

## R: fundamental skills for biologists

1-22 June 2022

This is an archive of questions asked via Slack during the workshop.

Question	Answer
<b>Session 1: Spreadsheets, organising data and first steps with R</b>	
Is it recommended that missing data is stated as NA as it is read as a non-value? What about using symbols instead, such as a dash that is used at times as well, would this cause issues?	Specifically in R - R knows how to handle NAs. Using a dash would cause a silent data mutation in R
In what cases is R better than python?	<p><b>Trainer A:</b> I'm a bit biased because I learnt R before I learnt a bit of python and I just find the R documentation a lot better overall for beginners than Python does. There's always this layer of assumed knowledge that I find frustrating whenever I'm trying to do something in Python.</p> <p>The plotting also I think is easier overall in R - python has 3 or more ecosystems for that iirc</p> <p><b>Trainer B:</b> Equally, the large R community means that there are a lot more pipelines available in R rather than python</p>
What can I do to see the console?	You might have minimised your console window - there is a minimise/maximise option in the pane
Is there a reason you would write in the console vs the script? Or is it an individual preference?	<p><b>Trainer A:</b> Use the console to try code out, script to record your code so you can re run it</p> <p><b>Trainer B:</b> Generally, I write in the script because then you have a record of the analysis that you are doing. I write in the</p>

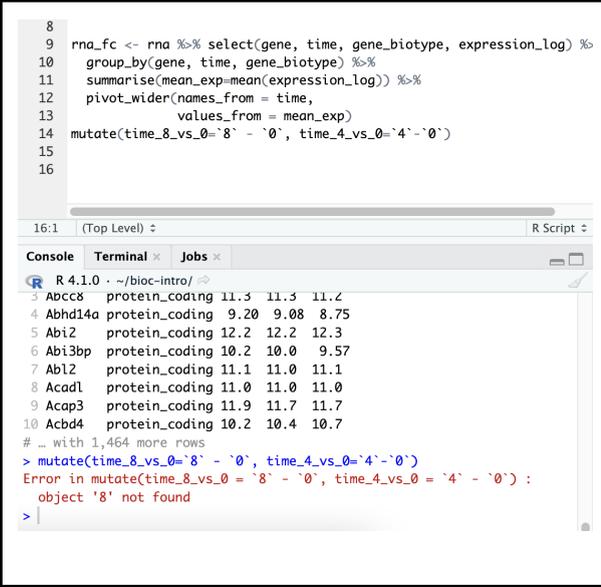
	<p>console when I have something that I want to try but don't necessarily think it'll be something that I want a record of ie, I'm quickly testing something.</p> <p><b>Trainer C:</b> Write in the console if it's something you know is not relevant to your analysis - for example, calling help functions. Otherwise even when testing things I'd personally put them in the script (so if they do work, I can go back and reuse or edit those commands). And also when you are doing once only actions like installing packages</p>
<p>What does the bioc-intro.Rproj file do?</p>	<p>The project creates a directory that will act as a starting point for running .R scripts. The project structure also allows R to store information of each session that you run, such as the objects that you create and the history of commands that you have run</p>
<p>Is there a way to move these files into our own personal folders OR is there a way to save into these folders upfront from the outset?</p>	<p>You can export the data. Click the checkbox for the folder you want to keep, click the <code>More</code> button for the dropdown menu, then click <code>Export</code></p>
<p>What makes data visualisation in R better than something like Prism? Example with your BarGraph Plot</p>	<p><b>Trainer A:</b> It is much much more customizable</p> <p><b>Trainer B:</b> I have not used Prism, but R is very flexible. Some packages in R can produce similar graphs in Prism, with but more complex graphs R can do a lot better</p> <p><b>Trainer C:</b> Will be able to reproduce much more easily/send to someone else who will be able to generate the same plot if given the data :)</p> <p><b>Trainer D:</b> huge variety and flexibility of plotting types. check out <a href="https://r-graph-gallery.com/">https://r-graph-gallery.com/</a> (includes example code for each)  The R Graph Gallery <a href="#">The R Graph Gallery – Help and inspiration for R charts</a> The R graph gallery displays hundreds of charts made with R, always providing the reproducible code.</p> <p><b>Trainer E:</b> It's easier to wrangle data into the format you need to plot in R e.g. to plot different subsets of the data</p>

