Methods of Applied Statistics I

STA2101H F LEC9101

Week 1

September 14 2022



← Tweet

The Guy Medal in Gold is awarded to Nancy Reid @reid_nancy for her pioneering work on higher-order approximate inference



E.E.A.DM. Can 10, 0000 Tuittan Mak Ann

Ś	Chrome File	Edit View	History	Bookmarks	Profiles	Tab Windov	v Help		₩ 5	•	🔊 O 省	88 0) 💌	* 💷	Ŷ	۹ 🛢	Wed	Sep 1
•	🛑 🔍 M. Inbox - nancy.reid1: x 📴 12 photo printing oj: x 🛓 Contact Annex Pho: x N Chen, Wei Faculty: x 💥 www.cs.cmu.edu/- x 💩 Yao Xie ISyE Geo: x 👩 Program 🛛 x +																	
÷	\rightarrow C $rac{1}{2}$ virt	ual.oxfordab	stracts.co	m /#/event/272	26/program											QÅ	$\stackrel{\circ}{\Box}$	* 1
€	E Royal Statistical Society 2022 International Conference																	
÷	Program Q Search program Timezone ~ Dates ~ Topic streams ~ Your bookmarks ~																	
	Titles	BST 13																
00	Topic Streams	08:3	0 09:00-10:00 Contribute	d: Infectious	09:00-10:00 Contributed:	Novel	09:00-10:00	09:00 Conti	-10:00 ributed: Causal	D	09:00-10:00 Contributed: Clin	ical Д	09:00-10: Contribu	.00 uted: Design o	<mark>,</mark> ۵	09:00-10:4 Contribu	10 ted: Repc	orting
			disease Topic strear Statistics	ns Applications of	applications sets Topic streams	and data	point process models Topic streams Environmental & Spatial Statistics	infere Topic Statis	ence streams Medical tics		Trials Topic streams Mec Statistics		experim studies Topic stru Theory	eents and		uncertain Topic stre Theory	nty ams Meth	ods &
		10:0																
		10:1	10:10-11:10 Keynote 2 (President Lecture - A Topic strear	Campion 's Invited) Adrian Raftery ms Plenary														
		11:1	11:10-11:40 Refreshme Topic strear	ents ns Break														
			11:40-13:00 The role of supporting ballet dane Topic strear Statistics	f statistics in a professional sers ns Applications of	11:40-13:00 Democratisa statistics in Topic streams Industry & Fina	Lation of GSK Business, ance	11:40-13:00 RSS Statistical Ambassadors' Showcase Topic streams Communicating & Teaching Statistics	11:40 Data - an i MLO Topic	-13:00 Science in Indus ntroduction to ps streams Data Scie	try D	11:40-13:00 Model selection discrimination fo environmental a spatial application Spatial Statistics	and or nd ons ronmental &	11:40-13: What is Topic stri Statistics	00 your estimanc eams Medical	וי <mark>וי</mark>	11:40-13:0 Papers fi Journals Cox regn developr perspect Topic stre Theory	10 rom the F i: 50 year: ession: ments and tives ams Meth	ISS s of d ods &



- 1. Course introduction: course details, evaluation, syllabus, people
- 2. Upcoming events of interest
- 3. Review of linear regression
- 4. In the news: \longrightarrow at the conference

Applied Statistics I

STA 2101F: Methods of Applied Statistics I Wednesday, 10am – 1 pm Eastern September 14 – December 7 2022 SF 3201

From the calendar:

This course will focus on principles and methods of applied statistical science. It is designed for MSc and PhD students in Statistics, and is required for the Applied Paper of the PhD comprehensive exams. The topics covered include: planning of studies, review of linear models, analysis of random and mixed effects models, model building and model selection, theory and methods for generalized linear models, and an introduction to nonparametric regression. Additional topics will be introduced as needed in the context of case studies in data analysis.

