Methods of Applied Statistics I

STA2101H F LEC9101

Week 2

September 22 2021

29th -> Wilson Hall (New College)

← Tweet



•••

My colleagues have done this fine empirical work on risk matrices - I really like the idea of moving away from a symmetric grid to give a feel for non-linearity. About time that this familiar tool was investigated!

Gabriel Recchia @mesotronium - Sep 15

Risk matrices are very familiar to those involved in health & safety and risk assessments – but do people actually understand them? We took a first look, published today: doi.org/10.1111/risa.1...

Show this thread



7:54 AM · Sep 20, 2021 · TweetDeck





- Office hours, Upcoming events, HW 2
 OH: Monday 7pm-8.30pm, Wednesday 4pm-5.30pm, Zoom until further notice
- 2. Steps in analysis; types of studies \vee
- 3. Linear Regression Part 2: recap, testing groups of variables, checking model assumptions, collinearity, p > n
- \rightarrow 4. In the News

1. Office hours, Upcoming events, HW 2

OH: Monday 7pm-8.30pm, Wednesday 4pm-5.30pm, Zoom until further notice

- 2. Steps in analysis; types of studies
- 3. Linear Regression Part 2: recap, testing groups of variables, checking model assumptions, collinearity, p > n
- 4. In the News

- SM Statistical Models by Davison
- LM-1,2 Linear Models with R by Faraway (1st and 2nd editions)
 LM (both)
- ELM-1,2 Extending the Linear Model with R by Faraway (1st and 2nd editions)
- CD Principles of Applied Statistics by Cox & Donnelly

Upcoming

- Learning Individualized Treatment Rule for a Target Population
- Thursday 3.30 Link

About Guanhua Chen



I am an Assistant Professor of Biostatistics and Medical Informatics at the University of Wisconsin-Madison. I got my Ph.D. from the University of North Carolina at Chapel Hill in 2014 under the direction of Dr. Michael Kosorok. Before joining UW, I was an Assistant Professor of Biostatistics at Vanderbilt University.

Research Interest

Develop statistical learning methods for clinical and biomedical research. In particular, I am interested in analyzing heterogeneous, high-dimensional-omics data (genome, microbiome) and electronic health record data to advance precision medicine. My current research is supported by PCORI and NSF grants.

Fridays at noon Toronto Data Workshop link

This term is a special series of talks featuring University of Toronto speakers on the relationship between data science and their other field of expertise.

Date	Speaker	Recording
Fri 24 September 2021, noon - 1pm	Karen Chapple, Geography, planning, cities	
Fri 1 October 2021, noon - 1pm	Special on 2021 Canadian Election	
Fri 8 October 2021, noon - 1pm	Fedor Dokshin, Sociology	

Applied Statistics I September 22 2021

Applied Statistics I

HW Question Week 2

STA2101F 2021

Due September 29 2021 11.59 pm

Homework to be submitted through Quercus

This question concerns the article "How people understand risk matrices..." by Sutherland et al., available on the course web site.

- (a) What is a risk matrix? Find an example of a risk matrix for the assessment of risk related to either COVID-19 or wildfire. Provide the reference (link) to the example and a snapshot of one or two of the published risk matrices.
- (b) The authors describe two randomized controlled experiments, in Sections 2 and 3, respectively. What are the units of analysis in these two experiments?
- (c) In the first experiment, the authors describe derived variables: a "basic knowledge score", a "risk comparison score", a "prioritization score", the results of a "matrix
- September 2pperference test", and a "total numeracy score". Which of these are treated as response variables and which are treated as explanatory variables?
 - (d) What 'treatments' were rendemly assigned to units in Experiment 12

Steps in Analysis

- understand the physical background
- understand the objective
- make sure you know what the client wants
- put the problem into statistical terms