	<p><b>Trainer F:</b> Also with the graphs generated in R you can then use the sister programs such as R shiny - where you can create interactive plots</p>
<p>I cannot see an output for any of these simple commands; eg 3 + 5 etc</p>	<p>You need to select 'Run' in the right-hand corner otherwise your code will not run into the console</p> <p>Or</p> <p>within the script, place your cursor at the line and then hit ctrl+enter / command + enter</p>
<p>Is there a point in writing "x=", given excluding it spits out the same result? e.g. round(digits=4, 3.1415926)</p>	<p>It's for making reading code easier but not strictly necessary.</p> <p>As long as you specify what the other arguments are, you won't have to specify 'x'</p>
<p>Can you use other script programs to write code e.g. VSCode then paste this into RStudio?</p>	<p>Yes, you can write your code in another program (or if it is already written in another format) and copy it over. Provided the language is R that you are copying over.</p> <p>Just note you would need to install the R packages you're using in both VStudio and RStudio if you wanted to run the code in both</p> <p><a href="https://code.visualstudio.com/docs/languages/r">https://code.visualstudio.com/docs/languages/r</a></p>
<p>Can you add a numeric value to a character vector and visa versa?</p>	<p>You can but it will change the data type - this is called coercion</p>
<p>How can we omit a number?</p>	<p><code>weight_g &lt;- weight_g[-1]</code> will delete the first entry</p>
<p>What does double vector mean?</p>	<p>A 'double precision floating point number'. <a href="https://en.wikipedia.org/wiki/Double-precision_floating-point_format">https://en.wikipedia.org/wiki/Double-precision_floating-point_format</a>. Basically a computer term for a value with some number of decimal places</p>
<p>Rule of thumb: round brackets - indicate a function and will need to have arguments within the brackets. e.g., round(5.56831, digits =2)</p>	

<p>square brackets - are for sub-setting objects ( I like to think of orientating within box or drawer); where you specify the location of the element you're searching for e.g., <code>molecules[2]</code></p>	
<p><b>Session 2: Manipulating and analysing data with dplyr</b></p>	
<p>When I have a data spreadsheet with colour and italic or bold text, will it confuse R or will the computer just ignore these characteristics? Or in other words: Can I structure my table for optical reasons for me to see things more easily (not store data in colour or format) or does that cause problems in R?</p>	<p>In general, it's just ignored. If you save as a .csv, the colour information is stripped out. If it's saved as an .xlsx file, the colour is there, but note that you need to use a special package to read these (e.g. readxl).</p>
<p>Tab complete is not prompting automatically.</p>	<p>That normally indicates that you are not in the right directory</p> <p>Check/set your working directory</p> <p>Might not be the issue but: you can only tab-autocomplete paths when they are within quote marks.</p> <p>e.g. <code>rna &lt;- read.csv("course-data/&lt;hit tab&gt;")</code></p> <p>You can also try just typing <code>read_csv("")</code>. Then with your cursor in between the quotation marks, hit tab. It should create a drop down menu of what is available</p>
<p>Can you count unique genes by the column header name rather than column number?</p>	<p><b>Absolutely!</b></p> <pre>length(unique(rna[,c("gene")]))</pre> <p>This is called subsetting with column names</p>
<p>What is actually displayed if we only did <code>rna[1]</code> instead of <code>rna[,1]</code>?</p>	<p><code>rna[1]</code> will show the column 1 in data.frame format whereas <code>rna[,1]</code> will show the contents of the column 1 as a vector</p>
<p>Is it always rows first and columns second in factors and the other way round in matrix?</p>	<p>In the <code>matrix()</code> function it doesn't matter if you give <code>ncol=</code> first or <code>nrow=</code> as you've named the arguments</p>
<p>Can you subset a list the same way you do a data frame?</p>	<p>First you have to specify the 'element' you want to subset. e.g. to get the 4th row of the</p>

	<p>cars element (which is the 4th in the list), you would use <code>l[[4]][4, ]</code></p> <p>or to get the first two elements, it's a more familiar subset command: <code>l[c(1, 2)]</code></p>
<p>Can we save lists as xlsx file?</p>	<p><b>Trainer A:</b> Thought it is always possible to combine all elements of a list in a dataframe with the <code>bind_rows()</code> function and export that as a xlsx file</p> <p>I recommend it only for specific reasons though, because the structure of your initial list will completely change</p> <p><b>Trainer B:</b> I find it useful sometimes to save a list of data.frames into a multi-tabbed excel file, but you need to use a special package (maybe xlsx??) to do so. Usually .tab or .tsv is simpler</p>
<p>Is filtering case sensitive? male (didn't work) vs Male</p>	<p>It is case sensitive when it comes to the filtering, but when I compare <code>select(rna, starts_with("T"))</code> and <code>select(rna, starts_with("t"))</code>, I get the same output</p> <p>The select helpers are not case sensitive by default, you change that by adding <code>ignore.case=FALSE</code></p> <p><a href="https://tidyselect.r-lib.org/reference/starts_with.html">https://tidyselect.r-lib.org/reference/starts_with.html</a></p>
<p>So mutate always adds the new column to the end of the data you start with? sounds dangerous. What if you (by mistake) run the same mutate command again? (error? extra columns? overwrite?)</p>	<p><b>Trainer A:</b> At this point of the exercise we haven't change our initial object yet</p> <p><b>Trainer B:</b> Mutate will change an existing column if you ask it to, but it only does it on the copy of the data you're working on.</p> <p><b>Trainer C:</b> Also good with the 'pipes' you can see in one place all the changes you have made and you will not make a new object until you assign it again...</p> <p><b>Follow up question:</b> So if I want to update the data object, I have to pipe it back to that object name?</p>

