# Methods of Applied Statistics I

# STA2101H F LEC9101

Week 3

September 28 2022

The Economist = Menu Weekly

≡ Menu Weekly edition Q Search ✓

#### Science & technology | Peer review

# An influential academic safeguard is distorted by status bias

To those that have, more shall be given



Sep 14th 2022

📮 Save 🔍 < Share 🌰 Give



- 1. Upcoming events
- 2. Comments re HW
- 3. Linear Regression Part 3: recap, checking model assumptions, collinearity, model-building, p > n
- 4. In the News

# Upcoming

- September 29: CANSSI Ontario Research Day
   Schedule and Registration
- Distinguished Lecture Series in Statistical Sciences
- Xihong Lin, Harvard U
   Details and Registration
- September 29 3.30 89 Chestnut Street, 3rd Floor Lessons learned from the COVID-19 Pandemic: a statistician's reflection
- September 30 3.30 UY9014 Ensemble methods for testing a global null hypothesis
- September 30 1.00 Zoom data\_4\_lyf Toronto Data workshop
   "How the NFL blocks black coaches"

Applied Statistics I September 28 2022



2022 DLSS: Xihong Lin

Professor, Department of BiostatisticsCoordinating Director, Program in Quantitative Genomics; Harvard T.H. Chan School of Public Health; Professor of Statistics, Department of Statistics, Harvard University

# ... upcoming

October 3 3:30 Data Science ARES online

James Zou, Stanford

"AI for clinical trials and clinical trials for AI"

**Register here** 



#### STA2101F 2022

#### Due September 21 2022 11.59 pm

#### Homework to be submitted through Quercus

You can submit this HW in Word, Latex, or R Markdown, but in future please use R Markdown. If you are using Word or Latex with a R script for the computational work, then this R script should be provided as an Appendix. In the document itself you would just include properly formatted output.

You are welcome to discuss questions with others, but the solutions and code must be written independently. Any R output that is included in a solution should be formatted as part of the discussion (i.e. not cut and pasted from the Console).

The dataset vafer concerns a study on semiconductors. You can get more information about the data with ?vafer; you will first need library(faraway);data(vafer), and possibly install.packages("faraway"). The questions below are adapted from LM Ch.3.

- (a) Fit the linear model resist ~ x1 + x2 + x3 + x4. Extract the X matrix using the model.matrix function. How have the levels of the factors been coded? Level '-' has been coded 0, level '+' coded 1.
- (b) Compute the correlation between the columns of the X matrix. Why are there some missing values? The R output tells you the standard error of the intercept column is 0, so it seems likely that dividing by 0 in the formula for correlation is the problem. It's slightly more subtle, R will give Inf if the numerator is not 0 (try 5/0 for example), but gives NaH for 0/0, and corr(X, 1, 1, X, 2, 1), for example, returns NA. However corr(X) gives 1 for the correlation between the intercept and itself. It somehow recognizes that the numerator and denominator are equal, and that seems to take precedence over other conventions. Which is why it's good to study statistical computing.
- (c) What difference in resistance is expected when moving from the low to the high level of x1? The estimated difference in resistance is 25.8 units. Note that it is not necessary to add "all other variables held fixed", because of (d).
- (d) Refit the model without x4 and examine the regression coefficients and standard errors. What stayed the same and what changed? How is this related to the correlation matrix of X? The coefficients on x1, x2, x3 are unchanged, as the X<sup>2</sup>X matrix is diagonal. The estimated standard errors by the coefficients are slightly larger, because the residual to the standard errors by 2022.

