

Math 6040/7260 Linear Models

Mon/Wed/Fri 9:00am - 9:50am

Instructor: Dr. Xiang Ji, xji4@tulane.edu

1 Lecture 1: Jan 18

Today

- Introduction
- Introduce yourself
- Course logistics

What is this course about?

The term “linear models” describes a wide class of methods for the statistical analysis of multivariate data. The underlying theory is grounded in linear algebra and multivariate statistics, but applications range from biological research to public policy. The objective of this course is to provide a solid introduction to both the theory and practice of linear models, combining mathematical concepts with realistic examples.

Prerequisite

- **Must:** Introduction to Probability (Math 3070/6070), Mathematical Statistics (Math 3080/6080)
- **Good to have:** Scientific Computation II

A hierarchy of linear models

- The linear mean model:

$$\underset{n \times 1}{\mathbf{y}} = \underset{n \times p}{\mathbf{X}} \underset{p \times 1}{\boldsymbol{\beta}} + \underset{n \times 1}{\boldsymbol{\epsilon}}$$

where $\mathbf{E}(\boldsymbol{\epsilon}) = \mathbf{0}$. Only assumption is that errors have mean 0.

- Gauss-Markov model:

$$\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\epsilon}$$

where $\mathbf{E}(\boldsymbol{\epsilon}) = \mathbf{0}$ and $\mathbf{Var}(\boldsymbol{\epsilon}) = \sigma^2 \mathbf{I}$. Uncorrelated errors with constant variance.

- Aitken model or general linear model:

$$\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\epsilon}$$

where $\mathbf{E}(\boldsymbol{\epsilon}) = \mathbf{0}$ and $\mathbf{Var}(\boldsymbol{\epsilon}) = \sigma^2 \mathbf{V}$. \mathbf{V} is fixed and known.

- Variance components models: $\mathbf{y} \sim N(\mathbf{X}\beta, \sigma_1^2\mathbf{V}_1 + \sigma_2^2\mathbf{V}_2 + \dots + \sigma_r^2\mathbf{V}_r)$ with $\mathbf{V}_1, \mathbf{V}_2, \dots, \mathbf{V}_r$ known.
- General mixed linear Model:

$$\mathbf{y} = \mathbf{X}\beta + \epsilon$$

where $\mathbf{E}(\epsilon) = \mathbf{0}$ and $\mathbf{Var}(\epsilon) = \Sigma(\theta)$.

- Generalized linear models (GLMs). Logistic regression, probit regression, log-linear model (Poisson regression), ... Note the difference from the general linear model. GLMs are generalization of the *concept* of linear models. They are covered in Math 7360 - Data Analysis class (<https://tulane-math-7360-2021.github.io/>).

Syllabus

Check course website frequently for updates and announcements.

<https://tulane-math-7260-2023.github.io/>

HW submission

Through Github with demo on Friday class.

Presentations

Let me know your pick by the end of Friday (01/20/2023).

Last year (spring 2022) comments

1. Experience in this course

- I really enjoyed this class.
- The lecture is dry. Expect Dr. Ji to read through the lecture notes until someone asks a question. Also, having a background in statistics gives context to the lecture notes.

Response: I did emphasize on pre-requisite last year, but it still seemed not enough.

- I think this course was well laid out. Even though I did take it as a graduation requirement, I ended up enjoying the course. I also think the work load is manageable and professor Xi does provide all the tools necessary to succeed in the course. It can be overwhelming at first but with a little time and effort you can get the hang of the material.

Response: Mind the typos, students.

2. Strong aspects of this course

- Professor Ji is a great professor.

- Dr. Ji provides an inhuman amount of course material to help supplement learning, it was extremely helpful to have labs (answered and unanswered) as well as homework keys posted.

Response: There will be labs again.

- Notes are very structured and the professor is nice.
- The strongest aspect of this course is the homework. They are a great way to interact with and learn the material. The problems can seem challenging, but are doable with a little effort. Also, he publishes the keys afterword so you can check your work and see what the reasoning behind the answer is. This coupled with the lab sessions are a great way to prepare for the exams.

3. There will be an internal mid-term-ish evaluation for this course. Will remember to go over them.

Last last year (spring 2021) comments

1. Experience in this course

- Overall, I had a pretty good experience in this course. It moved quickly, but that is expected from this level of course. Sometimes it was hard to stay engaged with the lectures and to really absorb the course material. Because the lectures moved so fast, I really appreciated how Professor made the full notes available at the time of the lecture. I would have liked if there were a few more examples with the notes, as sometimes the homework felt disjoint from the notes.

Response: I will try to move slower this semester. I will start lab sessions earlier too.

- The professor is an extremely intelligent, kind, and understanding professor. He prioritizes in making sure that we understand the material and seeing how the material can be applied. His lecture notes were a godsend because the texts could be a bit ambiguous at times but he elucidated the material in such a comprehensible manner.

Response: I will try to fix the left-over typos.

- Mentioned in class from other students/internal evaluation, conveying the mathematical concepts through the presentation is not a good idea to follow the class in real-time. Prepared presentation can give rise to a distraction on what we have been going over.

Response: I am still delivering this class in hybrid-mode. I found the presentations fit online teaching better. I think the difficulty might be caused by (1) fast moving lecture (2) I only realized the need of reviewing basic concepts of probability almost a quarter into the semester...

- I found the setup of the course not very engaging. Additionally, many of the class notes came directly from the additional sources with no additional information or

explanation, which I found to be not very helpful.

Response: I actually like them. I was the guinea pig to test them.

- Easily help us to understand the main course, and the notes and details are great.

Response: There will be notes.

- Moves very quickly and can be hard to keep up with. Sometimes instructions are unclear.

Response: I will try to slow down.

- Both the instructor and the TA were helpful. It was hard to follow along in class though.

Response: We don't have TA this time. Make use of the office hour. And I have to say, it needs effort to ace in this class.

2. Strong aspects of this course

- Having the lecture notes and labs available was very helpful. Professor was also always very nice and accommodating, and willing to meet with me when I needed help. He also always responded to student feedback, if we asked for an extra day or two on the homework or something like that.

Response: Here is an example of correctly using the office hours.

- His lecture notes and the lab sessions.

Response: They will be there again.

- Lab session is necessarily required to this class. A lot of computations in the class would be done by computer due to the complexity, and students are expected to handle with the computer programming properly at a desired level. The course can be an introduction to the statistical computation, which does not exist in the mathematics department.

Response: Hmm, there is a course Math 7360 Data Analysis that focuses more on the computational side.

- I appreciated the homework reviews in class and felt these helped clarify the material.

Response: Of course, the reviews will be there again. The purpose of the course is for you to learn.

- Grading was easy which made up for the rigor.

Response: Don't rely on this...

- Really appreciate that Professor Xiang made such a neat and tidy notes for us. It is really helpful for me to review. And notes have a great interaction with us, Professor Xiang also leaves some questions to help us think about the logic behind.

Response: Well, Xiang is my first name. Please call me Prof. X.

- Prof. Xiang was highly organized and wanted his students to understand the course content more than he made them worry about grades. I learned a lot

about Linear Models and feel confident applying the course content professionally and academically. I wish most of the Math department had his teaching style and implemented his course documents and organization structure. Prof. Xiang made the course content in class digestible and if I needed to review the material I could easily find it through his course notes and textbook. I wish I could say the same about my other courses.

Response: Hmm, I like Prof. X. better.

- I really appreciated the emphasis on learning. It allowed for most students to take it at the pace that was good for them.

Response: Please don't let your score rely on this comment.

3. There will be an internal mid-term-ish evaluation for this course. Will remember to go over them.

2 Lecture 2:Jan 20

Last time

- Introduction
- Course logistics

Today

- Reply to the “Presentation Dates” thread on Canvas by the end of today.
- Git

What is git?

Git is currently the most popular system for version control according to [Google Trend](#). Git was initially designed and developed by [Linus Torvalds](#) in 2005 for Linux kernel development. Git is the British English slang for unpleasant person.

Why using git?

- [GitHub](#) is becoming a de facto central repository for open source development.
- **Advertise** yourself through GitHub (e.g., host a free personal webpage on GitHub).
- a skill that employers look for (according to [this AmStat article](#)).

Git workflow

Figure 2.1 shows its basic workflow.

What do I need to use Git?

- A **Git server** enabling multi-person collaboration through a centralized repository.
- A **Git client** on your own machine.
 - Linux: Git client program is shipped with many Linux distributions, e.g., Ubuntu and CentOS. If not, install using a package manager, e.g., `yum install git` on CentOS.
 - Mac: follow instructions at <https://www.atlassian.com/git/tutorials/install-git>.
 - Windows: Git for Windows at <https://gitforwindows.org> (GUI) aka Git Bash.
- Do **not** totally rely on GUI or IDE. Learn to use Git on command line, which is needed for cluster and cloud computing.

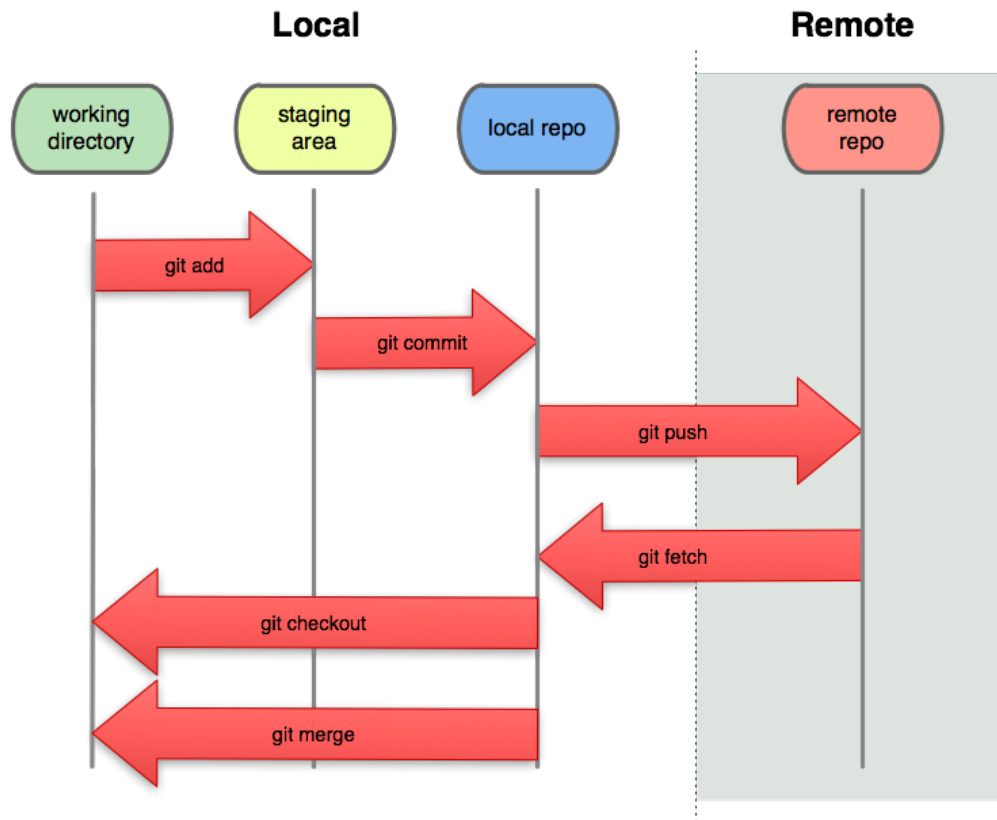


Figure 2.1

Git survival commands

- `git pull` synchronize local Git directory with remote repository.
- Modify files in local working directory.
- `git add FILES` add snapshots to staging area
- `git commit -m "message"` store snapshots permanently to (**local**) Git repository
- `git push` push commits to remote repository.

Git basic usage

Working with your local copy.

- `git pull` : update local Git repository with remote repository (fetch + merge).
- `git log FILENAME` : display the current status of working directory.
- `git diff` : show differences (by default difference from the most recent commit).
- `git add file1 file2 ...` : add file(s) to the staging area.

- `git commit` : commit changes in staging area to Git directory.
- `git push` : publish commits in local Git repository to remote repository.
- `git reset --soft HEAD 1` : undo the last commit.
- `git checkout FILENAME` : go back to the last commit, discarding all changes made.
- `git rm FILENAME` : remove files from git control.

Git demonstration

Show how to create a private git repository for HW and Exam submissions.

On [GitHub](#)

- Obtain [student developer pack](#).
- Create a private repository `math-6040-2023-spring` (please substitute 6040 by 7260 if you are taking the graduate level). Add `xji3` as your collaborators with write permission ([instruction](#)).

On your local machine:

- clone the repository: please refer to [this webpage](#) with instructions for your operating system.
- enter the folder: `cd math-6040-2023-spring`.
- after finishing the rest of the questions, save your file inside your git repository folder `math-6040-2023-spring` with name `hw1.pdf` (for example). Please make it human-readable.
- now using git commands to stage this change: `git add hw1.pdf`
- commit: `git commit -m "hw1 submission"` (remember to replace the quotation mark)
- push to remote server: `git push`
- tag version hw1: `git tag hw1` and push: `git push --tags`.

Take a look at the tags on GitHub ([instructions](#)).

When submitting your hw, please email your instructor (xji4@tulane.edu) a link to your tag ([instructions](#)).

3 Lecture 3:Jan 23

Last time

- Git

Today

- HW1 posted
- Linear algebra: vector and vector space, rank of a matrix
- Column space and Nullspace (JM Appendix A)

Notations

$$\underset{n \times 1}{\mathbf{y}} = \underset{n \times p}{\mathbf{X}} \underset{p \times 1}{\boldsymbol{\beta}} + \underset{n \times 1}{\boldsymbol{\epsilon}}$$

$$\begin{pmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{pmatrix} = \begin{pmatrix} x_{11} & x_{12} & \dots & x_{1p} \\ x_{21} & x_{22} & \dots & x_{2p} \\ \vdots & \vdots & & \vdots \\ x_{n1} & x_{n2} & \dots & x_{np} \end{pmatrix} \begin{pmatrix} \beta_1 \\ \beta_2 \\ \vdots \\ \beta_p \end{pmatrix} + \begin{pmatrix} \epsilon_1 \\ \epsilon_2 \\ \vdots \\ \epsilon_n \end{pmatrix}$$

- All vectors are column vector
- Write dimensions underneath as in $\underset{n \times p}{\mathbf{X}}$ or as $\mathbf{X} \in \mathbb{R}^{n \times p}$
- Bold upper-case letters for Matrices. Bold lower-case letters for Vectors.

Vector and vector space

(from JM Appendix A)

- A set of vectors $\mathbf{x}_1, \dots, \mathbf{x}_n$ are *linearly dependent* if there exist coefficients c_j for $j = 1, 2, \dots, n$ such that $\sum_{j=1}^n c_j \mathbf{x}_j = \mathbf{0}$ and $\|\mathbf{c}\|_2 = \sum_{j=1}^n c_j^2 > 0$. They are *linearly independent* if $\sum_{j=1}^n c_j \mathbf{x}_j = \mathbf{0}$ implies (i.e. \implies) $c_j = 0$ for all j .
- Two vectors are *orthogonal* to each other, written $\mathbf{x} \perp \mathbf{y}$, if their inner product is 0, that is $\mathbf{x}^T \mathbf{y} = \mathbf{y}^T \mathbf{x} = \sum_j x_j y_j = 0$.
- A set of vectors $\mathbf{x}^{(1)}, \mathbf{x}^{(2)}, \dots, \mathbf{x}^{(n)}$ are mutually orthogonal iff (i.e. \iff) $\mathbf{x}^{(i)T} \mathbf{x}^{(j)} = 0$ for $\forall i \neq j$.
- The most common set of vectors that are mutually orthogonal are the *elementary* vectors $\mathbf{e}^{(1)}, \mathbf{e}^{(2)}, \dots, \mathbf{e}^{(n)}$, which are all zero, except for one element equal to 1, so that $\mathbf{e}_i^{(i)} = 1$ and $\mathbf{e}_j^{(i)} = 0, \forall j \neq i$.

- A *vector space* \mathcal{S} is a set of vectors that are closed under addition and scalar multiplication, that is
 - if $\mathbf{x}^{(1)}$ and $\mathbf{x}^{(2)}$ are in \mathcal{S} , then $c_1\mathbf{x}^{(1)} + c_2\mathbf{x}^{(2)}$ is in \mathcal{S} .
- A vector space \mathcal{S} is *generated* or *spanned* by a set of vectors $\mathbf{x}^{(1)}, \mathbf{x}^{(2)}, \dots, \mathbf{x}^{(n)}$, written as $\mathcal{S} = \text{span}\{\mathbf{x}^{(1)}, \mathbf{x}^{(2)}, \dots, \mathbf{x}^{(n)}\}$, if any vector \mathbf{x} in the vector space is a linear combination of $\mathbf{x}_i, i = 1, 2, \dots, n$.
- A set of linearly independent vectors that generate or span a space \mathcal{S} is called a *basis* of \mathcal{S} .

Example A.1

Let

$$\mathbf{x}^{(1)} = \begin{bmatrix} 1 \\ 1 \\ 1 \\ 1 \end{bmatrix}, \mathbf{x}^{(2)} = \begin{bmatrix} 1 \\ 2 \\ 3 \\ 4 \end{bmatrix}, \text{ and } \mathbf{x}^{(3)} = \begin{bmatrix} -3 \\ -1 \\ 1 \\ 3 \end{bmatrix}.$$

Then $\mathbf{x}^{(1)}$ and $\mathbf{x}^{(2)}$ are linearly independent, but $\mathbf{x}^{(1)}, \mathbf{x}^{(2)}$, and $\mathbf{x}^{(3)}$ are linearly dependent since $5\mathbf{x}^{(1)} - 2\mathbf{x}^{(2)} + \mathbf{x}^{(3)} = \mathbf{0}$

Rank

Some matrix concepts arise from viewing columns or rows of the matrix as vectors. Assume $\mathbf{A} \in \mathbb{R}^{m \times n}$.

- $\text{rank}(\mathbf{A})$ is the maximum number of linearly independent rows or columns of a matrix.
- $\text{rank}(\mathbf{A}) \leq \min\{m, n\}$.
- A matrix is *full rank* if $\text{rank}(\mathbf{A}) = \min\{m, n\}$. It is *full row rank* if $\text{rank}(\mathbf{A}) = m$. It is *full column rank* if $\text{rank}(\mathbf{A}) = n$.
- a square matrix $\mathbf{A} \in \mathbb{R}^{n \times n}$ is *singular* if $\text{rank}(\mathbf{A}) < n$ and *non-singular* if $\text{rank}(\mathbf{A}) = n$.
- $\text{rank}(\mathbf{A}) = \text{rank}(\mathbf{A}^T) = \text{rank}(\mathbf{A}^T \mathbf{A}) = \text{rank}(\mathbf{A} \mathbf{A}^T)$. (Show this in HW.)
- $\text{rank}(\mathbf{A} \mathbf{B}) \leq \min\{\text{rank}(\mathbf{A}), \text{rank}(\mathbf{B})\}$. (Hint: Columns of $\mathbf{A} \mathbf{B}$ are spanned by columns of \mathbf{A} and rows of $\mathbf{A} \mathbf{B}$ are spanned by rows of \mathbf{B} .)
- if $\mathbf{A} \mathbf{x} = \mathbf{0}_m$ for some $\mathbf{x} \neq \mathbf{0}_n$, then $\text{rank}(\mathbf{A}) \leq n - 1$.

Column space

Definition: The column space of a matrix, denoted by $\mathcal{C}(\mathbf{A})$ is the vector space spanned by the columns of the matrix, that is,

$$\mathcal{C}(\mathbf{A}) = \{\mathbf{x} : \text{there exists a vector } \mathbf{c} \text{ such that } \mathbf{x} = \mathbf{A} \mathbf{c}\}.$$

This means that if $\mathbf{x} \in \mathcal{C}(\mathbf{A})$, we can find coefficients c_j such that

$$\mathbf{x} = \sum_j c_j \mathbf{a}^{(j)}$$

where $\mathbf{a}^{(j)} = \mathbf{A}_{\cdot j}$ denotes the j^{th} column of matrix \mathbf{A} .

- The column space of a matrix consists of all vectors formed by multiplying that matrix by any vector.
- The number of basis vectors for $\mathcal{C}(\mathbf{A})$ is then the number of linearly independent columns of the matrix \mathbf{A} , and so, $\dim(\mathcal{C}(\mathbf{A})) = \text{rank}(\mathbf{A})$.
- The dimension of a space is the number of vectors in its basis.

Example A.2

Let $\mathbf{A} = \begin{bmatrix} 1 & 1 & -3 \\ 1 & 2 & -1 \\ 1 & 3 & 1 \\ 1 & 4 & 3 \end{bmatrix}$ and $\mathbf{c} = \begin{bmatrix} 5 \\ 4 \\ 3 \end{bmatrix}$. Show that \mathbf{Ac} is a linear combination of columns in \mathbf{A} .

solution:

$$\mathbf{Ac} = \begin{bmatrix} 1 \times 5 + 1 \times 4 + (-3) \times 3 \\ 1 \times 5 + 2 \times 4 + (-1) \times 3 \\ 1 \times 5 + 3 \times 4 + 1 \times 3 \\ 1 \times 5 + 4 \times 4 + 3 \times 3 \end{bmatrix} = \begin{bmatrix} 0 \\ 10 \\ 20 \\ 30 \end{bmatrix}.$$

You could recognize that

$$\mathbf{Ac} = 5 \times \begin{bmatrix} 1 \\ 1 \\ 1 \\ 1 \end{bmatrix} + 4 \times \begin{bmatrix} 1 \\ 2 \\ 3 \\ 4 \end{bmatrix} + 3 \times \begin{bmatrix} -3 \\ -1 \\ 1 \\ 3 \end{bmatrix} = 5\mathbf{a}^{(1)} + 4\mathbf{a}^{(2)} + 3\mathbf{a}^{(3)} = \begin{bmatrix} 0 \\ 10 \\ 20 \\ 30 \end{bmatrix}.$$

Result A.1

$$\text{rank}(\mathbf{AB}) \leq \min(\text{rank}(\mathbf{A}), \text{rank}(\mathbf{B})).$$

proof: Each column of \mathbf{AB} is a linear combination of columns of \mathbf{A} (i.e. $(\mathbf{AB})_{\cdot j} = \mathbf{A}\mathbf{b}^{(j)}$), so the number of linearly independent columns of \mathbf{AB} cannot be greater than that of \mathbf{A} . Similarly, $\text{rank}(\mathbf{AB}) = \text{rank}(\mathbf{B}^T \mathbf{A}^T)$, the same argument gives $\text{rank}(\mathbf{B}^T)$ as an upper bound.

Result A.2

- (a) If $\mathbf{A} = \mathbf{BC}$, then $\mathcal{C}(\mathbf{A}) \subseteq \mathcal{C}(\mathbf{B})$.
- (b) If $\mathcal{C}(\mathbf{A}) \subseteq \mathcal{C}(\mathbf{B})$, then there exists a matrix \mathbf{C} such that $\mathbf{A} = \mathbf{BC}$.

proof: For (a), any vector $\mathbf{x} \in \mathcal{C}(\mathbf{A})$ can be written as $\mathbf{x} = \mathbf{Ad} = \mathbf{B}(\mathbf{Cd})$.

For (b), $\mathbf{A}_{\cdot j} \in \mathcal{C}(\mathbf{B})$, so that there exists a vector $\mathbf{c}^{(j)}$ such that $\mathbf{A}_{\cdot j} = \mathbf{B}\mathbf{c}^{(j)}$. The matrix $\mathbf{C} = (\mathbf{c}^{(1)}, \mathbf{c}^{(2)}, \dots, \mathbf{c}^{(n)})$ satisfies that $\mathbf{A} = \mathbf{BC}$.

4 Lecture 4: Jan 25

Last time

- Linear algebra: vector and vector space, rank of a matrix
- Column space (JM Appendix A)

Today

- Nullspace
- Probability review

Null space

Definition: The null space of a matrix, denoted by $\mathcal{N}(\mathbf{A})$, is $\mathcal{N}(\mathbf{A}) = \{\mathbf{y} : \mathbf{A}\mathbf{y} = \mathbf{0}\}$.

Result A.3

If \mathbf{A} has full-column rank, then $\mathcal{N}(\mathbf{A}) = \{\mathbf{0}\}$.

proof: Matrix \mathbf{A} has full-column rank means its columns are linearly independent, which means that $\mathbf{A}\mathbf{c} = \mathbf{0}$ implies $\mathbf{c} = \mathbf{0}$.

Theorem A.1

Assume $\mathbf{A} \in \mathbb{R}^{m \times n}$, then $\dim(\mathcal{C}(\mathbf{A})) = r$ and $\dim(\mathcal{N}(\mathbf{A})) = n - r$, where $r = \text{rank}(\mathbf{A})$.

See JM Appendix Theorem A.1 for the proof.

proof: Denote $\dim(\mathcal{N}(\mathbf{A}))$ by k , to be determined, and construct a set of basis vectors for $\mathcal{N}(\mathbf{A}) : \{\mathbf{u}^{(1)}, \mathbf{u}^{(2)}, \dots, \mathbf{u}^{(k)}\}$, so that $\mathbf{A}\mathbf{u}^{(i)} = \mathbf{0}$, for $i = 1, 2, \dots, k$. Now, construct a basis for \mathbb{R}^n by adding the vectors $\{\mathbf{u}^{(k+1)}, \dots, \mathbf{u}^{(n)}\}$, which are not in $\mathcal{N}(\mathbf{A})$. Clearly, $\mathbf{A}\mathbf{u}^{(i)} \in \mathcal{C}(\mathbf{A})$ for $i = k+1, \dots, n$, and so the span of these vectors form a subspace of $\mathcal{C}(\mathbf{A})$. These vectors $\{\mathbf{A}\mathbf{u}^{(i)}, i = k+1, \dots, n\}$ are also linearly independent from the following argument: suppose $\sum_{i=k+1}^n c_i \mathbf{A}\mathbf{u}^{(i)} = \mathbf{0}$; then $\sum_{i=k+1}^n c_i \mathbf{A}\mathbf{u}^{(i)} = \mathbf{A} [\sum_{i=k+1}^n c_i \mathbf{u}^{(i)}] = \mathbf{0}$, and hence $\sum_{i=k+1}^n c_i \mathbf{u}^{(i)}$ is a vector in $\mathcal{N}(\mathbf{A})$. Therefore, there exist b_i such that $\sum_{i=k+1}^n c_i \mathbf{u}^{(i)} = \sum_{i=1}^k b_i \mathbf{u}^{(i)}$, or $\sum_{i=1}^k b_i \mathbf{u}^{(i)} - \sum_{i=k+1}^n c_i \mathbf{u}^{(i)} = \mathbf{0}$. Since $\{\mathbf{u}^{(i)}\}$ form a basis for \mathbb{R}^n , c_i must all be zero. Therefore $\mathbf{A}\mathbf{u}^{(i)}, i = k+1, \dots, n$ are linearly independent. At this point, since $\text{span}\{\mathbf{A}\mathbf{u}^{(k+1)}, \dots, \mathbf{A}\mathbf{u}^{(n)}\} \subseteq \mathcal{C}(\mathbf{A})$, $\dim(\mathcal{C}(\mathbf{A}))$ is at least $n - k$. Suppose there is a vector \mathbf{y} that is in $\mathcal{C}(\mathbf{A})$, but not in the span; then there exists $\mathbf{u}^{(n+1)}$ so that $\mathbf{y} = \mathbf{A}\mathbf{u}^{(n+1)}$ and $\mathbf{u}^{(n+1)}$ is linearly independent of $\{\mathbf{u}^{(k+1)}, \dots, \mathbf{u}^{(n)}\}$ (and clearly not in $\mathcal{N}(\mathbf{A})$), making $n+1$ linearly independent vectors in \mathbb{R}^n . Since that is not possible, the span is equal to $\mathcal{C}(\mathbf{A})$ and $\dim(\mathcal{C}(\mathbf{A})) = n - k = r = \text{rank}(\mathbf{A})$, so that $k = \dim(\mathcal{N}(\mathbf{A})) = n - r$.

Interpretation: “dimension of column space + dimension of null space = # columns”

Mis-Interpretation: Columns space and null space are orthogonal complement to each other.

They are of different orders in general! Next result gives the correct statement.

Reference:

- Statistical Inference, 2nd Edition, by George Casella & Roger L. Berger
- [Review of Probability Theory](#) by Arian Maleki and Tom Do

Probability theory review

A few basic elements to define a probability on a set:

- **Sample space** S is the set that contains all possible outcomes of a particular experiment.
- An **event** is any collection of possible outcomes of an experiment, that is, any subset of S (including S itself).
- Event operations
 1. Union: The union of A and B , written $A \cup B$, is the set of elements that belong to either A or B or both:

$$A \cup B = \{x : x \in A \text{ or } x \in B\}$$

2. Intersection: The intersection of A and B , written $A \cap B$, is the set of elements that belong to both A and B :

$$A \cap B = \{x : x \in A \text{ and } x \in B\}$$

3. Complementation: The complement of A , written as A^c , is the set of all elements that are not in A :

$$A^c = \{x : x \notin A\}.$$

- **Sigma algebra (or Borel field):** A collection of subsets of S is called a sigma algebra (or Borel field), denoted by \mathcal{B} , if it satisfies the following three properties:
 1. $\emptyset \in \mathcal{B}$ (the empty set is an element of \mathcal{B})
 2. If $A \in \mathcal{B}$, then $A^c \in \mathcal{B}$ (\mathcal{B} is closed under complementation).
 3. If $A_1, A_2, \dots \in \mathcal{B}$, then $\cup_{i=1}^{\infty} A_i \in \mathcal{B}$ (\mathcal{B} is closed under countable unions).
- **Axioms of probability:** Given a sample space S and an associated sigma algebra \mathcal{B} , a *probability function* is a function $\Pr()$ with domain \mathcal{B} that satisfies
 1. $\Pr(A) \geq 0$ for all $A \in \mathcal{B}$
 2. $\Pr(S) = 1$.
 3. If $A_1, A_2, \dots \in \mathcal{B}$ are pairwise disjoint, then $\Pr(\cup_{i=1}^{\infty} A_i) = \sum_{i=1}^{\infty} \Pr(A_i)$.

Properties:

If $\Pr()$ is a *probability function* and A and B are any sets in \mathcal{B} , then

- $\Pr(\emptyset) = 0$, where \emptyset is the empty set
Proof: $1 = \Pr(S) = \Pr(S \cup \emptyset)$
- $\Pr(A) \leq 1$
Proof: see below and remember $\Pr(A^c) \geq 0$
- $\Pr(A^c) = 1 - \Pr(A)$
Proof: $1 = \Pr(S) = \Pr(A \cup A^c) = \Pr(A) + \Pr(A^c)$
- $\Pr(B \cap A^c) = \Pr(B) - \Pr(A \cap B)$
Proof: $B = \{B \cap A\} \cup \{B \cap A^c\}$
- $\Pr(A \cup B) = \Pr(A) + \Pr(B) - \Pr(A \cap B)$
Proof: $A \cup B = A \cup \{B \cap A^c\}$ and use the above property.
- $\Pr(A \cup B) = \Pr(A) + \Pr(B \cap A^c) = \Pr(A) + \Pr(B) - \Pr(A \cap B)$
- If $A \subset B$, then $\Pr(A) \leq \Pr(B)$.
Proof: If $A \subset B$, then $A \cap B = A$ and use $\Pr(B \cap A^c) = \Pr(B) - \Pr(A \cap B)$.

Conditional probability

Definition: If A and B are events in S , and $\Pr(B) > 0$, then the conditional probability of A given B , written $\Pr(A|B)$, is

$$\Pr(A|B) = \frac{\Pr(A \cap B)}{\Pr(B)}$$

Note that what happens in the conditional probability calculation is that B becomes the sample space: $\Pr(B|B) = 1$, in other words, $\Pr(A|B)$ is the probability measure of the event A after observing the occurrence of event B .

Definition: Two events A and B are statistically independent if $\Pr(A \cap B) = \Pr(A) \Pr(B)$. When A and B are independent events, then $\Pr(A|B) = \Pr(A)$ and the following pairs are also independent

- A and B^c
proof:

$$\begin{aligned}\Pr(A \cap B^c) &= \Pr(A) - \Pr(A \cap B) \\ &= \Pr(A) - \Pr(A) \Pr(B) \\ &= \Pr(A)(1 - \Pr(B)) \\ &= \Pr(A) \Pr(B^c)\end{aligned}$$

- A^c and B
- A^c and B^c

Random variables

Definition: A random variable is a function from a sample space S into the real numbers.