	<p><b>Trainer D:</b> Yes and it would not change your underlying raw data, so you can always start over</p> <p><b>Trainer E:</b> If you run the same mutate command again (and you've saved the output as an object) it will replace the column as it has the same name</p> <p><b>Follow up question:</b> so if I just leave a pipe hanging, that version of the data just disappears? it's not held as an object?</p> <p><b>Trainers B and C:</b> Yes ! Just displays to screen</p> <p><b>Trainer D:</b> it is like running a command without assigning it to an object</p>
<p><b>Session 3: Data visualisation</b></p>	
<p>I got this error.</p> <pre>&gt; ggplot(data=rna) Error in .Call.graphics(C_palette2, .Call(C_palette2, NULL)) :   invalid graphics state</pre>	<p>Try restarting RStudio</p>
<p>Can you easily modify upper and lower boundaries (i.e. 0- 50 000 instead of 100 000 on x axis in this example)?</p>	<p>You can also use</p> <pre>scale_x_continuous(limits = c(0, 50000))</pre>
<p>What does this message mean?</p> <pre>&gt; rna_fc &lt;- rna %&gt;% select(gene, time, gene_biotype, expression_log) %&gt;% + group_by(gene, time, gene_biotype) %&gt;% + summarize(mean_exp=mean(expression_log) ) `summarise()` has grouped output by 'gene', 'time'. You can override using the `.groups` argument.</pre>	<p>This is a standard warning from using <code>group_by</code> (it's about different ways of ungrouping data later)</p> <p>When we use <code>group_by()</code>, the output will retain the same grouping structure. This can be avoided if we use <code>%&gt;% ungroup()</code> at the end, OR, use <code>group_by(gene, time) %&gt;% summarize(mean_exp = mean(expression), .groups = 'drop')</code></p>
<p>I have this error message coming up:</p>	<p>You're missing the pipe after <code>pivot_wider</code></p>

<pre> 8 9 rna_fc &lt;- rna %&gt;% select(gene, time, gene_biotype, expression_log) %&gt;% 10   group_by(gene, time, gene_biotype) %&gt;% 11   summarise(mean_exp=mean(expression_log)) %&gt;% 12   pivot_wider(names_from = time, 13               values_from = mean_exp) 14   mutate(time_8_vs_0='8' - '0', time_4_vs_0='4'-'0') 15 16 </pre> 	
<p>Can we input hexadecimals for the colours?</p>	<p>Yes, see here  <a href="http://www.cookbook-r.com/Graphs/Colors_(ggplot2)/#hexadecimal-color-code-chart">http://www.cookbook-r.com/Graphs/Colors_(ggplot2)/#hexadecimal-color-code-chart</a></p> <p>There's also a table of standard R colours with names (like 'blue' :  <a href="http://www.stat.columbia.edu/~tzheng/files/Rcolor.pdf">http://www.stat.columbia.edu/~tzheng/files/Rcolor.pdf</a></p>
<p>What is the alpha?</p>	<p>Transparency: 1 is opaque, 0 fully transparent</p> <p><a href="https://ggplot2.tidyverse.org/reference/aes_colour_fill_alpha.html#:~:text=Alpha%20refers%20to%20the%20opacity.corresponding%20to%20more%20transparent%20colors">https://ggplot2.tidyverse.org/reference/aes_colour_fill_alpha.html#:~:text=Alpha%20refers%20to%20the%20opacity.corresponding%20to%20more%20transparent%20colors</a></p>
<p>How do you get the sample names to actually align with each boxplot or column for that matter?</p>	<p>Try adding:</p> <pre> theme(axis.text.x = element_text(angle = 90, hjust = 0.5, vjust = 0.5)) </pre> <p>e.g. usually googling 'align axis labels ggplot' can find something (I always forget this)  <a href="https://stackoverflow.com/questions/37488075/align-axis-label-on-the-right-with-ggplot2">https://stackoverflow.com/questions/37488075/align-axis-label-on-the-right-with-ggplot2</a></p>
<p>I got an error</p> <pre> &gt; ggplot(data=mean_exp_by_time, +         mapping=aes(x=time, y=mean_exp, group=gene)) + geom_line() Error in ggplot(data = mean_exp_by_time, mapping = aes(x = time, y = mean_exp, : object 'mean_exp_by_time' not found </pre>	<p>This is usually due to a mismatch between the object name in your code and what you have actually called the object. Check what you have called it in your environment in RStudio and change your code as needed.</p>

