

TweetKell

Análisis de los datos privados de una cuenta de Twitter con Haskell

Universidad de Sevilla
Ingeniería Informática Tecnologías Informáticas
Programación Declarativa - Tercer curso

Juan Arteaga Carmona
Herrera, Sevilla, España
`JuanArteaga@andalu30.me`

4 de septiembre de 2019

Índice

1	Introducción	3
2	Expectativas	3
3	Criterios de evaluación	3
4	Elección del tipo de trabajo	4
5	Haskell	4
6	Datos de Twitter	4
7	Implementación	4
7.1	Estructura del código	4
7.2	Aplicación de los conceptos y ajuste a los criterios de evaluación	5
8	Uso de la aplicación	8
8.1	Utilizando el archivo compilado	9
8.2	Utilizando el código fuente	9
9	Posibles mejoras	10
9.1	Analizar más información	10
9.2	Paralelización y computación distribuida	10

Índice de figuras

1	Menú principal de la aplicación	9
2	Pantalla de confirmación de la ruta del archivo a analizar	9
3	Distintos tipos de análisis que se pueden realizar sobre el archivo <i>Tweet.js</i>	10
4	Pantalla de resultados del análisis	10

1. Introducción

A continuación se presenta la documentación realizada para la entrega final de la convocatoria de Septiembre del trabajo de curso de la asignatura Programación Declarativa del grado de Ingeniería Informática - Tecnologías Informáticas de la Universidad de Sevilla.

2. Expectativas

Al comienzo del trabajo se pretendía obtener una versión final que cumpliera con los requisitos expuestos por nuestro profesor, los cuales se pueden observar en el apartado 3.

3. Criterios de evaluación

A continuación se muestran los criterios de evaluación a los que se someterá este trabajo, de acuerdo con la información que se encontraba en el aula virtual de la asignatura a día 2 de septiembre de 2019.

- La fecha de entrega del trabajo es el día 8 de septiembre hasta media noche. La defensa del trabajo tendrá lugar la semana siguiente, del 9 al 13 (se anunciará por este canal las fechas).
- El trabajo debe ser un programa de temática "libre" (dentro de los ofertados en la convocatoria de febrero, o los que se haya hablado con los profesores previamente), totalmente funcional. Para la superación del trabajo, éste debe contener como mínimo:
 - Dos aplicaciones de cada concepto básico de programación funcional visto en cada tema del primer bloque de la asignatura (temas 2,4,5,6,7,8). Es decir: al menos usar 2 funciones básicas de haskell, definir 2 funciones recursivas, 2 funciones por patrones, 2 usos de guardas, 2 usos de listas por comprensión, 2 usos de orden superior, 2 usos de tipados, etc.
 - Creación de un módulo
 - Creación de un tipo de datos nuevo y uso de éste.
 - Uso de al menos dos de los conceptos avanzados y librerías vistos en los temas 3, 9, 11, 12, 13, 14 o 15 (deben ser de temas distintos).
- La entrega debe contener, en un archivo comprimido (zip, máximo 50 mb):
 - El código fuente Haskell (incluyendo un pequeño readme explicando el uso del programa).
 - Un documento (formato libre) explicando la estructura del código y cómo se ha aplicado cada uno de los conceptos.

4. Elección del tipo de trabajo

Se decidió que el trabajo fuese uno completamente diferente al que se propuso en la primera convocatoria [1] ya que, tal y como se explicó en las conclusiones de este, elaborar un juego con Haskell y CodeWord no parecía ser una tarea sencilla ni interesante.

Así pues, este trabajo se centrara en la temática de Data Science y análisis de datos, con especial énfasis en la recolección de los datos a través de la librería Aeson [2] y en la interactividad del programa como se verá en el apartado 7.2.

5. Haskell

Para este proyecto utilizaremos en lenguaje de programación Haskell [3]. Haskell es un lenguaje de programación puramente funcional creado por Paul Hodax en 1990 [4].

6. Datos de Twitter

Al utilizar esta aplicación podremos extraer y consultar los datos de una cuenta de Twitter, para esto es necesario que dispongamos de los datos en un formato que nuestro PC pueda entender.

Para conseguir los datos deberemos de hacer a Twitter una petición para que nos permita descargarlos. Esta petición se hace a través de los ajustes de privacidad en nuestro perfil y, dependiendo del tamaño de nuestra cuenta, en unos minutos se nos permitirá acceder a ellos. [5]

Una vez descargados, tendremos que descomprimirlos y ya podremos utilizar la aplicación.

7. Implementación

7.1. Estructura del código

El código se encuentra estructurado en varios módulos que, además de ser un requisito para la entrega del trabajo, permiten una mayor comodidad y simplicidad a la hora de desarrollar la aplicación. A continuación se describen los módulos:

- `TweeKell.hs`

Archivo principal del programa, contiene una función `main` que llama al menú principal del módulo `MenuInteractivo`.