Experiment	Random variable
Toss two dice	$X = \text{sum of the numbers}$
Toss a coin 25 times	$X = \text{number of heads in 25 tosses}$
Apply different amounts of fertilizer to corn plants	$X = \text{yield/acre}$

Suppose we have a sample space

$$S = \{s_1, \dots, s_n\}$$

with a probability function \Pr and we define a random variable X with range $\mathcal{X} = \{x_1, \dots, x_m\}$. We can define a probability function \Pr_X on \mathcal{X} in the following way. We will observe $X = x_i$ if and only if the outcome of the random experiment is an $s_j \in S$ such that $X(s_j) = x_i$. Thus,

$$\Pr_X(X = x_i) = \Pr(\{s_j \in S : X(s_j) = x_i\}).$$

We will simply write $\Pr(X = x_i)$ rather than $\Pr_X(X = x_i)$.

A note on notation: Random variables are often denoted with uppercase letters and the realized values of the variables (or its range) are denoted by corresponding lowercase letters.

Distribution functions

Definition: The cumulative distribution function or cdf of a random variable (r.v.) X , denoted by $F_X(x)$ is defined by

$$F_X(x) = \Pr(X \leq x), \text{ for all } x.$$

The function $F(x)$ is a cdf if and only if the following three conditions hold:

1. $\lim_{x \rightarrow \infty} F(x) = 1$.
2. $F(x)$ is a nondecreasing function of x .
3. $F(x)$ is right-continuous; that is, for every number x_0 , $\lim_{x \downarrow x_0} F(x) = F(x_0)$.

Definition: A random variable X is continuous if $F(x)$ is a continuous function of x . A random variable X is discrete if $F(x)$ is a step function of x .

The following two statements are equivalent:

1. The random variables X and Y are identically distributed.
2. $F_X(x) = F_Y(x)$ for every x .

Density and mass functions

Definition: The probability mass function (pmf) of a discrete random variable X is given by

$$f_X(x) = \Pr(X = x) \text{ for all } x.$$

Example (Geometric probabilities) For the geometric distribution, we have the pmf

$$f_X(x) = \Pr(X = x) = \begin{cases} p(1-p)^{x-1} & \text{for } x = 1, 2, \dots \\ 0 & \text{otherwise.} \end{cases}$$

Definition: The probability density function or pdf, $f_X(x)$, of a continuous random variable X is the function that satisfies

$$F_X(x) = \int_{-\infty}^x f_X(t)dt \quad \text{for all } x.$$

A note on notation: The expression “ X has a distribution given by $F_X(x)$ ” is abbreviated symbolically by “ $X \sim F_X(x)$ ”, where we read the symbol “ \sim ” as “is distributed as”.

Example (Logistic distribution) For the logistic distribution, we have

$$F_X(x) = \frac{1}{1 + e^{-x}}$$

and, hence,

$$f_X(x) = \frac{d}{dx}F_X(x) = \frac{e^{-x}}{(1 + e^{-x})^2}.$$

A function $f_X(x)$ is a pdf (or pmf) of a random variable X if and only if

1. $f_X(x) \geq 0$ for all x
2. $\sum_x f_X(x) = 1$ (pmf) or $\int_{-\infty}^{\infty} f_X(x)dx = 1$ (pdf).

Expectations

The expected value, or expectation, of a random variable is merely its average value, where we speak of “average” value as one that is weighted according to the probability distribution.

Definition: The expected value or mean of a random variable $g(X)$, denoted by $\mathbf{E}(g(X))$, is

$$\mathbf{E}(g(X)) = \begin{cases} \int_{-\infty}^{\infty} g(x)f_X(x)dx & \text{if } X \text{ is continuous} \\ \sum_{x \in \mathcal{X}} g(x)f_X(x) = \sum_{x \in \mathcal{X}} g(x)\Pr(X = x) & \text{if } X \text{ is discrete,} \end{cases}$$

Exponential mean

Suppose $X \sim \text{Exp}(\lambda)$ distribution, that is, it has pdf given by

$$f_X(x) = \frac{1}{\lambda}e^{-x/\lambda}, \quad 0 \leq x < \infty, \quad \lambda > 0$$

Then $\mathbf{E}(X)$ is:

$$\begin{aligned}\mathbf{E}(X) &= \int_0^\infty \frac{1}{\lambda} x e^{-x/\lambda} dx \\ &= -x e^{-x/\lambda} \Big|_0^\infty + \int_0^\infty e^{-x/\lambda} dx \\ &= \int_0^\infty e^{-x/\lambda} dx = \lambda\end{aligned}$$

Binomial mean

IF X has binomial distribution, i.e. $X \sim \text{binomial}(n, p)$, its pmf is given by

$$\Pr(X = x) = \binom{n}{x} p^x (1-p)^{n-x}, \quad x = 0, 1, \dots, n,$$

where n is a positive integer, $0 \leq p \leq 1$, and for every fixed pair n and p the pmf sums to 1. The expected value of a binomial random variable is then given by

$$\mathbf{E}(X) = \sum_{x=0}^n x \binom{n}{x} p^x (1-p)^{n-x}$$

Now, use the identity $x \binom{n}{x} = n \binom{n-1}{x-1}$ to derive the Expected value.

$$\begin{aligned}\mathbf{E}(X) &= \sum_{x=1}^n x \binom{n}{x} p^x (1-p)^{n-x} \\ &= \sum_{x=1}^n n \binom{n-1}{x-1} p^x (1-p)^{n-x} \\ &= \sum_{y=0}^{n-1} n \binom{n-1}{y} p^{y+1} (1-p)^{n-(y+1)} \\ &= np \sum_{y=0}^{n-1} \binom{n-1}{y} p^y (1-p)^{n-1-y} \\ &= np,\end{aligned}$$

since the last summation must be 1, being the sum over all possible values of a $\text{binomial}(n-1, p)$ pmf.

properties:

Let X be a random variable and let a, b and c be constants. Then for any functions $g_1(x)$ and $g_2(x)$ whose expectations exist,

$$1. \quad \mathbf{E}(a \cdot g_1(X) + b \cdot g_2(X) + c) = a\mathbf{E}(g_1(X)) + b\mathbf{E}(g_2(X)) + c.$$

2. If $g_1(x) \geq 0$ for all x , then $\mathbf{E}(g_1(X)) \geq 0$.
3. If $g_1(x) \geq g_2(x)$ for all x , then $\mathbf{E}(g_1(X)) \geq \mathbf{E}(g_2(X))$.
4. If $a \leq g_1(x) \leq b$ for all x , then $a \leq \mathbf{E}(g_1(X)) \leq b$.

Moments

The various moments of a distribution are an important class of expectations.

Definition: For each integer n , the n^{th} moment of X (or $F_X(x)$), μ'_n , is

$$\mu'_n = \mathbf{E}(X^n).$$

The n^{th} central moment of X , μ_n , is

$$\mu_n = \mathbf{E}((X - \mu)^n),$$

where $\mu = \mu'_1 = \mathbf{E}(X)$.

Variance

Definition: The variance of a random variable X is its second central moment, $\mathbf{Var}(X) = \mathbf{E}((X - EX)^2)$. The positive square root of $\mathbf{Var}(X)$ is the standard deviation of X .

Exponential variance

Let X have the exponential(λ) distribution, $X \sim \text{Exp}(\lambda)$. Then the variance of X is

$$\begin{aligned} \mathbf{Var}(X) &= \mathbf{E}((X - EX)^2) = \mathbf{E}((X - \lambda)^2) \\ &= \int_0^\infty (x - \lambda)^2 \frac{1}{\lambda} e^{-x/\lambda} dx \\ &= \int_0^\infty (x^2 - 2x\lambda + \lambda^2) \frac{1}{\lambda} e^{-x/\lambda} dx \\ &= \lambda^2. \end{aligned}$$

properties

1. $\mathbf{Var}(aX + b) = a^2 \mathbf{Var}(X)$.

proof:

$$\begin{aligned} \mathbf{Var}(aX + b) &= \mathbf{E}(((aX + b) - \mathbf{E}(aX + b))^2) \\ &= \mathbf{E}((aX - aEX)^2) \\ &= a^2 \mathbf{E}((X - EX)^2) \\ &= a^2 \mathbf{Var}(X) \end{aligned}$$

$$2. \text{Var}(X) = \mathbf{E}(X^2) - (\mathbf{E}(X))^2.$$

proof:

$$\begin{aligned} \text{Var}(X) &= \mathbf{E}(X - EX)^2 \\ &= \mathbf{E}(X^2 - 2X\mathbf{E}(X) + (\mathbf{E}(X))^2) \\ &= \mathbf{E}(X^2) - 2\mathbf{E}(X)\mathbf{E}(X) + (\mathbf{E}(X))^2 \\ &= \mathbf{E}(X^2) - (\mathbf{E}(X))^2 \end{aligned}$$

Moment generating function

Definition: Let X be a random variable with cdf F_X . The moment generating function or mgf of X (or F_X), denoted by $M_X(t)$, is

$$M_X(t) = \mathbf{E}(e^{tX}),$$

provided that the expectation exists for t in some neighborhood of 0. That is, there exists an $h > 0$ such that for all t in $-h < t < h$, $\mathbf{E}(e^{tX})$ exists. If the expectation does not exist in a neighborhood of 0, we say that the moment generating function does not exist.

Property: If X has mgf $M_X(t)$, then

$$\mathbf{E}(X^n) = M_X^{(n)}(0),$$

where we define

$$M_X^{(n)}(0) = \left. \frac{d^n}{dt^n} M_X(t) \right|_{t=0}.$$

7 Lecture 7: Feb 6

Last time

- R basics

Today

- Simple Linear Regression (JF Chapter 5)

Simple linear regression

Figure 7.1 shows Davis's data on the measured and reported weight in kilograms of 101 women who were engaged in regular exercise.

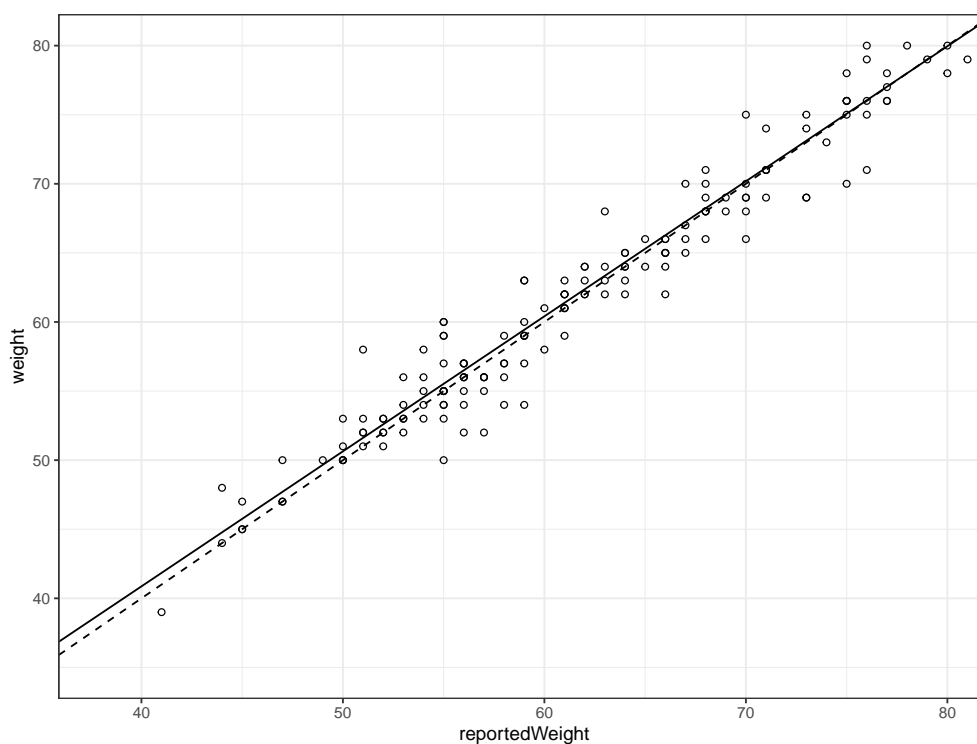


Figure 7.1: Scatterplot of Davis's data on the measured and reported weight of 101 women. The dashed line gives $y = x$.

It's reasonable to assume that the relationship between measured and reported weight appears to be linear. Denote:

- measured weight by y_i : **response variable** or **dependent variable**
- reported weight by x_i : **predictor variable** or **independent variable**
- intercept: β_0

- slope: β_1
- residual/error term ϵ_i .

Then the simple linear regression model writes:

$$y_i = \beta_0 + \beta_1 x_i + \epsilon_i.$$

For given $(\hat{\beta}_0, \hat{\beta}_1)$ values, the *fitted value* or *predicted value* for observation i is:

$$\hat{y}_i = \hat{\beta}_0 + \hat{\beta}_1 x_i.$$

Therefore, the residual is

$$\hat{\epsilon}_i = y_i - \hat{y}_i$$

Fitting a linear model

Choose the “best” values for β_0, β_1 such that

$$SS[E] = \sum_1^n \left(y_i - (\hat{\beta}_0 + \hat{\beta}_1 x_i) \right)^2 = \sum_1^n (y_i - \hat{y}_i)^2 = \sum_1^n \epsilon_i^2$$

is minimized. These are **least squares** (LS) estimates:

$$\begin{aligned}\hat{\beta}_0 &= \bar{y} - \hat{\beta}_1 \bar{x} \\ \hat{\beta}_1 &= \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{\sum (x_i - \bar{x})^2}.\end{aligned}$$

Definition: The line satisfying the equation

$$y = \hat{\beta}_0 + \hat{\beta}_1 x$$

is called the linear regression of y on x which is also called the least squares line.

For Davis’s data, we have

$$\begin{aligned}n &= 101 \\ \bar{y} &= \frac{5780}{101} = 57.228 \\ \bar{x} &= \frac{5731}{101} = 56.743 \\ \sum (x_i - \bar{x})(y_i - \bar{y}) &= 4435.9 \\ \sum (x_i - \bar{x})^2 &= 4539.3,\end{aligned}$$

so that

$$\begin{aligned}\hat{\beta}_1 &= \frac{4435.9}{4539.3} = 0.97722 \\ \hat{\beta}_0 &= 57.228 - 0.97722 \times 56.743 = 1.7776\end{aligned}$$

Least squares estimates

The simple linear regression (SLR) model writes:

$$y_i = \beta_0 + \beta_1 x_i + \epsilon_i.$$

The least squares estimates minimizes the sum of squared error (SSE) which is

$$SS[E] = \sum_1^n \left(y_i - (\hat{\beta}_0 + \hat{\beta}_1 x_i) \right)^2 = \sum_1^n (y_i - \hat{y}_i)^2 = \sum_1^n \hat{\epsilon}_i^2.$$

The **least squares** (LS) estimates (in vector form):

$$\hat{\beta}_{ls} = \begin{pmatrix} \hat{\beta}_0 \\ \hat{\beta}_1 \end{pmatrix} = \begin{pmatrix} \bar{y} - \hat{\beta}_1 \bar{x} \\ \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{\sum (x_i - \bar{x})^2} \end{pmatrix}.$$

Definition: The line satisfying the equation

$$y = \hat{\beta}_0 + \hat{\beta}_1 x$$

is called the linear regression of y on x which is also called the least squares line.

SLR Model in Matrix Form

$$\begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix} = \begin{bmatrix} \beta_0 + \beta_1 x_1 \\ \beta_0 + \beta_1 x_2 \\ \vdots \\ \beta_0 + \beta_1 x_n \end{bmatrix} + \begin{bmatrix} \epsilon_1 \\ \epsilon_2 \\ \vdots \\ \epsilon_n \end{bmatrix}$$

$$\begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix} = \begin{bmatrix} 1 & x_1 \\ 1 & x_2 \\ \vdots & \vdots \\ 1 & x_n \end{bmatrix} \begin{bmatrix} \beta_0 \\ \beta_1 \end{bmatrix} + \begin{bmatrix} \epsilon_1 \\ \epsilon_2 \\ \vdots \\ \epsilon_n \end{bmatrix}$$

Jargons

- \mathbf{X} is called the *design matrix*
- $\boldsymbol{\beta}$ is the vector of parameters
- $\boldsymbol{\epsilon}$ is the error vector
- \mathbf{Y} is the response vector.

The Design Matrix

$$\mathbf{X}_{n \times 2} = \begin{bmatrix} 1 & x_1 \\ 1 & x_2 \\ \vdots & \vdots \\ 1 & x_n \end{bmatrix}$$

Vector of Parameters

$$\boldsymbol{\beta}_{2 \times 1} = \begin{bmatrix} \beta_0 \\ \beta_1 \end{bmatrix}$$

Vector of Error terms

$$\boldsymbol{\epsilon}_{n \times 1} = \begin{bmatrix} \epsilon_1 \\ \epsilon_2 \\ \vdots \\ \epsilon_n \end{bmatrix}$$

Vector of Responses

$$\mathbf{Y}_{n \times 1} = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix}$$

Gramian Matrix

$$\mathbf{X}^T \mathbf{X} = \begin{bmatrix} n & \sum_i x_i \\ \sum_i x_i & \sum_i x_i^2 \end{bmatrix}$$

Therefore, we have

$$\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\epsilon}.$$

Assume the Gramian matrix has full rank (which actually should be the case, why?), we want to show that

$$\hat{\boldsymbol{\beta}}_{ls} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{Y}.$$

The inverse of the Gramian matrix is

$$(\mathbf{X}^T \mathbf{X})^{-1} = \frac{1}{n \sum_i (x_i - \bar{x})^2} \begin{bmatrix} \sum_i x_i^2 & -\sum_i x_i \\ -\sum_i x_i & n \end{bmatrix}$$

Now we have

$$\begin{aligned}
\hat{\beta}_{ls} &= (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{Y} \\
&= \frac{1}{n \sum_i (x_i - \bar{x})^2} \begin{bmatrix} \sum_i x_i^2 & -\sum_i x_i \\ -\sum_i x_i & n \end{bmatrix} \begin{bmatrix} \mathbf{1}_n^T \\ \mathbf{x}^T \end{bmatrix} \mathbf{y} \\
&= \frac{1}{n \sum_i (x_i - \bar{x})^2} \begin{bmatrix} \sum_i x_i^2 & -\sum_i x_i \\ -\sum_i x_i & n \end{bmatrix} \begin{bmatrix} \sum_i y_i \\ \sum_i x_i y_i \end{bmatrix} \\
&= \frac{1}{n \sum_i (x_i - \bar{x})^2} \begin{bmatrix} (\sum_i x_i^2)(\sum_i y_i) - (\sum_i x_i)(\sum_i x_i y_i) \\ n \sum_i x_i y_i - (\sum_i x_i)(\sum_i y_i) \end{bmatrix} \\
&= \begin{bmatrix} \bar{y} - \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{\sum (x_i - \bar{x})^2} \bar{x} \\ \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{\sum (x_i - \bar{x})^2} \end{bmatrix}
\end{aligned}$$

9 Lecture 9: Feb 8

Last time

- Introduction of simple linear regression

Today

- Simple correlation

From last lecture: Assume the Gramian matrix has full rank (which actually should be the case, why?)

$$\mathbf{X}^T \mathbf{X} = \begin{bmatrix} n & \sum_i x_i \\ \sum_i x_i & \sum_i x_i^2 \end{bmatrix}$$

Proof: By Cauchy-Schwarz inequality, we have

$$n \sum_i x_i^2 \geq (\sum_i x_i)^2$$

where the equality holds only if all x_i are equal.

Some properties:

- (a) $\sum x_i \hat{\epsilon}_i = 0$.
- (b) $\sum \hat{y}_i \hat{\epsilon}_i = 0$ (HW1).

Proof: For (a), we look at

$$\begin{aligned} & \mathbf{X}^T \hat{\epsilon} \\ &= \mathbf{X}^T (\mathbf{Y} - \mathbf{X} \hat{\beta}) \\ &= \mathbf{X}^T [\mathbf{Y} - \mathbf{X} (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{Y}] \\ &= \mathbf{X}^T \mathbf{Y} - \mathbf{X}^T \mathbf{X} (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{Y} \\ &= \mathbf{X}^T \mathbf{Y} - \mathbf{X}^T \mathbf{Y} \\ &= \mathbf{0} \end{aligned}$$

Other quantities in Matrix Form

Fitted values

$$\hat{\mathbf{Y}} = \begin{bmatrix} \hat{y}_1 \\ \hat{y}_2 \\ \vdots \\ \hat{y}_n \end{bmatrix} = \begin{bmatrix} \hat{\beta}_0 + \hat{\beta}_1 x_1 \\ \hat{\beta}_0 + \hat{\beta}_1 x_2 \\ \vdots \\ \hat{\beta}_0 + \hat{\beta}_1 x_n \end{bmatrix} = \begin{bmatrix} 1 & x_1 \\ 1 & x_2 \\ \vdots & \vdots \\ 1 & x_n \end{bmatrix} \begin{bmatrix} \hat{\beta}_0 \\ \hat{\beta}_1 \end{bmatrix} = \mathbf{X} \hat{\beta}$$

Hat matrix

$$\begin{aligned}\hat{\mathbf{Y}} &= \mathbf{X}\hat{\boldsymbol{\beta}} \\ \hat{\mathbf{Y}} &= \mathbf{X}(\mathbf{X}^T\mathbf{X})^{-1}\mathbf{X}^T\mathbf{Y} \\ \hat{\mathbf{Y}} &= \mathbf{H}\mathbf{Y}\end{aligned}$$

where $\mathbf{H} = \mathbf{X}(\mathbf{X}^T\mathbf{X})^{-1}\mathbf{X}^T$ is called “hat matrix” because it turns \mathbf{Y} into $\hat{\mathbf{Y}}$.

Simple correlation

Having calculated the least squares line, it is of interest to determine how closely the line fits the scatter of points. There are many ways of answering it. The standard deviation of the residuals, S_E , often called the *standard error of the regression* or the *residue standard error*, provides one sort of answer. Because of estimation considerations, the variance of the residuals is defined using *degrees of freedom* $n - 2$:

$$S_\epsilon^2 = \frac{\sum \hat{\epsilon}_i^2}{n - 2}.$$

The residual standard error is,

$$S_\epsilon = \sqrt{\frac{\sum \hat{\epsilon}_i^2}{n - 2}}$$

For the Davis’s data, the sum of squared residuals is $\sum \hat{\epsilon}_i^2 = 418.87$, and thus the standard error of the regression is

$$S_\epsilon = \sqrt{\frac{418.87}{101 - 2}} = 2.0569\text{kg}.$$

On average, using the least-squares regression line to predict measured weight from reported weight results in an error of about 2 kg.

Sum of squares:

- Total sum of squares (TSS) for Y: $\text{TSS} = \sum (y_i - \bar{y})^2$
- Residual sum of squares (RSS): $\text{RSS} = \sum (y_i - \hat{y}_i)^2$
- regression sum of squares (RegSS): $\text{RegSS} = \text{TSS} - \text{RSS} = \sum (\hat{y}_i - \bar{y})^2$
- $\text{RegSS} + \text{RSS} = \text{TSS}$

Sample correlation coefficient

Definition: The sample correlation coefficient r_{xy} of the paired data $(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)$ is defined by

$$r_{xy} = \frac{\sum (x_i - \bar{x})(y_i - \bar{y}) / (n - 1)}{\sqrt{\sum (x_i - \bar{x})^2 / (n - 1) \times \sum (y_i - \bar{y})^2 / (n - 1)}} = \frac{s_{xy}}{s_x s_y}$$

s_{xy} is called the sample covariance of x and y :

$$s_{xy} = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{n - 1}$$

$s_x = \sqrt{\sum (x_i - \bar{x})^2 / (n - 1)}$ and $s_y = \sqrt{\sum (y_i - \bar{y})^2 / (n - 1)}$ are, respectively, the sample standard deviations of X and Y .

Some properties of r_{xy} :

- r_{xy} is a measure of the linear association between x and y in a dataset.
- correlation coefficients are always between -1 and 1 :

$$-1 \leq r_{xy} \leq 1$$

- The closer r_{xy} is to 1 , the stronger the positive linear association between x and y
- The closer r_{xy} is to -1 , the stronger the negative linear association between x and y
- The bigger $|r_{xy}|$, the stronger the linear association
- If $|r_{xy}| = 1$, then x and y are said to be perfectly correlated.
- $\hat{\beta}_1 = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{\sum (x_i - \bar{x})^2} = \frac{s_{xy}}{s_x^2} = r_{xy} \frac{s_y}{s_x}$

11 Lecture 11: Feb 13

Last time

- Introduction of simple linear regression

Today

- HW2 posted
- The statistical model of the SLR (JF chapter 6)
- Properties of the Least-Squares estimator
- Inference of SLR model

R-square

The ratio of RegSS to TSS is called the *coefficient of determination*, or sometimes, simply “r-square”. it represents the proportion of variation observed in the response variable y which can be “explained” by its linear association with x .

- In simple linear regression, “r-square” is in fact equal to r_{xy}^2 . (But this isn’t the case in multiple regression.)
- It is also equal to the squared correlation between y_i and \hat{y}_i . (This is the case in multiple regression.)

For Davis’s regression of measured on reported weight:

$$\text{TSS} = 4753.8$$

$$\text{RSS} = 418.87$$

$$\text{RegSS} = 4334.9$$

Thus,

$$r^2 = \frac{4334.9}{4753.8} = 1 - \frac{418.87}{4753.8} = 0.9119$$

The statistical model of Simple Linear Regression

Standard statistical inference in simple regression is based on a *statistical model* that describes the population or process that is sampled:

$$y_i = \beta_0 + \beta_1 x_i + \epsilon_i$$

where the coefficients β_0 and β_1 are the *population regression parameters*. The data are randomly sampled from some population of interest.

- y_i is the value of the response variable

- x_i is the explanatory variable
- ϵ_i represents the aggregated omitted causes of y (i.e., the causes of y beyond the explanatory variable), other explanatory variables that could have been included in the regression model, measurement error in y , and whatever component of y is inherently random.

Key assumptions of SLR

The key assumptions of the SLR model concern the behavior of the errors, equivalently, the distribution of y conditional on x :

- *Linearity.* The expectation of the error given the value of x is 0: $\mathbf{E}(\epsilon) \equiv \mathbf{E}(\epsilon|x_i) = 0$. And equivalently, the expected value of the response variable is a linear function of the explanatory variable: $\mu_i \equiv \mathbf{E}(y_i) \equiv \mathbf{E}(y_i|x_i) = \mathbf{E}(\beta_0 + \beta_1 x_i + \epsilon_i|x_i) = \beta_0 + \beta_1 x_i$.
- *Constant variance.* The variance of the errors is the same regardless of the value of x : $\mathbf{Var}(\epsilon|x_i) = \sigma_\epsilon^2$. The constant error variance implies constant conditional variance of y on given x : $\mathbf{Var}(y|x_i) = \mathbf{E}((y_i - \mu_i)^2) = \mathbf{E}((y_i - \beta_0 - \beta_1 x_i)^2) = \mathbf{E}(\epsilon_i^2) = \sigma_\epsilon^2$. (Question: why the last equal sign?)
- *Normality.* The errors are independent identically distributed with Normal distribution with mean 0 and variance σ_ϵ^2 . Write as $\epsilon_i \stackrel{iid}{\sim} N(0, \sigma_\epsilon^2)$. Equivalently, the conditional distribution of the response variable is normal: $y_i \stackrel{iid}{\sim} N(\beta_0 + \beta_1 x_i, \sigma_\epsilon^2)$.
- *Independence.* The observations are sampled independently.
- *Fixed X , or X measured without error and independent of the error.*
 - For experimental research where X values are under direct control of the researcher (i.e. X 's are fixed). If the experiment were replicated, then the values of X would remain the same.
 - For research where X values are sampled, we assume the explanatory variable is measured without error and the explanatory variable and the error are independent in the population from which the sample is drawn.
- *X is not invariant.* X 's can not be all the same.

Figure 11.1 shows the assumptions of linearity, constant variance, and normality in SLR model.



Figure 11.1: The assumptions of linearity, constant variance, and normality in simple regression. The graph shows the conditional population distributions $\Pr(Y|x)$ of Y for several values of the explanatory variable X , labeled as x_1, x_2, \dots, x_5 . The conditional means of Y given x are denoted μ_1, \dots, μ_5 .

Properties of the Least-Squares estimator

Under the strong assumptions of the simple regression model, the sample least squares coefficients $\hat{\beta}_{ls}$ have several desirable properties as estimators of the population regression coefficients β_0 and β_1 :

- The least-squares intercept and slope are *linear estimators*, in the sense that they are linear functions of the observations y_i .

Proof:

method (a) $\hat{\beta} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{Y}$

method (b) $\hat{\beta}_1 = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{\sum (x_i - \bar{x})^2} = \frac{\sum (x_i - \bar{x})y_i}{\sum (x_i - \bar{x})^2} - \frac{\sum (x_i - \bar{x})\bar{y}}{\sum (x_i - \bar{x})^2} = \sum \frac{(x_i - \bar{x})}{\sum (x_i - \bar{x})^2} y_i = \sum k_i y_i$ where

$$k_i = \frac{(x_i - \bar{x})}{\sum (x_i - \bar{x})^2}$$

and $\hat{\beta}_0 = \bar{y} - \hat{\beta}_1 \bar{x}$

- The sample least-squares coefficients are *unbiased estimators* of the population regression coefficients:

$$\mathbf{E}(\hat{\beta}_0) = \beta_0$$

$$\mathbf{E}(\hat{\beta}_1) = \beta_1$$

Proof:

method (a) $\mathbf{E}(\hat{\beta}) = \mathbf{E}((\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{Y}) = \mathbf{E}((\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{X} \beta) = \beta$. (note: $\mathbf{E}(Y) = \mathbf{E}(\mathbf{X}\beta + \epsilon) = \mathbf{E}(\mathbf{X}\beta) + \mathbf{E}(\epsilon) = \mathbf{X}\beta$)

method (b) recall that $\hat{\beta}_1 = \sum k_i y_i$ where $k_i = \frac{(x_i - \bar{x})}{\sum (x_i - \bar{x})^2}$. First, we want to show

1. $\sum k_i = 0$
2. $\sum k_i x_i = 1$

They are actually quite easy: $\sum k_i = \sum_i \frac{(x_i - \bar{x})}{\sum_j (x_j - \bar{x})^2} = \frac{(\sum_i x_i) - n\bar{x}}{\sum_j (x_j - \bar{x})^2} = 0$, and $\sum k_i x_i = \sum_i \frac{(x_i - \bar{x})x_i}{\sum_j (x_j - \bar{x})^2} = \frac{(\sum_i x_i^2) - \bar{x}(\sum_i x_i)}{\sum_j (x_j - \bar{x})^2} = \frac{(\sum_i x_i^2) - n\bar{x}^2}{\sum_j (x_j - \bar{x})^2} = 1$.

Now $\mathbf{E}(\hat{\beta}_1) = \mathbf{E}(\sum k_i y_i) = \sum [k_i \mathbf{E}(y_i)] = \sum [k_i(\beta_0 + \beta_1 x_i)] = \beta_0 \sum k_i + \beta_1 \sum (k_i x_i) = \beta_1$, and $\mathbf{E}(\hat{\beta}_0) = \mathbf{E}(\bar{y} - \hat{\beta}_1 \bar{x}) = \mathbf{E}(\bar{y}) - \bar{x} \mathbf{E}(\hat{\beta}_1) = \mathbf{E}(\frac{1}{n} \sum y_i) - \bar{x} \beta_1 = \frac{1}{n} [\sum \mathbf{E}(y_i)] - \bar{x} \beta_1 = \frac{1}{n} \sum [\beta_0 + x_i \beta_1] - \bar{x} \beta_1 = \beta_0$

- Both $\hat{\beta}_0$ and $\hat{\beta}_1$ have simple sampling variances:

$$\text{Var}(\hat{\beta}_0) = \frac{\sigma_\epsilon^2 \sum x_i^2}{n \sum (x_i - \bar{x})^2}$$

$$\text{Var}(\hat{\beta}_1) = \frac{\sigma_\epsilon^2}{\sum (x_i - \bar{x})^2}$$

Proof:

$\text{Var}(\hat{\beta}_1) = \text{Var}(\sum k_i y_i) = \sum k_i^2 \text{Var}(y_i) = \sigma_\epsilon^2 \sum k_i^2 = \sigma_\epsilon^2 \frac{\sum_i (x_i - \bar{x})^2}{[\sum_j (x_j - \bar{x})^2]^2} = \frac{\sigma_\epsilon^2}{\sum (x_i - \bar{x})^2}$, and

$\text{Var}(\hat{\beta}_0) = \text{Var}(\bar{y} - \hat{\beta}_1 \bar{x}) = \text{Var}(\bar{y}) + (\bar{x})^2 \text{Var}(\hat{\beta}_1) - 2\bar{x} \text{Cov}(\bar{Y}, \hat{\beta}_1)$.

