

GitHub热门项目综合分析报告

乔奕博 10242140474

1. 定义问题

1.1 分析背景 ✨

GitHub作为全球最大的代码托管平台，其项目趋势反映了技术发展的动态。

2015-2024年间，GitHub见证了传统开发工具的成熟和AI、云原生等技术的兴起；

2025年，随着AI Agents、MCP生态系统、音视频生成等技术的快速发展，GitHub生态再次发生显著变化。

本项目旨在结合历史数据和最新趋势，全面分析GitHub平台上项目的演变规律，为开发者和技术决策者提供参考。

1.2 问题陈述 💡

核心问题：如何准确分析GitHub项目的长期发展趋势和短期热点变化，以及两者之间的关联？

具体问题：

- 从2015年开始的过去十年间，经典项目的影响力格局如何？
- 2025年GitHub上最热门的技术领域是什么？
- 如何定义和识别不同时期的'流行项目'？
- 月度项目增长趋势在不同时期有何差异？
- 不同技术领域的项目表现随时间如何变化？
-

1.3 分析目标 🔍

- 识别过去十年间的经典项目和2025年新兴热点项目
- 分析跨时期项目增长趋势和月度变化
- 定义科学的“流行项目”标准
- 生成直观的可视化图表和综合报告
- 构建可复用的数据分析流程，支持多时期数据整合

2. 获取数据

2.1 数据来源 📄

本项目使用了两种主要数据来源：

数据类型	来源	路径	包含内容
外部指标	"OpenRank杯" OpenSODA 开放数据集	top_300_metrics	全球Top 300项目的详细指标
项目类型分析	分析脚本生成	final_report/ project_type_analysis.csv	项目类型分布和指标
月度趋势数据	GitHub爬虫	github_new_trend/data_2025/ monthly_trends_*.json	2025年1-12月的项目趋势数据
年度报告	分析脚本生成	github_new_trend/data_2025/ github_2025_yearly_report.json	2025年年度分析报告

2.2 数据获取方法 📁

- **GitHub爬虫**：使用 github_2025_yearly_collector.py 脚本，通过GitHub API 获取2025年每月的热门项目数据，包括项目基本信息、星数、语言、topics等
- **外部数据集成**：设计了 utils/config.py 配置管理工具，支持从外部目录加载 top_300_metrics开源数据集
- **数据整合**：将月度数据汇总为年度报告，便于后续分析

