## Trabajo Práctico Final, Progr. 2

Ivan Campos Wainer UNR

FCEIA.

## Resumen del objetivo

Proceso de selección de conceptos

Antes de seleccionar mi estrategia final para abarcar este trabajo, lleve adelante un proceso de búsqueda de distintas formas para alcanzar el objetivo

### **Primer concepto**

A partir de la obtención del archivo y su posterior separación a través de un split, conseguí una lista con las palabras del mismo. A partir de esto, llevé a cabo la creación de un diccionario en el que las claves eran cada palabra del archivo, y su valor, una lista que contiene las palabras que la anteceden y preceden (teniendo en cuenta que puede ser mas de una opción).

Por ejemplo, la frase "el amor verdadero y el agua es azul"

```
{
    "el": ["amor", "y", "agua"]
}
```

La gran desventaja que encontré, es que al usar un diccionario de N claves (N = cantidad de palabras), y una lista compuesta por valores de uno como longitud mínima, se vuelve compleja de recorrer y quedan M valores (en el mejor de los casos, M=N) que se vuelven innecesarios.

## **Segundo Concepto**

Buscando otra alternativa, intenté asignar una probabilidad a cada palabra del texto.

Para llevarlo a cabo, guardé en un diccionario cada palabra como una clave, pero asignándole un porcentaje como valor, y tomando como probabilidad la cantidad de

veces que la misma aparecía dentro del texto.

Por ejemplo, así quedaría la siguiente oración "el amor despues del amor":

```
{
  "el": 20,
  "amor": 40,
  "despues": 20,
  "del": 20
}
```

A diferencia de las ideas anteriores, esta posibilidad tiene como valor un número (float) y no posee una lista de N elementos. Sin embargo, al igual que la idea inicial, tiene como desventaja que debe guardar todas las palabras del texto dentro de una estructura, implicando elementos no deseados.

Finalmente, como ninguna de estas ideas me pareció totalmente adecuada para solucionar la consigna, planteé la siguiente idea.

## **Concepto elegido**

Por últmo, separé las palabras que estaban alrededor del guión bajo de la frase y busqué en la letra de la canción los siguientes puntos para guardarlas en un diccionario junto a la probabilidad de aparición:

- La palabra que la precede, en caso de que el guión se encuentre detrás
- La palabra que la antecede, en caso de que el guión se encuentre delante

Por ejemplo:

#### Letra:

El amor, después del amor tal vez

Se parezca a este rayo de sol

Y ahora que busque y ahora que encontre

El perfume que lleva el dolor

#### Frase:

Y ahora que \_ y ahora que encontre

#### Diccionario:

```
{
    "busque": 2,
    "encontre": 1
}
```

#### Resumen parseo de texto en C:

Para el parseo en C, lo primero que hice fue crear el archivo de texto que tiene en su interior los nombres de todos los archivos contenidos en la carpeta del artista separadas por un '\n'. Luego, por cada una de las canciones, apliqué la función add\_text\_to\_output que, como explico en el código, se encarga de iterar por cada caracter contenido en el archivo y "empujarlos" a uno nuevo, formateando el texto. Para finalizar, el programa ejecuta el codigo en Python.

#### Resumen sustitución de barra baja y archivo de salida en Python:

Primero, el programa lee el archivo del texto formateado, la letra del artista y las frases. A continuación, formatea las palabras que incluye cada una, actuando como un auxiliar que sustituye las mayúsculas y elimina los puntos y comas de las oraciones. Después, separa las palabras que rodean al guión y llama a la función más importante del programa: add words que verifica 5 casos:

- Cuando el guión bajo se encuentra entre las dos palabras coincidentes
   Como vemos en el ejemplo, el guión bajo se encuentra entre dos palabras presentes en la canción: "que" e "y". Y es por eso que se agrega al diccionario con +100 de chances de ser el que sustituya el mismo.
- Cuando el guión bajo está primero en la frase
   Se puede dar que el guión se encuentre al principio de la frase, por lo que buscará unicamente en la letra aquellas palabras que están detrás de la que precede el guión.
- Cuando el guión bajo esta a lo último de la frase
   Se puede dar que el guión se encuentre al final de la frase, por lo que buscará unicamente en la letra aquellas palabras que están delante de la que antecede el guión.

Cuando la palabra está atrás o adelante del guión

El diccionario contiene las palabras y su posibilidad de sustituir el guión bajo con la que se encuentre adelante y atrás del guión mientra recorre el texto.

Aclaración: Es importante que este caso se verifique luego de los desarrollados anteriormente.

Cuando se necesite aumentar el radio

En algunos casos, las palabras que rodean el guión, no se encuentran en el texto. Por lo tanto, se debe aumenta el radio y se busca en i elementos detras y delante del mismo, siendo i = palabras distancia del guión.

