

Reproducible Clinical Data Analysis with R and RStudio

Session 4
Reproducible Reports
April 5, 2018

Introduction	Early: 3/27, 4 PM Late: 3/27, 5:30 PM
Getting Data	Early: 3/28, 4 PM Late: 3/28, 5:30 PM
Exploring Data	Early: 4/3, 4 PM Late: 4/2, 5:30 PM
Reproducible Reports	Early: 4/5, 4 PM Late: 4/5, 5:30 PM
Course Project Presentations	TBA

R Markdown and R Notebooks



R Markdown

Reference Guide
Learn more about R Markdown at rmarkdown.rstudio.com

Learn more about Interactive Docs at shiny.rstudio.com/articles

Contents:

1. **Markdown Syntax**
2. Knitr chunk options
3. Pandoc options

Syntax

Plain text

End a line with two spaces
to start a new paragraph.

italics and _italics_

bold and __bold__

superscript^{^2}

~~strikethrough~~

[link](www.rstudio.com)

Header 1

Header 2

Header 3

Header 4

Becomes

Plain text

End a line with two spaces to start a new paragraph.

italics and *italics*

bold and **bold**

superscript²

strikethrough

link

Header 1

Header 2

Header 3

~/code/rcda_mgh_2018/mar_2018 - RStudio

esr-reference-ranges.Rmd

Preview Insert Run

```
1: ---  
2: title: "ESR Reference Range Study"  
3: output: html_notebook  
4: ---  
5:  
6: ```{r setup, include=FALSE}  
7: # Load required packages  
8: library(tidyverse)  
9: l  
10: l ESR Reference Range Study  
11: l Chunk 1: setup  
12: l Background and Objectives  
13: # MGH  
14: ARUP Laboratories  
15: Quest Diagnostics  
16: LabCorp  
17: Bakeman's ABCs of Interpretive Laboratory Data  
18: Data Acquisition  
19: T Chunk 2: connect-datamart  
20: p username()  
21: s password()  
1:1 # ESR Reference Range Study
```

Environment History Connections Git

Import Dataset List Global Environment

Environment is empty

Files Plots Packages Help Viewer

R Markdown

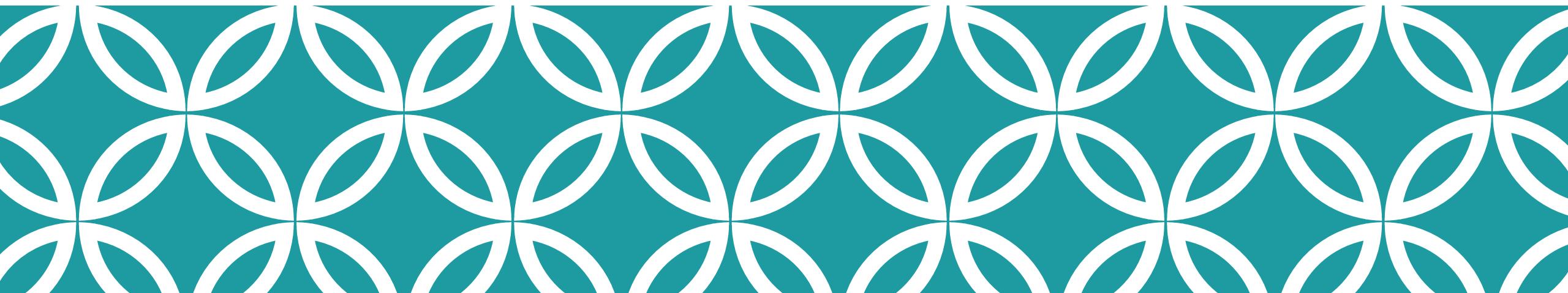
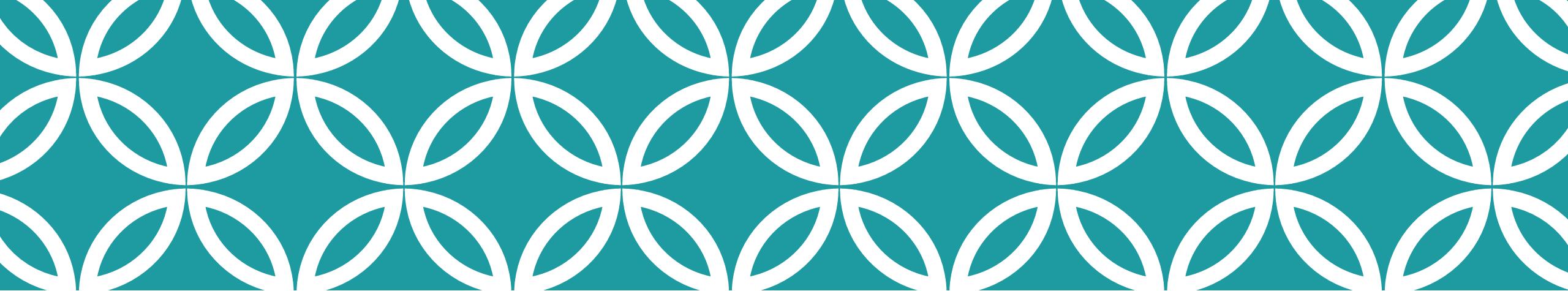
Console

Your Turn #1

Create a new R Notebook. Save it on your Desktop as Practice.Rmd and preview it.

Now delete all the contents below the header and try to write the R Markdown that recapitulates the handout titled “R Notebook Practice”.