2.3 数据规模

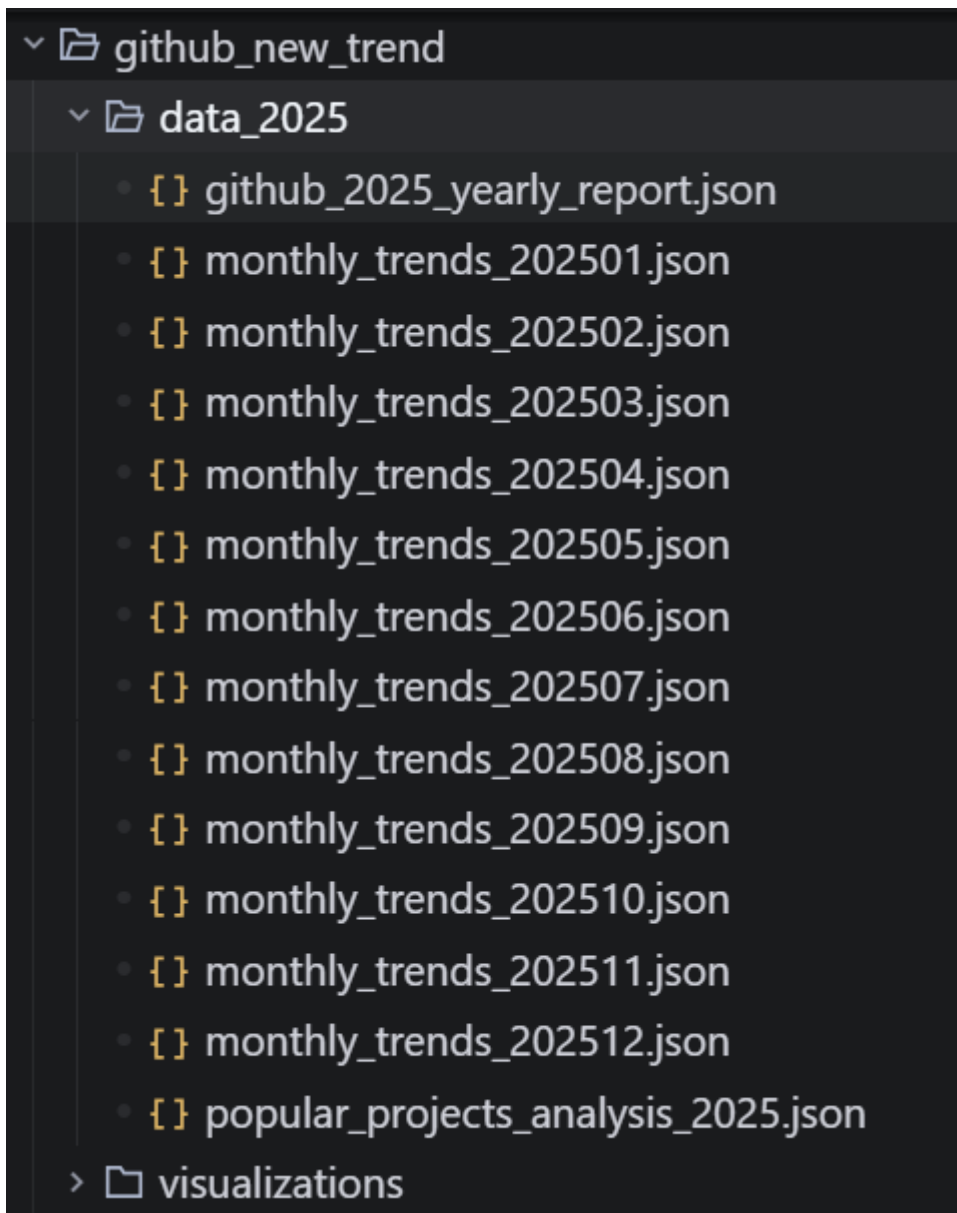
- 2015-2023年数据集（如下图）
 - 数据来源: GitHub Top 300 Metrics
 - 项目数量: 121个 (2015-2023)
 - 组织数量: 91个
 - 主要指标: OpenRank、月度新增Stars、活跃度



OpenSODA 挑战赛开放数据集正式发布！

随着更多的队伍开始通过“[初赛知识问答](#)”进入复赛，为了帮助各参赛队伍更好的设计复赛作品方案，我们特为此整理一批覆盖 2 类赛题（6 类子赛题）的开放数据集及其相关说明，供大家使用。后续，组委会还会继续为进入决赛阶段的队伍，提供更加系统完整的数据支持服务。希望大家取得好的成绩！

- 2025年爬取数据
 - 数据来源: GitHub API爬取
 - 爬取时间: 2026年1月
 - 月度数据: 12个月



3. 准备数据

3.1 数据清洗 ☒

- 处理不同数据类型：实现了 `get_stars_count` 函数，安全处理数字、字符串等不同类型的星数数据
- 修复数据格式：统一JSON和CSV数据的字段命名和格式
- 处理缺失值：对于缺失的描述、语言等字段，设置默认值或忽略
- 移除重复数据：检查并移除重复的项目记录

3.2 数据转换

- 统一时间格式：将月度数据的时间格式统一，便于时间序列分析
- 标准化字段名称：将不同来源的数据字段名称标准化，便于整合

3.3 数据验证

- 检查数据完整性：验证月度数据是否包含所有必要字段
 - 验证数据一致性：检查同一项目在不同月份的数据是否一致
 - 异常值处理：识别并处理异常高的星数或增长率数据，分析其原因
-

4. 探索数据

4.1 数据探索

- 统计组织数量和项目分布
- 查看Microsoft等典型组织的项目
- 分析VSCode等热门项目的指标文件结构

基于OpenRank数据集的2015-2024年项目分析

项目影响力排名:

- 第一名: microsoft/vscode (OpenRank: 1954.1)
- 第二名: pytorch/pytorch (OpenRank: 1328.6)
- 第三名: microsoft/winget-pkgs (OpenRank: 1054.3)

核心指标统计:

- 平均OpenRank: 441.9
- 最大OpenRank: 1954.1

2025年项目趋势分析

- 月度项目数量: 每月约100个热门项目

2025年语言统计

- Python: 420个项目, 平均Stars: 5035.4
- TypeScript: 263个项目, 平均Stars: 4958.2

- JavaScript: 88个项目, 平均Stars: 3908.8
- Go: 60个项目, 平均Stars: 5254.1
- Rust: 43个项目, 平均Stars: 5074.1

2025年热门项目:

- system-prompts-and-models-of-ai-tools (105,627 stars)
- DeepSeek-R1 (91,642 stars)
- gemini-cli (89,677 stars)

4.2 趋势分析 🔮

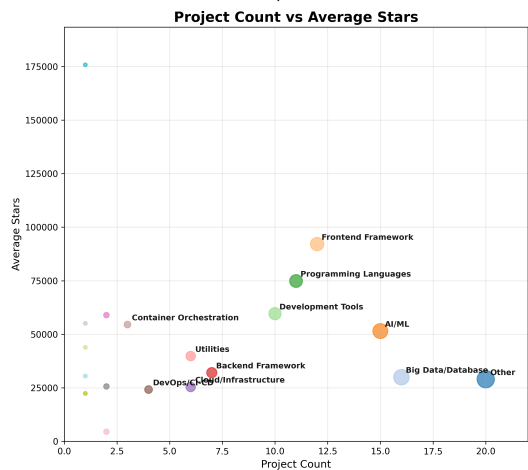
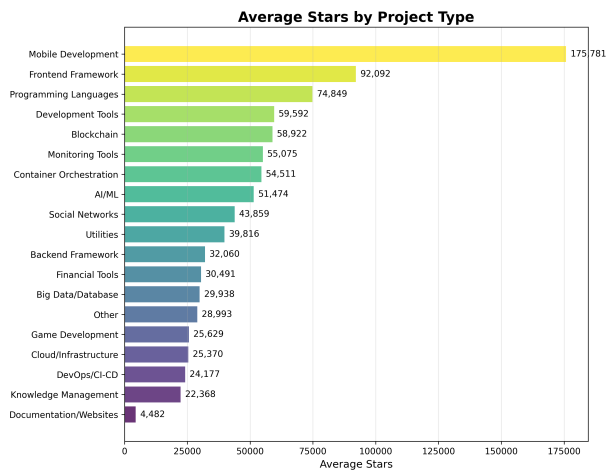
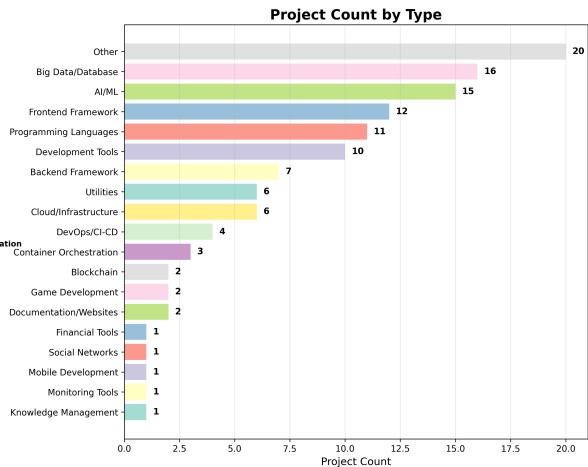
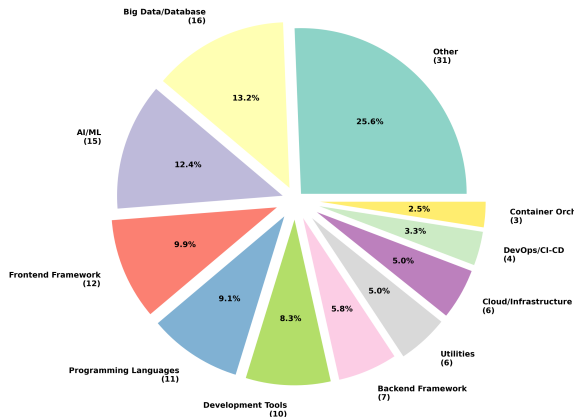
- 月度趋势：分析每月流行项目数量、总星数、平均星数的变化
- 技术领域分布：统计不同编程语言和技术领域的项目数量和星数
- 项目增长模式：识别快速增长的项目和稳定发展的项目

