

Math 6040/7260 Linear Models

Mon/Wed/Fri 11:00am - 11:50am

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

1 Lecture 1: Jan 26

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
- **Good to have:** Mathematical Statistics, 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: $\mathbf{y} \sim N(\mathbf{X}\boldsymbol{\beta}, \sigma_1^2 \mathbf{V}_1 + \sigma_2^2 \mathbf{V}_2 + \cdots + \sigma_r^2 \mathbf{V}_r)$ with $\mathbf{V}_1, \mathbf{V}_2, \dots, \mathbf{V}_r$ known.

- General mixed linear Model:

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

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

- 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-7260-2022.github.io/>

HW submission

Through Github with demo on Friday class.

Presentations

Let me know your pick by the end of Friday (01/28/2022).

Last year 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 28

Last time

- Introduction
- Course logistics

Today

- Reply to the “Presentation Dates” thread on Canvas by the end of today.
- Introduce yourself (remind remote students to record a short video)
 - basic info (name, department, year, ...)
 - why taking this course
- Git
- Linear algebra: vector and vector space, rank of a matrix (maybe)

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>.



Figure 2.1

– 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.

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-6040-2022-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-6040-2022-spring`.
- after finishing the rest of the questions, save your file inside your git repository folder `math-6040-2022-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 25

Last time

- Git

Today

- 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.)
- $\text{rank}(\mathbf{AB}) \leq \min\{\text{rank}(\mathbf{A}), \text{rank}(\mathbf{B})\}$. (Hint: Columns of \mathbf{AB} are spanned by columns of \mathbf{A} and rows of \mathbf{AB} are spanned by rows of \mathbf{B} .)
- if $\mathbf{Ax} = \mathbf{0}_m$ 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,

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

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 $\mathbf{A} = \mathbf{BC}$.

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.

4 Lecture 4: Feb 2

Last time

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

Today

- Probability review

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

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

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

proof:

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

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

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

6 Lecture 6: Feb 7

Last time

- Probability review

Today

- R basics
- Probability review
- Basic statistical concepts (PVR Chapter 1 - 2)
- Simple Linear Regression (JF Chapter 5)

Some common random variables

Discrete random variables

- $X \sim \text{Bernoulli}(p)$ (where $0 \leq p \leq 1$):

$$\Pr(x) = \begin{cases} p & \text{if } x = 1 \\ 1 - p & \text{if } x = 0 \end{cases}$$

- $X \sim \text{Binomial}(n, p)$ (where $0 \leq p \leq 1$):

$$\Pr(x) = \binom{n}{x} p^x (1 - p)^{n-x}$$

- $X \sim \text{Geometric}(p)$ (where $0 \leq p \leq 1$):

$$\Pr(x) = p(1 - p)^{x-1}$$

- $X \sim \text{Poisson}(\lambda)$ (where $\lambda > 0$):

$$\Pr(x) = e^{-\lambda} \frac{\lambda^x}{x!}$$

Continuous random variables

- $X \sim \text{Uniform}(a, b)$ (where $a < b$):

$$f(x) = \begin{cases} \frac{1}{b-a} & \text{if } a \leq x \leq b \\ 0 & \text{otherwise} \end{cases}$$

- $X \sim \text{Exponential}(\lambda)$ (where $\lambda > 0$):

$$f(x) = \begin{cases} \lambda e^{-\lambda x} & \text{if } x \geq 0 \\ 0 & \text{otherwise} \end{cases}$$

- $X \sim \text{Normal}(\mu, \sigma^2)$:

$$f(x) = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{1}{2\sigma^2}(x-\mu)^2}$$

The following table provides a summary of some of the properties of these distributions.

Distribution	PDF or PMF	Mean	Variance
<i>Bernoulli</i> (p)	$\begin{cases} p & \text{if } x = 1 \\ 1 - p & \text{if } x = 0 \end{cases}$	p	$p(1 - p)$
<i>Binomial</i> (n, p)	$\binom{n}{x} p^x (1 - p)^{n-x}$, for $0 \leq x \leq n$	np	$np(1 - p)$
<i>Geometric</i> (p)	$p(1 - p)^{x-1}$, for $k = 1, 2, \dots$	$\frac{1}{p}$	$\frac{1-p}{p^2}$
<i>Poisson</i> (λ)	$e^{-\lambda} \frac{\lambda^x}{x!}$, for $k = 1, 2, \dots$	λ	λ
<i>Uniform</i> (a, b)	$\frac{1}{b-a} I(a \leq x \leq b)$	$\frac{a+b}{2}$	$\frac{(b-a)^2}{12}$
<i>Gaussian</i> (μ, σ^2)	$\frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{1}{2\sigma^2}(x-\mu)^2}$	μ	σ^2
<i>Exponential</i> (λ)	$\lambda e^{-\lambda x} I(x \geq 0)$	$\frac{1}{\lambda}$	$\frac{1}{\lambda^2}$

