

# TURNING A TWITTER HASHTAG INTO A WORD CLOUD

Jennifer Pannell  
Ecology postdoc, AUT

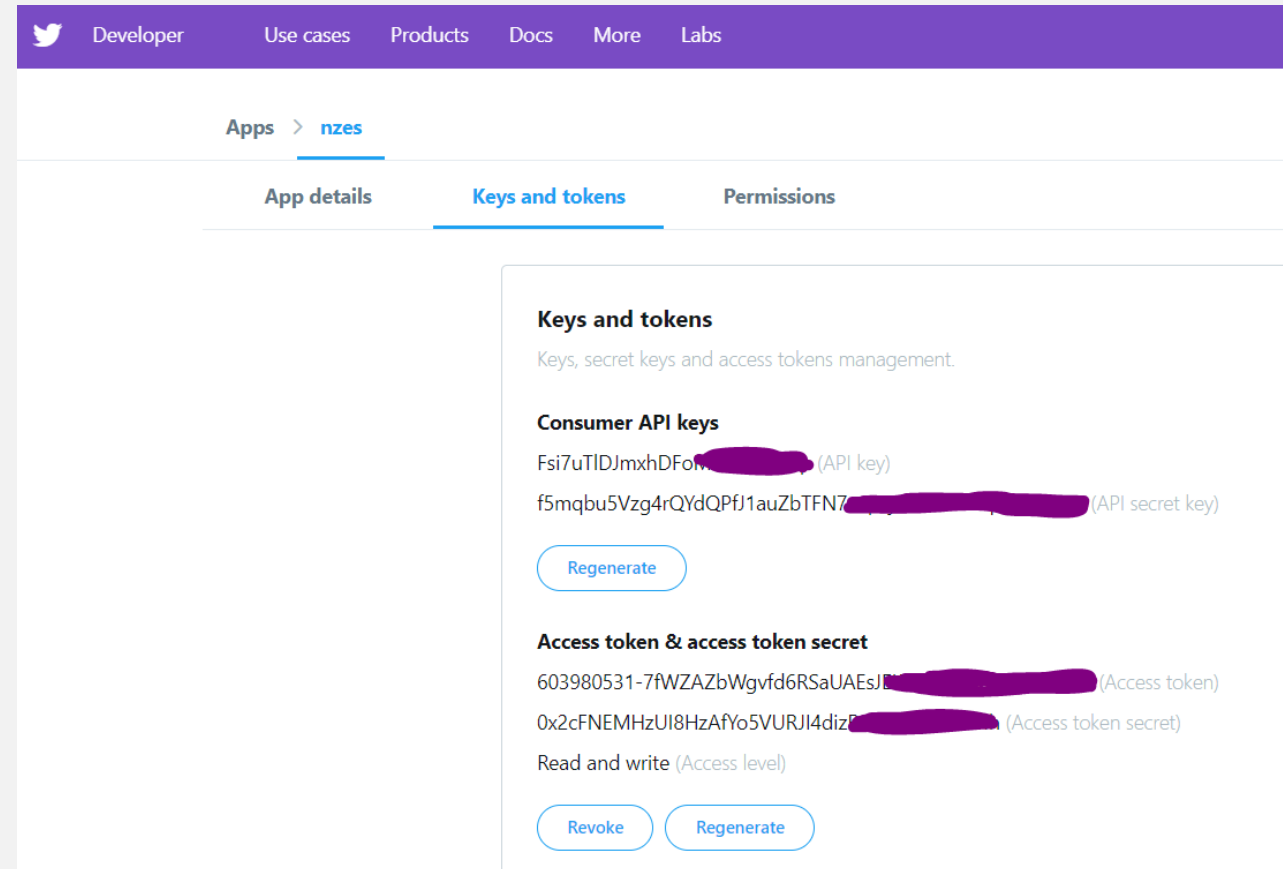


@jennypannell



# STEP I: GET TWITTER CREDENTIALS

- Go to <https://developer.twitter.com/en/apps>
- Create a blank app
- Copy consumer API keys, access token & secret



The screenshot shows the Twitter Developer Portal interface. At the top, there is a navigation bar with the Twitter logo and links for 'Developer', 'Use cases', 'Products', 'Docs', 'More', and 'Labs'. Below this, the breadcrumb 'Apps > nzes' is visible. The main content area has three tabs: 'App details', 'Keys and tokens' (which is selected), and 'Permissions'. Under the 'Keys and tokens' tab, there is a section titled 'Keys and tokens' with the subtitle 'Keys, secret keys and access tokens management.' Below this, there are two sections: 'Consumer API keys' and 'Access token & access token secret'. The 'Consumer API keys' section shows an 'API key' (Fsi7uTIDJmxhDFo...) and an 'API secret key' (f5mqbu5Vzg4rQYdQPfJ1auZbTFN7...), both of which are partially redacted with purple boxes. A 'Regenerate' button is located below these keys. The 'Access token & access token secret' section shows an 'Access token' (603980531-7fWZAZbWgvfd6RSaUAesJ...) and an 'Access token secret' (0x2cFNEMHzUI8HzAfYo5VURJI4diz...), also partially redacted. Below these, the access level is listed as 'Read and write'. 'Revoke' and 'Regenerate' buttons are located at the bottom of this section.

