

## Pins for Sharing Clinical R Assets

Phil Bowsher, RStudio PBC;  
Ryan Johnson, RStudio PBC

### ABSTRACT

The pins package helps statistical programmers publish and share data sets, models, and other assets across the organization. Programmers can pin objects to a variety of “boards” and others can download and access the data/objects from a central location. The goal of this paper is to introduce statistical programmers to the pins package and how to share data sets across the organizations in a more reproducible way than legacy methods like emailing data. The main pins site for reference is here:

<https://pins.rstudio.com/>

### INTRODUCTION

The pins package was created in early 2019 with the primary purpose of publishing and sharing R objects. Often workflows are restricted to local data that collaborators do not have access to. Other times, the data has grown stale and needs to be refreshed. In other situations, data are in complex relational databases that can be tricky to query.

The above examples complicate data access across projects and teams. Moreover, managing access control can be difficult in situations involving sensitive data. The pins package can help remedy these situations by providing an accessible place to access data for various workflows. Moreover, **processes such as ETL (Extract, Transform, Load) jobs** can benefit from pins integration to ensure data is refreshed on an interval specified by the data owner.

This paper is for statistical programmers that are building clinical routines and have wondered how best to approach datasets in workflows.

### UNDERSTANDING PINS AND BOARDS

Pins can be used to publish R objects in a variety of formats (RDS, CSV, arrow/feather, JSON, etc.). Small data sets are the most commonly pinned object, such as ephemeral data or reference tables, not suited for databases, but still require a physical “home” that workflows can point to. These homes are known as “boards.” Boards have a variety of storage backends which are detailed here: <https://pins.rstudio.com/reference/index.html#boards>

### EXAMPLE WORKFLOW

A common workflow is as follows:

An R developer has processed raw data resulting in a polished dataset that is ready to be accessed by other workflows (eg., a clinical table pipeline or Shiny application). Here is a sample pins workflow:

#### **Load pins package**

This will make the functions of the pins package available to users.

#### **Create the board**

This will create a location to store the object.

#### **Write data to the board as a pin**

This process will transfer the object from RStudio to your board.

Then others at the organization can access the pin like this:

#### **Load pins package**

This will make the functions of the pins package available to users.

### Identify the board

This will tell R the location and name of the board we want to use that has the pin.

### Read the shared data

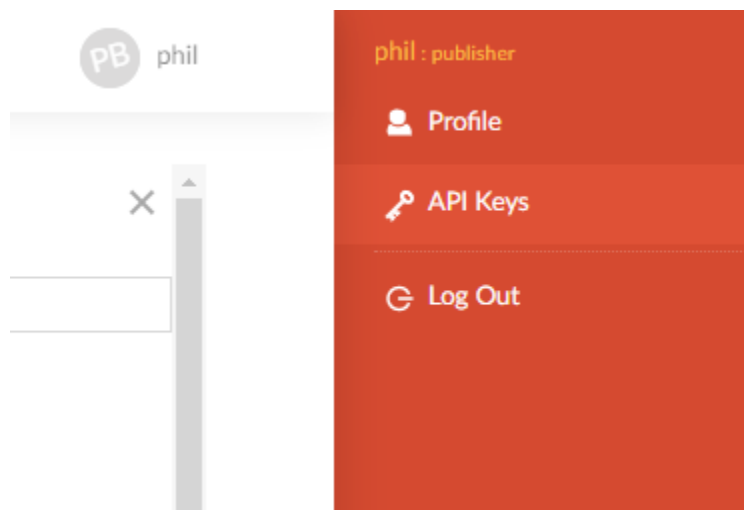
This will bring the pin into the RStudio connections tab. Users can also create their own datasets or objects as a copy for their own work.