Applied Statistics The estimated standard errors of the coefficients are slightly larger, because the residual sum of significant significant straightly larger, because the residual sum of significant signific

#### bonus, rmd, soln's

# Linear regression recap

• Analysis of variance:  $y^{\mathsf{T}}y = (y - X\hat{\beta})^{\mathsf{T}}(y - X\hat{\beta}) + \hat{\beta}^{\mathsf{T}}X^{\mathsf{T}}X\hat{\beta}$ 

Source	DF	SS	MS		
Regression	р — 1	SS <sub>REG</sub>	$\textit{RegMS} = \textit{SS}_{\textit{REG}}/(p-1)$		
Residual	n – p	RSS	ResMS = RSS/(n-p)		
Total (corrected)	n — 1	TSS			
$F = rac{RegMS}{ResMS} \sim F_{p-1,n-p}$ under					

• regression SS can be further partitioned

depends on the order

### ... Linear regression recap

- same principle can be used to test for sets of variables
- or for testing any linear constraint on  $\beta$

 $A\beta = c$ 

• numerator degrees of freedom for F-statistic depend on the rank of A

$$F_{1,
u}\equiv t_{
u}^2$$

• sometimes only an F-test can be used to assess the effect of an explanatory variable when?

.

# Q on Piazza

$$(y - \bar{y}\mathbf{1})^{\mathrm{T}}(y - \bar{y}\mathbf{1}) = (y - X\hat{\beta})^{\mathrm{T}}(y - X\hat{\beta}) + \hat{\beta}^{\mathrm{T}}(X^{\mathrm{T}}X)\hat{\beta} - n\bar{y}^{2}$$
$$\sum_{i=1}^{n} (y_{i} - \bar{y})^{2} = \sum_{i=1}^{n} (y_{i} - x_{i}^{\mathrm{T}}\hat{\beta})^{2} + \hat{\beta}_{2}^{\mathrm{T}}(X_{2}^{\mathrm{T}}X_{2})\hat{\beta}_{2}$$

### **Factor variables**

- F-tests are used when the columns to be removed form a group
- if a covariate is a factor, i.e. categorical, then 1m 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

 prostate\$gleason\_factor <- factor(prostate\$gleason) levels(prostate\$gleason\_factor)
 [1] "6" "7" "8" "9" model\_fac <- lm(lpsa ~ .-gleason, data=prostate)</li>

### ... factor variables

```
model_fac <- lm(lpsa ~ .-gleason, data=prostate)</pre>
 sumary(model_fac)
 Estimate Std. Error t value Pr(>|t|)
(Intercept)
               0.91328
                         0.84084 1.09
                                         0.2804
lcavol
               0.56999
                        0.09010
                                   6.33 1.1e-08
                         0.16961 2.76
                                         0.0070
lweight
               0.46879
              -0.02175
                         0.01136 -1.91
                                         0.0589
age
lbph
               0.09968
                         0.05898 1.69
                                         0.0946
svi
               0.74588
                         0.24740 3.01
                                         0.0034
                                  -1.31
                                         0.1941
lcp
              -0.12511
                         0.09559
               0.00499
                         0.00467 1.07
                                         0.2885
pgg45
                         0.21942 1.22
                                         0.2259
gleason_factor7 0.26761
gleason_factor8 0.49682
                         0.76927 0.65
                                         0.5201
gleason_factor9 -0.05621
                         0.50020
                                  -0.11
                                         0.9108
```

n = 97, p = 11, Residual SE = 0.70, R-Squared = 0.67

Applied Statistics I September 28 2022

model\_nog <- lm(lpsa ~ . - gleason - gleason\_factor, data = prostate)</pre>

anova(model\_fac, model\_nog) # compare two models

```
Analysis of Variance Table
```

```
Model 1: lpsa ~ (lcavol + lweight + age + lbph + svi + lcp + gleason +
    pgg45 + gleason_factor) - gleason - gleason_factor
Model 2: lpsa ~ (lcavol + lweight + age + lbph + svi + lcp + gleason +
    pgg45 + gleason_factor) - gleason
    Res.Df RSS Df Sum of Sq F Pr(>F)
    1    89 44.2
    2    86 42.7 3    1.48 0.99    0.4
Applied Statistics | September 28 2022
```

### **Model checking**

plot(model1)