Prerequisites: ECO374H1/ECO375H1/STA302H1 (regression); STA305H1 (design of studies)

September 14 2022 Course Delivery:

On Sontombor 14, the along will be delivered online at the scheduled

Applied Statistics I

STA 2101F: Methods of Applied Statistics I Wednesday, 10am – 1 pm Eastern September 14 – December 7 2022 SF 3201

From the calendar:

This course will focus on principles and methods of applied statistical science. It is designed for MSc and PhD students in Statistics, and is required for the Applied Paper of the PhD comprehensive exams. The topics covered include: planning of studies, review of linear models, analysis of random and mixed effects models, model building and model selection, theory and methods for generalized linear models, and an introduction to nonparametric regression. Additional topics will be introduced as needed in the context of case studies in data analysis.

Prerequisites: ECO374H1/ECO375H1/STA302H1 (regression); STA305H1 (design of studies)

September 14 2022 Course Delivery:

On Sontombor 14, the along will be delivered online at the scheduled

- Grading
- Academic Integrity
- Computing
- References
 Modules



• Contact

Use Piazza for course questions; email for personal questions



- 1. Course introduction: technical issues, course details, evaluation, syllabus
- 2. Upcoming events of interest
- 3. Review of linear regression
- 4. In the news: \longrightarrow at the conference

Upcoming events

- Thursdays 3.30 Departmental Seminar
- Mondays 3.30 Data Science and Applied Research Seminar
- Fridays 12.00 Toronto Data Workshop
- Special:
 - September 29: CANSSI Ontario Research Day
 - September 29, 30 3.30: Distinguished Lecture Series in Statistical Sciences link



link

2022 DLSS: Xihong Lin

link

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

Sep 29 (3:30-4:30 pm): Lessons Learned from the COVID-19 Pandemic: A Statistician's Reflection Sep 30 (3:30-4:30 pm): Ensemble Methods for Testing a Global Null Hypothesis



- 1. Course introduction: course details, evaluation, syllabus
- 2. Upcoming events of interest
- 3. Review of linear regression
- 4. Steps in analysis
- 5. In the news: \longrightarrow at the conference

• Model:

$$Y = X \beta + \epsilon$$

• Model:

$$Y_{n \times 1} = X_{n \times p} \beta_{p \times 1} + \epsilon_{n \times 1}$$

Review of Linear Regression

• Model:

$$Y_{n \times 1} = X_{n \times p} \beta_{p \times 1} + \epsilon_{n \times 1}$$

• Equivalently:

 $y_i =$

• Model:

$$Y_{n \times 1} = X_{n \times p} \beta_{p \times 1} + \epsilon_{n \times 1}$$

• Equivalently:

 $y_i =$

- Standard Assumptions
 - y_i independent equivalently ϵ_i independent

•
$$\mathbb{E}(\epsilon_i) = 0$$

- $var(\epsilon_i) = \sigma^2$
- x_i known, β to be estimated

y is often called response why? constant

 x_i often called explanatory variables

v is often called response

x; often called explanatory variables

Model:

$$Y_{n \times 1} = X_{n \times p} \beta_{p \times 1} + \epsilon_{n \times 1}$$

• Equivalently:

 $y_i =$

- Standard Assumptions
 - y_i independent equivalently ϵ_i independent
 - $\mathbb{E}(\epsilon_i) = 0$
 - $var(\epsilon_i) = \sigma^2$
 - x_i known, β to be estimated
- More concisely:

 $\mathbb{E}(Y \mid X) =$, $var(Y \mid X) =$

1??

whv?

constant

Nice big equation:

$$\begin{pmatrix} y_1 \\ \vdots \\ y_n \end{pmatrix} = \begin{pmatrix} x_{11} & \dots & x_{1p} \\ \vdots & \vdots & \vdots \\ x_{n1} & \dots & x_{np} \end{pmatrix} \begin{pmatrix} & & \\ & \vdots \\ & & \end{pmatrix} + \begin{pmatrix} & & \\ & \vdots \\ & & \end{pmatrix}$$

Nice big equation:

$$\begin{pmatrix} y_1 \\ \vdots \\ y_n \end{pmatrix} = \begin{pmatrix} x_{11} & \dots & x_{1p} \\ \vdots & \vdots & \vdots \\ x_{n1} & \dots & x_{np} \end{pmatrix} \begin{pmatrix} & \vdots \\ & \end{pmatrix} + \begin{pmatrix} & \vdots \\ & \end{pmatrix}$$