## UNDERSTANDING API KEYS AND ACCESS CONTROLS WITH RSTUDIO CONNECT

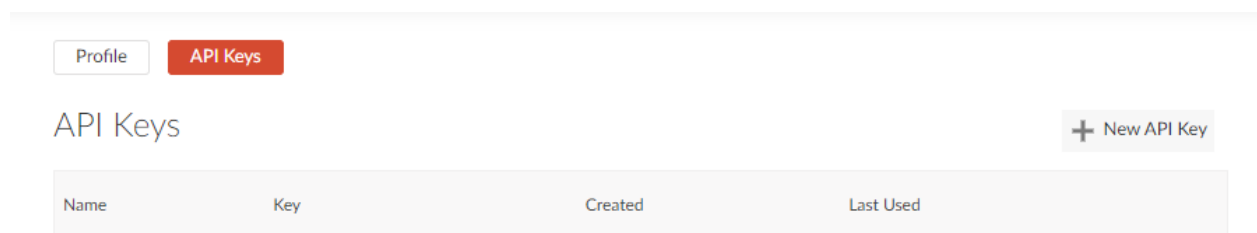
RStudio Connect can serve as a board for pins and provide an additional layer of access control. For users that wish to access a secure pin on RStudio Connect, a user specific API key will need to be generated as described below. Please ensure you treat API Keys like a password. API keys enable you to send requests to systems as if you were logged in manually. See the API Keys documentation for more details:

<https://docs.rstudio.com/connect/user/api-keys/>

Click on your name in the top right corner of RStudio Connect. Then click on API Keys:



Then click on New API Key:



Then create a name for your key.

## Create API Key



Enter a name for the new API Key

OK

Make sure you save your key in a safe place.

An important package to manage API keys is the `usethis` R package:

<https://usethis.r-lib.org/>

It is advised to place your API token as an environment variable in your ``.Renvi`` file. The `usethis` package makes it easy to edit this file for this purpose.

The code would look like this:

```
usethis::edit_r_environ()
```

Then enter in the file and save:

```
CONNECT_API_KEY="your-api-key"
```

```
CONNECT_SERVER="https://rstudio-connect-server"
```

These API keys allow users to programmatically access content on RStudio Connect and use the Connect Server API.

Below we will show a real workflow in R.

### ADVERSE EVENTS EXAMPLE

The FDA maintains an API that houses a number of high-value, high priority and scalable structured datasets, including adverse events, drug product labeling, and recall enforcement reports.

<https://open.fda.gov/>

The R package, `openfda`, provides convenient access to the OpenFDA API. Data in the API are often changing as new data are sent from hospitals, doctors, and other organizations.

Below we will work to create a reproducible workflow that creates a master dataset that is pinned to RStudio Connect. Then we will create a Shiny app and dashboard that will use the pinned data.

All of the examples below are available here:

<https://github.com/philbowsher/RMD-Shiny-session-2020-4-21/tree/master/RStudio-Connect-session-2019-12-06/shiny-days-master/Pins>

## CREATING THE LIVING DATASET AND ETL WORKFLOW

Below we will create the ETL job that we can run weekly to check for new data. We will schedule the ETL job to update our pinned data set each week so that all other processes and artifacts that depend on the data will always have the updated results.

For this example, we will use a R Notebook. It is important to either use a R Notebook, R Markdown or Quarto file as we will want the report to be executable by our production system, in this example it will be RStudio Connect.

[https://raw.githubusercontent.com/philbowsher/RMD-Shiny-session-2020-4-21/master/RStudio-Connect-session-2019-12-06/shiny-days-master/Pins/02\\_building\\_blocks.Rmd](https://raw.githubusercontent.com/philbowsher/RMD-Shiny-session-2020-4-21/master/RStudio-Connect-session-2019-12-06/shiny-days-master/Pins/02_building_blocks.Rmd)

Our ETL job will have the following structure:

1. Load the packages
2. Create helper functions to use the openFDA API and package to query adverse events data
3. Using our openFDA helper functions, we can pull adverse event data by gender for a specific drug and wrangle as needed
4. Publish our master data set as a pin to our board
5. Create some plots to test that our workflow is working

The pins code looks like this:

```
## Publish Pin
```{r}
library(pins)



# depends on CONNECT_SERVER and CONNECT_API_KEY variable
board <- board_rsconnect()

board %>% pin_write(adverse, "adverse", type = "csv")
```
```

For Connecting to your board in the code above, you will need to update your code with your API key like this:

```
board <- board_rsconnect(auth = "envvar", server =
"http://ec2-18-216-69-78.us-east-2.compute.amazonaws.com/rsconnect/", key =
Sys.getenv("VETIVER_API"))
```