Now,

$$\text{Var}(\bar{y}) = \text{Var}\left(\frac{1}{n} \sum_{i=1}^n y_i\right) = \frac{1}{n^2} \sum_{i=1}^n \text{Var}(y_i) = \frac{\sigma^2}{n},$$

$$\text{Var}(\hat{\beta}_1) = \frac{\sigma_\epsilon^2}{\sum (x_i - \bar{x})^2},$$

and

$$\begin{aligned} \text{Cov}(\bar{Y}, \hat{\beta}_1) &= \text{Cov}\left\{\frac{1}{n} \sum_{i=1}^n Y_i, \frac{\sum_{j=1}^n (x_j - \bar{x}) Y_j}{\sum_{i=1}^n (x_i - \bar{x})^2}\right\} \\ &= \frac{1}{n} \frac{1}{\sum_{i=1}^n (x_i - \bar{x})^2} \text{Cov}\left\{\sum_{i=1}^n Y_i, \sum_{j=1}^n (x_j - \bar{x}) Y_j\right\} \\ &= \frac{1}{n \sum_{i=1}^n (x_i - \bar{x})^2} \sum_{i=1}^n (x_j - \bar{x}) \sum_{j=1}^n \text{Cov}(Y_i, Y_j) \\ &= \frac{1}{n \sum_{i=1}^n (x_i - \bar{x})^2} \sum_{i=1}^n (x_j - \bar{x}) \sigma^2 \\ &= 0. \end{aligned}$$

Finally,

$$\begin{aligned}\text{Var}(\hat{\beta}_0) &= \frac{\sigma^2}{n} + \frac{\sigma^2 \bar{x}^2}{\sum_{i=1}^n (x_i - \bar{x})^2} \\ &= \frac{\sigma^2}{n \sum_{i=1}^n (x_i - \bar{x})^2} \left\{ \sum_{i=1}^n (x_i - \bar{x})^2 + n \bar{x}^2 \right\} \\ &= \frac{\sigma^2 \sum_{i=1}^n x_i^2}{n \sum_{i=1}^n (x_i - \bar{x})^2}.\end{aligned}$$

- Rewrite the formula for $\text{Var}(\hat{\beta}_1) = \frac{\sigma_\epsilon^2}{(n-1)S_X^2}$, we see that the sampling variance of the slope estimate will be small when
 - The error variance σ_ϵ^2 is small
 - The sample size n is large
 - The explanatory-variable values are spread out (i.e. have a large variance, S_X^2)

- (Gauss-Markov theorem) Under the assumptions of linearity, constant variance, and independence, the least-squares estimators are BLUE (Best Linear Unbiased Estimator), that is they have the smallest sampling variance and are unbiased. (show this)

Proof:

Let $\tilde{\beta}_1$ be another linear unbiased estimator such that $\tilde{\beta}_1 = \sum c_i y_i$. For $\tilde{\beta}_1$ is still unbiased as above, $\mathbf{E}(\tilde{\beta}_1) = \beta_0 \sum c_i + \beta_1 \sum c_i x_i = \beta_1$ for all β_1 , we have $\sum c_i = 0$ and $\sum c_i x_i = 1$.

$$\mathbf{Var}(\tilde{\beta}_1) = \sigma_\epsilon^2 \sum c_i^2$$

Let $c_i = k_i + d_i$, then

$$\begin{aligned}\mathbf{Var}(\tilde{\beta}_1) &= \sigma_\epsilon^2 \sum (k_i + d_i)^2 \\ &= \sigma_\epsilon^2 \left[\sum k_i^2 + \sum d_i^2 + 2 \sum k_i d_i \right] \\ &= \mathbf{Var}(\hat{\beta}_1) + \sigma_\epsilon^2 \sum d_i^2 + 2\sigma_\epsilon^2 \sum k_i d_i\end{aligned}$$

Now we show the last term is 0 to finish the proof.

$$\begin{aligned}\sum k_i d_i &= \sum k_i (c_i - k_i) = \sum c_i k_i - \sum k_i^2 \\ &= \sum_i \left[c_i \frac{x_i - \bar{x}}{\sum_j (x_j - \bar{x})^2} \right] - \frac{1}{\sum_i (x_i - \bar{x})^2} \\ &= 0\end{aligned}$$

- Under the full suite of assumptions, the least-squares coefficients $\hat{\beta}_0$ and $\hat{\beta}_1$ are the maximum-likelihood estimators of β_0 and β_1 . (show this)

Proof:

The log likelihood under the full suite of assumptions is $\ell = -\log \left[(2\pi)^{\frac{n}{2}} \sigma_\epsilon^n \right] - \frac{1}{2\sigma_\epsilon^2} (\mathbf{Y} - \mathbf{X}\beta)^T (\mathbf{Y} - \mathbf{X}\beta)$. Maximizing the likelihood is equivalent as minimizing $(\mathbf{Y} - \mathbf{X}\beta)^T (\mathbf{Y} - \mathbf{X}\beta) = \epsilon^T \epsilon$ which is the SSE.

- Under the assumption of normality, the least-squares coefficients are themselves normally distributed. Summing up,

$$\hat{\beta}_0 \sim N\left(\beta_0, \frac{\sigma_\epsilon^2 \sum x_i^2}{n \sum (x_i - \bar{x})^2}\right)$$

$$\hat{\beta}_1 \sim N\left(\beta_1, \frac{\sigma_\epsilon^2}{\sum (x_i - \bar{x})^2}\right)$$

12 Lecture 12 Feb 15

Last time

- Properties of the Least-Squares estimator

Today

- HW1 due Feb 17
- Inference of SLR model
- Multiple linear regression

Properties of the Least-Squares estimator

Under the strong assumptions of the simple regression model, the sample least squares coefficients $\hat{\beta}_{ls}$ have several desirable properties as estimators of the population regression coefficients β_0 and β_1 :

- The least-squares intercept and slope are *linear estimators*, in the sense that they are linear functions of the observations y_i .

Proof:

method (a) $\hat{\beta} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{Y}$

method (b) $\hat{\beta}_1 = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{\sum (x_i - \bar{x})^2} = \frac{\sum (x_i - \bar{x})y_i}{\sum (x_i - \bar{x})^2} - \frac{\sum (x_i - \bar{x})\bar{y}}{\sum (x_i - \bar{x})^2} = \sum \frac{(x_i - \bar{x})}{\sum (x_i - \bar{x})^2} y_i = \sum k_i y_i$ where

$$k_i = \frac{(x_i - \bar{x})}{\sum (x_i - \bar{x})^2}$$

and $\hat{\beta}_0 = \bar{y} - \hat{\beta}_1 \bar{x}$

- The sample least-squares coefficients are *unbiased estimators* of the population regression coefficients:

$$\mathbf{E}(\hat{\beta}_0) = \beta_0$$

$$\mathbf{E}(\hat{\beta}_1) = \beta_1$$

Proof:

method (a) $\mathbf{E}(\hat{\beta}) = \mathbf{E}((\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{Y}) = \mathbf{E}((\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{X} \beta) = \beta$. (note: $\mathbf{E}(Y) = \mathbf{E}(\mathbf{X}\beta + \epsilon) = \mathbf{E}(\mathbf{X}\beta) + \mathbf{E}(\epsilon) = \mathbf{X}\beta$)

method (b) recall that $\hat{\beta}_1 = \sum k_i y_i$ where $k_i = \frac{(x_i - \bar{x})}{\sum (x_i - \bar{x})^2}$. First, we want to show

$$1. \sum k_i = 0$$

$$2. \sum k_i x_i = 1$$

They are actually quite easy: $\sum k_i = \sum_i \frac{(x_i - \bar{x})}{\sum_j (x_j - \bar{x})^2} = \frac{(\sum_i x_i) - n\bar{x}}{\sum_j (x_j - \bar{x})^2} = 0$, and $\sum k_i x_i = \sum_i \frac{(x_i - \bar{x})x_i}{\sum_j (x_j - \bar{x})^2} = \frac{(\sum_i x_i^2) - \bar{x}(\sum_i x_i)}{\sum_j (x_j - \bar{x})^2} = \frac{(\sum_i x_i^2) - n\bar{x}^2}{\sum_j (x_j - \bar{x})^2} = 1$.

Now $\mathbf{E}(\hat{\beta}_1) = \mathbf{E}(\sum k_i y_i) = \sum [k_i \mathbf{E}(y_i)] = \sum [k_i (\beta_0 + \beta_1 x_i)] = \beta_0 \sum k_i + \beta_1 \sum (k_i x_i) =$

$$\beta_1, \text{ and } \mathbf{E}(\hat{\beta}_0) = \mathbf{E}(\bar{y} - \hat{\beta}_1 \bar{x}) = \mathbf{E}(\bar{y}) - \bar{x} \mathbf{E}(\hat{\beta}_1) = \mathbf{E}\left(\frac{1}{n} \sum y_i\right) - \bar{x} \beta_1 = \frac{1}{n} [\sum \mathbf{E}(y_i)] - \bar{x} \beta_1 = \frac{1}{n} \sum [\beta_0 + x_i \beta_1] - \bar{x} \beta_1 = \beta_0$$

- Both $\hat{\beta}_0$ and $\hat{\beta}_1$ have simple sampling variances:

$$\text{Var}(\hat{\beta}_0) = \frac{\sigma_\epsilon^2 \sum x_i^2}{n \sum (x_i - \bar{x})^2}$$

$$\text{Var}(\hat{\beta}_1) = \frac{\sigma_\epsilon^2}{\sum (x_i - \bar{x})^2}$$

Proof:

$$\text{Var}(\hat{\beta}_1) = \text{Var}(\sum k_i y_i) = \sum k_i^2 \text{Var}(y_i) = \sigma_\epsilon^2 \sum k_i^2 = \sigma_\epsilon^2 \frac{\sum_i (x_i - \bar{x})^2}{[\sum_j (x_j - \bar{x})^2]^2} = \frac{\sigma_\epsilon^2}{\sum (x_i - \bar{x})^2}, \text{ and}$$

$$\text{Var}(\hat{\beta}_0) = \text{Var}(\bar{y} - \hat{\beta}_1 \bar{x}) = \text{Var}(\bar{y}) + (\bar{x})^2 \text{Var}(\hat{\beta}_1) - 2\bar{x} \text{Cov}(\bar{Y}, \hat{\beta}_1).$$

Now,

$$\text{Var}(\bar{y}) = \text{Var}\left(\frac{1}{n} \sum_{i=1}^n y_i\right) = \frac{1}{n^2} \sum_{i=1}^n \text{Var}(y_i) = \frac{\sigma^2}{n},$$

$$\text{Var}(\hat{\beta}_1) = \frac{\sigma_\epsilon^2}{\sum (x_i - \bar{x})^2},$$

and

$$\begin{aligned} \text{Cov}(\bar{Y}, \hat{\beta}_1) &= \text{Cov}\left\{\frac{1}{n} \sum_{i=1}^n Y_i, \frac{\sum_{j=1}^n (x_j - \bar{x}) Y_j}{\sum_{i=1}^n (x_i - \bar{x})^2}\right\} \\ &= \frac{1}{n} \frac{1}{\sum_{i=1}^n (x_i - \bar{x})^2} \text{Cov}\left\{\sum_{i=1}^n Y_i, \sum_{j=1}^n (x_j - \bar{x}) Y_j\right\} \\ &= \frac{1}{n \sum_{i=1}^n (x_i - \bar{x})^2} \sum_{i=1}^n (x_j - \bar{x}) \sum_{j=1}^n \text{Cov}(Y_i, Y_j) \\ &= \frac{1}{n \sum_{i=1}^n (x_i - \bar{x})^2} \sum_{i=1}^n (x_j - \bar{x}) \sigma^2 \\ &= 0. \end{aligned}$$

Finally,

$$\begin{aligned} \text{Var}(\hat{\beta}_0) &= \frac{\sigma^2}{n} + \frac{\sigma^2 \bar{x}^2}{\sum_{i=1}^n (x_i - \bar{x})^2} \\ &= \frac{\sigma^2}{n \sum_{i=1}^n (x_i - \bar{x})^2} \left\{ \sum_{i=1}^n (x_i - \bar{x})^2 + n \bar{x}^2 \right\} \\ &= \frac{\sigma^2 \sum_{i=1}^n x_i^2}{n \sum_{i=1}^n (x_i - \bar{x})^2}. \end{aligned}$$

- Rewrite the formula for $\text{Var}(\hat{\beta}_1) = \frac{\sigma_\epsilon^2}{(n-1)S_X^2}$, we see that the sampling variance of the slope estimate will be small when

- The error variance σ_ϵ^2 is small
- The sample size n is large
- The explanatory-variable values are spread out (i.e. have a large variance, S_X^2)
- (Gauss-Markov theorem) Under the assumptions of linearity, constant variance, and independence, the least-squares estimators are BLUE (Best Linear Unbiased Estimator), that is they have the smallest sampling variance and are unbiased. (show this)

Proof:

Let $\tilde{\beta}_1$ be another linear unbiased estimator such that $\tilde{\beta}_1 = \sum c_i y_i$. For $\tilde{\beta}_1$ is still unbiased as above, $\mathbf{E}(\tilde{\beta}_1) = \beta_0 \sum c_i + \beta_1 \sum c_i x_i = \beta_1$ for all β_1 , we have $\sum c_i = 0$ and $\sum c_i x_i = 1$.

$$\mathbf{Var}(\tilde{\beta}_1) = \sigma_\epsilon^2 \sum c_i^2$$

Let $c_i = k_i + d_i$, then

$$\begin{aligned} \mathbf{Var}(\tilde{\beta}_1) &= \sigma_\epsilon^2 \sum (k_i + d_i)^2 \\ &= \sigma_\epsilon^2 \left[\sum k_i^2 + \sum d_i^2 + 2 \sum k_i d_i \right] \\ &= \mathbf{Var}(\hat{\beta}_1) + \sigma_\epsilon^2 \sum d_i^2 + 2\sigma_\epsilon^2 \sum k_i d_i \end{aligned}$$

Now we show the last term is 0 to finish the proof.

$$\begin{aligned} \sum k_i d_i &= \sum k_i (c_i - k_i) = \sum c_i k_i - \sum k_i^2 \\ &= \sum_i \left[c_i \frac{x_i - \bar{x}}{\sum_j (x_j - \bar{x})^2} \right] - \frac{1}{\sum_i (x_i - \bar{x})^2} \\ &= 0 \end{aligned}$$

- Under the full suite of assumptions, the least-squares coefficients $\hat{\beta}_0$ and $\hat{\beta}_1$ are the maximum-likelihood estimators of β_0 and β_1 . (show this)

Proof:

The log likelihood under the full suite of assumptions is $\ell = -\log \left[(2\pi)^{\frac{n}{2}} \sigma_\epsilon^n \right] - \frac{1}{2\sigma_\epsilon^2} (\mathbf{Y} - \mathbf{X}\beta)^T (\mathbf{Y} - \mathbf{X}\beta)$. Maximizing the likelihood is equivalent as minimizing $(\mathbf{Y} - \mathbf{X}\beta)^T (\mathbf{Y} - \mathbf{X}\beta) = \epsilon^T \epsilon$ which is the SSE.

- Under the assumption of normality, the least-squares coefficients are themselves normally distributed. Summing up,

$$\begin{aligned} \hat{\beta}_0 &\sim N\left(\beta_0, \frac{\sigma_\epsilon^2 \sum x_i^2}{n \sum (x_i - \bar{x})^2}\right) \\ \hat{\beta}_1 &\sim N\left(\beta_1, \frac{\sigma_\epsilon^2}{\sum (x_i - \bar{x})^2}\right) \end{aligned}$$

Statistical inference of the SLR model

Now we have the distribution of $\hat{\beta}_0$ and $\hat{\beta}_1$

$$\begin{aligned}\hat{\beta}_0 &\sim N(\beta_0, \frac{\sigma_\epsilon^2 \sum x_i^2}{n \sum (x_i - \bar{x})^2}) \\ \hat{\beta}_1 &\sim N(\beta_1, \frac{\sigma_\epsilon^2}{\sum (x_i - \bar{x})^2}).\end{aligned}$$

However, σ_ϵ is never known in practice. Instead, an *unbiased* estimator of σ_ϵ^2 is given by

$$\hat{\sigma}_\epsilon^2 = MS[E] = \frac{SS[E]}{n-2}.$$

show that $\mathbf{E}(\sum (y_i - \hat{y}_i)^2) = \sigma_\epsilon^2(n-2)$.

Proof:

$$MS[E] = \frac{\sum (y_i - \hat{y}_i)^2}{n-2},$$

we want to show $\mathbf{E}(\sum (y_i - \hat{y}_i)^2) = \sigma_\epsilon^2(n-2)$.

LHS: $\mathbf{E}(\sum (y_i - \hat{y}_i)^2) = \sum_i [\mathbf{E}(y_i - \hat{y}_i)^2]$

and $\mathbf{E}[(y_i - \hat{y}_i)^2] = \text{Var}(y_i - \hat{y}_i) + [\mathbf{E}(y_i - \hat{y}_i)]^2 = \text{Var}(y_i - \hat{y}_i) = \text{Var}(y_i) + \text{Var}(\hat{y}_i) - 2\text{cov}(y_i, \hat{y}_i)$

$$\text{Var}(y_i) = \sigma_\epsilon^2$$

$$\text{Var}(\hat{y}_i) = \text{Var}(\bar{y} + \hat{\beta}_1(x_i - \bar{x}))$$

$$= \text{Var}(\bar{y}) + (x_i - \bar{x})^2 \text{Var}(\hat{\beta}_1) + 2(x_i - \bar{x}) \text{Cov}(\bar{y}, \hat{\beta}_1)$$

$$\text{Cov}(\bar{y}, \hat{\beta}_1) = \text{Cov}(\bar{y}, \sum k_i y_i)$$

$$= \sum_i \text{Cov}(\bar{y}, k_i y_i)$$

$$= \sum_i \frac{k_i}{n} \text{Var}(y_i)$$

$$= \frac{1}{n} \sum k_i$$

$$= 0$$

$$\therefore \text{Var}(\hat{y}_i) = \text{Var}(\bar{y}) + (x_i - \bar{x})^2 \text{Var}(\hat{\beta}_1)$$

$$= \frac{1}{n} \sigma_\epsilon^2 + \frac{\sigma_\epsilon^2 (x_i - \bar{x})^2}{\sum (x_i - \bar{x})^2}$$

$$= \sigma_\epsilon^2 \left[\frac{1}{n} + \frac{(x_i - \bar{x})^2}{\sum (x_i - \bar{x})^2} \right]$$

Now, we derive the last term $\text{cov}(y_i, \hat{y}_i)$:

$$\begin{aligned}
\text{cov}(y_i, \hat{y}_i) &= \text{cov}(y_i, \bar{y} + \hat{\beta}_1(x_i - \bar{x})) \\
&= \text{cov}(y_i, \frac{1}{n} \sum_j y_j + (x_i - \bar{x}) \sum_j k_j y_j) \\
&= \text{cov}(y_i, \sum_j \left[\frac{1}{n} + (x_i - \bar{x}) k_j \right] y_j) \\
&= \sigma_\epsilon^2 \left[\frac{1}{n} + (x_i - \bar{x}) k_i \right] \\
&= \sigma_\epsilon^2 \left[\frac{1}{n} + \frac{(x_i - \bar{x})^2}{\sum (x_i - \bar{x})^2} \right]
\end{aligned}$$

Therefore, we have for i th residue

$$\begin{aligned}
\text{Var}(y_i - \hat{y}_i) &= \text{Var}(y_i) + \text{Var}(\hat{y}_i) - 2\text{cov}(y_i, \hat{y}_i) \\
&= \sigma_\epsilon^2 + \sigma_\epsilon^2 \left[\frac{1}{n} + \frac{(x_i - \bar{x})^2}{\sum (x_i - \bar{x})^2} \right] - 2\sigma_\epsilon^2 \left[\frac{1}{n} + \frac{(x_i - \bar{x})^2}{\sum (x_i - \bar{x})^2} \right] \\
&= \sigma_\epsilon^2 \left[1 - \frac{1}{n} - \frac{(x_i - \bar{x})^2}{\sum (x_i - \bar{x})^2} \right].
\end{aligned}$$

And finally, sum over i we get

$$\sum_i \text{Var}(y_i - \hat{y}_i) = \sigma_\epsilon^2 \sum_i \left[1 - \frac{1}{n} - \frac{(x_i - \bar{x})^2}{\sum (x_i - \bar{x})^2} \right] = (n - 2)\sigma_\epsilon^2$$

Confidence intervals

Now we substitute $\hat{\sigma}_\epsilon^2$ into the distribution of $\hat{\beta}_0$ and $\hat{\beta}_1$

$$\begin{aligned}
\hat{\beta}_1 &\sim N(\beta_1, \frac{\sigma_\epsilon^2}{\sum (x_i - \bar{x})^2}) \\
\hat{\beta}_0 &\sim N(\beta_0, \frac{\sigma_\epsilon^2 \sum x_i^2}{n \sum (x_i - \bar{x})^2})
\end{aligned}$$

to get the estimated standard errors:

$$\begin{aligned}
\widehat{SE}(\hat{\beta}_1) &= \sqrt{\frac{MS[E]}{\sum (x_i - \bar{x})^2}} \\
\widehat{SE}(\hat{\beta}_0) &= \sqrt{MS[E] \left(\frac{1}{n} + \frac{\bar{x}^2}{\sum (x_i - \bar{x})^2} \right)}
\end{aligned}$$

And the $100(1 - \alpha)\%$ confidence intervals for β_1 and β_0 are given by

$$\hat{\beta}_1 \pm t(n - 2, \alpha/2) \sqrt{\frac{MS[E]}{S_{xx}}}$$

$$\hat{\beta}_0 \pm t(n-2, \alpha/2) \sqrt{MS[E] \left(\frac{1}{n} + \frac{\bar{x}^2}{S_{xx}} \right)}$$

where $S_{xx} = \sum (x_i - \bar{x})^2$

Confidence interval for $\mathbf{E}(Y|X = x_0)$

The conditional mean $\mathbf{E}(Y|X = x_0)$ can be estimated by evaluating the regression function $\mu(x_0)$ at the estimates $\hat{\beta}_0, \hat{\beta}_1$. The conditional variance of the expression isn't too difficult (already shown):

$$\text{Var}(\hat{\beta}_0 + \hat{\beta}_1 x_0 | X = x_0) = \sigma^2 \left(\frac{1}{n} + \frac{(x_0 - \bar{x})^2}{S_{xx}} \right)$$

This leads to a confidence interval of the form

$$\hat{\beta}_0 + \hat{\beta}_1 x_0 \pm t(n-2, \alpha/2) \sqrt{MS[E] \left(\frac{1}{n} + \frac{(x_0 - \bar{x})^2}{S_{xx}} \right)}$$

Prediction interval

Often, prediction of the response variable Y for a given value, say x_0 , of the independent variable of interest. In order to make statements about future values of Y , we need to take into account

- the sampling distribution of $\hat{\beta}_0$ and $\hat{\beta}_1$
- the randomness of a future value Y .

We have seen the predicted value of Y based on the linear regression is given by $\hat{Y}_0 = \hat{\beta}_0 + \hat{\beta}_1 x_0$.

The 95% prediction interval has the form

$$\hat{Y}_0 \pm t(n-2, \alpha/2) \sqrt{MS[E] \left(1 + \frac{1}{n} + \frac{(x_0 - \bar{x})^2}{S_{xx}} \right)}.$$

Hypothesis test

To test the hypothesis $H_0 : \beta_1 = \beta_{slope}$ that the population slope is equal to a specific value β_{slope} (most commonly, the null hypothesis has $\beta_{slope} = 0$), we calculate the test statistic (T -statistics) with $df = n - 2$

$$t_0 = \frac{\hat{\beta}_1 - \beta_{slope}}{\widehat{SE}(\hat{\beta}_1)} \sim t_{n-2}$$

Some questions to answer using regression analysis:

1. What is the meaning, in words, of β_1 ?

Answer: β_1 is the population slope parameter of the SLR model that represents the

amount of increase in the mean of the response variable with a unit increase of the explanatory variable.

2. True/False: (a) β_1 is a statistic (b) β_1 is a parameter (c) β_1 is unknown.

Answer: (a) False (b) True (C) True. In reality, the true population parameters are almost never known. However, in simulation studies, we do know them.

3. True/False: (a) $\hat{\beta}_1$ is a statistic (b) $\hat{\beta}_1$ is a parameter (c) $\hat{\beta}_1$ is unknown

Answer: (a) True (b) False (C) False. $\hat{\beta}_1$ is an estimate of the population parameter β_1 .

4. Is $\hat{\beta}_1 = \beta_1$?

Answer: No. However, $\mathbf{E}(\hat{\beta}_1) = \beta_1$

Multiple linear regression

JF 5.2+6.2

Multiple linear regression - an example

An example on the prestige, education, and income levels of 45 U.S. occupations (Duncan's data):

	income	education	prestige
accountant	62	86	82
pilot	72	76	83
architect	75	92	90
author	55	90	76
chemist	64	86	90
minister	21	84	87
professor	64	93	93
dentist	80	100	90
reporter	67	87	52
engineer	72	86	88
lawyer	76	98	89
teacher	48	91	73

“prestige” represents the percentage of respondents in a survey who rated an occupation as “good” or “excellent” in prestige, “education” represents the percentage of incumbents in the occupation in the 1950 U.S. Census who were high school graduates, and “income” represents the percentage of occupational incumbents who earned incomes in excess of \$3,500.

Using the `pairs` command in R, we can look at the pairwise scatter plot between the three variables as in Figure 12.1.

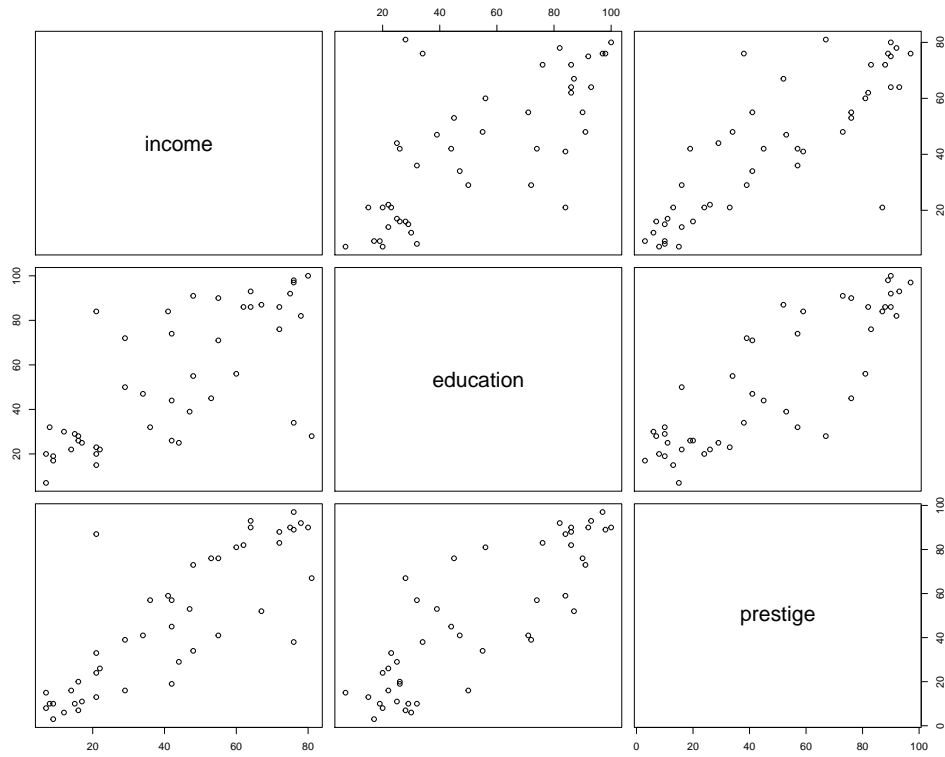


Figure 12.1: Scatterplot matrix for occupational prestige, level of education, and level of income of 45 U.S. occupations in 1950.

Consider a regression model for the “prestige” of occupation i , Y_i , in which the mean of Y_i is a linear function of two predictor variables $X_{i1} = \text{income}$, $X_{i2} = \text{education}$ for occupations $i = 1, 2, \dots, 45$:

$$Y = \beta_0 + \beta_1 \text{income} + \beta_2 \text{education} + \text{error}$$

or

$$Y_i = \beta_0 + \beta_1 X_{i1} + \beta_2 X_{i2} + \epsilon_i$$

or

$$Y_1 = \beta_0 + \beta_1 X_{11} + \beta_2 X_{12} + \epsilon_1$$

$$Y_2 = \beta_0 + \beta_1 X_{21} + \beta_2 X_{22} + \epsilon_2$$

$$\vdots = \vdots$$

$$Y_{45} = \beta_0 + \beta_1 X_{45,1} + \beta_2 X_{45,2} + \epsilon_{45}$$

A multiple linear regression (MLR) model w/ p independent variables

Let p independent variables be denoted by x_1, \dots, x_p .

- Observed values of p independent variables for i^{th} subject from sample denoted by x_{i1}, \dots, x_{ip}

- response variable for i^{th} subject denoted by Y_i
- For $i = 1, \dots, n$, MLR model for Y_i :

$$Y_i = \beta_0 + \beta_1 x_{i1} + \beta_2 x_{i2} + \dots + \beta_p x_{ip} + \epsilon_i$$

- As in SLR, $\epsilon_1, \dots, \epsilon_n \stackrel{iid}{\sim} N(0, \sigma^2)$

Least squares estimates of regression parameters minimize $SS[E]$:

$$SS[E] = \sum_{i=1}^n (y_i - \beta_0 - \beta_1 x_{i1} - \dots - \beta_p x_{ip})^2$$

$$\hat{\sigma}^2 = \frac{SS[E]}{n-p-1}$$

Interpretations of regression parameters:

- σ^2 is unknown error variance parameter
- $\beta_0, \beta_1, \dots, \beta_p$ are $p + 1$ unknown regression parameters:
 - β_0 : average response when $x_1 = x_2 = \dots = x_p = 0$
 - β_i is called a partial slope for x_i . Represents mean change in y per unit increase in x_i *with all other independent variables held fixed*.

12 Lecture 14 Feb 22

Last time

- Inference of SLR model

Today

- Confidence intervals for SLR
- Multiple linear regression

Confidence intervals

Now we substitute $\hat{\sigma}_\epsilon^2$ into the distribution of $\hat{\beta}_0$ and $\hat{\beta}_1$

$$\begin{aligned}\hat{\beta}_1 &\sim N(\beta_1, \frac{\sigma_\epsilon^2}{\sum(x_i - \bar{x})^2}) \\ \hat{\beta}_0 &\sim N(\beta_0, \frac{\sigma_\epsilon^2 \sum x_i^2}{n \sum(x_i - \bar{x})^2})\end{aligned}$$

to get the estimated standard errors:

$$\begin{aligned}\widehat{SE}(\hat{\beta}_1) &= \sqrt{\frac{MS[E]}{\sum(x_i - \bar{x})^2}} \\ \widehat{SE}(\hat{\beta}_0) &= \sqrt{MS[E] \left(\frac{1}{n} + \frac{\bar{x}^2}{\sum(x_i - \bar{x})^2} \right)}\end{aligned}$$

And the $100(1 - \alpha)\%$ confidence intervals for β_1 and β_0 are given by

$$\begin{aligned}\hat{\beta}_1 \pm t(n - 2, \alpha/2) \sqrt{\frac{MS[E]}{S_{xx}}} \\ \hat{\beta}_0 \pm t(n - 2, \alpha/2) \sqrt{MS[E] \left(\frac{1}{n} + \frac{\bar{x}^2}{S_{xx}} \right)}\end{aligned}$$

where $S_{xx} = \sum(x_i - \bar{x})^2$

Confidence interval for $\mathbf{E}(Y|X = x_0)$

The conditional mean $\mathbf{E}(Y|X = x_0)$ can be estimated by evaluating the regression function $\mu(x_0)$ at the estimates $\hat{\beta}_0, \hat{\beta}_1$. The conditional variance of the expression isn't too difficult (already shown):

$$\text{Var}(\hat{\beta}_0 + \hat{\beta}_1 x_0 | X = x_0) = \sigma^2 \left(\frac{1}{n} + \frac{(x_0 - \bar{x})^2}{S_{xx}} \right)$$

This leads to a confidence interval of the form

$$\hat{\beta}_0 + \hat{\beta}_1 x_0 \pm t(n - 2, \alpha/2) \sqrt{MS[E] \left(\frac{1}{n} + \frac{(x_0 - \bar{x})^2}{S_{xx}} \right)}$$

Prediction interval

Often, prediction of the response variable Y for a given value, say x_0 , of the independent variable of interest. In order to make statements about future values of Y , we need to take into account

- the sampling distribution of $\hat{\beta}_0$ and $\hat{\beta}_1$
- the randomness of a future value Y .