**Applied Statistics I** 



Applied Statistics I

#### Model assumptions

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



- 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:  $\hat{\epsilon}_i = y_i \hat{y}_i$
- $\operatorname{Var}(\hat{\epsilon}) = \sigma^2(I H), \quad \operatorname{Var}(y_i \hat{y}_i) = \sigma^2(1 h_{ii})$   $\circ < h_{ii} < 1, \Sigma h_{ii} = p$
- i.e. don't all have the same variance
- hat matrix  $H = X(X^{T}X)^{-1}X^{T}$   $Hy = X(X^{T}X)^{-1}X^{T}y = X\hat{\beta} = \hat{y}$
- standardized residuals:  $r_i = rac{\hat{\epsilon}_i}{\tilde{\sigma}(1-h_{ii})^{1/2}}$  approx var 1
- Cook's distance  $C_i = \frac{(\hat{y} \hat{y}_{-i})^{\mathrm{T}}(\hat{y} \hat{y}_{-i})}{p\tilde{\sigma}^2} = \frac{r_i^2 h_{ii}}{p(1 h_{ii})}$

measure of influence

high leverage or high residual

#### Applied Statistics I September 28 2022

- standard diagnostics check for non-constant variance, influential observations
- and for normality of residuals

using qqnorm

- assumption of independence across *i* may be more important
- but more difficult to assess
- exception: observations collected over time LM-2, §6.1.3, LM-1 §4.1.3

# Aside on normal plots



Applied Statistics I

### ... Aside

```
library(ggplot2); library(nullabor); library(tidyverse)
df5_frame <- data.frame(x = rt(30, df = 5))
lineup_df5_data <- lineup(
   method = null_dist("x", dist = "norm", params = list(mean = 0, sd = 1)),
   true = df5_frame, n=12)</pre>
```

```
lineup_df5_data %>%
ggplot(aes(sample = x)) +
geom_qq_line() +
geom_qq() +
facet_wrap(~ .sample)
```

- Model  $y = X\beta + \epsilon$ , alternatively,
- $E(y \mid X) = X\beta$ ,  $Var(Y \mid X) = \sigma^2 I$
- plots of y against each column of x can be helpful
- for(i in 1:8){plot(prostate[,i],prostate[,9]... }
- added variable plots can be more helpful
- plot residuals from y on  $X_{-j}$  against residuals from  $x_j$  on  $X_{-j}$

partial regression plots slope of this line is  $\hat{eta}_j$ 

### **Prostate data**

**Applied Statistics I** 





Applied Statistics I

Figure 4.13 Partial regression (left) and partial residual (right) September 28 2022 plots for the savings data.



Applied StatisticsFigure & Aber Introducing another dimension to diagnostic plots. Shape is used denote the status variable on the left while faceting is used on the right.

# Collinearity

- simple model  $y_i = \beta_0 + \beta_1 x_{1i} + \beta_2 x_{2i} + \epsilon_i$ , i = 1, ..., n
- + if  $x_1 \perp x_2$ , then interpretation of  $\beta_1$  and  $\beta_2$  clear
- + if  $x_1 = x_2$  then  $\beta_1$  and  $\beta_2$  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<sup>T</sup>X 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

ridge, Lasso

### Three tasks related to linear regression

• Estimation of  $\beta$ , and estimation of its standard error – for inference about  $\mathbb{E}(y \mid x)$ 

alternatively comparing sub-models using *F*-tests

• Prediction of  $y_+$ , say, given a new vector of explanatory variables  $x_+$ 

LM-2 Ch.4, LM-1 §3.5, SM §8.3.2

 Model Selection: which explanatory variables do we need for prediction or inference?

These same questions arise in other models such as logistic regression, analysis of survival data, and so on, but the generic linear model is often a good starting point

• Prediction:  $y_+ = x_+^T \beta + \epsilon$ ;  $\hat{y}_+ = x_+^T \hat{\beta}$ ;  $\operatorname{var}(\hat{y}_+) = \sigma^2 x_+ (X^T X)^{-1} x_+$ 

assuming ...

error in expected response different from

prediction error  $\mathbb{E}(\mathbf{y}_+ - \hat{\mathbf{y}}_+)^2 = \sigma^2 + \operatorname{var}(\hat{\mathbf{y}}_+)$ 

