

PRÁCTICA 5: TF-IDF

RECUPERACIÓN DE INFORMACIÓN
LICENCIATURA EN CIENCIAS DE LA COMPUTACIÓN
PRIMAVERA 2022

GUSTAVO IVÁN MOLINA REBOLLEDO
Benemérita Universidad Autónoma de Puebla

ABSTRACT. A Haskell tool for CCOS 264 (Information retrieval) class. The full code, input and output is included in the project repository: https://github.com/ivanmoreau/tf_idf

2 de marzo de 2022

I. INTRODUCCIÓN

En esta práctica el objetivo es obtener consultas a partir de la información conseguida en la práctica anterior. El programa debe

1. Programar el pesado de términos usando TF/IDF utilizando el corpus de la práctica de pre-procesamiento.
2. Guardar el pesado obtenido.
3. Crear 5 consultas y mostrar el pesado.

II. IMPLEMENTACIÓN

Para implementar esta herramienta utilizo el lenguaje de programación Haskell, así como, en especial, la librería Parsec.

Este programa consta de tres partes (o archivos) en los que se dividen los distintos funcionamientos. El sistema de búsqueda y de matriz se encuentran en TFIDF.hs.

Las demás partes son reutilizadas de las prácticas anteriores como librerías (sistema de Query, por ejemplo).

Main.hs sólo se dedica a manejar el IO, e interconectar todo el funcionamiento.

III. TF/IDF

Esta sección es Literate Haskell. El mismo código es parte del ejecutable y parte del documento.

Licencia:

```
1 {- |
2 Copyright: (c) 2022 Iván Molina Rebolledo
3
4 SPDX-License-Identifier: GPL-3.0-only
5 Maintainer: Iván Molina Rebolledo <ivanmolinarebolledo@gmail.com>
6
7 See README for more info
8 -}
```

Boilerplate:

```

1 {-# LANGUAGE OverloadedStrings #-}
2
3 module TFIDF ( genMatrix, search ) where
4
5 import           Data.List (delete, nub, elemIndices)
6 import qualified Data.Map as Dm
7 import           Data.Map (Map, fromList, insert, (!))
8 import           Data.Serialize.Text ()
9 import           Data.Text (Text, pack, replace, splitOn)
10 import           Query
11
12 type MatrixW = Map Text (Map Int Double)

```

Lista de palabras a ignorar:

```

1 rone :: [Text]
2 rone = ["$user tweeted ", " $ht", "$ht ", " $user", "$user "]

```

Es necesario limpiar todas las líneas del corpus (eliminar una lista de palabras)

```

1 fixWords :: Text → [Text] → Text
2 fixWords xs [] = xs
3 fixWords xs (l : ls) = fixWords (replace l (pack "") xs) ls

```

con el fin de obtener una matriz limpia con las palabras.

```

1 getWords :: Text → MatrixW
2 getWords l =
3   fromList
4     ( map
5       (\x → (x, fromList [(0, 0)]))
6       (Data.List.delete "" (nub (splitOn " " (replace "\n" " " l))))
7     )

```

Se ha de inicializar la matriz con el número de documentos:

```

1 emptymatrix :: MatrixW → Int → MatrixW
2 emptymatrix m c = Dm.map (\_ → fromList (map (\x → (x, 0)) [0 .. c]))
   ↪ m

```

En este código se establece la existencia de una palabra para el documento n ; establece o actualiza la frecuencia.

```

1 aa :: MatrixW → Text → Int → Double → MatrixW
2 aa m t i val = let r = m ! ( t ) in insert t (insert i val r) m

```

Se define la variante del peso TF, idf y el tf-idf:

```

1 tf :: Text → Text → Double
2 tf i t = let o = Prelude.length $ elemIndices i $ splitOn " " t in
3   if o > 0 then
4     1.0 + logBase 2.0 (fromIntegral o)
5   else
6     0
7
8 findcol :: Text → Text → Bool
9 findcol term text = (>0) $ Prelude.length $ filter (==term) $ splitOn "
   ↪ " text
10 findtex :: Text → [Text] → Int
11 findtex term textl = Prelude.length $ filter (==True) $ map (findcol
   ↪ term) textl
12
13 idf :: Text -- Termino.
14   → [Text] -- [Lineas].
15   → Double
16 idf t m = let n = findtex t m in
17   logBase 2.0 $ (fromIntegral $ Prelude.length m) / (fromIntegral n)
18
19 tfidf :: Text → Text → [Text] → Double
20 tfidf term line lines_ = (tf term line) * (idf term lines_)