Statistics: its objectives and scope (PVR Chapter 1)

We will use the word *statistics* in a broader sense:

Statistics refers to a body of scientific principles and methodologies that are useful for obtaining information about a phenomenon or a large collection of items. Statistical methods are techniques for using limited amounts of information to arrive at conclusions – called statistical inferences – about the phenomenon or the collection of items of interest.

Population and sample population

A *population* (sometimes referred to as a statistical population) is a collection (or aggregate) of measurements about which an inference is desired.

Example: An investigator is interested in evaluating the relationship between age, blood sugar level, and blood cholesterol level of insulin-dependent diabetics who are on a special experimental diet. The investigator wants to answer the following questions, among others:

1. How does the blood cholesterol level change with age and blood sugar level?
2. Are higher cholesterol levels associated with higher sugar levels?
3. Do older diabetics tend to have higher sugar and cholesterol levels?

What is the population of interest in this example? The population of interest is a collection of measurements – each of which consists of three values (age, blood sugar level, and blood cholesterol level) – for an insulin-dependent diabetic who is on the experimental diet.

Not that, in statistics, a measurement is one of the elements that form the population. In certain populations, each measurement may consist of several values. Populations in which each measurement

- is a single value are called *univariate* populations
- contains more than one value is called a *multivariate* population.

Sample and sample size

A *sample* consists of a finite number of measurements chosen from a population. The number of measurements in a sample is called the *sample size*.

Example: Answers to the questions about associations between the age, blood sugar level, and blood cholesterol level of diabetics can be based on measurements made on a sample of, say, $n = 40$ treated insulin-dependent diabetics. Such a sample is a collection of 40 measurements, each of which consists of three values: the age, blood sugar level, and blood cholesterol level of a treated patient.

Statistical components of a research study

A typical research study consists of three stages. The statistical techniques useful in these three stages are commonly known as *statistical methods in research*, and can be divided into three groups:

1. Methods for designing the research study
2. Methods for organizing and summarizing data
3. Methods for making inferences

In this course, we focus on the third stage that is to use the information in the samples to make conclusions about populations (i.e. making inferences). The key statistical issue in such inferences is their accuracy.

Example: Suppose that the average indoor radiation level in a sample of 15 homes built on reclaimed phosphate mine lands is 0.032 WL (working level is a historical unit of concentration of radioactive decay products of radon). Then 0.032 WL could be regarded as an estimate of the average indoor level in all homes built on reclaimed phosphate mine lands.

- How accurate is this estimate?
- Suppose there are a total of 4000 homes built on reclaimed lands, is our small sample representative of the whole population?
- If our sample could be regarded as representative of the population, it would be reasonable to expect that the difference between the estimated value of 0.032 WL and the true mean radiation level for all homes will be small, but how do we get/estimate the actual magnitude of this difference?

The natural question is whether it is possible to assess, with reasonable certainty, the magnitude of the error in our estimate. For example, can we say, with a reasonable degree of confidence, that the average level for the population of all homes will be within 0.001 WL of the average value calculated from sample homes?

Types of populations (PVR Chapter 2.2)

Statistical populations can be classified into categories depending upon the characteristics of the measurements contained in them.

- Univariate and multivariate populations
 - In a *univariate* population, each measurement consists of a single value
 - In a *multivariate* population, measurements consist of more than one value
- Real and conceptual populations
 - The population of 4000 indoor radon levels is a real population
 - The population of digestibility values for sheep fed June-harvested Pensacola Bahia grass is a conceptual population.
- Finite and infinite populations
 - A population may contain only a finite number of measurements, as in the case of the population of indoor radon levels of 4000 homes

- A population may have infinitely many measurements, as in the case of a conceptual population of potential digestibility measurements, in which every value in the interval $[0\%, 100\%]$ is a possible value of a measurement in the population.
- Quantitative and qualitative populations
 - A measurement is said to be *quantitative* if its value can be interpreted on a natural and meaningful scale
 - A measurement is *qualitative* if its value serves the sole purpose of identifying an object or a characteristic. The value of a qualitative measurement has no numerical implications.
- Discrete and continuous populations
 - A population is said to be *discrete* if the distinct values of the measurements contained in it can be arranged in a sequence.
 - A continuous population consists of measurements that take all the values in one or more intervals of a real line.