- "analyses should be as simple as possible, but no simpler"
- What variables should we keep in the model ?
- Hierarchical models: some models have a natural hierarchy: polynomials, factorial structure, auto-regressive, sinusoidal, ...
- in these models the 'highest' level of the hierarchy is removed first
- e.g.  $y = \beta_0 + \beta_1 x + \beta_2 x^2 + \epsilon$  should \*not\* be simplified to  $y = \beta_0 + \beta_2 x^2 + \epsilon$
- e.g. if interaction terms are included, then main effects and other 2nd-order terms also need to be included:  $y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \beta_{12} x_1 x_2 + \beta_{11} x_1^2 + \beta_{22} x_2^2 + \epsilon$
- \*not\*  $y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \beta_{12} x_1 x_2 + \epsilon$  unless x = 0/1
- $y = \beta_0 + \beta_1 \sin(2\pi x) + \beta_2 \cos(2\pi x) + \beta_3 \sin(4\pi x) + \beta_4 \cos(4\pi x) + \epsilon$
- $y_t = \beta_0 + \alpha y_{t-1} + \epsilon$   $y_t = \beta_0 + \alpha_1 y_{t-1} + \alpha_2 y_{t-2} \epsilon$  \*not\*  $y_t = \beta_0 + \alpha_2 y_{t-2} + \epsilon$

### ... Model Selection

- testing procedures: forward selection, backward selection, stepwise selection
- it is quite common to fit all explanatory variables, and then drop if p > 0.05
- if estimates and estimated standard errors don't change very much, may be okay
- if estimates and estimated standard errors change a lot, cause for concern
- · if estimates change sign, points to possibly extreme confounding

```
step(model1)
Start: AIC=-58.32
lpsa ~ lcavol + lweight + age + lbph + svi + lcp + gleason +
    pgg45
```

	Df	Sum of Sq	RSS	AIC
- gleason	1	0.0412	44.204	-60.231
- pgg45	1	0.5258	44.689	-59.174
- lcp	1	0.6740	44.837	-58.853
<none></none>			44.163	-58.322
- age	1	1.5503	45.713	-56.975
Applied Statistics I	1	September 2 1.6835	8 <u>2022</u> 45.847	-56.693
1		0 5064	17 740	F0 740

#### step(model1)

```
...
Step: AIC=-61.37
lpsa ~ lcavol + lweight + age + lbph + svi
```

	Df	Sum of Sq	RSS	AIC
<none></none>			45.526	-61.374
- age	1	0.9592	46.485	-61.352
- lbph	1	1.8568	47.382	-59.497
- lweight	1	3.2251	48.751	-56.735
- svi	1	5.9517	51.477	-51.456
- lcavol	1	28.7665	74.292	-15.871

#### Call:

lm(formula = lpsa ~ lcavol + lweight + age + lbph + svi, data = prostate)

#### Coefficients:

(Intercept)	(Intercept) lcavol	lweight	age	lbph	svi
0.95100	0.56561	0.42369	-0.01489	0.11184	0.72095

### ... Model Selection

.

.

.

.

- Criterion-based procedures
- AIC, BIC, Mallows  $C_p$ ,  $R_a^2$

 $AIC = n \log(RSS/n) + 2p$ 

 $BIC = n \log(RSS/n) + \log(n)p$ 

$$C_p = RSS_p/\tilde{\sigma}^2 + 2p - n$$

$$R_a^2 = 1 - rac{ ilde{\sigma}_{model}^2}{ ext{TSS}/(n-1)}$$

Applied Statistics I September 28 2022 Sion  $AIC_c$  which may be better than AIC for linear models

most widely used RSS: residual sum of squares

#### In the News



E Menu Weekly edition Q Search ✓

#### Science & technology | Peer review

# An influential academic safeguard is distorted by status bias

To those that have, more shall be given



Sep 14th 2022

📮 Save 🛛 < Share 🏻 🌰 Give

KARL-FRANZENS-UNIVERSITÄT GRAZ UNIVERSITY OF GRAZ School of Business, Economics and Social Sciences



