

TROUBLESHOOTING R

Isolating Issues and Asking Questions



Erik Kusch

erik.kusch@au.dk

Section for Ecoinformatics & Biodiversity

Center for Biodiversity and Dynamics in a Changing World (BIOCHANGE)

Aarhus University

17/02/2021

1 Dealing with Problems

- Problem-Sources
- Approaching Problems

2 Reproducibiliy

- What & Why
- How

3 Reporting Issues and Asking for Help

- Minimal Working Examples
- Minimality
- Completeness
- Reproducibility

1 Dealing with Problems

- Problem-Sources
- Approaching Problems

2 Reproducibiliy

- What & Why
- How

3 Reporting Issues and Asking for Help

- Minimal Working Examples
- Minimality
- Completeness
- Reproducibility

So your analysis is misbehaving...

Start by **investigating common error sources:**

Code Errors

Mis-Specified Statistics

Unexpected Findings

So your analysis is misbehaving...

Start by **investigating common error sources:**

Code Errors

Mis-Specified Statistics

Unexpected Findings

So your analysis is misbehaving...

Start by **investigating common error sources:**

Code Errors

- Capitalisation
- Mismatched...
 - Object types
 - Value modes
 - Object dimensions
- Directory structures
- Non-self-contained scripts

Mis-Specified Statistics

Unexpected Findings

So your analysis is misbehaving...

Start by **investigating common error sources:**

Code Errors

- Capitalisation
- Mismatched...
 - Object types
 - Value modes
 - Object dimensions
- Directory structures
- Non-self-contained scripts

Mis-Specified Statistics

- Assumptions not met
- Erroneous outliers
- Mis-interpretation of outputs
- Incorrect data handling

Unexpected Findings

So your analysis is misbehaving...

Start by **investigating common error sources:**

Code Errors

- Capitalisation
- Mismatched...
 - Object types
 - Value modes
 - Object dimensions
- Directory structures
- Non-self-contained scripts

Mis-Specified Statistics

- Assumptions not met
- Erroneous outliers
- Mis-interpretation of outputs
- Incorrect data handling

Unexpected Findings

"[...] in an experimental framework, surprising results are significant results, rather than signs of failure."

~ Curtin & Parker, 2014

So your analysis is misbehaving...

Start by **investigating common error sources:**

Code Errors

- Capitalisation
- Mismatched...
 - Object types
 - Value modes
 - Object dimensions
- Directory structures
- Non-self-contained scripts

Mis-Specified Statistics

- Assumptions not met
- Erroneous outliers
- Mis-interpretation of outputs
- Incorrect data handling

Unexpected Findings

"[...] in an experimental framework, surprising results are significant results, rather than signs of failure."

~ Curtin & Parker, 2014

How do we resolve these?

Errors and Warnings

Error

R: "I can't and won't. You won't force me. I quit. I am done."

```
> fit[5,100,]
Error in fit[5, 100, ] : subscript out of bounds
```

The problem it caused

The code that 'threw' the error

Message defined as error

Warning

R: "I can, but you might not like the output better look at that, my guy."

```
Warning message:
Removed 1 rows containing missing values (geom_path).
```

The behaviour being warned about

The declaration of a warning

The piece of code that caused the warning

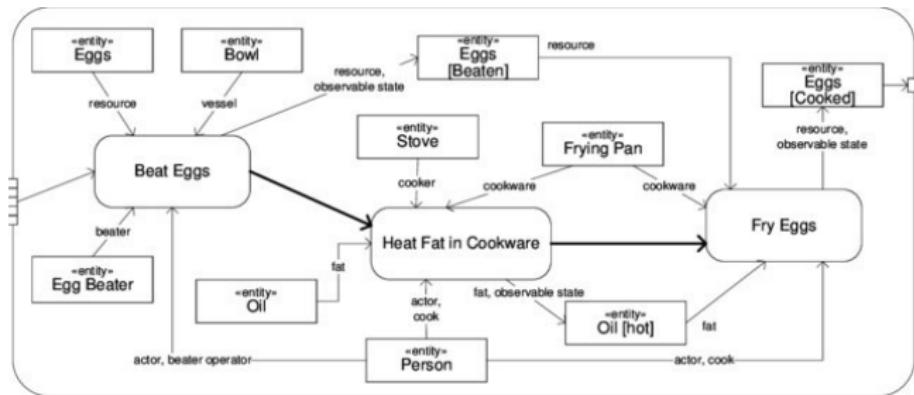
<http://rex-analytics.com/decoding-error-messages-r/>