4.3 可视化探索 📊

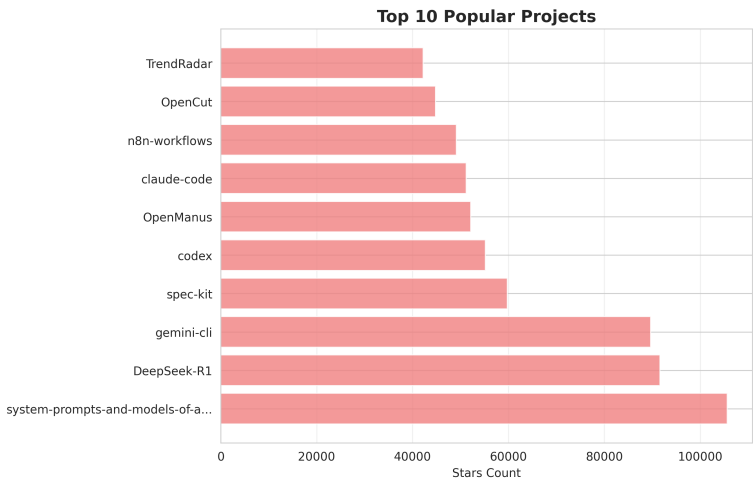
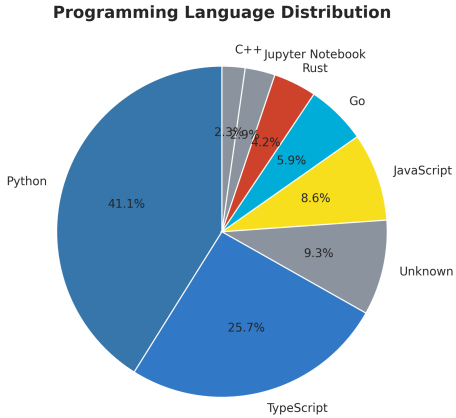
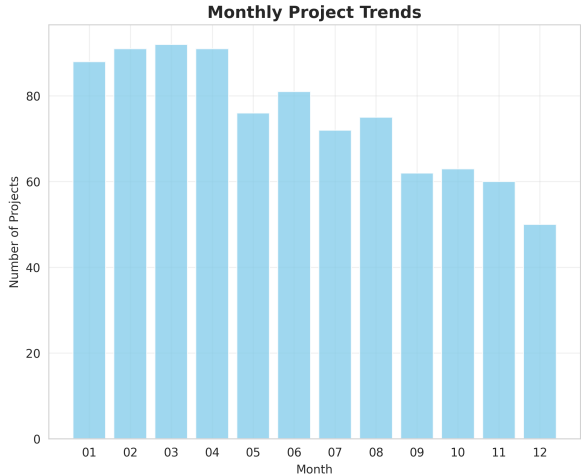
- 饼图：展示项目类型分布
- 柱状图：展示热门项目排名、技术领域比较
- 折线图：展示月度趋势变化

GitHub Project Type Analysis Dashboard

Project Type Distribution (Simplified View - Major Categories Only)



GitHub 2025 Project Trends Analysis Dashboard



Statistics Summary

Data Collection Statistics:

Total Projects: 1,200
Total Stars: 5,800,073
Programming Languages: 37
Covered Months: 12