- `MenuInteractivo.hs`

Módulo encargado de generar los menús interactivos del programa además de recibir las instrucciones del usuario y llamar a la lógica correcta.

- `TiposDatos.hs`

Módulo encargado de definir los tipos de datos que se utilizan para decodificar los archivos JSON y mostrar la información.

- **TiposDatos2.hs**

Igual que TiposDatos, con la única diferencia de que este contiene las definiciones de los tipos de datos que producen alguna colisión en caso de que se encontrasen en TiposDatos.

- **ParsersArchivos.hs**

Módulo encargado de parsear los archivos JSON de Twitter y de crear los objetos de los tipos que se definen en TiposDatos y TiposDatos2.

- **AnalisisSimples.hs**

Módulo encargado de recibir el parseado realizado por ParsersArchivos y analizar y mostrar la información que el usuario ha solicitado para todos los tipos de análisis excepto para los que dependen de Tweets.

- **AnalisisTweets.hs**

Igual que AnalisisSimples pero centrado en los que dependen del tipo Tweet.

7.2. Aplicación de los conceptos y ajuste a los criterios de evaluación

- **Aplicaciones de conceptos básicos de programación funcional.**

- 2 funciones básicas de Haskell

El uso de funciones básicas de Haskell se puede observar por todo el trabajo, por ejemplo, se puede observar en dos de las funciones mas simples y cortas:

```
1 modPhoneFile :: String -> String
2 modPhoneFile string = "{" ++ (drop 50 $ reverse $ drop 4 $ reverse string)
3
4 modAccFile :: String -> String
5 modAccFile string = drop 45 $ reverse $ drop 3 $ reverse string
```

- 2 funciones recursivas:

A continuación podemos ver dos funciones recursivas de este trabajo, una de ellas se encarga de imprimir los textos de una lista de Tweets y otra se encarga de contar el número de apariciones de los distintos idiomas de los tweets.

```

1 imprimeTextosTweets :: [String] -> IO()
2 imprimeTextosTweets [] = putStrLn "" --Base
3 imprimeTextosTweets (x:ls) = do --Recursiva
4     putStrLn "\n->"
5     print x
6     imprimeTextosTweets ls
7
8 cuentaIdiomas :: Foldable t => [t a] -> [Int] -> [Int]
9 cuentaIdiomas [] ac = ac
10 cuentaIdiomas (x:xs) ac = do
11     let newac = ac ++ [(length x)]
12     cuentaIdiomas xs newac

```

- 2 funciones por patrones

Las dos funciones recursivas anteriores están, claramente, definidas por patrones. Sin embargo no son las únicas. Un ejemplo de otra función sería el siguiente:

```

1 tipoAnalisis :: String -> [String]
2 tipoAnalisis "account.js" = ["1: Ver informacion", "2: Abrir en navegador"]
3 tipoAnalisis "profile.js" = ["1: Ver informacion", "8: Ver icono de perfil"]
4 tipoAnalisis "verified.js" = ["1: Ver informacion"]
5 tipoAnalisis "phone-number.js" = ["1: Ver informacion"]
6 tipoAnalisis "tweet.js" = ["3: Ver tweets", "4: Ver Retweets", "5: Tweet mas retweeteado",
7                             "6: Tweet con más MG", "7: Idiomas mas utilizados"]

```

- 2 usos de guardas

En cuanto a funciones con guardas, dos ejemplos podrían ser:

```

1 verInformacion :: String -> String -> IO()
2 verInformacion ruta archivo
3     | archivo == "account.js" = verInformacionAccount (ruta ++ "/" ++ archivo)
4     | archivo == "profile.js" = verInformacionProfile (ruta ++ "/" ++ archivo)
5     | archivo == "verified.js" = verInformacionVerified (ruta ++ "/" ++ archivo)
6     | archivo == "phone-number.js" = verInformacionPhone (ruta ++ "/" ++ archivo)
7
8 abrirNavegador :: String -> String -> IO()
9 abrirNavegador ruta archivo
10    | archivo == "account.js" = abrirNavegadorAccount (ruta ++ "/" ++ archivo)
11    | archivo == "profile.js" = abrirNavegadorAccount (ruta ++ "/" ++ archivo)

```

- 2 usos de tipados

En cuanto al uso de tipados, debido a que Haskell es un lenguaje fuertemente tipado se deben de usar sin excepción. En las funciones de ejemplo anteriores se puede ver como, además, se escribe el tipo de la función para que el lenguaje siquiera tenga que inferirlo.

■ Creación de al menos un módulo

Un ejemplo de la creación de un módulo en este trabajo sería el módulo `TiposDatos`, que se encarga de definir parte de los distintos tipos de datos que se usan en el trabajo.

