

Introductory Review

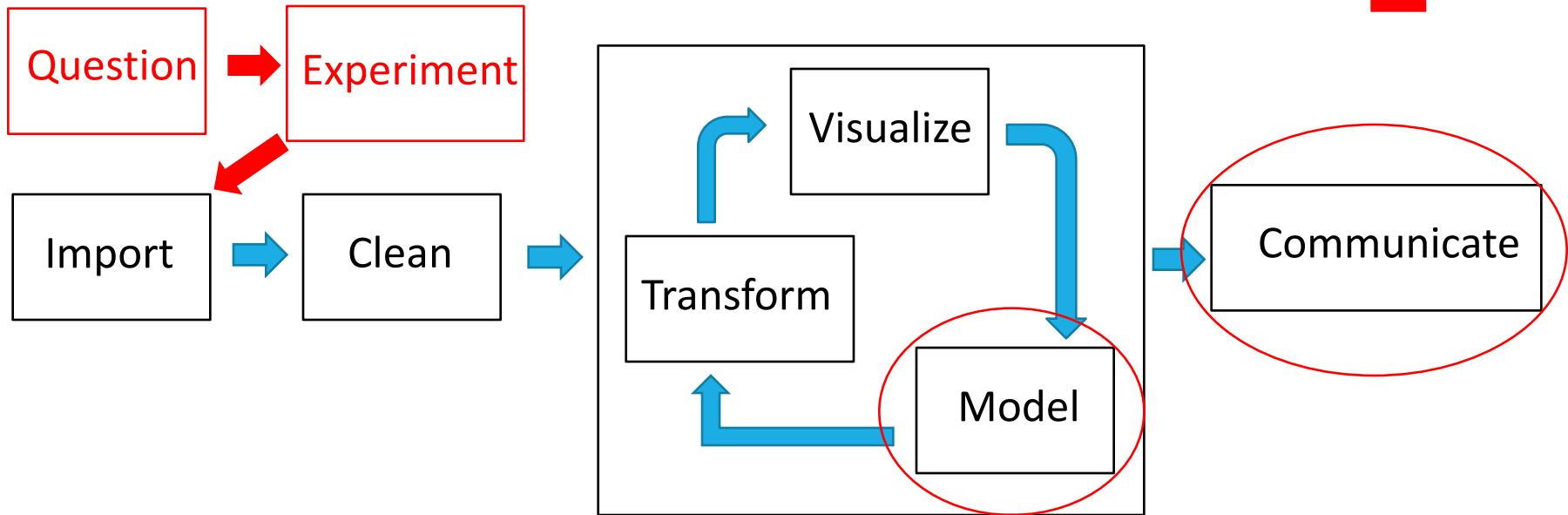
SAS

MEASURES OF CENTRALITY

CONFIDENCE INTERVALS

HYPOTHESIS TESTING

Statistical flow chart



- Most statistics classes emphasize “Model” and “Communicate”
- It is important to keep in mind the entire process
- Actually, the flow chart is really a loop

SAS

Getting Started with SAS

SAS is a powerful tool for analyzing data

To get the most out of this class, you should become conversant in SAS

You can access SAS in at least three ways:

- A local copy installed on your computer
- A remote copy hosted by the “Citrix” SMU server
- A remote copy hosted by SAS (SAS onDemand)