Anatomy of a Reproducible Report

Header

starts and ends
with 3 dashes

```
1 ---  
2 title: "ESR Reference Range Study"  
3 output: html_notebook  
4 ---
```

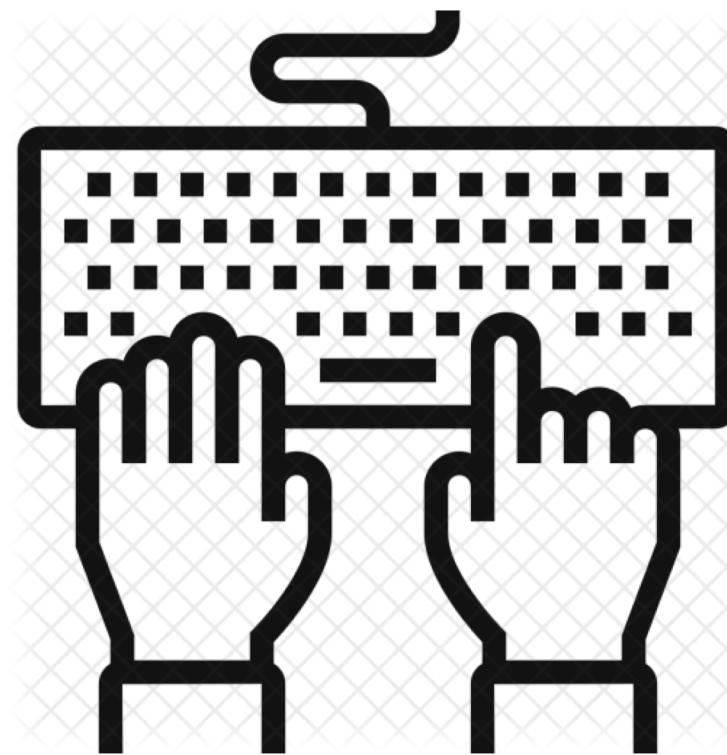
title:
Document title

output:
must say html_notebook

“Copy/Paste” Code



“Write Yourself” Code



“Setup” Chunk

chunk name
(optional)

chunk option:
don't show code in rendered document

comments start
with #

```
6 - ````{r setup, include=FALSE}  
7 # Load required packages  
8 library(tidyverse)  
9 library(odbc)  
10 library(lubridate)  
11 ...
```



for database
access

for dealing
with dates

Background and Objectives



Data Acquisition



Data Exploration and Clean-Up



Visualization and Modeling



Summary

Background and Objectives

- What is the question to be addressed by the project?
- How could addressing this question change patient management or further generalizable knowledge?
- What data resources could be used?
- What are the *first 1-2 key graphs or statistics* that should be obtained?

Background and Objectives



Data Acquisition



Data Exploration and Clean-Up



Visualization and Modeling

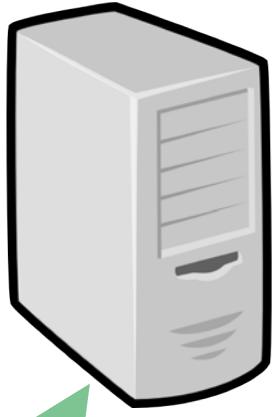


Summary

Data Acquisition

- **Narrative**
 - Database sources
 - Database server address and database table name
 - Data file sources
 - Origin, with references
 - File name and format
 - Both
 - General properties of the data set
 - Code book for variables
- **Code**
 - Database sources
 - Connection and query
 - Don't hard-code credentials!
 - Data dump
 - Disconnect
 - Data file sources
 - Code to read data file

Database Management System



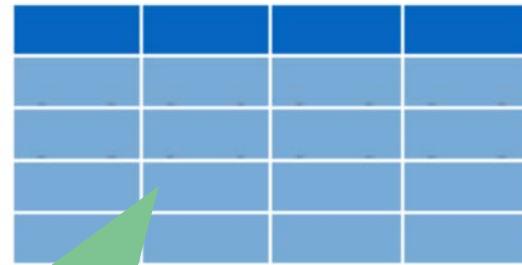
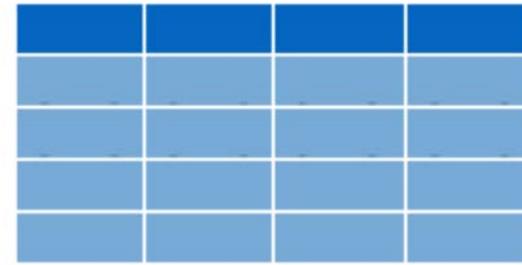
PHSSQL2057
("Datamart")

Database



MGHLABUTIL

Database Tables



MGHLABUTIL_LabResults

Connect to Datamart

functions to ask login credentials

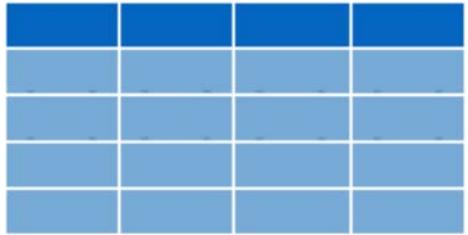
connection handle

```
62 - ````{r connect-datamart}
63 - username <- function() {
64   .rs.api.showPrompt(title = "Username",
65                       message = "Please enter your Partners User Name")
66 }
67
68 - password <- function() {
69   .rs.api.askForPassword(prompt = "Please enter your Partners Password")
70 }
71
72 con <- dbConnect(odbc(),
73                   driver = "{FreeTDS}",
74                   dsn = "PHSSQL2057",
75                   uid = str_c("PARTNERS\\\", username()),
76                   pwd = password())
77
78 labresults_tbl <- tbl(con, "MGHLABUTIL_LabResults")
79 ````
```



pointer to database table

Data frame



filter rows by a condition

Select one or
more columns

Database Table

Your Turn #2

Connect to the Datamart by running the code chunk labeled connect-datamart.

Find the MGHLABUTIL database and within it, find the MGHLABUTIL_LabResults table.

Preview it by clicking the symbol that looks like a small table.