Now your pinned dataset will show up in RStudio Connect like this:



Studio Connect

Content

People

Documentation

## Content



### adverse: a pinned 60 x 4 data frame


gt

Type: Pin

Last Deployed: Today

After clicking on the data set in RStudio Connect, viewers will see:

Content / adverse: a pinned 60 x 4 data frame

gt/adverse

**Last updated:** 2022-05-20 16:47:00 • **Format:** csv • **API:** v1  
Download data: [adverse.csv](#)  
▶ Raw metadata

Code

```
library(pins)
board <- board_rsconnect("envvar", server = "http://ec2-18-216-69-78.us-east-2.compute.amazonaws.com")
pin_read(board, "gt/adverse")
```

Preview (up to 100 rows)

| term                  | count | drug    | gender |
|-----------------------|-------|---------|--------|
| FATIGUE               | 10134 | Aspirin | male   |
| DYSPOEA               | 9459  | Aspirin | male   |
| DIARRHOEA             | 7942  | Aspirin | male   |
| DIZZINESS             | 7811  | Aspirin | male   |
| MYOCARDIAL INFARCTION | 7791  | Aspirin | male   |
| DRUG INEFFECTIVE      | 7622  | Aspirin | male   |
| NAUSEA                | 7223  | Aspirin | male   |

Info
Access
Runtime
Schedule
Tags
Vars
Logs

Sharing settings

- ☒ Anyone - no login required
- ☐ All users - login required
- ☐ Specific users or groups

