Building a package that fits into an evolving ecosystem

OCRUG

Emil Hvitfeldt

2019-11-19

Overview

- Open Source Software Development

- My journey

- What I learned

This talk is based on anecdotes, I dearly hope that they generalize!

As a new developer it can be hard to find problems that are:

- Easy enough for you do
- Prominent enough that they are worth solving

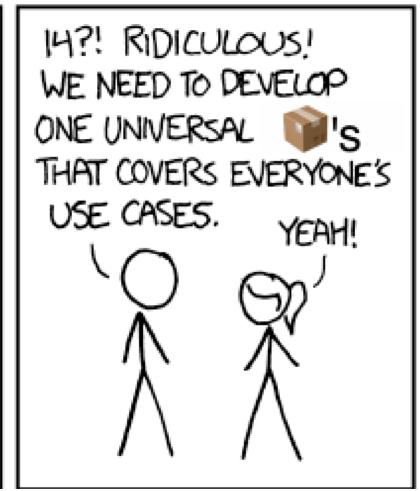
My advice

- Keep list of ideas
- Google early and often about implementations

Working on a implentation can still be fruitful even if it doesn't end up on CRAN.

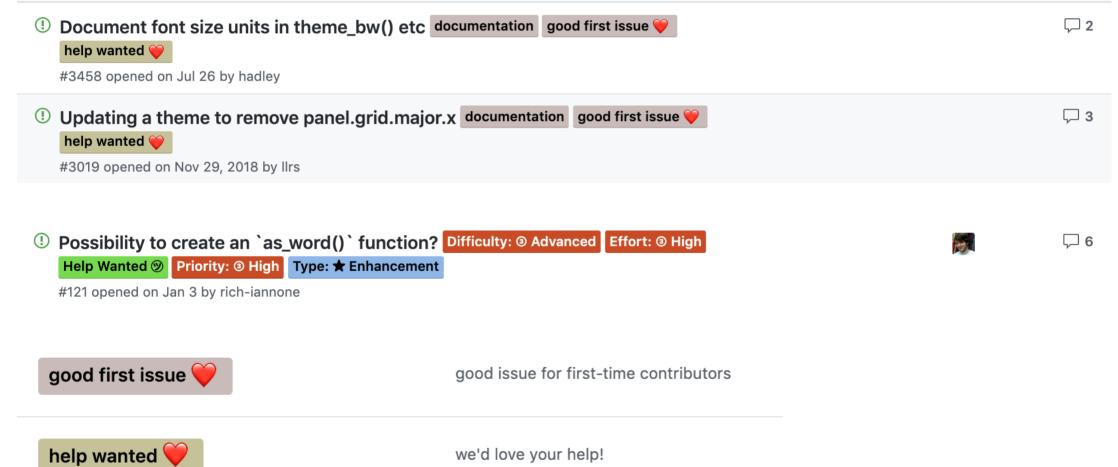
HOW PROLIFERATE: (SEE: A/C CHARGERS, CHARACTER ENCODINGS, INSTANT MESSAGING, ETC.)

