
title:
“In-
tro-
duc-
tion
to R
and
RStu-
dio”
au-
thor:
“Vic-
tor H
Tor-
res”
out-
put:
pdf_document:
de-
fault
html_document:
in-
cludes:
in_header:
header.html
css:
./lab.css
high-
light:
pyg-
ments
theme:
cerulean
toc:
true
toc_float:
true
edi-
tor_options:
chunk_output_type:
con-
sole

The
RStu-
dio
Inter-
face

The goal of this lab is to introduce you to R and RStudio, which you'll be using throughout the course both to learn the statistical concepts discussed in the course and to analyze real data and come to informed conclusions. To clarify which is which: R is the name of the programming language.

As
the
labs
progress,
you
are
en-
cour-
aged
to ex-
plore
be-
yond
what
the
labs
dic-
tate;
a
will-
ing-
ness
to ex-
peri-
ment
will
make
you a
much
bet-
ter
pro-
gram-
mer.
Be-
fore
we
get to
that
stage,
how-
ever,
you
need
to
build
some
basic
flu-
ency
in R.
To-
day
we
begin
with
the
fun-
da-
men-

Go
ahead
and
launch
RStudio.
You
should
see a
window
that
looks
like
the
image
shown
below.

The panel on the lower left is where the action happens. It's called the *console*. Every time you launch RStudio, it will have the same text at the top of the console telling you the version of R that you're running. Below that information is the *prompt*. As its name suggests, this prompt is really

The
panel
in the
upper
right
con-
tains
your
*envi-
ron-
ment*
as
well
as a
his-
tory
of the
com-
mands
that
you've
previ-
ously
en-
tered.

Any
plots
that
you
gen-
erate
will
show
up in
the
panel
in the
lower
right
cor-
ner.
This
is
also
where
you
can
browse
your
files,
ac-
cess
help,
man-
age
pack-
ages,
etc.

R
Pack-
ages

R is
an
open-
source
pro-
gram-
ming
lan-
guage,
mean-
ing
that
users
can
con-
tribute
pack-
ages
that
make
our
lives
eas-
ier,
and
we
can
use
them
for
free.
For
this
lab,
and
many
oth-
ers in
the
fu-
ture,
we
will
use
the
fol-
low-
ing R
pack-
ages:

- The
suite
of
tidy-
verse
pack-
ages:
for
data
wran-
gling
and
data
visu-
aliza-
tion -
open-
in-
tro:
for
data
and
cus-
tom
func-
tions
with
the
Open-
Intro
re-
sources

If these packages are not already available in your R environment, install them by typing the following three lines of code into the console of your RStudio session, pressing the enter/return key after each one.

Note that you can check to see which packages (and which versions) are installed by inspecting the *Packages* tab in the lower right panel of RStudio.

You
may
need
to
select
a
server
from
which
to
down-
load;
any
of
them
will
work.
Next,
you
need
to
load
these
pack-
ages
in
your
work-
ing
envi-
ron-
ment.
We
do
this
with
the
library
func-
tion.
Run
the
fol-
low-
ing
three
lines
in
your
con-
sole.

```
r
library(tidyverse)
library(openintro)
You
only
need
to in-
stall
pack-
ages
once,
but
you
need
to
load
them
each
time
you
re-
launch
RStu-
dio.
```

The
Tidy-
verse
pack-
ages
share
com-
mon
philoso-
phies
and
are
de-
signed
to
work
to-
gether.
You
can
find
more
about
the
pack-
ages
in the
tidy-
verse
at
tidy-
verse.org.

Cre-
ating
a re-
pro-
ducible
lab
re-
port

We
will
be
using
R
Mark-
down
to
cre-
ate
re-
pro-
ducible
lab
re-
ports.
See
the
fol-
low-
ing
videos
de-
scrib-
ing
why
and
how:
**Why
use
R
Mark-
down
for
Lab
Re-
ports?
Using
R
Mark-
down
for
Lab
Re-
ports
in
RStu-
dio**

This
file
(with
the
.Rmd
file
ex-
ten-
sion)
will
serve
as the
lab
re-
port.
You
can
just
type
your
an-
swers
in
this
docu-
ment
in-
stead
of
creat-
ing a
sepa-
rate
docu-
ment.