Or, if you prefer:

$$y_i = x_{i_1}\beta_1 + x_{i_2}\beta_2 + \cdots + x_{i_p}\beta_p + \epsilon_i, \quad \epsilon_i \sim i = 1, \dots, n$$

Nice big equation:

$$\begin{pmatrix} y_1 \\ \vdots \\ y_n \end{pmatrix} = \begin{pmatrix} x_{11} & \dots & x_{1p} \\ \vdots & \vdots & \vdots \\ x_{n1} & \dots & x_{np} \end{pmatrix} \begin{pmatrix} & \vdots \\ & \end{pmatrix} + \begin{pmatrix} & \vdots \\ & \end{pmatrix}$$

Or, if you prefer:

$$y_i = x_{i1}\beta_1 + x_{i2}\beta_2 + \dots + x_{ip}\beta_p + \epsilon_i, \quad \epsilon_i \sim i = 1, \dots, n$$

Or, if you prefer:

 $\mathbb{E}(\mathbf{y}_i \mid \mathbf{x}_i) = \mathbf{x}_i^{\mathrm{T}} \boldsymbol{\beta}, \qquad \quad \text{var}(\mathbf{y}_i \mid \mathbf{x}_i) = \sigma^2, \qquad \quad i = 1, \dots, n$

y_i independent

e.g.?

• often not completely clear: *X* might be fixed by design, or measured on each individual

e.g.?

- often not completely clear: X might be fixed by design, or measured on each individual
- If measured, then should we consider its distribution? E.g. should our model be $(y_i, x_i^T) \sim ??$ some (p+1)-dimensional distribution

- often not completely clear: X might be fixed by design, or measured on each individual
- If measured, then should we consider its distribution? E.g. should our model be $(y_i, x_i^T) \sim ??$ some (p + 1)-dimensional distribution
- Almost always in regression settings we condition on X, as on previous slide

ancillary statistic

e.g.?

- often not completely clear: X might be fixed by design, or measured on each individual
- If measured, then should we consider its distribution? E.g. should our model be $(y_i, x_i^T) \sim ??$ some (p + 1)-dimensional distribution
- Almost always in regression settings we condition on X, as on previous slide

ancillary statistic

e.g.?

• often not emphasized: interpretation of β_i

- often not completely clear: X might be fixed by design, or measured on each individual
- If measured, then should we consider its distribution? E.g. should our model be $(y_i, x_i^T) \sim ??$ some (p + 1)-dimensional distribution
- Almost always in regression settings we condition on X, as on previous slide

ancillary statistic

e.g.?

- often not emphasized: interpretation of β_j
 - version 1: effect on the expected response of a unit change in *j*th explanatory variable, all other variables held fixed

- often not completely clear: X might be fixed by design, or measured on each individual
- If measured, then should we consider its distribution? E.g. should our model be $(y_i, x_i^T) \sim ??$ some (p + 1)-dimensional distribution
- Almost always in regression settings we condition on X, as on previous slide

ancillary statistic

e.g.?

- often not emphasized: interpretation of β_i
 - version 1: effect on the expected response of a unit change in *j*th explanatory variable,

all other variables held fixed

version 2:

$$eta_j = rac{\partial \mathbb{E}(\mathbf{y}_i \mid \mathbf{x}_{ij})}{\partial \mathbf{x}_{ij}}$$

 $\frac{\partial \mathbb{E}(y \mid x_j)}{\partial x_j}$ notation ambiguous, see CD §6.5.2

• Definition

$$\hat{eta}_{LS} := \min_{eta} \sum_{i=1}^n (y_i - x_i^{\mathrm{\scriptscriptstyle T}}eta)^2$$

• Definition

$$\hat{\beta}_{LS} := \min_{\beta} \sum_{i=1}^{n} (y_i - x_i^{\mathrm{T}} \beta)^2$$

• Equivalently,

• Definition

$$\hat{\beta}_{LS} := \min_{\beta} \sum_{i=1}^{n} (y_i - x_i^{\mathrm{T}} \beta)^2$$

- Equivalently,
- Equivalently,

$$\hat{\beta}_{LS} :=$$

L2 distance

• Definition

$$\hat{\beta}_{LS} := \min_{\beta} \sum_{i=1}^{n} (y_i - x_i^{\mathrm{T}} \beta)^2$$

- Equivalently,
- Equivalently,

$$\hat{\beta}_{LS} :=$$

L2 distance

- Equivalently, $\hat{\beta}_{LS}$ is the solution of the score equation

 $X^{\mathrm{T}}(y - X\beta) = 0$

?how?