There are 71 columns in this database table. Pair up and discuss which specific columns you might want to pull to examine ESR reference ranges.



Code Book

Retrieve all ESR values from 2017, along with the following columns:

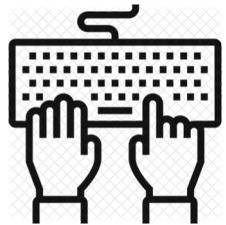
1. **CollAge.** Patient age at time of collection.
2. **CollectDateTime.** Date and time of collection.
3. **PtName.** Patient name.
4. **PtNumber.** Patient MRN (MGH).
5. **PtSex.** Patient sex.
6. **Result.** ESR result value.
7. **TstOrderName.** Test order name.
8. **UserQAFlags.** This column marks “high” and “low” results.

Database Query

You can treat a database table pointer like a data frame!

collect():
pull all results
(By default, max. 1000)

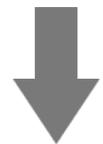
```
96  ```{r query}
97  esr <- labresults_tbl %>%
98    select(CollAge,
99           CollectDateTime,
100          PtName,
101          PtNumber,
102          PtSex,
103          Result,
104          TstOrderName,
105          UserQAFlags) %>%
106    filter(TstOrderName == "ESR",
107            year(CollectDateTime) == "2017") %>%
108    arrange(CollectDateTime) %>%
109    collect()
110  ```
```



Save Data Dump

str_c(...):
join multiple strings
into one string

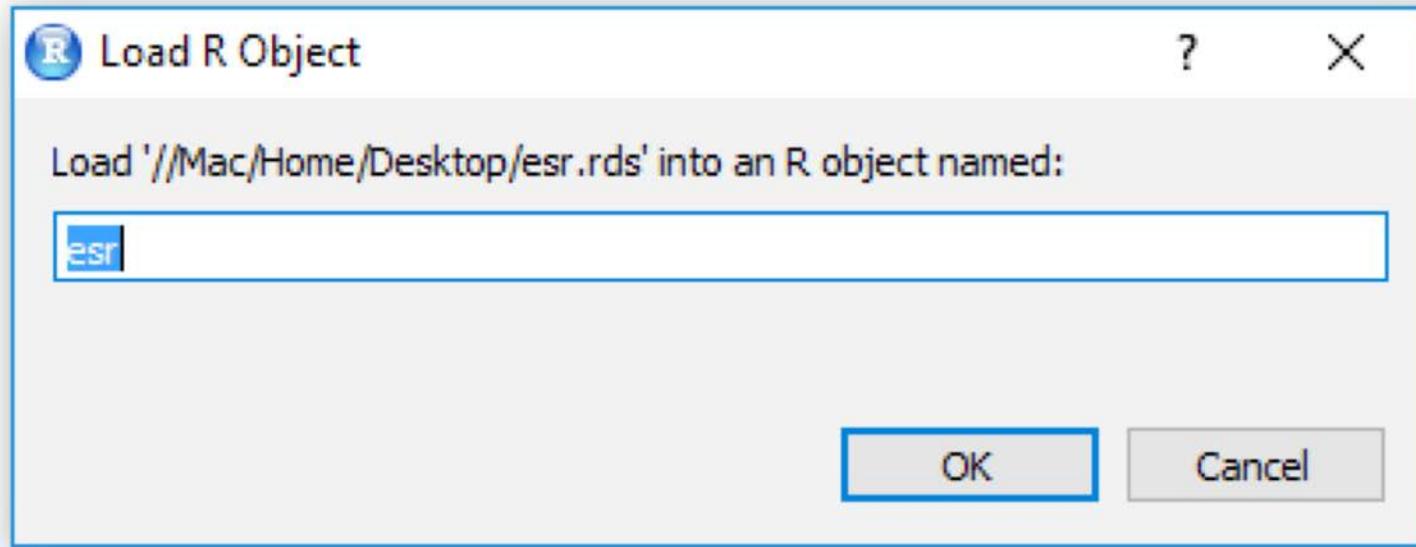
```
110 + ```{r save-rds}  
111   esr %>%  
112     write_rds(str_c("esr_data_", today(), ".rds"))  
113   ````
```



esr_data_2018-04-05.rds



Load Data Dump (if necessary)



Disconnect from Database

```
117 ~ ````{r}  
118   dbDisconnect(con)  
119   ````
```



Background and Objectives



Data Acquisition



Data Exploration and Clean-Up



Visualization and Modeling



Summary

Data Exploration and Clean-Up

- Convert messy data to tidy
- Check for missing values
- Verify (and if necessary, convert) data types
 - Numeric variables should be `integer` or `numeric`
 - Date-time variables should be `Date` or `POSIXct`
 - Categorical variables should be `factor`
 - All others should be `character`
- Visualize distributions
 - Numeric variables: `histogram`
 - Categorical variables: `numerical summary`
 - Character variables: look at a small sample
- Is the data as expected?
 - Range and shape of distribution
 - Number and breakdown of categories
 - Properties of character variables
- Explain and/or fix discrepancies

Convert Messy Data to Tidy

A data set is **tidy** if:

CollAge	PtNumber	PtSex	Result
7	5143567	M	22
42	3459254	F	5
19	2332467	F	5
80	3445732	M	89
41	7245673	F	12

1. Each **variable** is in its own **column**
2. Each **case** is in its own **row**
3. Each **value** is in its own **cell**

Missing values

Count the number of `NA`s (missing values) in each column of `esr`.

Hide

```
esr %>%
  map_df(function(x) sum(is.na(x))) %>%
  gather(column, NAs)
```



column	NAs
<chr>	<int>
CollAge	0
CollectDateTime	0
PtName	0
PtNumber	0
PtSex	0
Result	0
TstOrderName	0
UserQAFlags	0
8 rows	

None of the columns of `esr` contain any missing values.