Steps in Analysis

- understand the physical background
- understand the objective
- make sure you know what the client wants
- put the problem into statistical terms
- How were the data collected:
 - are the data observational or experimental? etc.
 - is there nonresponse
 - are there missing values
 - how are the data coded
 - what are the units of measurement
 - beware of data entry errors

- slow

LM < plats, summery stat.

date 'cleaning' 'corpentry'

CD §1.2

- start with a scientific question
- assess how data could shed light on this
- plan data collection
- consider of sources of variation and how careful planning can minimize their impact

- start with a scientific question
- assess how data could shed light on this
- plan data collection
- consider of sources of variation and how careful planning can minimize their impact
- develop strategies for data analysis: modelling, computation, methods of analysis
- assess the properties of the methods and their impact on the question at hand



- start with a scientific question
- assess how data could shed light on this
- plan data collection
- consider of sources of variation and how careful planning can minimize their impact
- develop strategies for data analysis: modelling, computation, methods of analysis
- assess the properties of the methods and their impact on the question at hand
- communicate the results: accurately

but not pessimistically

CD §1.2

- start with a scientific question
- assess how data could shed light on this
- plan data collection
- consider of sources of variation and how careful planning can minimize their impact
- develop strategies for data analysis: modelling, computation, methods of analysis
- assess the properties of the methods and their impact on the question at hand
- communicate the results: accurately
- visualization strategies, conveyance of uncertainties

but not pessimistically



choice of material/individuals to study – "units of analysis"

CD Ch. 1

Design and Analysis

- choice of material/individuals to study "units of analysis" cifies failes
- "For studies of a new phenomenon it will usually be best to examine situations in which the phenomenon is likely to appear in the most striking form, even if this is in some sense artificial"

- choice of material/individuals to study "units of analysis"
- "For studies of a new phenomenon it will usually be best to examine situations in which the phenomenon is likely to appear in the most striking form, even if this is in some sense artificial"
- statistical analysis needs to take account of the design (even if statistician enters the project at the analysis stage)

- choice of material/individuals to study "units of analysis"
- "For studies of a new phenomenon it will usually be best to examine situations in which the phenomenon is likely to appear in the most striking form, even if this is in some sense artificial"
- statistical analysis needs to take account of the design (even if statistician enters the project at the analysis stage)
- need to be clear at the design stage about broad features of the statistical analysis

 more publicly convincing and "reduces the possibility that the data cannot be satisfactorily analysed"



- choice of material/individuals to study "units of analysis"
- "For studies of a new phenomenon it will usually be best to examine situations in which the phenomenon is likely to appear in the most striking form, even if this is in some sense artificial"
- statistical analysis needs to take account of the design (even if statistician enters the project at the analysis stage)
- need to be clear at the design stage about broad features of the statistical analysis

 more publicly convincing and "reduces the possibility that the data cannot be satisfactorily analysed"
- "it is unrealistic and indeed potentially dangerous to follow an initial plan unswervingly ... it may be a crucial part of the analysis to clarify the research objectives"

Experimental and observational studies

experiment is a study in which all key elements are under the control of the investigator
 Freatments "

CD Ch. 1

Experimental and observational studies

experiment is a study in which all key elements are under the control of the investigator
assip diet (Mediterrane) Je
in an observational study key elements cannot be manipulated by the investigator.
Smoking Unp cancer (mmes)

CD Ch. 1

- experiment is a study in which all key elements are under the control of the investigator
- in an observational study key elements cannot be manipulated by the investigator.
- "It often, however, aids the interpretation of an observation study to consider the question: what would have been done in a comparable experiment?"

- experiment is a study in which all key elements are under the control of the investigator
- in an observational study key elements cannot be manipulated by the investigator.
- "It often, however, aids the interpretation of an observation study to consider the question: what would have been done in a comparable experiment?"
- Example: hormone replacement therapy and heart disease

- experiment is a study in which all key elements are under the control of the investigator
- in an observational study key elements cannot be manipulated by the investigator.
- "It often, however, aids the interpretation of an observation study to consider the question: what would have been done in a comparable experiment?"
- Example: hormone replacement therapy and heart disease
- observational study strong and statistically significant reduction in heart disease among women taking hormone replacement therapy

- experiment is a study in which all key elements are under the control of the investigator
- in an observational study key elements cannot be manipulated by the investigator.

