

PRÁCTICA 1 - ANÁLISIS LÉXICO

MÓDULO I – LENGUAJE ESTRUCTURADO

PROCESAMIENTO DE LENGUAJE ESCRITO

GRADO EN CIENCIA E INGENIERÍA DE DATOS

2022/2023

PROFESORES: CARLOS DAFONTE – carlos.dafonte@udc.es
DANIEL GARABATO – daniel.garabato@udc.es

DEPARTAMENTO DE CIENCIAS DE LA COMPUTACIÓN Y TECNOLOGÍAS DE LA INFORMACIÓN
UNIVERSIDADE DA CORUÑA

Diseño e implementación mediante la herramienta **SLY** de un analizador léxico para el tratamiento de mensajes de *log* centralizados del servicio de autenticación de una máquina Linux.

La entrada del analizador será un fichero de texto plano como el que se adjunta junto con este enunciado (**auth_example.log**). Cada línea representa un mensaje de *log*, compuesto por los siguientes campos:

- **Fecha:** Compuesto por una secuencia de letras que representan el mes, seguidas de una secuencia de dígitos que representa el día (separados por un espacio o un tabulador).
- **Nombre de la máquina:** Secuencia de caracteres alfanuméricos que representan el nombre del equipo en que se registró el evento.
- **Servicio:** Secuencia de caracteres que identifica el proceso o servicio que registra el evento. Nos centraremos en el servicio SSH: **sshd[%PID%]**.
- **Mensaje:** Encapsula el tipo de evento que se ha registrado, junto con información adicional. Trataremos especialmente aquellos asociados a un usuario autenticado correctamente (“**Accepted password for...**”), una contraseña errónea para un usuario (“**Failed password for..**”), o un usuario no autorizado (“**Invalid user...**” y “**Failed password for invalid user**”). Cabe destacar que existen muchos otros tipos de eventos que, no obstante, no consideraremos para el desarrollo de la práctica y simplemente catalogaremos como “otros”.

El objetivo de la práctica consiste en implementar los siguientes contadores:

- Número total de eventos.
- Número de eventos por mes y quincena (primera o segunda).
- Número de eventos según la franja horaria: mañana (08:00h – 15:59h), tarde (16:00 – 23:59h) y noche (00:00h – 07:59h).
- Número de eventos por máquina.
- Según el tipo de evento identificado:
 - Usuarios autenticados correctamente (“**Accepted password for...**”).
 - Contraseña errónea (“**Failed password for...**”).
 - Usuario no autorizado (“**Invalid user...**” y “**Failed password for invalid user...**”)

Se obtendrá:

- Número total de ocurrencias.
- Número de eventos por máquina.
- Número de eventos por usuario.

- Número de eventos por clase (A, B o C) y tipo (privada o pública) de IP.
- Número total de ocurrencias para los demás tipos de eventos (otros).

Deben definirse adecuadamente cada uno de los **tokens** a identificar, especialmente aquellos que se necesiten contabilizar de algún modo. Aquellos elementos que no requieran tratamiento podrán ignorarse mediante expresiones regulares, pero no pueden dar lugar a errores de procesamiento.

Junto con el enunciado, se proporciona un esqueleto en python (`p1_base.py`) del que tendréis que partir y cuya sección `__main__` **NO debéis modificar**. En ella se lee íntegramente el fichero de *log* proporcionado como entrada por el sistema estándar. Así mismo, procesa la entrada con el analizador léxico e invoca el método `print_output` incluido en el esqueleto, que debéis completar adecuadamente para mostrar por pantalla el valor de los contadores arriba indicados, siguiendo este formato:

```
1  #contadores_generales
2  total_eventos,x
3  total_aceptados,x
4  total_fallidos,x
5  total_no_autorizados,x
6  total_otros,x
7  #eventos_por_fecha
8  mes_1,primera,x
9  mes_1,segunda,x
10 mes_2,segunda,x
11 [...]
12 #eventos_por_hora
13 manana,x
14 tarde,x
15 noche,x
16 #eventos_por_maquina
17 maquina_1,x
18 maquina_2,x
19 [...]
20 #eventos_aceptados_por_maquina
21 maquina_1,x
22 maquina_2,x
23 [...]
24 #eventos_aceptados_por_usuario
25 usuario_1,x
26 usuario_2,x
27 [...]
28 #eventos_aceptados_por_ip
29 clase_a,privada,x
30 clase_a,pública,x
31 clase_b,privada,x
32 clase_b,pública,x
```

```

33 clase_c,privada,x
34 clase_c,pública,x
35 #eventos_fallidos_por_maquina
36 maquina_1,x
37 maquina_2,x
38 [...]
39 #eventos_fallidos_por_usuario
40 usuario_1,x
41 usuario_2,x
42 [...]
43 #eventos_fallidos_por_ip
44 clase_a,privada,x
45 clase_a,pública,x
46 clase_b,privada,x
47 clase_b,pública,x
48 clase_c,privada,x
49 clase_c,pública,x
50 #eventos_no_autorizados_por_maquina
51 maquina_1,x
52 maquina_2,x
53 [...]
54 #eventos_no_autorizados_por_usuario
55 usuario_1,x
56 usuario_2,x
57 [...]
58 #eventos_no_autorizados_por_ip
59 clase_a,privada,x
60 clase_a,pública,x
61 clase_b,privada,x
62 clase_b,pública,x
63 clase_c,privada,x
64 clase_c,pública,x

```

Obviamente, los valores *x* tendrán que corresponderse con el número de ocurrencias identificadas para cada caso. Los nombres de los meses, así como los nombres de máquinas y de usuarios tendrán que corresponderse con los que figuren en los registros. Las cabeceras (líneas precedidas por #) deben ser las que se indican, mientras que la secuencia [...] se ha utilizado a modo de ejemplo para indicar una sucesión de elementos (los que sean). Si para algún caso no se detectasen ocurrencias, la cabecera debe aparecer igualmente, pero sin sucesión de elementos (pasando directamente a la siguiente cabecera). Si se omite algún elemento en las listas, se asumirá que tiene 0 ocurrencias.