Going
for-
ward
you
should
re-
frain
from
typ-
ing
your
code
di-
rectly
in the
con-
sole,
and
in-
stead
type
any
code
(final
cor-
rect
an-
swer,
or
any-
thing
you're
just
try-
ing
out)
in the
R
Mark-
down
file
and
run
the
chunk
using
either
the
Run
but-
ton
on
the
chunk
(green
side-
ways
trian-
gle)
or by

Dr. Arbuthnot's
Baptism
Records
To
get
started,
let's
take
a
peek
at the
data.
r
data('arbuthnot',
package='openintro')
You
can
run
the
com-
mand
by

-
click-
ing
on
the
green
arrow
at the
top
right
of the
code
chunk
in the
R
Mark-
down
(Rmd)
file,
or -
putting
your
cur-
sor
on
this
line,
and
click-
ing
the
Run
but-
ton
on
the
upper
right
cor-
ner of
the
pane,
or -
hold-
ing
Ctrl-Shift-Enter,
or -
typ-
ing
the
code
in the
con-
sole.

This
com-
mand
in-
structs
R to
load
some
data:
the
Ar-
buth-
not
bap-
tism
counts
for
boys
and
girls.
You
should
see
that
the
envi-
ron-
ment
area
in the
upper
right-
hand
cor-
ner of
the
RStu-
dio
win-
dow
now
lists a
data
set
called
arbuthnot
that
has
82
ob-
serva-
tions
on 3
vari-
ables.
As
you
inter-
act
with

The
Ar-
buth-
not
data
set
refers
to the
work
of
Dr. John
Ar-
buth-
not,
an
18th
cen-
tury
physi-
cian,
writer,
and
math-
e-
mati-
cian.
He
was
inter-
ested
in the
ratio
of
new-
born
boys
to
new-
born
girls,
so he
gath-
ered
the
bap-
tism
records
for
chil-
dren
born
in
Lon-
don
for
every
year
from
1629
to

—
r
arbutnot

```
## #  
A  
tibble:  
82 x  
3 ##  
year  
boys  
girls  
##  
<int>  
<int>  
<int>  
##  
1  
1629  
5218  
4683  
##  
2  
1630  
4858  
4457  
##  
3  
1631  
4422  
4102  
##  
4  
1632  
4994  
4590  
##  
5  
1633  
5158  
4839  
##  
6  
1634  
5035  
4820  
##  
7  
1635  
5106  
4928  
##  
8  
1636  
4917  
4605  
##  
9  
1637  
4708  
4457  
##  
10  
1638
```

However,
print-
ing
the
whole
dataset
in the
con-
sole is
not
that
use-
ful.
One
ad-
van-
tage
of
RStu-
dio is
that
it
comes
with
a
built-
in
data
viewer.
Click
on
the
name
`arbuthnot`
in the
*Envi-
ron-
ment*
pane
(up-
per
right
win-
dow)
that
lists
the
ob-
jects
in
your
envi-
ron-
ment.
This
will
bring
up an
alter-

What
you
should
see
are
four
columns
of
num-
bers,
each
row
repre-
sent-
ing a
differ-
ent
year:
the
first
entry
in
each
row is
sim-
ply
the
row
num-
ber
(an
index
we
can
use to
ac-
cess
the
data
from
indi-
vid-
ual
years
if we
want),
the
sec-
ond is
the
year,
and
the
third
and
fourth
are
the
num-
bers

Note that the row numbers in the first column are not part of Arbutnot's data. R adds them as part of its print-out to help you make visual comparisons. You can think of them as the index that you see on the left side of a spreadsheet. In fact, the comparisons to a spreadsheet will

You
can
see
the
di-
men-
sions
of
this
data
frame
as
well
as the
names
of the
vari-
ables
and
the
first
few
ob-
serva-
tions
by
typ-
ing:
`r`
`glimpse(arbuttnot)`

```
##
Rows:
82
##
Columns:
3 ##
$
year
<int>
1629,
1630,
1631,
1632,
1633,
1634,
1635,
1636,
1637,
1638,
1639~
## $
boys
<int>
5218,
4858,
4422,
4994,
5158,
5035,
5106,
4917,
4703,
5359,
5366~
## $
girls
<int>
4683,
4457,
4102,
4590,
4839,
4820,
4928,
4605,
4457,
4952,
4784~
```

