

Stata: a short history viewed through epidemiology

Bianca L De Stavola

UCL Great Ormond Street Institute of Child Health

b.destavola@ucl.ac.uk

UK Stata Conference, 12-13 September 2024

A personal perspective



► This talk is a personal reflection on 35+ years of applied research in epidemiology

► Aims:

- Pay tribute to influential contributors
- Share some highlights
- Offer reflections

Overview



Some history
Before Stata
The 1990s
The 2000s
The 2010s
The 2020s





Fortran and the first statistical software



```
DIMENSION KA(10,10), KB(10,10), KC(10,10)
EQUIVALENCE (KB(1,1), KC(1,1))
      DEFINE FILE 1(10,10,U,N), 2(10,10,U,L)
      READ (2,100) (KA(1,M),M=1,10)
      WRITE (1'1) (KA(1,H), H=1,10)
      READ (1'1) (KA(1,M),M-1,10)
      00 30 (=1.10
      DO 30 J-1,10
     KC(1,J) = KA(J.1)
     00 40 1=1,10
WRITE (2'1) (KB(1,M),N=1,10)
      DO 50 1=1,10
      READ (1'1) (KA(1,M),M=1,10)
      READ (2'1) (KB(1,M).M=1,10)
      WRITE (1,200) (KA(1,M),M=1,10), (KB(1,M),M=1,10)
100
     FORMAT (10/3)
700 FORMAT (10X, 1015, 10X, 1015)
```



Fortran and the first statistical software



```
DIMENSION KA(10,10), KB(10,10), KC(10,10)
EQUIVALENCE (KB(1,1), KC(1,1))
OLE INE FILE 1(10,10,U,M), Z(10,10,U,L)

DO 10 :=1,10

READ (Z,100) (KA(1,M),M=1,10)

DO 20 :=1,10

DO 30 :=1,10

REC(1,0) = KA(J,1)

DO 30 :=1,10

DO 50 :=1,10

READ (1') (KA(1,M),M=1,10)

READ (1') (KA(1,M),M=1,10)

READ (2') (KA(1,M),M=1,10)

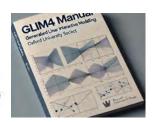
SO WELL (2') (KB(1,M),M=1,10)

FORMAT (1013)

TO FORMAT (1013)

FORMAT (1013)

FORMAT (1013)
```





Fortran and the first statistical software



```
DIMENSION KA(10,10), KB(10,10), KC(10,10)

EQUIVALENCE (KB(1,1), KC(1,1))

OLE INFE FILE 1(10,10,U,N), Z(10,10,U,L)

DO 10 :=1,10

READ (Z,100) (KA(1,N),M=1,10)

DO 20 :=1,10

DO 30 :=1,10

DO 40 :=1,10

MRITE (2*1) (KB(1,M),M=1,10)

DO 50 :=1,10

DO 50 :=1,10

READ (1*1) (KR(1,M),M=1,10)

READ (2*1) (KR(1,M),M=1,10)

READ (2*1) (KR(1,M),M=1,10)

READ (2*1) (KR(1,M),M=1,10)

FORMAT (10)3)

FORMAT (10)3

FORMAT (10)3
```







The 1990s



► Michael Hills and David Clayton

The 1990s Michael Hills and David Clayton













Acknowledgments

The original version of etrate was written by David Clayton (retired) of the Cambridge Institute for Medicial Research and Michael Hills (1934–2021) of the London School of Hygiene and Tropical Medicine.

Acknowledgments

step1st] and atjoin are extensions of lexts by David Clayton (retired) of the Cambridge Institute for Medical Research and Michael Hills (1934-2021) of the London School of Hygiene and Tropical Medicine (Clayton and Hills 1995). The original step1st and styoin commands were written by Jeroen Weesie of the Department of Sociology at Utrecht University, The Netherlands (Weesie 1998a, 1998b), as was the revised attap1st command.

Acknowledgments

We thank David Clayton (retired) of the Cambridge Institute for Medical Research and Michael Hills (1934—2021) of the London School of Hygiene and Tropical Medicine, who wrote the original versions of mhodds and tabodds.