Who can view or change this Pin

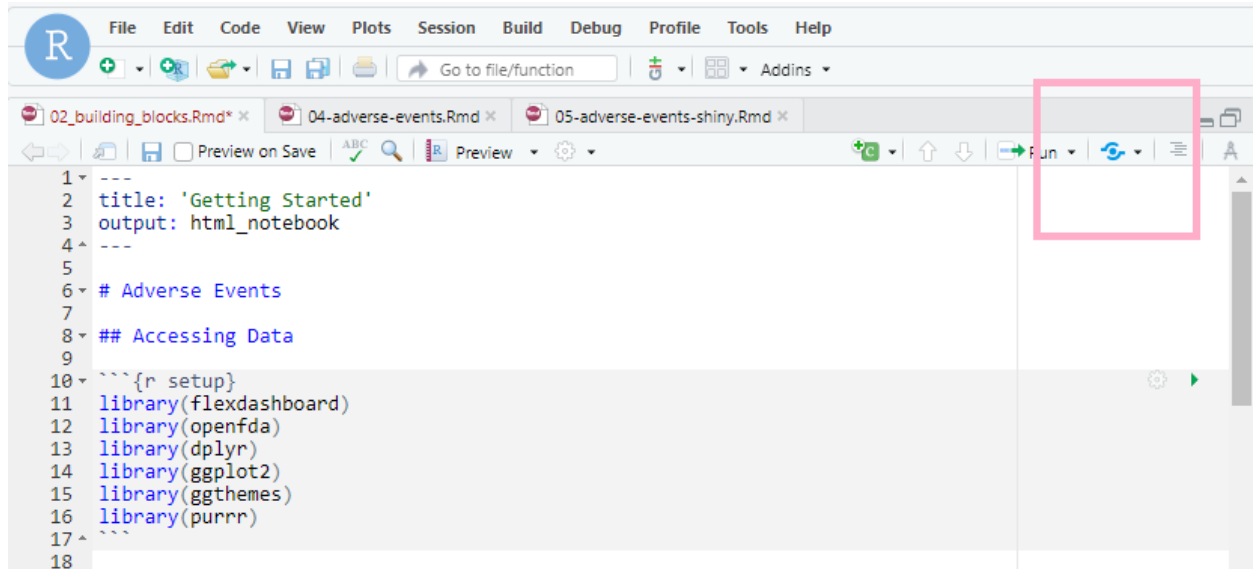
gt gt

Content URL

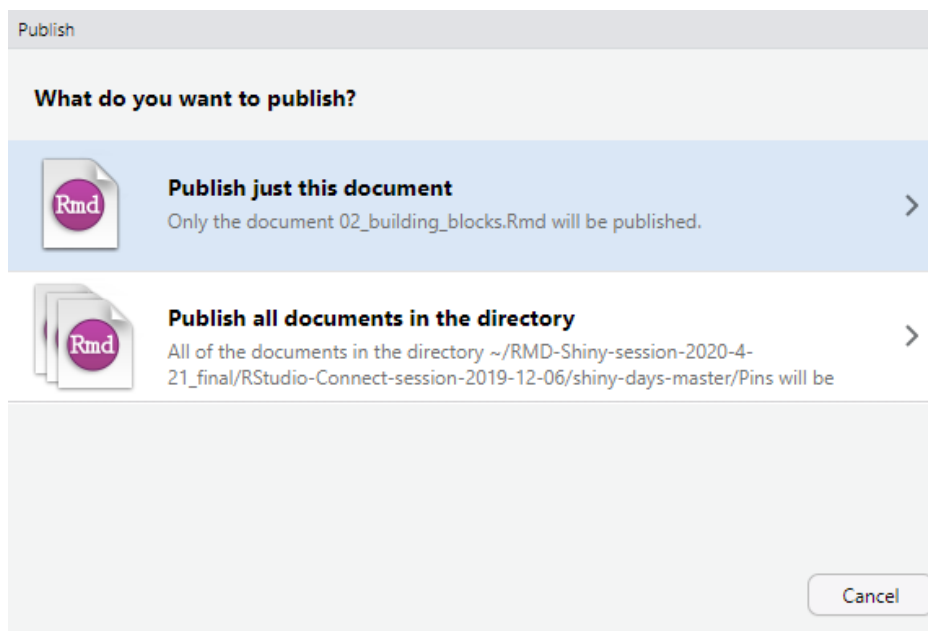
Customize

http://ec2-18-216-69-78.us-east-2.comput... Copy

Now that the data prep is working and our data is pinned properly, we can publish the R Notebook ETL to the RStudio Connect Production Environment. This can be done by clicking the blue publish button at the top right of the file like this:



Be sure to publish the document and any other files needed. In our case, it will just be the one ETL job.




If it is your first time publishing, there will be a series of steps for authenticating into the RStudio Connect Production server. After those steps, you can publish like this:

Publish

Back Publish

Publish Files From:

☒  02\_building\_blocks.Rmd

Add More...

Publish From Account: [Add New Account](#)

gt: ec2-18-216-69-78.us-east-2.compute.amazonaws.com

philbowsher: beta.rstudioconnect.com

Title:

02\_building\_blocks

Publish Cancel

After “publish” is clicked, RStudio Workbench will communicate with RStudio Connect and pass along the required information for rebuilding a sandbox where the ETL job will live in Production. In Production, the ETL job will have the required packages and dependencies to operate as needed:

```

Console Terminal x Render x Deploy x Launcher x
.../Pins/02_building_blocks.Rmd

Preparing to deploy document...DONE
Uploading bundle for document: 10...DONE
Deploying bundle: 53 for document: 10 ...
[Connect] Building R Markdown document...
[Connect] Bundle created with R version 4.0.2 is compatible with environment Local with R version 4.0.2
from /opt/R/4.0.2/bin/R
[Connect] Bundle requested R version 4.0.2; using /opt/R/4.0.2/bin/R which has version 4.0.2
[Connect] 2022/05/20 17:00:18.327768797 Using user agent string: 'RStudio R (4.0.2 x86_64-pc-linux-gnu x
86_64 linux-gnu)'
[Connect] 2022/05/20 17:00:18.770866320 Running on host: ip-10-8-11-176
[Connect] 2022/05/20 17:00:18.787107776 Linux distribution: Ubuntu 18.04.6 LTS (bionic)
[Connect] 2022/05/20 17:00:18.787121026 R version: 4.0.2
[Connect] 2022/05/20 17:00:18.787153894 # Validating R library read / write permissions -----
-----
[Connect] 2022/05/20 17:00:18.787155200 Using R library for packrat bootstrap: /opt/rstudio-connect/mnt/
R/4.0.2
[Connect] 2022/05/20 17:00:18.787164978 # Validating managed packrat installation -----
-----
[Connect] 2022/05/20 17:00:18.787165926 Vendored packrat archive: /opt/rstudio-connect/ext/R/packrat_0.
5.0-31_6a5fcdcfcc31a7b1b5e89ecff6c0916a77648e0d.tar.gz
[Connect] 2022/05/20 17:00:18.798992131 Vendored packrat SHA: 6a5fcdcfcc31a7b1b5e89ecff6c0916a77648e0d
[Connect] 2022/05/20 17:00:18.801804649 Managed packrat SHA: 6a5fcdcfcc31a7b1b5e89ecff6c0916a77648e0d
[Connect] 2022/05/20 17:00:18.802986501 Managed packrat version: 0.5.0.31
[Connect] 2022/05/20 17:00:18.803809334 Managed packrat is up-to-date.
[Connect] 2022/05/20 17:00:18.804203303 # Validating packrat cache read / write permissions -----
-----