PACKAGES REQUIRED: twitteR

## STEP 2: SCRAPE THE HASHTAG

Provide credentials and connect to twitter

Search parameters, and max no. tweets to return

```
25
26 ## Get the data
27 {r connect to twitter}
28
29 consumer_key <- 'Fsi7uTlDJmxhDFoMamPRIPfHp'
30 consumer_secret <- 'f5mqbu5Vzg4rQYdQPfJ1auZbTFN7vepEjL8DoxF9nqKwh2vK75'
31 access_token <- '603980531-7fwAZzbwgvfd6RSauAesJEi00XkwPoRys2CdFVkm'
32 access_token_secret <- '0x2cFNEMHZUI8HzAfYo5VURJI4dizRsQBijlFjzfhzoHh'
33
34 setup_twitter_oauth(consumer_key, consumer_secret, access_token, access_token_secret)
35
36
37 [1] "Using direct authentication"
38
39 {r download tweets}
40 kind_twitter <- searchTwitter("#kindnessInScience",since='2017-01-01',n=5000)
41 kind_twitter_df <- twListToDF(kind_twitter) # Convert to data frame
42
43 # only returns tweets from last 7 days due to the API the package uses
44
45 rm(kind_twitter, kind_twitter_df, consumer_key, consumer_secret, access_token, access_token_secret)
46
47 5000 tweets were requested but the API can only return 20
```



PACKAGES REQUIRED: twitteR

## STEP 2: SCRAPE THE HASHTAG

Provide credentials and connect to twitter

Search parameters, and max no. tweets to return

- Issue: twitter only allows searching back as far as 7 days
- Python alternative using screen scraping: <https://github.com/taspinar/twitterscraper>

