

ITAM Text Miner

Felipe Gerard

2 de noviembre de 2015

0. El problema por resolver

1. Orquestación

Es importante que un sistema de análisis de esta escala sea robusto y confiable. Para ello se optó por utilizar el orquestador `luigi`, que fue desarrollado por Spotify en el lenguaje `Python`. Este esquema tiene varias ventajas sobre los scripts simples, algunas de las cuales se enumeran a continuación:

1. *Modularidad*: El código se segmenta en pasos bien definidos, lo que facilita su mantenimiento y expansión.
2. *Robustez*: Se requiere que todos los pasos serialicen sus resultados a disco, lo que hace que el sistema sea resistente a fallos y que no tenga que repetir todo el proceso después de un error.
3. *Idempotencia*: Si se programa correctamente, los procesos sólo corren una vez. Incluso si se corre el proceso una segunda vez, sólo corre lo que no haya sido corrido antes.

Obtener estas ventajas es relativamente simple con `luigi`. Lo único que se requiere es que los pasos tengan entradas y salidas bien definidas; un proceso únicamente requiere sus dependencias para correr y debe arrojar los resultados que necesite la siguiente sección.

`luigi` corre nativamente en `Python`, por lo que se puede hacer cualquier cosa que se pueda hacer en `Python`. Esto incluye llamar códigos en `shell`, `R` y muchos otros lenguajes. La mayoría de los procesos se hicieron directamente en `Python` por simplicidad, pero algunos se dejaron en otros lenguajes por diversas razones de conveniencia.

a. Extracción de texto

En esta sección definimos a grandes rasgos los diversos componentes del *pipeline* de datos, desde los archivos en PDF crudos hasta las salidas necesarias para mostrar en el producto final.

Pegado de PDFs en archivos de 50 MB

-> Petri

Extracción de texto a partir de PDFs

Input: PDFs crudos, en el formato descrito abajo.

Output: Libros en formato texto tal cual fue extraído, archivos JSON con extractos de los libros.

Paquetes relevantes: `pdfminer`

El proceso empieza con una carpeta general en la que deben estar todos los PDFs. Adentro de ella debe haber una carpeta por libro que contenga los PDFs de las diversas hojas, por ejemplo: `pdf/libro_de_arte/hoja_i.pdf`

Se extraen los textos utilizando el paquete `pdfminer` de `Python` y se juntan los resultados de todas las hojas en un solo archivo por libro. Los libros en versión texto se meten en una carpeta según su idioma. Adicionalmente, se guardan dos archivos de metadatos, uno con los idiomas registrados y otro con una relación entre los libros y sus idiomas correspondientes.

Por practicidad se optó por extraer los párrafos representativos en este mismo proceso. Para ello se tomaron 3 secciones de aproximadamente 500 letras al 10%, 50% y 90% de avance del libro aproximadamente. El criterio tiene cierta flexibilidad e intenta encontrar párrafos completos (que empiecen en salto de línea y mayúscula), aunque esto en general es poco preciso por la imperfección del OCR efectuado al momento de digitalizar los libros.

Nota 1: Esta extracción se considera que tiene nivel de limpieza “nulo”, por lo que los resultados de la extracción se meten a la carpeta `raw` dentro de la carpeta de textos.

Nota 2: Por razones internas del funcionamiento de `luigi`, es imposible checar de antemano a qué carpeta debe ir cada libro, ya que no se puede saber su idioma antes de procesarlo. Para darle la vuelta a esta limitación, se generó una carpeta con un archivo de metadatos por libro, `libro.meta`, que sirve para que `luigi` pueda checar el árbol de dependencias.

Limpieza de textos

Input: Libros en formato texto crudo.

Output: Libros en formato texto limpio, según la especificación.

Paquetes relevantes: `nltk`, `unicodedata`

Según el nivel de limpieza elegido, se genera una estructura similar a la de los textos crudos (`raw`). Los dos niveles son incrementalmente más “limpios”:

- Limpieza básica (`clean`):
 - Se quitan los acentos.
 - Se quitan los saltos de página.
 - Se pasa todo a minúsculas.
 - Se quitan los caracteres especiales que queden después de quitar los acentos.
 - Se quitan palabras cortas (de 3 o menos caracteres).
 - Se quitan palabras con algún carácter repetido 3 o más veces.
- Limpieza avanzada (`stopwords`):
 - Todo lo de la limpieza básica.
 - Se quitan las palabras poco informativas (`stopwords`).

Los dos pasos anteriores se generan de manera incremental, lo que significa que si se genera uno más avanzado, siempre se generan todos los anteriores. Cada nivel de limpieza genera una estructura similar a la generada en el paso de extracción (i.e. a la carpeta `raw`), pero con nombre y características según sea el caso.

Nota: Nuevamente se utilizó un esquema de dependencias artificiales como en la extracción.

b. Minería de textos

Detección de idioma

Vectorización de textos

Latent Dirichlet Analysis: Clusters por tópicos

→ Lechuga

Latent Semantic Indexing: Documentos similares

c. Paquete