CollAge

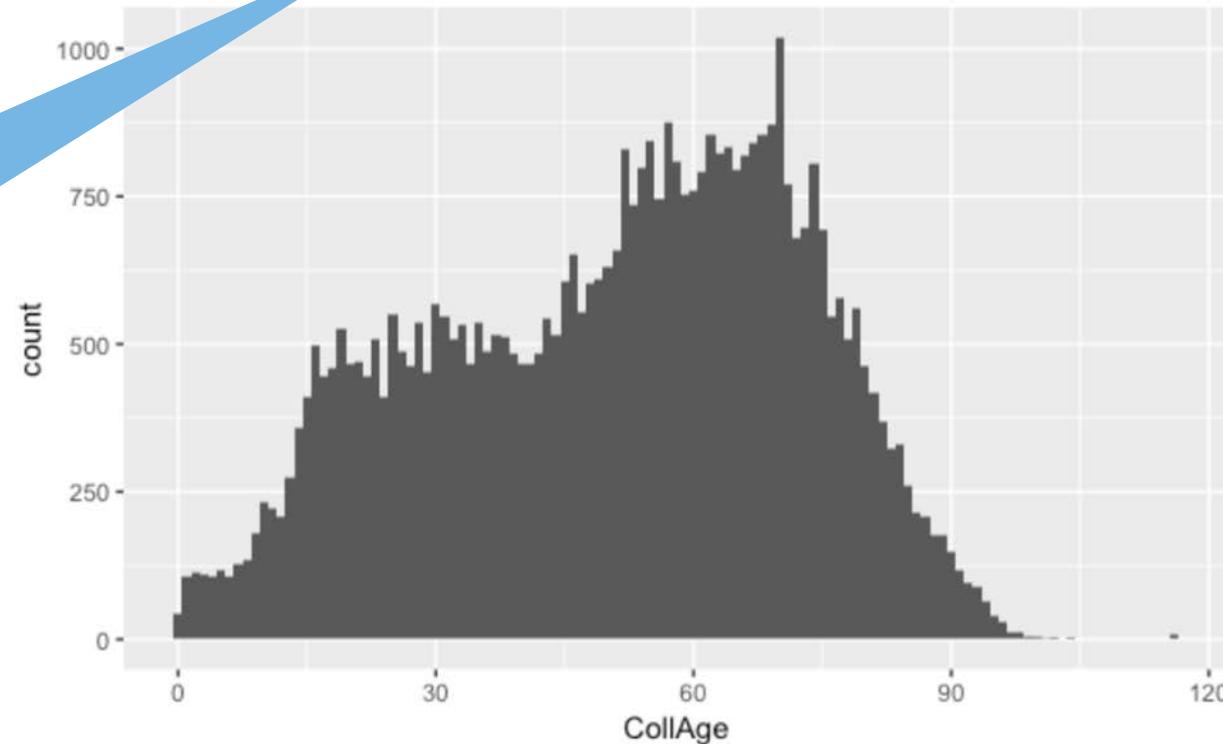
CollAge , the age at collection, is an **integer** column, which is appropriate.

```
ggplot(data = esr) +  
  geom_histogram(mapping = aes(x = CollAge),  
    binwidth = 1)
```

Hide



binwidth:
width of bins.
Experiment!



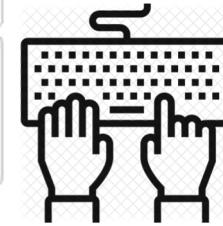
The distribution of CollAge is bimodal, with peaks at around 30 and at around 70 years.

CollectDateTime

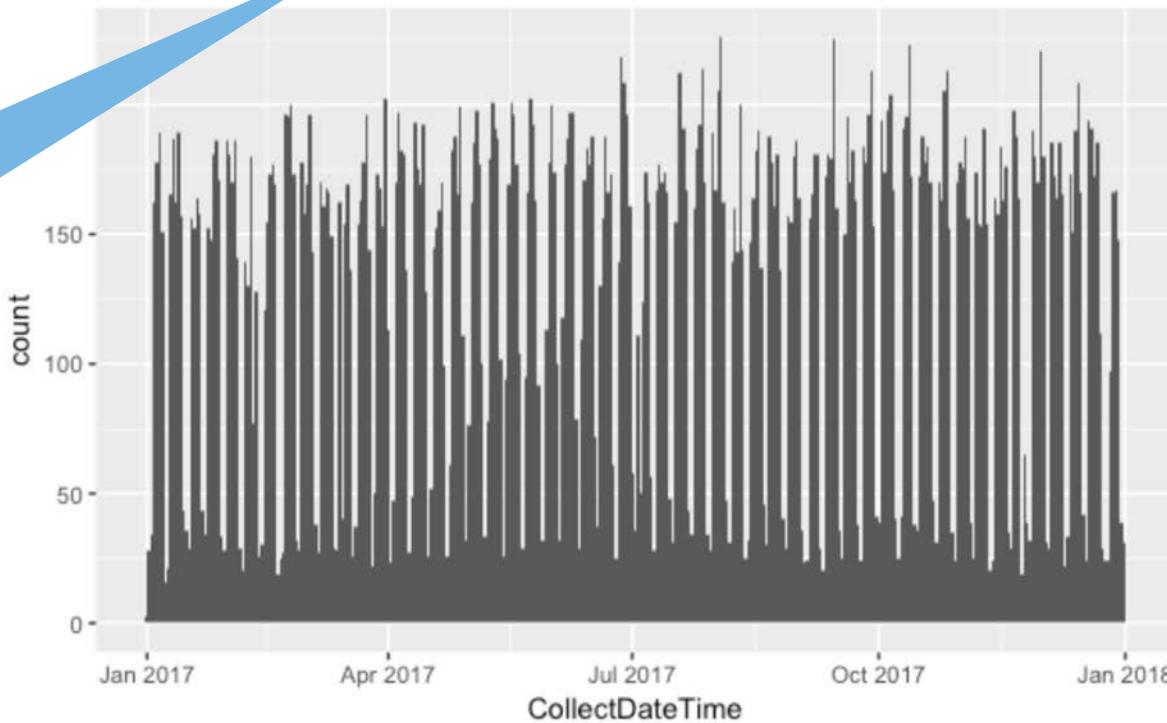
CollectDateTime is a **POSIXct** (datetime) column, which is appropriate.

