

Raw Logs Processing

Nitish Rangarajan

November 22, 2017

Preprocess the logs into a csv dataset

The raw dataset was collected as a log from the servers. The log has two types of arrays namely Persistent_Storage and RoundInfo. [3] The avatar's history is stored in the Persistent_Storage. Each element in this array contains the information about the avatar during the observed timeframe. They contain 11 fields which are separated by commas. They are dummy, query_time, query_sequence_number, avatar_ID, guild, level, race, character_class, zone, dummy and dummy.

Create a regular expression that would process the logs.

```

import csv
import re
import os          # This is for os.listdir
import os.path     # This is for the other dir stuff.
import string      # Maybe for directory name cycling etc.
import time        # For timing how long it takes.

# VARIABLES
max_char = 0       # Tracks char ID for error testing.
max_guild = 0      # Tracks guild ID for error testing.
mfc = 0            # Tracks how many little files were counted.
write_data_loc = 'C:/Users/nrangara/Downloads/WorldOfWarcraft/output/'
                  # Adjust to your dir as needed.
write_data_file = 'wowah_data.csv'
write_data_filename = write_data_loc + write_data_file
the_dir = 'C:/Users/nrangara/Downloads/WorldOfWarcraft/WoWAH/' # This is
                  # where the WoWAH folders are located, adjust as needed. Have them in their own subdir.

#REGEX
line_re = re.compile(r'^.*"[\d+],\s(.*)\s?(\d+)\s?(\d+)\s?(\d*)\s?(\d*)\s?([A-Z].*)\s?([A-Z].*)\s?([A-Z].*)".*$')
#
#           dummy   time1   seq2   3char   4guild   5level   6race   7charc
lass      8zone
#line_re = re.compile(r'^.*"[\d+],\s(.*)\s?(\d+)\s?(\d+)\s?(\d*)\s?(\d*)\s?([A-Z][0-9].*)\s?([A-Z][0-9].*)\s?([A-Z][0-9].*)".*$')
#
#           dummy   time1   seq2   3char   4guild   5level   6race   7cha
rclass   8zone
#line_re = re.compile(r'^.*"[\d+],\s(.*)\s?(\d+)\s?(\d+)\s?(\d*)\s?(\d*)\s?([A-Z][0-9].*)\s?([A-Z][0-9].*)\s?([A-Z][0-9].*)".*$')
# REGEX NOTES
# groups: 1=timestamp, 3=avatarID, 4=guild.
# [1] = "0, 03/30/06 23:59:49, 1,10772, , 1, Orc, Warrior, Orgrimmar, , 0",
#      "0, 01/10/09 00:03:50, 1,55517, , 3, Orc, Warlock, Orgrimmar, WARLOCK, 0", -- [1]
#      "0, 01/10/09 00:04:10, 5,4002,1, 75, Orc, Hunter, Zul'Gurub, HUNTER, 0", -- [26]
#      "0, 01/10/09 00:04:10, 5,78122,342, 80, Orc, Hunter, The Storm Peaks, HUNTER, 0", -- [3
2]
#      "0, 01/10/09 00:08:04, 51,64635,161, 80, Blood Elf, Paladin, The Obsidian Sanctum, PALAD
IN, 0", -- [447]
# dummy, query time, query sequence number, avatar ID, guild, level, race, class, zone, dummy, d
ummy

def get_subdirs(the_folder):
    this_list = []
    this_list = os.listdir(the_folder)
    print ('From get_subdirs, a list is: ', this_list) # Printing for error control.
    for item in this_list:
        if item.startswith('.'):
            this_list.remove(item)
    return(this_list)
# End of get_subdirs
# '.DS_Store'

def get_file_list(the_folder): # Yes these two are the same, just diff names.
    this_list = []

