

Lecture Notes on Neural Networks: Part 2B, Import Into R, Trim Rows, Encode Easy Columns

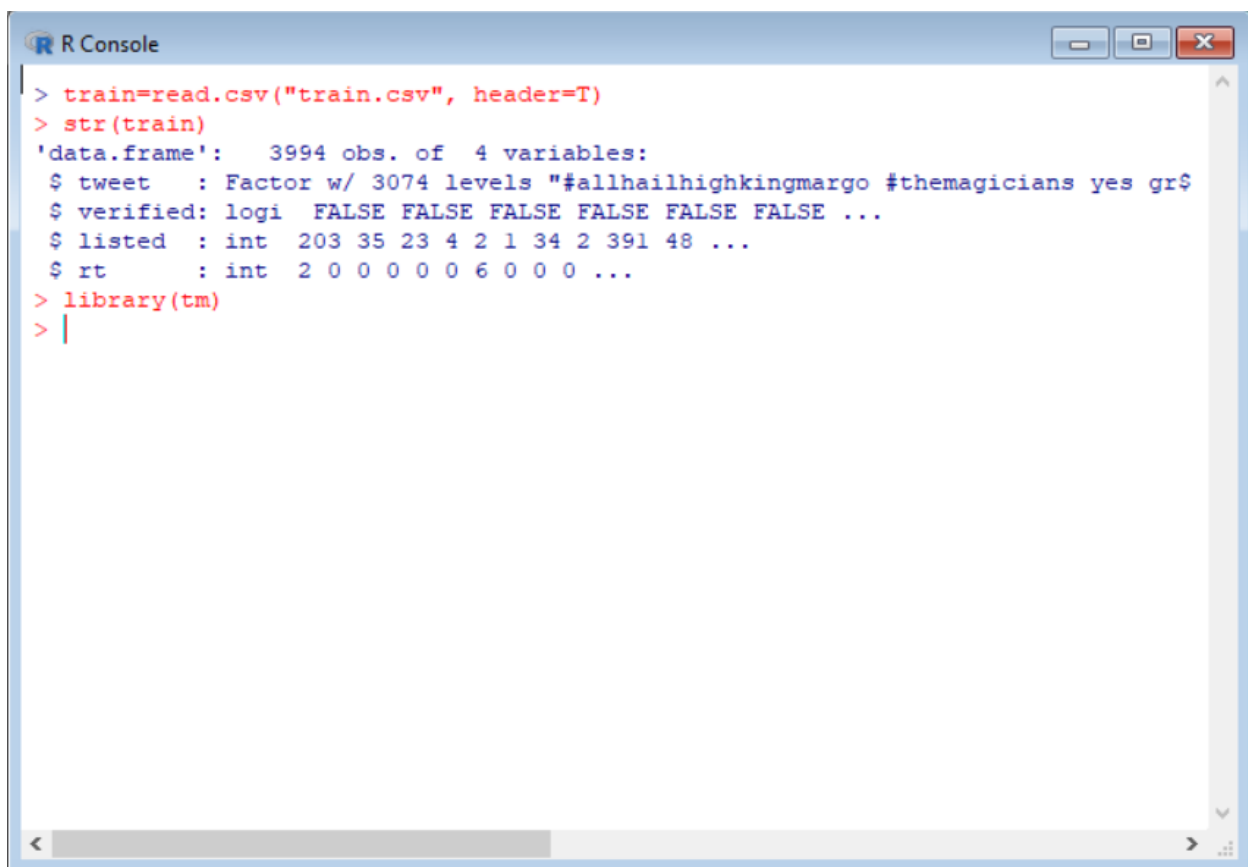
Nick V. Flor, nickflor@unm.edu

In the previous tutorial, we did some data cleaning, and copied into a new sheet just the columns of data we would use as inputs (tweet, listed, verified) as well as output (rt). We saved the file as *train.csv*.

In this tutorial, we'll import the file into R, then encode the “easy” columns. The easy columns are any column other than the tweet column—because converting text into numbers is not easy.

