Methods of Applied Statistics I

STA2101H F LEC9101

Week 1

September 15 2021



Ten million reasons to vaccinate the world



- 1. Course introduction: technical issues, course details, evaluation, syllabus, people
- 2. Upcoming events of interest
- 3. Review of linear regression
- 4. In the news: excess deaths
- 5. Computing: RStudio, RMarkdown

Technical Issues

- If **you** are having technical difficulties
 - If possible, send **me** a message in chat
 - Try leaving the class and re-joining
 - Try switching to Chrome if you are using something else
 - Don't panic, the lecture is being recorded and both the recording and the slides will be posted
- If **Prof** is having technical difficulties
 - Check the chat to see if there's any information there
 - If I've disappeared completely, give me 15 minutes before closing the call
 - Look for an announcement on Quercus
 - Don't panic, Prof, you'll figure it out

Applied Statistics I

STA 2101F: Methods of Applied Statistics I Wednesday, 10am – 1 pm Eastern September 15 – December 8 2021

Updated September 14

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 15 2021 Course Delivery:

On Contembor 15 and 22, the class will be delivered online at the scheduled time

Applied Statistics I

STA 2101F: Methods of Applied Statistics I Wednesday, 10am – 1 pm Eastern September 15 – December 8 2021

Updated September 14

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 15 2021 Course Delivery:

On Contembor 15 and 22 the class will be delivered online at the scheduled time.

Course Description

• Course Delivery

Piazza, Notifications

- Grading
- Academic Integrity
- Computing



- References
 Modules
- Contact

Use Piazza for course questions; email for personal questions

Applied Statistics I September 15 2021

- $\bullet \text{ about me} \longrightarrow$
- TA: Ruoyong Xu

- Please turn on your camera to introduce yourself
- Name, program, current location (city)





- 1. Course introduction: technical issues, course details, evaluation, syllabus, people
- 2. Upcoming events of interest
- 3. Review of linear regression
- 4. Steps in analysis
- 5. In the news: excess deaths

Upcoming events 1

About Lucy Gao



Lucy is an Assistant Professor in the Department of Statistics and Actuarial Science at the University of Waterloo. She received her PhD in Biostatistics from the University of Washington. Her research interests are in statistical learning, selective inference, and experiment design.

- Weekly Department Seminar Series
- Selective Inference on Trees
- via Zoom

Sep 16 15.30 EDT

Register here

 $|-\zeta$



- Launch of Data Sciences Institute, U of T
- Speakers:
 - Jennifer Chayes, UC Berkeley
 - Andrew Gelman, Columbia
 - Rob Tibshirani, Stanford



(י)

Upcoming events 3

Distinguished Lecture Series in Statistical Sciences

Register here



- September 20, 2021, 4-5pm Eastern Bayesian Modelling and Analysis of Challenging Data: Making New Sources of Data Trustworthy
- September 21, 2021, 4-5pm Eastern Bayesian Modelling and Analysis of Challenging Data: Identifying the Intrinsic Dimension of High-Dimensional Data



- 1. Course introduction: technical issues, people, course details, evaluation, syllabus
- 2. Upcoming events of interest
- 3. Review of linear regression
- 4. Steps in analysis
- 5. In the news: excess deaths

• Model:

• Model:

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

Applied Statistics I September 15 2021

• Model:

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

• Equivalently:

$$y_i = \chi_i^{\dagger} \not \xi + \varepsilon_i \qquad i = 1, ..., n$$



y is often called response

x_i often called explanatory variables

why?

constant

• Model:

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

• Equivalently:

$$y_i = \chi_i \beta f z_i$$
 $i = 1, ..., n$

- 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:

$$\mathbb{E}(Y \mid X) = X\beta , \quad \text{var}(Y \mid X) = \sigma^2 \mathcal{I}_{y} \in \text{id. under}$$

$$I ??$$
September 15 2021

Applied Statistics I

Nice big equation:

$$\begin{pmatrix} y_{1} \\ \vdots \\ \vdots \\ y_{n} \end{pmatrix} = \begin{pmatrix} x_{11} & \dots & x_{1p} \\ \vdots & \vdots & \vdots \\ \vdots & \vdots & \vdots \\ x_{n1} & \dots & x_{np} \end{pmatrix} \begin{pmatrix} \beta_{i} \\ \vdots \\ \beta_{p} \end{pmatrix} + \begin{pmatrix} \xi_{1} \\ \vdots \\ \xi_{n} \\ \xi_{n} \end{pmatrix} \qquad \underbrace{Y = \chi \beta \neq \xi}_{\xi_{n}} \xi$$