Simple linear regression

Figure 6.1 shows Davis's data on the measured and reported weight in kilograms of 101 women who were engaged in regular exercise.



Figure 6.1: 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

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

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

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

7 Lecture 7: Feb 9

Last time

- SLR in Matrix Form

Today

- Simple correlation
- The statistical model of the SLR (JF chapter 6)
- Properties of the Least-Squares estimator

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

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

Other quantities in Matrix Form

Fitted values

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

Hat matrix

$$\hat{\mathbf{Y}} = \mathbf{X}\hat{\boldsymbol{\beta}}$$

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

$$\hat{\mathbf{Y}} = \mathbf{H}\mathbf{Y}$$

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

Davis’s data example

For Davis’s data, we have

$$n = 101$$

$$\bar{y} = \frac{5780}{101} = 57.228$$

$$\bar{x} = \frac{5731}{101} = 56.743$$

$$\sum (x_i - \bar{x})(y_i - \bar{y}) = 4435.9$$

$$\sum (x_i - \bar{x})^2 = 4539.3,$$

so that

$$\hat{\beta}_1 = \frac{4435.9}{4539.3} = 0.97722$$

$$\hat{\beta}_0 = 57.228 - 0.97722 \times 56.743 = 1.7776$$

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 $y = x$. The solid line gives the least squares line $y = \hat{\beta}_0 + \hat{\beta}_1 x$.

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

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.
- $\hat{\beta}_1 = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{\sum (x_i - \bar{x})^2} = \frac{s_{xy}}{s_x^2} = r_{xy} \frac{s_y}{s_x}$

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



Figure 7.2: 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 .

9 Lecture 9: Feb 14

Last time

- Simple correlation
- The statistical model of the SLR (JF chapter 6)

Today

- Properties of the Least-Squares estimator
- Inference of SLR model

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

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

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 $H_0 : \beta_1 = \beta_{slope_0}$ that the population slope is equal to a specific value β_{slope_0} (most commonly, the null hypothesis has $\beta_{slope_0} = 0$), we calculate the test statistic (T -statistics) with $df = n - 2$

$$t_0 = \frac{\hat{\beta}_1 - \beta_{slope_0}}{\widehat{SE}(\hat{\beta}_1)} \sim t_{n-2}$$

10 Lecture 10: Feb 16

Last time

- Properties of the Least-Squares estimator

Today

- HW1 due Feb 18
- Inference of SLR model
- Multiple linear regression

Properties of the Least-Squares estimator

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

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

Statistical inference of the SLR model

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

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

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

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

$$\begin{aligned} \text{Var}(y_i) &= \sigma_\epsilon^2 \\ \text{Var}(\hat{y}_i) &= \text{Var}(\bar{y} + \hat{\beta}_1(x_i - \bar{x})) \\ &= \text{Var}(\bar{y}) + (x_i - \bar{x})^2 \text{Var}(\hat{\beta}_1) + 2(x_i - \bar{x}) \text{Cov}(\bar{y}, \hat{\beta}_1) \\ \text{Cov}(\bar{y}, \hat{\beta}_1) &= \text{Cov}(\bar{y}, \sum k_i y_i) \\ &= \sum_i \text{Cov}(\bar{y}, k_i y_i) \\ &= \sum_i \frac{k_i}{n} \text{Var}(y_i) \\ &= \frac{1}{n} \sum k_i \\ &= 0 \\ \therefore \text{Var}(\hat{y}_i) &= \text{Var}(\bar{y}) + (x_i - \bar{x})^2 \text{Var}(\hat{\beta}_1) \\ &= \frac{1}{n} \sigma_\epsilon^2 + \frac{\sigma_\epsilon^2 (x_i - \bar{x})^2}{\sum (x_i - \bar{x})^2} \\ &= \sigma_\epsilon^2 \left[\frac{1}{n} + \frac{(x_i - \bar{x})^2}{\sum (x_i - \bar{x})^2} \right] \end{aligned}$$

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

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

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 $H_0 : \beta_1 = \beta_{slope_0}$ that the population slope is equal to a specific value β_{slope_0} (most commonly, the null hypothesis has $\beta_{slope_0} = 0$), we calculate the test statistic (T -statistics) with $df = n - 2$

$$t_0 = \frac{\hat{\beta}_1 - \beta_{slope_0}}{\widehat{SE}(\hat{\beta}_1)} \sim t_{n-2}$$

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

12 Lecture 12: Feb 21

Last time

- Inference of SLR model

Today

- remember to send me the link to HW1 tag
- Lab 2 review
- Multiple linear regression

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 $SS[E]$:

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

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

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*.