CD Ch. 1

Subgroup Covander HET &

- "It often, however, aids the interpretation of an observation study to consider the question: what would have been done in a comparable experiment?"
- Example: hormone replacement therapy and heart disease
- observational study strong and statistically significant reduction in heart disease among women taking hormone replacement therapy
- women's health study (JAMA, 2002, p.321) statistically significant increase in risk among women randomized to hormone replacement therapy

ORIGINAL ARTICLE

Observational Study of Hydroxychloroquine in Hospitalized Patients with Covid-19

Joshua Geleris, M.D., Yifei Sun, Ph.D., Jonathan Platt, Ph.D., Jason Zucker, M.D., Matthew Baldwin, M.D., George Hripcsak, M.D., Angelena Labella, M.D., Daniel K. Manson, M.D., Christine Kubin, Pharm.D., R. Graham Barr, M.D., Dr.P.H., Magdalena E. Sobieszczyk, M.D., M.P.H., and Neil W. Schluger, M.D.

Article Figures/Media	Metrics	June 18, 2020 N Engl Med 2020: 382:2411-2418	
14 References 300 Citing Articles		DOI: 10.1056/NEJMoa2012410 Chinese Translation 中文翻译	
Abstract			June 2020

ORIGINAL ARTICLE

A Randomized Trial of Hydroxychloroquine as Postexposure Prophylaxis for Covid-19

David R. Boulware, M.D., M.P.H., Matthew F. Pullen, M.D., Ananta S. Bangdiwala, M.S., Katelyn A. Pastick, B.Sc., Sarah M. Lofgren, M.D., Elizabeth C. Okafor, B.Sc., Caleb P. Skipper, M.D., Alanna A. Nascene, B.A., Melanie R. Nicol, Pharm.D., Ph.D., Mahsa Abassi, D.O., M.P.H., Nicole W. Engen, M.S., Matthew P. Cheng, M.D., et al.

Article	Figures/Media	Metrics	August 6, 2020	
18 Referenc	es 128 Citing Articles Letters 11 Comments		DOI: 10.1056/NEJMoa2016638	August 2020

Applied Statistics I September 22 2021

Hydroxychloroquine



Articles

RETRACTED: Hydroxychloroquine or chloroquine with or without a macrolide for treatment of COVID-19: a multinational registry analysis

Prof Mandeep R Mehra MD ^a lpha 🖾, Sapan S Desai MD ^b, Prof Frank Ruschitzka MD ^c, Amit N Patel MD ^d, ^e

retracted

... hydroxychloroquine

Coch	Trusted evidence. Informed decisions. Better health.	Search	Q
Our evidence	About us Join Cochrane News and jobs	Cochrane Libra	ry 🕨
	Coronavirus (COVID-19) resource	'S	
schrane	Is chloroquine or hydroxychloroquine us COVID-19, or in preventing infection in p the virus?	seful in treating people witl people who have been expo	sed to

Kerievs !!

Cochrane Reviews

weta-analysis "Study of studies"

Applied Statistics I

September 22 2021

... hydroxychloroquine

Cochrane Reviews 🔻

Cochrane Database of Systematic Reviews Review - Intervention

Trials **•**

Chloroquine or hydroxychloroquine for prevention and treatment of COVID-19

Clinical Answers 🔻

About 🔻

Help 🔻

Shagteshwar Singh, Hannah Ryan, Tamara Kredo, Marty Chaplin, Tom Fletcher Authors' declarations of interest Version published: 12 February 2021 Version history https://doi.org/10.1002/14651858.CD013587.pub2

Collapse all Expand all

Hydroxychloroquine does not reduce deaths from COVID-19, and probably does not reduce the number of people needing mechanical ventilation.