Nice big equation:

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

Or, if you prefer:

refer:

$$y_{i} = x_{i1}\beta_{1} + x_{i2}\beta_{2} + \dots + x_{ip}\beta_{p} + \epsilon_{i}, \quad \epsilon_{i} = \frac{2}{5} \epsilon_{i} + \epsilon_{i}, \quad i = 1, \dots, n$$

$$f = 1, \dots, n$$

Nice big equation:

$$\begin{pmatrix} y_1 \\ \vdots \\ \vdots \\ y_n \end{pmatrix} = \begin{pmatrix} x_{11} & \cdots & x_{1p} \\ \vdots & \vdots & \vdots \\ \vdots & \vdots & \vdots \\ x_{n1} & \cdots & x_{np} \end{pmatrix} \begin{pmatrix} \beta_1 \\ \vdots \\ \beta_p \end{pmatrix} + \begin{pmatrix} \varepsilon_1 \\ \vdots \\ \vdots \\ \varepsilon_n \end{pmatrix}$$

Or, if you prefer:

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

Or, if you prefer:

$$\mathbb{E}(y_i \mid x_i) = x_i^{\mathrm{T}} \beta,$$

Systematic random $i = 1, \dots, n$
 $y_i \text{ independent}$

Applied Statistics I

September 15 2021

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
- Almost always in regression settings we condition on X, as on previous slide

ancillary statistic if its dist= down't dyp. on Box 5

e.g.?



September 15 2021

- 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
 - often not emphasized: interpretation of $\beta_j \begin{pmatrix} y & s \\ y & z \end{pmatrix} \neq z \end{pmatrix} \neq z \end{pmatrix}$ ancillary statistic
 - version 1: effect on the expected response of a unit change in *j*th explanatory variable, $M_{1} = \frac{1}{2} \sum_{i=1}^{n} \frac{$

$$y = \begin{cases} \chi_{i} + f_{i} \chi_{j} + f_{i} \chi_{j} + f_{i} \chi_{j} + f_{i} \chi_{j} \\ y = f_{i} \chi_{i} + f_{i} \chi_{j} + f_{i} \chi_{j$$

* correct" -> ??

• Definition





• Definition

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

• Equivalently,

Definition

- Equivalently,
- Equivalently,

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

$$\hat{\beta}_{LS} := \min(Y - X \beta) (Y - X \beta) (Y - Y \beta)$$

$$\hat{\beta}_{LS} := 12 \text{ distance}$$

unstante

• Definition

- Equivalently,
- Equivalently,



L2 distance

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

?how?

Definition

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

 $\widehat{\beta}_{LS} = (X^T \times) \hat{X}^T y$

- Equivalently,
- Equivalently,

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

L2 distance

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

$$X^{\mathrm{T}}(y - X\beta) = 0 \quad \longrightarrow \quad X^{\mathsf{T}}y = X^{\mathsf{T}}X \beta_{\mathsf{LS}}$$

?how?

15

If (XTX) - exists

check dimensions

Applied Statistics I September 15 2021

Solution

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

check dimensions

ASIDE: here and following all assume X is fixed

Applied Statistics I September 15 2021

Solution

• Expected value

$$\hat{\beta}_{LS} = (X^{T}X)^{-1}(X^{T}Y)$$
check dimensions
$$\mathbb{E}(\hat{\beta}_{LS}) = \mathbb{E}\{\{Y, Y\} = (X^{T}X)^{T}X^{T}EY = (X^{T}X)^{T}X^{T}X \neq \{Y\} = \{X^{T}X\}^{T}X^{T}X \neq \{Y\} = \{X^{T}X\}^{T}X \neq \{Y\} =$$

ASIDE: here and following all assume X is fixed

Solution

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

check dimensions

• Expected value

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

why?