13 Lecture 13: Feb 23

Last time

- Lab 2 review
- Multiple linear regression

Today

- HW1 review next week
- HW2 posted, due March 4th
- Inference of MLR
- more review on probability

Matrix formulation of MLR

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

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

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:

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

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 $(p + 1) \times 1$
- $\boldsymbol{\epsilon}$ denotes an error vector of size $n \times 1$

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

$$y_i = \mathbf{x}_{i\cdot} \boldsymbol{\beta} + \epsilon_i.$$

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

$$\begin{aligned}
\mathbf{X}^T \mathbf{X} &= \begin{bmatrix} 45 & 1884 & 2365 \\ 1884 & 105148 & 122197 \\ 2365 & 122197 & 163265 \end{bmatrix} \\
(\mathbf{X}^T \mathbf{X})^{-1} &= \begin{bmatrix} 0.10211 & -0.00085 & -0.00084 \\ -0.00085 & 0.00008 & -0.00005 \\ -0.00084 & -0.00005 & 0.00005 \end{bmatrix} \\
(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{Y} &= \begin{bmatrix} -6.0646629 \\ 0.5987328 \\ 0.5458339 \end{bmatrix} = ? \\
SS[E] = \boldsymbol{\epsilon}^T \boldsymbol{\epsilon} = (\mathbf{Y} - \hat{\mathbf{Y}})^T (\mathbf{Y} - \hat{\mathbf{Y}}) &= 7506.7 \\
MS[E] = \frac{SS[E]}{df} = \frac{7506.7}{45 - 2 - 1} &= 178.73 \\
\hat{\boldsymbol{\Sigma}} = MS[E](\mathbf{X}^T \mathbf{X})^{-1} &= \begin{bmatrix} 18.249481 & -0.151845008 & -0.150706025 \\ -0.151845 & 0.014320275 & -0.008518551 \\ -0.150706 & -0.008518551 & 0.009653582 \end{bmatrix}
\end{aligned}$$

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:

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

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.,

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

Or more generally, test the null hypothesis

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

where $0 \leq q_1 < q_2 < \cdots < q_k \leq p$ 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 13.1.

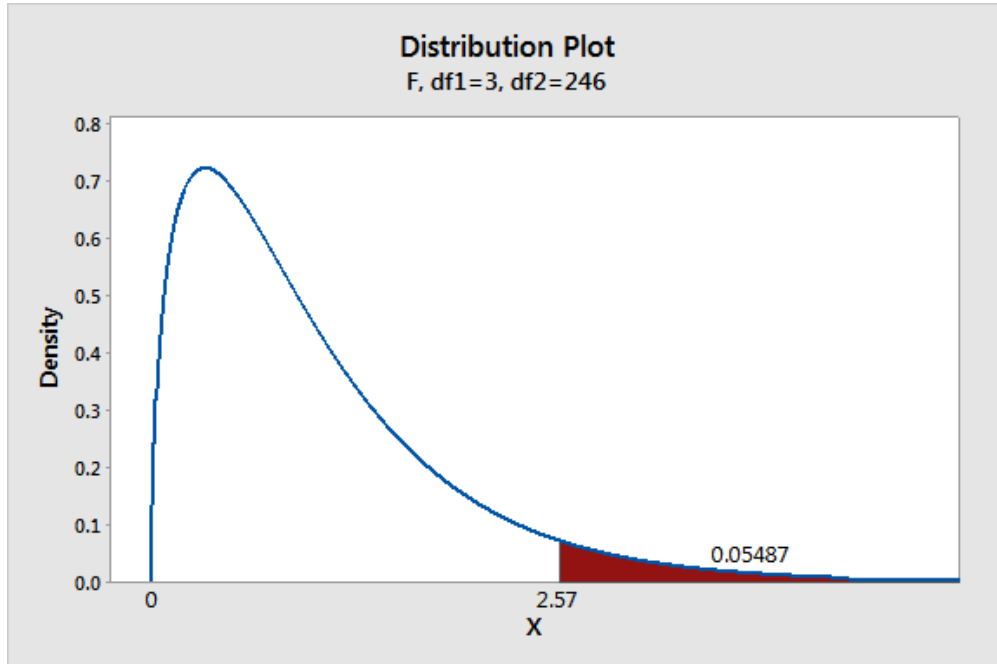


Figure 13.1: An example for p -value for F -statistic value 2.57 with an $F_{3,246}$ distribution

15 Lecture 15: March 2

Last time

- Inference of MLR

Today

- HW2 posted, due March 11th
- HW1 review
- more review on probability
- Dummy-Variable regression
- Interactions

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 $X^2 \equiv Z_1^2 + Z_2^2 + \dots + Z_k^2 \sim \chi_k^2$ (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 $\mathbf{Var}(\mathbf{X}) = \mathbf{Var}(\mathbf{Y}) + \mathbf{Var}(\mathbf{Z})$

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 15.1 (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.