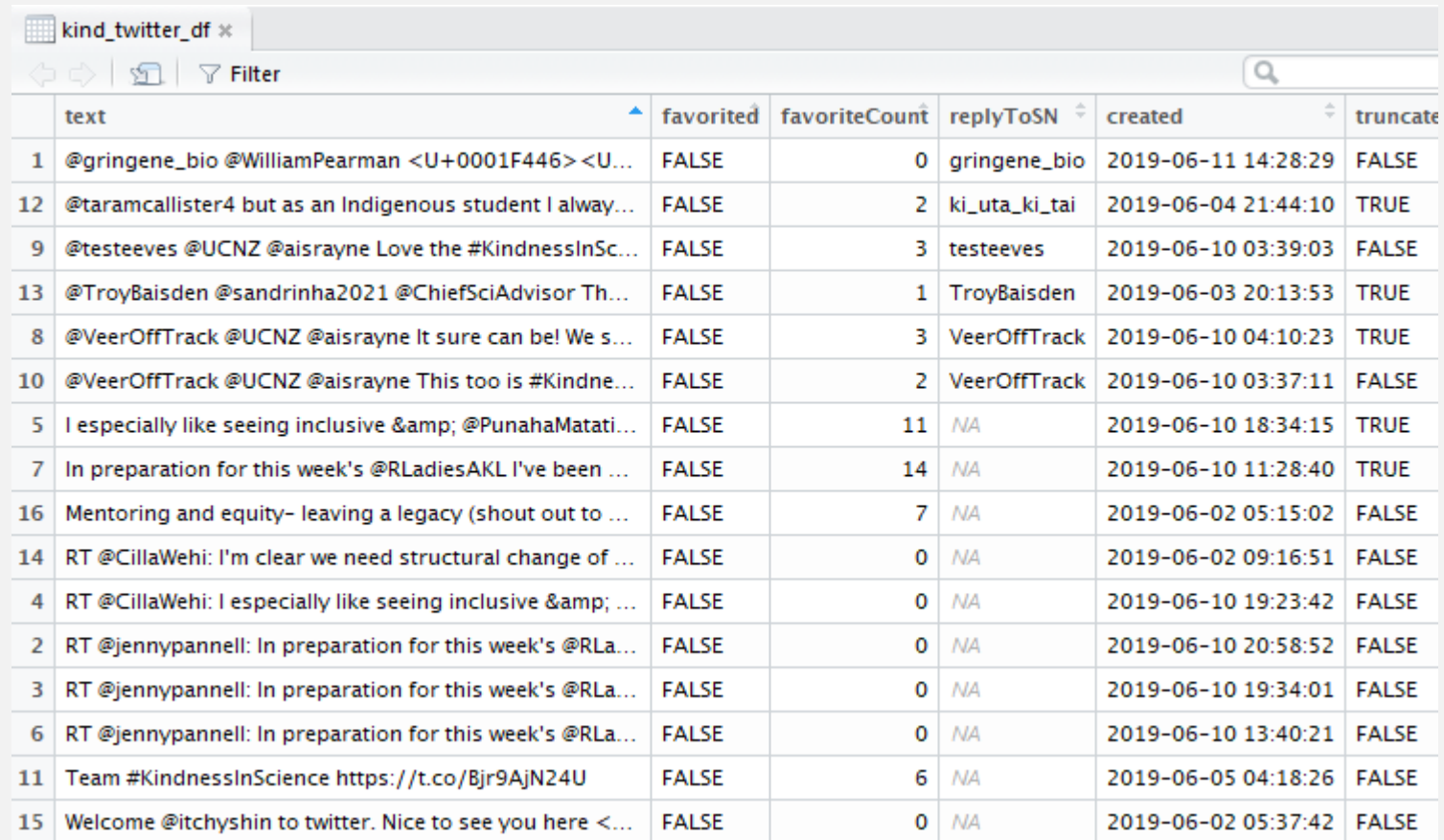
```
25
26 ## Get the data
27 {r connect to twitter}
28
29 consumer_key <- 'Fsi7uTlDjmxhDFoMamPRIPfHp'
30 consumer_secret <- 'f5mqbu5Vzg4rQYdQPfJ1auZbTFN7vepEjL8DoxF9nqKwh2vK75'
31 access_token <- '603980531-7fwAZAbwgVfd6RSauAesJEi00XkwPoRys2CdFVkm'
32 access_token_secret <- '0x2cFNEMHZUI8HzAfYo5VURJI4dizRsQBijlFjzfhzoHh'
33
34 setup_twitter_oauth(consumer_key, consumer_secret, access_token, access_token_secret)
35
36
37 [1] "Using direct authentication"
38
39 {r download tweets}
40 kind_twitter <- searchTwitter("#kindnessInScience",since='2017-01-01',n=5000)
41 kind_twitter_df <- twListToDF(kind_twitter) # Convert to data frame
42
43 # only returns tweets from last 7 days due to the API the package uses
44
45 rm(kind_twitter, kind_twitter_df, consumer_key, consumer_secret, access_token, access_token_secret)
46
47
48 5000 tweets were requested but the API can only return 20
```



@jennypannell

## PACKAGES REQUIRED: twitteR

# STEP 2: SCRAPE THE HASHTAG



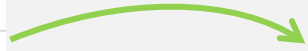
	text	favorited	favoriteCount	replyToSN	created	truncate
1	@gringene_bio @WilliamPearman <U+0001F446><U...	FALSE	0	gringene_bio	2019-06-11 14:28:29	FALSE
12	@taramcallister4 but as an Indigenous student I alway...	FALSE	2	ki_uta_ki_tai	2019-06-04 21:44:10	TRUE
9	@testeeves @UCNZ @aisrayne Love the #KindnessInSc...	FALSE	3	testeeves	2019-06-10 03:39:03	FALSE
13	@TroyBaisden @sandrinha2021 @ChiefSciAdvisor Th...	FALSE	1	TroyBaisden	2019-06-03 20:13:53	TRUE
8	@VeerOffTrack @UCNZ @aisrayne It sure can be! We s...	FALSE	3	VeerOffTrack	2019-06-10 04:10:23	TRUE
10	@VeerOffTrack @UCNZ @aisrayne This too is #Kindne...	FALSE	2	VeerOffTrack	2019-06-10 03:37:11	FALSE
5	I especially like seeing inclusive & @PunahaMatati...	FALSE	11	NA	2019-06-10 18:34:15	TRUE
7	In preparation for this week's @RLadiesAKL I've been ...	FALSE	14	NA	2019-06-10 11:28:40	TRUE
16	Mentoring and equity- leaving a legacy (shout out to ...	FALSE	7	NA	2019-06-02 05:15:02	FALSE
14	RT @CillaWehi: I'm clear we need structural change of ...	FALSE	0	NA	2019-06-02 09:16:51	FALSE
4	RT @CillaWehi: I especially like seeing inclusive & ...	FALSE	0	NA	2019-06-10 19:23:42	FALSE
2	RT @jennypannell: In preparation for this week's @RLa...	FALSE	0	NA	2019-06-10 20:58:52	FALSE
3	RT @jennypannell: In preparation for this week's @RLa...	FALSE	0	NA	2019-06-10 19:34:01	FALSE
6	RT @jennypannell: In preparation for this week's @RLa...	FALSE	0	NA	2019-06-10 13:40:21	FALSE
11	Team #KindnessInScience <a href="https://t.co/Bjr9AjN24U">https://t.co/Bjr9AjN24U</a>	FALSE	6	NA	2019-06-05 04:18:26	FALSE
15	Welcome @itchyshin to twitter. Nice to see you here <...	FALSE	0	NA	2019-06-02 05:37:42	FALSE