· Least squares estimates are unbiased

ASIDE: here and following all assume X is fixed

Applied Statistics I September 15 2021

• Expected value

• Solution

$$\hat{\beta}_{LS} = (X^{T}X)^{-1}(X^{T}y) \qquad \text{var}(A^{T}y) \qquad \text{$$

· Least squares estimates are unbiased

 \mathbb{E}

• Variance

really variance-covariance matrix

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

ASIDE: here and following all assume X is fixed

What about the normal distribution?

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

$$y_i \sim \mathcal{N}(x_i^T \beta, \sigma^2)$$

$$\hat{\beta}_{is} = (X^T x)^T x^T y = Ay \sim \mathcal{N}(\beta, \sigma^2 (x x)^T)$$
- 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

 $L(\beta,\sigma^{2};y) = \frac{1}{(2\pi\sigma^{2})^{n/2}} \exp\left\{-\frac{1}{(2\sigma^{2})^{(2}}(y-X\beta)^{T}(y-X\beta)\right\}, \quad (y-y)^{T} \mathcal{E}^{-1}(y+y)$ $\begin{pmatrix} \sigma & 0 \\ 0 & a \end{pmatrix} \sum_{i=1}^{n} \sum_{j=1}^{n} \sum_{j=1}^{n} \sum_{i=1}^{n} \sum_{j=1}^{n} \sum_{j=1}^{n} \sum_{i=1}^{n} \sum_{i=1}^{n} \sum_{i=1}^{n} \sum_{i=1}^{n} \sum_{i=1}^{n} \sum_{j=1}^{n} \sum_{i=1}^{n} \sum_{i=1}^{n} \sum_{i=1}^{n$ $\begin{pmatrix} y_{\sigma^{\perp}} \\ \vdots \\ y_{\sigma^{\alpha}} \end{pmatrix} \overline{\Sigma}^{r'}$

September 15 2021

52I

= losk it my

 $N_{\mu}(\mu, \Sigma) f(y; \mu, \Sigma)$

- 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

$$L(\beta, \sigma^{2}; y) = \frac{1}{(2\pi\sigma^{2})^{n/2}} \exp\left\{-\frac{1}{2\sigma^{2}}(y - X\beta)^{T}(y - X\beta)\right\},$$
• the log-likelihood function is
$$\ell(\beta, \sigma^{2}; y) = -\frac{n}{2}\log(\sigma^{2}) - \frac{1}{2\sigma^{2}}(y - X\beta)^{T}(y - X\beta),$$
• the log-likelihood function is
$$\ell(\beta, \sigma^{2}; y) = -\frac{n}{2}\log(\sigma^{2}) - \frac{1}{2\sigma^{2}}(y - X\beta)^{T}(y - X\beta),$$
• the log-likelihood function is
$$\ell(\beta, \sigma^{2}; y) = -\frac{n}{2}\log(\sigma^{2}) - \frac{1}{2\sigma^{2}}(y - X\beta)^{T}(y - X\beta),$$
• the log-likelihood function is
$$\ell(\beta, \sigma^{2}; y) = -\frac{n}{2}\log(\sigma^{2}) - \frac{1}{2\sigma^{2}}(y - X\beta)^{T}(y - X\beta),$$
• the log-likelihood function is
$$\ell(\beta, \sigma^{2}; y) = -\frac{n}{2}\log(\sigma^{2}) - \frac{1}{2\sigma^{2}}(y - X\beta)^{T}(y - X\beta),$$
• the log-likelihood function is
$$\ell(\beta, \sigma^{2}; y) = -\frac{n}{2}\log(\sigma^{2}) - \frac{1}{2\sigma^{2}}(y - X\beta)^{T}(y - X\beta),$$
• the log-likelihood function is
$$\ell(\beta, \sigma^{2}; y) = -\frac{n}{2}\log(\sigma^{2}) - \frac{1}{2\sigma^{2}}(y - X\beta)^{T}(y - X\beta),$$
• the log-likelihood function is
$$\ell(\beta, \sigma^{2}; y) = -\frac{n}{2}\log(\sigma^{2}) - \frac{1}{2\sigma^{2}}(y - X\beta)^{T}(y - X\beta),$$
• the log-likelihood function is
$$\ell(\beta, \sigma^{2}; y) = -\frac{n}{2}\log(\sigma^{2}) - \frac{1}{2\sigma^{2}}(y - X\beta)^{T}(y - X\beta),$$
• the log-likelihood function is
$$\ell(\beta, \sigma^{2}; y) = -\frac{n}{2}\log(\sigma^{2}) - \frac{1}{2\sigma^{2}}(y - X\beta)^{T}(y - X\beta),$$
• the log-likelihood function is
$$\ell(\beta, \sigma^{2}; y) = -\frac{n}{2}\log(\sigma^{2}) - \frac{1}{2\sigma^{2}}(y - X\beta)^{T}(y - X\beta),$$
• the log-likelihood function is
$$\ell(\beta, \sigma^{2}; y) = -\frac{n}{2}\log(\sigma^{2}) - \frac{1}{2\sigma^{2}}(y - X\beta)^{T}(y - X\beta),$$
• the log-likelihood function is
$$\ell(\beta, \sigma^{2}; y) = -\frac{n}{2}\log(\sigma^{2}) - \frac{1}{2\sigma^{2}}(y - X\beta)^{T}(y - X\beta),$$
• the log-likelihood function is
$$\ell(\beta, \sigma^{2}; y) = -\frac{n}{2}\log(\sigma^{2}) - \frac{n}{2}\log(\sigma^{2}) - \frac{n}{2}$$