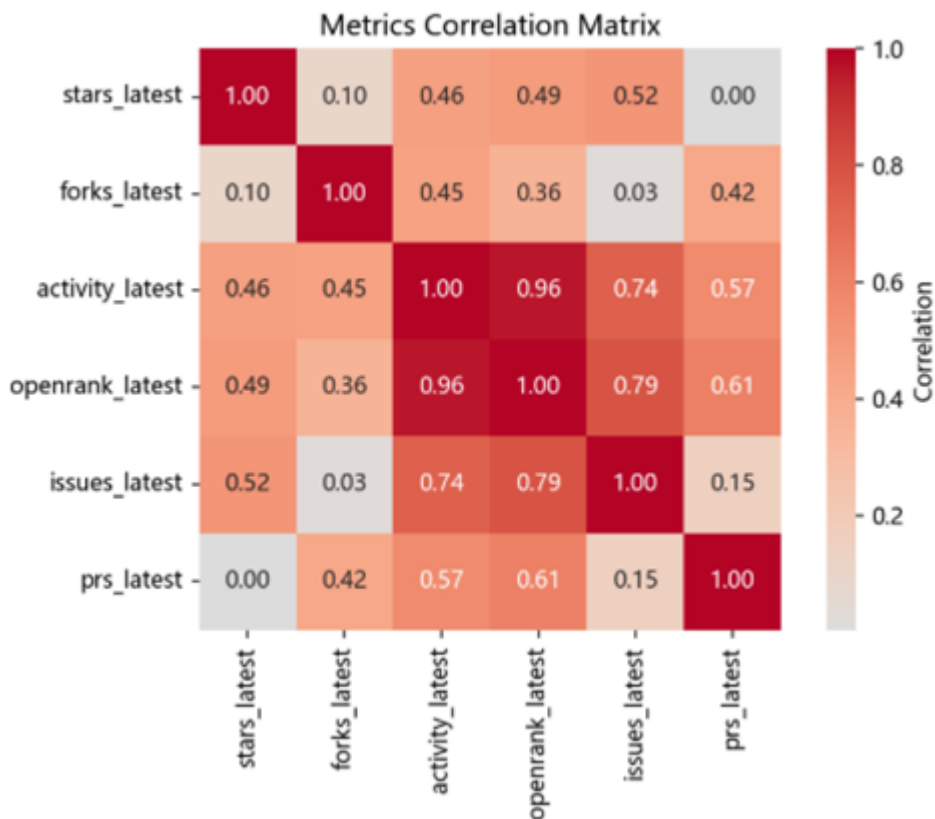
Highest Starred Project:
system-prompts-and-models-of-ai-tools: 105,627 stars

Data Time Range:
Jan 2025 - Dec 2025
+ Jan 2026 Latest Trends

Data Collection Success Rate:
12 Success
0 Failed

4.4 建模分析

OpenRank指标与项目Stars数,activity指标的相关性——使用Pandas的 .corr() 函数计算 反应相关度的系数



5. 交付与迭代

5.1 交付内容 📁

- 报告文件：生成 executive_summary.txt 、 growth_analysis.csv 等报告
- 可视化图表：生成饼图、柱状图、折线图等可视化结果
- 分析脚本：提供可复现的分析脚本和工具

5.2 问题与迭代 📋

多次优化迭代：

- 中文显示问题：修改字体设置，确保仪表板上的中文正确显示
- 数据路径问题：修正脚本中的数据目录路径，确保能找到有效数据
- 精细化项目领域分类：增加分类指标，扩充分类关键词，优化饼图显示
- 更新爬虫：增加对2025年数据多维度爬取

5.3 迭代效果 🎓

- 修复了所有已知错误

- 提高了分析效率和准确性
 - 增强了代码的可维护性和可扩展性
-

6. 结论与展望

6.1 主要结论 🔍

根据GitHub近十年热门项目的分析，得到以下一般性结论

- **OpenRank是评估项目影响力的最佳指标：**与项目活跃度高度相关（相关系数0.964）比单看Stars或活跃度更能全面反映项目在开源社区的综合影响力
- **组织影响力与项目数量和质量密切相关：**Microsoft以6802.69的总影响力位居首位，主要得益于12个高影响力项目（Top40项目的统计数据）- 相比之下，Google虽有高影响力项目，但项目数量较少，顶级组织往往拥有多个相互补充的高质量项目，形成生态系统效应
- **开源项目的成功模式：**顶级项目往往在多个指标上都表现优秀，形成良性循环，不同类型的项目（如开发工具、编程语言、框架）在各指标上的表现模式不同，持续的开发活跃度和社区参与是维持高影响力的关键因素