SITUATION:
THERE ARE
14 COMPETING
III'S

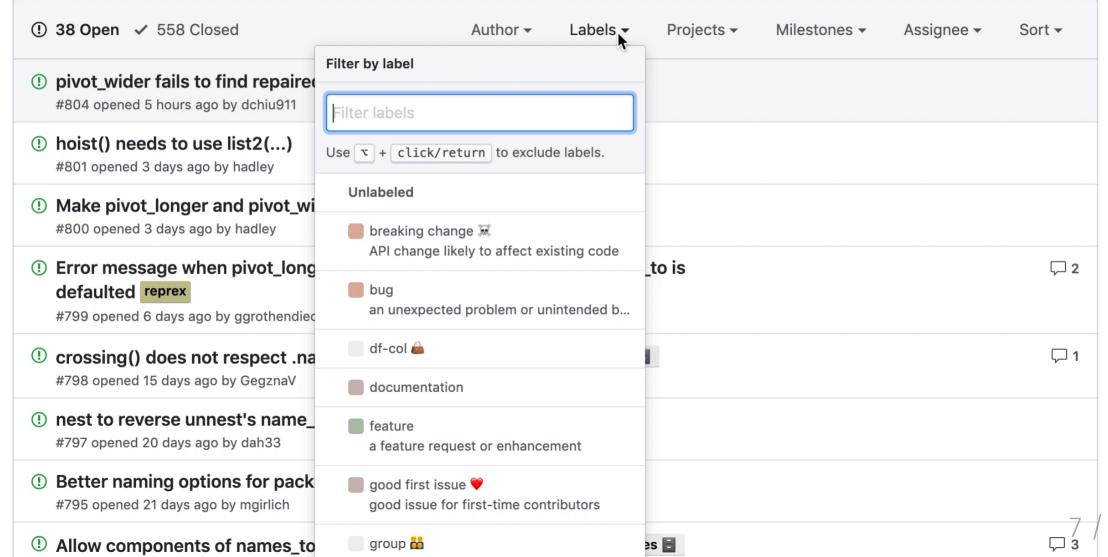




Look for Github tags



New issue



† Milestones 0

Ask before starting the work

Makes sure that:

- you are not doing the same work as someone else
- the work is wanted

Github Projects

Text Processing Recipes Closed

Updated on Dec 18, 2018

currently in development with the textrecipes package.

· seriousness: medium

· complexity: low

Add steps for stop words, stemming, tf-idf, n-grams and so on.



text recipe steps #192



topepo opened this issue on Sep 5, 2018 · 14 comments



topepo commented on Sep 5, 2018

Member



•••

For this project, the idea is to have steps that can be used to process text data (contained in a new package). I've made placeholders in that project for some obvious processing candidates.

@EmilHvitfeldt has volunteered to get started. Perhaps @juliasilge, @skeydan, and others might have some suggestions and opinions. I'd be happy to include tensorflow methods for text processing too.

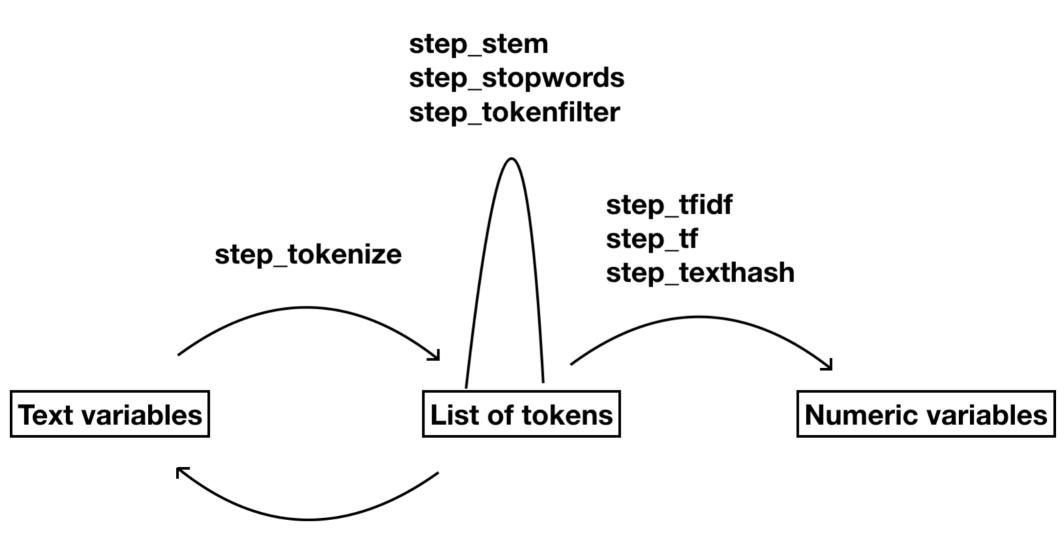
I can add anyone interested in helping to the project page. We can use this issue to discuss ideas and kick around implementation questions.