17

- 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

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

• 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

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

Applied Statistics I September 15 2021

• 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

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

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

why?

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

- maximum likelihood estimate of β is

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

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

• distribution of \hat{eta}_j is

why?

0

maximum likelihood estimate of β is

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

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

whv?

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

- distribution of $\hat{\beta}_i$ is
- maximum likelihood estimate of σ^2 is $\frac{1}{n}(y X\hat{\beta})^{\mathrm{T}}(y X\hat{\beta})$ $\hat{\beta}_j \pm 2\left[\sigma\left\{\left(X^{\mathrm{T}}X\right)^{-1}, j\right\}^{1/2} + 25\%\right] C T$ Statistics 1 September 15 2021

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

Applied

- maximum likelihood estimate of β is

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

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

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

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

$$(X^{T}X)^{-1}(\beta, \sigma^{2}(X^{T}X)^{-1})$$

$$(X^{T}X)^{-1}(\beta, \sigma^{2}(X^{T}X)^{-1})$$

18

• distribution of
$$\beta_j$$
 is

$$N(\beta_j, \sigma^2(X^TX)_{jj}^{-1}), \quad j = 1, ..., p \qquad usually used formula is a set of σ^2 is $\frac{1}{n}(y - X\hat{\beta})^T(y - X\hat{\beta})$ and $\sigma^2 = \frac{1}{n-p}(y - X\hat{\beta})^T(y - X\hat{\beta})$ for $\sigma^2 = \sigma^2$ and γ_j
• but we use
 $p < n$
 $j_n = (\int_{a_1}^{b_1} \int_{a_2}^{\sigma^2} \frac{1}{n-p}(y - X\hat{\beta})^T(y - X\hat{\beta}) \sum_{i=1}^{n-p} f_{i}(y - X\hat{\beta})^T(y - X\hat{\beta}) \sum_{i=1}^{n-p} f_{i}(y - X\hat{\beta}) \sum_{i=1}^$$$

Pause

 $\hat{\beta}_{i} \sim N(\beta_{i}, se_{j})$ ~ 95% C.I J (ander Nors=) B. ± 1.96. se; (1) I'm lost $\sqrt{var\beta_j} = \sqrt{\sigma^2 (X^T X)^2}$ (2) I'm good (3) I'm bored $\hat{\beta}_{j} \pm \pm \pm \frac{x^{2}}{n-1} = \sqrt{\hat{\sigma}^{2}} (x^{T}x)^{T}$ 22 but > 1.96

HW Question Week 1

STA2101F 2021

Due September 22 2021 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).

Applied Statistics: install.packages("if2007"). The questions below are adapted from LM Ch.3.