Breaking problems down

- Investigate objects and their contents
- Read function help files (especially the arguments and values sections)
- Simulate data if in doubt

Breaking problems down

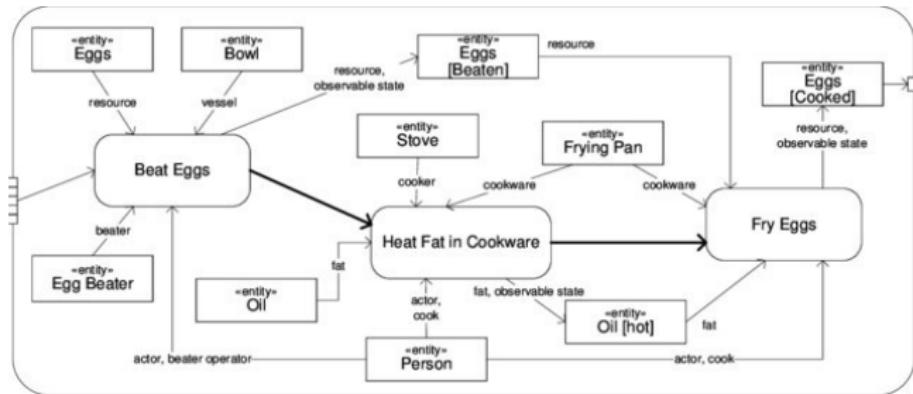
- Investigate objects and their contents
- Read function help files (especially the arguments and values sections)
- Simulate data if in doubt



Sousa et al., 2007

Breaking problems down

- Investigate objects and their contents
- Read function help files (especially the arguments and values sections)
- Simulate data if in doubt



Sousa et al., 2007

Your analyses hardly ever fail at all steps.

1 Dealing with Problems

- Problem-Sources
- Approaching Problems

2 Reproducibiliy

- What & Why
- How

3 Reporting Issues and Asking for Help

- Minimal Working Examples
- Minimality
- Completeness
- Reproducibility

What it is and why you need to care

Reproducibility ensures:

- Honest and true reporting
- Results can be scrutinised
- Problems can be traced and identified
- Co-authors and reviewers get less annoyed

Our Goal:

An analysis is **reproducible** when it can be executed by different researchers across operating systems with a **single button prompt** and return the same results for everyone.

Coding Styles

■ Code

■ *Consistency*

- The same result on every execution.
- The same naming scheme throughout.

■ *Self-contained* - Code runs in an empty \mathbb{R} environment to completion.

■ Documentation

■ *Commenting* - Informative comments for:

- What each line is doing/supposed to do.
- Why each line is doing what it does/is supposed to do.

■ *Header* - Information of what the script does, and who authored it.

■ *File-Versioning* - GitHub or an equivalent to trace when a script broke.

Packages

Never include forced package installation in a script.

Instead, do something like this:

```
install.load.package <- function(x) {  
  if (!require(x, character.only = TRUE))  
    install.packages(x, repos='http://cran.us.r-project.org')  
  require(x, character.only = TRUE)  
}  
  
package_vec <- c(  
  "rethinking", # for quadratic approximation of posteriors  
  "reshape2" # for data reformatting  
)  
  
sapply(package_vec, install.load.package)
```

Working Directories

Never include **hard-coded** directory indexing.

Make an Omlette out of Easter Eggs

- 1 Find the egg(s)/data
- 2 Open the egg(s)/data
- 3 Make the egg(s) into an omelette /
do your analysis

.... but... where are the eggs/the data?

Working Directories

Never include **hard-coded** directory indexing.

Make an Omlette out of Easter Eggs

- 1 Find the egg(s)/data
- 2 Open the egg(s)/data
- 3 Make the egg(s) into an omelette /
do your analysis
.... but... where are the eggs/the data?

Working Directories

Never include **hard-coded** directory indexing.

Make an Omlette out of Easter Eggs

- 1 Find the egg(s)/data
 - 2 Open the egg(s)/data
 - 3 Make the egg(s) into an omelette / do your analysis
- but... where are the eggs/the data?



Working Directories

Never include **hard-coded** directory indexing.

