

目录

项目说明

1. 变量设置
2. `get_filters` 函数
2. `load_data` 函数
3. 其他函数说明

项目说明

本次项目是依据文档信息来解决相应的问题，主要是看文档字符串了解函数作用，查看参数信息了解参数的数据类型以及参数的信息。特别注意⚠️脚本中的 `main` 函数和最后的 `if` 条件语句部分不用修改——可以进行优化，但是不要修改文件结构

1. 变量设置

变量使用，根据变量的作用来判断使用何种方式的变量——全局变量以及局部变量。例如 MONTHS 和 CITY_DATA 等可能在全局使用的可以在文件开始申明并赋值

```
1  #!/usr/bin/env python
2  # -*- coding:utf-8 -*-
3
4  import time
5  import pandas as pd
6  import numpy as np
7
8  CITY_DATA = {'chicago': 'chicago.csv',
9              'new york city': 'new_york_city.csv',
10             'washington': 'washington.csv'}
11
12 MONTHS = ['january', 'february', 'march', 'april', 'may', 'june']
```

2. `get_filters` 函数

根据文档字符串信息，确认需要筛选出来的哪些数据，三个结果需要一定限定条件才能输出，另一个需要解决的是用户交互问题——怎么保障用户输入是正确的值，筛选三个值都是一样的结构，这里隐藏优化点——使用函数。所以该问题需要考虑使用循环以及 `input` 函数来控制交互，使用条件语句去控制循环

```

1 city = input("Choose city:")
2 CITY_DATA = { 'chicago': 'chicago.csv',
3               'new york city': 'new_york_city.csv',
4               'washington': 'washington.csv' }
5
6 while True:
7     if city in CITY_DATA:
8         break
9     else:
10        print("Your value is valide")
11        city = input("Choose city:")

```

使用函数进行优化：

```

1 def option(conditions, tip_message, error_message):
2     value = input(tip_message)
3
4     while True:
5         if value in conditions:
6             break
7         else:
8             print(error_message)
9             value = input(tip_message)
10    return value
11
12
13 city = option(CITY_DATA, "Choose city:", "Your value is valide")
14 month = option(MONTHS, "Choose month:", "Your value is valide")

```

2. load_data 函数

需要根据前面 `get_filter` 函数得到的 city、month 和 day 来收集和筛选数据，其中还需要使用像一个相应的日期进行处理，输出相应的值

```

1 # 读取数据
2 df = pd.read_csv(file)
3
4 # 列表数据得到对应值—解决某个值在列表中的第一个索引值，在本例中列表是唯一值，所以可以筛选出对应的月份和周数的数字值
5 months = ["a", "b", "c", "d"]
6 month = 'a'
7 months.index(month)
8
9 # output
10 0
11
12 # 因为需要对应 dataframe 中对应的月份和周数，所以需要加 1

```

```

13 months.index(month) + 1
14 # output
15 1
16
17 # 数据筛选可以使用一般掩码筛选
18 df[df["month"] == 1]
19
20 # 对于多条件筛选的方式，月数是 1 周数是 3 的数据
21 df[(df["month"] == 1) and (df["day"] == 3)]

```

3. 其他函数说明

- 针对最常见的数值处理——包括最常见月份、周数、起始时间、起始和终止站点，需要筛选出数据中的众数

```

1 # 针对筛选对应列的众数
2 a = pd.DataFrame({'a': ['a', 'b', 'a', 'c', 'b', 'a', 'b', 'c', 'd',
3     'a'], 'b': ['c', 'c', 'd', 'a', 'b', 'a', 'a', 'b', 'a',
4     'd']})
5 a
6 # output
7
8     a  b
9 0  a  c
10 1  b  c
11 2  a  d
12 3  c  a
13 4  b  b
14 5  a  a
15 6  b  a
16 7  c  b
17 8  d  a
18 9  a  d
19
20 # 筛选 a 中 a 列的众数
21 In [23]: a['a'].mode()
22 Out[23]:
23 0      a
24 dtype: object
25
26 # 可以筛选出对应的行值
27 In [24]: a['a'].mode()[0]
28 Out[24]: 'a'

```

- 常见的起点站和终点站组合，这里的语句已经说明了是需要计算频繁的车程，可以想一下现实情况中怎么表现频繁线段——是否可以考虑为计算线段的两头端点的组合频繁。假设上海，北京，成都以及广州等城市之间的频繁车程：我们可以考虑为统计起始点和终点 北京-成都，北京-上海，成都-广州，成都-北京，成都-上海 等组合方式的统计数量就可以计算出来了。这里的解决方式，可以将两个端点数据构造为一个新的列再计算众数

```
1 In [15]: a['c'] = a['a'] + a['b']
2
3 In [16]: a
4 Out[16]:
5      a  b  c
6 0  a  c  ac
7 1  b  c  bc
8 2  a  d  ad
9 3  c  a  ca
10 4  b  b  bb
11 5  a  a  aa
12 6  b  a  ba
13 7  c  b  cb
14 8  d  a  da
15 9  a  d  ad
16
17 In [20]: a['c'].mode()[0]
18 Out[20]: 'ad'
```

- 总旅行和平均旅行时间，这个是需要通过聚合计算，也就是使用 `Series` 的 `sum` 和 `mean` 方法

```
1 In [28]: b = pd.Series([1,2, 3, 4, 5,6], name="trip")
2
3 In [29]: b.mean()
4 Out[29]: 3.5
5
6 In [30]: b.sum()
7 Out[30]: 21
```

- 分别统计用户信息的数量——用户的类型和用户性别的统计量，可以使用 `Series` 的 `value_counts` 方法进行统计

```
1 In [26]: a["a"].value_counts()
2 Out[26]:
3 a      4
4 b      3
5 c      2
6 d      1
7 Name: a, dtype: int64
```