```

Se llena la matriz por linea:

```

1 byLine :: [Text] → [Text] → MatrixW → Int → MatrixW
2 byLine [] _ m _ = m
3 byLine (l : ls) lines_ m c =
4   let words_ = filter (≠ "") $ splitOn " " l
5       w_ = map (\w → (w, tfidf w l lines_)) words_ in
6   byLine ls lines_ (f m (w_)) (c + 1)
7   where f m_ [] = m_
8         f m_ ((w, w_):xs) = f (aa m_ w c w_) xs

```

Se genera la matriz:

```

1 genMatrix :: Text → MatrixW
2 genMatrix t =
3   let ff = fixWords t rone
4       gw = getWords ff
5       ts = (splitOn "\n" ff)
6       em = emptymatrix gw (length ts)
7   in byLine ts ts em 0

```

Busca la palabra y regresa los índices de los documentos:

```

1 readd :: Text → MatrixW → [(Int, Double)]
2 readd t m = case Dm.lookup t m of
3   Just v → (Dm.toList (Dm.filter (>0) v))
4   Nothing → []

```

Se evalua el árbol de Query:

```

1 lmin :: [(Int, Double)] → [(Int, Double)] → [(Int, Double)]
2 lmin _ [] = []
3 lmin l0 (l1:ls1) = lmin (filter (\(x, _) → x ≠ (fst l1)) l0) ls1
4
5 lunion :: [(Int, Double)] → [(Int, Double)] → [(Int, Double)]
6 lunion l0 l1 = let n = filter (\(x, _) → x `notElem` (fst $ unzip l0))
7   → l1 in
8   l0 ++ n

```

```

9  lintersect :: [(Int, Double)] → [(Int, Double)] → [(Int, Double)]
10 lintersect l0 l1 = let n = filter (\(x,_) → x `elem` (fst $ unzip l0))
    → l1 in
11     n
12
13 pw :: Query → Int → MatrixW → [(Int, Double)]
14 pw (QAnd l r) t m = (pw l t m) `liintersect` (pw r t m)
15 pw (QOr l r) t m = (pw l t m) `lunion` (pw r t m)
16 pw (QNot e) t m = (zip [0..t] (replicate t 0)) `lmin` (pw e t m)
17 pw (QP v) _ m = read v m

```

Se realiza una consulta con una Query:

```

1 search :: Query → MatrixW → [(Int, Double)]
2 search q m = pw q (Dm.size (snd (Dm.elemAt 0 m))) m

```

IV. FUNCIONAMIENTO

Podemos usarlo de la siguiente forma:

```

ivanmolinarebolledo@Ivans-macOS tf_idf % stack run -- --help
The tfidf program

```

```
tfidf [COMMAND] ... [OPTIONS]
```

Common flags:

```

-? --help          Display help message
-V --version       Print version information

```

```
booleanmodel matrix [OPTIONS]
```

```

-f --from=ITEM
-t --to=ITEM

```

```
booleanmodel query [OPTIONS]
```

```
-o --ogfile=ITEM
-m --matrix=ITEM
-q --query=ITEM
ivanmolinarebolledo@Ivans-macOS tf_idf %
```

La única parte importante que es necesario resaltar es el orden de los parsers en cada una de las funciones "grandes". Este orden es importante porque nos permite distinguir "tweeted" de una palabra cualquiera, por ejemplo. Todo lo demás son combinaciones de parsers.