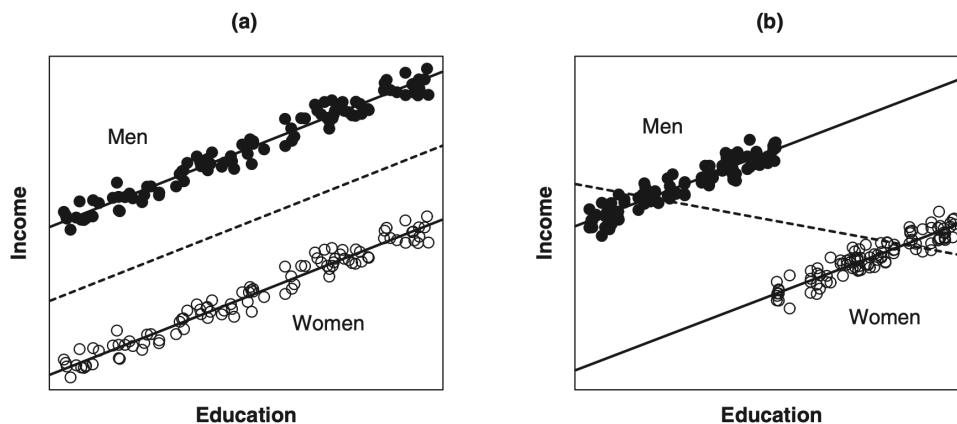


Figure 15.1: Idealized data representing the relationship between income and education for populations of men (filled circles) and women (open circles). In (a), there is no relationship between education and gender; in (b), women have a higher average level of education than men. In both (a) and (b), 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:

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

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 15.2.

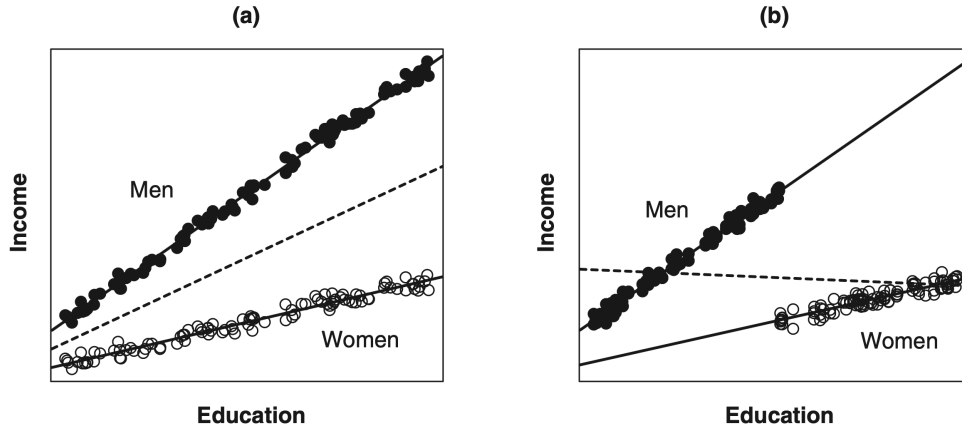


Figure 15.2: Idealized data representing the relationship between income and education for populations of men (filled circles) and women (open circles). In (a), there is no relationship between education and gender; in (b), women have a higher average level of education than men. In both (a) and (b), 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 15.2 (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:

Professional:	$Y_i =$	$(\beta_0 + \gamma_1) +$	$(\beta_1 + \delta_{11})X_{i1} +$	$(\beta_2 + \delta_{21})X_{i2} +$	ϵ_i
White collar:	$Y_i =$	$(\beta_0 + \gamma_2) +$	$(\beta_1 + \delta_{12})X_{i1} +$	$(\beta_2 + \delta_{22})X_{i2} +$	ϵ_i
Blue collar:	$Y_i =$	$\beta_0 +$	$\beta_1 X_{i1} +$	$\beta_2 X_{i2} +$	ϵ_i