■ Creación de un tipo de datos y uso de este

Utilizando el mismo archivo que en el apartado anterior, podemos ver varias definiciones de tipos de datos, por ejemplo uno de los mas interesantes sería el tipo de datos `Tweet`:

```
1 data Tweet = Tweet {
2     favorite_count :: String,
3     retweet_count  :: String,
4     created_at     :: String,
5     full_text      :: String,
6     lang           :: String
7 } deriving (Generic, Show)
8 instance FromJSON Tweet
```

En este tipo de datos recogemos 5 de entre todos los parámetros que nos encontramos en los archivos de Twitter.

■ Uso de dos conceptos avanzados

En cuanto a uso de conceptos avanzados, como este trabajo se centraba en Data Science y analisis de datos, uno de los mas importantes ha sido el procesamiento de archivos JSON con Aeson

- Procesamiento JSON con Aeson:

Para decodificar archivos JSON se ha usado la libreria Aeson [2], que nos ha permitido recoger los datos de los archivos sin problemas. Todo esto se ha realizado con las distintas funciones del archivo *ParsersArchivos.hs*. A continuación podemos ver un ejemplo de parseado de los Tweets.

```
1 parseTweets :: String -> IO [Tweet]
2 parseTweets path = do
3     preparaTweets path
4     file <- B.readFile (path++".mod")
5     let tweets = decode file :: Maybe [Tweet]
6     case tweets of
7         Nothing -> exitSuccess
8         Just tweets -> return tweets
```

- Programa interactivo:

En cuanto al segundo concepto avanzado, la primera idea fue utilizar la documentación de la asignatura para intentar implementar una paralelización

sencilla. Sin embargo, finalmente se decidió implementar un programa interactivo, ya que sería una solución mas sencilla y rápida, además de que el tamaño de los datos de una cuenta de Twitter, especialmente de la parte que analizamos nosotros, no es tan grande como para necesitar paralelismo o clustering.

Esta interactividad se hace con las funciones del módulo *MenuInteractivo.hs*, por ejemplo, aqui podemos ver el menú principal de la aplicación:

```
1 menuPrincipal :: IO()
2 menuPrincipal = do
3     clearScreen
4     putStrLn "Hola, bienvenido al wizard de Tweekell"
5     putStrLn "Menu principal:"
6     putStrLn "Pulse enter para comenzar o escriba alguna opción"
7
8     putStrLn "\ni: Información"
9     putStrLn "q: Salir"
10
11     option <- getLine
12     case option of
13         "i" -> menuInformacion
14         "q" -> exitSuccess
15         _   -> menu1
```

Como se puede observar, se espera la interacción del usuario y se responde correctamente dependiendo de lo elegido.

8. Uso de la aplicación

Para usar la aplicación se puede ejecutar el archivo compilado en cualquier sistema GNU/Linux o, si se prefiere, se puede utilizar el código fuente. Una vez ejecutado deberíamos de ver el menu principal de la aplicación tal y como muestra la figura 1.

A continuación, siguiendo los menus introduciremos la ruta al archivo que queremos analizar y comprobaremos que lo hemos introducido correctamente tal y como se ve en la figura 2

Tras esto, podremos seleccionar el tipo de análisis que queremos realizarle al archivo, un ejemplo que podemos ver en la figura 3 serían los tipos de analisis sobre el archivo *Tweet.js*

Finalmente, tras seleccionar el tipo de análisis que queremos ver nos encontraremos con una pantalla que nos presentará los resultados. Un ejemplo seria la pantalla de la figura 4