V. GRAMMAR

```
Query → a
a → o ("and" o)*
o → term ("or" term)*
term → "not" Query | "(" Query ")" | Identifier
Identifier → alphaNum (alphaNum | "'")*
```

VI. RESULTADOS

Se prueban las siguientes queries con resultados satisfactorios:

```
ivanmolinarebolledo@Ivans-macOS tf_idf % stack run -- query
→ -o="Tweets.txt" -m="mm.ths" -q="cama or museo"
[13] @LuisGcortez tweeted: Parece que falta algo de atención y
→ mantenimiento al museo de la ciudad por parte del @ICAdifusion
→ http://t.co/FLbC5SbcoF (W: 8.45532722030456)
[201] @JaAC9510 tweeted: Hoy me niego a levantarme de mi cama (W:
→ 7.455327220304562)
[16] @yosoyloagui tweeted: Tu y yo juntos en un lugar perfectollamado mi
→ cama *.* (W: 7.455327220304562)

ivanmolinarebolledo@Ivans-macOS tf_idf % stack run -- query
→ -o="Tweets.txt" -m="mm.ths" -q="(invito and not carnal) or novio"
```

[186] @Andresrgtz tweeted: Mi amiga se besó a una mujer muy sexi el fin
 → de semana no le digan a su novio (W: 6.870364719583405)
 [49] @KarlyyRG tweeted: Nada mejor que tener un novio super divertido
 → que en lugar de limitarte, te sigue el pedo. (W: 6.870364719583405)
 [42] @KarlyyRG tweeted: Lo que mas me encanta de mi novio es que me hace
 → morir de risa. (W: 6.870364719583405)

ivanmolinarebolledo@Ivans-macOS tf_idf % stack run -- query
 → -o="Tweets.txt" -m="mm.ths" -q="vida and riesgo or dios"
 [206] @alex_dsr01 tweeted: Ni pedo la vida es un riesgo -Fergus (W:
 → 8.45532722030456)
 [345] @CaballoNegroII tweeted: RT @MyRyCaR: Gracias por tu presencia en
 → mi vida ...un regalo de Dios ☺ teee amooo @CaballoNegroII
 → http://t.co/grQRS43NK1 (W: 6.133399125417199)
 [203] @abrilushiZ tweeted: RT @_equiswe: Si Dios me quita la vida antes
 → que a ti... (W: 6.133399125417199)

ivanmolinarebolledo@Ivans-macOS tf_idf % stack run -- query
 → -o="Tweets.txt" -m="mm.ths" -q="tortillas or hambre or feliz"
 [69] @PiernasAlLomo tweeted: @La_Montserrat a las tortillas? (W:
 → 8.45532722030456)
 [27] @IAN_BLACK26 tweeted: Mitad de quincena que comiencen los juegos
 → del hambre. (W: 7.455327220304562)
 [14] @mc_gabucha tweeted: Ahora entiendo por que se juntó la necesidad
 → con el hambre... ¡Que gente tan acomplejada! (W: 7.455327220304562)
 [325] @123anotaz tweeted: RT @CarlaMorrAgs: La mujer es una obra de arte
 → que ilumina los ojos de quien la mira. Feliz Dia de la Mujer
 → @CarlaMorrisonmx ♥ (W: 6.455327220304562)
 [265] @fantasmaluigui tweeted: Siempre recuerda: No tomes decisiones
 → cuando estés enojado, y no hagas promesas cuando estés feliz. (W:
 → 6.455327220304562)
 [236] @nichelopezc tweeted: Hoy ni la pizza me hace feliz carajo!! (W:
 → 6.455327220304562)

[109] @zubyelynl tweeted: Hola feliz domingo tengan tod@s hoy les
→ presento esta capa o poncho que se puede usar tanto en temporada de
→ frio ... <http://t.co/qh1xeinf22> (W: 6.455327220304562)

ivanmolinarebolledo@Ivans-macOS tf_idf % stack run -- query
→ -o="Tweets.txt" -m="mm.ths" -q="amiga and novio"

[186] @Andresrgtz tweeted: Mi amiga se besó a una mujer muy sexi el fin
→ de semana no le digan a su novio (W: 6.870364719583405)

VII. REFERENCIAS

parsec. (2022, February 01). Retrieved from <https://hackage.haskell.org/package/parsec>