- ► London School of Hygiene and Tropical Medicine
- European Education Program in Epidemiology in Florence



Ana Timberlake





- ► London School of Hygiene and Tropical Medicine
- European Education Program in Epidemiology in Florence





The 2000s



- ► Mixed effects models
- Missing data









gllamm - Generalized linear and latent mixed models

Description Remarks and examples References Also ser

Description

GLIAMM stands for generalized linear latent and mixed models, and gllamm is a Stata command for fitting such models written by Sophia Rabe-Hesketh (University of California-Berkeley) as part of joint work with Anders Skrondal (Norwegian Institute of Public Health) and Andrew Pickles (King's College London).

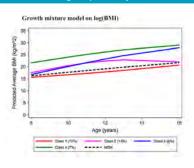






Extensions: latent and grouped trajectories



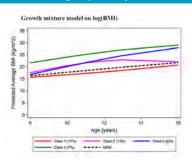


Using mixed and gllamm

[Herle et al. EJE 2021]

Extensions: latent and grouped trajectories

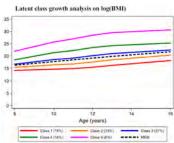




Using mixed and gllamm

[Herle et al. EJE 2021]

Using mixed and traj (Jones and Nagin, 2013)





- Increasing awareness of bias from ignoring missing data bias
- ► Rubin's Multiple Imputation approach and van Buuren's Multiple Imputation by Chained Equations were starting to gain traction



- Increasing awareness of bias from ignoring missing data bias
- ► Rubin's Multiple Imputation approach and van Buuren's Multiple Imputation by Chained Equations were starting to gain traction



The 2010s



► Causal inference



- ► The currently dominant approach in biostatistics and epidemiology relies on potential outcomes (POs) [Rubin, 1974; Robins, 1986; Pearl, 1995]
- ▶ Adopting this approach, we are concerned with questions formulated as contrasts of outcomes that would occur under hypothetical interventions on the exposure:

"Would the outcome of an individual differ if they had/not had that exposure?"

- Robins proposed solutions for estimation of POs*:
 - (a) inverse probability weighting (IPW) (of marginal structural models)
- (b) the g-computation formula
- (c) g-estimation (of structural nested models)
- ▶ teffects implements (a) and (b) for time-fixed exposures

De Stavola/Short history

^{*}Under assumptions of: no interference & consistency (i.e. SUTVA) and conditional atohangeability 🛢 🕨



- ► The currently dominant approach in biostatistics and epidemiology relies on potential outcomes (POs) [Rubin, 1974; Robins, 1986; Pearl, 1995]
- ▶ Adopting this approach, we are concerned with questions formulated as contrasts of outcomes that would occur under hypothetical interventions on the exposure:

"Would the outcome of an individual differ if they had/not had that exposure?"

- ▶ Robins proposed solutions for estimation of POs*:
 - (a) inverse probability weighting (IPW) (of marginal structural models)
 - (b) the g-computation formula
 - (c) g-estimation (of structural nested models)
- teffects implements (a) and (b) for time-fixed exposures

De Stavola/Short history 14/22

^{*}Under assumptions of: no interference & consistency (i.e. SUTVA) and conditional exchangeability 📱 🕨



- ► The currently dominant approach in biostatistics and epidemiology relies on potential outcomes (POs) [Rubin, 1974; Robins, 1986; Pearl, 1995]
- ▶ Adopting this approach, we are concerned with questions formulated as contrasts of outcomes that would occur under hypothetical interventions on the exposure:

"Would the outcome of an individual differ if they had/not had that exposure?"

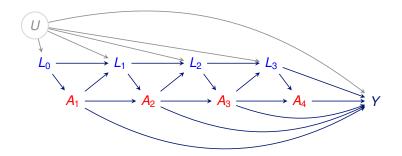
- ▶ Robins proposed solutions for estimation of POs*:
 - (a) inverse probability weighting (IPW) (of marginal structural models)
 - (b) the g-computation formula
 - (c) g-estimation (of structural nested models)
- ▶ teffects implements (a) and (b) for time-fixed exposures

De Stavola/Short history 14/22

^{*}Under assumptions of: no interference & consistency (i.e. SUTVA) and conditional exchangeability 🖹 🕒 🚆 🔗

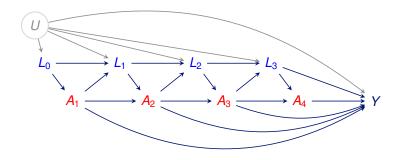


We often deal with scenarios with time-varying confounding of the effect of a time-varying exposure A by a time-varying confounder L:





We often deal with scenarios with time-varying confounding of the effect of a time-varying exposure A by a time-varying confounder L:



Here the total causal effect of A involves L_1 , L_2 , L_3 , although these are also confounders for A_2 , A_3 , A_4 : standard regression modelling does not work!