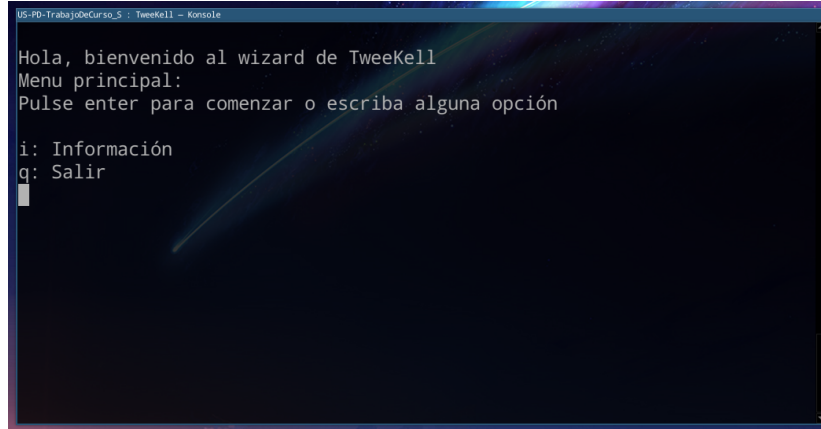


Figura 1: Menú principal de la aplicación

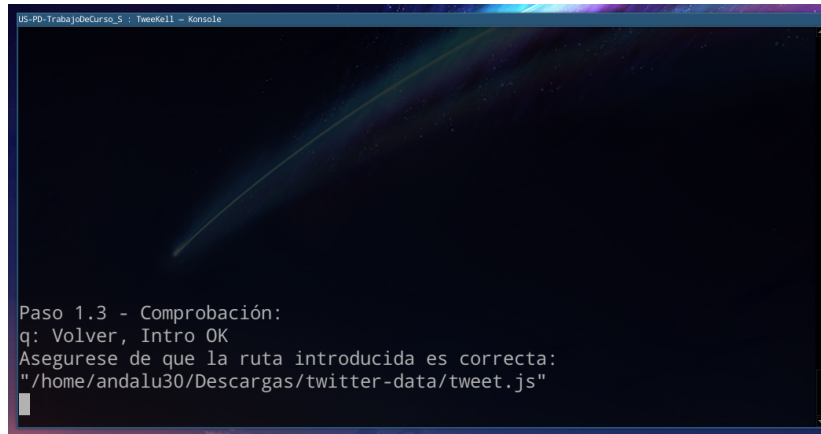


Figura 2: Pantalla de confirmación de la ruta del archivo a analizar

8.1. Utilizando el archivo compilado

El archivo compilado *Tweekell* se debe de poder ejecutar en cualquier sistema GNU/Linux sin problemas¹ mediante la ejecución del comando `./Tweekell`. Tras esto podremos ver el menu principal tal y como se ve en la figura 1

8.2. Utilizando el código fuente

Si se prefiere utilizar el código fuente para ejecutar la aplicación se podrá, por ejemplo, utilizar *ghci*. En este caso, ejecute *ghci* en una terminal, cargue el archivo *Tweekell.hs* con `:l Tweekell.hs` y ejecute *main*, debería de ver el menu principal.

¹Testeado en varios sistemas corriendo Fedora Workstation, Manjaro Linux, Ubuntu Server y Arch Linux

```
US-PD-TrabajoDeCurso_5 : TweetK11 - Konsole

Paso 2 - Tipo de análisis:
q: Volver, Intro OK

Para el archivo seleccionado existen los siguientes tipos de análisis:
3: Ver tweets
4: Ver Retweets
5: Tweet mas retweeteado
6: Tweet con más MG
7: Idiomas mas utilizados

Por favor seleccione el tipo de análisis que desea realizar:
█
```

Figura 3: Distintos tipos de análisis que se pueden realizar sobre el archivo *Tweet.js*

```
US-PD-TrabajoDeCurso_5 : bash - Konsole

El tweet con mas retweets es:
->"No se como me las apa\241o pero en todas las series que veo siempre t
iene que haber un piano por enmedio..."
con "5"
RTs
[andalu30@Jupiter US-PD-TrabajoDeCurso_5]$ █
```

Figura 4: Pantalla de resultados del análisis

9. Posibles mejoras

9.1. Analizar más información

Una de las posibles mejoras a este trabajo, sería la de analizar más información de la que se analiza actualmente. De entre los 45 archivos que Twitter nos ofrece, apenas analizamos 5 de ellos. Los archivos más interesantes de analizar serían posiblemente *like.js*, *ad-engagement.js* y *ad-impressions.js*.

9.2. Paralelización y computación distribuida

Otra posible mejora al trabajo sería la implementación de paralelismo en los análisis para que de esta forma podamos aprovechar la gran potencia de los PCs actuales.

Otra forma de escalar nuestra aplicación para que soporte grandes cantidades de

datos sería la de utilizar computación distribuida. Un ejemplo de esto sería Apache Spark [6], que nos permitiría crear clusters para poder analizar una gran cantidad de información.

Referencias

- [1] Juan Arteaga Carmona *Recreando un juego interactivo con Haskell y CodeWorld-API* <https://github.com/Andalu30/US-PD-TrabajoDeCurso> - Repositorio privado en GitHub
- [2] Bryan O’Sullivan. *Aeson: A fast Haskell JSON library* <https://github.com/bos/aeson> - Repositorio en GitHub
- [3] Paul Hudak, *Haskell* <https://www.haskell.org/> Web del lenguaje.
- [4] Hudak, Paul; Hughes, John; Peyton Jones, Simon; Wadler, Philip (2007), *A History of Haskell: Being Lazy with Class* ISBN 978-1-59593-766-7.
- [5] Twitter *Cómo acceder a tus datos de Twitter* <https://help.twitter.com/es/managing-your-account/accessing-your-twitter-data> Enlace a la web de asistencia de Twitter
- [6] The Apache Software Foundation *Apache Spark, Lightning-fast unified analytics engine* <https://spark.apache.org> Enlace a la web