• Definition

$$\hat{\beta}_{LS} := \min_{\beta} \sum_{i=1}^{n} (y_i - x_i^{\mathrm{T}} \beta)^2$$

- Equivalently,
- Equivalently,

$$\hat{\beta}_{LS} :=$$

 $\hat{\beta}_{is} =$

L2 distance

- Equivalently, $\hat{\beta}_{LS}$ is the solution of the score equation

$$X^{\mathrm{T}}(y - X\beta) = 0$$

?how?

12

Solution

Applied Statistics I September 14 2022

check dimensions

Solution

$$\hat{eta}_{LS} = (X^{ ext{ iny T}}X)^{-1}(X^{ ext{ iny T}}y)$$

check dimensions

ASIDE: here and following all assume X is fixed

Solution

$$\hat{\beta}_{LS} = (X^{\mathrm{\scriptscriptstyle T}}X)^{-1}(X^{\mathrm{\scriptscriptstyle T}}y)$$

check dimensions

• Expected value

 $\mathbb{E}(\hat{eta}_{LS}) =$

why?

ASIDE: here and following all assume X is fixed

Solution

$$\hat{\beta}_{LS} = (X^{\mathrm{\scriptscriptstyle T}}X)^{-1}(X^{\mathrm{\scriptscriptstyle T}}y)$$

check dimensions

• Expected value

 $\mathbb{E}(\hat{eta}_{LS}) =$

why?

· Least squares estimates are unbiased

ASIDE: here and following all assume X is fixed

Solution

$$\hat{\beta}_{LS} = (X^{\mathrm{\scriptscriptstyle T}}X)^{-1}(X^{\mathrm{\scriptscriptstyle T}}y)$$

check dimensions

• Expected value

$$\mathbb{E}(\hat{eta}_{LS}) =$$

why?

- Least squares estimates are unbiased
- Variance

really variance-covariance matrix

$$\operatorname{var}(\hat{\beta}_{LS}) = (X^{\mathrm{T}}X)^{-1}X^{\mathrm{T}}\operatorname{var}(y)X(X^{\mathrm{T}}X)^{-1} = (X^{\mathrm{T}}X)^{-1}X^{\mathrm{T}}\sigma^{2}IX(X^{\mathrm{T}}X)^{-1} = \sigma^{2}(X^{\mathrm{T}}X)^{-1}$$

ASIDE: here and following all assume X is fixed

SM §8.2

• If we further assume $\epsilon_i \sim N(0, \sigma^2)$ (and independent across *i*), then

- If we further assume $\epsilon_i \sim N(0, \sigma^2)$ (and independent across *i*), then
- $y \mid X \sim N(X\beta, \sigma^2 I)$, and

- If we further assume $\epsilon_i \sim N(0, \sigma^2)$ (and independent across *i*), then
- $y \mid X \sim N(X\beta, \sigma^2 I)$, and
- the likelihood function is

$$L(\beta,\sigma^2;\mathbf{y}) = \frac{1}{(2\pi\sigma^2)^{n/2}} \exp\left\{-\frac{1}{2\sigma^2}(\mathbf{y}-\mathbf{X}\beta)^{\mathsf{T}}(\mathbf{y}-\mathbf{X}\beta)\right\},\,$$

- If we further assume $\epsilon_i \sim N(0, \sigma^2)$ (and independent across *i*), then
- $y \mid X \sim N(X\beta, \sigma^2 I)$, and
- the likelihood function is

$$L(\beta,\sigma^2;\mathbf{y}) = \frac{1}{(2\pi\sigma^2)^{n/2}} \exp\left\{-\frac{1}{2\sigma^2}(\mathbf{y}-\mathbf{X}\beta)^{\mathsf{T}}(\mathbf{y}-\mathbf{X}\beta)\right\},\,$$