Remote Copy

<https://odamid.oda.sas.com/SASODARegistration/>

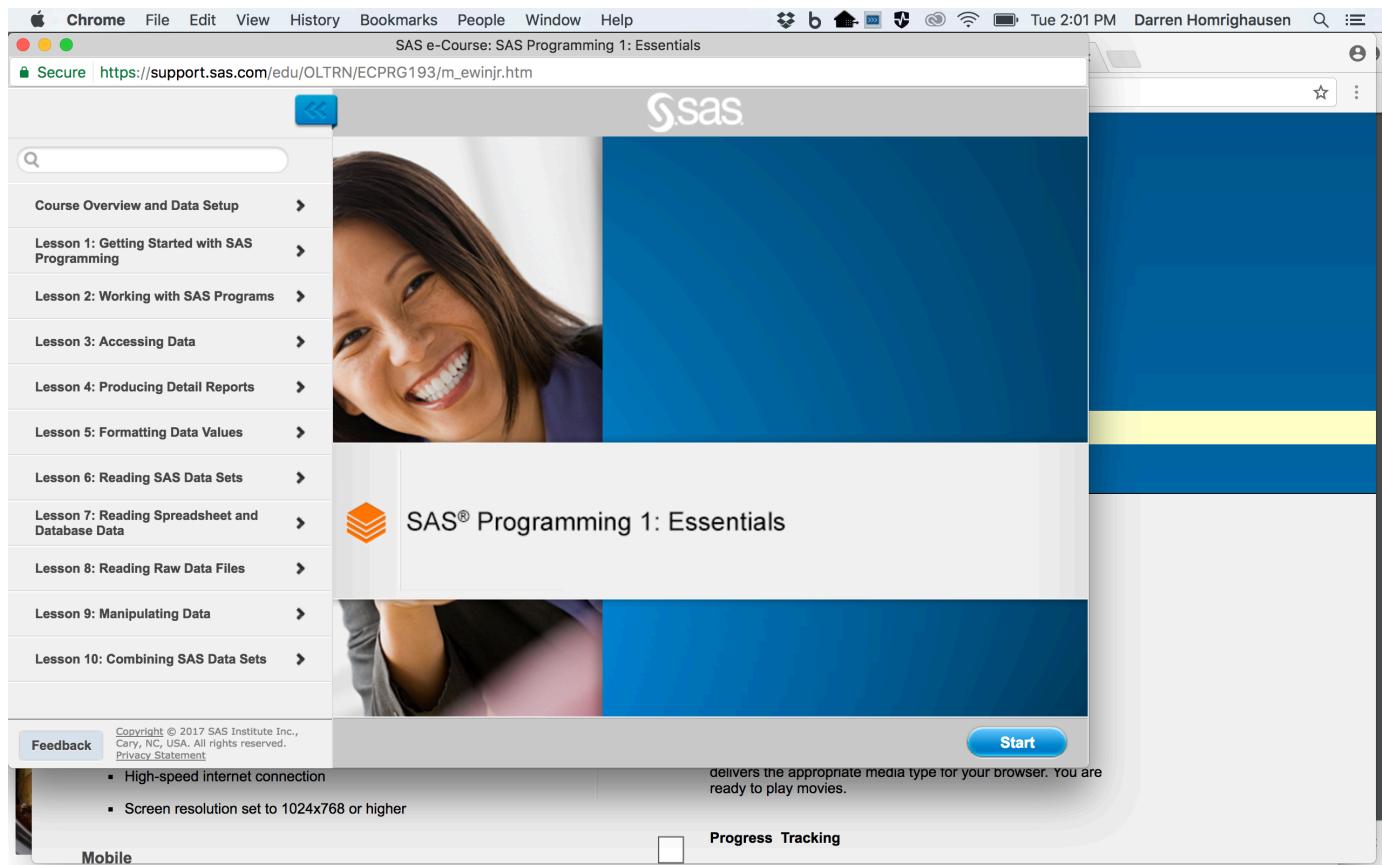
A confirmation email will be sent to the email address

After registering, you can join this class:

<https://odamid.oda.sas.com/SASODAControlCenter/enroll.html?enroll=79763c48-3d66-4faa-b802-14db272bc718>

Learning SAS

After registering for SAS onDemand, you get access to some SAS tutorials:



Learning SAS

Also, here are some links to brief introductory videos:

<http://video.sas.com/detail/videos/programming/video/4573016765001/writing-a-basic-sas-program?autoStart=true>

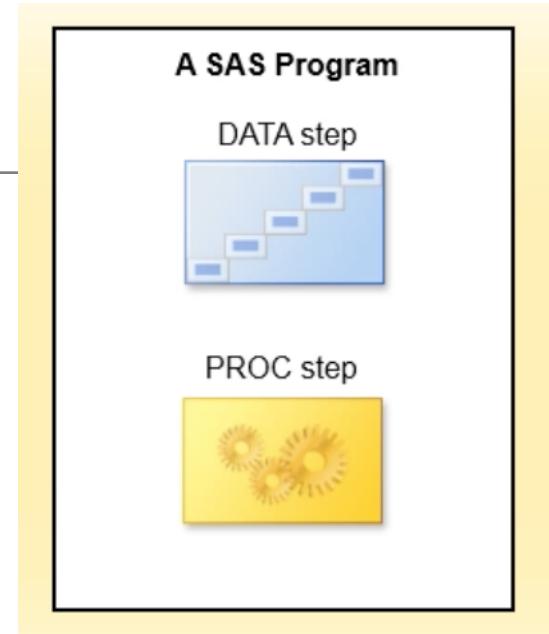
<http://video.sas.com/detail/videos/managing-data/video/4573016761001/filtering-a-sas-table-in-a-data-step?autoStart=true>