Fuera de add\_words, en caso de que ninguna palabra que rodee el guión bajo exista dentro de la letra de la canción, se le asigna una palabra random que se encuentre en ella.

Para finalizar, busca las palabras con las probabilidades más altas dentro del diccionario y sustituye la barra baja con las mismas.

## Serie de pasos

A continuación, desarrollé la serie de pasos que lleve adelante para hacer este trabajo

## **▼** Parseo de texto en C

1. Ejecución de generate text files name y apertura de archivo

```
/* Recibe como parametro el nombre del directorio al que debemos acceder
 2 dentro de Textos/ luego, ejecuta el comando ls y guarda en el directorio raíz, un archivo llamado
 3 text files name.txt *
    int generate_text_files_name(char name[])
5
6
7
         char command[100];
8
        sprintf(command, "./Textos/%s", name);
FILE *faux = fopen(command, "r");
10
11
         if (!faux)
12
13
             return -1;
14
15
16
17
        }
        fclose(faux);
18
        sprintf(command, "cd Textos/%s && ls > ../../text_files_name.txt", name);
19
20
21
         return 1;
22 }
```

Definición de la función

```
if (generate_text_files_name(argv[1]) == -1)
    {
        printf("No existe el directorio: Textos/%s\n", argv[1]);
        return -1;
    }
}
FILE *text_files_name = fopen("./text_files_name.txt", "r");
```

Llamada de la función y apertura del archivo en main



Ejemplo de salida de text\_files\_name.txt

2. Declaración de string auxiliar path y creación de archivo de salida, (Entradas/nombre.txt)

```
char path[100]; // String auxiliar para guardar rutas (se irá modificando a lo largo del programa)
sprintf(path, "./Entradas/%s.txt", argv[1]);
FILE *file_output = fopen(path, "w");
```

Aclaración: se abre en modo 'w' para sobrescribir el archivo del artista en caso de que se ejecute de nuevo el programa.

3. Declaración del bucle while por cada línea de text\_files\_name y ejecución de add\_text\_to\_output por cada archivo recorrido (almacenado en path)

```
1 while ((fscanf(text_files_name, "%s", path)) != EOF)
2 {
3     add_text_to_output(argv[1], path, file_output); // <---
4 }</pre>
```

```
void add_text_to_output(char *name, char *filetxt, FILE *file_output)
{
    // ...
}
```

Le pasa por parámetros el nombre del artista, el path, que sería el nombre del archivo recorrido y el file\_output (Entradas/)

• a. Declaración de variable auxiliar path y archivos

```
char path[100];
sprintf(path, "./Textos/%s/%s", name, filetxt);
FILE *file_input = fopen(path, "r");
```

 b. Declaración de ciclo while y variables auxiliares character, prev\_character

```
char character, prev_character;
while ((character = fgetc(file_input)) != EOF)
{
    //...
}
```

Al igual que antes, recorre el archivo de la letra (file\_input)

 c. Condiciones para el parseo (para poner carácter en el file\_output) y sustitución de caracteres

```
while ((character = fgetc(file_input)) != EOF)
3
       if (esMayus(character))
           fputc(character + 32, file_output);
4
       else if (character == '\n' && prev_character != '.')
5
           fputc(' ', file_output);
6
       else if (character != '.' && character != ',')
7
8
            fputc(character, file_output);
9
10
       }
11
12
        prev_character = character;
13 }
```

• d. Agregado de \n para evitar bugs y cierre de archivo

```
1 /* d. Agregado de \n para evitar bugs y cierre de archivo */
2 fputc('\n', file_output);
3
4 fclose(file_input);
```

4. Liberación de memoria, cierre de archivos y llamada al programa Python

```
1 fclose(file_output);
2 fclose(text_files_name);
3
4 sprintf(path, "python3 main.py %s", argv[1]);
5 system(path);
```

# ▼ Sustitución de barra baja y archivo de salida en Python

1. Declaración de la función main y argumentos

```
1 def main(name):
2  # ...
3
4 """ ARGUMENTS """
5 script, name = argv
6 main(name)
```

2. Apertura de archivo de los textos parseados; creación de lista

```
1 f_lyrics = open("./Entradas/" + name + ".txt")
2 lyrics = f_lyrics.read().replace("\n", " ").split(" ") # La lista contiene como elementos cada palabra del archivo
3 f_lyrics.close()
4
```

La lista contiene como elementos cada palabra del archivo

3. Declaración de f\_phrase y ejecución de read\_file (función auxiliar) y creación de archivo de salida

```
1 f_phrase = read_file('./Frases/' + name + '.txt')
2
3 file_output = open("./Salidas/" + name + ".txt", 'w+')
```