Make an Omlette out of Easter Eggs

- 1 Find the egg(s)/data
- 2 Open the egg(s)/data
- 3 Make the egg(s) into an omelette /
do your analysis
.... but... where are the eggs/the data?



Working Directories

Never include **hard-coded** directory indexing.

Make an Omlette out of Easter Eggs

- 1 Find the egg(s)/data
- 2 Open the egg(s)/data
- 3 Make the egg(s) into an omelette /
do your analysis
.... but... where are the eggs/the data?



Working Directories

Never include **hard-coded** directory indexing.

Make an Omlette out of Easter Eggs

- 1 Find the egg(s)/data
- 2 Open the egg(s)/data
- 3 Make the egg(s) into an omelette /
do your analysis
.... but... where are the eggs/the data?



But what if you want to hunt for Easter Eggs/data somewhere other than the garden at Achmore/your specific project folder?

Working Directories

Never include hard-coded directory indexing.

Soft-code your working directories from a base/project directory:

```
Dir.Base <- getwd() # to find the project folder  
Dir.Data <- file.path(Dir.Base, "Data") # index the data folder
```

Now we have our directories indexed without any hardcoded:

```
Dir.Base  
## [1] "D:/Documents/Teaching/Excursions-into-Biostatistics/Troubleshooting - Isolating Issues and Asking Questions"  
Dir.Data  
  
## [1] "D:/Documents/Teaching/Excursions-into-Biostatistics/Troubleshooting - Isolating Issues and Asking Questions/Da
```

Randomness

Always make random processes reproducible.

Imagine you survey the elevation of a random surface with paratroopers:

- Clouds obscure the landscape
- All flights are on the same route
- Troopers report elevation of their landing
- There is no drift of paratroopers
- Troopers jump at fixed intervals starting at a certain time after take-off



Randomness

Always make random processes reproducible.

Imagine you survey the elevation of a random surface with paratroopers:

- Clouds obscure the landscape
- All flights are on the same route
- Troopers report elevation of their landing
- There is no drift of paratroopers
- Troopers jump at fixed intervals starting at a certain time after take-off



Randomness

Always make random processes reproducible.

Imagine you survey the elevation of a random surface with paratroopers:

- Clouds obscure the landscape
- All flights are on the same route
- Troopers report elevation of their landing
- There is no drift of paratroopers
- Troopers jump at fixed intervals starting at a certain time after take-off

→ How can we get the same elevations from the troopers of two different flights?



Randomness

Always make random processes reproducible.

Imagine you survey the elevation of a random surface with paratroopers:

- Clouds obscure the landscape
- All flights are on the same route
- Troopers report elevation of their landing
- There is no drift of paratroopers
- Troopers jump at fixed intervals starting at a certain time after take-off



→ We let all flights commence their jump at the same time after take-off.

→ How can we get the same elevations from the troopers of two different flights?

Randomness

Always make random processes reproducible.

```
rnorm(5, 1, 0.2)
## [1] 0.8796 1.1960 1.2139 0.9883 0.9878
rnorm(5, 1, 0.2)
## [1] 1.1414 0.7078 1.3137 1.1270 1.0204
set.seed(42); rnorm(5, 1, 0.2)
## [1] 1.2742 0.8871 1.0726 1.1266 1.0809
set.seed(42); rnorm(5, 1, 0.2)
## [1] 1.2742 0.8871 1.0726 1.1266 1.0809
```

A seed needs to be set before **every** random process.

1 Dealing with Problems

- Problem-Sources
- Approaching Problems

2 Reproducibiliy

- What & Why
- How

3 Reporting Issues and Asking for Help

- Minimal Working Examples
- Minimality
- Completeness
- Reproducibility

What is a Minimal Working Example (MWE)?

A reproducible, simplified version of code producing a desired outcome.

- **Minimal.** Information and code is reduced as much as possible. This involves reducing:
 - Lines of R code
 - Number of R packages needed
 - Reliance on directory structures
 - Description of work, result, aim, and potential error
- **Complete.** All information, data, and code necessary to execute the MWE is presented.
- **Reproducible.** The MWE returns the same output on all operating systems when executed by any person. Ideally, this happens at the push of a single button.

Minimal Description

The Goal: I want to use a k-means clustering approach to identify biome classes across Alaska during the year 1982 using GIMMs NDVI 3g annual mean and seasonality values.