Hide

```
ggplot(data = esr) +  
  geom_histogram(mapping = aes(x = CollectDateTime),  
    bins = 365)
```



bins:
total number of
bins.



Weekly cycling of ESR test volume is apparent. The weekly ESR test volume was approximately constant throughout 2017, without clear peaks or troughs.

Your Turn #3

Do you agree that weekly ESR testing volume was constant?

Can you think of an easy graph to verify this idea?

Discuss with your group and then try to create the graph in the R Markdown document.



sample_n():
sample rows
from a data
frame

pull():
extract a
single column
and display
it compactly

PtName

PtName is a **character** column, which is appropriate.

```
set.seed(1)
esr %>%
  sample_n(10) %>%
  pull(PtName)
```

```
[1] "R...,Z..."
[3] "M...,G..."
[5] "B...,J..."
[7] "C...,W...E..."
[9] "D...,D..."
```

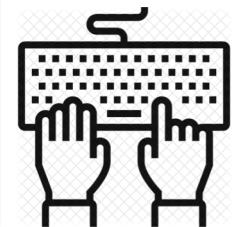
```
"C...,I..."
"...,S...,Z..."
"...,J..."
"C...,H..."
"D...,F..."
```

A random sample of 10 names is as expected.

set.seed():

set the random number generator
makes sure “random” results are reproducible

Hide



PtNumber

`PtNumber`, the MGH MRNs of the patients, is a **character** column, which is appropriate.

Note: if `PtNumber` were converted to integer, MRNs with leading zeroes would be altered.

[Hide](#)

```
set.seed(1)
esr %>%
  sample_n(10) %>%
  pull(PtNumber)
```



```
[1] " "
[8] " "
[9] " "
[10] " "
```

A random sample of 10 MRNs shows that all are seven digits long, as expected for MGH patients.

PtSex

PtSex is a **character** column, but since this is a categorial variable, we will convert it to a **factor**.

summary():
concise
statistical
summary

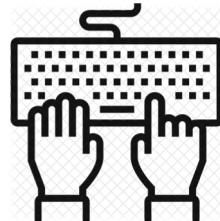
```
esr <- esr %>%
  mutate(PtSex = as_factor(PtSex))

esr %>%
  pull(PtSex) %>%
  summary()
```

F	M	U
26701	19857	17

as_factor():
convert to
factor

Overwrite existing
PtSex variable



The majority of ESR tests were performed on female (F) patients. Unexpectedly, in addition to M and F , there are some rows with PtSex denoted as U .

Overwrite
existing esr
data frame

```
esr %>%
  filter(PtSex == "U")
```

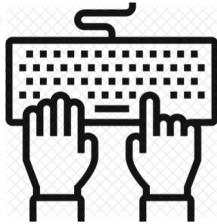
CollAge	CollectDateTime	PtName	PtNumber	PtSex
<int>	<S3: POSIXct>	<chr>	<chr>	<fctr>
35	2017-06-01 09:20:00	MGW CAP,ESR 01	QAPR-7475	U
39	2017-06-01 09:21:00	MGW CAP,ESR 02	QAPR-7476	U
33	2017-06-01 09:21:00	MGW CAP,ESR 03	QAPR-7477	U
35	2017-06-01 09:21:00	MGW CAP,ESR 04	QAPR-7478	U
35	2017-06-01 09:21:00	MGW CAP,ESR 05	QAPR-7479	U
35	2017-06-01 09:21:00	MGW CAP,ESR 06	QAPR-7480	U
35	2017-06-01 09:21:00	MGW CAP,ESR 07	QAPR-7481	U
35	2017-06-01 09:21:00	MGW CAP,ESR 08	QAPR-7482	U
35	2017-06-01 09:21:00	MGW CAP,ESR 09	QAPR-7483	U
35	2017-06-01 09:21:00	MGW CAP,ESR 10	QAPR-7484	U
35	2017-06-01 09:21:00	MGW CAP,ESR 11	QAPR-7485	U
35	2017-06-01 09:21:00	MGW CAP,ESR 12	QAPR-7486	U
35	2017-06-01 09:21:00	MGW CAP,ESR 13	QAPR-7487	U
35	2017-06-01 09:21:00	MGW CAP,ESR 14	QAPR-7488	U
35	2017-06-01 09:21:00	MGW CAP,ESR 15	QAPR-7489	U
35	2017-06-01 09:21:00	MGW CAP,ESR 16	QAPR-7490	U
35	2017-06-01 09:21:00	MGW CAP,ESR 17	QAPR-7491	U

1-10 of 17 rows | 1-5 of 8 columns

Previous 1 2 Next

It appears that the rows with PtSex equal to U belong to a small number of patients as well as CAP survey samples. Since these rows represent less than 0.1% of the whole data set, it will be safe to remove them.

```
esr <- esr %>%
  filter(PtSex != "U")
```



Conversion of Character to Numeric Types can create Missing Values!

`as.integer()`

`"123"`



`123`

`"abc"`



`NA`

get non-numeric values of Result, arranged by order of frequency

convert Result to integer and remove all rows where Result is now NA

Result

Result is a **character** column but should be an **integer**, so we will convert it.