<http://video.sas.com/detail/videos/managing-data/video/4573016759001/performing-conditional-logic-in-sas?autoStart=true>

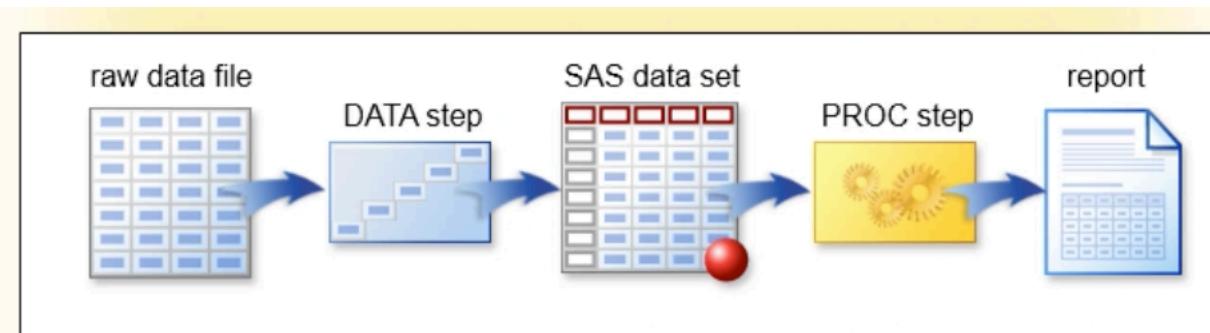
What does a typical SAS session look like?

