

Comparing Datasets

Sofia Bliss-Carrascosa

9/12/2022

```
library(tidyverse)
library(knitr)
```

V2: TESTING THE R-VERSION COMPARED TO ORIGINAL VERSION Loading R-version

```
v0_all_us_unsorted <- read_csv("allUS_unsorted.csv")
```

```
v1_all_us_polclaims <- read_csv("all_us_pol_claims_updated_sep15.csv")
```

```
v2_5parties <- v1_all_us_polclaims %>%
  filter(claimant_party == "Republican"|
         claimant_party == "Democratic"|
         claimant_party == "Independent"|
         claimant_party == "Libertarian"|
         claimant_party == "unknown_affiliation")
```

```
v2_5parties %>%
  count(publisher.site) %>%
  kable(caption = "All Claims in 5 Parties
            Sorted by Publisher")
```

Table 1: All Claims in 5 Parties Sorted by Publisher

publisher.site	n
cbsnews.com	167
checkyourfact.com	10
factcheck.org	1290
factcheck.thedispatch.com	58
newsweek.com	66
nytimes.com	464
politifact.com	4275
polygraph.info	3
poynter.org	10
thegazette.com	7
usatoday.com	19
vox.com	2
washingtonpost.com	1252

```
dim(v2_5parties)
```

```
## [1] 7623 13
```

Loading Grace Original Version

```
grace_v2 <- read_csv("grace_v2_CSVversion.csv")
```

```
grace_v2 %>%  
  count(publisher.site) %>%  
  kable(caption = "All Claims in 5 Parties  
          Sorted by Publisher")
```

Table 2: All Claims in 5 Parties Sorted by Publisher

<u>publisher.site</u>	<u>n</u>
cbsnews.com	167
checkyourfact.com	10
factcheck.org	1290
factcheck.thedispatch.com	58
newsweek.com	66
nytimes.com	464
politifact.com	4275
polygraph.info	3
poynter.org	10
thegazette.com	7
usatoday.com	19
vox.com	2
washingtonpost.com	1252

```
dim(grace_v2)
```

```
## [1] 7623 13
```

Data counts work out!

```
anti_join(v2_5parties, grace_v2)
```

```
## Joining, by = c("...1", "url", "title", "textualRating", "languageCode", "publisher.name", "publisher.site")
```

```
## # A tibble: 0 x 13  
## #   ... with 13 variables: ...1 <dbl>, url <chr>, title <chr>,  
## #   textualRating <chr>, languageCode <chr>, publisher.name <chr>,  
## #   publisher.site <chr>, reviewDate <dtm>, text <chr>, claimant <chr>,  
## #   claimDate <dtm>, claimant_party <chr>, reason <chr>
```

ALL MATCH!! WOOT!!

V3: TESTING THE R-VERSION COMPARED TO ORIGINAL VERSION Loading R-version

```
v3_deduped <- v2_5parties %>%
  distinct(url, text, title, .keep_all = TRUE)

v3_deduped %>%
  count(publisher.site) %>%
  kable(caption = "All Claims in 5 Parties
          Sorted by Publisher")
```

Table 3: All Claims in 5 Parties Sorted by Publisher

<u>publisher.site</u>	<u>n</u>
cbsnews.com	157
checkyourfact.com	10
factcheck.org	1238
factcheck.thedispatch.com	58
newsweek.com	63
nytimes.com	444
politifact.com	4176
polygraph.info	2
poynter.org	9
thegazette.com	7
usatoday.com	17
vox.com	1
washingtonpost.com	1179

```
dim(v3_deduped)
```

```
## [1] 7361 13
```

Loading Grace version

```
grace_v3 <- read_csv("gracededupedata.csv")

grace_v3 %>%
  count(publisher.site) %>%
  kable(caption = "All Claims in 5 Parties
          Sorted by Publisher")
```

Table 4: All Claims in 5 Parties Sorted by Publisher

<u>publisher.site</u>	<u>n</u>
cbsnews.com	157
checkyourfact.com	10
factcheck.org	1237
factcheck.thedispatch.com	58
newsweek.com	63
nytimes.com	444
politifact.com	4176
polygraph.info	2

<u>publisher.site</u>	<u>n</u>
poynter.org	9
thegazette.com	7
usatoday.com	17
vox.com	1
washingtonpost.com	1179

```
dim(grace_v3)
```

```
## [1] 7360 13
```

Data Counts are not the same: why?

```
anti_join(v3_deduped, grace_v3)
```

```
## Joining, by = c("...1", "url", "title", "textualRating", "languageCode", "publisher.name", "publisher.site")
```

```
## # A tibble: 14 x 13
```

```
##   ...1 url      title textualRating languageCode publisher.name publisher.site
##   <dbl> <chr>   <chr> <chr>          <chr>          <chr>          <chr>
## 1 17524 http://~ 'Brok~ FALSE          en          FactCheck.org factcheck.org
## 2 17855 http://~ 'Buy ~ Spins the Fa~ en          FactCheck.org factcheck.org
## 3 18059 http://~ CO2: ~ Some Pros; M~ en          FactCheck.org factcheck.org
## 4 16629 http://~ 'Deat~ Repeats Rura~ en          FactCheck.org factcheck.org
## 5 15616 https:~ Fact ~ FALSE          en          Newsweek       newsweek.com
## 6 5758  http://~ 'Song~ Pants on Fire en          PolitiFact     politifact.com
## 7 3341  http://~ 'Abol~ Needs Context en          PolitiFact     politifact.com
## 8 4296  https:~ 'Alwa~ FALSE          en          PolitiFact     politifact.com
## 9 3392  https:~ 8 tim~ 'Slip of the~ en          PolitiFact     politifact.com
## 10 5664  https:~ Fact~ Also means r~ en          PolitiFact     politifact.com
## 11 3078  https:~ 'Medi~ Accurate          en          PolitiFact     politifact.com
## 12 5616  https:~ 'We a~ FALSE          en          PolitiFact     politifact.com
## 13 15541 https:~ Trump~ FALSE          en          The New York ~ nytimes.com
## 14 14674 https:~ Analy~ Three Pinocc~ en          The Washingto~ washingtonpos~
## # ... with 6 more variables: reviewDate <dtm>, text <chr>, claimant <chr>,
## #   claimDate <dtm>, claimant_party <chr>, reason <chr>
```