Challenges

- Inheret messy structure of text
- steps don't have specified order

Design choice

Flexibily > speed



step_untokenize

```
step stem <-
 function(recipe,
          role = NA,
          trained = FALSE,
          columns = NULL,
          options = list(),
          custom stemmer = NULL,
          skip = FALSE,
          id = rand id("stem")
 ) {
   add step(
     recipe,
      step_stem_new(
       terms = ellipse_check(...),
       role = role,
       trained = trained,
        options = options,
        custom stemmer = custom stemmer,
        columns = columns,
        skip = skip,
        id = id
step_stem_new <-
  function(terms, role, trained, columns, options, custom_stemmer, skip, id) {
     subclass = "stem",
     terms = terms,
     role = role,
     trained = trained,
     columns = columns,
      options = options,
     custom_stemmer = custom_stemmer,
      skip = skip,
      id = id
prep.step_stem <- function(x, training, info = NULL, ...) {</pre>
 col_names <- terms_select(x$terms, info = info)</pre>
  check_list(training[, col_names])
  step_stem_new(
   terms = x$terms,
   role = x$role,
   trained = TRUE,
   columns = col_names,
   options = x$options,
   custom_stemmer = x$custom_stemmer,
   skip = x$skip,
    id = x$id
```

```
bake.step stem <- function(object, new data, ...) {</pre>
 col names <- object$columns
 stem_fun <- object$custom_stemmer %||%</pre>
   SnowballC::wordStem
  for (i in seq along(col names)) {
   stemmed text <- map(new data[, col names[i], drop = TRUE],
                        stem fun)
   new data[, col names[i]] <- tibble(stemmed text)</pre>
 new data <- factor to text(new data, col names)</pre>
 as tibble(new data)
print.step stem <-
 function(x, width = max(20, options()) width - 30), ...) {
   cat("Stemming for ", sep = "")
   printer(x$columns, x$terms, x$trained, width = width)
   invisible(x)
tidy.step_stem <- function(x, ...) {</pre>
 if (is trained(x)) {
   res <- tibble(terms = x$terms,
                  is custom stemmer = is.null(x$custom stemmer))
 } else {
   term names <- sel2char(x$terms)</pre>
   res <- tibble(terms = term names,
                  value = na chr)
 res$id <- x$id
 res
```

```
step stem <-
 function(recipe,
          role = NA,
          trained = FALSE,
          columns = NULL,
          options = list(),
          custom stemmer = NULL,
          skip = FALSE,
          id = rand id("stem")
 ) {
   add step(
     recipe,
      step_stem_new(
       terms = ellipse_check(...),
       role = role,
       trained = trained,
        options = options,
        custom stemmer = custom stemmer,
        columns = columns,
        skip = skip,
        id = id
step_stem_new <-
  function(terms, role, trained, columns, options, custom_stemmer, skip, id) {
     subclass = "stem",
     terms = terms,
     role = role,
     trained = trained,
     columns = columns,
      options = options,
     custom stemmer = custom stemmer,
      skip = skip,
      id = id
prep.step_stem <- function(x, training, info = NULL, ...) {</pre>
 col_names <- terms_select(x$terms, info = info)</pre>
  check_list(training[, col_names])
  step_stem_new(
   terms = x$terms,
   role = x$role,
   trained = TRUE,
   columns = col_names,
   options = x$options,
   custom_stemmer = x$custom_stemmer,
   skip = x$skip,
    id = x$id
```

```
bake.step stem <- function(object, new data, ...) {</pre>
 col names <- object$columns
  stem fun <- object$custom stemmer %||%
   SnowballC::wordStem
  for (i in seq along(col names)) {
   stemmed text <- map(new data[, col names[i], drop = TRUE],
                        stem fun)
   new data[, col names[i]] <- tibble(stemmed text)</pre>
 new data <- factor to text(new data, col names)</pre>
 as tibble(new data)
print.step stem <-</pre>
 function(x, width = max(20, options()) width - 30), ...) {
   cat("Stemming for ", sep = "")
   printer(x$columns, x$terms, x$trained, width = width)
   invisible(x)
tidy.step stem <- function(x, ...) {
 if (is trained(x)) {
   res <- tibble(terms = x$terms,
                  is custom stemmer = is.null(x$custom stemmer))
 } else {
   term names <- sel2char(x$terms)</pre>
   res <- tibble(terms = term names,
                  value = na chr)
 res$id <- x$id
 res
```