SAS overview

Every SAS session is some combination of DATA and PROC steps

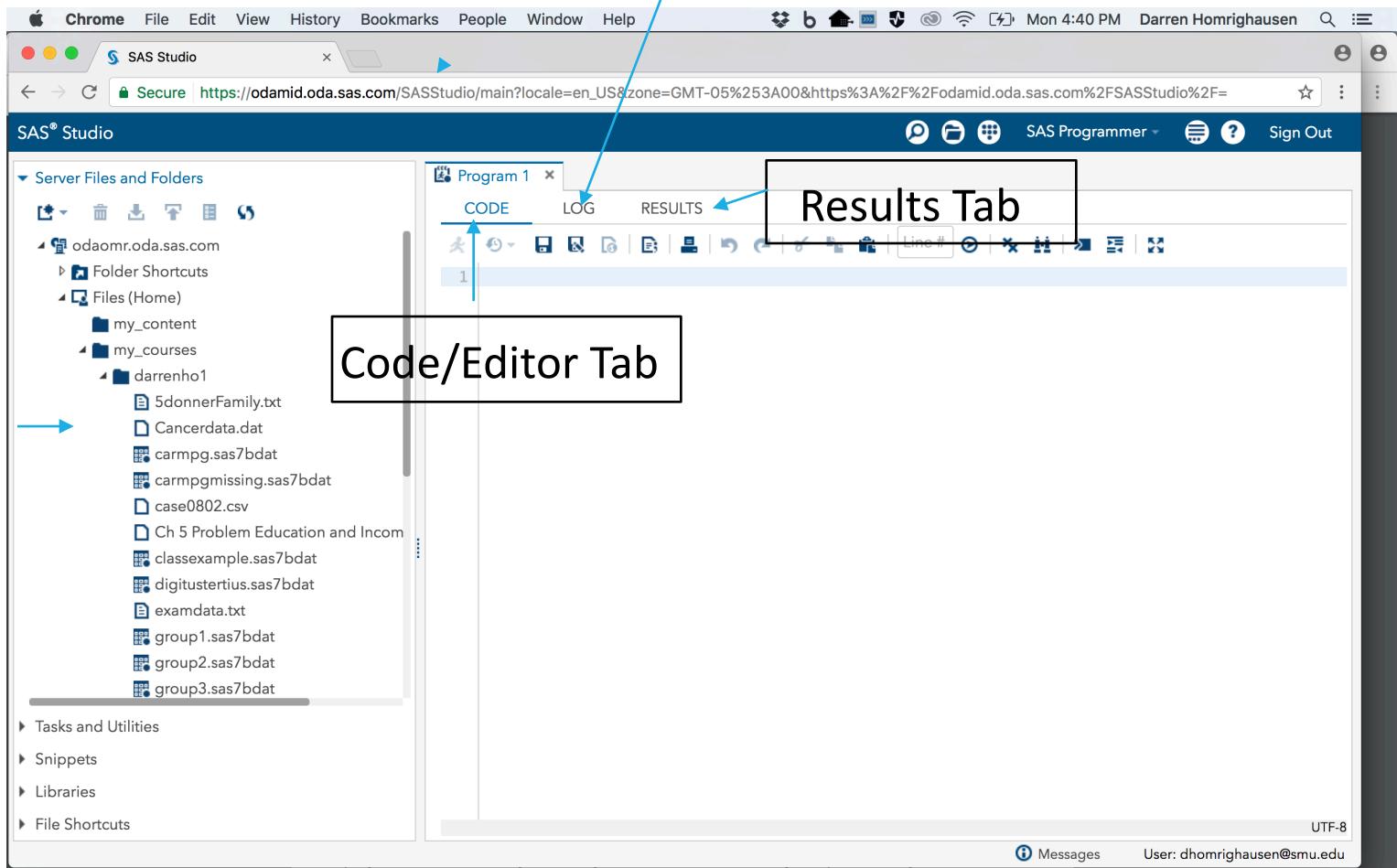


An example SAS session might look like:



SAS session

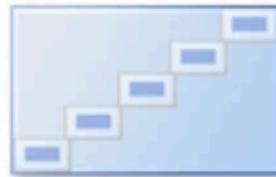
Log Tab



SAS session: Editor

A SAS Program

DATA step



PROC step

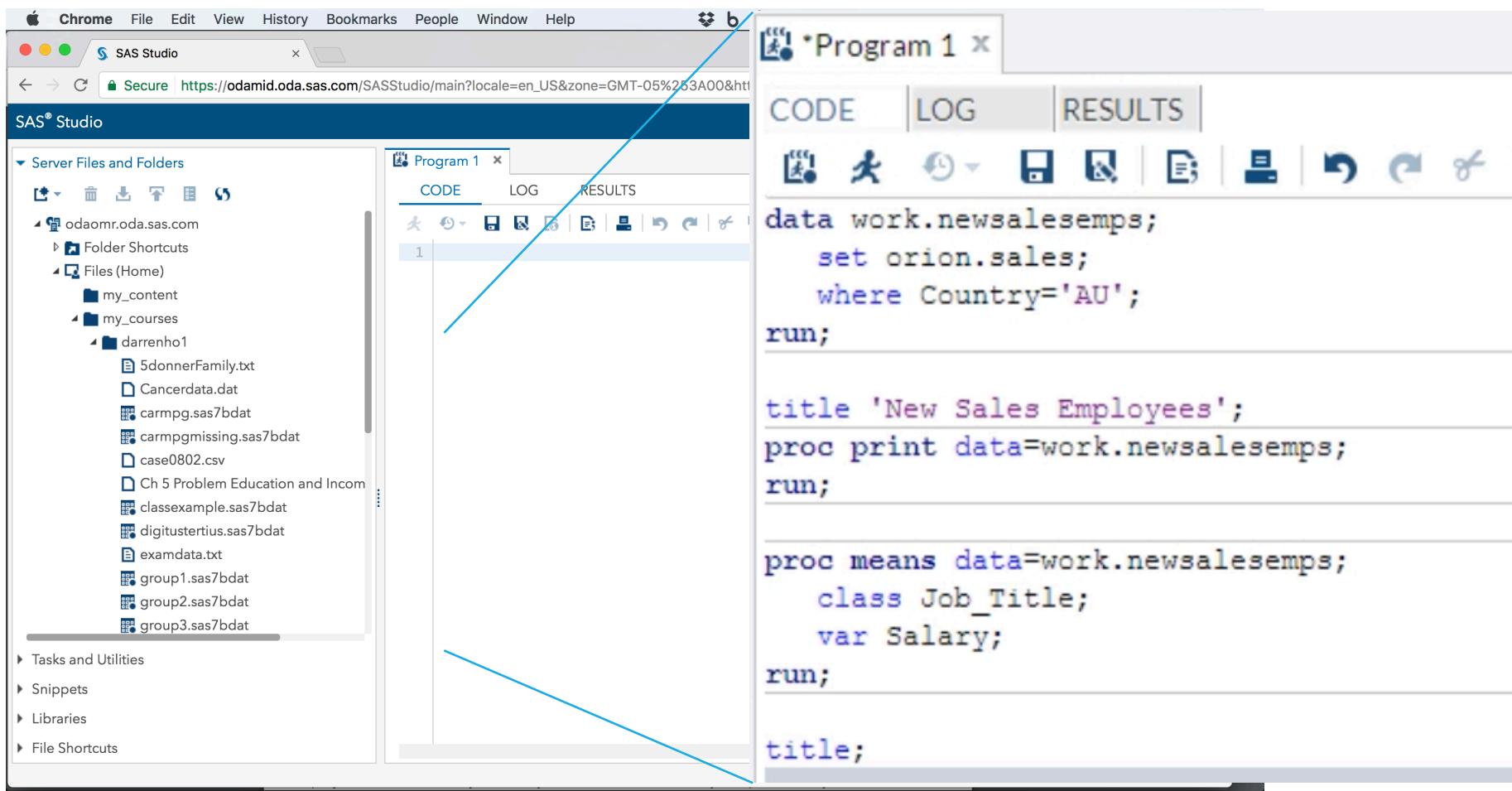


```
data work.newsalesemps;
  set orion.sales;
  where Country='AU';
run;

proc print data=work.newsalesemps;
run;

proc means data=work.newsalesemps;
  class Job_Title;
  var Salary;
run;
```