We have seen the predicted value of Y based on the linear regression is given by $\hat{Y}_0 = \hat{\beta}_0 + \hat{\beta}_1 x_0$.

The 95% prediction interval has the form

$$\hat{Y}_0 \pm t(n-2, \alpha/2) \sqrt{MS[E] \left(1 + \frac{1}{n} + \frac{(x_0 - \bar{x})^2}{S_{xx}} \right)}.$$

Hypothesis test

To test the hypothesis $H_0 : \beta_1 = \beta_{slope}$ that the population slope is equal to a specific value β_{slope} (most commonly, the null hypothesis has $\beta_{slope} = 0$), we calculate the test statistic (T -statistics) with $df = n - 2$

$$t_0 = \frac{\hat{\beta}_1 - \beta_{slope}}{\widehat{SE}(\hat{\beta}_1)} \sim t_{n-2}$$

Some questions to answer using regression analysis:

1. What is the meaning, in words, of β_1 ?
Answer: β_1 is the population slope parameter of the SLR model that represents the amount of increase in the mean of the response variable with a unit increase of the explanatory variable.
2. True/False: (a) β_1 is a statistic (b) β_1 is a parameter (c) β_1 is unknown.
Answer: (a) False (b) True (C) True. In reality, the true population parameters are almost never known. However, in simulation studies, we do know them.
3. True/False: (a) $\hat{\beta}_1$ is a statistic (b) $\hat{\beta}_1$ is a parameter (c) $\hat{\beta}_1$ is unknown
Answer: (a) True (b) False (C) False. $\hat{\beta}_1$ is an estimate of the population parameter β_1 .
4. Is $\hat{\beta}_1 = \beta_1$?
Answer: No. However, $\mathbf{E}(\hat{\beta}_1) = \beta_1$

Multiple linear regression

JF 5.2+6.2

Multiple linear regression - an example

An example on the prestige, education, and income levels of 45 U.S. occupations (Duncan's data):

	income	education	prestige
accountant	62	86	82
pilot	72	76	83
architect	75	92	90
author	55	90	76
chemist	64	86	90
minister	21	84	87
professor	64	93	93
dentist	80	100	90
reporter	67	87	52
engineer	72	86	88
lawyer	76	98	89
teacher	48	91	73

“prestige” represents the percentage of respondents in a survey who rated an occupation as “good” or “excellent” in prestige, “education” represents the percentage of incumbents in the occupation in the 1950 U.S. Census who were high school graduates, and “income” represents the percentage of occupational incumbents who earned incomes in excess of \$3,500.

Using the `pairs` command in R, we can look at the pairwise scatter plot between the three variables as in Figure 12.1.



Figure 12.1: Scatterplot matrix for occupational prestige, level of education, and level of income of 45 U.S. occupations in 1950.

Consider a regression model for the “prestige” of occupation i , Y_i , in which the mean of Y_i is a linear function of two predictor variables $X_{i1} = \text{income}$, $X_{i2} = \text{education}$ for occupations $i = 1, 2, \dots, 45$:

$$Y = \beta_0 + \beta_1 \text{income} + \beta_2 \text{education} + \text{error}$$

or

$$Y_i = \beta_0 + \beta_1 X_{i1} + \beta_2 X_{i2} + \epsilon_i$$

or

$$Y_1 = \beta_0 + \beta_1 X_{11} + \beta_2 X_{12} + \epsilon_1$$

$$Y_2 = \beta_0 + \beta_1 X_{21} + \beta_2 X_{22} + \epsilon_2$$

$$\vdots = \vdots$$

$$Y_{45} = \beta_0 + \beta_1 X_{45,1} + \beta_2 X_{45,2} + \epsilon_{45}$$

A multiple linear regression (MLR) model w/ p independent variables

Let p independent variables be denoted by x_1, \dots, x_p .

- Observed values of p independent variables for i^{th} subject from sample denoted by x_{i1}, \dots, x_{ip}

- response variable for i^{th} subject denoted by Y_i
- For $i = 1, \dots, n$, MLR model for Y_i :

$$Y_i = \beta_0 + \beta_1 x_{i1} + \beta_2 x_{i2} + \dots + \beta_p x_{ip} + \epsilon_i$$

- As in SLR, $\epsilon_1, \dots, \epsilon_n \stackrel{iid}{\sim} N(0, \sigma^2)$

Least squares estimates of regression parameters minimize $SS[E]$:

$$SS[E] = \sum_{i=1}^n (y_i - \beta_0 - \beta_1 x_{i1} - \dots - \beta_p x_{ip})^2$$

$$\hat{\sigma}^2 = \frac{SS[E]}{n-p-1}$$

Interpretations of regression parameters:

- σ^2 is unknown error variance parameter
- $\beta_0, \beta_1, \dots, \beta_p$ are $p + 1$ unknown regression parameters:
 - β_0 : average response when $x_1 = x_2 = \dots = x_p = 0$
 - β_i is called a partial slope for x_i . Represents mean change in y per unit increase in x_i *with all other independent variables held fixed*.

Matrix formulation of MLR

Let a vector for p observed independent variables for individual i be defined by

$$x_{i.} = (1, x_{i1}, x_{i2}, \dots, x_{ip}).$$

The MLR model for Y_1, \dots, Y_n is given by

$$\begin{aligned} Y_1 &= \beta_0 + \beta_1 X_{11} + \beta_2 X_{12} + \dots + \beta_p X_{1p} + \epsilon_1 \\ Y_2 &= \beta_0 + \beta_1 X_{21} + \beta_2 X_{22} + \dots + \beta_p X_{2p} + \epsilon_2 \\ &\vdots \\ Y_n &= \beta_0 + \beta_1 X_{n1} + \beta_2 X_{n2} + \dots + \beta_p X_{np} + \epsilon_n \end{aligned}$$

This system of n equations can be expressed using matrices:

$$\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\epsilon}$$

where

- \mathbf{Y} denotes a response vector of size $n \times 1$
- \mathbf{X} denotes a design matrix of size $n \times (p + 1)$
- $\boldsymbol{\beta}$ denotes a vector of regression parameters of size $(p + 1) \times 1$

- ϵ denotes an error vector of size $n \times 1$

Here, the error vector ϵ is assumed to follow a multivariate normal distribution with variance-covariance matrix $\sigma^2 \mathbf{I}_n$. For individual i ,

$$y_i = \mathbf{x}_i \boldsymbol{\beta} + \epsilon_i.$$

Some simplified expressions: (\mathbf{a} is a known $p \times 1$ vector)

$$\begin{aligned}\hat{\boldsymbol{\beta}} &= (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{Y} \\ \text{Var}(\hat{\boldsymbol{\beta}}) &= \sigma^2 (\mathbf{X}^T \mathbf{X})^{-1} \\ &= \boldsymbol{\Sigma} \\ \widehat{\text{Var}}(\hat{\boldsymbol{\beta}}) &= MS[E] (\mathbf{X}^T \mathbf{X})^{-1} \\ &= \hat{\boldsymbol{\Sigma}} \\ \widehat{\text{Var}}(\mathbf{a}^T \hat{\boldsymbol{\beta}}) &= \mathbf{a}^T \hat{\boldsymbol{\Sigma}} \mathbf{a}\end{aligned}$$

Question: what are the dimensions of each of these quantities?

- $(\mathbf{X}^T \mathbf{X})^{-1}$ may be verbalized as “x transposed x inverse”
- $\hat{\boldsymbol{\Sigma}}$ is the estimated variance-covariance matrix for the estimate of the regression parameter vector $\hat{\boldsymbol{\beta}}$
- \mathbf{X} is assumed to be of full *rank*.

Some more simplified expressions:

$$\begin{aligned}\hat{\mathbf{Y}} &= \mathbf{X} \hat{\boldsymbol{\beta}} \\ &= \mathbf{X} (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{Y} \\ &= \mathbf{H} \mathbf{Y} \\ \hat{\boldsymbol{\epsilon}} &= \mathbf{Y} - \hat{\mathbf{Y}} \\ &= \mathbf{Y} - \mathbf{X} \hat{\boldsymbol{\beta}} \\ &= (\mathbf{I} - \mathbf{H}) \mathbf{Y}\end{aligned}$$

- $\hat{\mathbf{Y}}$ is called the vector of fitted or predicted values
- $\mathbf{H} = \mathbf{X} (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T$ is called the hat matrix
- $\hat{\boldsymbol{\epsilon}}$ is the vector of residuals

For the Duncan's data example on income, education and prestige, with $p = 2$ independent variables and $n = 45$ observations,

$$\mathbf{X} = \begin{bmatrix} 1 & 62 & 86 \\ 1 & 72 & 76 \\ \vdots & \vdots & \vdots \\ 1 & 8 & 32 \end{bmatrix}$$

and

$$\mathbf{X}^T \mathbf{X} = \begin{bmatrix} 45 & 1884 & 2365 \\ 1884 & 105148 & 122197 \\ 2365 & 122197 & 163265 \end{bmatrix}$$

$$(\mathbf{X}^T \mathbf{X})^{-1} = \begin{bmatrix} 0.10211 & -0.00085 & -0.00084 \\ -0.00085 & 0.00008 & -0.00005 \\ -0.00084 & -0.00005 & 0.00005 \end{bmatrix}$$

$$(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{Y} = \begin{bmatrix} -6.0646629 \\ 0.5987328 \\ 0.5458339 \end{bmatrix} = ?$$

$$SS[E] = \boldsymbol{\epsilon}^T \boldsymbol{\epsilon} = (\mathbf{Y} - \hat{\mathbf{Y}})^T (\mathbf{Y} - \hat{\mathbf{Y}}) = 7506.7$$

$$MS[E] = \frac{SS[E]}{df} = \frac{7506.7}{45 - 2 - 1} = 178.73$$

$$\hat{\boldsymbol{\Sigma}} = MS[E](\mathbf{X}^T \mathbf{X})^{-1} = \begin{bmatrix} 18.249481 & -0.151845008 & -0.150706025 \\ -0.151845 & 0.014320275 & -0.008518551 \\ -0.150706 & -0.008518551 & 0.009653582 \end{bmatrix}$$

16 Lecture 16 Feb 27

Last time

- Confidence intervals for SLR
- Multiple linear regression

Today

- Multiple correlation
- More review
- Dummy variable regression

Multiple correlation, JF 5.2.3

The sums of squares in multiple regression are defined in the same manner as in SLR:

$$\begin{aligned}TSS &= \sum (Y_i - \bar{Y})^2 \\RegSS &= \sum (\hat{Y}_i - \bar{Y})^2 \\RSS &= \sum (Y_i - \hat{Y}_i)^2 = \sum \hat{\epsilon}_i^2\end{aligned}$$

Not surprisingly, we have a similar analysis of variance for the regression:

$$TSS = RegSS + RSS$$

The squared multiple correlation R^2 , representing the proportion of variation in the response variable captured by the regression, is defined in terms of the sums of squares:

$$R^2 = \frac{RegSS}{TSS} = 1 - \frac{RSS}{TSS}.$$

Because there are several slope coefficients, potentially with different signs, the *multiple correlation coefficient* is, by convention, the positive square root of R^2 . The multiple correlation is also interpretable as the simple correlation between the fitted and observed Y values, i.e. $r_{\hat{Y}Y}$.

Adjusted- R^2

Because the multiple correlation can only rise, never decline, when explanatory variables are added to the regression equation (HW1), investigators sometimes penalize the value of R^2 by a “correction” for degrees of freedom. The corrected (or “adjusted”) R^2 is defined as:

$$\begin{aligned}R_{adj}^2 &= 1 - \frac{\frac{RSS}{n-p-1}}{\frac{TSS}{n-1}} \\&= 1 - \left[\frac{(1 - R^2)(n - 1)}{n - p - 1} \right]\end{aligned}$$

Confidence intervals

Confidence intervals and hypothesis tests for individual coefficients closely follow the pattern of simple-regression analysis:

1. substitute an estimate of the error variance (MSE) for the unknown σ^2 into the variance term of $\hat{\beta}_i$
2. find the estimated standard error of a slope coefficient $\widehat{SE}(\hat{\beta}_i)$
3. $t = \frac{\hat{\beta}_i - \beta_i}{\widehat{SE}(\hat{\beta}_i)}$ follows a t -distribution with degrees of freedom as associated with SSE.

Therefore, we can construct the $100(1 - \alpha)\%$ confidence interval for a single slope parameter by (why?):

$$\hat{\beta}_i \pm t(n - p - 1, \alpha/2) \widehat{SE}(\hat{\beta}_i)$$

Hand-waving proof:

we know that $t = \frac{\hat{\beta}_i - \beta_i}{\widehat{SE}(\hat{\beta}_i)} \sim t_{n-p-1}$, such that

$$\begin{aligned} 1 - \alpha &= \Pr(-t_c < t < t_c) \\ &= \Pr\left(t_c < \frac{\hat{\beta}_i - \beta_i}{\widehat{SE}(\hat{\beta}_i)} < t_c\right) \\ &= \Pr\left(\hat{\beta}_i - t_c \cdot \widehat{SE}(\hat{\beta}_i) < \beta_i < \hat{\beta}_i + t_c \cdot \widehat{SE}(\hat{\beta}_i)\right) \end{aligned}$$

where $t_c = t(n - p - 1, \alpha/2)$ is the critical value.

Hypothesis tests

We first test the null hypothesis that all population regression slopes are 0:

$$H_0 : \beta_1 = \beta_2 = \dots = \beta_p = 0$$

The test statistics,

$$F = \frac{RegSS/p}{RSS/(n - p - 1)}$$

follows an F -distribution with p and $n - p - 1$ degrees of freedom.

We can also test a null hypothesis about a *subset* of the regression slopes, e.g.,

$$H_0 : \beta_1 = \beta_2 = \dots = \beta_q = 0.$$

Or more generally, test the null hypothesis

$$H_0 : \beta_{q_1} = \beta_{q_2} = \dots = \beta_{q_k} = 0$$

where $0 \leq q_1 < q_2 < \dots < q_k \leq p$ is a subset of k indices. To get the F -statistic for this case, we generally perform the following steps:

1. Fit the *full* (“unconstrained”) model, in other words, model that provides context for H_0 . Record SSR_{full} and the associated df_{full}
2. Fit the *reduced* (“constrained”) model, in other words, full model constrained by H_0 . Record SSR_{red} and the associated df_{red}
3. Calculate the F-statistic by

$$F = \frac{[SSR_{red} - SSR_{full}]/(df_{red} - df_{full})}{SSR_{full}/df_{full}}$$

4. Find p -value (the probability of observing an F-statistic that is at least as high as the value that we obtained) by consulting an F-distribution with numerator $df(ndf) = df_{red} - df_{full}$ and denominator $df(ddf) = df_{full}$. Notation: $F_{ndf,ddf}$, see Figure 16.1.

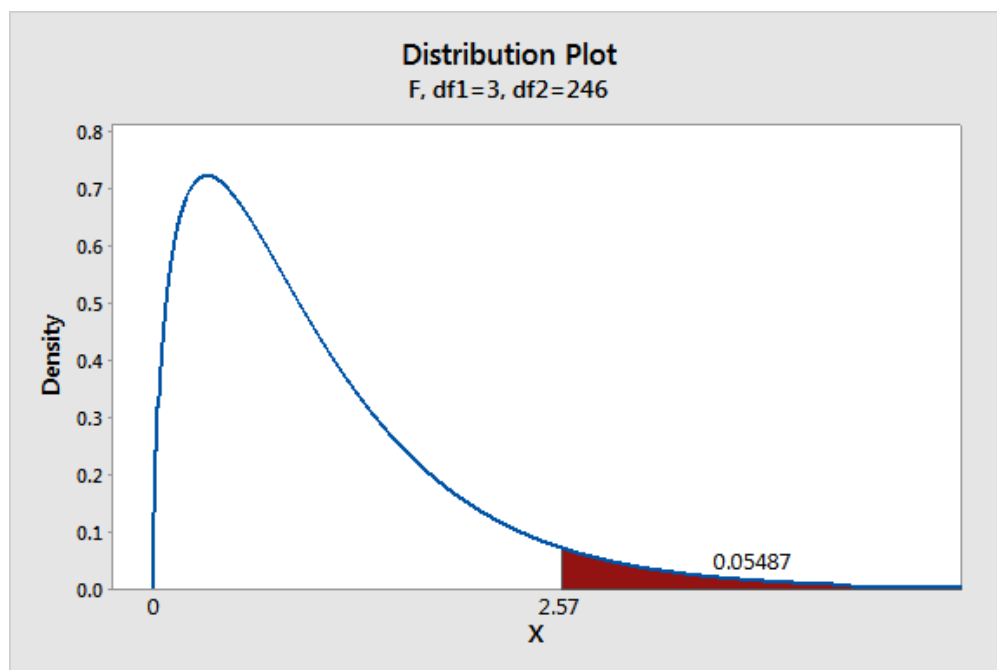


Figure 16.1: An example for p -value for F-statistic value 2.57 with an $F_{3,246}$ distribution

A little more background review

Reference:

- Statistical Inference, 2nd Edition, by George Casella & Roger L. Berger
- [Review of Probability Theory](#) by Arian Maleki and Tom Do

Chi-square, t-, and F-Distributions

Let $Z_1, Z_2, \dots, Z_k \stackrel{iid}{\sim} N(0, 1)$, then $X \equiv Z_1^2 + Z_2^2 + \dots + Z_k^2 \sim \chi_k^2$ (with k degrees of freedom).
If $X \sim \chi_k^2$

$$\begin{aligned}\mathbf{E}(X) &= k \\ \mathbf{Var}(X) &= 2k.\end{aligned}$$

Student's t versus χ^2

If $X_1, \dots, X_n \stackrel{iid}{\sim} N(\mu, \sigma^2)$, then

$$\frac{\bar{X} - \mu}{\sigma/\sqrt{n}} \sim N(0, 1).$$

When σ is unknown,

$$\frac{\bar{X} - \mu}{\hat{\sigma}/\sqrt{n}} \sim t_{n-1}, \quad \text{where } \hat{\sigma} = \sqrt{\frac{\sum (X_i - \bar{X})^2}{n-1}}.$$

Note that

$$\begin{aligned}\frac{\bar{X} - \mu}{\hat{\sigma}/\sqrt{n}} &= \frac{\bar{X} - \mu}{\sigma/\sqrt{n}} \cdot \frac{1}{\frac{\hat{\sigma}}{\sigma}} \\ &= Z \cdot \frac{1}{\sqrt{\frac{\sum (X_i - \bar{X})^2}{(n-1)\sigma^2}}} \\ &= \frac{Z}{\sqrt{\frac{\chi_{n-1}^2}{n-1}}}\end{aligned}$$

F versus χ^2

$$F_{ndf,ddf} \equiv \frac{\chi_{ndf}^2/ndf}{\chi_{ddf}^2/ddf}$$

t versus F

$$\begin{aligned}t_k &= \frac{Z}{\sqrt{\chi_k^2/k}} \\ &= \frac{\sqrt{\chi_1^2/1}}{\sqrt{\chi_k^2/k}} \\ &= \sqrt{F_{1,k}}\end{aligned}$$

or, in other words, $t_k^2 = F_{1,k}$

Random vectors and matrices

The cdf for random vector

$$\mathbf{Y} = \begin{bmatrix} Y_1 \\ Y_2 \\ \vdots \\ Y_n \end{bmatrix} \text{ is } F_{\mathbf{Y}}(\mathbf{y}) = \Pr(Y_1 \leq y_1, Y_2 \leq y_2, \dots, Y_n \leq y_n)$$

If a joint pdf exists, then $f_{\mathbf{Y}}(\mathbf{y}) = f_{\mathbf{Y}}(y_1, \dots, y_n)$ and

$$F_{\mathbf{Y}}(\mathbf{y}) = \int_{-\infty}^{y_1} \int_{-\infty}^{y_2} \dots \int_{-\infty}^{y_n} f_{\mathbf{Y}}(\mathbf{t}) d\mathbf{t}$$

Moments

$$\begin{aligned} \mathbf{E}(\mathbf{Y}) = \boldsymbol{\mu}_{\mathbf{Y}} &= \begin{bmatrix} E(Y_1) \\ E(Y_2) \\ \vdots \\ E(Y_n) \end{bmatrix} = \begin{bmatrix} \mu_1 \\ \mu_2 \\ \vdots \\ \mu_n \end{bmatrix} \\ \mathbf{Var}(\mathbf{Y}) &= \mathbf{E}((\mathbf{Y} - \boldsymbol{\mu}_{\mathbf{Y}})(\mathbf{Y} - \boldsymbol{\mu}_{\mathbf{Y}})^T) \\ &= \mathbf{E} \left(\begin{bmatrix} (Y_1 - \mu_1)^2 & (Y_1 - \mu_1)(Y_2 - \mu_2) & \dots \\ (Y_2 - \mu_2)(Y_1 - \mu_1) & (Y_2 - \mu_2)^2 & \dots \\ \dots & \dots & \dots \end{bmatrix} \right) \\ &= \mathbf{E}([(Y_i - \mu_i)(Y_j - \mu_j), i = 1, 2, \dots, n, j = 1, 2, \dots, n]) \\ &= (\sigma_{ij})_{i=1,2,\dots,n; j=1,2,\dots,n} \end{aligned}$$

where $\sigma_{ij} = Cov(Y_i, Y_j)$

Linear functions

Let $\mathbf{X} \in \mathbb{R}^{k \times 1}$, $\mathbf{Y} \in \mathbb{R}^{n \times 1}$ and $\mathbf{A} \in \mathbb{R}^{k \times 1}$, $\mathbf{B} \in \mathbb{R}^{k \times n}$ be non-random, then

$$\begin{aligned} \mathbf{X} &= \mathbf{A} + \mathbf{B} \mathbf{Y} \\ \mathbf{E}(\mathbf{X}) &= \mathbf{A} + \mathbf{B} \mathbf{E}(\mathbf{Y}) \\ \mathbf{Var}(\mathbf{X}) &= \mathbf{B} \mathbf{Var}(\mathbf{Y}) \mathbf{B}^T \end{aligned}$$

Sums of random vectors

$$\begin{aligned} \mathbf{X} &= \mathbf{Y} + \mathbf{Z} \\ \mathbf{E}(\mathbf{X}) &= \mathbf{E}(\mathbf{Y}) + \mathbf{E}(\mathbf{Z}) = \mathbf{E}(\mathbf{Y} + \mathbf{Z}) \end{aligned}$$

Note that there is no independence assumed above.

$$\mathbf{Var}(\mathbf{X}) = \mathbf{Var}(\mathbf{Y} + \mathbf{Z}) = \mathbf{Var}(\mathbf{Y}) + \mathbf{Var}(\mathbf{Z}) + Cov(\mathbf{Y}, \mathbf{Z}) + Cov(\mathbf{Z}, \mathbf{Y})$$

If \mathbf{Y}, \mathbf{Z} are uncorrelated, then $\mathbf{Var}(\mathbf{X}) = \mathbf{Var}(\mathbf{Y}) + \mathbf{Var}(\mathbf{Z})$

17 Lecture 17: March 1

Last time

- more review on probability

Today

- HW1 review
- Dummy-Variable regression
- Interactions

Dummy-variable regression

For categorical data (factor), we use dummy variable regression:

$$Y_i = \beta_0 + \beta_1 X_i + \beta_2 D_i + \epsilon_i$$

where D , called a dummy variable regressor or an indicator variable, is coded 1 for one level and 0 for all others,

$$D_i = \begin{cases} 1 & \text{for men} \\ 0 & \text{for women} \end{cases}.$$

Therefore, for women, the model becomes

$$Y_i = \beta_0 + \beta_1 X_i + \epsilon_i$$

and for men

$$Y_i = \beta_0 + \beta_1 X_i + \beta_2 + \epsilon_i = (\beta_0 + \beta_2) + \beta_1 X_i + \epsilon_i$$

For example, Figure 17.1 (a) and (b) represents two small (idealized) populations. In both cases, the within-gender regressions of income on education are parallel. Parallel regressions imply additive effects of education and gender on income: Holding education constant, the “effect” of gender is the vertical distance between the two regression lines, which, for parallel lines, is everywhere the same.



Figure 17.1: Idealized data representing the relationship between income and education for populations of men (filled circles) and women (open circles). In (a), there is no relationship between education and gender; in (b), women have a higher average level of education than men. In both (a) and (b), the within-gender (i.e., partial) regressions (solid lines) are parallel. In each graph, the overall (i.e. marginal) regression of income on education (ignoring gender) is given by the broken line. JF Figure 7.1.

Multi-level factor

We can model the effects of classification factors with m categories (levels) by using $m - 1$ indicator variables.

For example, the three-category occupational-type factor can be represented in the regression equation by introducing two dummy regressors:

Category	D_1	D_2
Professional and managerial	1	0
White collar	0	1
Blue collar	0	0

A model for the regression of prestige on income, education, and type of occupation is then

$$Y_i = \beta_0 + \beta_1 X_{i1} + \beta_2 X_{i2} + \gamma_1 D_{i1} + \gamma_2 D_{i2} + \epsilon_i$$

where X_1 is income and X_2 is education. This model describes three parallel regression planes, which can differ in their intercepts:

$$\begin{aligned} \text{Professional: } Y_i &= (\beta_0 + \gamma_1) + \beta_1 X_{i1} + \beta_2 X_{i2} + \epsilon_i \\ \text{White collar: } Y_i &= (\beta_0 + \gamma_2) + \beta_1 X_{i1} + \beta_2 X_{i2} + \epsilon_i \\ \text{Blue collar: } Y_i &= \beta_0 + \beta_1 X_{i1} + \beta_2 X_{i2} + \epsilon_i \end{aligned}$$

Therefore, the coefficient β_0 gives the intercept for blue-collar occupations; γ_1 represents the constant vertical difference between the parallel regression planes for professional and blue-collar occupations (fixing the values of education and income); and γ_2 represents the constant vertical distance between the regression planes for white-collar and blue-collar occupations (again, fixing education and income).

In the above prestige example, we chose “blue collar” as the baseline category. Sometimes, it is natural to pick a particular category as the baseline category, for example, the “control group” in an experiment. However, in most applications, the choice of a baseline category is entirely arbitrary.

Matrix representation

For the above prestige model

$$Y_i = \beta_0 + \beta_1 X_{i1} + \beta_2 X_{i2} + \gamma_1 D_{i1} + \gamma_2 D_{i2} + \epsilon_i$$

we have the design matrix \mathbf{X} as

$$\mathbf{X} = \begin{bmatrix} 1 & X_{11} & X_{12} & D_{11} & D_{12} \\ 1 & X_{21} & X_{22} & D_{21} & D_{22} \\ \vdots & \vdots & \vdots & \vdots & \vdots \\ 1 & X_{n1} & X_{n2} & D_{n1} & D_{n2} \end{bmatrix}$$

and the vector of coefficients β is

$$\beta = \begin{bmatrix} \beta_0 \\ \beta_1 \\ \beta_2 \\ \gamma_1 \\ \gamma_2 \end{bmatrix}$$

such that we have (again) the linear model in matrix form:

$$\mathbf{Y} = \mathbf{X}\beta + \epsilon$$

where $\epsilon_i \stackrel{iid}{\sim} N(0, \sigma^2)$, in other words, $\epsilon \sim N(\mathbf{0}, \sigma^2 \mathbf{I}_n)$.

Interactions

Two explanatory variables are said to interact in determining a response variable when the partial effect of one depends on the value of the other. Consider the hypothetical data shown in Figure 19.1.

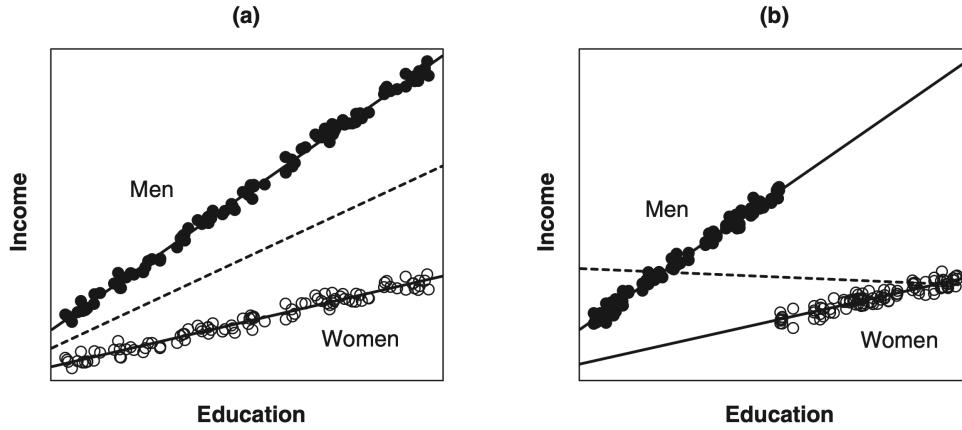


Figure 17.2: Idealized data representing the relationship between income and education for populations of men (filled circles) and women (open circles). In (a), there is no relationship between education and gender; in (b), women have a higher average level of education than men. In both (a) and (b), the within-gender (i.e., partial) regressions (solid lines) are not parallel. The slope for men is greater than the slope for women, and consequently education and gender interact in affecting income. In each graph, the overall regression of income on education (ignoring gender) is given by the broken line. JF Figure 7.7.

It is apparent in both Figure 19.1 (a) and (b) the within-gender regressions of income on education are not parallel: In both cases, the slope for men is larger than the slope for women.

Modeling interactions

We accommodate the interaction of education and gender by:

$$Y_i = \beta_0 + \beta_1 X_i + \beta_2 D_i + \beta_3 (X_i D_i) + \epsilon_i$$

where we introduce the interaction regressor XD into the regression equation. For women, the model becomes

$$\begin{aligned} Y_i &= \beta_0 + \beta_1 X_i + \beta_2 \cdot 0 + \beta_3 (X_i \cdot 0) + \epsilon_i \\ &= \beta_0 + \beta_1 X_i + \epsilon_i \end{aligned}$$

and for men

$$\begin{aligned} Y_i &= \beta_0 + \beta_1 X_i + \beta_2 \cdot 1 + \beta_3 (X_i \cdot 1) + \epsilon_i \\ &= (\beta_0 + \beta_2) + (\beta_1 + \beta_3) X_i + \epsilon_i \end{aligned}$$

The parameters β_0 and β_1 are, respectively, the intercept and slope for the regression of income on education among women (the baseline category for gender); β_2 gives the difference in intercepts between the male and female groups; and β_3 gives the difference in slopes between the two groups.

Usual guidance: Models that include an interaction between two predictors should also include the individual predictors by themselves regardless of the statistical significance of the associated β 's.

Test for the interaction

We can simply test the hypothesis $H_0 : \beta_3 = 0$ and construct the test statistic $t = \frac{\hat{\beta}_i - 0}{\widehat{SE}(\hat{\beta}_i)} \sim t_{n-4}$ ($p = 3$).

Interactions with multi-level factor

We can easily extend the method for modeling interactions by forming product regressors to multi-level factors, to several factors, and to several quantitative explanatory variables. Using the occupational prestige example, the occupational type could possibly interact both with income (X_1) and with education (X_2):

$$Y_i = \beta_0 + \beta_1 X_{i1} + \beta_2 X_{i2} + \gamma_1 D_{i1} + \gamma_2 D_{i2} \\ + \delta_{11} X_{i1} D_{i1} + \delta_{12} X_{i1} D_{i2} + \delta_{21} X_{i2} D_{i1} + \delta_{22} X_{i2} D_{i2} + \epsilon_i$$

The model therefore permits different intercepts and slopes for the three types of occupations:

$$\begin{array}{lll} \text{Professional:} & Y_i = & (\beta_0 + \gamma_1) + (\beta_1 + \delta_{11})X_{i1} + (\beta_2 + \delta_{21})X_{i2} + \epsilon_i \\ \text{White collar:} & Y_i = & (\beta_0 + \gamma_2) + (\beta_1 + \delta_{12})X_{i1} + (\beta_2 + \delta_{22})X_{i2} + \epsilon_i \\ \text{Blue collar:} & Y_i = & \beta_0 + \beta_1 X_{i1} + \beta_2 X_{i2} + \epsilon_i \end{array}$$

19 Lecture 19: March 6

Last time

- HW1 review
- Dummy-Variable regression

Today

- Interactions
- Git tag demo
- Lab session review

Interactions

Two explanatory variables are said to interact in determining a response variable when the partial effect of one depends on the value of the other. Consider the hypothetical data shown in Figure 19.1.