It is better practice to type this command into your console, since it is not necessary code to include in your solution file. This command should output the following

Rows:
 82
 Columns:
 3 \$
 year
 1629,
 1630,
 1631,
 1632,
 1633,
 1634,
 1635,
 1636,
 1637,
 1638,
 1639~
 \$
 boys
 5218,
 4858,
 4422,
 4994,
 5158,
 5035,
 5106,
 4917,
 4703,
 5359,
 5366~
 \$
 girls
 4683,
 4457,
 4102,
 4590,
 4839,
 4820,
 4928,
 4605,
 4457,
 4952,
 4784~

We
can
see
that
there
are
82
ob-
serva-
tions
and 3
vari-
ables
in
this
dataset.
The
vari-
able
names
are
year,
boys,
and
girls.
At
this
point,
you
might
notice
that
many
of the
com-
mands
in R
look
a lot
like
func-
tions
from
math
class;
that
is, in-
vok-
ing R
com-
mands
means
sup-
ply-
ing a
~~func-~~
tion
with
some
num

```
##  
Some  
Ex-  
plo-  
ration  
Let's  
start  
to ex-  
amine  
the  
data  
a  
little  
more  
closely.  
We  
can  
ac-  
cess  
the  
data  
in a  
single  
col-  
umn  
of a  
data  
frame  
sepa-  
rately  
using  
a  
com-  
mand  
like  
r  
arbutnnot$boys
```

##

[1]

5218

4858

4422

4994

5158

5035

5106

4917

4703

5359

5366

5518

5470

5460

4793

##

[16]

4107

4047

3768

3796

3363

3079

2890

3231

3220

3196

3441

3655

3668

3396

3157

##

[31]

3209

3724

4748

5216

5411

6041

5114

4678

5616

6073

6506

6278

6449

6443

6073

##

[46]

6113

6058

6552

6423

6568

6247

6548

This
com-
mand
will
only
show
the
num-
ber of
boys
bap-
tized
each
year.
The
dollar
sign
basi-
cally
says
“go
to the
data
frame
that
comes
be-
fore
me,
and
find
the
vari-
able
that
comes
after
me”.

1.
What
com-
mand
would
you
use to
ex-
tract
just
the
counts
of
girls
bap-
tized?
Try
it!
r
arbutnot\$girls

##

[1]

4683

4457

4102

4590

4839

4820

4928

4605

4457

4952

4784

5332

5200

4910

4617

##

[16]

3997

3919

3395

3536

3181

2746

2722

2840

2908

2959

3179

3349

3382

3289

3013

##

[31]

2781

3247

4107

4803

4881

5681

4858

4319

5322

5560

5829

5719

6061

6120

5822

##

[46]

5738

5717

5847

6203

6033

6041

6200

Notice that the way R has printed these data is different. When we looked at the complete data frame, we saw 82 rows, one on each line of the display. These data are no longer structured in a table with other variables, so they are displayed one right after another. Objects that print out in this way

```
###  
Data  
visu-  
aliza-  
tion  
R has  
some  
pow-  
erful  
func-  
tions  
for  
mak-  
ing  
graph-  
ics.  
We  
can  
cre-  
ate a  
sim-  
ple  
plot  
of the  
num-  
ber of  
girls  
bap-  
tized  
per  
year  
with  
the  
com-  
mand  
r  
ggplot(data  
=  
arbuthnot,  
aes(x  
=  
year,  
y =  
girls))  
+  
geom_point()
```



We use the `ggplot()` function to build plots. If you run the plotting code in your console, you should see the plot appear under the *Plots* tab of the lower right panel of RStudio. Notice that the command above again looks like a function, this time with arguments separated by commas.

With
`ggplot()`:

- The first argument is always the dataset.

- Next, you provide the variables from the dataset to be assigned to aesthetic elements of the plot, e.g. the x and the y axes.