The Stata Journal (2011) 11, Number 4, pp. 479-517

gformula: Estimating causal effects in the presence of time-varying confounding or mediation using the g-computation formula

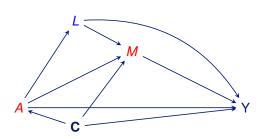
Rhian M. Daniel
Centre for Statistical Methodology
London School of Hygiene and Tropical Medicine
London, UK
rhian.daniel@ishtm.ac.uk

Bianca L. De Stavola Centre for Statistical Methodology London School of Hygiene and Tropical Medicine London, UK

Simon N. Cousens Centre for Statistical Methodology London School of Hygiene and Tropical Medicine London, UK



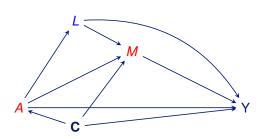




- gformula can be used to estimate natural and interventional effects
- ightharpoonup medeff (Hicks and Tingley, 2011) and paramed (Emsley and Liu, 2013)[†] can only be used when L is not an intermediate confounder

Now incorporated in version 18

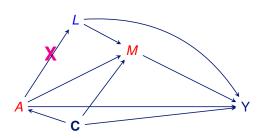




- ▶ gformula can be used to estimate natural and interventional effects
- ▶ medeff (Hicks and Tingley, 2011) and paramed (Emsley and Liu, 2013)[†] can only be used when L is not an intermediate confounder

Now incorporated in version 18





- ▶ gformula can be used to estimate natural and interventional effects
- ightharpoonup medeff (Hicks and Tingley, 2011) and paramed (Emsley and Liu, 2013)[†] can only be used when L is not an intermediate confounder



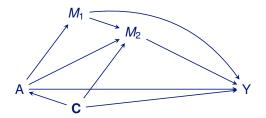
^TNow incorporated in version 18

The 2010s

Mediation: extensions to multiple mediators (in Stata)



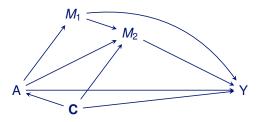
Vansteelandt & Daniel "Interventional effects for mediation analysis with multiple mediators", *Epidemiology* 2017



Mediation: extensions to multiple mediators (in Stata)



Vansteelandt & Daniel "Interventional effects for mediation analysis with multiple mediators", *Epidemiology* 2017



Micali *et al.* "Maternal Prepregnancy Weight Status and Adolescent Eating Disorder Behaviors", *Epidemiology* 2018

A: Prepregnancy maternal BMI

Y: Binge eating score at 13/14y

M₁: Childhood growth 8-12y

M2: Maternal food avoidance at 8y

Effect of Maternal overweight		
	Mean difference	95% CI
Total	0.25	0.18, 0.32
Direct	-0.02	-0.08, 0.05
Indirect via growth	0.28	0.23, 0.33
Indirect via environment	-0.02	-0.04, -0.01

The 2020s



- ► Administrative databases
- ▶ High-dimensional covariates



- ▶ Linked administrative data sources increasingly available for:
 - comparative effectiveness research
 - policy evaluations
- ▶ Recognition of biases potentially affecting such research:
 - Confounding and measurement error
 - Selection bias
 - Lack of positivity
 - Immortal time bias
 - High dimensionality
- Advantages in emulating the design principles of experimental studies to avoid some of these biases ("target trial emulation")



- ► Linked administrative data sources increasingly available for:
 - comparative effectiveness research
 - policy evaluations
- ▶ Recognition of biases potentially affecting such research:
 - Confounding and measurement error
 - Selection bias
 - Lack of positivity
 - Immortal time bias
 - High dimensionality
- ► Advantages in emulating the design principles of experimental studies to avoid some of these biases ("target trial emulation")