Before doing so, examine any rows that do *not* contain an integer value:

```
esr %>%
  filter(is.na(as.integer(Result))) %>%
  group_by(Result) %>%
  summarize(n = n()) %>%
  arrange(desc(n))
```

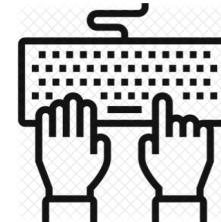
Result	n
<chr>	<int>
REFUS	1292
CREDIT	118
>140	14
QNS	12
CANCEL	6

1-10 of 19 rows Previous 1 2 Next

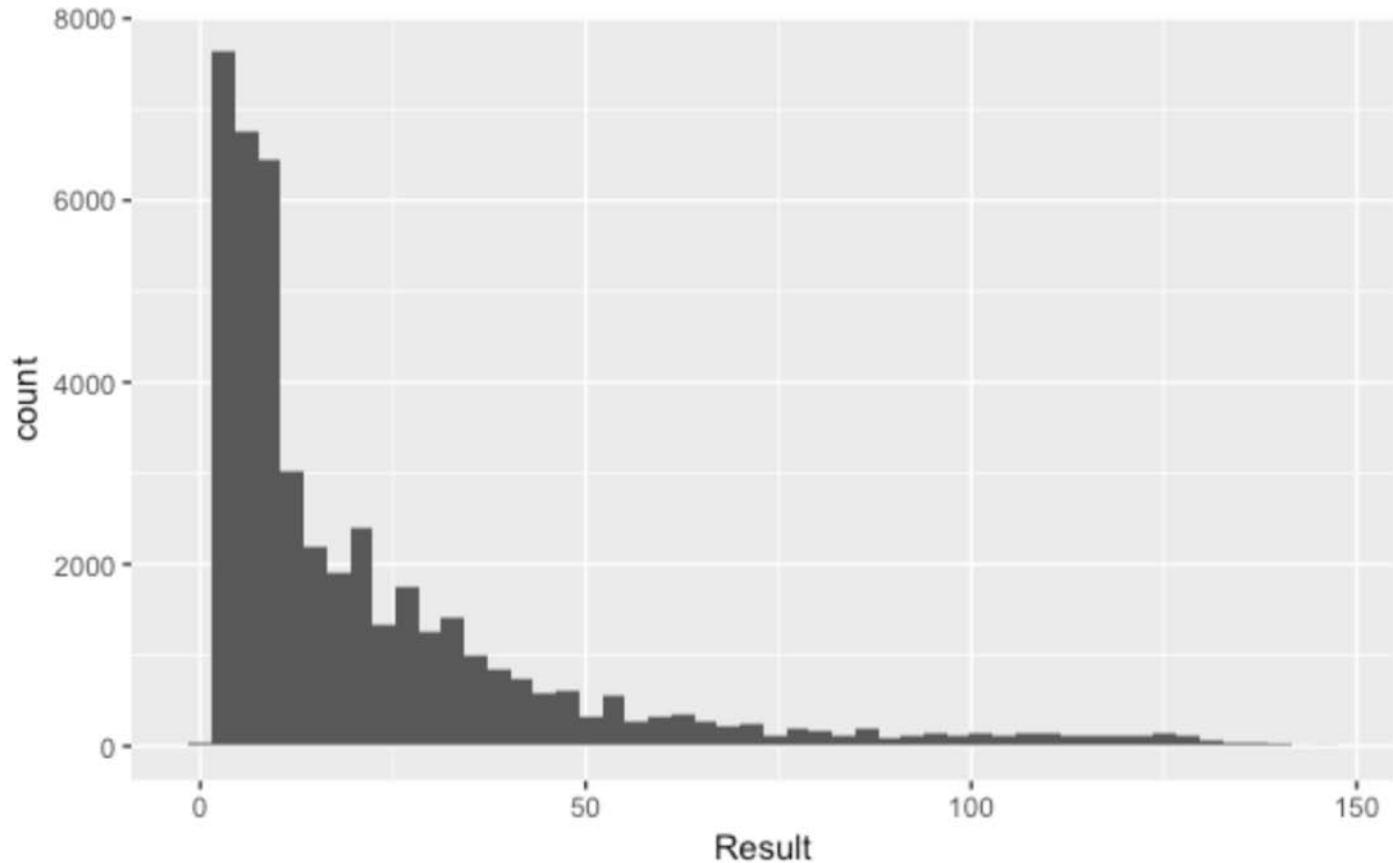
About 1400 rows (out of about 47,000) have a non-integer value in Result - about 3%. The large majority of them have a value of REFUS or CREDIT , indicating that the sample was not run or reported. These rows should be discarded.

Only a very small number of samples had excessive values (>140 , >144). It is acceptable to discard these rows as well.

```
esr <- esr %>%
  mutate(Result = as.integer(Result)) %>%
  drop_na()
```



```
ggplot(data = esr) +  
  geom_histogram(mapping = aes(Result),  
    bins = 50)
```



The distribution of ESR values appears to have a mode at around 10, with a long right tail, which represents the abnormal values.