- Finally, you use another layer, separated by a + to specify the geometric object for the plot. Since we want to41 scatter-plot,

For
in-
stance,
if you
wanted
to
visu-
alize
the
above
plot
using
a line
graph,
you
would
re-
place
`geom_point()`
with
`geom_line()`.
r
`ggplot(data
=
arbuthnot,
aes(x
=
year,
y =
girls))
+
geom_line()`



You
might
won-
der
how
you
are
sup-
posed
to
know
the
syn-
tax
for
the
`ggplot`
func-
tion.
Thank-
fully,
R
docu-
ments
all of
its
func-
tions
ex-
ten-
sively.
To
learn
what
a
func-
tion
does
and
its
argu-
ments
that
are
avail-
able
to
you,
just
type
in a
ques-
tion
mark
fol-
lowed
by
the
name
of the

Try
the
fol-
low-
ing in
your
con-
sole:
r
?ggplot
Notice
that
the
help
file
re-
places
the
plot
in the
lower
right
panel.
You
can
toggle
be-
tween
plots
and
help
files
using
the
tabs
at the
top of
that
panel.

1. Is there an apparent trend in the number of girls baptized over the years? How would you describe it? (To ensure that your lab report is comprehensive, be sure to include the code needed to make the plot as well as your written interpretation.)

There is a significant decrease of girls baptized from 1620 to 1660. After that year, the amount of girls baptized increase drastically until the 1700's, where we can see some reduction of the number of baptisms. The graphic below shows the changes on the amount of females bap

```
r
ggplot(data
=
arbuthnot,
aes(x
=
year,
y =
girls))
+
geom_line()
+
geom_point(
colour
=
'red'
)
```



```
###
R as
a big
calcu-
lator
```

Now,
sup-
pose
we
want
to
plot
the
total
num-
ber of
bap-
tisms.
To
com-
pute
this,
we
could
use
the
fact
that
R is
really
just a
big
calcu-
lator.
We
can
type
in
math-
emat-
ical
ex-
pres-
sions
like
`r`
5218
+
4683

[1]
9901

to see
 the
 total
 num-
 ber of
 bap-
 tisms
 in
 1629.
 We
 could
 re-
 peat
 this
 once
 for
 each
 year,
 but
 there
 is a
 faster
 way.
 If we
 add
 the
 vec-
 tor
 for
 bap-
 tisms
 for
 boys
 to
 that
 of
 girls,
 R will
 com-
 pute
 all
 sums
 si-
 mul-
 tane-
 ously.
 r
 arbutnnot\$boys
 +
 arbutnnot\$girls

[1]
9901
9315
8524
9584
9997
9855
10034
9522
9160
10311
10150
10850

[13]
10670
10370
9410
8104
7966
7163
7332
6544
5825
5612
6071
6128

[25]
6155
6620
7004
7050
6685
6170
5990
6971
8855
10019
10292
11722

[37]
9972
8997
10938
11633
12335
11997
12510
12563
11895
11851
11775
12399
40

[49]
12626
12601

What
you
will
see
are
82
num-
bers
(in
that
packed
dis-
play,
be-
cause
we
aren't
look-
ing at
a
data
frame
here),
each
one
repre-
sent-
ing
the
sum
we're
after.
Take
a
look
at a
few of
them
and
verify
that
they
are
right.

Adding
a new
vari-
able
to the
data
frame

We'll
be
using
this
new
vec-
tor to
gen-
erate
some
plots,
so
we'll
want
to
save
it as
a per-
ma-
nent
col-
umn
in our
data
frame.
r
arbuthnot
<-
arbuthnot
%>%
mutate(total
=
boys
+
girls)

The `%>%` operator is called the **pip-ing** operator. It takes the output of the previous expression and pipes it into the first argument of the function in the following one. To continue our analogy with mathematical functions, $\mathbf{x} \%>\% \mathbf{f}(y)$ is equivalent to $\mathbf{f}(\mathbf{x}, y)$.