• the log-likelihood function is

$$\ell(\beta,\sigma^2; \mathbf{y}) = -\frac{n}{2}\log(\sigma^2) - \frac{1}{2\sigma^2}(\mathbf{y} - \mathbf{X}\beta)^{\mathrm{T}}(\mathbf{y} - \mathbf{X}\beta),$$

constants in params don't matter

- If we further assume $\epsilon_i \sim N(o, \sigma^2)$ (and independent across *i*), then
- $y \mid X \sim N(X\beta, \sigma^2 I)$, and
- the likelihood function is

$$L(\beta,\sigma^2;\mathbf{y}) = \frac{1}{(2\pi\sigma^2)^{n/2}} \exp\left\{-\frac{1}{2\sigma^2}(\mathbf{y}-\mathbf{X}\beta)^{\mathsf{T}}(\mathbf{y}-\mathbf{X}\beta)\right\},\,$$

• the log-likelihood function is

$$\ell(\beta,\sigma^2;\mathbf{y}) = -\frac{n}{2}\log(\sigma^2) - \frac{1}{2\sigma^2}(\mathbf{y} - \mathbf{X}\beta)^{\mathrm{T}}(\mathbf{y} - \mathbf{X}\beta),$$

constants in params don't matter

• the maximum likelihood estimate of β is

$$\hat{\beta}_{ML} = (X^{\mathrm{T}}X)^{-1}X^{\mathrm{T}}y = \hat{\beta}_{LS}$$

- maximum likelihood estimate of β is

$$\hat{eta}_{\mathsf{ML}} = (X^{ ext{ iny{T}}}X)^{-1}X^{ ext{ iny{T}}}y = \hat{eta}_{\mathsf{LS}}$$

- maximum likelihood estimate of β is

$$\hat{\beta}_{ML} = (X^{\mathrm{\scriptscriptstyle T}}X)^{-1}X^{\mathrm{\scriptscriptstyle T}}y = \hat{\beta}_{LS}$$

- distribution of $\hat{\beta}$ is normal

why?

$$\hat{eta} \sim N_{p}(eta, \sigma^{2}(X^{\mathrm{\scriptscriptstyle T}}X)^{-1})$$

- maximum likelihood estimate of β is

$$\hat{\beta}_{ML} = (X^{\mathrm{\scriptscriptstyle T}}X)^{-1}X^{\mathrm{\scriptscriptstyle T}}y = \hat{\beta}_{LS}$$

- distribution of $\hat{\beta}$ is normal

 $\hat{\beta} \sim N_{p}(\beta, \sigma^{2}(X^{\mathrm{T}}X)^{-1})$

• distribution of $\hat{\beta}_j$ is

$$N(\beta_j, \sigma^2(X^{\mathrm{T}}X)_{jj}^{-1}), \quad j = 1, \dots, p$$

whv?

- maximum likelihood estimate of β is

$$\hat{\beta}_{ML} = (X^{\mathrm{\scriptscriptstyle T}}X)^{-1}X^{\mathrm{\scriptscriptstyle T}}y = \hat{\beta}_{LS}$$

- distribution of $\hat{\beta}$ is normal

 $\hat{\beta} \sim N_{p}(\beta, \sigma^{2}(X^{\mathrm{T}}X)^{-1})$

• distribution of $\hat{\beta}_j$ is

$$N(\beta_j, \sigma^2 (X^T X)_{jj}^{-1}), \quad j = 1, ..., p$$

aximum likelihood estimate of σ^2 is $\frac{1}{n}(y - X\hat{\beta})^T(y - X\hat{\beta})$

• m

why?

• maximum likelihood estimate of β is

$$\hat{\beta}_{ML} = (X^{\mathrm{\scriptscriptstyle T}}X)^{-1}X^{\mathrm{\scriptscriptstyle T}}y = \hat{\beta}_{LS}$$

- distribution of $\hat{\beta}$ is normal

 $\hat{\beta} \sim N_{p}(\beta, \sigma^{2}(X^{\mathrm{T}}X)^{-1})$

• distribution of $\hat{\beta}_j$ is

$$N(\beta_j, \sigma^2(X^{\mathrm{T}}X)_{jj}^{-1}), \quad j = 1, \dots, p$$

- maximum likelihood estimate of σ^2 is $\frac{1}{n}(y X\hat{\beta})^{\mathrm{T}}(y X\hat{\beta})$
- but we use

$$\tilde{\sigma}^2 = \frac{1}{n-p} (y - X\hat{\beta})^{\mathrm{T}} (y - X\hat{\beta})$$

Applied Statistics I September 14 2022

15

why?