```

After the packages are set up, the content will show in RStudio Connect:

Content / 02\_building\_blocks

Getting Started

Adverse Events

Accessing Data

Code

```
library(flexdashboard)
library(openfda)
library(dplyr)
```

Attaching package: 'dplyr'

The following objects are masked from 'package:stats':

Info

Access

Runtime

Schedule

Tags

Vars

Logs

Sharing settings

Anyone - no login required

All users - login required

Specific users or groups

Who can view or change this Document

gt

To the right, the user can set the required settings for the ETL job such as access and scheduling. Since we want our pinned dataset to be refreshed weekly, we can set the ETL job to run systematically to refresh our pinned data:

Info

Access

Runtime

Schedule

Tags

Vars

Lo

☒ Schedule output for default

Timezone

(GMT-04:00) America - Indianapolis

Start date & time

May 6, 2022 12 : 54 am pm

Reset to Local Time

Schedule type

Weekly

Run every 1 week.

Run every...

Sunday

Monday

Tuesday

Wednesday

Thursday

☒ Friday

Saturday

☒ Publish output after it is generated

☒ Send email after update

Owners are always notified unless they opt-out.

☒ Send to all collaborators

☐ Send to all viewers

Additional Recipients

8



This process will also send an email letting us know our ETL job ran successfully.

Please note that content like our ETL job above can also be deployed with github and or CI/CD. Please read this doc for more information.

<https://www.pharmasug.org/proceedings/2021/AD/PharmaSUG-2021-AD-205.pdf>

## CREATING ARTIFACTS AND CONTENT THAT USE OUR PINNED DATA

Now that our data are pinned to RStudio Connect, we can create reports, dashboards, Shiny apps, etc. that use the data. Below we will do that.

### **Flexdashboard R Markdown Dashboard**

Flexdashboard is a great package for creating dashboards:

<https://pkgs.rstudio.com/flexdashboard/articles/examples.html>

Below is a flexdashboard we created that uses our pinned data:

<https://github.com/philbowsher/RMD-Shiny-session-2020-4-21/blob/master/RStudio-Connect-session-2019-12-06/shiny-days-master/Pins/04-adverse-events.Rmd>