PACKAGES REQUIRED: tm, tidytext, stringr

## STEP 3: CLEAN UP TWEET DATA

	text	favorited	favoriteCount	replyToSN	created	truncate
1	@gringene_bio @WilliamPearman <U+0001F446> <U...	FALSE	0	gringene_bio	2019-06-11 14:28:29	FALSE
12	@taramcallister4 but as an Indigenous student I alway...	FALSE	2	ki_uta_ki_tai	2019-06-04 21:44:10	TRUE
9	@testeeves @UCNZ @aisrayne Love the #KindnessInSc...	FALSE	3	testeeves	2019-06-10 03:39:03	FALSE
13	@TroyBaisden @sandrinha2021 @ChiefSciAdvisor Th...	FALSE	1	TroyBaisden	2019-06-03 20:13:53	TRUE
8	@VeerOffTrack @UCNZ @aisrayne It sure can be! We s...	FALSE	3	VeerOffTrack	2019-06-10 04:10:23	TRUE
10	@VeerOffTrack @UCNZ @aisrayne This too is #Kindne...	FALSE	2	VeerOffTrack	2019-06-10 03:37:11	FALSE
5	I especially like seeing inclusive & @PunahaMatati...	FALSE	11	NA	2019-06-10 18:34:15	TRUE
7	In preparation for this week's @RLadiesAKL I've been ...	FALSE	14	NA	2019-06-10 11:28:40	TRUE
16	Mentoring and equity- leaving a legacy (shout out to ...	FALSE	7	NA	2019-06-02 05:15:02	FALSE
14	RT @CillaWehi: I'm clear we need structural change of ...	FALSE	0	NA	2019-06-02 09:16:51	FALSE
4	RT @CillaWehi: I especially like seeing inclusive & ...	FALSE	0	NA	2019-06-10 19:23:42	FALSE
2	RT @jennypannell: In preparation for this week's @RLa...	FALSE	0	NA	2019-06-10 20:58:52	FALSE
3	RT @jennypannell: In preparation for this week's @RLa...	FALSE	0	NA	2019-06-10 19:34:01	FALSE
6	RT @jennypannell: In preparation for this week's @RLa...	FALSE	0	NA	2019-06-10 13:40:21	FALSE
11	Team #KindnessInScience <a href="https://t.co/Bjr9AjN24U">https://t.co/Bjr9AjN24U</a>	FALSE	6	NA	2019-06-05 04:18:26	FALSE
15	Welcome @itchyshin to twitter. Nice to see you here <...	FALSE	0	NA	2019-06-02 05:37:42	FALSE



	word	n
1	kindnessinscience	242
2	status	97
3	twittercom	97
4	https	93
5	science	45
6	testeeves	23
7	pictwittercom	20
8	conserteam	16
9	kindness	15
10	sgalla	15
11	http	12
12	love	12
13	change	11
14	inclusive	11
15	research	11
16	cillawehi	10
17	diversity	9



PACKAGES REQUIRED: tm, tidytext, stringr

## STEP 3: CLEAN UP TWEET DATA

Break up tweets into vector of single words

```
3 ## Tidy up the text
4 Tweet text needs to be tidied by splitting into words, removing punctuation and numbers
  and words like "the", "and", etc.
5 ## {r} triangle data and tidy
6 # make sure the tweets are stored as character data
7 kind$text <- as.character(kind$text)
8
9 # before we split up, make sure any words that need to stay together are replaced
10 kind$text <- gsub("te reo", "tereo", kind$text)
11 kind$text <- gsub("ka rawe", "karawe", kind$text)
12 kind$text <- gsub("kia ora", "kiaora", kind$text)
13 kind$text <- gsub("kia kaha", "kiakaha", kind$text)
14
15 use the unnest_tokens function in tidytext to make the word frequency table
16 # make into data frame with 1 row per word
17 kindTable <- kind %>%
18   unnest_tokens(word, text)
```