Background and Objectives



Data Acquisition



Data Exploration and Clean-Up



Visualization and Modeling



Summary

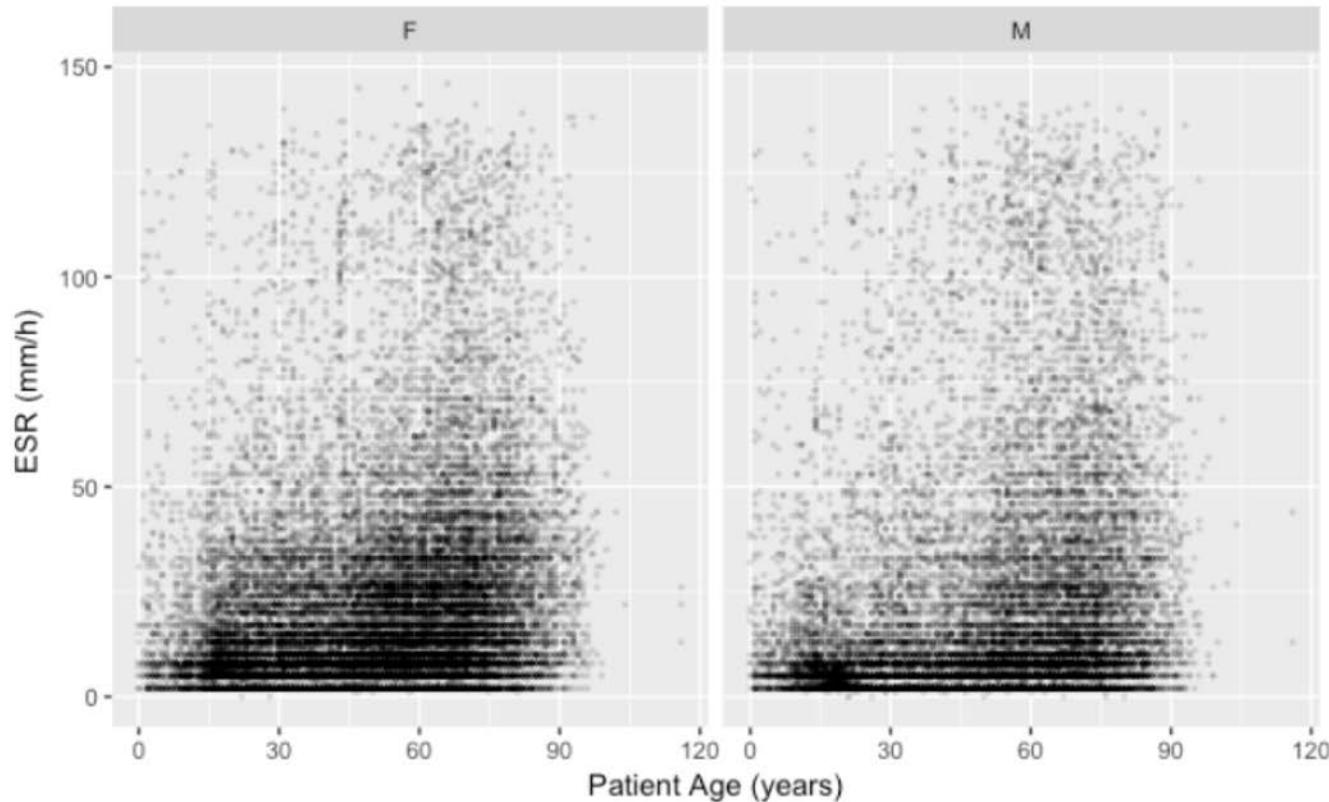
Visualization and Modeling

Iterative cycles of:

- Pose a question or hypothesis.
- Create a graph or statistical summary that addresses the question or tests the hypothesis.
- Comment on new insights gained and new problems raised.

The primary goal of this project is to test whether normal and abnormal ESR values can be separated into distinct populations, so that reference intervals can be set at the boundary. From the literature, I expect that ESR values increase with age and sex. To reduce these confounding effects, we will look at the relationship of ESR and age, broken down by sex.

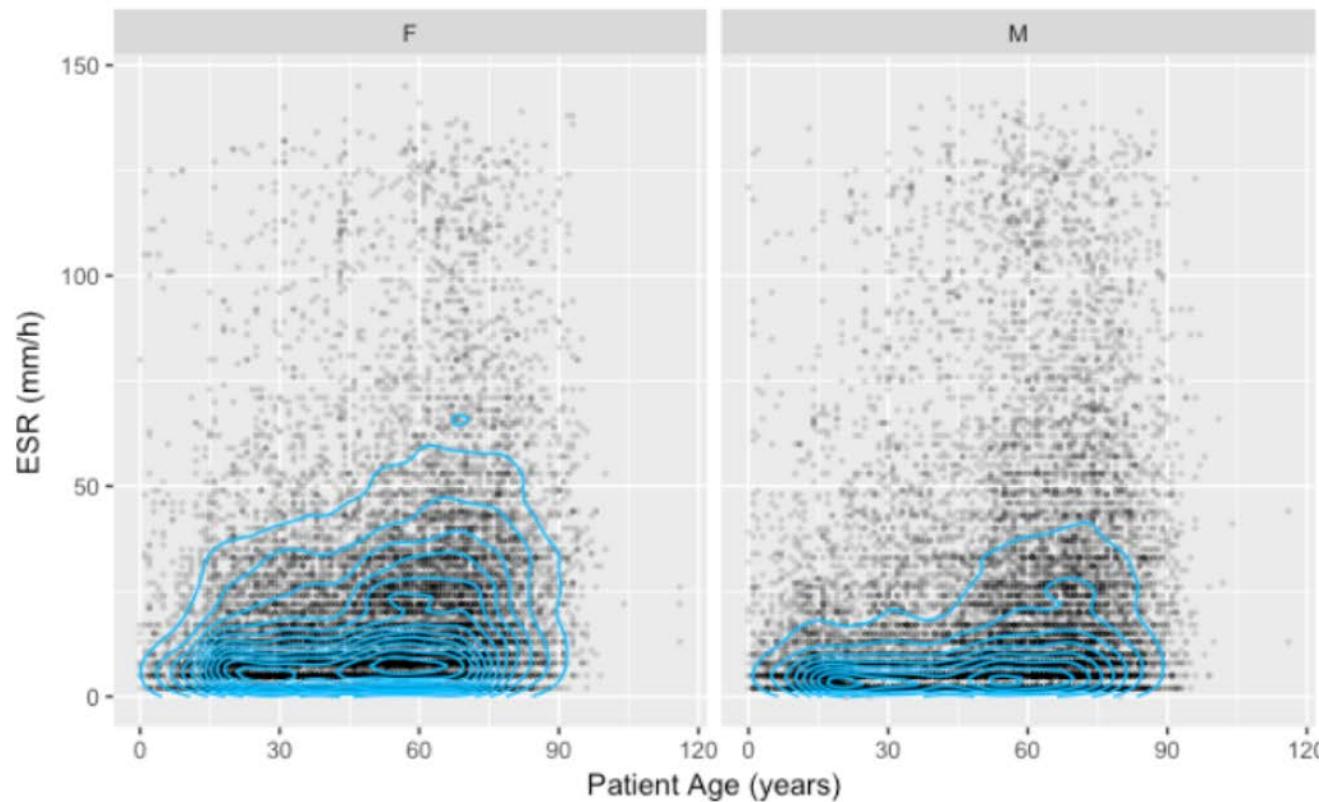
[Code](#)