Hydroxychloroquine caused more unwanted effects than a placebo treatment, though it did not appear to increase the number of serious unwanted effects.

We do not think new studies of hydroxychloroquine should be started for treatment of COVID

... hydroxychloroquine



Cochrane Database of Systematic Reviews Editorial

Contested effects and chaotic policies: the 2020 story of (hydroxy) chloroquine for treating COVID-19

Susan Gould, Susan L Norris Authors' declarations of interest Version published: 25 March 2021 https://doi.org/10.1002/14651858.ED000151 C

Ivermectin



... ivermectin



"Moderate-certainty evidence finds that large reductions in COVID-19 deaths are possible using ivermectin. ... The apparent safety and low cost suggest that ivermectin is likely to have a Applied Statistics I Septembel ghilfeant impact on the SARS-CoV-2 pandemic globally."

American J Therapeutics, July 2021

Linear regression recap

- generic form of linear regression, in matrix notation $y = X\beta + \epsilon$
- least squares estimate of β is $\hat{\beta} = (X^{T}X)^{-1}X^{T}y$ $\hat{\beta}$ has expected value β and variance-covariance matrix $\sigma^{2}(X^{T}X)$

ßr

Juxi

E(+1x)=xB

Linear regression recap

- generic form of linear regression, in matrix notation $y = X\beta + \epsilon$
- least squares estimate of β is $\hat{\beta} = (X^T X)^{-1} X^T y$
- $\hat{\beta}$ has expected value β and variance-covariance matrix $\sigma^2(X^T X)^{-1}$
- this is the maximum likelihood estimate if $\epsilon \sim \bigvee_{\beta} (0, \sigma^2 I)$ $\hat{\beta} \sim N(\beta, \sigma^2 (X^T X)^{-1})$

•
$$\tilde{\sigma}^2 = (y - X\hat{\beta})^{\mathrm{T}}(y - X\hat{\beta})/(n - p)$$

• leads to *t*-tests for individual components β_i

called s^2 in SM

and confidence intervals

Linear regression recap

- generic form of linear regression, in matrix notation $y = X\beta + \epsilon$
- least squares estimate of β is $\hat{\beta} = (X^T X)^{-1} X^T Y$
- $\hat{\beta}$ has expected value β and variance-covariance matrix $\sigma^2(X^T X)^{-1}$
- this is the maximum likelihood estimate if $\epsilon \sim N(0, \sigma^2 I)$
- $\hat{\beta} \sim N(\beta, \sigma^2(X^{\mathrm{T}}X)^{-1})$

•
$$\tilde{\sigma}^2 = (y - X\hat{\beta})^{\mathrm{T}}(y - X\hat{\beta})/(n - p)$$

- leads to *t*-tests for individual components β_i
- X is an n × p matrix of explanatory variables, which may be the highly "number"
 measured in the sample (SM Fx 8.2)

 - fixed by design (SM Ex 8.4),
 - introduced to make the model more flexible (SM Ex 8.2) 1/x f x , x , x]^{n R, model.matrix}
 - X often called the design matrix (IJALM

September 22 2021 Applied Statistics I

SM - Davison

called s^2 in SM

and confidence intervals

15

Aside: Lazy Notation



• "where we hope there is no confusion"



Aside: ancillarity

- in a statistical model (i.e. likelihood function) that factorizes, in a way that separates the parameters, there are strong theoretical (and often practical) reasons for using only the relevant factor for inference
- in our example, even if (y, X) have a (p + 1)-dimensional distribution (maybe even jointly multivariate normal), it would typically be the case that β , which by definition is $\mathbb{E}(y \mid X)$, is not a parameter of the distribution of X. $E(\gamma \mid X ; \beta) \quad \text{vor}(\gamma \mid X; \sigma^2)$