PART B. IMPORT INTO R, TRIM ROWS

Step 1. Import into R



```
> train=read.csv("train.csv", header=T)
> str(train)
'data.frame':  3994 obs. of  4 variables:
 $ tweet   : Factor w/ 3074 levels "#allhailhighkingmargo #themagicians yes gr$
 $ verified: logi  FALSE FALSE FALSE FALSE FALSE FALSE ...
 $ listed  : int   203 35 23 4 2 1 34 2 391 48 ...
 $ rt      : int    2 0 0 0 0 0 6 0 0 0 ...
> library(tm)
> |
```

Explanation:

I used *read.csv* to import into R, *str* to check the structure, and then *library* to load the text mining package (tm). Remember, we did an *install.packages("tm")*, early in the semester during the text-mining lectures.

Bottom Line: *train* contains our spreadsheet.

(continued on the next page)

Step 2. Trim Rows and Create a New Data Frame

```
R Console
> noRTs=grep("^rt ", train$tweet, invert=T)
> train1=data.frame(
+ tweet=train$tweet[noRTs],
+ verified=train$verified[noRTs],
+ listed=train$listed[noRTs],
+ rt=train$rt[noRTs])
> str(train1)
'data.frame':  2657 obs. of  4 variables:
 $ tweet   : Factor w/ 3074 levels "#allhailhighkingmargo #themagicians yes gr$
 $ verified: logi  FALSE FALSE FALSE FALSE FALSE FALSE ...
 $ listed  : int   203 35 23 4 2 1 2 391 48 56 ...
 $ rt      : int    2 0 0 0 0 0 0 0 0 0 ...
> |
```

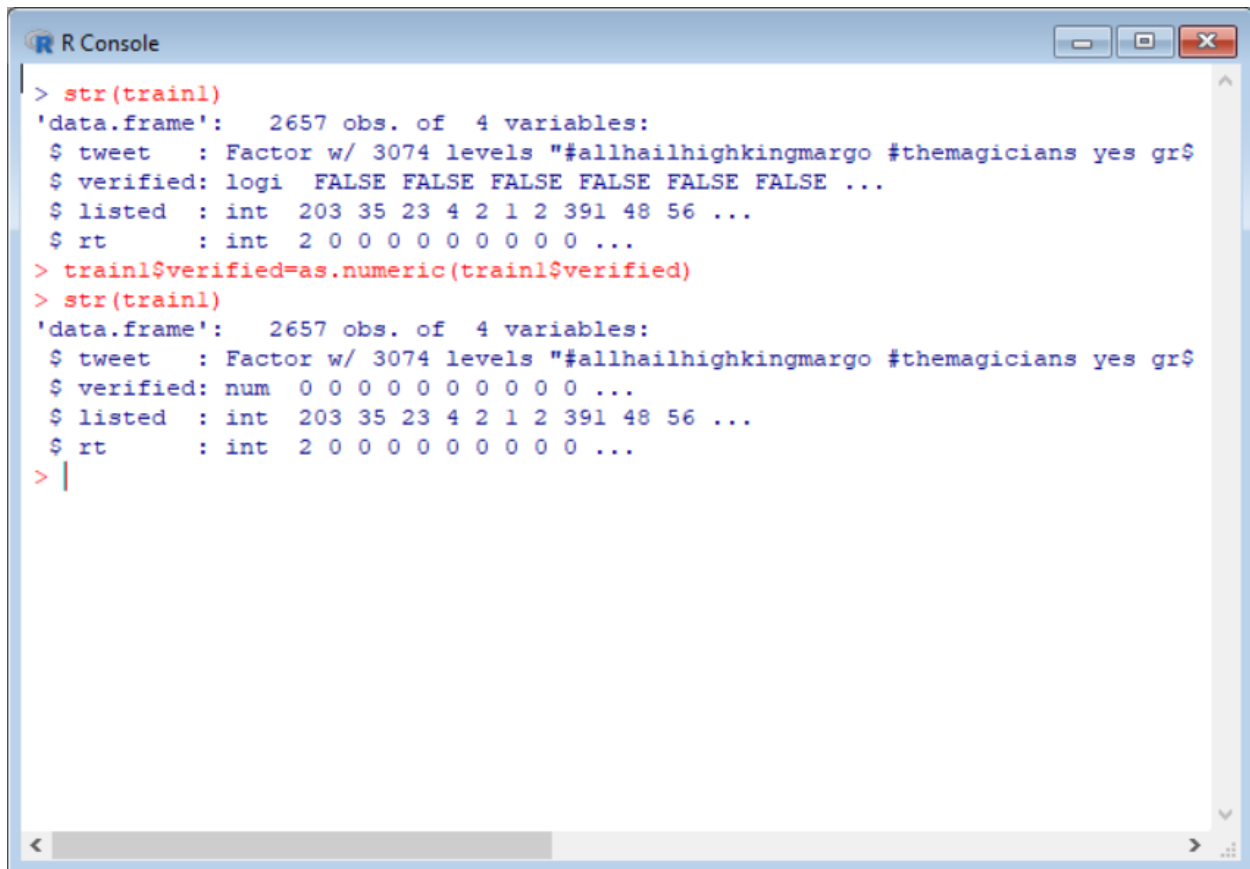
Explanation:

We do not want to train our network on retweets. A retweet always starts with the letter *rt* ("*rt*" as a regular expression); the *grep* returns the row number of any tweet which does not (*invert=T*) start with an *rt*. I store these row numbers in the variable *noRTs*.

Finally, I create a new *data.frame* named *train1*, with the same column headings as *train* (*tweet*, *verified*, *listed*, *rt*) but the [*noRTs*] means the tweets are not retweets.

Bottom line: *train1* is *train* without the retweets.

Step 3A. Encode Booleans (TRUE/FALSE) as 1/0



```
> str(train1)
'data.frame':  2657 obs. of  4 variables:
 $ tweet   : Factor w/ 3074 levels "#allhailhighkingmargo #themagicians yes gr$
 $ verified: logi  FALSE FALSE FALSE FALSE FALSE FALSE ...
 $ listed  : int   203 35 23 4 2 1 2 391 48 56 ...
 $ rt      : int    2 0 0 0 0 0 0 0 0 0 ...
> train1$verified=as.numeric(train1$verified)
> str(train1)
'data.frame':  2657 obs. of  4 variables:
 $ tweet   : Factor w/ 3074 levels "#allhailhighkingmargo #themagicians yes gr$
 $ verified: num    0 0 0 0 0 0 0 0 0 0 ...
 $ listed  : int   203 35 23 4 2 1 2 391 48 56 ...
 $ rt      : int    2 0 0 0 0 0 0 0 0 0 ...
> |
```

Explanation:

The function *as.numeric* coerces TRUE/FALSE to 1/0. I call *str* before and after *as.numeric* just to show you the effects.

Bottom Line: Whenever you want to convert a column of Booleans to 1's and 0's, use *as.numeric*.

(continued)

Step 3B. Encode the listed column as 1 if high, 0 otherwise

```
R Console
> llevel=mean(train1$listed)+sd(train1$listed)
> llevel
[1] 175.7202
> gtl=which(train1$listed>=llevel)
> lt1=which(train1$listed<llevel)
> train1$listed[gtl]=1
> train1$listed[lt1]=0
> |
>
```

Explanation:

What does it mean for a user to have a high level of being listed? Instead of arbitrarily picking a cut-off, I decided to define one standard deviation (*sd*) from the *mean* as the cut-off. I stored this cut-off in *llevel*. Thus, any users listed more than 175.72 times has a high listing (1), anything less than that is a low listing level (0).

I used the *which* function to return the listed row numbers greater than *llevel* (*gtl*). I also used the *which* function to return the listed row numbers less than *llevel* (*ltl*).

Finally, I set all rows greater than *llevel* to 1, and all rows less than *llevel* to 0.

Bottom Line: the listed values for all rows greater than *llevel* are now 1 (vs say 200), otherwise 0.