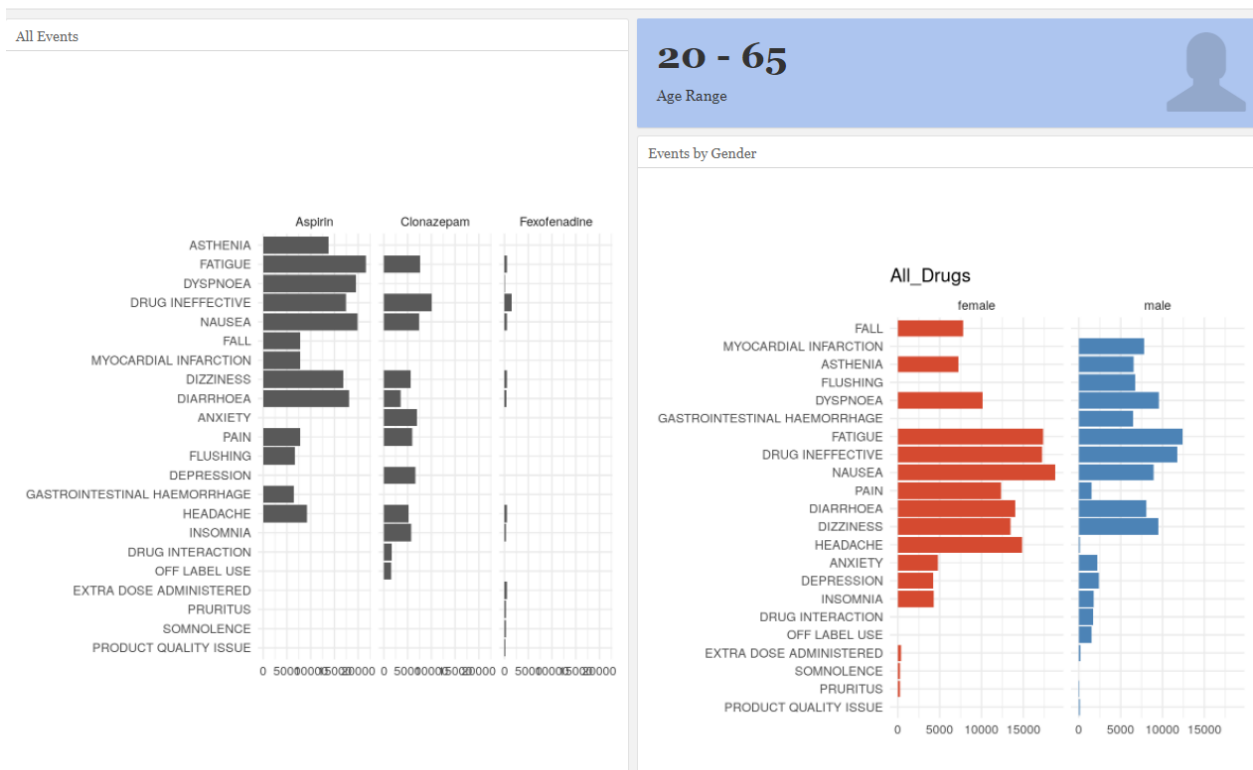
Now that the data are in RStudio Connect, we do not need to include the data wrangling routine in our dashboard code as it is now handled by our published ETL job. For example, the pink section blow is no longer required:

```
38
39 ~~~{r}
40 #this to be replaced by pins pull
41
42 age <- create_age(20,65)
43
44 # jnk <- capture.output(male <- get_adverse("1", params$drug, age))
45 # if (!is.null(male)) {
46 #   male$gender <- 'male'
47 # }
48
49 # jnk <- capture.output(female <- get_adverse("2", params$drug, age))
50 # if (!is.null(female)) {
51 #   female$gender <- 'female'
52 # }
53
54 # adverse <- rbind(male, female)
55
56 # adverse <- pin_get("Adverse", board = "rsconnect") %>% group_by(term) %>% summarise(count =
57   sum(count))
58
59 library(pins)
60 board <- board_rsconnect("envvar", server =
61   "http://ec2-18-216-69-78.us-east-2.compute.amazonaws.com/rsconnect", key =
62   Sys.getenv("VETIVER_API"))
63 adverse <- pin_read(board, "gt/adverse")
64
65 ~~~
```

This is important because now we have separated the often time-consuming data prep to a separate process outside of our artifact.

