

项目：可视化电影数据

第一步：清理数据和选择变量

- 清理数据：
 - 1、将数据文件读入 python，删除杂行数据
`movies = movies[movies.iloc[:,21].isnull() & movies.iloc[:,22].isnull()]`
 - 2、去掉重复行
`movies = movies.drop_duplicates()`
`movies.to_csv('movies1.csv')`
 - 3、在表格中通过 `release_year` 筛选非年份的数据行，去掉错乱行
 - 4、通过 excel 分列功能将 `genres` 列分列
- 选择变量：

根据需要解答的问题，选定了'id', 'original_title', 'keywords', 'genres', 'production_companies', 'release_year', 'vote_count', 'vote_average', 'budget_adj', 'revenue_adj' 这些变量，删除其他列，将数据文件导入 Tableau。

第二步：问题

- 回答下列问题，引用你在线可视化结果去支持你的答案：
 - 问题 1：电影类型是如何随着时代变化而变化的？

视图显示：

1. 喜剧、剧情、惊悚、动作类题材随时代变化增长很快，战争、西部等题材电影几乎没有增长。
2. 喜剧、剧情、动作题材始终排在前几名，依然是最常见的电影题材

- 问题 2：环球影业和派拉蒙影业的电影之间数据指标有什么区别？

视图显示：

1. 随着时间的推移，环球影业和派拉蒙影业制作的电影数均呈增长趋势，1988 年之后，环球影业增长量超过派拉蒙影业
2. 收入和成本具有一定的正相关，大成本电影 Titanic 收入高达 25 亿，环球制作的 E.T.和 Jaws 小成本高票房收入，性价比高。

- 问题 3：和非小说改编的电影相比，基于小说改编的电影表现得怎么样？

从 1960 年到 2015 年，上映的电影总量增长迅速，IP 改编电影数增长较慢，IP 改编电影占比随时间降低。

从类型分布上看，改编电影以惊悚、浪漫、奇幻、剧情题材为主；大体来看，电影预算与电影收入呈正相关，IP 改编对电影收入助益不大。

○ 问题 4：叫好又叫座的电影与叫好不叫座的电影有什么区别？

我提出这个问题，是因为想知道有些电影能够得到高评分，票房却并不乐观的原因。

从可视化图中可以看出：

1. 评分人数与票房有一定联系，票房不高，很可能是因为传播度不够。
2. 题材与票房也有一定关系，像冒险、动作、科幻等题材就非常受影院用户的欢迎。

当然，还有很多因素可能影响电影叫好却不叫座，譬如演员人气、前期宣发，这部分还需要其他数据补充，才能进一步研究。

第三步：可视化

链接：

https://public.tableau.com/views/movies2/Q4Story?:embed=y&:display_count=yes&publish=yes