SAS session: Editor/Code



The screenshot shows a SAS Studio interface running in a Chrome browser. The left sidebar contains navigation links for 'Server Files and Folders', 'Tasks and Utilities', 'Snippets', 'Libraries', and 'File Shortcuts'. The main area features a 'Program 1' window with three tabs: 'CODE', 'LOG', and 'RESULTS'. The 'CODE' tab is active, displaying the following SAS code:

```
data work.newsalesemps;
  set orion.sales;
  where Country='AU';
run;

title 'New Sales Employees';
proc print data=work.newsalesemps;
run;

proc means data=work.newsalesemps;
  class Job_Title;
  var Salary;
run;

title;
```

SAS session: Log

The screenshot shows the SAS Studio interface and a Chrome browser window. The SAS Studio window has a sidebar with 'Server Files and Folders' containing 'odaomr.oda.sas.com' and a 'Program 1' panel with tabs for 'CODE', 'LOG' (selected), and 'RESULTS'. The 'LOG' tab displays a log of SAS code execution. The browser window shows a secure connection to https://odamid.oda.sas.com/SASStudio/main?locale=en_US&zone=GMT-05%253A00&https%3A%2F%2F.

```
1      OPTIONS NONOTES NOSTIMER NOSOURCE NOS
38     ;
39     data work.newsalesemps;
40       set orion.sales;
41       where Country='AU';
42     run;

NOTE: The data set WORK.NEWSALESEMP has 63 obse
NOTE: DATA statement used (Total process time):
      real time          0.03 seconds
      cpu time           0.00 seconds

43
44       title 'New Sales Employees';
45       proc print data=work.newsalesemps;
```

