

实训4：日志文件解析与分析

• ***目标***：掌握正则表达式和数据聚合。

• ***任务***：解析服务器日志，统计IP和URL访问频率。

• ***详细要求***：

o 读取日志文件（server.log），格式如192.168.1.1 - [2025-05-01 10:00:00] "GET /page1.html ...".

o 使用正则表达式提取IP、时间、URL。

o 统计每个IP的访问次数和每个URL的访问频率。

o 输出前10高频IP和URL到log_stats.txt，格式为“IP/URL: 次数”。

o 验证日志格式，跳过无效行并记录日志。

o 添加时间过滤功能（如仅分析某天的数据）。

o 保存分析结果，支持追加写入。

技能：正则表达式、Counter、文件I/O、日志记录。

```
1  # -*- coding: utf-8 -*-
2  import re
3  from collections import Counter
4  from datetime import datetime
5  import logging
6
7  # 设置日志记录
8  logging.basicConfig(filename='log_parser_debug.log', level=logging.INFO,
9                      format='%(asctime)s - %(message)s')
10
11 # 正则模式：提取 IP、时间、URL
12 LOG_PATTERN = re.compile(r'(?P<ip>\d+\.\d+\.\d+\.\d+) - \[(?P<time>
13 [\d\-:\s]+)\] "\w+ (?P<url>/[^\s]*)')
14
15 def parse_log(file_path, filter_date=None):
16     ip_counter = Counter()
17     url_counter = Counter()
18     total_lines = 0
19     valid_lines = 0
20
21     with open(file_path, 'r', encoding='utf-8') as f:
22         for line in f:
23             total_lines += 1
24             match = LOG_PATTERN.search(line)
25             if not match:
26                 logging.warning(f"无效日志行: {line.strip()}")
27                 continue
28
29             ip = match.group('ip')
30             time_str = match.group('time')
31             url = match.group('url')
32
33             # 时间过滤
34             if filter_date:
```

```
33         log_date = datetime.strptime(time_str, '%Y-%m-%d
%H:%M:%S').date()
34         if str(log_date) != filter_date:
35             continue
36
37         ip_counter[ip] += 1
38         url_counter[url] += 1
39         valid_lines += 1
40
41     logging.info(f"处理完成, 共 {total_lines} 行, 有效日志行 {valid_lines} 行")
42     return ip_counter.most_common(10), url_counter.most_common(10)
43
44 def save_stats(ip_stats, url_stats, output_file):
45     with open(output_file, 'a', encoding='utf-8') as f:
46         f.write('Top 10 IPs:\n')
47         for ip, count in ip_stats:
48             f.write(f'{ip}: {count}\n')
49         f.write('\nTop 10 URLs:\n')
50         for url, count in url_stats:
51             f.write(f'{url}: {count}\n')
52         f.write('\n')
53
54 if __name__ == "__main__":
55     log_file = 'server.log'
56     output_file = 'log_stats.txt'
57     date_filter = input("请输入要分析的日期 (格式 YYYY-MM-DD), 留空表示不过滤:
)").strip()
58     if date_filter == '':
59         date_filter = None
60
61     ip_stats, url_stats = parse_log(log_file, date_filter)
62     save_stats(ip_stats, url_stats, output_file)
63     print("分析完成, 结果已保存至 log_stats.txt")
64
```