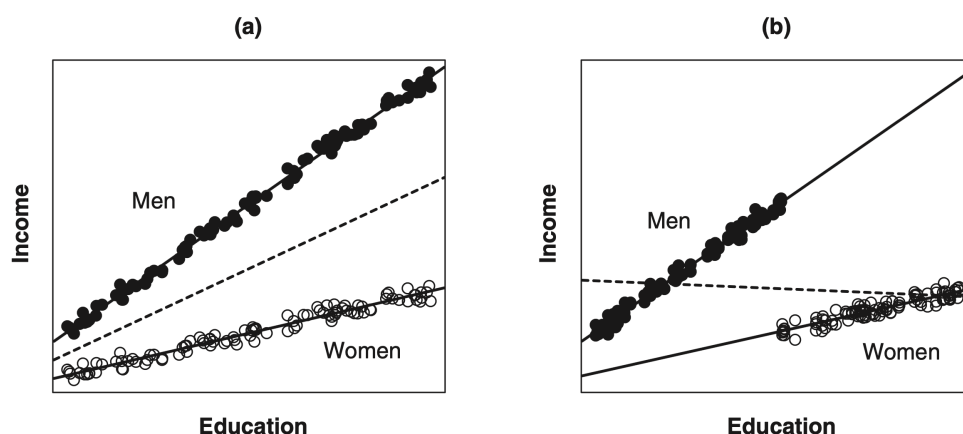


Figure 19.1: Idealized data representing the relationship between income and education for populations of men (filled circles) and women (open circles). In (a), there is no relationship between education and gender; in (b), women have a higher average level of education than men. In both (a) and (b), the within-gender (i.e., partial) regressions (solid lines) are not parallel. The slope for men is greater than the slope for women, and consequently education and gender interact in affecting income. In each graph, the overall regression of income on education (ignoring gender) is given by the broken line. JF Figure 7.7.

It is apparent in both Figure 19.1 (a) and (b) the within-gender regressions of income on education are not parallel: In both cases, the slope for men is larger than the slope for women.

Modeling interactions

We accommodate the interaction of education and gender by:

$$Y_i = \beta_0 + \beta_1 X_i + \beta_2 D_i + \beta_3 (X_i D_i) + \epsilon_i$$

where we introduce the interaction regressor XD into the regression equation. For women, the model becomes

$$\begin{aligned} Y_i &= \beta_0 + \beta_1 X_i + \beta_2 \cdot 0 + \beta_3 (X_i \cdot 0) + \epsilon_i \\ &= \beta_0 + \beta_1 X_i + \epsilon_i \end{aligned}$$

and for men

$$\begin{aligned} Y_i &= \beta_0 + \beta_1 X_i + \beta_2 \cdot 1 + \beta_3 (X_i \cdot 1) + \epsilon_i \\ &= (\beta_0 + \beta_2) + (\beta_1 + \beta_3) X_i + \epsilon_i \end{aligned}$$

The parameters β_0 and β_1 are, respectively, the intercept and slope for the regression of income on education among women (the baseline category for gender); β_2 gives the difference in intercepts between the male and female groups; and β_3 gives the difference in slopes between the two groups.

Usual guidance: Models that include an interaction between two predictors should also include the individual predictors by themselves regardless of the statistical significance of the associated β 's.

Test for the interaction

We can simply test the hypothesis $H_0 : \beta_3 = 0$ and construct the test statistic $t = \frac{\hat{\beta}_3 - 0}{\widehat{SE}(\hat{\beta}_3)} \sim t_{n-4}$ ($p = 3$).

Interactions with multi-level factor

We can easily extend the method for modeling interactions by forming product regressors to multi-level factors, to several factors, and to several quantitative explanatory variables. Using the occupational prestige example, the occupational type could possibly interact both with income (X_1) and with education (X_2):

$$\begin{aligned} Y_i &= \beta_0 + \beta_1 X_{i1} + \beta_2 X_{i2} + \gamma_1 D_{i1} + \gamma_2 D_{i2} \\ &\quad + \delta_{11} X_{i1} D_{i1} + \delta_{12} X_{i1} D_{i2} + \delta_{21} X_{i2} D_{i1} + \delta_{22} X_{i2} D_{i2} + \epsilon_i \end{aligned}$$

The model therefore permits different intercepts and slopes for the three types of occupations:

$$\begin{array}{lll} \text{Professional:} & Y_i = & (\beta_0 + \gamma_1) + (\beta_1 + \delta_{11})X_{i1} + (\beta_2 + \delta_{21})X_{i2} + \epsilon_i \\ \text{White collar:} & Y_i = & (\beta_0 + \gamma_2) + (\beta_1 + \delta_{12})X_{i1} + (\beta_2 + \delta_{22})X_{i2} + \epsilon_i \\ \text{Blue collar:} & Y_i = & \beta_0 + \beta_1 X_{i1} + \beta_2 X_{i2} + \epsilon_i \end{array}$$

20 Lecture 20: March 13

Last time

- Dummy-Variable regression (JF chapter 7)
- Interactions

Today

- Anonymous internal midterm evaluations on canvas
- Unusual and influential data (JF chapter 11)

Unusual and influential data

Linear models make strong assumptions about the structure of data, assumptions that often do not hold in applications. The method of least squares can be very sensitive to the structure of the data and may be markedly influenced by one or a few unusual observations.

Outliers

In simple regression analysis, an outlier is an observation whose response-variable value is *conditionally* unusual *given* the value of the explanatory variable: see Figure 20.1.

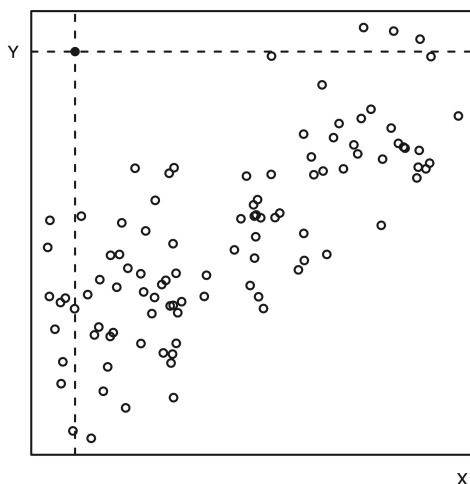


Figure 20.1: The black point is a regression outlier because it combines a relatively large value of Y with a relatively small value of X , even though neither its X -value nor its Y -value is unusual individually. Because of the positive relationship between Y and X , points with small X -values also tend to have small Y -values, and thus the black point is far from other points with similar X -values. JF Figure 11.1.

Unusual data are problematic in linear models fit by least squares because they can unduly

influence the results of the analysis. Their presence may be a signal that the model fails to capture important characteristics of the data.

Figure 20.2 illustrates some distinctions for the simple-regression model $Y = \beta_0 + \beta_1 X + \epsilon$.

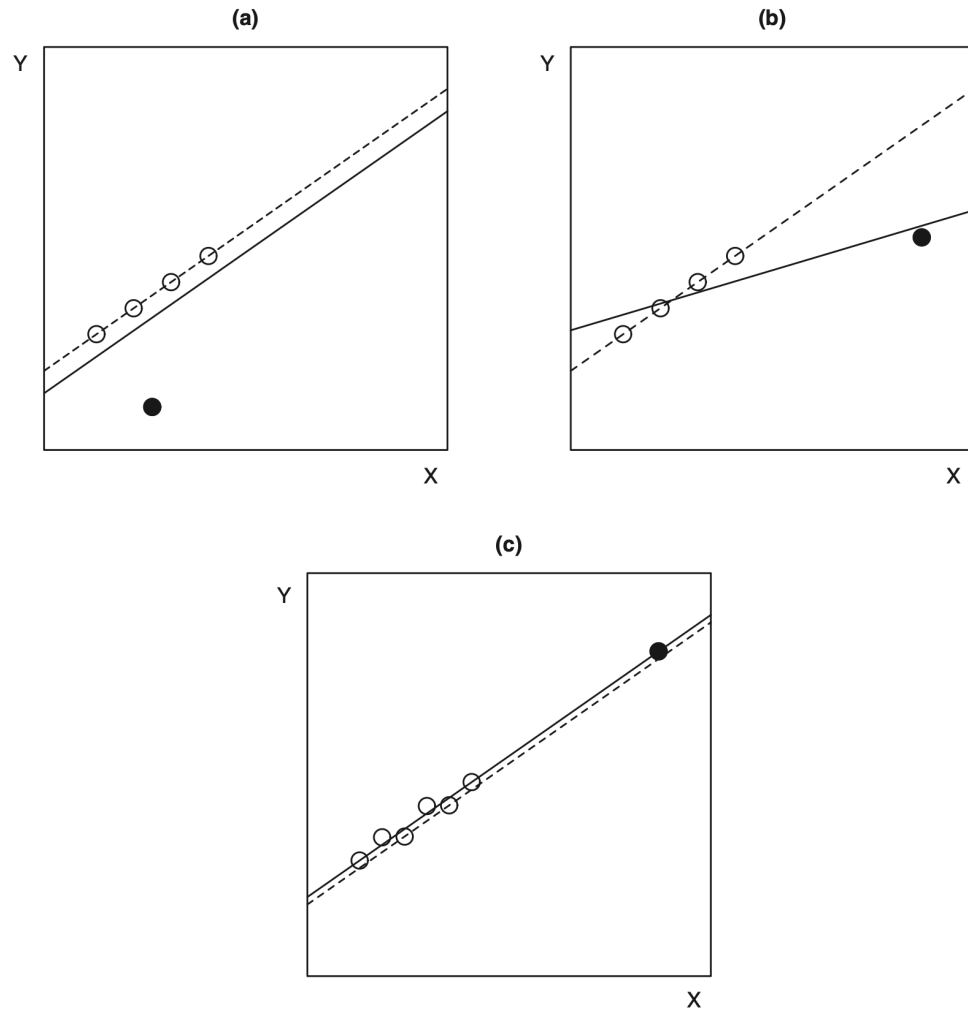


Figure 20.2: Leverage and influence in simple regression. In each graph, the solid line gives the least-squares regression for all the data, while the broken line gives the least-squares regression with the unusual data point (the black circle) omitted. (a) An outlier near the mean of X has low leverage and little influence on the regression coefficients. (b) An outlier far from the mean of X has high leverage and substantial influence on the regression coefficients. (c) A high-leverage observation in line with the rest of the data does not influence the regression coefficients. In panel (c), the two regression lines are separated slightly for visual effect but are, in fact, coincident JF Figure 11.2.

Some qualitative distinctions between outliers and high leverage observations:

- An outlier is a data point whose response Y does not follow the general trend of the rest of the data.

- A data point has high leverage if it has “extreme” predictor X values:
 - With a single predictor, an extreme X value is simply one that is particularly high or low.
 - With multiple predictors, extreme X values may be particularly high or low for one or more predictors, or may be “unusual” combinations of predictor values .

And the influence of a data point is the combination of leverage and discrepancy (“outlyingness”) though the following heuristic formula:

$$\text{Influence on coefficients} = \text{Leverage} \times \text{Discrepancy}.$$

Assessing leverage: hat-values

The hat-value h_i is a common measure of leverage in regression. They are named because it is possible to express the fitted values \hat{Y}_j (“Y-hat”) in terms of the observed values Y_i :

$$\hat{Y}_j = h_{1j}Y_1 + h_{2j}Y_2 + \cdots + h_{jj}Y_j + \cdots + h_{nj}Y_n = \sum_{i=1}^n h_{ij}Y_i.$$

The weight h_{ij} captures the contribution of observation Y_i to the fitted value \hat{Y}_j : If h_{ij} is large, then the i th observation can have a considerable impact on the j th fitted value. With the least square solutions, for the fitted values:

$$\hat{\mathbf{Y}} = \mathbf{X}\boldsymbol{\beta} = \mathbf{X}(\mathbf{X}^T\mathbf{X})^{-1}\mathbf{X}^T\mathbf{Y}$$

we (already) get the hat matrix:

$$\mathbf{H} = \mathbf{X}(\mathbf{X}^T\mathbf{X})^{-1}\mathbf{X}^T$$

Properties:

- (idempotent) $\mathbf{H} = \mathbf{H}\mathbf{H}$
- $h_i \equiv h_{ii} = \sum_{j=1}^n h_{ij}^2$
- $\frac{1}{n} \leq h_i \leq 1$ ([a proof](#) by Mohammad Mohammadi)
- $\bar{h} = (p+1)/n$

In the case of SLR, the hat-values are:

$$h_i = \frac{1}{n} + \frac{(X_i - \bar{X})^2}{\sum_{j=1}^n (X_j - \bar{X})^2}$$

Detecting outliers: studentized residuals

The variance of the residuals ($\hat{\epsilon}_i = Y_i - \hat{Y}_i$) do not have equal variances (even if the errors ϵ_i have equal variances):

$$\text{Var}(\hat{\epsilon}) = \text{Var}(\mathbf{Y} - \mathbf{X}\hat{\boldsymbol{\beta}}) = \text{Var}[(\mathbf{I} - \mathbf{H})\mathbf{Y}] = (\mathbf{I} - \mathbf{H})\text{Var}(\mathbf{Y})(\mathbf{I} - \mathbf{H}) = \sigma^2(\mathbf{I} - \mathbf{H})$$

so that for $\hat{\epsilon}_i$,

$$\text{Var}(\hat{\epsilon}_i) = \sigma^2(1 - h_i).$$

High-leverage observations tend to have small residuals (in other words, these observations can pull the regression surface toward them).

The standardized residual (sometimes called internally studentized residual)

$$\hat{\epsilon}'_i \equiv \frac{\hat{\epsilon}_i}{\hat{\sigma}\sqrt{1 - h_i}},$$

however, does not follow a t -distribution, because the numerator and denominator are not independent.

Suppose, we refit the model deleting the i th observation, obtaining an estimate $\hat{\sigma}_{(-i)}$ of σ that is based on the remaining $n - 1$ observations. Then the studentized residual (sometimes called externally studentized residual)

$$\hat{\epsilon}^*_i \equiv \frac{\hat{\epsilon}_i}{\hat{\sigma}_{(-i)}\sqrt{1 - h_i}}$$

has an independent numerator and denominator and follows a t -distribution with $n - p - 2$ degrees of freedom.

The studentized and the standardized residuals have the following relationship (Beckman and Trussell, 1974):

$$\hat{\epsilon}^*_i = \hat{\epsilon}'_i \sqrt{\frac{n - p - 2}{n - p - 1 - \hat{\epsilon}'_i{}^2}}$$

For large n ,

$$\hat{\epsilon}^*_i \approx \hat{\epsilon}'_i \approx \frac{\hat{\epsilon}}{\hat{\sigma}}$$

Test for outlier

It is of our interest to pick the studentized residual $\hat{\epsilon}^*_{max}$ with the largest absolute value among $\hat{\epsilon}^*_1, \hat{\epsilon}^*_2, \dots, \hat{\epsilon}^*_n$ to test for outlier. However, by doing so, we are effectively picking the biggest of n test statistics such that it is not legitimate simply to use t_{n-p-2} to find a p -value. We need a correction on the p -value because of multiple-comparisons.

Suppose that we have $p' = \Pr(t_{n-p-2} > |\hat{\epsilon}^*_{max}|)$, the p -value before correction. Then the Bonferroni adjusted p -value is $p = np'$.

Measuring influence

Influence on the regression coefficients combines leverage and discrepancy. The most direct measure of influence simply expresses the impact on each coefficient of deleting each observation in turn:

$$D_{ij} = \hat{\beta}_j - \tilde{\beta}_{j(-i)} \quad \text{for } i = 1, \dots, n \text{ and } j = 0, 1, \dots, p$$

where $\hat{\beta}_j$ are the least-squares coefficients calculated for all the data, and the $\tilde{\beta}_{j(-i)}$ are the least-squares coefficients calculated with the i th observation omitted. To assist in interpretation, it is useful to scale the D_{ij} by (deleted) coefficient standard errors:

$$D_{ij}^* = \frac{D_{ij}}{\widehat{SE}_{(-i)}(\tilde{\beta}_{j(-i)})}$$

Following Belsley, Kuh, and Welsh (1980), the D_{ij} are often termed DFBETA_{ij} , and D_{ij}^* are called DFBETAS_{ij} . One problem associated with using D_{ij} or D_{ij}^* is their large number: $n(p+1)$ of each.

Cook's distance calculated as

$$D_i = \frac{\sum_{j=1}^n (\tilde{y}_{j(-i)} - \hat{y}_j)^2}{(p+1)\hat{\sigma}^2} = \frac{\hat{\epsilon}_i'^2}{p+1} \times \frac{h_i}{1-h_i}$$

In effect, the first term in the formula for Cook's D is a measure of discrepancy, and the second is a measure of leverage. We look for values of D_i that stand out from the rest.

A similar measure suggested by Belsley et al. (1980)

$$\text{DFFITS}_i = \hat{\epsilon}_i^* \frac{h_i}{1-h_i}$$

Except for unusual data configurations, Cook's $D_i \approx \text{DFFITS}_i^2/(p+1)$.

Numerical cutoffs (suggested)

Diagnostic statistic	Cutoff value
h_i	$2\bar{h} = \frac{2(p+1)}{n}$, ($3\bar{h}$ for small sample)
D_{ij}^*	$ D_{ij}^* > 1$ or 2 ($2/\sqrt{n}$ for large samples)
Cook's D_i	$D_i > \frac{4}{n-p-1}$
DFFITS	$ \text{DFFITS}_i > 2\sqrt{\frac{p+1}{n-p-1}}$

21 Lecture 21: March 20

Last time

- Unusual and influential data (JF chapter 11)
- Midterm review
- Influence (JF chapter 11)

Today

- Added-variable plots
- Should unusual data be discarded
- Diagnosing non-normality, non-constant error variance, and nonlinearity (JF chapter 12)

Added-variable plots

Unlike the case of SLR, the scatterplot with the response variable and one predictor gives only the marginal effect in MLR. Instead, the added-variable plot (also called a partial-regression plot or a partial-regression leverage plot) gives a graphical inspection over each dimension.

Let $\tilde{Y}_i^{(1)}$ represent the residuals from the least-squares regression of Y on all the X s except X_1 , in other words, the residuals from the following fitted regression equation:

$$Y_i = \tilde{\beta}_0^{(1)} + \tilde{\beta}_2^{(1)} X_{i2} + \cdots + \tilde{\beta}_p^{(1)} X_{ip} + \tilde{Y}_i^{(1)}$$

where the parenthetical superscript (1) indicates the omission of X_1 from the right-hand side of the regression equation. Likewise, $\check{X}_i^{(1)}$ is the residual from the least-squares regression of X_1 on all the other X s:

$$X_{i1} = \check{\beta}_0^{(1)} + \check{\beta}_2^{(1)} X_{i2} + \cdots + \check{\beta}_p^{(1)} X_{ip} + \check{X}_i^{(1)}$$

Then, the residuals $\tilde{Y}_i^{(1)}$ and $\check{X}_i^{(1)}$ have the following interesting properties:

1. The slope from the least-squares regression of $\tilde{Y}_i^{(1)}$ on $\check{X}_i^{(1)}$ is simply the least-squares slope $\hat{\beta}_1$ from the *full* multiple regression.
2. The residuals from the simple regression of $\tilde{Y}_i^{(1)}$ on $\check{X}_i^{(1)}$ are the same as those from the full regression, that is

$$\tilde{Y}_i^{(1)} = \hat{\beta}_1 \check{X}_i^{(1)} + \hat{\epsilon}_i$$

3. The variation of $\check{X}_i^{(1)}$ is the *conditional variation* of X_1 holding the other X s constant.

Figure 21.1 shows that the conditional variation is smaller than its marginal variation – much smaller when X_1 is strongly collinear with other X s,

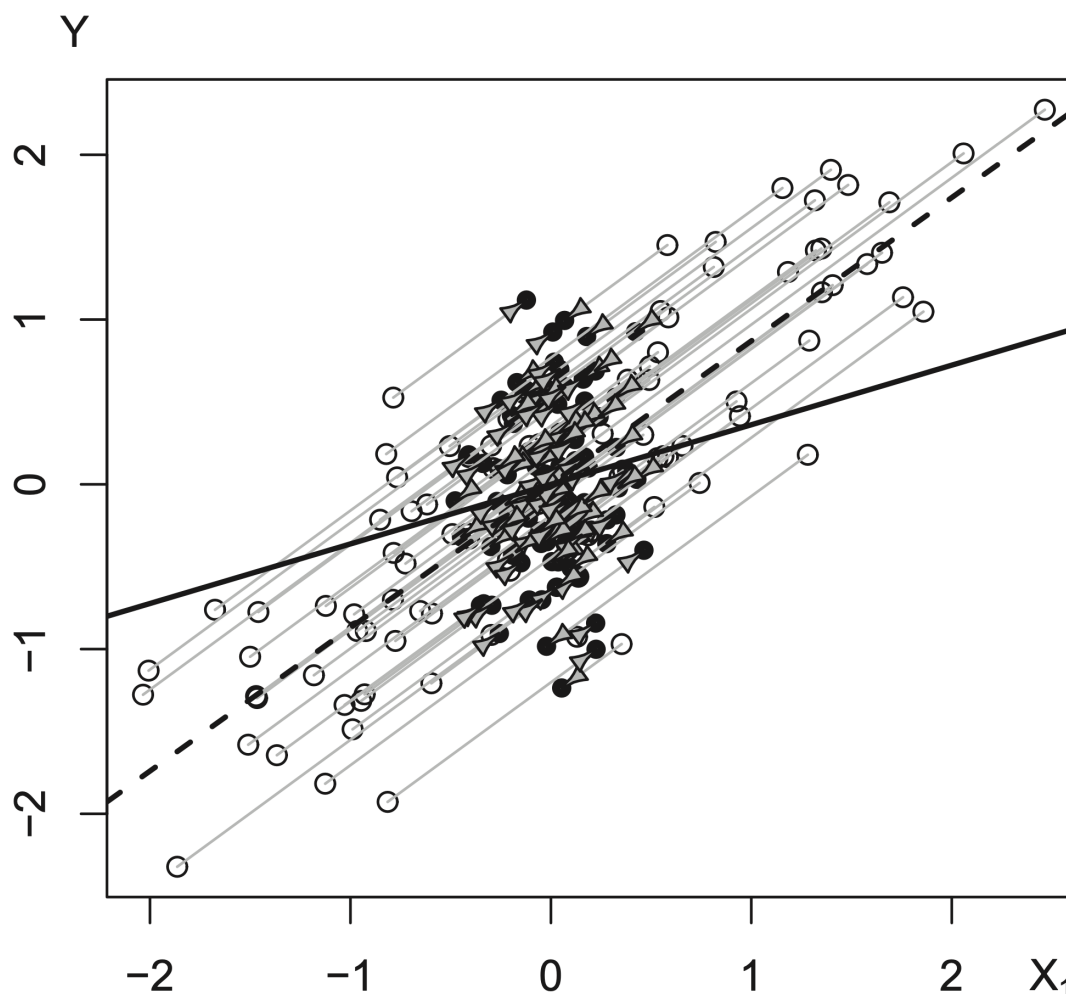


Figure 21.1: The marginal scatterplot (open circles) for Y and X_1 superimposed on the added-variable plot (filled circles) for X_1 in the regression of Y on X_1 and X_2 . The variables Y and X_1 are centered at their means to facilitate the comparison of the two sets of points. The arrows show how the points in the marginal scatterplot map into those in the AV plot. In this contrived data set, X_1 and X_2 are highly correlated ($r_{12} = 0.98$), and so the conditional variation in X_1 (represented by the horizontal spread of the filled points) is much less than its marginal variation (represented by the horizontal spread of the open points). The broken line gives the slope of the marginal regression of Y on X_1 alone, while the solid line gives the slope $\hat{\beta}_1$ of X_1 in the MLR of Y on both X s. JF Figure 11.9.

Figure 21.2 illustrates the added-variable plots using the Duncan's data.

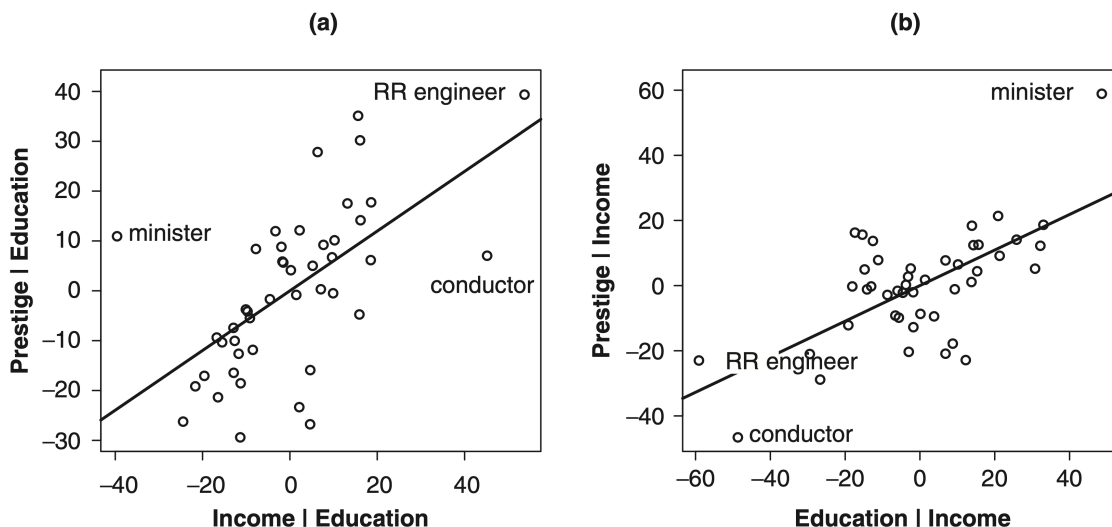


Figure 21.2: Added-variable plots for Duncan’s regression of occupational prestige on the (a) income and (b) education levels of 45 US occupations in 1950. Three unusual observations, *ministers*, *conductors*, and *railroadengineers*, are identified on the plots. The added-variable plot for the intercept $\hat{\beta}_0$ is not shown. JF Figure 11.10.

The added-variable plot for income in Figure 21.2(a) reveals three observations that exert substantial leverage on the income coefficient:

- *minister*, whose income is unusually low given the educational level of the occupation
- *conductor*, whose income is unusually high given education
- *railroad engineer*, whose income is relatively high given education.

Remember that the horizontal variable in this added-variable plot is the residual from the regression of income on education, and thus values far from 0 in this direction are for occupations with incomes that are unusually high or low given their levels of education.

Should unusual data be discarded?

In practice, although problematic data should not be ignored, they also should not be deleted automatically and without reflection:

- It is important to investigate *why* an observation is unusual. Truly “bad” data (e.g., an error in data entry) can often be corrected or, if correction is not possible, thrown away. When a discrepant data point is correct, we may be able to understand why the observation is unusual. For Duncan’s data, for example, it makes sense that ministers enjoy prestige not accounted for by the income and educational levels of the occupation and for a reason not shared by other occupations. In a case like this, where an outlying observation has characteristics that render it unique, we may choose to set it aside from the rest of the data.

- Alternatively, outliers, high-leverage points, or influential data may motivate model respecification, and the pattern of unusual data may suggest the introduction of additional explanatory variables. We noticed, for example, that both conductors and railroad engineers had high leverage in Duncan's regression because these occupations combined relatively high income with relatively low education. Perhaps this combination of characteristics is due to a high level of unionization of these occupations in 1950, when the data were collected. If so, and if we can ascertain the levels of unionization of all of the occupations, we could enter this as an explanatory variable, perhaps shedding further light on the process determining occupational prestige.
- Except in clear-cut cases, we are justifiably reluctant to delete observations or to re-specify the model to accommodate unusual data. Some researchers reasonably adopt alternative estimation strategies, such as robust regression, which continuously down-weights outlying data rather than simply discarding them. Because these methods assign zero or very small weight to highly discrepant data, however, the result is generally not very different from careful application of least squares, and, indeed, robust-regression weights can be used to identify outliers.
- Finally, in large samples, unusual data substantially alter the results only in extreme instances. Identifying unusual observations in a large sample, therefore, should be regarded more as an opportunity to learn something about the data not captured by the model that we have fit, rather than as an occasion to reestimate the model with the unusual observations removed.

Non-normally distributed errors

The assumption of normally distributed errors is almost always arbitrary. Nevertheless, the central limit theorem ensures that, under very broad conditions, inference based on the least-squares estimator is approximately valid in all but small samples. Why concern about non-normal errors?

- For some types of error distributions, particularly those with heavy tails, the efficiency of least-squares estimation decreases markedly.
- Highly skewed error distributions, aside from their propensity to generate outliers in the direction of the skew, compromise the interpretation of the least-squares fit. This fit is a conditional mean (of Y given the X s), and the mean is not a good measure of the center of a highly skewed distribution.
- A multimodal error distribution suggests that omission of one or more discrete explanatory variables that divide the data naturally into groups. An examination of the distribution of the residuals may motivate respecification of the model.

Note: The skewness α_3 is defined as $\alpha_3 \equiv \frac{\mu_3}{(\mu_2)^{3/2}}$ where μ_n denotes the n th central moment of a random variable X . The skewness measures the lack of symmetry in the pdf.

Quantile-comparison plot, JF 3.1.3

Quantile-comparison plots are useful for comparing an empirical sample distribution with a theoretical distribution, such as the normal distribution.

Let $P(x)$ represent the theoretical cumulative distribution function (cdf) with which we want to compare the data, that is $P(x) = \Pr(X \leq x)$. The quantile-comparison plot is constructed by:

1. Order the data values from smallest to largest, $X_{(1)}, X_{(2)}, \dots, X_{(n)}$. The $X_{(i)}$ are called the order statistics of the sample.
2. By convention, the cumulative proportion of the data “below” $X_{(i)}$ is given by

$$P_i = \frac{i - \frac{1}{2}}{n}$$

3. Use the inverse of the cdf to find the value z_i corresponding to the cumulative probability P_i , that is

$$z_i = P^{-1}\left(\frac{i - \frac{1}{2}}{n}\right)$$

4. Plot the z_i as horizontal coordinates against the $X_{(i)}$ as vertical coordinates. If X is sampled from the distribution P , then $X_{(i)} \approx z_i$.
 - if the distributions are identical except for location, then the plot is approximately linear with nonzero intercept, $X_{(i)} \approx \mu + z_i$
 - if the distributions are identical except for scale, then the plot is approximately linear with a slope different from 1, $X_{(i)} \approx \sigma z_i$
 - if the distributions differ both in location and scale but have the same shape, then $X_{(i)} \approx \mu + \sigma z_i$
5. It is often helpful to place a comparison line on the plot to facilitate the perception of departures from linearity. For a normal quantile-comparison plot (comparing the distribution of the data with the standard normal distribution), we can alternatively use the median as a robust estimator of μ and the interquartile range/1.39 as a robust estimator of σ .
6. We expect some departure from linearity because of sampling variation. It therefore assists interpretation to display the expected degree of sampling error in the plot. The standard error of the order statistic $X_{(i)}$ is

$$\text{SE}(X_{(i)}) = \frac{\hat{\sigma}}{p(z_i)} \sqrt{\frac{P_i(1 - P_i)}{n}}$$

where $p(z_i)$ is the probability density function, pdf, corresponding to the CDF $P(z)$. The values along the fitted line are given by $\hat{X}_{(i)} = \hat{\mu} + \hat{\sigma} z_i$. An approximate 95% confidence “envelope” around the fitted line is, therefore,

$$\hat{X}_{(i)} \pm 2 \times \text{SE}(X_{(i)})$$

- Figure 21.1 plots a sample of $n = 100$ observations from a normal distribution with mean $\mu = 50$ and standard deviation $\sigma = 10$.

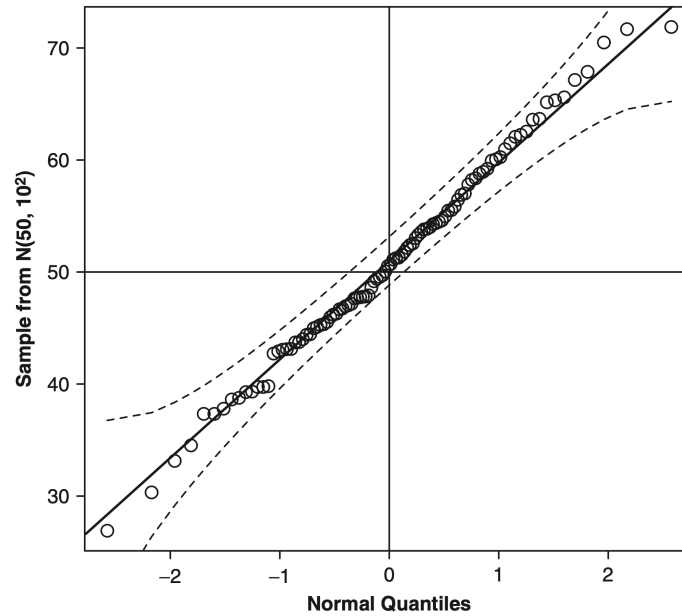


Figure 21.3: Normal quantile-comparison plot for a sample of 100 observations drawn from a normal distribution with mean 50 and standard deviation 10. The fitted line is through the quantiles of the distribution, the broken lines give a pointwise 95% confidence interval around the fit. JF Figure 3.8.