Method(s) & Material(s):

- Data:
 - GIMMs NDVI 3g (9x9km index of vegetation performance)
 - Natural Vector Shapefiles (shapefile-collection of states and provinces)
- Method:
 - K-Means-Clustering implemented via the `mclust` package

The Problem: Call to `mclustBIC()` produces the following error:

```
Error in mvnXII(data = data, prior = prior) : NA/NaN/Inf  
in foreign function call (arg 1)
```

Minimal Data - Before Reducing to MWE

The Data on our end before handing it off:

Data	
Alaska_Shp	Large SpatialPolygonsDataFrame (642.3 kB)
All1982_ras	Large RasterStack (279660 elements, 2.3 MB)
gimms_annual	Large RasterLayer (139830 elements, 1.1 MB)
gimms_raster	Large RasterBrick (3355920 elements, 26.9 MB)
gimms_raster_mvc	Large RasterBrick (1677960 elements, 13.4 MB)
gimms_seasonality	Large RasterLayer (139830 elements, 1.1 MB)
maxi	Large RasterLayer (139830 elements, 1.1 MB)
mini	Large RasterLayer (139830 elements, 1.1 MB)
Shapes	Large SpatialPolygonsDataFrame (4594 elements, 54.1 MB)
vals1982_mat	num [1:139830, 1:2] NA NA NA NA NA NA NA NA ...

values	
Dir_Base	"D:/Documents/Teaching/Excursions-into-Biostatistics/Troub...
Dir_Data	"D:/Documents/Teaching/Excursions-into-Biostatistics/Troub...
Dir_Plots	"D:/Documents/Teaching/Excursions-into-Biostatistics/Troub...
gimms_files	chr [1:2] "D:/Documents/Teaching/Excursions-into-Biostatistics/...
indices	num [1:24] 1 1 2 2 3 3 4 4 5 5 ...
Negatives	int [1:102595] 264 266 292 878 898 905 1439 1452 1453 1476...
Position	1665L

1982Mean.nc	16/02/2021 15:25	NC File	554 KB
1982Season.nc	16/02/2021 15:25	NC File	554 KB
ndv3g.geo_v1_1982_0106.nc4	09/03/2020 12:55	NC4 File	437,492 KB
ndv3g.geo_v1_1982_0712.nc4	09/03/2020 12:47	NC4 File	437,492 KB
ne_10m_admin_1_states_provinces.cpg	09/03/2020 14:05	CPG File	1 KB
ne_10m_admin_1_states_provinces.dbf	09/03/2020 14:05	DBF File	11,327 KB
ne_10m_admin_1_states_provinces.prj	09/03/2020 14:05	PRJ File	1 KB
ne_10m_admin_1_states_provinces.READ...	09/03/2020 14:05	Firefox HTML Doc...	29 KB
ne_10m_admin_1_states_provinces.shp	09/03/2020 14:06	SHP File	20,506 KB
ne_10m_admin_1_states_provinces.shx	09/03/2020 14:06	SHX File	36 KB
ne_10m_admin_1_states_provinces.VERS...	09/03/2020 14:05	Text Document	1 KB
Shapes.zip	09/03/2020 14:05	WinRAR ZIP archive	14,259 KB

Data needs to be **reduced** to the important parts (All1982_ras, in this case). We can export an R environment with the `save.image()` function.

The Data for the MWE that we hand over:

Data	
All1982_ras	Large RasterStack (279660 elements, 2.3 MB)

Data.rar	16/02/2021 15:18	WinRAR archive	181,813 KB
Workspace.RData	16/02/2021 15:16	R Workspace	294 KB