Llamada a la función desde el main

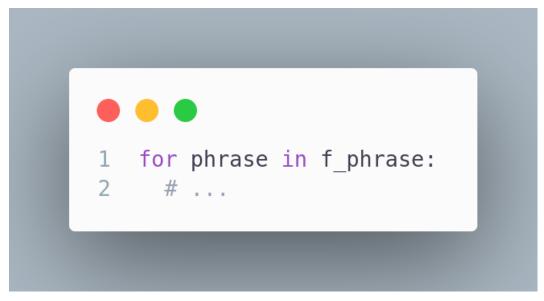
```
""" Se encarga de obtener el archivo que contiene las frases del artista
     para luego iterar en cada salto de linea:
      - guardar su contenido en forma de lista (word_line_aux)
- eliminar las comillas simples '' que quedan de sobra en word_line_aux
- guarda word_line_aux en content
     por ultimo, reotrna content, que es una lista que contiene listas de cada palabra contenida en linea del archivo frase """
 8 def read_file(file_name):
      file = open(file_name, 'r')
lines_of_text = file.readlines()
10
12
13
14
15
16
17
18
19
20
21
22
23
       content = []
       for line_of_text in lines_of_text:
          word_line_aux = line_of_text.replace('\n', '').split(' ')
         # Eliminar '' de la lista
if '' in word_line_aux:
            word_line_aux.remove('')
          content.append(word_line_aux) # Guarda ['Te', '_,', 'unos', 'chinos', 'en', 'madrid'] dentro de la lista
       file.close()
        return content
```

Función auxiliar read\_file

```
""" Frases/Fito_Paez.txt """
2 Te _, fumabas unos chinos en madrid
3 nadie puede y nadie debe vivir _ amor
4
5 f_phrase = [
6     ['Te', '_,', 'fumabas', 'unos', 'chinos', 'en', 'madrid'],
7     ['nadie', 'puede', 'y', 'nadie', 'debe', 'vivir', '_', 'amor']
8 ]
```

Ejemplo de cual sería el resultado de f\_phrase con ese archivo .txt

#### 4. Iteración por cada frase en lista de frases



Recordemos que f\_phrase es una lista de listas. Entonces phrase representa una lista de una frase

5. Declaración de phraseAux (frase formateada para trabajar con ella), ejecución de format\_text\_phrase

```
format_text_phrase: List(String) -> None | List(String)
format text_phrase se encarga de formatear la frase, eliminado comas, puntos
después e independientemente de '_'. También, convierte las mayusculas en minusculas.
"""
def format_text_phrase(phrase): # Retorna NONE en caso de que no haya '_'
phraseAux = ' '.join(phrase).replace("_,", "_").replace("_.", "_").replace(".", "").lower().split(" ")
return phraseAux
```

Definición de la función

```
for phrase in f_phrase:

phraseAux = format_text_phrase(phrase) # Utilizo un auxiliar formateado para mantener las frases originales de f_phrase

if '_' in phraseAux:

# ...
```

Verifica si existe '\_', debido a que puede haber frases que no contengan \_

• a. Declaración de words isolate y ejecución de isolate words

```
"""
2 isolate_words: List(String), Int -> List(String)
3 isolate_words se encarga de "aislar" las palabras en un radio n, que rodean el guion bajo para retornar una lista
4 que contiene a las mismas. Como maximo, pueden haber dos palabras aisaldas.
5 """
6 def isolate words(phrase, radius):
7 words_isolate = []
8 i = 0
9
10 while i < len(phrase) and phrase[i] != '_':
11 if i+radius < len(phrase) and phrase[i+radius] == '_':
12 words_isolate.append(phrase[i])
13 i += 1
14
15 if i+radius < len(phrase) and phrase[i] == '_': # En caso, el guión se encuentre como primer caracter
16 words_isolate.append(phrase[i+radius])
17
18 return words_isolate</pre>
```

Definición de la función isolate\_words

```
1 if '_' in phraseAux:
2 words_isolate = isolate_words(phraseAux, 1)
```

Devuelve una lista como máximo de dos elementos, que contiene las words que rodean el guion bajo