#### Nobel and novice: Author prominence affects peer review

Jürgen Huber, Sabiou Inoua, Rudolf Kerschbamer, Christian König-Kersting, Stefan Palan, Vernon L. Smith

Working Paper 2022-01 August 16, 2022

#### link

Applied Statistics I September 28 2022

# A randomized experiment

KARL-FRANZENS-UNIVERSITÄT GRAZ UNIVERSITY OF GRAZ School of Business, Economics and Social Sciences

#### Nobel and novice: Author prominence affects peer review

Jürgen Huber, Sabiou Inoua, Rudolf Kerschbamer, Christian König-Kersting, Stefan Palan, Vernon L. Smith

Working Paper 2022-01 August 16, 2022

**Applied Statistics I** 

#### Table 1: Invitations

	Low (LL)	Anonymized (AL, AA, AH)	High (HH)	Total
Invitations sent	781	2011	507	3299
Responses received	610	1591	410	2611
Invitations accepted	174	489	158	821
Acceptance rate	28.52%	30.74%	38.54%	31.44%
Anon. vs. Low		p = 0.3243		
Anon. vs. High		p = 0.0031		
Low vs. High		p = 0.0011		

Number of review invitations sent, number of replies received (declined or accepted), number of invitations accepted, fraction of invitations accepted when the review invitation listed the low prominence author (condition LL), no corresponding author (AL, AA, AH), or the high prominence author (HH). Two-sided Fisher's exact tests of invitation responses between conditions.



Figure 1: Recommendation percentages by condition. L stands for the relatively unknown author, A stands for anonymized and H stands for the highly prominent author. In conditions AL and AH, the invitation email is anonymized, but the respective corresponding author's name appears on the manuscript, while in AA both the invitation and the paper are anonymized. The tests are pairwise, two-sided Mann-Whitney U tests.

- common objectives
- to avoid systematic error, that is distortion in the conclusions arising from sources that do not cancel out in the long run
- to reduce the non-systematic (random) error to a reasonable level by replication and other techniques
- to estimate realistically the likely uncertainty in the final conclusions
- to ensure that the scale of effort is appropriate

# ... design of studies

- we concentrate largely on the careful analysis of individual studies
- in most situations synthesis of information from different investigations is needed
- but even there the quality of individual studies remains important
- examples include overviews (such as the Cochrane reviews)
- in some areas new investigations can be set up and completed relatively quickly; design of individual studies may then be less important

# ... design of studies

- formulation of a plan of analysis
- establish and document that proposed data are capable of addressing the research questions of concern
- main configurations of answers likely to be obtained should be set out
- · level of detail depends on the context
- even if pre-specified methods must be used, it is crucial not to limit analysis
- planned analysis may be technically inappropriate
- more controversially, data may suggest new research questions or replacement of objectives
- · latter will require confirmatory studies

# Unit of study and analysis

- smallest subdivision of experimental material that may be assigned to a treatment context: Expt
- Example: RCT unit may be a patient, or a patient-month (in crossover trial)
- Example: public health intervention unit is often a community/school/...
- split plot experiments have two classes of units of study and analysis
- in investigations that are not randomized, it may be helpful to consider what the primary unit of analysis would have been, had a randomized experiment been feasible
- the unit of analysis may not be the unit of interpretation ecological bias systematic difference between impact of *x* at different levels of aggregation
- on the whole, limited detail is needed in examining the variation within the unit of study

# Types of observational studies

- · secondary analysis of data collected for another purpose
- estimation of a some feature of a defined population (could in principle be found exactly)
- tracking across time of such features
- study of a relationship between features, where individuals may be examined
  - at a single time point
  - · at several time points for different individuals
  - at different time points for the same individual
- experiment: investigator has complete control over treatment assignment
- census
- meta-analysis: statistical assessment of a collection of studies on the same topic

- "distortion in the conclusions arising from irrelevant sources that do not cancel out in the long run"
- can arise through systematic aspects of, for example, a measuring process, or the spatial or temporal arrangement of units
- this can often be avoided by design, or adjustment in analysis
- can arise by the entry of personal judgement into some aspect of the data collection process
- this can often be avoided by randomization and blinding