Minimal Code

Original Script:

```

8+ ## PREABLE #####
9 rm(list=ls()) # clearing the entire environment
10
11 dir_base <- getwd() # identifying the current directory
12 dir_base <- file.path(dir_base, "data") # generating the folder path for data folder
13 dir.create(dir_base) # creating the folder
14 dir_plots <- paste(dir_base, "plots", sep="/") # generating the folder path for figures folder
15 dir.create(dir_plots) # creating the Figures folder
16
17 # install.packages("glims") # For glims NOt data download
18 # install.packages("raster") # for spatial raster format
19 library(raster)
20 # install.packages("mclust") # for mclust<clustering
21 library(mclust)
22 # install.packages("gridfs") # For shapefiles
23 library(gridfs)
24 # install.packages("ncdf4")
25 library(ncdf4)
26 # install.packages("ncdf4")
27 library(ncdf4)
28
29 ## DOWNLOADING DATA #####
30 ## CITIES DATA
31 glims_files <- downloadFiles(x = as.Date("1982-01-01"), # download from January 1982
32                               url = "http://www.naturalearthdata.com/http://www.naturalearthdata.com/download/10m/cultural/ne_10m_admin_1_states_provinces.zip",
33                               destfile = paste0(dir_data, "Shapes.zip", sep = "/")) # destination file
34 unzip(Shapes, exdir = dir_data, setNullDir = TRUE) # which file to unztp
35 exdir = dir_data
36
37 ## READING DATA #####
38 ## CITIES DATA
39 glims_files <- file.path(dir_data, list.files(path = dir_data, pattern = "ndvi3g.geo.vrt")) # identify the downloaded glims
40
41 glims_raster <- rasterizeGlims(x = glims_files, # the data we CASSICING
42                                remove_header = TRUE) # we don't need the header of the data
43
44 glims_raster # some information about the raster stack we created here
45
46 ## CROPPING DATA
47 Shapes <- readShapefile(dir_data, # where to look for the file
48                        "ne_10m_admin_1_states_provinces", # the file name
49                        verbose = TRUE) # print out what's going on of the loaded data
50 Position <- which(SHAPENAME == "Alaska") # Find the SHAPENAME name that is "Alaska" in our shapefiles
51 Alaska_Shp <- Shapes[Position,] # extract the alaska shapefile
52 plot(Alaska_Shp) # plot the Alaska shapefile
53 ext(Alaska_Shp) # extract crop coordinates the super-eastern coordinates
54 Alaska_Shp <- crop(Alaska_Shp, # what to crop
55                     extent(-130, -130, 55, 71)) # which extent to crop to
56 plot(Alaska_Shp) # visualizing the cropped product
57
58 ## HANDLING DATA #####
59 ## CROPPING
60 # crop
61 glims_raster <- crop(glims_raster, # data to be cropped
62                      ext(Alaska_Shp)) # extent of the Alaska extent
63 # mask (this keeps only the land parts from our shapefile)
64 glims_raster <- mask(glims_raster, # data to be masked
65                      Alaska_Shp) # cropped Alaska shapefile
66
67 ## MONTHLY AGGREGATES
68 # monthly indices
69 Indices <- monthlyIndices(glims_files) # generate month indices from the data
70 glims_raster_mvc <- monthlyComposite(glims_raster, # the data
71                                       Indices) # indices of the indices
72
73 glims_raster_mvc # some information about our monthly composite raster stack
74 negatives <- which(values(glims_raster_mvc) < 0) # identify all negative values
75 # values(glims_raster_mvc)[negatives] <- 0 # set threshold for barren land (negative)
76 glims_raster_mvc
77 plot(glims_raster_mvc)
78
79 ## ANNUAL AGGREGATES
80

```

MWE Script:

```

8+ ## PREABLE #####
9 rm(list=ls()) # clearing the entire environment
10
11 install.Load.package <- function(x) {
12   if (!require(x, character.only = TRUE))
13     install.packages(x, repos = "http://cran.us.r-project.org")
14   require(x, character.only = TRUE)
15 }
16 package_vec <- c(
17   "raster", # for handling raster format
18   "mclust" # for k-means <clustering
19 )
20 sapply(package_vec, install.Load.package)
21
22 ## LOADING DATA #####
23 load("workspace.Rdata")
24
25 ## CLUSTER ANALYSIS #####
26
27 ## EXTRACTING DATA
28 vals1982_mat <- readRDS("vals1982.rds")
29 rownames(vals1982_mat) <- 1:dim(vals1982_mat)[1] # column to index raster cell number

```

We removed code for:

- Download of data
- Cropping and masking of data
- Calculation of compound metrics

Complete Reporting

We also supply:

- Archive of our raw data directory
- Script which contains all the code we removed from full script when creating MWE script

We may also want to provide:

- sessionInfo()
- rstudioapi::versionInfo()
- Metadata if needed

MWE Script:

```

1 #-----#
2 #-----# This is a minimal R script for reporting issues and asking for help
3 #-----# It includes a download of the raw data, a reproducible
4 #-----# environment, and masking of negative values in the monthly composite
5 #-----# raster stack
6 #-----# Author: Kjetil R. Aasen
7 #-----# Date: 2018-01-18
8 #-----#
9
10
11 ##-----#
12 ##-----# clearing the entire environment
13 install.packages <- function(x) {
14   install.packages(x, character.only = TRUE)
15   install.packages(x, repos = "https://cran.us.r-project.org")
16   require(x, character.only = TRUE)
17 }
18 package_vec <- c(
19   "gridExtra", # handling easier formats
20   "mclust", # for k-means clustering
21 )
22 sapply(package_vec, install.packages)
23
24 ##-----#
25 ##-----# LOADING DATA #####
26
27 gimes_files <- downloadGIMES(x = as.Date("1982-01-01"), # downloaded from January 1982
28                               y = as.Date("1982-12-31"), # downloaded to December 1982
29                               dir = dir_out, # where download in data folder
30                               quiet = FALSE) # show download progress
31
32 ##-----#
33 ##-----# SHAPESFILE
34 download.FILES <- "https://www.naturalearthdata.com/http://www.naturalearthdata.com/download/10m/cultural"
35 gis_shp <- file.path(dir_out, "10m_admin_1_states_provinces.shp") # destination file
36 destfile <- paste0(dir_out, "/shapes.zip", sep = "/") # which file to unzip
37 unzip(paste0(dir_out, "/shapes.zip", sep = "/"), destfile) # which file to unzip
38 extdir <- dir_out # where to unzip it
39
40 ##-----#
41 ##-----# READING DATA #####
42
43 gimes_files <- file.path(dir_out, "10m_files/path = dir_out, pattern = "nh1314_glo.v1")) # Identify the downloaded gimes
44 files
45 gimes_raster <- rasterizeLwgs(x = gimes_files, # the data we rasterize
46                                remove_header = TRUE) # we don't need the header of the data
47
48 gimes_raster # some information about the raster stack we created here
49
50 ##-----#
51 ##-----# SHAPES
52 shapes <- readOGR(dir_out, "ne_10m_admin_1_states_provinces", # the file name
53                   verbose = FALSE) # we don't want an overview of the loaded data
54 position <- which(shapes$NAME == "Alaska") # we want the Alaska shapefile
55 Alaska_Shp <- shapes[position] # we want the Alaska shapefile
56 plot(Alaska_Shp) # plot it for inspection
57 Alaska_Shp <- crop(Alaska_Shp, # what to crop
58                     extent(-190, -110, 55, 75)) # which extent to crop to
59 plot(Alaska_Shp) # visualizing the cropped product
60
61 ##-----#
62 ##-----# MASKING DATA #####
63
64 ##-----# CROP
65 gimes_raster <- crop(gimes_raster, # data to be cropped
66                      extent(Alaska_Shp)) # extent of Alaska
67 # mask (this keeps all negative values, full from our shapefile)
68 gimes_raster <- mask(gimes_raster, # cropped data
69                      Alaska_Shp) # cropped Alaska shapefile
70
71 ##-----#
72 ##-----# MONTHLY INDICES
73 Indices <- monthlyIndices(gimes_files) # generate month indices from the data
74 gimes_raster_mvc <- monthlyComposite(gimes_raster, # the data
75                                       Indices) # Indices # the indices
76
77 ##-----#
78 ##-----# GIMES RASTER MVC # some information about our monthly composite raster stack
79 negatives <- which(values(gimes_raster_mvc) < 0) # identify all negative values
80 gimes_raster_mvc[negatives] <- 0 # set threshold for barren land (hence)
81 gimes_raster_mvc
82 plot(gimes_raster_mvc)
83
84 ##-----#
85 ##-----# AGGREGATE
86 gimes_raster_mvc <- calc(gimes_raster_mvc, fun = mean)
87

```

Reproducible Issues

We test our MWE script:

- In a new, empty directory by sourcing the script
- For completeness by running it on an empty R environment and by adding `rm(list=ls())` to the start of the script

Our script still produces the same error?

→ We ship it off to someone whom we hope can help us.

When you can't find an example of R code to steal from stack exchange



<https://knowyourmeme.com/photos/1297214-impossible-perhaps-the-archives-are-incomplete>