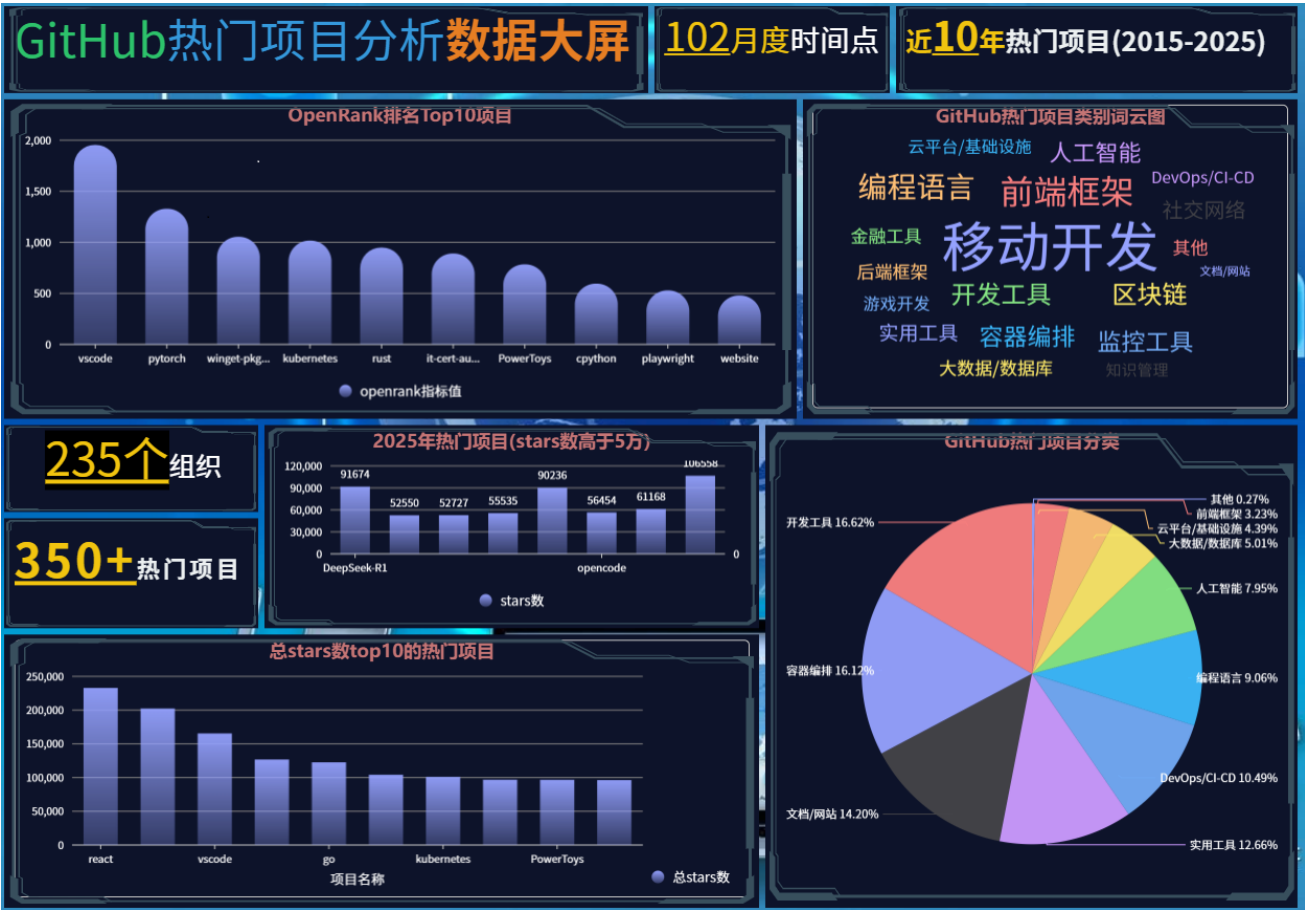
分析2025年GitHub项目的分析，有以下新趋势

- **AI Agents与MCP生态系统成为主导趋势：**2025年AI Agents项目（如DeepSeek-R1、sim、open-r1）占据GitHub热门项目榜前列，MCP（Model Context Protocol）相关项目（如context7、github-mcp-server）快速发展，形成完整的AI代理生态系统，与传统项目相比，AI Agents项目的Stars增长速度更快，社区参与度更高
- **音视频生成技术全面爆发：**2025年文本到语音（如index-tts、csm）、视频生成（如Wan2.1）技术成为新热点，这些项目专注于生成高质量、可控的音视频内容，满足多样化的应用场景需求，持续的模型优化和社区贡献是维持高影响力的关键
- **开源AI模型普及：**OpenAI、Google等公司纷纷开源大型语言模型和AI工具（如gpt-oss、qwen-code），推动了AI技术的民主化发展，开源模型的出现降低了AI应用开发的门槛，促进了AI技术的广泛应用

6.2 项目亮点 🚀

- 完整的数据科学流程：从问题定义到交付迭代
- 模块化的代码架构：便于维护和扩展
- 灵活的配置管理：支持不同环境和数据来源
- 健壮的数据处理：能处理不同格式和质量的数据

