Methods of Applied Statistics I

STA2101H F LEC9101

Week 3

September 29 2021





1. Upcoming events, HW 3

Office Hour Monday Oct 4 7pm 8.30pm \longrightarrow Tuesday Oct 5 7pm - 8 pm, Zoom

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

Upcoming

- Non-Central Squared Copulas: Properties and Applications
- Thursday 3.30 Link

About Bouchra Nasri



Dr. Nasri is Assistant Professor in Statistics at the Scholl of Public Health of Univesité de Montréal. Her research interests are dependence modelling, time series, and more recently spatial modelling. The main applications targeted by her research projects are related to climate change, public health and infectious diseases modelling. Dr. Nasri is an associate director of the new infectious diseases network OMNI-RÉUNIS.

Friday Oct 1 Toronto Data Workshop Zoom link

Toronto Data Workshop this Friday at noon (Toronto time) focuses on the recent Canadian election, with presentations from

- · Professor David Andrews on elections forecasting;
- · Professor Daniel Rubenson on the Canadian Election Study;
- · Johnson Vo on his model of the 2021 election; and
- Eric Zhu, Brian Diep, Ashely (Jing Yuan) Zhang, Kristin (Xi Yu Huang), and Tanvir Hyder on their model of the 2021 election.

Link: <u>https://utoronto.zoom.us/j/84277066292</u> Meeting ID: 842 7706 6292 Passcode: data_4_lyf September 29 2021

Applied Statistics I

HW 3

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 _{REG}	$\textit{RegMS} = \textit{SS}_{\textit{REG}}/(p-1)$
Residual	n – p	RSS	ResMS = RSS/(n-p)
Total (corrected)	n — 1	TSS	
F	$\overline{r} = \frac{Reg}{Res}$	$rac{MS}{MS}\simF$	_{p-1,n-p} under

• regression SS can be further partitioned

depends on the order

Analysis of variance:

anova(model1) Analysis of Variance Table

Response: lpsa

	\mathtt{Df}	Sum Sq	Mean Sq	F value	Pr(>F)	
lcavol	1	69.003	69.003	137.4962	< 2.2e-16	***
lweight	1	5.949	5.949	11.8531	0.0008832	***
age	1	0.420	0.420	0.8369	0.3627958	
lbph	1	1.069	1.069	2.1302	0.1479839	
svi	1	5.952	5.952	11.8594	0.0008806	***
lcp	1	0.129	0.129	0.2576	0.6130533	
gleason	1	0.708	0.708	1.4098	0.2382837	
pgg45	1	0.526	0.526	1.0476	0.3088604	
Residuals	88	44.163	0.502			

```
> summary(model1)
Call:
lm(formula = lpsa ~ ., data = prostate)
Coefficients:
```

	Estimate	Std. Error	t value	Pr(>
(Intercept)	0.669337	1.296387	0.516	0.60
lcavol	0.587022	0.087920	6.677	2.11e
lweight	0.454467	0.170012	2.673	0.00
age	-0.019637	0.011173	-1.758	0.08
lbph	0.107054	0.058449	1.832	0.07
svi	0.766157	0.244309	3.136	0.00
lcp	-0.105474	0.091013	-1.159	0.24
gleason	0.045142	0.157465	0.287	0.77
pgg45	0.004525	0.004421	1.024	0.30

Residual standard error: 0.7084 on 88 degrees of

Applied Statistics I September 29 2021

Analysis of variance:

anova(model1) Analysis of Variance Table

Response: lpsa

	$\mathtt{D}\mathtt{f}$	Sum Sq	Mean Sq	F value	Pr(>F)	
lcavol	1	69.003	69.003	137.4962	< 2.2e-16	***
lweight	1	5.949	5.949	11.8531	0.0008832	***
age	1	0.420	0.420	0.8369	0.3627958	
lbph	1	1.069	1.069	2.1302	0.1479839	
svi	1	5.952	5.952	11.8594	0.0008806	***
lcp	1	0.129	0.129	0.2576	0.6130533	
gleason	1	0.708	0.708	1.4098	0.2382837	
pgg45	1	0.526	0.526	1.0476	0.3088604	
Residuals	88	44.163	0.502			

> drop1(model1)
Single term deletions

1

lcp

pgg45

Model: lpsa ~ lcavol + lweight + age + lbph + svi + lo pgg45 Df Sum of Sq RSS ATC <none> 44.163 -58.322 lcavol 1 22.3721 66.535 -20.567 lweight 1 3.5861 47.749 -52.749 age 1 1.5503 45.713 -56.975 lbph 1 1.6835 45.847 -56.693 svi 1 4.9355 49.099 -50.046

0.6740 44.837 -58.853

1 0.5258 44.689 -59.174

gleason 1 0.0412 44.204 -60.231

... Linear regression recap

.