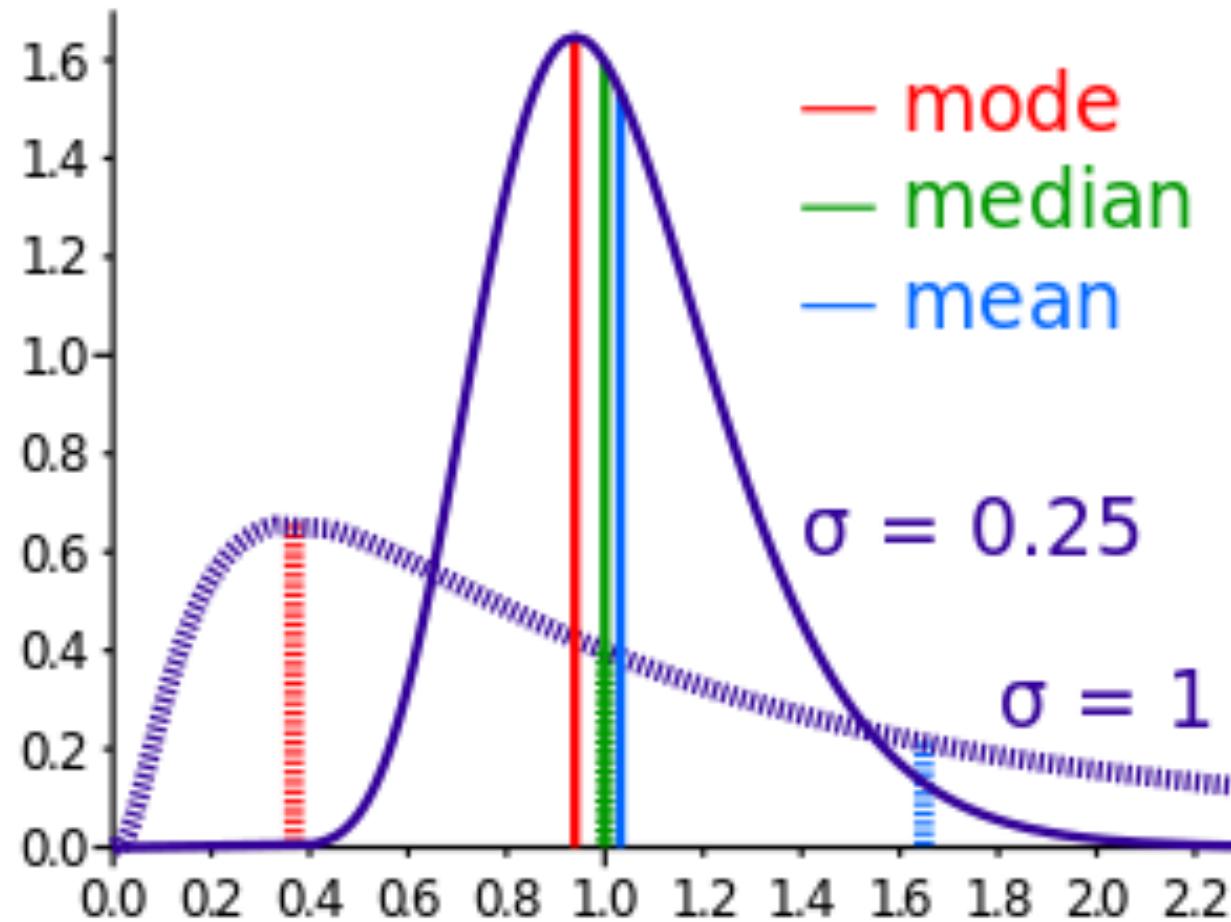
SAS session: Results

The screenshot shows a SAS Studio interface running in a Chrome browser window. The left sidebar displays 'Server Files and Folders' containing various SAS files and datasets. A central panel shows 'Program 1' with tabs for CODE, LOG, and RESULTS. The RESULTS tab is active, displaying a table titled 'New Sales Employee'. The table lists 16 rows of data with columns: OBS, Employee_ID, First_Name, Last_Name, Gender, and Salary.

OBS	Employee_ID	First_Name	Last_Name	Gender	Salary
1	120102	Tom	Zhou	M	108255
2	120103	Wilson	Dawes	M	87975
3	120121	Irenie	Elvish	F	26600
4	120122	Christina	Ngan	F	27475
5	120123	Kimiko	Hotstone	F	26190
6	120124	Lucian	Daymond	M	26480
7	120125	Fong	Hofmeister	M	32040
8	120126	Satyakam	Denny	M	26780
9	120127	Sharryn	Clarkson	F	28100
10	120128	Monica	Kletschkus	F	30890
11	120129	Alvin	Roebuck	M	30070
12	120130	Kevin	Lyon	M	26955
13	120131	Marinus	Surawski	M	26910
14	120132	Fancine	Kaiser	F	28525
15	120133	Petrea	Soltau	F	27440
16	120134	Sian	Shannan	M	28015

Measures of centrality

Mean, Median, and Mode



Estimating the population mean

Assumptions

We will return to these later...

The sample mean

If Y_1, Y_2, \dots, Y_n are the possible values and y_1, y_2, \dots, y_n are observed, then

$$\bar{Y} = \frac{(Y_1 + Y_2 + \dots + Y_n)}{n} \quad \text{and} \quad \bar{y} = \frac{(y_1 + y_2 + \dots + y_n)}{n}$$

The idea: \bar{Y} “approx. equals” μ

What is the definition of “approx. equals”?