6.4 数据大屏 📊



7. 技术栈与工具

类别	工具/库	用途
数据处理	pandas, numpy	数据清洗、转换、分析
可视化	matplotlib, seaborn	生成统计图表
配置管理	json, pathlib	管理配置文件和路径
爬虫	GitHub API	收集GitHub数据
开发语言	Python 3.13	主要开发语言
项目管理	config.json, requirements.txt	依赖管理和配置

8. 代码结构

根目录:

文件/文件夹	功能
original_data	OpenRank开源数据集
data_2025	2025年热门项目数据, 编程语言数据, 月度趋势数据
scripts	清洗, 处理, 可视化2024年以前GitHub热门项目数据
github_new_trend	收集, 分析, 可视化2025年GitHub热门项目数据
output	记录统计结果的csv文件, 生成的图表与仪表板
utils	配置管理工具, 处理路径和配置项
config.json	项目配置文件
README.md	项目简要介绍
requirements.txt	该项目的Python依赖列表

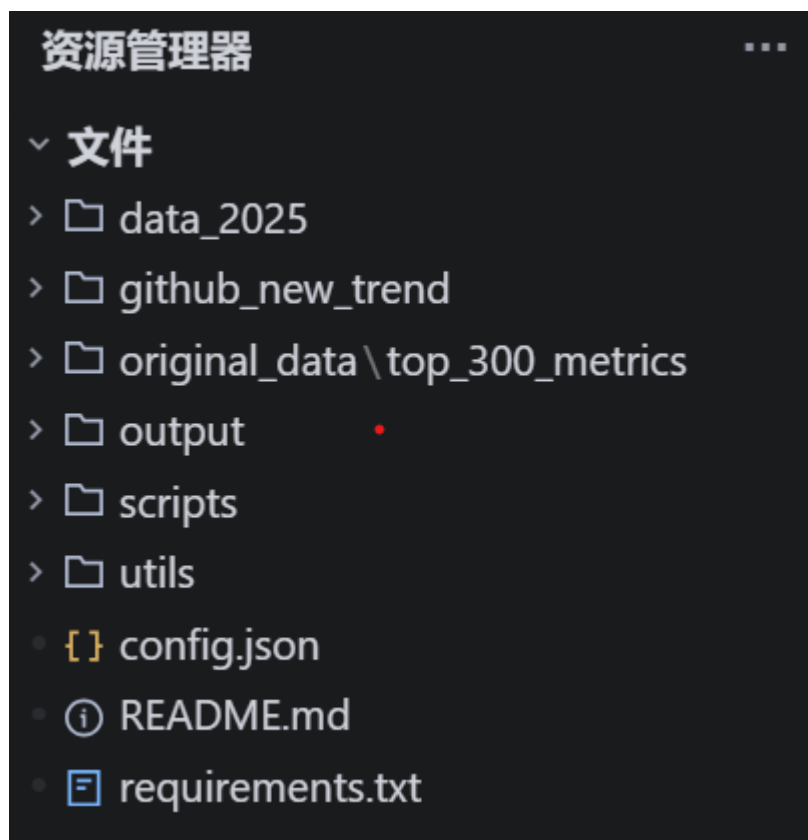
output/-目录: 分析结果输出

文件/文件夹	功能
2015-2024_charts	2015-2024年传统热门项目分析
organization_analysis	热门项目所在公司/组织分析
project_type_summary	热门项目领域分析
top_50_reports	按不同指标筛选top50项目
visualizations_2025	2025最新热门项目可视化图表
complete_total_stars_ranking.csv	完整的总stars数排名
metrics_correlation.csv	各指标关联度分析

github_new_trend/-目录

文件名	功能
github_2025_yearly_collector.py	从github上爬取2025年数据
github_popular_projects_analyzer.py	数据处理分析
github_data_visualizer.py	数据可视化

文件根目录如下：



9. 总结

本项目完整实现了数据科学的六个步骤，从定义问题到交付迭代，构建了一个健壮的GitHub项目分析系统。通过多次迭代优化，最终生成了近10年GitHub热门项目统计报告与2025年GitHub项目最新趋势分析报告。

项目的亮点在于模块化的设计、灵活的配置管理和健壮的数据处理能力，能够适应不同的数据来源和分析需求。未来可以进一步扩展功能，引入更先进的分析方法和可视化工具，提供更深入的技术趋势洞察。