There is no clear separation of normal and abnormal values. In addition, there are more data points on the female panel, which makes this graph difficult to interpret.

Will a separation become apparent when equal numbers of points are shown on each graph, and a density contour plot is overlaid?

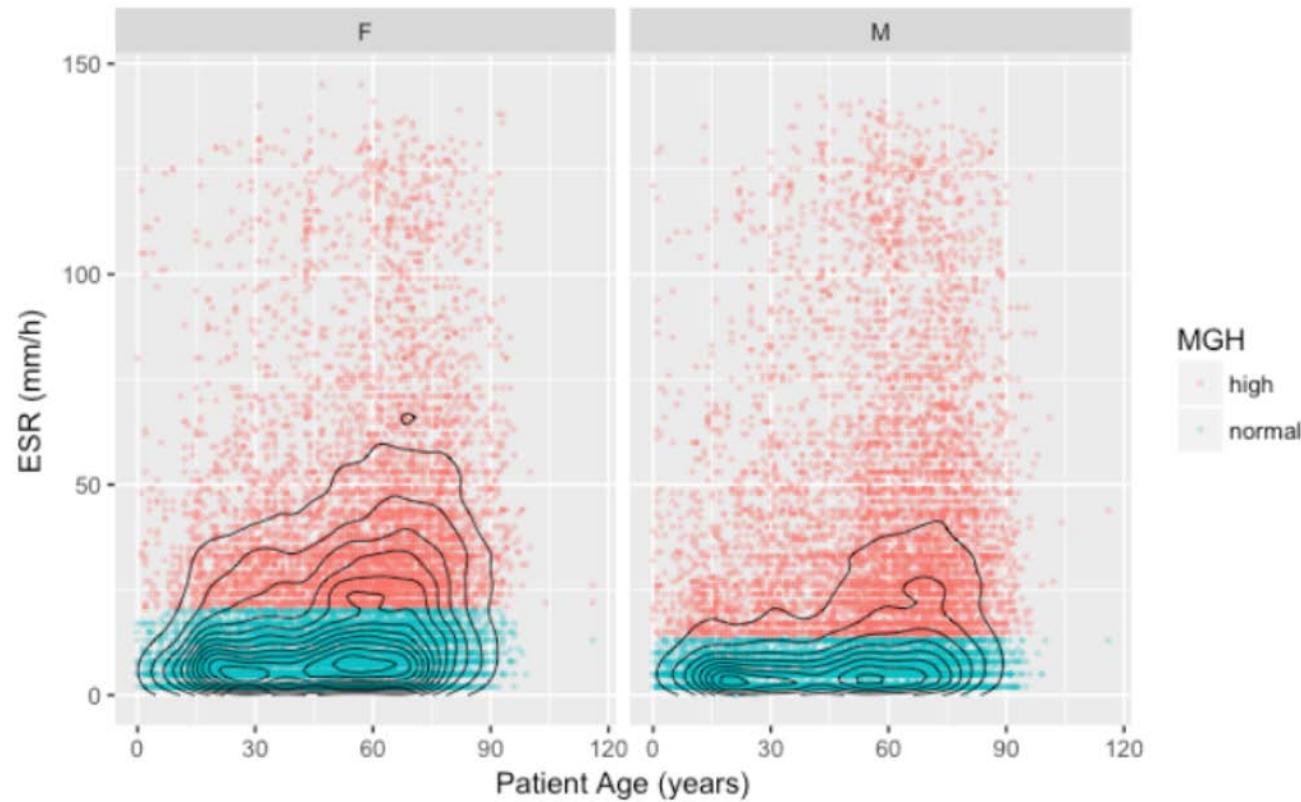
Code



Subsampling and density contour plotting did not reveal separable populations. This graph clearly shows that female patients tended to have higher ESR levels than male patients, and that older patients tended to have higher ESR levels than younger patients. However, this effect was not universal, and many old patients of both genders had low ESR values.

Do the MGH reference ranges appear to appropriately capture patients with normal ESRs?
Abnormal ESRs?

Code



MGH's normal reference range appears to capture the majority of normal ESR results. It does not capture the tendency of older individuals to have higher ESR values. However, this tendency is not universal, and it is uncertain whether high ESR values in older individuals are physiologically normal.

Background and Objectives



Data Acquisition



Data Exploration and Clean-Up



Visualization and Modeling



Summary

Summary

- Briefly restate background, objectives, and data acquisition strategy
- Restate each important result
- Explain how the new insights will be used to improve patient care

Your Turn #4

How would you summarize the ESR report?

Discuss with your group mates.