The plotted points are reasonably linear and stay within the rough 95% confidence envelope.

- Figure 21.2 plots a sample of $n = 100$ observations from the positively skewed chi-square distribution with 2 degrees of freedom. The positive skew of the data is reflected in points that lie *above* the comparison line in both tails of the distribution. (In contrast, the tails of negatively skewed data would lie *below* the comparison line.)

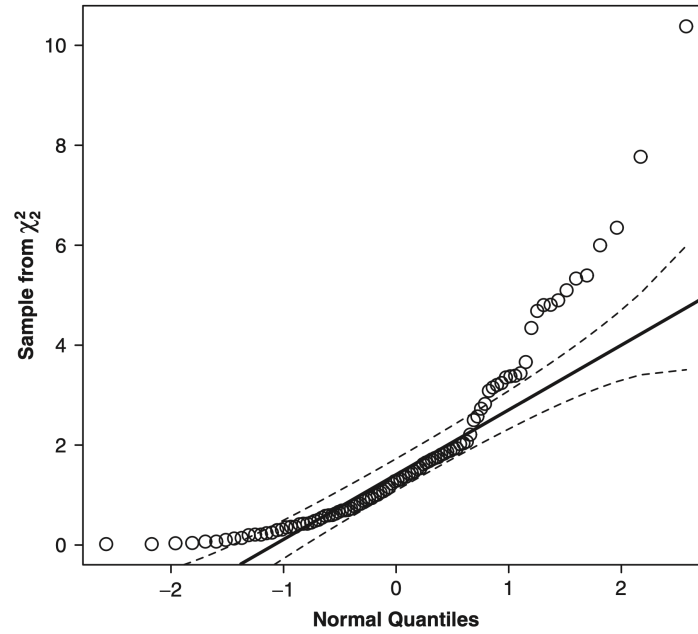


Figure 21.4: Normal quantile-comparison plot for a sample of 100 observations drawn from the positively skewed chi-square distribution with 2 degrees of freedom. JF Figure 3.9.

- Figure 21.3 plots a sample of $n = 100$ observations from the heavy-tailed t distribution with 2 degrees of freedom. In this case, values in the upper tail lie above the corresponding normal quantiles, the values in the lower tail below the corresponding normal quantiles.

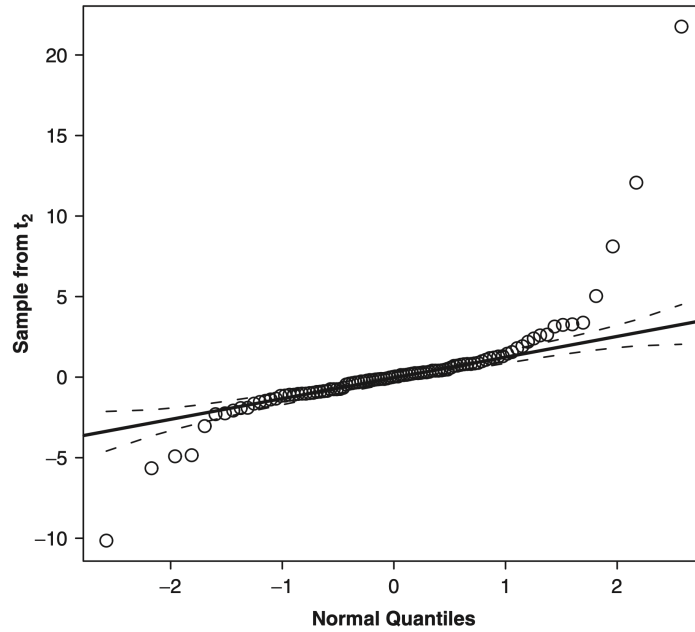


Figure 21.5: Normal quantile-comparison plot for a sample of 100 observations drawn from heavy-tailed t -distribution with 2 degrees of freedom. JF Figure 3.10.

- Figure 21.4 shows the normal quantile-comparison plot for the distribution of infant mortality. The positive skew of the distribution is readily apparent.

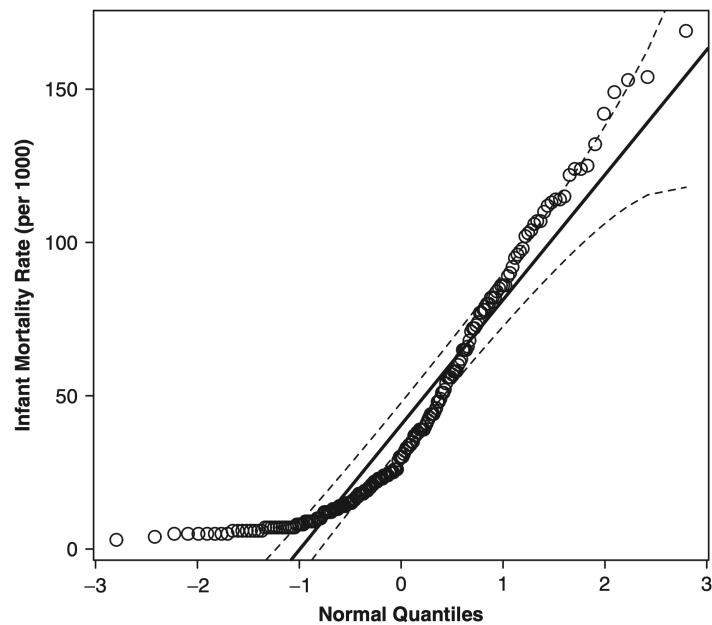


Figure 21.6: Normal quantile-comparison plot for the distribution of infant mortality. Note the positive skew. JF Figure 3.11.

Nonconstant error variance

One of the assumptions of the regression model is that the variation of the response variable around the regression surface (the error variance) is everywhere the same:

$$\text{Var}(\epsilon) = \text{Var}(Y|x_1, \dots, x_p) = \sigma_\epsilon^2$$

Constant error variance is often termed homoscedasticity, and similarly, nonconstant error variance is termed heteroscedasticity. We detect nonconstant error variances through graphical methods.

Residual plots

Because the least square residuals have unequal variance even when the constant variance assumption is correct:

$$\text{Var}(\hat{\epsilon}_i) = \sigma^2(1 - h_i).$$

It is preferable to plot studentized residuals against fitted values. A pattern of changing spread is often more easily discerned in a plot of absolute studentized residuals, $|\hat{\epsilon}_i^*|$, or squared studentized residuals, $\hat{\epsilon}_i^{*2}$, against \hat{Y} . If the values of \hat{Y} are all positive, then we can plot $\log|\hat{\epsilon}_i^*|$ against $\log \hat{Y}$. Figure 21.5 shows a plot of studentized residuals against fitted values and spread-level plot of studentized residuals, several points with negative fitted values were omitted.

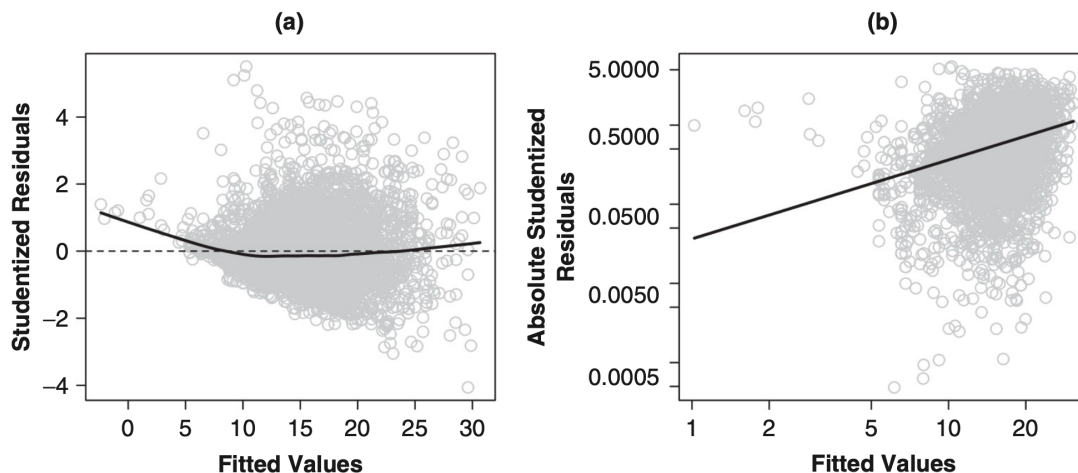


Figure 21.7: (a) Plot of studentized residuals versus fitted values and (b) spread-level plot for studentized residuals. JF Figure 12.3.

It is apparent from both graphs that the residual spread tends to increase with the level of the response, suggesting a violation of constant error variance assumption.

21 Lecture 23: March 22

Last time

- Added-variable plots
- Should unusual data be discarded
- Diagnosing non-normality, non-constant error variance, and nonlinearity (JF chapter 12)

Today

- lecture on Friday
- Diagnosing non-normality, non-constant error variance, and nonlinearity (JF chapter 12)

Quantile-comparison plot, JF 3.1.3

- Figure 21.1 plots a sample of $n = 100$ observations from a normal distribution with mean $\mu = 50$ and standard deviation $\sigma = 10$.

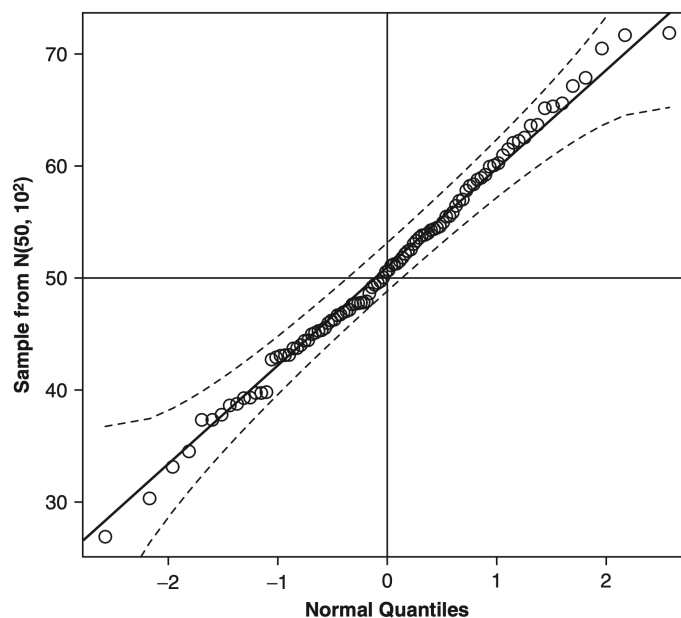


Figure 21.1: Normal quantile-comparison plot for a sample of 100 observations drawn from a normal distribution with mean 50 and standard deviation 10. The fitted line is through the quantiles of the distribution, the broken lines give a pointwise 95% confidence interval around the fit. JF Figure 3.8.

The plotted points are reasonably linear and stay within the rough 95% confidence envelope.

- Figure 21.2 plots a sample of $n = 100$ observations from the positively skewed chi-square distribution with 2 degrees of freedom. The positive skew of the data is reflected in points that lie *above* the comparison line in both tails of the distribution. (In contrast, the tails of negatively skewed data would lie *below* the comparison line.)

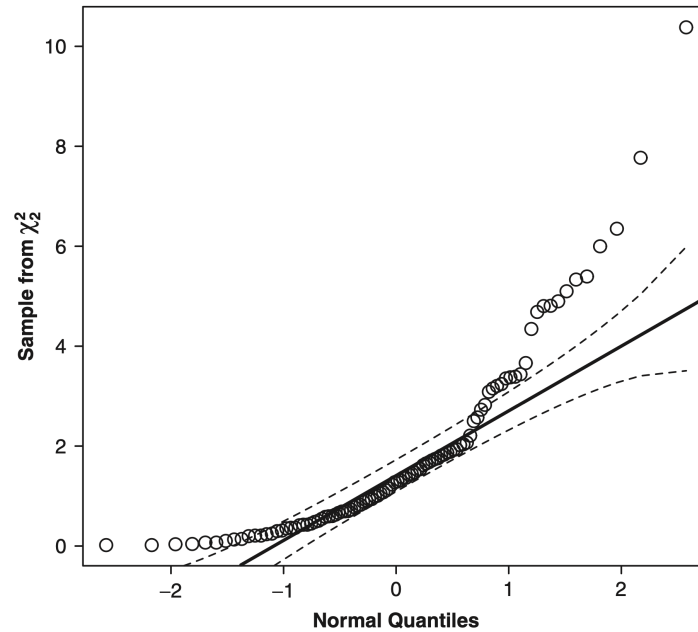


Figure 21.2: Normal quantile-comparison plot for a sample of 100 observations drawn from the positively skewed chi-square distribution with 2 degrees of freedom. JF Figure 3.9.

- Figure 21.3 plots a sample of $n = 100$ observations from the heavy-tailed t distribution with 2 degrees of freedom. In this case, values in the upper tail lie above the corresponding normal quantiles, the values in the lower tail below the corresponding normal quantiles.

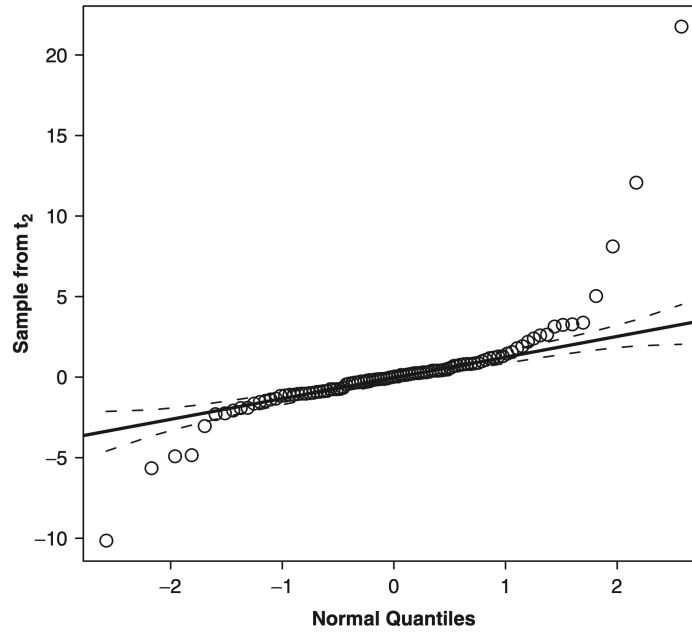


Figure 21.3: Normal quantile-comparison plot for a sample of 100 observations drawn from heavy-tailed t -distribution with 2 degrees of freedom. JF Figure 3.10.

- Figure 21.4 shows the normal quantile-comparison plot for the distribution of infant mortality. The positive skew of the distribution is readily apparent.

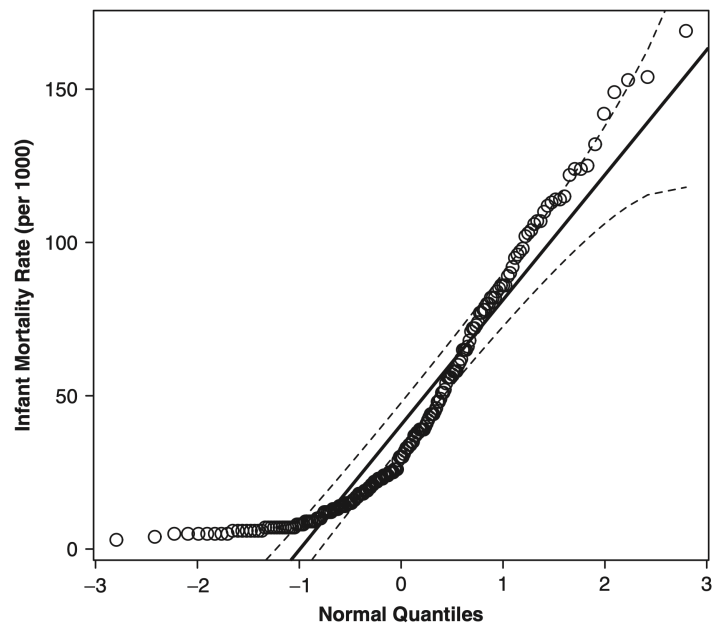


Figure 21.4: Normal quantile-comparison plot for the distribution of infant mortality. Note the positive skew. JF Figure 3.11.

Nonconstant error variance

One of the assumptions of the regression model is that the variation of the response variable around the regression surface (the error variance) is everywhere the same:

$$\text{Var}(\epsilon) = \text{Var}(Y|x_1, \dots, x_p) = \sigma_\epsilon^2$$

Constant error variance is often termed homoscedasticity, and similarly, nonconstant error variance is termed heteroscedasticity. We detect nonconstant error variances through graphical methods.

Residual plots

Because the least square residuals have unequal variance even when the constant variance assumption is correct:

$$\text{Var}(\hat{\epsilon}_i) = \sigma^2(1 - h_i).$$

It is preferable to plot studentized residuals against fitted values. A pattern of changing spread is often more easily discerned in a plot of absolute studentized residuals, $|\hat{\epsilon}_i^*|$, or squared studentized residuals, $\hat{\epsilon}_i^{*2}$, against \hat{Y} . If the values of \hat{Y} are all positive, then we can plot $\log|\hat{\epsilon}_i^*|$ against $\log \hat{Y}$. Figure 21.5 shows a plot of studentized residuals against fitted values and spread-level plot of studentized residuals, several points with negative fitted values were omitted.

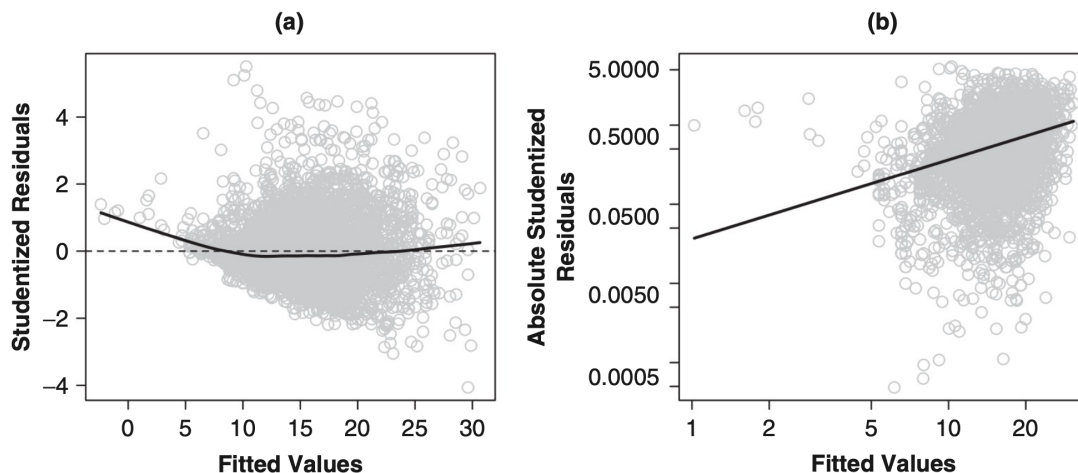


Figure 21.5: (a) Plot of studentized residuals versus fitted values and (b) spread-level plot for studentized residuals. JF Figure 12.3.

It is apparent from both graphs that the residual spread tends to increase with the level of the response, suggesting a violation of constant error variance assumption.

Weighted-least-squares estimation

Weighted-least-squares (WLS) regression provides an alternative approach to estimation in the presence of nonconstant error variance. Suppose that the errors from the linear regression model $\mathbf{Y} = \mathbf{X}\beta + \epsilon$ are independent and normally distributed, with zero means but *different* variances: $\epsilon_i \sim N(0, \sigma_i^2)$. Suppose further that the variances of the errors are known up to a constant of proportionality σ_ϵ^2 , so that $\sigma_i^2 = \sigma_\epsilon^2/w_i^2$. Then the likelihood for the model is

$$L(\beta, \sigma_\epsilon^2) = \frac{1}{(2\pi)^{n/2}|\mathbf{\Sigma}|^{1/2}} \exp \left[-\frac{1}{2}(\mathbf{Y} - \mathbf{X}\beta)^T \mathbf{\Sigma}^{-1}(\mathbf{Y} - \mathbf{X}\beta) \right]$$

where $\mathbf{\Sigma}$ is the covariance matrix of the errors,

$$\mathbf{\Sigma} = \sigma_\epsilon^2 \times \text{diag}\{1/w_1^2, \dots, 1/w_n^2\} \equiv \sigma_\epsilon^2 \mathbf{W}^{-1}$$

The maximum-likelihood estimators of β and σ_ϵ^2 are then

$$\begin{aligned} \hat{\beta} &= (\mathbf{X}^T \mathbf{W} \mathbf{X})^{-1} \mathbf{X}^T \mathbf{W} \mathbf{Y} \\ \hat{\sigma}_\epsilon^2 &= \frac{\sum (w_i \hat{\epsilon}_i)^2}{n} \end{aligned}$$

Correcting OLS standard errors for nonconstant variance

The covariance matrix of the ordinary-least-squares (OLS) estimator is

$$\begin{aligned} \mathbf{Var}(\hat{\beta}) &= (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{Var}(\mathbf{Y}) \mathbf{X} (\mathbf{X}^T \mathbf{X})^{-1} \\ &= \sigma_\epsilon^2 (\mathbf{X}^T \mathbf{X})^{-1} \end{aligned}$$

under the standard assumptions, including the assumption of constant error variance, $\mathbf{Var}(\mathbf{Y}) = \sigma_\epsilon^2 \mathbf{I}_n$. If, however, the errors are heteroscedastic but independent then $\mathbf{\Sigma} \equiv \mathbf{Var}(\mathbf{Y}) = \text{diag}\{\sigma_1^2, \dots, \sigma_n^2\}$, and

$$\mathbf{Var}(\hat{\beta}) = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{\Sigma} \mathbf{X} (\mathbf{X}^T \mathbf{X})^{-1}$$

White (1980) shows that the following is a consistent estimator of $\mathbf{Var}(\hat{\beta})$

$$\tilde{\mathbf{Var}}(\hat{\beta}) = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \hat{\mathbf{\Sigma}} \mathbf{X} (\mathbf{X}^T \mathbf{X})^{-1}$$

with $\hat{\mathbf{\Sigma}} = \text{diag}\{\hat{\sigma}_1^2, \dots, \hat{\sigma}_n^2\}$, where $\hat{\sigma}_i^2$ is the OLS residual for observation i .

Subsequent work suggested small modifications to White's coefficient-variance estimator, and in particular simulation studies by Long and Ervin (2000) support the use of

$$\tilde{\mathbf{Var}}^*(\hat{\beta}) = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \hat{\mathbf{\Sigma}}^* \mathbf{X} (\mathbf{X}^T \mathbf{X})^{-1}$$

where $\hat{\mathbf{\Sigma}}^* = \text{diag}\{\hat{\sigma}_i^2/(1 - h_i)^2\}$ and h_i is the hat-value associated with observation i . In large samples, where h_i is small, the distinction between $\tilde{\mathbf{Var}}(\hat{\beta})$ and $\tilde{\mathbf{Var}}^*(\hat{\beta})$ essentially disappears.

A rough *rule* is that nonconstant error variance seriously degrades the least-squares estimator only when the ratio of the largest to smallest variance is about 10 or more (or, more conservatively, about 4 or more).

Nonlinearity

If $\mathbf{E}(\mathbf{Y}|\mathbf{X})$ is not linear in \mathbf{X} (in other words, $\mathbf{E}(\epsilon|\mathbf{X}) \neq 0$ for some x), $\hat{\beta}$ may be biased and inconsistent. Usually we employ “linearity by default” but we should try to make sure this is appropriate: **detect** non-linearities and **model** them accurately.

Lowess smoother, JF 2.3

We can employ local averaging plots to help with diagnostics. Lowess method is in many respects similar to local-averaging smoothers, except that instead of computing an average Y -value within the neighborhood of a focal x , the lowess smoother computes a *fitted* value based on a locally weighted least-squares line, giving more weight to observations in the neighborhood that are close to the focal x than to those relatively far away. The name “lowess” is an acronym for *locally weighted scatterplot smoother* and is sometimes rendered as *loess*, for *local regression*.

Component-plus-residual plots

Added-variable plots, introduced for detecting influential data, can reveal nonlinearity. However, the added-variable plots are not always useful for locating a transformation:

- The added-variable plot adjusts X_j for the other X s.
- The *unadjusted* X_j is transformed in respecifying the model.

Moreover, Cook (1998, Section 14.5) shows that added-variable plots are biased toward linearity when the correlations among the explanatory variables are large.

Component-plus-residual plots (also called *partial-residual plots*) are often an effective alternative. The component-plus-residual plots are not as suitable as added-variable plots for revealing leverage and influence, though. The component-plus-residual plots are constructed by

1. Compute residuals from full regression:

$$\hat{\epsilon}_i = Y_i - \hat{Y}_i$$

2. Compute “linear component” of the partial relationship:

$$C_i = \hat{\beta}_j X_{ij}$$

3. Add linear component to residual to get partial residual for the j th explanatory variable

$$\hat{\epsilon}_i^{(j)} = \hat{\epsilon}_i + C_i = \hat{\epsilon}_i + \hat{\beta}_j X_{ij}$$

4. Plot $\hat{\epsilon}_i^{(j)}$ against $X_{.j}$

Figure 21.6 shows the component-plus-residual plots for the regression of log wages on variables (age, education and sex) of the 1994 wave of Statistics Canada's Survey of Labour and Income Dynamics (SLID) data. The SLID data set includes 3997 employed individuals who were between 16 and 65 years of age and who resided in Ontario.

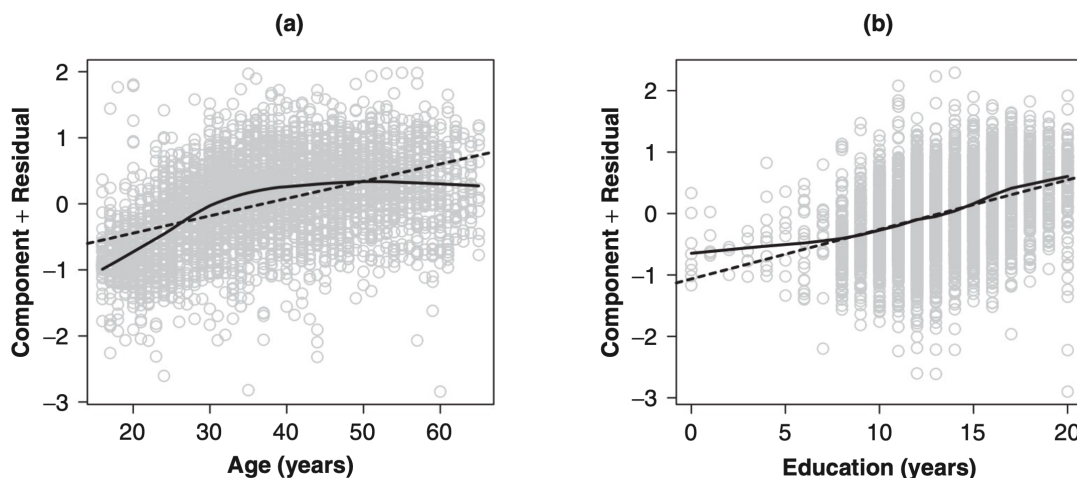


Figure 21.6: Component-plus-residual plots for age and education in SLID regression of log wages on these variables and sex. The solid lines are for lowess smooths with spans of 0.4, and the broken lines are for linear least-squares fits. JF Figure 12.6.

Data transformation

The family of powers and Roots, JF 4.1

A particularly useful group of transformations is the “family” of powers and roots:

$$X \rightarrow X^p$$

where the arrow indicates that we intend to replace X with the transformed variable X^p . If p is negative, then the transformation is an inverse power. For example, $X^{-1} = 1/X$. If p is a fraction, then the transformation represents a root. For example, $X^{1/3} = \sqrt[3]{X}$.

It is more convenient to define the family of power transformations in a slightly more complex manner, called the Box-Cox family of transformations (introduced in a seminal paper on transformations by Box & Cox, 1964):

$$X \rightarrow X^{(p)} = \frac{X^p - 1}{p}$$

Because $X^{(p)}$ is a linear function of X^p , the two transformations have the same essential effect on the data, but, as is apparent in Figure 21.7

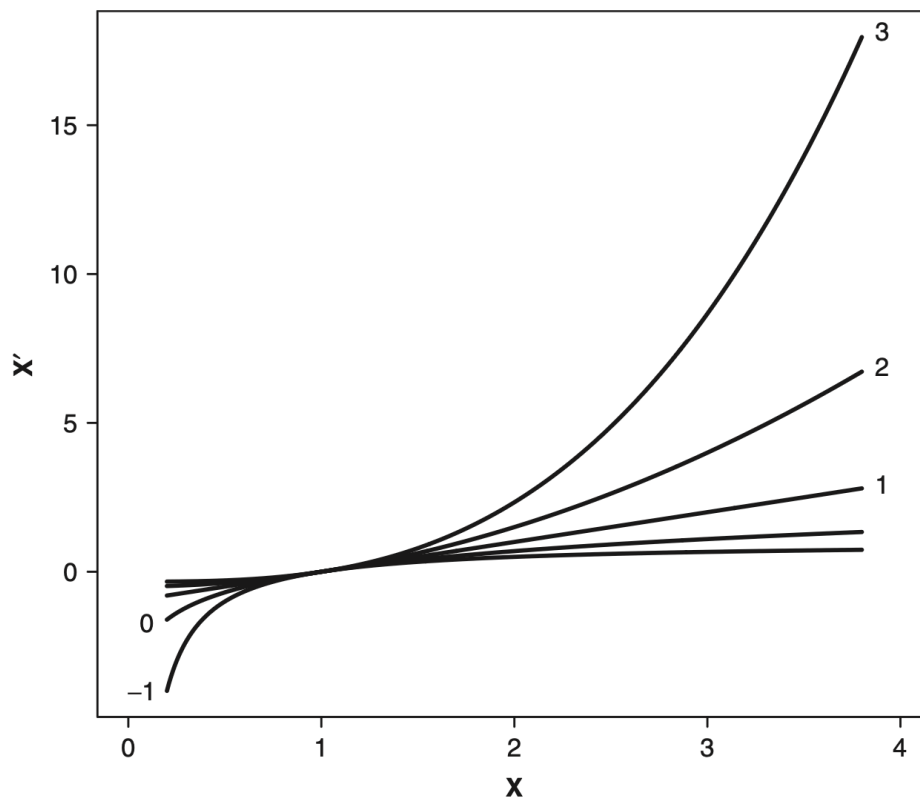


Figure 21.7: The Box-Cox family of power transformations X' of X . The curve labeled p is the transformation $X^{(p)}$, that is $(X^p - 1)/p$; $X^{(0)}$ is $\log_e(X)$. JF Figure 4.1.

- Dividing by p preserves the direction of X , which otherwise would be reversed when p is negative.
- The transformations $X^{(p)}$ are “matched” above $X = 1$ both in level and in slope:
 1. $1^{(p)} = 0$, for all values of p
 2. each transformation has a slope of 1 at $X = 1$.
- Descending the “ladder” of powers and roots towards $X^{(-1)}$ compresses the large values of X and spreads out the small ones. Ascending the ladder of powers and roots towards $X^{(2)}$ has the opposite effect. As p moves further from $p = 1$ (i.e., no transformation) in either direction, the transformation grows more powerful, increasingly “bending” the data.
- The power transformation X^0 is useless because it changes all values to 1, but we can think of the log transformation as a kind of “zeroth” power:

$$\lim_{p \rightarrow 0} \frac{X^p - 1}{p} = \log_e X$$

and by convention, $X^{(0)} \equiv \log_e X$.

Box-Cox transformation of Y

Box and Cox (1964) suggested a power transformation of Y with the object of normalizing the error distribution, stabilizing the error variance, and straightening the relationship of Y to the X s. The general Box-Cox model is

$$Y_i^{(\lambda)} = \beta_0 + \beta_1 X_{i1} + \cdots + \beta_p X_{ip} + \epsilon_i$$

where $\epsilon_i \stackrel{iid}{\sim} N(0, \sigma_\epsilon^2)$, and

$$Y_i^{(\lambda)} = \begin{cases} \frac{Y_i^\lambda - 1}{\lambda} & \text{for } \lambda \neq 0 \\ \log_e Y_i & \text{for } \lambda = 0 \end{cases}$$

Note: in statistics, \log_e is often written as \log .

