double debiased machine learning. *Econometrics Journal* **21**, C1–C68. doi: [10.1111/ectj.12097](https://doi.org/10.1111/ectj.12097) Neyman orthogonality
- Cox, D.R. (1961). Tests of separate families of hypotheses. *4th Berkeley Symposium* **1**, 105–123.

References ii

- Cox, D.R. (1962). Further results on tests of separate families of hypotheses. *J. R. Statist. Soc. B* **24**, 406–424.
- Cox, D.R. (1990). Role of models in statistical analysis. *Statist. Sci.* **5**, 169–174.
- Cox, D.R. and Donnelly, C.A. (2011). *Principles of Applied Statistics*. Cambridge University Press.
- Davison, A.C. (2003). *Statistical Models*. Cambridge University Press.
- Freedman, D.A. (2005). *Statistical Models*. Cambridge University Press.
- Genton, M.G., Padoan, S.A. and Sang, H. (2015). Multivariate max-stable processes. *Biometrika* **102**, 215–230.
- Hector, E. (2023). Fused mean structure learning in data integration with dependence. *Canad. J. Statist.*

References iii

- Huber, P.J. (1967). The behaviour of maximum likelihood estimates under non-standard conditions. *5th Berkeley Symposium* **1**, 221–233.
- Jorgensen, B. & Knudsen, S.J. (2004). Parameter orthogonality and bias adjustment for estimating functions. *Scand. J. Statist.* **31**, 93–114.
- Lehmann, E.L. (1990). Model specification: the views of Fisher and Neyman, and later developments. *Statist. Sci.* **5**, 160–168.
- McCullagh, P. (2002). What is a statistical model? (with discussion). *Ann. Statist.* **30**, 1225–1310.
- Ning, Y. & Liu, H. (2017). A general theory of hypothesis tests and confidence regions for sparse high-dimensional models. *Annals of Statistics* **45**, 158–195. decorrelated score
- Qu, A.m Lindsay, B.G. and Li, B. (2000). Improving generalised estimating equations using quadratic inference functions. *Biometrika* **87**, 823–836.

References iv

- Sartori, N., Severini, T.A. & Marras, E. (2010). An alternative specification of generalized linear mixed models. *Comp. Stat. Data. Anal.* **54**, 575–584. <https://arxiv.org/abs/2402.05708>
- Vansteelandt, S. & Dukes, O. (2022). Assumption-lean inference for generalized linear models (with discussion). *J. R. Statist. Soc. B* **84**, 657–685. focus on estimand
- White, H. (1982). Maximum likelihood estimation of misspecified models. *Econometrica* **50**, 1–25.
- Zhao, Q., Sur, P. and Candès, E. (2022). The asymptotic distribution of the MLE in high-dimensional logistic models: Arbitrary covariance. *Bernoulli* **28**, 1835–1861.