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

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

- numerator degrees of freedom for F-statistic depend on the rank of A
- sometimes only an F-test can be used to assess the effect of an explanatory variable

7

 $A\beta = c$



• §3.3: permutation test - doesn't rely on normal assumption

• §3.5: confidence intervals for β_j

and regions for $(\beta_j, \beta_{j'})$ confint; ellipse see Fig 3.2

• §3.6: bootstrap inference for β_j

resample $\hat{\epsilon}_i$

Model checking

SM 8.6, LM-2 Ch. 6, LM-1 Ch. 4

plot(model1)

Applied Statistics I



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

10

Applied Statistics I

Model assumptions



11

- 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: $\hat{\epsilon}_i = \mathbf{y}_i \hat{\mathbf{y}}_i$
- $\operatorname{Var}(\hat{\epsilon}) = \sigma^2(I H), \quad \operatorname{Var}(y_j \hat{y}_j) = \sigma^2(1 h_{jj})$ $\circ < h_{jj} < 1, \Sigma h_{jj} = p$
- 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 29 2021

- 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

15

```
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



18



Applied Statistics I

Figure 4.13 Partial regression (left) and partial residual (right) September 29 2021 plots for the savings data.

19



Applied StatisticsFigure 6cHaber Introducing another dimension to diagnostic plots. Shape is used denote the status variable on the left while faceting is used on the right.

Read Chapter 6 of LM-2 or Chapter 4 of LM-1, replicating the results

Read Section 8.6 of SM, working through the algebra

PhD, Stats

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 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 ridge, Lasso more next week

Three tasks related to linear regression

• Estimation of β , 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	45.847	-56.693
1		0 5004	47 740	50 740

step(model1)

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

		Df	Sum of	Sq	RSS	AIC	
<1	none>				45.526	-61.374	
-	age	1	0.95	592	46.485	-61.352	
-	lbph	1	1.8	568	47.382	-59.497	
-	lweight	1	3.22	251	48.751	-56.735	
-	svi	1	5.98	517	51.477	-51.456	
-	lcavol	1	28.76	665	74.292	-15.871	

Call:

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

Coefficients:

(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, R²_a

most widely used *RSS*: residual sum of squares

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

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

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

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

- SM has yet another version AIC_c which may be better than AIC for linear models
- C_p and R_a^2 are only useful for linear models; AIC and BIC more general



- "In "The First Political Order: How Sex Shapes Governance and National Security Worldwide", Ms Hudson, Ms Bowen and Ms Nielsen rank 176 countries on a scale of o to 16 for what they call the "patrilineal/fraternal syndrome". This is a composite of such things as unequal treatment of women in family law and property rights, early marriage for girls, patrilocal marriage, polygamy, bride price, son preference, violence against women and social attitudes towards it"
- "Ms Hudson and her co-authors tested the relationship between their patrilineal syndrome and violent political instability. They ran various regressions on their 176 countries, controlling for other things that might foster conflict, such as ethnic and religious strife, colonial history ..."

- "They did not prove that the syndrome caused instability that would require either longitudinal data that have not yet been collected or natural experiments that are virtually impossible with whole countries"
- "But they found a strong statistical link. The syndrome explained three-quarters of the variation in a country's score on the Fragile States index compiled by the Fund for Peace, a think-tank in Washington."
- Book website
- Blog

- "Examining approximately 176 nations, we examined whether national outcomes such as conflict, terrorism, poverty, and so forth, were significantly associated with a subordinative first political order, while controlling for background factors such as level of urbanization, levels of ethnic fractionalization, colonial history, and so forth.
- "Holding these characteristics constant, is that subordinative order strongly related to national outcomes? In all we examined 122 national outcome measures related to conflict, stability, governance, prosperity, health, demographics, education, environmental preservation, and social progress.
- "Across all 122 outcome variables, the subordination of women was both significant and the explanatory factor with the largest or second largest effect size over 70% of the time.

- 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