(1) I'm lost

(2) I'm good

(3) I'm bored

HW Question Week 1

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

(a) Fit the linear model regist α x1 + x2 + x3 + x4. Extract the X matrix using the

Inference

• If you really like likelihood theory, the expected Fisher information is SM §8.2.3

$$\mathcal{I}(\beta,\sigma^2) = \begin{pmatrix} \sigma^{-2} X^{\mathrm{T}} X & \mathbf{0} \\ \mathbf{0} & \frac{1}{2} \mathbf{n} \sigma^{-4} \end{pmatrix}$$

 \mathcal{I}^{-1} gives (asymptotic) variance of MLE

Inference

• If you really like likelihood theory, the expected Fisher information is SM §8.2.3

$$\mathcal{I}(\beta,\sigma^2) = \begin{pmatrix} \sigma^{-2} X^{\mathrm{T}} X & \mathbf{0} \\ \mathbf{0} & \frac{1}{2} \mathbf{n} \sigma^{-4} \end{pmatrix}$$

 \mathcal{I}^{-1} gives (asymptotic) variance of MLE

• but just using previous slide we have

$$rac{\hat{eta}_j - eta_j}{\sigma[\{(X^{ ext{ iny T}}X)^{-1}\}_{jj}\}]^{1/2}} \sim N(0, 1)$$

Inference

• If you really like likelihood theory, the expected Fisher information is SM §8.2.3

$$\mathcal{I}(\beta,\sigma^2) = \begin{pmatrix} \sigma^{-2} X^{\mathrm{T}} X & \mathbf{0} \\ \mathbf{0} & \frac{1}{2} \mathbf{n} \sigma^{-4} \end{pmatrix}$$

 \mathcal{I}^{-1} gives (asymptotic) variance of MLE

• but just using previous slide we have

$$rac{\hat{eta}_j - eta_j}{\sigma[\{(X^{ ext{ iny T}}X)^{-1}\}_{jj}\}]^{1/2}} \sim N(\mathsf{O},\mathsf{1})$$

and

$$rac{\hat{eta}_j - eta_j}{ ilde{\sigma}[\{(X^{ ext{ iny X}})^{-1}\}_{jj}\}]^{1/2}} \sim \mathsf{t}_{n-p}$$

Example

install.packages("faraway")
library(faraway)
data(prostate)
head(prostate)

Example

```
install.packages("faraway")
library(faraway)
data(prostate)
head(prostate)
```

```
model1 <- lm(lpsa ~ ., data = prostate)</pre>
   summary(model1)
   Coefficients:
               Estimate Std. Error t value Pr(>|t|)
   (Intercept)
               0.669337 1.296387 0.516 0.60693
   lcavol
               0.587022 0.087920 6.677 2.11e-09 ***
   lweight
               0.454467 0.170012 2.673 0.00896 **
Appliageatistics | Sept0m0496372 0.011173 -1.758 0.08229 .
  lhnh
               0 107054
                          0 050//0 1 020 0 070/0
```

Example

```
summary(model1)
Coefficients:
```

	Estimate	Std. Error	t value	Pr(> t)	
(Intercept)	0.669337	1.296387	0.516	0.60693	
lcavol	0.587022	0.087920	6.677	2.11e-09	***
lweight	0.454467	0.170012	2.673	0.00896	**
age	-0.019637	0.011173	-1.758	0.08229	
lbph	0.107054	0.058449	1.832	0.07040	
svi	0.766157	0.244309	3.136	0.00233	**
lcp	-0.105474	0.091013	-1.159	0.24964	
gleason	0.045142	0.157465	0.287	0.77503	
pgg45	0.004525	0.004421	1.024	0.30886	

Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1 Applied Statistics | September 14 2022

		322 views
-	Women in Statistics and Data Science	
- Pr	Follow @WomenInStat 18.7K followers	♀ tì ♡ 🎽
Jul 23rd 2020	0, <u>14 tweets, 4 min read</u>	
		☐ Bookmark 🛛 🖾 Save as PDF 🛛 + My Authors

Today, we're going to play a game I'm calling "IT'S JUST A LINEAR MODEL" (IJALM).