 b. Declaración de dict (diccionario con las palabras y probabilidad de que sustituya el '\_') y ejecución add\_words

```
*** add words: List(String), List(String) - List(String) -> dict

2 add words recibe como parametro la letra, la lista de las palabras aisladas que

3 rodean el '' (retorna isolate words) y la lista que contiene las frases formateadas.

4 Sus varios ifs sirven para verificar caoss unicos en el que el programa surifa un bug. Como

cuando el '_' es el primer elemento, o el ultimo; también cuando se necesita aumentar el radio.

7 Retorna un diccionario, sus claves serán las posibles palabras encontradas y sus valores sus posibilidades.

8 claraccion: En caso de estar entre dos palabras conocidas en el texto, se le suma en vez de una, dos chances, debido

9 a que es mas probable***

1 def add words(lyrics, words_isolate, f_phrase):

1 dict = {}

1 for i in range(0, len(lyrics)):

4 # Caso en el que el Guion baje se encuentre entre las dos palabras coincidentes

1 if len(words isolate) == 2 and i+2 < len(lyrics)-1 and lyrics[i] == words_isolate[0] and lyrics[i+2] == words_isolate[1]:

2 dict[lyrics[i+1]] = dict_eqt(lyrics[1+1], 0) + 2

3 # Caso en que el Guion bajo esté primero

2 elif [aphrase[0] == '' and lyrics[i] == words_isolate[0]:

3 dict[lyrics[i-1]] = dict_eqt(lyrics[i+1], 0) + 1

2 # Caso en que el Guion bajo esté primero

2 elif [phrase[0] == '' and lyrics[i] == words_isolate[0]:

3 dict[lyrics[i-1]] = dict_eqt(lyrics[i+1], 0) + 1

2 # Caso en que el Guion bajo esté primero

2 elif [phrase[0] == '' and lyrics[i] == words_isolate[0]:

3 dict[lyrics[i-1]] = dict_eqt(lyrics[i+1], 0) + 1

2 # Caso en que el Guion bajo esté primero

2 elif [phrase[0] == '' and lyrics[i] == words_isolate[0]:

3 dict[lyrics[i-1]] = dict_eqt(lyrics[i+1], 0) + 1

2 # Caso en que el Guion bajo esté primero

2 elif [phrase[0] == '' and lyrics[i] == words_isolate[0]:

3 dict[lyrics[i-1]] = dict_eqt(lyrics[i+1], 0) + 1

3 # Caso en que el Guion bajo esté primero

2 elif [phrase[0] == '' and lyrics[i] == words_isolate[0]:

3 dict[lyrics[i-1]] = dict_eqt(lyrics[i+1], 0) + 1

3 # Aumenta el indice de la palabra (key ded dict1 que se encu
```

```
1 if '_' in phraseAux:
2 words_isolate = isolate_words(phraseAux, 1)
3
4 dict = add_words(lyrics, words_isolate, phraseAux) # <--</pre>
```

Agrega palabras que se encuentran entre los dos elementos, o el elemento, en caso de que el guion bajo se encuentre ultimo o primero

 c. Caso de que las palabras que rodean el guion bajo, no existan, aumenta el radio

```
1  i = 2
2  while i < len(phrase)+1 and len(dict) == 0:
3   try:
4    dict = add_words(lyrics, isolate_words(phraseAux, i), phraseAux)
5    i+=1
6   except:
7   dict[lyrics[randint(0, len(lyrics))]]</pre>
```

Se pueden dar varios casos, en el que la palabra que antecede o precede el guion, no existan en la letra. Por lo tanto, se aumenta el radio y se busca la palabra que está detrás de la palabra que antecede el '\_'; el radio se puede aumentar N cantidad de veces (siendo N=longitud de la frase + 1). Es importante aclarar, que se busca generar una excepción en caso de que ninguna de las palabras que están en la frase existan en la letra; por lo tanto, se le asignará una palabra random (debido a que tiene la misma probabilidad de sustituir el '\_' como cualquier otra)

 c. Llamada a la función print\_phrase\_in\_output\_file encargada de "colocar" frase en archivo de salida con sus correspondiente guion bajo sustituido

```
print_phrase_in_output_file: List(String), dict, FILE -> None
Recibe una lista de la frase sin modificar, el diccionario con la clave de las palabras mas
probables y su valor de probabilidad y el archivo de salida donde serán escrita la frase.
Crea la variable keys max val, que obtiene una lista con las CLAVES del diccionario con mayores valores.
Luego, obtiene una clave random de la lista keys max_val; esto se debe a que si hay palabras que tienen
el mayor valor del diccionario, eligirá una random porque tiene la misma probabilidad de ser la que sustituya
el guino bajo.
Finalmente, sustituye el '_' con la palabra seleccionada, y la escribe en el archivo de salida

def print phrase in output file(phrase, dict, file_output):
Resy max val = get_keys max_val(dict)
random_word = keys_max_val(fandint(0, len(keys_max_val) - 1)] # Caso en el que las probabilidades de las palabras en el dicionario sean iguales
final_phrase = ''__join(phrase).replace('_', random_word)

file_output.write(final_phrase)
file_output.write(final_phrase)
file_output.write('\n')
```

Definición de la función print\_phrase\_in\_output\_file

```
1 print_phrase_in_output_file(phrase, dict, file_output)
```

Llamada de la función dentro del for (fuera del if)