 $f_{x}(x; x) \in no$ Boro²

Aside: ancillarity

- in a statistical model (i.e. likelihood function) that factorizes, in a way that separates the parameters, there are strong theoretical (and often practical) reasons for using only the relevant factor for inference
- in our example, even if (y, X) have a (p + 1)-dimensional distribution (maybe even jointly multivariate normal), it would typically be the case that β , which by definition is $\mathbb{E}(y \mid X)$, is not a parameter of the distribution of X.
- So we would have something like

$$f(\mathbf{y}, \mathbf{X}; \boldsymbol{\beta}, \sigma^2, \theta) \neq f_1(\mathbf{y} \mid \mathbf{X}; \boldsymbol{\beta}, \sigma^2) f_2(\mathbf{X} \mid \theta)$$

and we base our inference for β and σ² (which are the parameters of interest, by assumption) on f₁(y | X; β, σ²)

SM Ex.12.4; Cox 2006 §4.3.1

n: -x)2

Aside: ancillarity

- in a statistical model (i.e. likelihood function) that factorizes, in a way that separates the parameters, there are strong theoretical (and often practical) reasons for using only the relevant factor for inference
- in our example, even if (y, X) have a (p + 1)-dimensional distribution (maybe even jointly multivariate normal), it would typically be the case that β , which by definition is $\mathbb{E}(y \mid X)$, is not a parameter of the distribution of X. 50 + (2, x:
- So we would have something like

 $f(\mathbf{y}, \mathbf{X}; \boldsymbol{\beta}, \sigma^2, \theta) = f_1(\mathbf{y} \mid \mathbf{X}; \boldsymbol{\beta}, \sigma^2) f_2(\mathbf{X} \mid \theta) \quad \text{var} \mathbf{\beta}_{\mathbf{x}} =$ and we base our inference for β and σ^2 (which are the parameters of interest, by assumption) $\hat{\beta} = (X^T x)^{-1} X^T y$ on $f_1(\mathbf{y} \mid \mathbf{X}; \boldsymbol{\beta}, \sigma^2)$

- in technical terms X is ancillary for (β, σ^2)
- If we did decide to include the variability in X as part of our analysis, our inference about β would be much less precise, and needlessly so, because we are worrying about explanatory variable values that were not seen in our data set.
Comparing Models

SM 8.5, LM 3.1,2

• residual sum of squares
$$SS(\hat{\beta}) = (y - \chi\hat{\beta})^{T}(y - \chi\hat{\beta}) = (y - \hat{j})^{T}(y - \hat{j})$$

• Decomposition of variance $g^{T}y = (y - \hat{j} + \hat{j})^{T}(y - \hat{j} + \hat{j})$
• Decomposition of variance $y^{T}y = (y - \hat{j} + \hat{j})^{T}(y - \hat{j} + \hat{j})$
• Typically first column of X is $(1, ..., 1)^{T}$, so $y = \beta_{0} + X_{2}\beta_{2} + \epsilon$, say; then decomposition becomes
Total SS $I(y_{1} - \hat{y})^{T} = Z(y_{1} - x_{2}\hat{i}\hat{\beta}_{(2})^{T} + \hat{\beta}_{(2)}X_{2}\hat{x}_{1}\hat{\beta}_{(2)}) = X(xx)xy$
(connected $(y - y - \hat{j})^{T}(y - y - \hat{j})$
($y - X_{2}\hat{\beta}_{(2)})^{T}(y - \hat{y}_{2})$
Applied Statistics September 22 2021 $(y - X_{2}\hat{\beta}_{(2)})^{T}(y - \hat{y}_{2}\hat{\beta}_{(2)})$





- same argument can be derived for comparing submodels
- for example, testing $(\beta_2, \beta_3, \beta_4) = (0, 0, 0)$

- same argument can be derived for comparing submodels
- for example, testing $(\beta_2, \beta_3, \beta_4) = (0, 0, 0)$
- fit full model $\longrightarrow RSS_{full}$; fit reduced model $\longrightarrow RSS_{red}$