Remove unwanted words, punctuation, numbers, etc

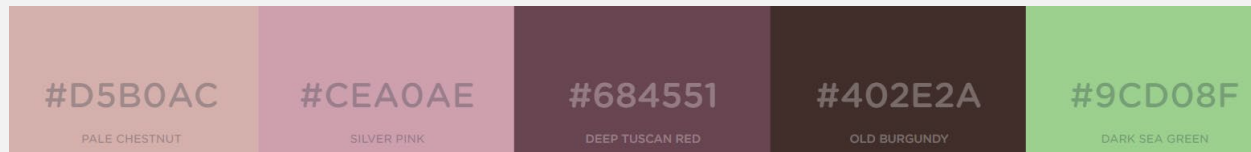
```
16 # remove numbers and punctuation
17 kindTable$word <- removeNumbers(kindTable$word) # remove numbers & punctuation
18 kindTable$word <- removePunctuation(kindTable$word)
19 kindTable <- kindTable[!(is.na(kindTable$word) | kindTable$word==""), ] # remove
  blank/whitespace words
20
21 # remove stop words - aka typically very common words such as "the", "of" etc
22 # now 7770 words
23 data(stop_words)
24
25 kindTable <- kindTable %>%
26   anti_join(stop_words)
27
28 rm(stop_words)
29
30 # add word count to table
31 kindTable <- kindTable %>%
32   dplyr::count(word, sort=TRUE)
```

Count how often each word occurs



PACKAGES REQUIRED: wordcloud2 (OPTIONAL: htmlwidgets, webshot)

## STEP 4: MAKE WORD CLOUD!

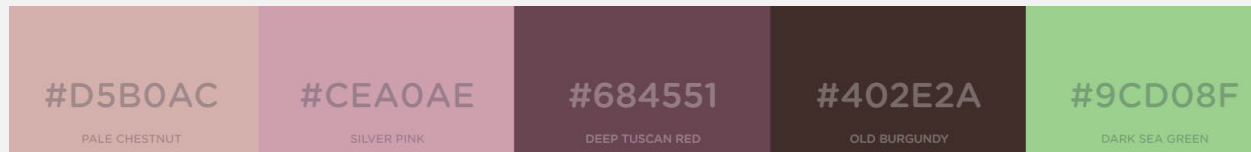


```
23
24 # Subset the data for the word cloud to a manageable size
25 ```{r subset for word cloud}
26 # Too many for word cloud to cope with, choose only words that occur more than three
  times
27 clouddat<-subset(kindTable, n>2)
28 ...
29
30 #Create the word cloud
31 Use a custom palette or built in, I used online palette generators
32 ```{r create the word cloud}
33 # Create Palette
34 kindPalette <- c("#D5B0AC", "#CEA0EA", "#684551", "#402E2A", "#9CD08F")
35
36 mycloud<-wordcloud2(clouddat, color=rep_len(kindPalette, nrow(clouddat)), size=0.5)
37 mycloud
38 # save it in html
39 savewidget(mycld,"kis.html",selfcontained = F)
40 # save as a pdf
41 webshot::install_phantomjs()
42 webshot("kis.html","kis.pdf", delay =20, vwidth = 1500, vheight=1000)
43 ...
44
```



PACKAGES REQUIRED: wordcloud2 (OPTIONAL: htmlwidgets, webshot)

## STEP 4: MAKE WORD CLOUD!



- I used a palette from <https://colors.co/app> but you can use built-in colours
- Issues: wordcloud2 is glitchy - difficult to customise cloud and save output.

```
23
24 # Subset the data for the word cloud to a manageable size
25 ```{r subset for word cloud}
26 # Too many for word cloud to cope with, choose only words that occur more than three
  times
27 clouddat<-subset(kindTable, n>2)
28 ...
29
30 #Create the word cloud
31 Use a custom palette or built in, I used online palette generators
32 ```{r create the word cloud}
33 # Create Palette
34 kindPalette <- c("#D5B0AC", "#CEA0EA", "#684551", "#402E2A", "#9CD08F")
35
36 mycloud<-wordcloud2(clouddat, color=rep_len(kindPalette, nrow(clouddat)), size=0.5)
37 mycloud
38 # save it in html
39 savewidget(mycLOUD,"kis.html",selfcontained = F)
40 # save as a pdf
41 webshot::install_phantomjs()
42 webshot("kis.html","kis.pdf", delay =20, vwidth = 1500, vheight=1000)
43
44 ...
45
```