(a) Fit the linear model regist x = x1 + x2 + x3 + x4 Extract the V matrix using the



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

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

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

Inference

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

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

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

• but just using previous slide we have

$$\frac{\hat{\beta}_j - \beta_j}{\sigma[\{(X^{\mathrm{\scriptscriptstyle T}}X)^{-1}\}_{jj}\}]^{1/2}} \sim N(\mathsf{O},\mathsf{1})$$

Inference

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

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

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

• but just using previous slide we have

See also Sep152021.Rmd

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

See also Sep152021.Rmd

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

Example

LM Exercise 2.4

summary(mode	11) 🔨	\checkmark			
Coefficients	: (Y)	Ser			
	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 .	t a co
lbph	0.107054	0.058449	1.832	0.07040 .	97-9 = 88
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	at set
pgg45	0.004525	0.004421	1.024	0.30886	Pj - ser c
Q: : C 1	0 (1 0 001 (1 0 01 0		0 1 () 1

Applied Statistics I September 15 2021

Signif. codes: $0 \ '***'$ $0.001 \ '**'$ $0.01 \ '*'$ $0.05 \ '.'$ $0.1 \ ' \ '1$ ied Statistics ISeptember 15 2021Confint (?)

			322 views
-	Women in Statistics and Data Science		
- Pr	Follow @WomeninStat 18.7K followers		1 🔿 🎽
Jul 23rd 2020), 14 tweets, 4 min read		
		🗍 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.

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





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

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



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

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

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



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

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

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



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

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

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

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

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



SM Example 8.5

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

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

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



for

1st column of *X*?

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

• $y_i = \gamma_0 x_{1i}^{\gamma_1} x_{2i}^{\gamma_2} \eta_{i}, \quad \eta_i \sim \text{positive r.v.} \quad \text{ff bog } \gamma_i \sim b_f \gamma_0 + \vartheta_i b_f \gamma_i; \quad \text{SM Example 8.5}$
• $y_i = \varphi_0 + \sum_{k=1}^{K} \varphi_k s_k(x_i) + \epsilon_i \quad K=3 \quad \text{Horn:} \quad \text{Smoothing splines, e.g.}$
Applied Statistics I September 15 2021 $\gamma \sim S(x_i, 3)$ $\gamma \sim S(x_i, 3)$

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

 $\longrightarrow \texttt{Sep152021.Rmd}$

Applied Statistics I September 15 2021

The linear model

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

 $\longrightarrow \texttt{Sep152021.Rmd}$

Applied Statistics I September 15 2021

The linear model

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

 $\epsilon \sim$

 $\longrightarrow \texttt{Sep152021.Rmd}$



- 1. Course introduction: technical issues, course details, evaluation, syllabus, people
- 2. Upcoming events of interest
- 3. Review of linear regression



5. In the news: excess deaths



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

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

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

- 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

- 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

• visualization strategies, conveyance of uncertainties

In the news



Applied Statistics I Septemb

September 15 2021

... in the news


Economist updates

value lies between 9.5m and 18.6m additional deaths.





Global statistical modelling done by The Economist estimates that the true number of those who died in Kenya as a result of the covid-19 pandemic is between 19,000 and 110,000, versus an official death toll of 4,746.



Applied Statistics I

September 15 2021 Seconomist.com Why the Economist's excess death model is misleading . Gordon Shotwell

Why the Economist's excess death model is misleading

Sep 7, 2021
10 min read

9/9/2021

The Economist has published a model which estimates that Kenyans are only detecting 4-25% of the true deaths which can be attributed to Covid. I think this is a good opportunity to learn about why many machine learning models are problematic. I'm going to talk about this particular model, but I should note that I've only spent about ten hours looking at this problem and I'm sure the authors of this model are smart thoughtful people who don't mean to mislead. That said, I think it's an excellent example of how machine learning models can lend a sheen of credibility to things that are basically unsupported assertions. When someone says that their model says something, most people assume that means that it's supporting that thing with hard data when it's often just making unsupported assertions. It's possible that the authors of this model have sound reasons about why they can make global excess death predictions based on a small unrepresentative sample of countries, but even so I think these observations are helpful for figuring out which models you should trust.

What got me started thinking about this subject was this tweet by one of the writers at The Economist suggesting that Kenya was radically undercounting deaths which have resulted from the Covid-19 pandemic.

Adrian Blomfield 🤣 @adrianblomfield



Applied Statistics I

September 15 2021 Global statistical modelling done by The Economist estimates that the true number of those who died in Kenva



Applied Statistics I S



Source: Johns Hopkins University CSSE COVID-19 Data – Last updated 14 September, 16:04 (London time) OurWorldInData.org/coronavirus • CC BY

Applied Statistics I September 15 2021