It works like this: I name a model for a quantitative response Y, and then you guess whether or not IJALM.

•
$$\mathbf{y}_i = \beta_0 + \beta_1 \mathbf{x}_i + \epsilon_i$$
, $i = 1, \dots, n$

•
$$y_i = \beta_0 + \beta_1 x_i + \epsilon_i$$
, $i = 1, \ldots, n$

1st column of X?

• $y_i = \beta_0 + \beta_1 x_i + \beta_2 x_i^2 + \beta_3 x_i^3 + \beta_4 x_i^4 + \beta_5 x_i^5 \epsilon_i$

•
$$y_i = \beta_0 + \beta_1 x_i + \epsilon_i$$
, $i = 1, \ldots, n$

- $y_i = \beta_0 + \beta_1 x_i + \beta_2 x_i^2 + \beta_3 x_i^3 + \beta_4 x_i^4 + \beta_5 x_i^5 \epsilon_i$
- $y_i = \beta_0 \pm \beta_1 + \epsilon_i$

•
$$y_i = \beta_0 + \beta_1 x_i + \epsilon_i$$
, $i = 1, \ldots, n$

- $y_i = \beta_0 + \beta_1 x_i + \beta_2 x_i^2 + \beta_3 x_i^3 + \beta_4 x_i^4 + \beta_5 x_i^5 \epsilon_i$
- $y_i = \beta_0 \pm \beta_1 + \epsilon_i$
- $y_i = \beta_0 + \beta_1 \sin(x_i) + \beta_2 \cos(x_i) + \epsilon_i$



•
$$y_i = \beta_0 + \beta_1 x_i + \epsilon_i$$
, $i = 1, \ldots, n$

- $y_i = \beta_0 + \beta_1 x_i + \beta_2 x_i^2 + \beta_3 x_i^3 + \beta_4 x_i^4 + \beta_5 x_i^5 \epsilon_i$
- $y_i = \beta_0 \pm \beta_1 + \epsilon_i$
- $y_i = \beta_0 + \beta_1 \sin(x_i) + \beta_2 \cos(x_i) + \epsilon_i$
- $y_i = \gamma_0 x_{1i}^{\gamma_1} x_{2i}^{\gamma_2} \eta_i$, $\eta_i \sim \text{positive r.v.}$



•
$$y_i = \beta_0 + \beta_1 x_i + \epsilon_i$$
, $i = 1, \ldots, n$

•
$$y_i = \beta_0 + \beta_1 x_i + \beta_2 x_i^2 + \beta_3 x_i^3 + \beta_4 x_i^4 + \beta_5 x_i^5 \epsilon_i$$

•
$$y_i = \beta_0 \pm \beta_1 + \epsilon_i$$

•
$$y_i = \beta_0 + \beta_1 \sin(x_i) + \beta_2 \cos(x_i) + \epsilon_i$$

• $y_i = \gamma_0 x_{1i}^{\gamma_1} x_{2i}^{\gamma_2} \eta_i$, $\eta_i \sim \text{positive r.v.}$

•
$$y_i = \varphi_0 + \sum_{k=1}^{K} \varphi_k \mathbf{s}_k(\mathbf{x}_i) + \epsilon_i$$



e.g smoothing splines

i + c, i - 1 n

Applied Statistics I September 14 2022

• expected value $\mathbb{E}(y) =$ linear in β

 $\longrightarrow \texttt{Sep142022.Rmd}$

The linear model

- expected value $\mathbb{E}(y) =$ linear in β
- measured with additive error $y = \mathbb{E}(y) + \epsilon$, $\epsilon \sim -$

 $\longrightarrow \texttt{Sep142022.Rmd}$

The linear model

- expected value $\mathbb{E}(y) =$ linear in β
- measured with additive error $y = \mathbb{E}(y) + \epsilon$, $\epsilon \sim -$
- generalizations

 $\epsilon \sim$

 $\longrightarrow \texttt{Sep142022.Rmd}$



- 1. Course introduction: technical issues, course details, evaluation, syllabus, people
- 2. Upcoming events of interest
- 3. Review of linear regression
- 4. In the news: \longrightarrow at the conference