- The mean of \bar{Y} equals μ
- $\text{Variance}(\bar{Y}) = \frac{\sigma^2}{n} \rightarrow \text{SD}(\bar{Y}) = \frac{\sigma}{\sqrt{n}}$ (*note: this goes to 0 as n increases and is the law of large numbers*)

So, $\bar{Y} = \mu + \text{remainder}$,

where $|\text{remainder}|$ is “unlikely” to exceed (multiple) $\cdot \frac{\sigma}{\sqrt{n}}$

How do we quantify “unlikely”? Via the central limit theorem (CLT)!

The central limit theorem CLT)

The CLT is a refinement of the law of large numbers

It states that not only does \bar{Y} eventually converge to μ , but the distribution of \bar{Y} “eventually” looks like normal/Gaussian distribution
 (“eventually”: commonly taken to be $n>30$ and obs. are independent)

In notation, we write this as $\bar{Y} \sim N(\mu, \frac{\sigma^2}{n})$

So, $\bar{Y} = \mu + \text{remainder} \leftrightarrow \bar{Y} - \mu = \text{remainder}$

Hence, remainder $\sim N(0, \frac{\sigma^2}{n})$

Confidence interval for the population mean

Suppose we want to estimate μ

A reasonable estimator is $\hat{\mu} = \bar{Y}$

What if we additionally want to find an interval which we believe contains μ ?

If that interval is symmetric around \bar{Y} , it would look like $[\bar{Y} - E, \bar{Y} + E]$

Our task: Quantify how “likely” it is that: $\bar{Y} - E \leq \mu \leq \bar{Y} + E$

For that, we use the result: remainder $\sim N(0, \frac{\sigma^2}{n})$

Hence, choose: $E = (\text{multiple}) \cdot \frac{\sigma}{\sqrt{n}}$

where **(multiple)** is chosen to quantify “likely”...

Confidence interval for the population mean

Reminder:

$$\bar{Y} \sim N\left(\mu, \frac{\sigma^2}{n}\right), \quad \bar{Y} = \mu + \text{remainder}, \quad \text{remainder} \sim N\left(0, \frac{\sigma^2}{n}\right), \quad E = (\text{multiple}) \cdot \frac{\sigma}{\sqrt{n}}$$

→ set **(multiple)** to be a quantile of the standard normal/Gaussian distribution
(we notate the normal distribution as $Z \sim N(0,1)$)

z_α is the α^{th} quantile of a normal provided: Probability($Z < z_\alpha$) = $1 - \alpha$

Claim: If we set $E = z_{\alpha/2} \cdot \frac{\sigma}{\sqrt{n}}$ then $[\bar{y} - E, \bar{y} + E]$ is a $(1 - \alpha)100\%$ confidence interval (CI)

(Proof sketch)

$$\begin{aligned} \text{Probability}(\bar{Y} - E \leq \mu \leq \bar{Y} + E) &= \text{Probability}(-E \leq \text{remainder} \leq E) \\ &= \text{Probability}(-z_{\alpha/2} \leq Z \leq z_{\alpha/2}) \end{aligned}$$

Example:

Claim: If we set $E = z_{\alpha/2} \cdot \frac{\sigma}{\sqrt{n}}$ then $[\bar{y} - E, \bar{y} + E]$ is a $(1 - \alpha)100\%$ confidence interval (CI) for μ

E is often referred to as the margin of error

Find a 99% confidence interval for μ if $\sigma^2 = 1.44$ and a sample size of 49

$$E = z_{\alpha/2} \cdot \frac{\sigma}{\sqrt{n}} = 2.575 \cdot \frac{1.2}{\sqrt{49}} = 0.4414 \quad (\text{note: this isn't a CI})$$



Hypothesis testing



- THE POLYCEPHALIC TWIN OF CONFIDENCE INTERVALS

Let's look at the mean

- The mean μ is unknown
- Suppose we presume $\mu = \mu_0$ (say, $\mu_0 = 10$)
- We calculate \bar{y} (say, $\bar{y} = 13$)
- Is this sufficient evidence that $\mu \neq \mu_0$?