—

*“Take
the
arbutnnot
dataset
and
pipe
it
into
the
mutate
func-
tion.
Mu-
tate
the
arbutnnot
data
set by
creat-
ing a
new
vari-
able
called
total
that
is the
sum
of the
vari-
ables
called
boys
and
girls.
Then
as-
sign
the
re-
sult-
ing
dataset
to the
object
called
arbutnnot,
i.e. over-
write
the
old
arbutnnot
dataset
with
the
new
one
con-
tain-
ing*

This
is
equiv-
alent
to
going
through
each
row
and
adding
up
the
boys
and
girls
counts
for
that
year
and
record-
ing
that
value
in a
new
col-
umn
called
total.

You'll see that there is now a new column called **total** that has been tacked onto the data frame. The special symbol `<-` performs an *assignment*, taking the output of one line of code and saving it into an object in your environment. In this case, you already have

You
can
make
a line
plot
of the
total
num-
ber of
bap-
tisms
per
year
with
the
com-
mand
r
ggplot(data
=
arbutnnot,
aes(x
=
year,
y =
total))
+
geom_line()



Similarly
 to
 you
 we
 com-
 puted
 the
 total
 num-
 ber of
 births,
 you
 can
 com-
 pute
 the
 ratio
 of the
 num-
 ber of
 boys
 to the
 num-
 ber of
 girls
 bap-
 tized
 in
 1629
 with
 r
 5218
 $/$
 4683
 $##$
 $[1]$
 1.114243
 or
 you
 can
 act
 on
 the
 com-
 plete
 columns
 with
 the
 ex-
 pres-
 sion

```

_____
r
arbutnnot
<-
arbutnnot
%>%
mutate(boy_to_girl_ratio
=
boys
/
girls)
You
can
also
com-
pute
the
pro-
por-
tion
of
new-
borns
that
are
boys
in
1629
r
5218
/
(5218
+
4683)
##
[1]
0.5270175

```

or
you
can
com-
pute
this
for all
years
si-
mul-
tane-
ously
and
ap-
pend
it to
the
dataset
r
arbuthnot
<-
arbuthnot
%>%
mutate(boy_ratio
=
boys
/
total)
Note
that
we
are
using
the
new
total
vari-
able
we
cre-
ated
ear-
lier in
our
calcu-
la-
tions.

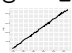
3.
Now,
gen-
erate
a plot
of the
pro-
por-
tion
of
boys
born
over
time.
What
do
you
see?

With
the
new
vari-
able
cre-
ated(total)
I am
able
to
gen-
erate
a lin-
ear
chart,
show-
ing
the
pro-
por-
tion
of
boys
born
over
time
dis-
play-
ing a
in-
crease
pat-
tern
over
the
years.

```

r
arbutnnot
<-
arbutnnot
%>%
mutate(total
=
boys
+
girls)
ggplot(data
=
arbutnnot,
aes(x
=
boys
, y
=
total))
+
geom_line()

```



Finally,
in ad-
dition
to
sim-
ple
math-
emat-
ical
oper-
ators
like
sub-
trac-
tion
and
divi-
sion,
you
can
ask R
to
make
com-
par-
isons
like
greater
than,
 $>$,
less
than,
 $<$,
and
equal-
ity,
 $==$.
For
ex-
am-
ple,
we
can
ask if
the
num-
ber of
births
of
boys
out-
num-
ber
that
of
girls
in
each
year
with

```
r
arbutnnot
<-
arbutnnot
%>%
mutate(more_boys
=
boys
>
girls)
```

This
com-
mand
adds
a new
vari-
able
to the
arbutnot
data
frame
con-
tain-
ing
the
val-
ues of
either
TRUE
if
that
year
had
more
boys
than
girls,
or
FALSE
if
that
year
did
not
(the
an-
swer
may
sur-
prise
you).
This
vari-
able
con-
tains
a dif-
ferent
kind
of
data
than
we
have
en-
countered
so
far.
All

More
Prac-
tice

In the
previ-
ous
few
pages,
you
recre-
ated
some
of the
dis-
plays
and
pre-
limi-
nary
anal-
ysis
of Ar-
buth-
not's
bap-
tism
data.
Your
as-
sign-
ment
in-
volves
re-
peat-
ing
these
steps,
but
for
present
day
birth
records
in the
United
States.
The
data
are
stored
in a
data
frame
called
present.

```
r  
data('present',  
package='openintro')
```