The bake step

the traceback

```
> library(textrecipes)
> library(recipes)
> xxx <- recipe(~ essay0, data = okc_text) %>%
   step_tokenize(essay0) %>%
   step_stem(essay0) %>%
   prep() %>%
+ juice()
Error in bake.step_stem(x$steps[[i]], new_data = trainina) : test
> traceback()
12: stop("test") at stem.R#151
11: bake.step_stem(x$steps[[i]], new_data = training)
10: bake(x$steps[[i]], new_data = training)
9: prep.recipe(.)
8: prep(.)
7: function_list[[i]](value)
6: freduce(value, `_function_list`)
5: `_fseq`(`_lhs`)
4: eval(quote(`_fseq`(`_lhs`)), env, env)
3: eval(quote(`_fseq`(`_lhs`)), env, env)
2: withVisible(eval(quote(`_fseq`(`_lhs`)), env, env))
1: recipe(~essay0, data = okc_text) %>% step_tokenize(essay0) %>%
       step_stem(essay0) %>% prep() %>% juice()
```

```
I call prep()
prep() calls prep.recipe()
prep.recipe() calls bake() in a loop
bake() calls bake.step_stem()
```

quite a few levels deep.

Browser to the rescue

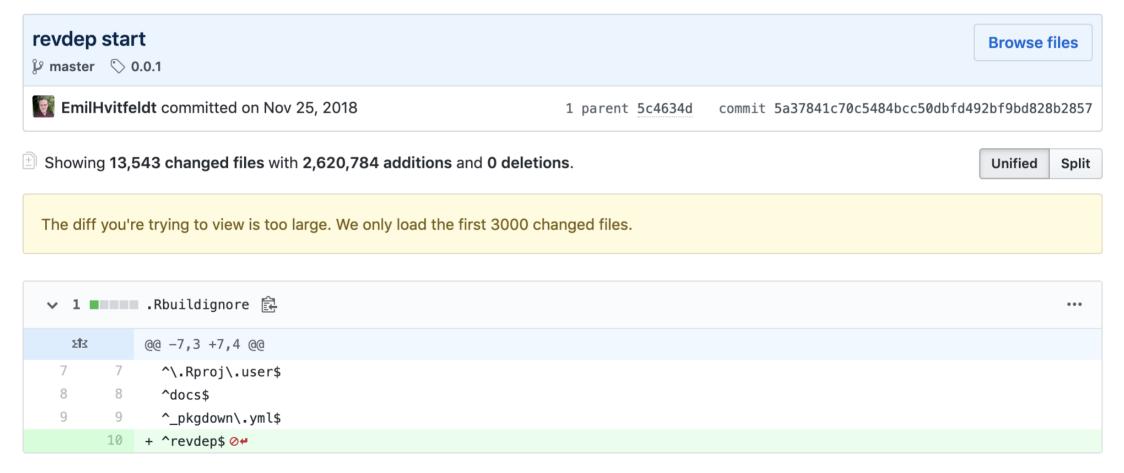
Plenty of follow up

Don't reinvent the wheel

textrecipes stands on the shoulders of

- recipe (obviously)
- tokenizers
- SnowballC
- stopwords
- text2vec
- textfeatures

My biggest git mistake



The reward - joined my first organization



Joined the tidymodels organization

on Nov 18



tidymodels

The reward - CRAN releases

The reward - Post on tidyverse.org

Thank you!

- 2 EmilHvitfeldt
- @ @ Emil_Hvitfeldt
 - 2 emilhvitfeldt
- ? www.hvitfeldt.me

Slides created via the R package xaringan.