•

- same argument can be derived for comparing submodels
- for example, testing $(\beta_2, \beta_3, \beta_4) = (0, 0, 0)$
- fit full model $\longrightarrow RSS_{full}$; fit reduced model $\longrightarrow RSS_{red}$

$$F = \frac{(RSS_{red} - RSS_{full})/(p-q)}{RSS_{full}/(n-p)}$$

- same argument can be derived for comparing submodels
- for example, testing $(\beta_2, \beta_3, \beta_4) = (0, 0, 0)$
- fit full model $\longrightarrow RSS_{full}$; fit reduced model $\longrightarrow RSS_{red}$

$$F = \frac{(RSS_{red} - RSS_{full})/(p-q)}{RSS_{full}/(n-p)}$$
• see LM 3.1, SM §8.2 (p.367) for connection to likelihood ratio

test

- same argument can be derived for comparing submodels
- for example, testing $(\beta_2, \beta_3, \beta_4) = (0, 0, 0)$



97.9

• when would we want to do this?

		n = 97	9 cov. + (3.	(mil
head(prostat	e)				(response
# lcavol	lweight	age lbph	svi lcp	gleason pgg45	lpsa (ach)
1 -0.5798185	2.7695	50 -1.386294	0 -1.38629	6- 0	-0.43078 bp (pso)
2 -0.9942523	3.3196	58 -1.386294	0 -1.38629	6 - 0	-0.16252
3 -0.5108256	2.6912	74 -1.386294	0 -1.38629	7 – 20	-0.16252
4 -1.2039728	3.2828	58 -1.386294	0 -1.38629	60	-0.16252
5 0.7514161	3.4324	62 -1.386294	0 -1.38629	6~ 0	0.37156
6 -1.0498221	3.2288	50 -1.386294	0 -1.38629	6 - 0	0.76547
			4	$\land \land \land$	t
model1 <- lm	(lpsa ~ .	., data = pros	state)		0

> summary(model1)

Coefficients:



```
model3 <- lm(lpsa ~ lcavol + lweight + age + lbph + svi, data = prostate)
anova(model3,model1)
Analysis of Variance Table
Model 1: lpsa ~ lcavol + lweight + age + lbph + svi
Model 3: lpsa ~ lcavol + lweight + age + lbph + svi + lcp + gleason +
                                                      ( RSSfel - RSSml)/2
   pgg45
 Res.Df
           RSS Df Sum of Sq F Pr(>F)
     91 45.526
     88 44.163 (3)
                    1.3625 0.905 0.4421
                                                              does this make sense?
                                                  Ha= Ka= Ka= Ka
```

Factor variables

Applied S¹

- F-tests are used when the columns to be removed form a group
- if a covariate is a factor, i.e. categorical, then lm will construct a set of dummy variables as part of the model matrix
- these variables should either all be in, or all be out

in most cases

Levels 6,7,89

- F-tests are used when the columns to be removed form a group
- if a covariate is a factor, i.e. categorical, then lm will construct a set of dummy variables as part of the model matrix
- these variables should either all be in, or all be out

in most cases



```
model4 <- lm(lpsa ~ .-gleason, data=prostate)
summary(model4)</pre>
```

> Coefficients:

	Estimate	Std. Error	t value	Pr(> t)		
(Intercept)	0.913282	0.840838	1.086	0.28044		
lcavol	0.569988	0.090100	6.326	1.09e-08	***	
lweight	0.468791	0.169610	2.764	0.00699	**	
age	-0.021749	0.011361	-1.914	0.05890	•	
lbph	0.099685	0.058984	1.690	0.09464	•	
svi	0.745879	0.247398	3.015	0.00338	**	
lcp	-0.125112	0.095591	-1.309	0.19408		
pgg45	0.004990	0.004672	1.068	0.28848	•	
gleason_factor7	0.267607	0.219419	1.220	0.22595	つし	
gleason_factor8	0.496820	0.769267	0.646	0.52011	3	pt. est. all related
gleason_factor9	-0.056215	0.500196	-0.112	0.91078	1-1	T o
and the second second second						m laca Score.