How likely is this observed value? Can we say “beyond a reasonable doubt” that this observed value is too extreme?

How much variation should we expect?

The mean of \bar{Y} is μ

$$\text{The Variance}(\bar{Y}) = \frac{\sigma^2}{n} \rightarrow \text{SD}(\bar{Y}) = \frac{\sigma}{\sqrt{n}}$$

Therefore, if we center and rescale:

$\rightarrow Z = \frac{\bar{Y} - \mu}{(\sigma/\sqrt{n})}$ has mean 0 and variance 1, by CLT: $Z \sim N(0,1)$
(assuming n is large enough and obs. are independent)

Large values of $|Z|$ are “unlikely”: $\text{Probability}(|Z| > z_{\alpha/2}) = \alpha$

Of course, computing this Z requires knowledge of μ **which is unknown**

What we **do** know is the hypothesized value μ_0

How much variation should we expect?

Let's look at $Z_0 = \frac{\bar{Y} - \mu_0}{(\sigma/\sqrt{n})}$ and compute $z_0 = \frac{\bar{y} - \mu_0}{(\sigma/\sqrt{n})}$ instead

As before, computing $\text{Probability}(|Z_0| > z_0)$ requires knowledge of μ
(after all, we showed that the distribution of \bar{Y} depends on μ)

Instead, we compute a related quantity: $\text{Probability}(|Z_0| > z_0 | \mu = \mu_0)$
(this is a conditional probability under the hypothesis that $\mu = \mu_0$)

Important: If $\mu = \mu_0$, then we of course know μ as we know μ_0

Hence, we can compute $\text{Probability}(|Z_0| > z_0 | \mu = \mu_0)$ using the CLT:

→ $Z_0 | \mu = \mu_0$ is approximately $N(0,1)$

Hypothesis Testing: Definitions

By convention, we define the following key ideas

- $H_0: \mu = \mu_0$ a “null hypothesis”
- $H_a: \mu \neq \mu_0$ an “alternative hypothesis”
- Probability($|Z_0| > z_0 | \mu = \mu_0$) as the “p-value”
(sometimes alternatively written: Probability($|Z_0| > z_0 | H_0$))

Massively important: the p-value is not Probability($\mu = \mu_0 | |Z_0| > z_0$)

Implication: the p-value...

...is not the conditional probability of the null hypothesis given $|Z_0| > z_0$

...is the conditional probability of getting a larger observation given H_0

If this probability is “small enough”, we have evidence to **reject H_0**

Hypothesis testing: An analogy

Suppose we are on a jury at a trial

Our verdict will either be guilty or innocent after seeing evidence

The doctrine of “innocent until proven guilty” demands evidence in excess of “reasonable doubt”

Statistics allows us to quantify both “innocent until proven guilty” and “beyond reasonable doubt” for a wide array of scientific problems

- “innocent until proven guilty” \leftrightarrow hypothesis testing
- “innocence” \leftrightarrow null hypothesis
- “guilty” \leftrightarrow alternative hypothesis
- “beyond reasonable doubt” \leftrightarrow small p-value

Confidence intervals and hypothesis tests



Suppose for a given problem we compute $z_0 = 1.96$

How likely is it that we would draw a $|Z_0| > z_0$, given $\mu = \mu_0$?

→ Probability($(|Z_0| > z_0) | \mu = \mu_0$) = 0.05



Let's compare this to a $(1 - \alpha)100\% = 95\%$ confidence interval

Reminder:

z_α is the α^{th} quantile of a normal provided: Probability($Z < z_\alpha$) = $1 - \alpha$

→ 0.025^{th} quantile is the value $z_{0.025}$ such that

Probability($|Z| > z_{0.025}$) = 0.05

When writing $E = z_{\alpha/2} \cdot \frac{\sigma}{\sqrt{n}}$ the multiplier $z_{0.025} = 1.96$, a main ingredient for a 95% confidence interval

The connection



100(1 – α)% confidence intervals will include μ_0
if and only if

$$\text{Probability}(|Z_0| > z_0 | \mu = \mu_0) \geq \alpha$$

Stated another way:

μ_0 not being in a 100(1 – α)% confidence interval is equivalent to
rejecting $H_0: \mu = \mu_0$ with a p-value less than α

(There are some caveats here. If Z takes on too few values or if we are doing “one sided hypothesis tests”, some adjustments are necessary)