To find the minimum and maximum values of columns, you can use the functions `min` and `max` within a `summarize()` call, which you will learn more about in the following lab. Here's an example of how to find the minimum and maximum amount of boy births in a year:

```

r
arbutnot
%>%
summarize(min
=
min(boys),
max
=
max(boys))
## #
A
tibble:
1 x
2 ##
min
max
##
<int>
<int>
## 1
2890
8426
1.
What
years
are
in-
cluded
in
this
data
set?
What
are
the
di-
men-
sions
of the
data
frame?
What
are
the
vari-
able
(col-
umn)
names?

```

The years included in this data ranges from 1940 to 2002. The dimensions of the data frame are 63 observations and 3 variables. There are three variables: Year, Boys, and Girls

```
r
data('present',
package='openintro')
present
%>%
summarize(min
=
min(boys),
max
=
max(boys))
```

```
## #
A
tibble:
  1 x
  2 ##
min
max
##
<dbl>
<dbl>
## 1
1211684
2186274
1.
How
do
these
counts
com-
pare
to
Ar-
buth-
not's?
Are
they
of a
simi-
lar
mag-
ni-
tude?
```

The
range
of
years
are
big-
ger
on
the
Ar-
buth-
not
dataframe,
also
the
dim-
men-
sions
are
big-
ger,
it
con-
tains
more
rows,
r
arbuthnot
<-
present
%>%
summarize(min
=
min(boys),
max
=
max(boys))


1.
Make
a plot
that
dis-
plays
the
pro-
por-
tion
of
boys
born
over
time.
What
do
you
see?
Does
Ar-
buth-
not's
ob-
serva-
tion
about
boys
being
born
in
greater
pro-
por-
tion
than
girls
hold
up in
the
U.S.?
In-
clude
the
plot
in
your
re-
sponse.
Hint:
You
should
be
able
to
reverse
your
code
from
Ever

The
Pro-
por-
tion
of
boys
be-
ing
born
in
greater
pro-
por-
tion
than
girls
in
the
U.S
changes
from
the
first
data
frame
to
the
sec-
ond
one,
an
in-
crease
of
girls
is
evi-
dent
in
the
charts.

```

r
present
<-
present
%>%
mutate(total
=
boys
+
girls)
ggplot(data
=
present,
aes(x
=
boys
, y
=
total))
+
geom_line()

```



1. In what year did we see the most total number of births in the U.S.? *Hint:* First calculate the totals and save it as a new variable. Then, sort your dataset in descending order based on the total column. You can do this interactively in the data viewer by clicking on the ar79 rows next to the vari

```
r
present
%>%
arrange(desc(total))
```

The
most
total
num-
ber
of
births
in
the
U.S
was
in
1961,
with
a
total
of
4268326.

These
data
come
from
re-
ports
by
the
Cen-
ters
for
Dis-
ease
Con-
trol.
You
can
learn
more
about
them
by
bring-
ing
up
the
help
file
using
the
com-
mand
?present.

Re-
sources
for
learn-
ing R
and
work-
ing in
RStu-
dio

That was a short introduction to R and RStudio, but we will provide you with more functions and a more complete sense of the language as the course progresses.

In
this
course
we
will
be
using
the
suite
of R
pack-
ages
from
the
**tidy-
verse**.

The
book
R For
Data
Sci-
ence
by
Grole-
mund
and
Wick-
ham
is a
fan-
tastic
re-
source
for
data
anal-
ysis
in R
with
the
tidy-
verse.

If you
are
googling
for R
code,
make
sure to
also
in-
clude
these
pack-
age
names
in
your
search

These
cheat-
sheets
may
come
in
handy
through-
out
the
semester:
-
RMark-
down
cheat-
sheet
-
Data
trans-
for-
ma-
tion
cheat-
sheet
-
Data
visu-
aliza-
tion
cheat-
sheet
-
More
cheat-
sheets

Note that some of the code on these cheatsheets may be too advanced for this course. However the majority of it will become useful throughout the semester.