Applied Statistics I September 22 2021

> anova(model1,model4)
Analysis of Variance Table

```
Model 1: lpsa ~ (lcavol + lweight + age + lbph + svi + lcp + gleason +
   pgg45 + gleason_factor) - gleason
Model 2: lpsa ~ lcavol + lweight + age + lbph + svi + lcp + glasson +
   pgg45
 Res.Df RSS Df Sum of Sq F Pr(>F)
1 86 42.724
2
 88 44.163 -2 -1.4392 1.4485 0.2406
     89 44.204 -3 -1.4805 .9934 .3999
```

- with designed experiments, covariates are often factors set at pre-determined levels
- see, e.g. Example 8.4 in SM

also Ch 14 in LM-2; Ch 13 in LM-1

- with designed experiments, covariates are often factors set at pre-determined levels
- see, e.g. Example 8.4 in SM

also Ch 14 in LM-2; Ch 13 in LM-1

- if the design is perfectly balanced, then X has orthogonal columns, and X^TX is diagonal
- so $\hat{\beta}_j$'s are uncorrelated, and hence independent (under normality assumption)

Applied Statistics I

- with designed experiments, covariates are often factors set at pre-determined levels
- see, e.g. Example 8.4 in SM

also Ch 14 in LM-2; Ch 13 in LM-1

- if the design is perfectly balanced, then X has orthogonal columns, and X^TX is diagonal
- so $\hat{\beta}_j$'s are uncorrelated, and hence independent (under normality assumption)
- more generally we might have $X^{T}X$ block diagonal, e.g.

W

September 22 202

27

importance?

• assumptions on errors: $\epsilon_i \sim_{i.i.d.} N(0, \sigma^2)$

• assumptions on errors: $\epsilon_i \sim_{i.i.d.} N(0, \sigma^2)$ on structure $\mathbb{E}(y \mid X) = X\beta$

- assumptions on errors: $\epsilon_i \sim_{i.i.d.} N(o, \sigma^2)$
- normality; constant variance; independent

on structure $\mathbb{E}(y \mid X) = X\beta$

- assumptions on errors: $\epsilon_i \sim_{i.i.d.} N(o, \sigma^2)$
- normality; constant variance; independent

plot(model1)

on structure $\mathbb{E}(y \mid X) = X\beta$

Model checking

• assumptions on errors: $\epsilon_i \sim_{i.i.d.} N(0, \sigma^2)$

on structure $\mathbb{E}(y \mid X) = X\beta$

normality; constant variance; independent





• residuals: $\hat{\epsilon}_i =$

- residuals: $\hat{\epsilon}_i =$
- $Var(\hat{\epsilon}) =$

- residuals: $\hat{\epsilon}_i =$
- $Var(\hat{\epsilon}) =$
- i.e. don't all have the same variance

- residuals: $\hat{\epsilon}_i =$
- $Var(\hat{\epsilon}) =$
- i.e. don't all have the same variance
- hat matrix *H* =

- residuals: $\hat{\epsilon}_i =$
- $Var(\hat{\epsilon}) =$
- i.e. don't all have the same variance
- hat matrix *H* =
- standardized residuals: $r_i =$

- residuals: $\hat{\epsilon}_i =$
- $Var(\hat{\epsilon}) =$
- i.e. don't all have the same variance
- hat matrix *H* =
- standardized residuals: $r_i =$
- Cook's distance $C_i =$

- residuals: $\hat{\epsilon}_i =$
- $Var(\hat{\epsilon}) =$
- i.e. don't all have the same variance
- hat matrix H =
- standardized residuals: $r_i =$
- Cook's distance $C_i =$

https://data.library.virginia.edu/diagnostic-plots/

- simple model $y_i = \beta_0 + \beta_1 x_{1i} + \beta_2 x_{2i} + \epsilon_i$, $i = 1, \dots n$
- if $x_1 \perp x_2$, then interpretation of β_1 and β_2 clear
- if $x_1 = x_2$ then β_1 and β_2 not separately identifiable