This process can also be repeated for anyone at the organization that is permitted to access to our pinned data. Below is the dashboard created by using pinned data:

## Adverse Events



```

35 ~ ```{r inputs}
36
37 #This to be replaced by pins pull
38
39 library(pins)
40 board <- board_register_rsconnect(auth = "envvar", account = "gt", server =
  "http://ec2-18-216-69-78.us-east-2.compute.amazonaws.com/rsconnect/",
41                                   key = Sys.getenv("VETIVER_API"))
42 adverseBase <- pin_read(board, "gt/adverse")
43
44 drugs <- unique(adverseBase$drug)
45
46
47 shiny::selectInput("sel_name", "Brand Name Drug", choices = drugs, selected = drugs[0])
48
49 ~ adverse <- reactive({
50   adverseBase %>%
51     filter(drug == input$sel_name)
52 ~ })
53
54 ~ ```

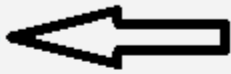
```

Now our pinned dataset has been converted to a Shiny reactive data frame and can be used like this:

```

59 ~ ### All Events
60
61 ~ ```{r}
62 ~ renderPlot({
63   req(adverse())
64   adverse() %>%
65   group_by(term) %>%
66   summarise(count = sum(count)) %>%
67   ggplot() +
68     geom_bar(aes(reorder(term,count), count), stat = 'identity') +
69     coord_flip() +
70     labs(
71       title = input$sel_name,
72       x = NULL,
73       y = NULL
74     ) +
75     theme_minimal()
76 ~ })
77
78 ~ ```

```



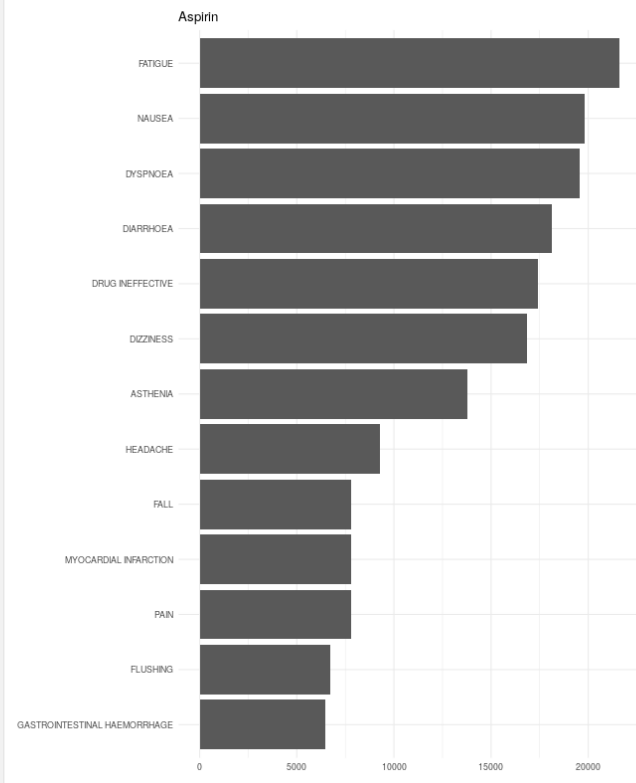
The Shiny app in Production that pulls in the pinned data is below:

## Adverse Events Phil

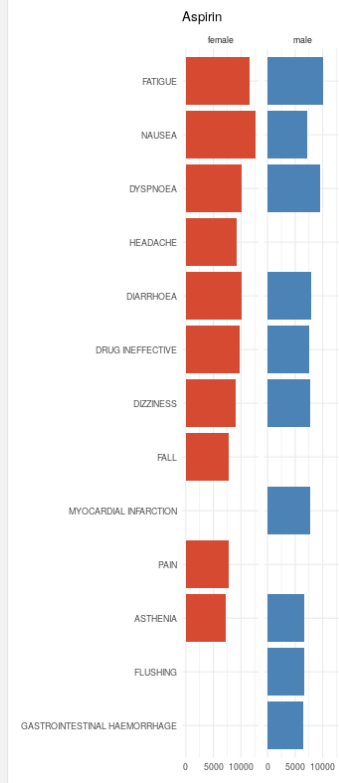
Brand Name Drug

Aspirin

All Events



Events by Gender



## CONTACT & SUMMARY

The information above highlights the exciting pins R package and example use cases, and how well established tools like these can help modernize clinical processes for reporting via reports, dashboards and Shiny applications.

Phil Bowsher

[phil@rstudio.com](mailto:phil@rstudio.com)

<https://github.com/philbowsher>