Does it matter how you spell colour? color?	No, it doesn't matter. ggplot is US/UK English agnostic and you can even mix and match.
Why when I write <code>facet_wrap(~gene)</code> , it does not show any changes in the plot area?	You need to include the + symbol on the line before this code
Why this time around the tilde is after sex and previously it was before gene?	It's related to the way that the plot displays a variable.  See also: <a href="https://benwhalley.github.io/just-enough-r/ggplot-details.html">https://benwhalley.github.io/just-enough-r/ggplot-details.html</a>
I am too deep in the woods. Why do I keep getting this message "  <code>geom_path: Each group consists of only one observation. Do you need to adjust the group aesthetic?</code> "	You might be missing a group or colour call
<b>Session 4: SummarizedExperiment and Bioconductor</b>	
For SummarizedExperiment:  I got this error  <pre>&gt; se &lt;- readRDS("course-data/data/GSE96870/GSE96870/se2.rds") Error in gzfile(file, "rb") : cannot open the connection In addition: Warning message: In gzfile(file, "rb") : cannot open compressed file 'course-data/data/GSE96870/GSE96870/se2.rds', probable reason 'No such file or directory'</pre>	Make sure you are implementing your command from the right level of the directory. Check your directory using <code>getwd()</code>  It might help running the full pathway (changing user to match your directory structure): <pre>se &lt;- readRDS("/home/user/Bioc_Intro/course-data/data/GSE96870/se2.rds")</pre>
I'm getting this error:  <pre>se_created &lt;- SummarizedExperiment(assays = SimpleList(counts=count_matrix) , colData = sample_metadata, rowRanges = gene_metadata)</pre>	try <code>rowData</code> instead of <code>rowRanges</code>  for the <code>SummarizedExperiment()</code> <code>rowRanges=</code> you need to supply a ranges (Granges/GenomicRanges) object. If you made one with <code>gene_metadata &lt;- rowRanges(se)</code> and gave that to <code>rowRanges=</code> that would work and that's what I did in the material. But in the demo I accidentally did <code>gene_metadata &lt;-</code>

<pre>Error in validObject(.Object) :   invalid class "RangedSummarizedExperiment" object: 1: invalid object for slot "rowRanges" in class "RangedSummarizedExperiment": got class "DFrame", should be or extend class "GenomicRanges_OR_GRangesList" invalid class "RangedSummarizedExperiment" object: 2:   'x@assays' is not parallel to 'x'</pre>	<p>rowData(se) which has no ranges 🤔 and it's a dataframe, and that you give to SummarizedExperiment() rowData=</p>
<p>Once tidySummarizedExperiment is loaded, is the dataset automatically viewed as a tibble?</p>	<p>Yes. You can change whether you want this functionality by changing options("restore_SummarizedExperiment_show" = TRUE)</p>
<p>I can't find .sample</p> <pre>&gt; se %&gt;% group_by(.samples) %&gt;% + summarise(total_counts = sum(counts)) tidySummarizedExperiment says: A data frame is returned for independent data analysis. Error in `dplyr::group_by()` : ! Must group by variables found in `.data`. ✘ Column `.samples` is not found.</pre>	<p>You made a typo. Instead of .samples try .sample</p>
<p>What about proteomics annotations? I could not find any in bioconductor</p>	<p>What about this? <a href="https://www.bioconductor.org/packages/release/data/experiment/vignettes/RforProteomics/inst/doc/RforProteomics.html#7_Annotation">https://www.bioconductor.org/packages/release/data/experiment/vignettes/RforProteomics/inst/doc/RforProteomics.html#7_Annotation</a></p> <p>I just googled proteomics annotation bioconductor 😊 but looking at the url I think it's not a package it's experiment data so in that section of bioconductor, see here <a href="https://www.bioconductor.org/packages/release/data/experiment/">https://www.bioconductor.org/packages/release/data/experiment/</a></p>