For a particular choice of λ , the conditional maximized log-likelihood (see JF 12.5.1 p.324 footnote 55) is

$$\begin{aligned} \log_e L(\beta_0, \beta_1, \dots, \beta_p, \sigma_\epsilon^2 | \lambda) &= -\frac{n}{2}(1 + \log_e 2\pi) \\ &\quad - \frac{n}{2} \log_e \hat{\sigma}_\epsilon^2(\lambda) + (\lambda - 1) \sum_{i=1}^n \log_e Y_i \end{aligned}$$

where $\hat{\sigma}_\epsilon^2(\lambda) = \sum \hat{\epsilon}_i^2(\lambda)/n$ and where $\hat{\epsilon}_i(\lambda)$ are the residuals from the least-squares regression of $Y^{(\lambda)}$ on X s. The least-squares coefficients from this regression are the maximum-likelihood estimates of β s conditional on the values of λ .

A simple procedure for finding the maximum-likelihood estimator $\hat{\lambda}$ is to evaluate the maximized $\log_e L$ (called the profile log-likelihood) for a range of values of λ . To test: $H_0 : \lambda = 1$, calculated the likelihood-ratio statistic

$$G_0^2 = -2[\log_e L(\lambda = 1) - \log_e L(\lambda = \hat{\lambda})]$$

which is asymptotically distributed as χ_1^2 with one degree of freedom under H_0 . A 95% confidence interval for λ includes those values for which

$$\log_e L(\lambda) > \log_e L(\lambda = \hat{\lambda}) - 1.92$$

The number 1.92 comes from $\frac{1}{2}\chi_{1,0.05}^2 = 0.5 \times 1.96^2$.

Figure 21.1 shows a plot of the profile log-likelihood against λ for the original SLID regression of composite hourly wages on sex, age, and education. The maximum-likelihood estimate of λ is $\hat{\lambda} = 0.09$, and a 95% confidence interval runs from 0.04 to 0.13. Although 0 is outside of the CI (confidence interval), it is essentially the same as log transformation of wages (the correlation between log wages and wages^{0.09} is 0.9996).

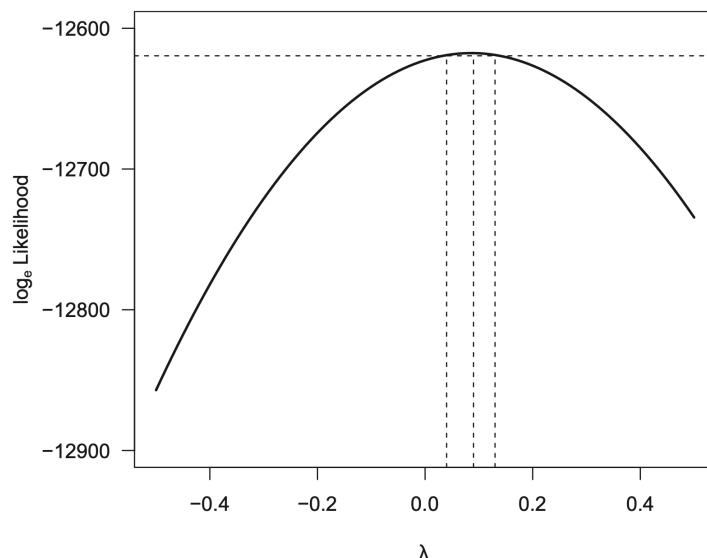


Figure 21.8: Box-Cox transformations for the SLID regression of wages on sex, age, and education. The maximized (profile) log-likelihood is plotted against the transformation parameter λ . The intersection of the line near the top of the graph with the profile log-likelihood curve marks off a 95% confidence interval for λ . The maximum of the log-likelihood corresponds to the MLE of λ . JF Figure 12.14.

Box-Tidwell transformation of X s

Now, consider the model

$$Y_i = \beta_0 + \beta_1 X_{i1}^{\gamma_1} + \cdots + \beta_p X_{ip}^{\gamma_p} + \epsilon_i$$

where the errors are independently distributed as $\epsilon_i \stackrel{iid}{\sim} N(0, \sigma_\epsilon^2)$ and all the X_{ij} are positive.

The parameters of this model ($\beta_0, \beta_1, \dots, \beta_p, \gamma_1, \dots, \gamma_p$, and σ_ϵ^2) could be estimated by general nonlinear least squares. Box and Tidwell (1962) suggested the following computationally more efficient procedure (also yields a constructed-variable diagnostic):

1. Regress Y on X_1, \dots, X_p , obtaining $\hat{\beta}_0, \hat{\beta}_1, \dots, \hat{\beta}_p$. (“Regress A on B s” is the same as “fitting the linear regression model with A as the response variable and B s as the explanatory variables”.)
2. Regress Y on X_1, \dots, X_p and the constructed variables $X_1 \log_e X_1, \dots, X_p \log_e X_p$ (again, by fitting the model of $Y = \beta_0 + \beta_1 X_1 + \cdots + \beta_p X_p + \delta_1 X_1 \log_e X_1 + \cdots + \delta_p X_p \log_e X_p + \epsilon_i$) to obtain $\tilde{\beta}_0, \tilde{\beta}_1, \dots, \tilde{\beta}_p, \tilde{\delta}_1, \dots, \tilde{\delta}_p$. In general $\hat{\beta}_i \neq \tilde{\beta}_i$. (The constructed variables result from the first-order Taylor-series approximation to $X_j^{\gamma_j}$ evaluated at $\gamma_j = 1$: $X_j^{\gamma_j} \approx X_j + (\gamma_j - 1)X_j \log_e X_j$.)
3. The constructed variable $X_j \log_e X_j$ can be used to assess the need for a transformation of X_j by testing the null hypothesis $H_0 : \delta_j = 0$. Added-variable plots for the constructed variables are useful for assessing leverage and influence on the decision to transform the X s.

4. A preliminary estimate of the transformation parameter γ_j (not the MLE) is

$$\tilde{\gamma}_j = 1 + \frac{\tilde{\delta}_j}{\hat{\beta}_j}$$

where $\tilde{\delta}_j$ is from step 2 and $\hat{\beta}_j$ is from step 1.

Polynomial regression

A machinery of multiple regression to fit non-linear relationships between predictor(s) and response.

- Linear: $y = \beta_0 + \beta_1 x + \epsilon$
- Quadratic: $y = \beta_0 + \beta_1 x + \beta_2 x^2 + \epsilon$
- Cubic: $y = \beta_0 + \beta_1 x + \beta_2 x^2 + \beta_3 x^3 + \epsilon$
- k^{th} order polynomial: $y = \beta_0 + \beta_1 x + \beta_2 x^2 + \cdots + \beta_k x^k + \epsilon$

Question:

Does quadratic model provide a significantly better fit than linear model?

Solution: Test $H_0 : \beta_2 = 0$ vs. $H_a : \beta_2 \neq 0$.

Alternatively, compare the corresponding adjusted- R^2 values.

21 Lecture 24: March 27

Last time

- Diagnosing non-normality, non-constant error variance, and nonlinearity (JF chapter 12)

Today

- lecture on Friday
- HW3 deadline extended to April 21th
- Data transformation
- Collinearity (JF chapter 13, RD 8.3.2)
- Principal component analysis (JF 13.1.1, RD 8.3.4)

Box-Cox transformation of Y

Box and Cox (1964) suggested a power transformation of Y with the object of normalizing the error distribution, stabilizing the error variance, and straightening the relationship of Y to the X s. The general Box-Cox model is

$$Y_i^{(\lambda)} = \beta_0 + \beta_1 X_{i1} + \cdots + \beta_p X_{ip} + \epsilon_i$$

where $\epsilon_i \stackrel{iid}{\sim} N(0, \sigma_\epsilon^2)$, and

$$Y_i^{(\lambda)} = \begin{cases} \frac{Y_i^\lambda - 1}{\lambda} & \text{for } \lambda \neq 0 \\ \log_e Y_i & \text{for } \lambda = 0 \end{cases}$$

Note: in statistics, \log_e is often written as \log .

For a particular choice of λ , the conditional maximized log-likelihood (see JF 12.5.1 p.324 footnote 55) is

$$\begin{aligned} \log_e L(\beta_0, \beta_1, \dots, \beta_p, \sigma_\epsilon^2 | \lambda) &= -\frac{n}{2}(1 + \log_e 2\pi) \\ &\quad - \frac{n}{2} \log_e \hat{\sigma}_\epsilon^2(\lambda) + (\lambda - 1) \sum_{i=1}^n \log_e Y_i \end{aligned}$$

where $\hat{\sigma}_\epsilon^2(\lambda) = \sum \hat{\epsilon}_i^2(\lambda)/n$ and where $\hat{\epsilon}_i(\lambda)$ are the residuals from the least-squares regression of $Y^{(\lambda)}$ on X s. The least-squares coefficients from this regression are the maximum-likelihood estimates of β s conditional on the values of λ .

A simple procedure for finding the maximum-likelihood estimator $\hat{\lambda}$ is to evaluate the maximized $\log_e L$ (called the profile log-likelihood) for a range of values of λ . To test: $H_0 : \lambda = 1$, calculated the likelihood-ratio statistic

$$G_0^2 = -2[\log_e L(\lambda = 1) - \log_e L(\lambda = \hat{\lambda})]$$

which is asymptotically distributed as χ_1^2 with one degree of freedom under H_0 . A 95% confidence interval for λ includes those values for which

$$\log_e L(\lambda) > \log_e L(\lambda = \hat{\lambda}) - 1.92$$

The number 1.92 comes from $\frac{1}{2}\chi_{1,0.05}^2 = 0.5 \times 1.96^2$.

Figure 21.1 shows a plot of the profile log-likelihood against λ for the original SLID regression of composite hourly wages on sex, age, and education. The maximum-likelihood estimate of λ is $\hat{\lambda} = 0.09$, and a 95% confidence interval runs from 0.04 to 0.13. Although 0 is outside of the CI (confidence interval), it is essentially the same as log transformation of wages (the correlation between log wages and wages^{0.09} is 0.9996).

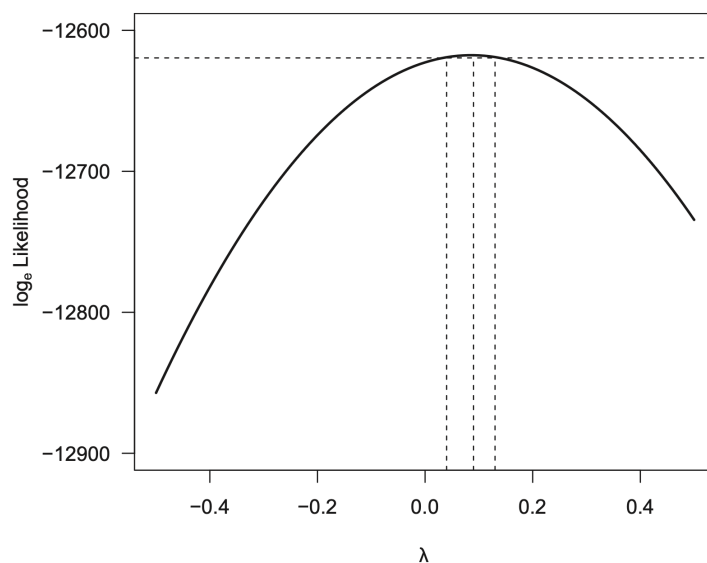


Figure 21.1: Box-Cox transformations for the SLID regression of wages on sex, age, and education. The maximized (profile) log-likelihood is plotted against the transformation parameter λ . The intersection of the line near the top of the graph with the profile log-likelihood curve marks off a 95% confidence interval for λ . The maximum of the log-likelihood corresponds to the MLE of λ . JF Figure 12.14.

Box-Tidwell transformation of X s

Now, consider the model

$$Y_i = \beta_0 + \beta_1 X_{i1}^{\gamma_1} + \cdots + \beta_p X_{ip}^{\gamma_p} + \epsilon_i$$

where the errors are independently distributed as $\epsilon_i \stackrel{iid}{\sim} N(0, \sigma_\epsilon^2)$ and all the X_{ij} are positive.

The parameters of this model ($\beta_0, \beta_1, \dots, \beta_p, \gamma_1, \dots, \gamma_p$, and σ_ϵ^2) could be estimated by general nonlinear least squares. Box and Tidwell (1962) suggested the following computationally more efficient procedure (also yields a constructed-variable diagnostic):

1. Regress Y on X_1, \dots, X_p , obtaining $\hat{\beta}_0, \hat{\beta}_1, \dots, \hat{\beta}_p$. (“Regress A on B s” is the same as “fitting the linear regression model with A as the response variable and B s as the explanatory variables”.)
2. Regress Y on X_1, \dots, X_p and the constructed variables $X_1 \log_e X_1, \dots, X_p \log_e X_p$ (again, by fitting the model of $Y = \beta_0 + \beta_1 X_1 + \dots + \beta_p X_p + \delta_1 X_1 \log_e X_1 + \dots + \delta_p X_p \log_e X_p + \epsilon_i$) to obtain $\tilde{\beta}_0, \tilde{\beta}_1, \dots, \tilde{\beta}_p, \tilde{\delta}_1, \dots, \tilde{\delta}_p$. In general $\hat{\beta}_i \neq \tilde{\beta}_i$. (The constructed variables result from the first-order Taylor-series approximation to $X_j^{\gamma_j}$ evaluated at $\gamma_j = 1$: $X_j^{\gamma_j} \approx X_1 + (\gamma_j - 1)X_1 \log_e X_1$.)
3. The constructed variable $X_j \log_e X_j$ can be used to assess the need for a transformation of X_j by testing the null hypothesis $H_0 : \delta_j = 0$. Added-variable plots for the constructed variables are useful for assessing leverage and influence on the decision to transform the X s.
4. A preliminary estimate of the transformation parameter γ_j (not the MLE) is

$$\tilde{\gamma}_j = 1 + \frac{\tilde{\delta}_j}{\hat{\beta}_j}$$

where $\tilde{\delta}_j$ is from step 2 and $\hat{\beta}_j$ is from step 1.

Polynomial regression

A machinery of multiple regression to fit non-linear relationships between predictor(s) and response.

- Linear: $y = \beta_0 + \beta_1 x + \epsilon$
- Quadratic: $y = \beta_0 + \beta_1 x + \beta_2 x^2 + \epsilon$
- Cubic: $y = \beta_0 + \beta_1 x + \beta_2 x^2 + \beta_3 x^3 + \epsilon$
- k^{th} order polynomial: $y = \beta_0 + \beta_1 x + \beta_2 x^2 + \dots + \beta_k x^k + \epsilon$

Question:

Does quadratic model provide a significantly better fit than linear model?

Solution: Test $H_0 : \beta_2 = 0$ vs. $H_a : \beta_2 \neq 0$.

Alternatively, compare the corresponding adjusted- R^2 values.

Additional reference

- “A First Course in Linear Model Theory” by Nalini Ravishanker and Kipak K. Dey
- [Lecture notes](#) by Cedric Ginestet

Collinearity

In linear model

$$\begin{aligned} \mathbf{Y} &= \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\epsilon} \\ \boldsymbol{\epsilon} &\sim \mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{I}_n) \end{aligned}$$

Collinearity (or multicollinearity) exists when there is “near-dependency” between the columns of the design matrix \mathbf{X} .

- Two or more columns.
- In other words, high correlation between explanatory variables.
- the data/model pair is ill-conditioned when $\mathbf{X}^T\mathbf{X}$ is nearly singular.

Perfect collinearity leads to rank-deficiency in \mathbf{X} such that $\mathbf{X}^T\mathbf{X}$ is singular. In the case of perfect collinearity, two or more columns are linear-dependent.

An example of perfect collinearity

$$Y_i = \beta_0 + \beta_1 X_{i1} + \beta_2 X_{i2} + \beta_3 X_{i3} + \beta_4 X_{i4} + \beta_5 X_{i5} + \epsilon_i$$

Consider the case, where

- Y_i represents the amount of sales.
- $X_{i1}, X_{i2}, \dots, X_{i4}$ are categorical that represent the quarter in which the sample is collected: $X_{ij} = \mathbf{1}(\text{sample } i \text{ collected in quarter } j)$.
- X_{i5} represents expense spent in advertising.

The dummy variable trap $X_{i4} = 1 - X_{i1} - X_{i2} - X_{i3}$. Recall that we need $m - 1$ dummy variables for m categories.

An example of high correlation between predictors

$$Y_i = \beta_0 + \beta_1 X_{i1} + \beta_2 X_{i2} + \epsilon_i$$

Consider the case, where

- Y_i represents the salary of individual i .
- X_{i1} represents the age of individual i .
- X_{i2} represents the experience of individual i .

How to interpret β_1 ?

We expect high correlation between age and experience.

Problems caused by multicollinearity

1. large standard errors of the regression coefficients
 - small associated t-statistics
 - conclusion that truly useful explanatory variables are insignificant in explaining the regression

2. the sign of regression coefficients may be the opposite of what a mechanistic understanding of the problem would suggest
3. deleting a column of the predictor matrix will cause large changes in the coefficient estimates for other variables

However, multicollinearity does **not** greatly affect the **predicted values**.

25 Lecture 25: March 29

Last time

- Data transformation
- Collinearity (JF chapter 13, RD 8.3.2)

Today

- lecture on Friday
- Collinearity (JF chapter 13, RD 8.3.2)
- Principal component analysis (JF 13.1.1, RD 8.3.4)
- Biased estimation:
 - Ridge Regression
 - Lasso Regression
- Model selection

Signs and detections of multicollinearity

Some signs for multicollinearity:

1. Simple correlation between a pair of predictors exceeds 0.9 or R^2 .
2. High value of the multiple correlation coefficient with some high partial correlations between the explanatory variables.
3. Large F -statistics with some small t -statistics for individual regression coefficients

Some approaches for detecting multicollinearity:

1. Pairwise correlations among the explanatory variables
2. Variance inflation factor
3. Condition number

Variance inflation factor

For a multiple linear regression with k explanatory variables. We can regress X_j on the $(k - 1)$ other explanatory variables and denote R_j as the coefficient of determination.

Then the variance inflation factor (VIF) is defined as

$$\text{VIF}_j = \frac{1}{1 - R_j^2}$$

- $\text{VIF}_j \in [1, +\infty)$

- A suggested threshold is 10
- May use the averaged $\overline{\text{VIF}} = \sum_{j=1}^k \text{VIF}_j / k$.

Condition index and condition number

We first scale the design matrix \mathbf{X} into column-equilibrated predictor matrix \mathbf{X}_E such that $\{X_E\}_{ij} = X_{ij} / \sqrt{\mathbf{X}_j^T \mathbf{X}_j}$.

Let $\mathbf{X}_E = \mathbf{U}\mathbf{D}\mathbf{V}^T$ be the singular-value decomposition (SVD) of the $n \times p$ matrix \mathbf{X}_E where $\mathbf{U}^T \mathbf{U} = \mathbf{V}^T \mathbf{V} = \mathbf{I}_p$ and $\mathbf{D} = \text{diag}(d_1, d_2, \dots, d_p)$ is a diagonal matrix with $d_j \geq 0$.

The j^{th} condition index is defined as

$$\eta(\mathbf{X}_E) = d_{\max} / d_j, \quad j = 1, 2, \dots, p$$

The condition number is defined as

$$C = d_{\max} / d_{\min}$$

$$C \geq 1, \quad d_{\max} = \max_{1 \leq j \leq p} d_j \quad \text{and} \quad d_{\min} = \min_{1 \leq j \leq p} d_j$$

Some properties of the condition number

- Large condition number indicates evidence of multicollinearity
- Typical cutoff values, 10, 15 to 30.

Some problems with the condition number

- practitioners have different opinions of whether \mathbf{X} should be centered around their means for SVD.
 - centering may remove nonessential ill conditioning, e.g. $\text{Cor}(X, X^2)$
 - centering may mask the role of the constant term in any underlying near-dependencies
- the degree of multicollinearity with dummy variables may be influenced by the choice of reference category
- condition number is affected by the scale of the \mathbf{X} measurements
 - By scaling down any column of \mathbf{X} , the condition number can be made arbitrarily large
 - Known as *artificial ill-conditioning*
 - The condition number of the scaled matrix \mathbf{X}_E is also referred to as the *scaled condition number*

Recall that $\mathbf{X}_E = \mathbf{U}\mathbf{D}\mathbf{V}^T$ is the singular-value decomposition (SVD) of \mathbf{X}_E , where $\mathbf{U}^T\mathbf{U} = \mathbf{V}^T\mathbf{V} = \mathbf{I}_p$ and $\mathbf{D} = \text{diag}(d_1, d_2, \dots, d_p)$ is a diagonal matrix with $d_j \geq 0$.

Then

$$\begin{aligned}\mathbf{X}_E^T\mathbf{X}_E &= \mathbf{V}\mathbf{D}\mathbf{U}^T\mathbf{U}\mathbf{D}\mathbf{V}^T \\ &= \mathbf{V}\mathbf{D}^2\mathbf{V}^T\end{aligned}$$

is the spectral decomposition of the Gramian matrix $\mathbf{X}_E^T\mathbf{X}_E$ with $\{d_j^2\}$ being the eigenvalues and \mathbf{V} being the corresponding eigen vector matrix. This relationship links the condition numbers to the eigen values of the Gramian matrix.

Variance decomposition method

The variance-covariance matrix of the coefficient

$$\begin{aligned}\text{Cov}(\hat{\boldsymbol{\beta}}) &= \sigma^2(\mathbf{X}_E^T\mathbf{X}_E)^{-1} \\ &= \sigma^2\mathbf{V}\mathbf{D}^{-2}\mathbf{V}^T\end{aligned}$$

Its j^{th} diagonal element is the estimated variance of the j^{th} coefficient, $\hat{\beta}_j$. Then

$$\text{Var}(\hat{\beta}_j) = \sigma^2 \sum_{h=1}^p \frac{v_{jh}^2}{d_h^2}$$

- Let $q_{jh} = \frac{v_{jh}^2}{d_h^2}$ and $q_j = \sum_{h=1}^p q_{jh}$.
- The variance decomposition proportion is $\pi_{jh} = q_{jh}/q_j$.
- π_{jh} denotes the proportion of the variance of the j^{th} regression coefficient associated with the h^{th} component of its decomposition.
- The variance decomposition proportion matrix is $\boldsymbol{\Pi} = \{\pi_{jh}\}$.

In practice, it is suggested to combine condition index and proportions of variance for multicollinearity diagnostic. Identify multicollinearity if

- Two or more elements in the j^{th} row of matrix $\boldsymbol{\Pi}$ are relatively large
- And its associated condition index η_j is large too

Principal Components

The method of principal components, introduced by Karl Pearson (1901) and Harold Hotelling (1933), provides a useful representation of the correlational structure of a set of variables. Some advantages of the principal component analysis include

Condition	Proportions of variance			
Index	$Var(\hat{\beta}_1)$	$Var(\hat{\beta}_2)$...	$Var(\hat{\beta}_3)$
η_1	π_{11}	π_{12}	...	π_{1p}
η_2	π_{21}	π_{22}	...	π_{2p}
\vdots	\vdots	\vdots		\vdots
η_p	π_{p1}	π_{p2}	...	π_{pp}

Table 1: Table of condition index and proportions of variance

- more unified
- linear transformation of the original predictors into a new set of orthogonal predictors
- the new orthogonal predictors are called principal components

Principal components regression is an approach that inspects the sample data (\mathbf{Y}, \mathbf{X}) for directions of variability and uses this information to reduce the dimensionality of the estimation problem. The procedure is based on the observation that every linear regression model can be restated in terms of a set of orthogonal predictor variables, which are constructed as linear combinations of the original variables. The new orthogonal variables are called the principal components of the original variables.

Let $\mathbf{X}^T \mathbf{X} = \mathbf{Q} \mathbf{\Delta} \mathbf{Q}^T$ denote the spectral decomposition of $\mathbf{X}^T \mathbf{X}$, where $\mathbf{\Delta} = \text{diag}\{\lambda_1, \dots, \lambda_p\}$ is a diagonal matrix consisting of the (real) eigenvalues of $\mathbf{X}^T \mathbf{X}$, with $\lambda_1 \geq \dots \geq \lambda_p$ and $\mathbf{Q} = (\mathbf{q}_1, \dots, \mathbf{q}_p)$ denotes the matrix whose columns are the orthogonal eigenvectors of $\mathbf{X}^T \mathbf{X}$ corresponding to the ordered eigenvalues. Consider the transformation

$$\mathbf{Y} = \mathbf{X} \mathbf{Q} \mathbf{Q}^T \boldsymbol{\beta} + \boldsymbol{\epsilon} = \mathbf{Z} \boldsymbol{\theta} + \boldsymbol{\epsilon},$$

where $\mathbf{Z} = \mathbf{X} \mathbf{Q}$, and $\boldsymbol{\theta} = \mathbf{Q}^T \boldsymbol{\beta}$.

The elements of $\boldsymbol{\theta}$ are known as the regression parameters of the principal components. The matrix $\mathbf{Z} = \{\mathbf{z}_1, \dots, \mathbf{z}_p\}$ is called the matrix of principal components of $\mathbf{X}^T \mathbf{X}$. $\mathbf{z}_j = \mathbf{X} \mathbf{q}_j$ is the j th principal component of $\mathbf{X}^T \mathbf{X}$ and $\mathbf{z}_j^T \mathbf{z}_j = \lambda_j$, the j th largest eigenvalue of $\mathbf{X}^T \mathbf{X}$.

Principal components regression consists of deleting one or more of the variables \mathbf{z}_j (which correspond to small values of λ_j), and using OLS estimation on the resulting reduced regression model.

Derivation under standardized predictors, JF 13.1.1

Consider the vectors of standardized predictors, $\mathbf{x}_1^*, \mathbf{x}_2^*, \dots, \mathbf{x}_p^*$ (obtained by subtracting the mean and divided by standard deviation of the original predictor vectors). Because the principal components are linear combinations of the original predictors, we write the first principal component as

$$\begin{aligned} \mathbf{w}_1 &= A_{11} \mathbf{x}_1^* + A_{21} \mathbf{x}_2^* + \dots + A_{p1} \mathbf{x}_p^* \\ &= \mathbf{X}^* \mathbf{a}_1 \end{aligned}$$

The variance of the first component becomes

$$\begin{aligned} S_{w_1}^2 &= \frac{1}{n-1} \mathbf{w}_1^T \mathbf{w}_1 \\ &= \frac{1}{n-1} \mathbf{a}_1^T \mathbf{X}^{*T} \mathbf{X}^* \mathbf{a}_1 \\ &= \mathbf{a}_1^T \mathbf{R}_{XX} \mathbf{a}_1 \end{aligned}$$

where $\mathbf{R}_{XX} = \frac{1}{n-1} \mathbf{X}^{*T} \mathbf{X}^*$. We want to maximize $S_{w_1}^2$ under the normalizing constraint $\mathbf{a}_1^T \mathbf{a}_1 = 1$ (otherwise $S_{w_1}^2$ can be arbitrarily large by inflating \mathbf{a}_1). Consider

$$F_1 \equiv \mathbf{a}^T \mathbf{R}_{XX} \mathbf{a}_1 - L_1(\mathbf{a}_1^T \mathbf{a}_1 - 1)$$

where L_1 is a Lagrange multiplier. By differentiating this equation with respect to \mathbf{a}_1 and L_1 ,

$$\begin{aligned} \frac{\partial F_1}{\partial \mathbf{a}_1} &= 2\mathbf{R}_{XX} \mathbf{a}_1 - 2L_1 \mathbf{a}_1 \\ \frac{\partial F_1}{\partial L_1} &= -(\mathbf{a}_1^T \mathbf{a}_1 - 1) \end{aligned}$$

Setting the partial derivatives to 0 produces

$$\begin{aligned} (\mathbf{R}_{XX} - L_1 \mathbf{I}_p) \mathbf{a}_1 &= \mathbf{0} \\ \mathbf{a}_1^T \mathbf{a}_1 &= 1 \end{aligned}$$

From the first equation, we see that L_1 is an eigenvalue of \mathbf{R}_{XX} such that $\mathbf{R}_{XX} \mathbf{a}_1 = L_1 \mathbf{a}_1$ such that

$$S_{w_1}^2 = \mathbf{a}_1^T \mathbf{R}_{XX} \mathbf{a}_1 = L_1 \mathbf{a}_1^T \mathbf{a}_1 = L_1$$

To maximize $S_{w_1}^2$, we only need to pick the largest eigenvalue of \mathbf{R}_{XX} .

Additional reference

- “A First Course in Linear Model Theory” by Nalini Ravishanker and Kipak K. Dey
- [Lecture notes](#) by Cedric Ginestet

Ridge Regression

Ridge regression and the Lasso regression are two forms of regularized regression. These methods can be used to alleviate the consequences of multicollinearity.

1. When variables are highly correlated, a large coefficient in one variable may be alleviated by a large coefficient in another variable, which is negatively correlated to the former.
2. Regularization imposes an upper threshold on the values taken by the coefficients, thereby producing a more parsimonious solution, and a set of coefficients with smaller variance.

Constrained optimization

Ridge regression is motivated by a constrained minimization problem, which can be formulated as

$$\begin{aligned}\hat{\boldsymbol{\beta}}^{ridge} &= \arg \min_{\boldsymbol{\beta} \in \mathbb{R}^p} \sum_{i=1}^n (y_i - \mathbf{x}_i^T \boldsymbol{\beta})^2 \\ \text{subject to } \|\boldsymbol{\beta}\|_2^2 &= \sum_{j=1}^p \beta_j^2 \leq t\end{aligned}$$

for $t \geq 0$.

Use a Lagrange multiplier, we can rewrite the formula as

$$\hat{\boldsymbol{\beta}}^{ridge} = \arg \min_{\boldsymbol{\beta} \in \mathbb{R}^p} \left\{ \sum_{i=1}^n (y_i - \mathbf{x}_i^T \boldsymbol{\beta})^2 + \lambda \sum_{j=1}^p \beta_j^2 \right\}$$

for $\lambda \geq 0$ and where there is a one-to-one correspondence between t and λ . λ is an arbitrary constant usually referred to as the “ridge constant”.

Analytical solutions

The ridge-regression estimator has analytical solution

$$\hat{\boldsymbol{\beta}}^{ridge} = (\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I})^{-1} \mathbf{X}^T \mathbf{Y}$$

This is obtained by differentiating the objective function with respect to $\boldsymbol{\beta}$ and set it to 0:

$$\begin{aligned}\frac{\partial}{\partial \boldsymbol{\beta}} \{ (\mathbf{Y} - \mathbf{X}\boldsymbol{\beta})^T (\mathbf{Y} - \mathbf{X}\boldsymbol{\beta}) + \lambda \boldsymbol{\beta}^T \boldsymbol{\beta} \} \\ = 2(\mathbf{X}^T \mathbf{X})\boldsymbol{\beta} - 2\mathbf{X}^T \mathbf{Y} + 2\lambda \boldsymbol{\beta} \\ = 0\end{aligned}$$

Therefore,

$$(\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I})\boldsymbol{\beta} = \mathbf{X}^T \mathbf{Y}$$

Since we are adding a positive constant to the diagonal of $\mathbf{X}^T \mathbf{X}$, we are, in general, producing an invertible matrix, $\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I}$ even if $\mathbf{X}^T \mathbf{X}$ is singular. Historically, this particular aspect of ridge regression was the main motivation behind the adoption of this particular extension of OLS theory.

The ridge regression estimator is related to the classical OLS estimator, $\hat{\boldsymbol{\beta}}^{OLS}$, in the following manner

$$\hat{\boldsymbol{\beta}}^{ridge} = [\mathbf{I} + \lambda(\mathbf{X}^T \mathbf{X})^{-1}]^{-1} \hat{\boldsymbol{\beta}}^{OLS},$$

assuming $\mathbf{X}^T \mathbf{X}$ is non-singular. This relationship can be verified by applying the definition of $\hat{\boldsymbol{\beta}}^{OLS}$,

$$\begin{aligned}\hat{\boldsymbol{\beta}}^{ridge} &= [\mathbf{I} + \lambda(\mathbf{X}^T \mathbf{X})^{-1}]^{-1} (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{Y} \\ &= (\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I})^{-1} \mathbf{X}^T \mathbf{Y}\end{aligned}$$

using the fact $\mathbf{B}^{-1}\mathbf{A}^{-1} = (\mathbf{AB})^{-1}$.