- simple model $y_i = \beta_0 + \beta_1 x_{1i} + \beta_2 x_{2i} + \epsilon_i$, $i = 1, \dots n$
- if $x_1 \perp x_2$, then interpretation of β_1 and β_2 clear
- if $x_1 = x_2$ then β_1 and β_2 not separately identifiable
- usually we're somewhere in between, at least in observational studies
- may be very difficult to dis-entangle effects of correlated covariates

- simple model $y_i = \beta_0 + \beta_1 x_{1i} + \beta_2 x_{2i} + \epsilon_i$, $i = 1, \dots n$
- if $x_1 \perp x_2$, then interpretation of β_1 and β_2 clear
- if $x_1 = x_2$ then β_1 and β_2 not separately identifiable
- usually we're somewhere in between, at least in observational studies
- may be very difficult to dis-entangle effects of correlated covariates
- example: health effects of air pollution
- measurable increase in mortality on high-pollution days
- measurable increase in mortality on high-temperature days
- high temperatures and high levels of pollutants tend to co-occur

- simple model $y_i = \beta_0 + \beta_1 x_{1i} + \beta_2 x_{2i} + \epsilon_i$, $i = 1, \dots n$
- if $x_1 \perp x_2$, then interpretation of β_1 and β_2 clear
- if $x_1 = x_2$ then β_1 and β_2 not separately identifiable
- usually we're somewhere in between, at least in observational studies
- may be very difficult to dis-entangle effects of correlated covariates
- example: health effects of air pollution
- measurable increase in mortality on high-pollution days
- measurable increase in mortality on high-temperature days
- high temperatures and high levels of pollutants tend to co-occur +++
- mathematically, X^TX is nearly singular, or at least ill-conditioned, so calculation of its inverse is subject to numerical errors
- if p > n then $X^{T}X$ not invertible, no LS solution

On the 'Breast cancer epidemic' poster

A poster and presentation at the RSS conference featured flawed analysis.



🧸 Anthony B. Masters 6 days ago · 4 min read \star



Applied Statistics I

First: Correlation is not causation



Figure 5. England & Wales. Cumulated Cohort Incidence for Birth cohorts of Women. Breast Cancer within ages 50-54 and cumulated cohort abortion rates. Correlation coefficient 0.86.


Second, there is no discussion of other well-known risk factors for breast cancer. Instead of calling upon extensive research about breast cancer, the authors cite themselves. Large prospective cohort studies do not imply heightened risks of breast cancer after abortion. Second, there is no discussion of other well-known risk factors for breast cancer. Instead of calling upon extensive research about breast cancer, the authors cite themselves. Large prospective cohort studies do not imply heightened risks of breast cancer after abortion.

Third, there is no exploration of absolute and relative risks. There is an established increased risk of breast cancer using contraceptive pills. Breast cancer is rare among young women. Increased risk during this time means a low number of additional cases. Second, there is no discussion of other well-known risk factors for breast cancer. Instead of calling upon extensive research about breast cancer, the authors cite themselves. Large prospective cohort studies do not imply heightened risks of breast cancer after abortion.

Third, there is no exploration of absolute and relative risks. There is an established increased risk of breast cancer using contraceptive pills. Breast cancer is rare among young women. Increased risk during this time means a low number of additional cases.

Fourth, methods and sources are unclear

In the News



Jonas Schöley @jschoeley

Life-expectancy drop in 2020 compared to previous years. Another look at the results of our forthcoming IJE paper. @OxfordDemSci @CPop_SDU #rstats Code: github.com/jschoeley/de0a... pic.twitter.com/vSe4xtvJ2X 2021-09-15, 3:54 AM

link