```

```

this_list = os.listdir(the_folder)
for item in this_list:
    if item.startswith('.'):
        this_list.remove(item)
return(this_list)
# End of get_file_list

def parse_and_write(file, output_file):
    for line in file:
        # Oh the first "line" is a hard return???
        # print 'A line is: ', line
        data = line_re.match(line)
        #print (line)
        if data is not None:
            #print ("Matched the regex.")
            timestamp = data.group(1)
            char = data.group(3)
            level = data.group(5)
            race = data.group(6)
            charclass = data.group(7)
            zone = data.group(8)

            if data.group(4) is not(''):
                guild = data.group(4)
            else:
                guild = '-1' # Note there are some missing values, i.e. errant -1.
            #print (timestamp) # This is so you can keep track of where it is. Max Jan 2009 I
IRC.

            new_line = char + ',' + level + ',' + race + ',' + charclass + ',' + zone + ',' + g
uild + ',' + timestamp + '\n'
            #new_line = char + ',' + guild + ',' + timestamp + '\n'
            output_file.write(new_line)

        #else:
            #print ("Didn't match the regex.")

# End of parse_and_write
# Note the two diff formats they use in the files, it changes partway through:
# [1] = "0, 03/30/06 23:59:49, 1,10772, , 1, Orc, Warrior, Orgrimmar, , 0",
# "0, 01/10/09 00:03:50, 1,55517, , 3, Orc, Warlock, Orgrimmar, WARLOCK, 0", -- [1]
# dummy, query time, query sequence number, avatar ID, guild, level, race, class, zone, dummy, d
ummy

def read_tree(output_file):
    global the_dir
    months_folders = get_subdirs(the_dir) # This is why the subdirs should be in their own
location that you set in the vars section up top.
    for folder in months_folders: # Run isdir(dir) first, try/
except. Make sure no funny folders/dirs.
        folder = the_dir + folder # Expands the folder name to
the long version.
        day_folders = get_subdirs(folder)
        for day_folder in day_folders:

```

```

day_folder = folder + '/' + day_folder
file_list = get_file_list(day_folder)
for file in file_list:
    try:
        file = day_folder + '/' + file
        with open(file, 'r', encoding='utf8') as f:
            this_file = f.readlines()
            # Should read the whole
            # file as a string?
            parse_and_write(this_file, output_file)
    except IOError:
        print ('Error opening hoped for data-text file,', str(file), ', reason: ', IOError)
# End of read_tree

def main():
    #open write file here
    output_file = open(write_data_filename, 'a', encoding='utf8')    # 'a' is very important, it
    # appends the new data to the big file.
    fieldnames = ('char, level, race, charclass, zone, dummy1, dummy2, guild, timestamp\n')
    output_file.write(fieldnames)
    start_time = time.time()
    read_tree(output_file)
    #close write file here
    output_file.close()
    spent_time = time.time() - start_time
    mins_spent = int(spent_time / 60)
    secs_remainder = int(spent_time % 60)
    print ('Time of process: ', mins_spent, ':', secs_remainder)    # 13m:42s on iMac. Also 14
    # m:39s another time.

#    print 'Files scanned (or tried), ', mfc    # 138,084
#    print 'Max Chars: ', max_char    # They say 91,065 ">= 1" but it starts at 0, my
#    count says: 91064 + 1 = 91,065.
#    print 'Max Guilds: ', max_guild    # They say "An integer within [1, 513]" but no
#    since they start at 0. 512 + 1 = 513.
# End of main

# Main call

main()

```

[illegible]

[illegible]

[illegible]

```
8-11-22', '2008-11-23', '2008-11-24', '2008-11-25', '2008-11-26', '2008-11-27', '2008-11-28', '2008-11-29', '2008-11-30', '2008-12-01', '2008-12-02', '2008-12-03', '2008-12-04', '2008-12-05', '2008-12-06', '2008-12-07', '2008-12-08', '2008-12-09', '2008-12-10', '2008-12-11', '2008-12-12', '2008-12-13', '2008-12-14', '2008-12-15', '2008-12-16', '2008-12-17', '2008-12-18', '2008-12-19', '2008-12-20', '2008-12-21', '2008-12-22', '2008-12-23', '2008-12-24', '2008-12-25', '2008-12-26', '2008-12-27', '2008-12-28', '2008-12-29', '2008-12-30', '2008-12-31']  
## From get_subdirs, a list is: ['2009-01-01', '2009-01-02', '2009-01-03', '2009-01-04', '2009-01-05', '2009-01-06', '2009-01-07', '2009-01-08', '2009-01-09', '2009-01-10']  
## Time of process: 8 : 55
```