Moreover, when \mathbf{X} is composed of orthonormal variables, such that $\mathbf{X}^T\mathbf{X} = \mathbf{I}_p$, it then follows that

$$\hat{\boldsymbol{\beta}}^{ridge} = \frac{1}{1 + \lambda} \hat{\boldsymbol{\beta}}^{OLS}$$

Bias and variance of ridge estimator

Ridge estimation produces a biased estimator of the true parameter $\boldsymbol{\beta}$. With the definition of $\hat{\boldsymbol{\beta}}^{ridge}$ and the model assumption $\mathbf{E}(\mathbf{Y}|\mathbf{X}) = \mathbf{X}\boldsymbol{\beta}$, we obtain,

$$\begin{aligned} \mathbf{E}(\hat{\boldsymbol{\beta}}^{ridge}|\mathbf{X}) &= (\mathbf{X}^T\mathbf{X} + \lambda\mathbf{I})^{-1}\mathbf{X}^T\mathbf{X}\boldsymbol{\beta} \\ &= (\mathbf{X}^T\mathbf{X} + \lambda\mathbf{I})^{-1}(\mathbf{X}^T\mathbf{X} + \lambda\mathbf{I} - \lambda\mathbf{I})\boldsymbol{\beta} \\ &= \boldsymbol{\beta} - \lambda(\mathbf{X}^T\mathbf{X} + \lambda\mathbf{I})^{-1}\boldsymbol{\beta} \end{aligned}$$

where the bias of the ridge estimator is proportional to λ . The variance of the ridge estimator is

$$\mathbf{Var}(\hat{\boldsymbol{\beta}}^{ridge}|\mathbf{X}) = \sigma^2(\mathbf{X}^T\mathbf{X} + \lambda\mathbf{I})^{-1}(\mathbf{X}^T\mathbf{X})(\mathbf{X}^T\mathbf{X} + \lambda\mathbf{I})^{-1}.$$

When λ increases, the inverted term $(\mathbf{X}^T\mathbf{X} + \lambda\mathbf{I})^{-1}$ is increasingly dominated by $\lambda\mathbf{I}$. The variance of the ridge estimator, therefore, is a decreasing function of λ . This result is intuitively reasonable because the estimator itself is driven toward $\mathbf{0}$.

Variance-bias tradeoff

The mean-squared error of an estimator can be decomposed into the sum of its squared bias and sampling variance.

$$\begin{aligned} MSE(\hat{\theta}) &= \mathbf{E}((\hat{\theta} - \theta)^2) = \mathbf{E}(\hat{\theta}^2) + \theta^2 - 2\theta\mathbf{E}(\hat{\theta}) \\ Bias^2(\hat{\theta}) &= [\mathbf{E}(\hat{\theta}) - \theta]^2 = \mathbf{E}^2(\hat{\theta}) + \theta^2 - 2\theta\mathbf{E}(\hat{\theta}) \\ Var(\hat{\theta}) &= \mathbf{E}(\hat{\theta}^2) - \mathbf{E}^2(\hat{\theta}) \end{aligned}$$

Therefore

$$MSE(\hat{\theta}) = Bias^2(\hat{\theta}) + Var(\hat{\theta})$$

The essential idea here is to trade a small amount of bias in the coefficient estimates for a large reduction in coefficient sampling variance. Hoerl and Kennard (1970) prove that it is always possible to choose a positive value of the ridge constant λ so that the mean-squared error of the ridge estimator is less than the mean-squared error of the least-squares estimator. These ideas are illustrated heuristically in Figure 25.1

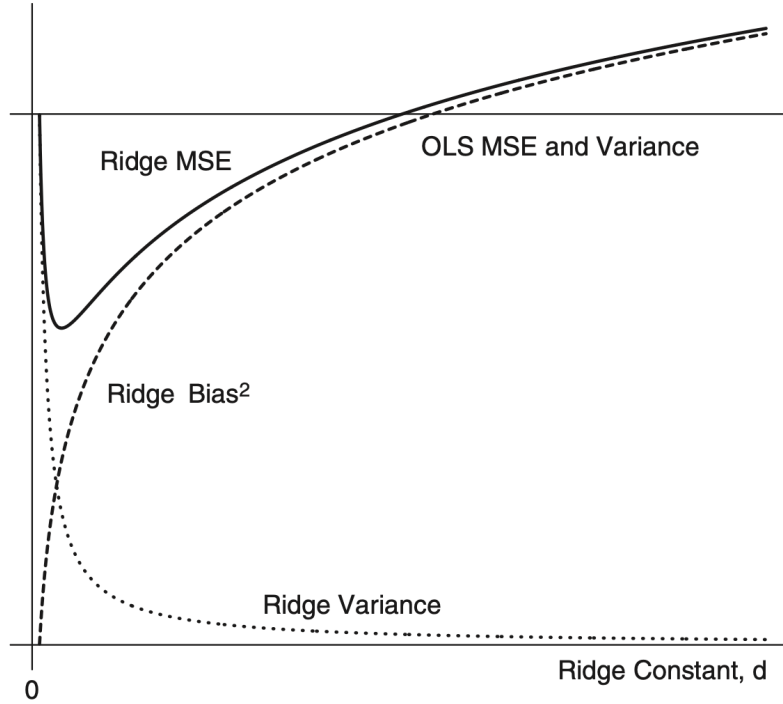


Figure 25.1: Trade-off of bias and against variance for the ridge-regression estimator. The horizontal line gives the variance of the least-squares (OLS) estimator; because the OLS estimator is unbiased, its variance and mean-squared error are the same. The broken line shows the squared bias of the ridge estimator as an increasing function of the ridge constant d (i.e. λ in our notes). The dotted line shows the variance of the ridge estimator. The mean-squared error (MSE) of the ridge estimator, given by the heavier solid line, is the sum of its variance and squared bias. For some values of d , the MSE error of the ridge estimator is below the variance of the OLS estimator. JF Figure 13.9.

Lasso regression

We have seen that ridge regression essentially re-scales the OLS estimates. The lasso, by contrast, tries to produce a *sparse* solution, in the sense that several of the slope parameters will be set to zero.

Constrained optimization

Different from the L_2 penalty for ridge regression, the Lasso regression employs L_1 -penalty.

$$\hat{\beta}^{lasso} = \arg \min_{\beta \in \mathbb{R}^p} \sum_{i=1}^n (y_i - \mathbf{x}_i^T \beta)^2$$

$$\text{subject to } \|\beta\|_1 = \sum_{j=1}^p |\beta_j| \leq t$$

for $t \geq 0$; which can again be re-formulated using the Lagrangian for the L_1 -penalty,

$$\hat{\boldsymbol{\beta}}^{lasso} = \arg \min_{\boldsymbol{\beta} \in \mathbb{R}^p} \left\{ \sum_{i=1}^n (y_i - \mathbf{x}_i^T \boldsymbol{\beta})^2 + \lambda \sum_{j=1}^p |\beta_j| \right\}$$

where $\lambda > 0$ and, as before, there exists a one-to-one correspondence between t and λ .

Parameter estimation

Contrary to ridge regression, the Lasso does not have a closed-form solution. The L_1 -penalty makes the solution non-linear in y_i 's. The above constrained minimization is a quadratic programming problem, for which many solvers exist.

Choice of Hyperparameters

Regularization parameter

The choice of λ in both ridge and lasso regressions is more of an art than a science. This parameter can be constructed as a complexity parameter, since as λ increases, less and less effective parameters are likely to be included in both ridge and lasso regressions. Therefore, one can adopt a model selection perspective and compare different choices of λ using cross-validation or an information criterion. That is, the value of λ should be chosen adaptively, in order to minimize an estimate of the expected prediction error (as in cross-validation), for instance, which is well approximated by AIC. We will discuss model selection in more detail later.

Bayesian perspective

The penalty terms in ridge and lasso regression can also be justified, using a Bayesian framework, whereby these terms arise as a result of the specification of a particular prior distribution on the vector of slope parameters.

1. The use of an L_2 -penalty in multiple regression is analogous to the choice of a Normal prior on the β_j 's, in Bayesian statistics.

$$\begin{aligned} y_i &\stackrel{iid}{\sim} \mathcal{N}(\beta_0 + \mathbf{x}_i^T \boldsymbol{\beta}, \sigma^2), \quad i = 1, \dots, n \\ \beta_j &\stackrel{iid}{\sim} \mathcal{N}(0, \tau^2), \quad j = 1, \dots, p \end{aligned}$$

2. Similarly, the use of an L_1 -penalty in multiple regression is analogous to the choice of a Laplace prior on the β_j 's, such that

$$\beta_j \stackrel{iid}{\sim} \text{Laplace}(0, \tau^2), \quad j = 1, \dots, p$$

In both cases, the value of the hyperparameter, τ^2 , will be inversely proportional to the choice of the particular value for λ . For ridge regression, λ is exactly equal to the shrinkage parameter of the hierarchical model, $\lambda = \sigma^2/\tau^2$.

Model selection

Model selection is conceptually simplest when our goal is *prediction* – that is, the development of a regression model that will predict new data as accurately as possible. However, prediction is not often the only desirable characteristic in a statistical model that model interpretation, data summary and explanations are also desired. We discuss several criteria for selecting among m competing statistical models $\mathcal{M} = \{M_1, M_2, \dots, M_m\}$ for n observations of a response variable Y and associated predictors X s.

Adjusted- R^2

The squared multiple correlation “corrected” (or “adjusted”) for degrees of freedom is intuitively reasonable criterion for comparing linear-regression models with different numbers of parameters. Suppose model M_j is one of the models under consideration. If M_j has s_j regression coefficients (including the regression constant) and is fit to a data set with n observations, then the adjusted- R^2 for the model is

$$R_{adj,j}^2 = 1 - \frac{n-1}{n-s_j} \times \frac{RSS_j}{TSS}$$

Models with relatively large numbers of parameters are penalized for their lack of parsimony. The model with the highest adjusted- R^2 value is selected as the best model. Beyond this intuitive rationale, however, there is no deep justification for using R_{adj}^2 as a model selection criterion.

Cross-validation and generalized cross-validation

The key idea in cross-validation (more accurately, leave-one-out cross-validation) is to omit the i th observation to obtain an estimate of $E(Y|x_i)$ based on the other observations as $\hat{Y}_{-i}^{(j)}$ for model M_j . Omitting the i th observation makes the fitted value $\hat{Y}_{-i}^{(j)}$ independent of the observed value Y_i . The cross-validation criterion for model M_j is

$$CV_j \equiv \frac{\sum_{i=1}^n \left[\hat{Y}_{-i}^{(j)} - Y_i \right]^2}{n}$$

We prefer the model with the smallest value of CV_j .

In linear least-squares regression, there are efficient procedures for computing the leave-one-out fitted values $\hat{Y}_{-i}^{(j)}$ that do not require literally refitting the model (recall the discussions of standardized residuals). However, in other applications, leave-one-out cross-validation can be computationally expensive (that requires literally refitting the model n times).

An alternative is to divide the data into a relatively small number of subsets of roughly equal size and to fit the model omitting one subset at a time, obtaining fitted values for all observations in the omitted subset. This method is termed as K -fold cross-validation where K is the number of subsets. The cross-validation criterion is defined the same way as before.

An alternative criterion is to approximate CV by the generalized cross-validation criterion

$$GCV_j \equiv \frac{n \times RSS_j}{df_{res_j}^2}$$

which however is less popular given the increasing computational power we have in the modern era.

AIC and BIC

The Akaike information criterion (AIC) and the Bayesian information criterion (BIC) are also popular model selection criteria. Both are members of a more general family of *penalized* model-fit statistics (in the form of “*IC”), applicable to regression models fit by maximum likelihood, that take the form

$$*IC_j = -2 \log_e L(\hat{\theta}_j) + cs_j$$

where $L(\hat{\theta}_j)$ is the maximized likelihood under model M_j ; $\hat{\theta}_j$ is the vector of parameters of the model (including, for example, regression coefficients and an error variance); s_j is the number of parameters in $\hat{\theta}_j$; and c is a constant that differs from one model selection criterion to another. The first term, $-2 \log_e L(\hat{\theta}_j)$, is the residual deviance under the model; for a linear model with normal errors, it is simply the residual sum of squares.

The model with the smallest *IC is the one that receives most support from the data (the selected model). The AIC and BIC are defined as follows:

$$\begin{aligned} AIC_j &\equiv -2 \log_e L(\hat{\theta}_j) + 2s_j \\ BIC_j &\equiv -2 \log_e L(\hat{\theta}_j) + s_j \log_e(n) \end{aligned}$$

The lack-of-parsimony penalty for the BIC grows with the sample size, while that for the AIC does not. When $n \geq 8$ the penalty for the BIC is larger than that for the AIC resulting in BIC tends to nominate models with fewer parameters. Both AIC and BIC are based on deeper statistical considerations, please refer to JF 22.1 sections **A closer look at the AIC** and **A closer look at the BIC** for more details.

Sequential procedures

Besides the ranking systems above, there is another class loosely defined as sequential procedures for model selection.

1. Forward selection
2. Backwards elimination
3. Stepwise selection

Forward selection :

1. Choose a threshold significance level for adding predictors, “SLENTRY” (SL stands for significance level). For example, $SLENTRY = 0.10$.
2. Initialize with $y = \beta_0 + \epsilon$.
3. Form a set of candidate models that differ from the working model by addition of one new predictor
4. Do any of the added predictors have $p - value \leq SLENTRY$?
 - Yes: add predictor with smallest p -value to working model + repeat steps 3 to 4.
 - No: stop. Final model = working model.

Backwards elimination

1. Choose threshold level for removing predictors. For example, $SLSTAY = 0.05$.
2. Initialize with most general model (biggest possible): $y = \beta_0 + \beta_1 x_1 + \dots + \epsilon$.
3. Form a set of candidate models that differ from working model by deletion of one term
4. Do any $p - value > SLSTAY$ (from fitting the current working model)?
 - Yes: remove the term with largest p -value and repeat steps 3 and 4.
 - No: stop. Final model = working model.

Stepwise Alternate forwards + backwards steps. Initialize with $y = \beta_0 + \epsilon$. Stop when consecutive forward + backward steps do not change working model. ($SLENTRY \leq SLSTAY$)

Some examples

- [Model selection by AIC](#)
- [Model selection by AIC and Lasso](#)

26 Lecture 26: April 10

Last time

- Collinearity (JF chapter 13, RD 8.3.2)
- Principal component analysis (JF 13.1.1, RD 8.3.4)
- Biased estimation:
 - Ridge Regression

Today

- Biased estimation:
 - Lasso Regression
- Model selection
- Analysis of Variance (JF chapter 8)
 - one-way anova
 - two-way anova

Lasso regression

We have seen that ridge regression essentially re-scales the OLS estimates. The lasso, by contrast, tries to produce a *sparse* solution, in the sense that several of the slope parameters will be set to zero.

Constrained optimization

Different from the L_2 penalty for ridge regression, the Lasso regression employs L_1 -penalty.

$$\begin{aligned}\hat{\boldsymbol{\beta}}^{lasso} &= \arg \min_{\boldsymbol{\beta} \in \mathbb{R}^p} \sum_{i=1}^n (y_i - \mathbf{x}_i^T \boldsymbol{\beta})^2 \\ \text{subject to } \|\boldsymbol{\beta}\|_1 &= \sum_{j=1}^p |\beta_j| \leq t\end{aligned}$$

for $t \geq 0$; which can again be re-formulated using the Lagrangian for the L_1 -penalty,

$$\hat{\boldsymbol{\beta}}^{lasso} = \arg \min_{\boldsymbol{\beta} \in \mathbb{R}^p} \left\{ \sum_{i=1}^n (y_i - \mathbf{x}_i^T \boldsymbol{\beta})^2 + \lambda \sum_{j=1}^p |\beta_j| \right\}$$

where $\lambda > 0$ and, as before, there exists a one-to-one correspondence between t and λ .

Parameter estimation

Contrary to ridge regression, the Lasso does not have a closed-form solution. The L_1 -penalty makes the solution non-linear in y_i 's. The above constrained minimization is a quadratic programming problem, for which many solvers exist.

Choice of Hyperparameters

Regularization parameter

The choice of λ in both ridge and lasso regressions is more of an art than a science. This parameter can be constructed as a complexity parameter, since as λ increases, less and less effective parameters are likely to be included in both ridge and lasso regressions. Therefore, one can adopt a model selection perspective and compare different choices of λ using cross-validation or an information criterion. That is, the value of λ should be chosen adaptively, in order to minimize an estimate of the expected prediction error (as in cross-validation), for instance, which is well approximated by AIC. We will discuss model selection in more detail later.

Bayesian perspective

The penalty terms in ridge and lasso regression can also be justified, using a Bayesian framework, whereby these terms arise as a result of the specification of a particular prior distribution on the vector of slope parameters.

1. The use of an L_2 -penalty in multiple regression is analogous to the choice of a Normal prior on the β_j 's, in Bayesian statistics.

$$\begin{aligned} y_i &\stackrel{iid}{\sim} \mathcal{N}(\beta_0 + \mathbf{x}_i^T \boldsymbol{\beta}, \sigma^2), \quad i = 1, \dots, n \\ \beta_j &\stackrel{iid}{\sim} \mathcal{N}(0, \tau^2), \quad j = 1, \dots, p \end{aligned}$$

2. Similarly, the use of an L_1 -penalty in multiple regression is analogous to the choice of a Laplace prior on the β_j 's, such that

$$\beta_j \stackrel{iid}{\sim} \text{Laplace}(0, \tau^2), \quad j = 1, \dots, p$$

In both cases, the value of the hyperparameter, τ^2 , will be inversely proportional to the choice of the particular value for λ . For ridge regression, λ is exactly equal to the shrinkage parameter of the hierarchical model, $\lambda = \sigma^2/\tau^2$.

Model selection

Model selection is conceptually simplest when our goal is *prediction* – that is, the development of a regression model that will predict new data as accurately as possible. However, prediction is not often the only desirable characteristic in a statistical model that model interpretation, data summary and explanations are also desired. We discuss several criteria for selecting among m competing statistical models $\mathcal{M} = \{M_1, M_2, \dots, M_m\}$ for n observations of a response variable Y and associated predictors X s.

Adjusted- R^2

The squared multiple correlation “corrected” (or “adjusted”) for degrees of freedom is intuitively reasonable criterion for comparing linear-regression models with different numbers of parameters. Suppose model M_j is one of the models under consideration. If M_j has s_j regression coefficients (including the regression constant) and is fit to a data set with n observations, then the adjusted- R^2 for the model is

$$R_{adj,j}^2 = 1 - \frac{n-1}{n-s_j} \times \frac{RSS_j}{TSS}$$

Models with relatively large numbers of parameters are penalized for their lack of parsimony. The model with the highest adjusted- R^2 value is selected as the best model. Beyond this intuitive rationale, however, there is no deep justification for using R_{adj}^2 as a model selection criterion.

Cross-validation and generalized cross-validation

The key idea in cross-validation (more accurately, leave-one-out cross-validation) is to omit the i th observation to obtain an estimate of $E(Y|x_i)$ based on the other observations as $\hat{Y}_{-i}^{(j)}$ for model M_j . Omitting the i th observation makes the fitted value $\hat{Y}_{-i}^{(j)}$ independent of the observed value Y_i . The cross-validation criterion for model M_j is

$$CV_j \equiv \frac{\sum_{i=1}^n \left[\hat{Y}_{-i}^{(j)} - Y_i \right]^2}{n}$$

We prefer the model with the smallest value of CV_j .

In linear least-squares regression, there are efficient procedures for computing the leave-one-out fitted values $\hat{Y}_{-i}^{(j)}$ that do not require literally refitting the model (recall the discussions of standardized residuals). However, in other applications, leave-one-out cross-validation can be computationally expensive (that requires literally refitting the model n times).

An alternative is to divide the data into a relatively small number of subsets of roughly equal size and to fit the model omitting one subset at a time, obtaining fitted values for all observations in the omitted subset. This method is termed as K -fold cross-validation where K is the number of subsets. The cross-validation criterion is defined the same way as before.

An alternative criterion is to approximate CV by the generalized cross-validation criterion

$$GCV_j \equiv \frac{n \times RSS_j}{df_{res_j}^2}$$

which however is less popular given the increasing computational power we have in the modern era.

AIC and BIC

The Akaike information criterion (AIC) and the Bayesian information criterion (BIC) are also popular model selection criteria. Both are members of a more general family of *penalized* model-fit statistics (in the form of “*IC”), applicable to regression models fit by maximum likelihood, that take the form

$$*IC_j = -2 \log_e L(\hat{\theta}_j) + cs_j$$

where $L(\hat{\theta}_j)$ is the maximized likelihood under model M_j ; $\hat{\theta}_j$ is the vector of parameters of the model (including, for example, regression coefficients and an error variance); s_j is the number of parameters in $\hat{\theta}_j$; and c is a constant that differs from one model selection criterion to another. The first term, $-2 \log_e L(\hat{\theta}_j)$, is the residual deviance under the model; for a linear model with normal errors, it is simply the residual sum of squares.

The model with the smallest *IC is the one that receives most support from the data (the selected model). The AIC and BIC are defined as follows:

$$\begin{aligned} AIC_j &\equiv -2 \log_e L(\hat{\theta}_j) + 2s_j \\ BIC_j &\equiv -2 \log_e L(\hat{\theta}_j) + s_j \log_e(n) \end{aligned}$$

The lack-of-parsimony penalty for the BIC grows with the sample size, while that for the AIC does not. When $n \geq 8$ the penalty for the BIC is larger than that for the AIC resulting in BIC tends to nominate models with fewer parameters. Both AIC and BIC are based on deeper statistical considerations, please refer to JF 22.1 sections **A closer look at the AIC** and **A closer look at the BIC** for more details.

Sequential procedures

Besides the ranking systems above, there is another class loosely defined as sequential procedures for model selection.

1. Forward selection
2. Backwards elimination
3. Stepwise selection

Forward selection :

1. Choose a threshold significance level for adding predictors, “SLENTRY” (SL stands for significance level). For example, $SLENTRY = 0.10$.
2. Initialize with $y = \beta_0 + \epsilon$.
3. Form a set of candidate models that differ from the working model by addition of one new predictor
4. Do any of the added predictors have $p - value \leq SLENTRY$?
 - Yes: add predictor with smallest p -value to working model + repeat steps 3 to 4.
 - No: stop. Final model = working model.

Backwards elimination

1. Choose threshold level for removing predictors. For example, $SLSTAY = 0.05$.
2. Initialize with most general model (biggest possible): $y = \beta_0 + \beta_1 x_1 + \dots + \epsilon$.
3. Form a set of candidate models that differ from working model by deletion of one term
4. Do any $p - value > SLSTAY$ (from fitting the current working model)?
 - Yes: remove the term with largest p -value and repeat steps 3 and 4.
 - No: stop. Final model = working model.

Stepwise Alternate forwards + backwards steps. Initialize with $y = \beta_0 + \epsilon$. Stop when consecutive forward + backward steps do not change working model. ($SLENTRY \leq SLSTAY$)

Some examples

- [Model selection by AIC](#)
- [Model selection by AIC and Lasso](#)

Additional reference

[Course notes](#) by Dr. Jason Osborne.

Analysis of Variance

The term analysis of variance is used to describe the partition of the response-variable sum of squares into “explained” and “unexplained” components, noting that this decomposition applies generally to linear models. For historical reasons, analysis of variance (abbreviated ANOVA) also refers to procedures for fitting and testing linear models in which the explanatory variables are categorical.

One-way ANOVA

Suppose that there are *no* quantitative explanatory variables, but only a single factor (categorical data). For example, for a three-category classification, we have the model

$$Y_i = \alpha + \gamma_1 D_{i1} + \gamma_2 D_{i2} + \epsilon_i \quad (3)$$

employing the following coding for the dummy regressors:

Group	D_1	D_2
1	1	0
2	0	1
3	0	0

The expectation of the response variable in each group (i.e. in each category or level of the factor) is the population group mean, denoted by μ_j for the j th group. Equation ?? produces the following relationship between group means and model parameters:

$$\text{Group 1: } E(Y_i | D_{i1} = 1, D_{i2} = 0) = \alpha + \gamma_1 \times 1 + \gamma_2 \times 0 = \alpha + \gamma_1$$

$$\text{Group 2: } E(Y_i | D_{i1} = 0, D_{i2} = 1) = \alpha + \gamma_1 \times 0 + \gamma_2 \times 1 = \alpha + \gamma_2$$

$$\text{Group 3: } E(Y_i | D_{i1} = 0, D_{i2} = 0) = \alpha + \gamma_1 \times 0 + \gamma_2 \times 0 = \alpha$$

There are three parameters (α , γ_1 and γ_2) and three group means, so we can solve uniquely for the parameters in terms of the group means:

$$\hat{\alpha} = \mu_3$$

$$\hat{\gamma}_1 = \mu_1 - \mu_3$$

$$\hat{\gamma}_2 = \mu_2 - \mu_3$$

Not surprisingly, α represents the mean of the baseline category (Group 3) and that γ_1 and γ_2 captures differences between the other group means and the mean of the baseline category.

notations

Because observations are partitioned according to groups, it is convenient to let Y_{jk} denote the k th observation within the j th of m groups. The number of observations in the j th group is n_j , and the total number of observations is $n = \sum_{j=1}^m n_j$. Let $\mu_j \equiv E(Y_{jk})$ be the population mean in group j .

The one-way ANOVA model is

$$Y_{jk} = \mu + \alpha_j + \epsilon_{jk}$$

where μ represents the general level of response variable in the population; α_j represents the effect on the response variable of membership in the j th group; ϵ_{jk} is an error variable that follows the usual linear-model assumptions: $\epsilon_{jk} \stackrel{iid}{\sim} \mathcal{N}(0, \sigma^2)$.

By taking expectations, we have

$$\mu_j = \mu + \alpha_j$$

The parameters of the model are, therefore, underdetermined, for there are $m + 1$ parameters (including μ) but only m population group means (recall the dummy variable trap introduced in collinearity). To produce easily interpretable parameters and that estimates and generalizes usefully to more complex models, we impose the sum-to-zero constraint

$$\sum_{j=1}^m \alpha_j = 0$$

With the sum-to-zero constraint, we solve for the parameters

$$\begin{aligned}\hat{\mu} &= \frac{\sum \mu_j}{m} \\ \hat{\alpha}_j &= \mu_j - \mu\end{aligned}$$

The fitted Y values are the group means for the one-way ANOVA model:

$$\hat{Y}_{jk} = \hat{\mu} + \hat{\alpha}_j$$

and the regression and residual sums of squares therefore take particularly simple forms in one-way ANOVA:

$$\begin{aligned}RegSS &= \sum_{j=1}^m \sum_{k=1}^{n_j} (\hat{Y}_{jk} - \bar{Y})^2 = \sum_{j=1}^m n_j (\bar{Y}_j - \bar{Y})^2 \\ RSS &= \sum_{j=1}^m \sum_{k=1}^{n_j} (Y_{jk} - \hat{Y}_{jk})^2 = \sum_{j=1}^m \sum_{k=1}^{n_j} (Y_{jk} - \bar{Y}_j)^2\end{aligned}$$

and can be presented in an ANOVA table.

Table 2: General one-way ANOVA table

Source	Sum of Squares	df	Mean Square	F	H_0
Groups	$\sum n_j (\bar{Y}_j - \bar{Y})^2$	$m - 1$	$\frac{RegSS}{m-1}$	$\frac{RegMS}{RMS}$	$\alpha_1 = \dots = \alpha_m = 0$
Residuals	$\sum \sum (Y_{jk} - \bar{Y}_j)^2$	$n - m$	$\frac{RSS}{n-m}$		
Total	$\sum \sum (Y_{jk} - \bar{Y})^2$	$n - 1$			

Sometimes, the column of Source can also be denoted with Treatments (for Groups) and Error (for Residuals). And a balanced one-way ANOVA model has the same number of observations in one group (or treatment), in other words, $n_1 = \dots = n_m = \frac{n}{m}$.

one-way ANOVA example

The following data come from study investigating binding fraction for several antibiotics using $n = 20$ bovine serum samples:

Antibiotic	Binding Percentage	Sample mean
Penicillin G	29.6 24.3 28.5 32.0	28.6
Tetracyclin	27.3 32.6 30.8 34.8	31.4
Streptomycin	5.8 6.2 11.0 8.3	7.8
Erythromycin	21.6 17.4 18.3 19	19.1
Chloramphenicol	29.2 32.8 25.0 24.2	27.8

Question: Are the population means for these 5 treatments plausibly equal?

Answer:

One model parameterizes antibiotic effects as differences from mean

$$Y_{jk} = \mu + \alpha_j + \epsilon_{jk}$$

for $j = 1, \dots, 5$ and $k = 1, \dots, 4$, where $\epsilon_{jk} \stackrel{iid}{\sim} \mathcal{N}(0, \sigma^2)$ errors.

Unknown parameters

1. μ - overall population mean (average of 5 treatment population means)
2. α_j - difference between (population) mean for treatment j and μ
3. σ^2 - (population) variance of binding fraction for a given antibiotic

To test $H_0 : \alpha_1 = \alpha_2 = \dots = \alpha_5 = 0$, we just carry out one-way ANOVA:

Source	Sum of Squares	df	Mean Square	F
Groups	1481	4	370	41
Residuals	136	15	9	
Total	1617	19		

Compared to $F(0.05, 4, 15) = 3.06$, we have $F = 41 > 3.06$.

Conclusion: we reject the null hypothesis of all population means for 5 treatment being equal at 0.05 significance level.

How do we obtain standard errors of parameter estimates? (HW)

Two-Way ANOVA

The inclusion of a second factor permits us to model and test partial relationships, as well as to introduce interactions. Let's take a look at the patterns of relationship that can occur when a quantitative response variable is classified by two factors.

Patterns of Means in the two-way classification

Consider the following table:

	C_1	C_2	\dots	C_c	
R_1	μ_{11}	μ_{12}	\dots	μ_{1c}	$\mu_{1\cdot}$
R_2	μ_{21}	μ_{22}	\dots	μ_{2c}	$\mu_{2\cdot}$
\vdots	\vdots	\vdots		\vdots	\vdots
R_r	μ_{r1}	μ_{r2}	\dots	μ_{rc}	$\mu_{r\cdot}$
	$\mu_{\cdot 1}$	$\mu_{\cdot 2}$	\dots	$\mu_{\cdot c}$	$\mu_{\cdot\cdot}$

The factors, R and C (for “rows” and “columns” of the table of means), have r and c categories, respectively. The factor categories are denoted R_j and C_k . Within each cell of the design - that is, for each combination of categories $\{R_j, C_k\}$ of the two factors - there is a population cell mean μ_{jk} for the response variable. Extending the dot notation, we have

$$\mu_{j\cdot} \equiv \frac{\sum_{k=1}^c \mu_{jk}}{c}$$

is the marginal mean of the response variable in row j .

$$\mu_{\cdot k} \equiv \frac{\sum_{j=1}^r \mu_{jk}}{r}$$

is the marginal mean in column k . And

$$\mu_{\cdot\cdot} \equiv \frac{\sum_j \sum_k \mu_{jk}}{r \times c}$$

is the grand mean.

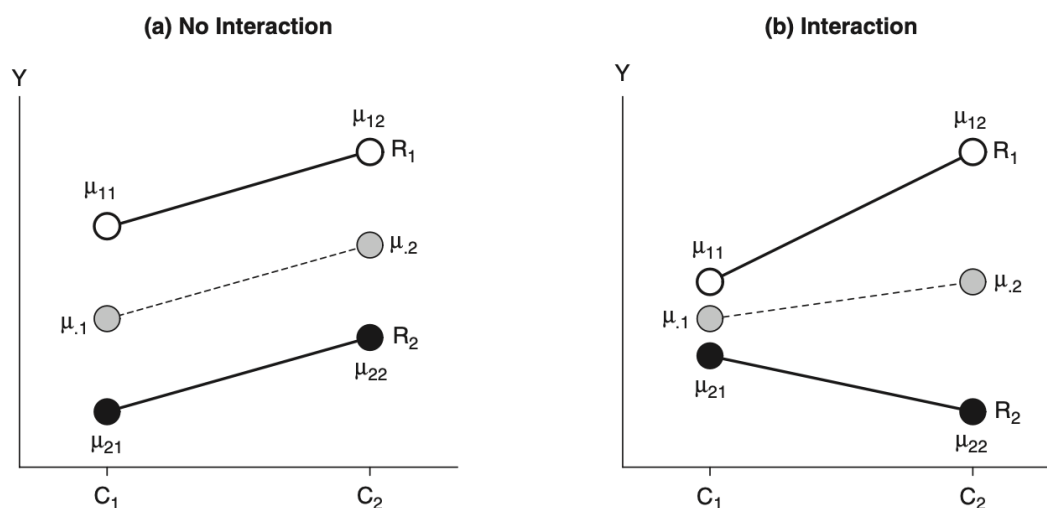


Figure 26.1: Interaction in the two-way classification. In (a), the parallel profiles of means (given by the white and black circles connected by solid lines) indicate that R and C do not interact in affecting Y . The R -effect – that is, the difference between the two profiles – is the same at both C_1 and C_2 . Likewise, the C -effect – that is, the rise in the line from C_1 to C_2 – is the same for both profiles. In (b), the R -effect differs at the two categories of C , and the C -effect differs at the two categories of R : R and C interact in affecting Y . In both graphs, the column marginal means $\mu_{.1}$ and $\mu_{.2}$ are shown as averages of the cell means in each column (represented by the gray circles connected by broken lines). JF Figure 8.2.