- Linked administrative data sources increasingly available for:
 - comparative effectiveness research
 - policy evaluations
- ▶ Recognition of biases potentially affecting such research:
 - Confounding and measurement error
 - Selection bias
 - Lack of positivity
 - Immortal time bias
 - High dimensionality
- ► Advantages in emulating the design principles of experimental studies to avoid some of these biases ("target trial emulation")



- ► Linked administrative data sources increasingly available for:
 - comparative effectiveness research
 - policy evaluations
- ▶ Recognition of biases potentially affecting such research:
 - Confounding and measurement error
 - Selection bias
 - Lack of positivity
 - Immortal time bias
 - High dimensionality
- ► Advantages in emulating the design principles of experimental studies to avoid some of these biases ("target trial emulation")

The 2020s

Linked administrative data: Nguyen et al. 2024



- ▶ Background: Special educational needs (SEN) provision is designed to help pupils with additional educational, behavioural or health needs
- ► Aim: assess the impact of SEN provision on an educational outcomes during primary education for children with a certain congenital abnormality
- ▶ Data: ECHILD, linked educational and health records across England
- ▶ Results with/without (correct) lasso selection (using telasso)[‡]:

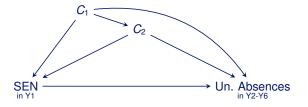
[‡]As developed by Chernozhukov (2018); Code to be deposited in GitHub < □ ▶ < 점 ▶ < 필 ▶ ◀ 필 ▶ ▼ 및 ◆ 의 및 ◆

The 2020s

Linked administrative data: Nguyen et al. 2024



- ▶ Background: Special educational needs (SEN) provision is designed to help pupils with additional educational, behavioural or health needs
- ► Aim: assess the impact of SEN provision on an educational outcomes during primary education for children with a certain congenital abnormality
- Data: ECHILD, linked educational and health records across England



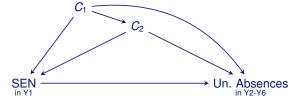
Results with/without (correct) lasso selection (using telasso)[‡]:

De Stavola/Short history

TAs developed by Chernozhukov (2018); Code to be deposited in GitHub ∢ □ ▶ ∢ ₺ ▶ ∢ ₺ ▶ ▼ ₺ ♥ ९ €



- ▶ Background: Special educational needs (SEN) provision is designed to help pupils with additional educational, behavioural or health needs
- ► Aim: assess the impact of SEN provision on an educational outcomes during primary education for children with a certain congenital abnormality
- Data: ECHILD, linked educational and health records across England

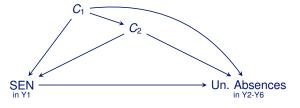


Results with/without (correct) lasso selection (using telasso)[‡]:

^{*}As developed by Chernozhukov (2018); Code to be deposited in GitHub « 🗆 ト « 🗇 ト « 🗏 ト 🧸 🖹 🔻 🔗 🥄



- ▶ Background: Special educational needs (SEN) provision is designed to help pupils with additional educational, behavioural or health needs
- ► Aim: assess the impact of SEN provision on an educational outcomes during primary education for children with a certain congenital abnormality
- Data: ECHILD, linked educational and health records across England



Results with/without (correct) lasso selection (using telasso)[‡]:

Effect of SEN in Y1		
	Rate Ratio	95% CI
Crude	1.22	1.11, 1.34
IPW	0.86	0.76, 0.97
G-computation	0.98	0.86, 1.09
AIPW-lasso with int.	0.80	0.66, 0.95

Final thoughts ...



Final thoughts ...



Positives

- Wonderful Stata community
- ► Cross-pollination with econometricians
- ► Results increasingly reproducible

Final thoughts ...



Positives

- Wonderful Stata community
- ► Cross-pollination with econometricians
- ► Results increasingly reproducible

Future challenges

► Access to Stata within secure environments: only via Google Notebooks and/or Python

Final thoughts



Positives

- Wonderful Stata community
- ► Cross-pollination with econometricians
- ► Results increasingly reproducible

Future challenges

► Access to Stata within secure environments: only via Google Notebooks and/or Python

Thank you for listening!