



数据分析方向

Python 完整学习路线

从零基础到独立分析的全流程指南

零基础友好

8-12周

面试就绪

适用人群：零基础或有少量编程经验的转行者

总时长预估：8-12 周（每天 2-3 小时业余学习）

最终目标：独立完成数据提取、清洗、分析、可视化全流程

目 录

概述

阶段零：心态与环境准备（1-2 天）

阶段一：Python 核心语法（10-14 天）

阶段二：统计学基础思维（5-7 天）

阶段三：SQL 数据提取（7-10 天）

阶段四：数据处理核心库（14-18 天）

阶段五：数据可视化（5-7 天）

阶段六：实战项目（10-14 天）

阶段七：进阶补充（持续学习）

学习节奏建议

数据分析方向 Python 完整学习路线

适用人群：零基础或有少量编程经验的转行者

总时长预估：8-12 周（按每天 2-3 小时业余学习计算）

最终目标：能独立完成数据提取、清洗、分析、可视化全流程，具备面试数据分析岗位的基本能力

阶段零：心态与环境准备（1-2 天）

认知校准

在开始前，请理解数据分析师的工作本质是用数据回答业务问题，而非单纯写代码。Python 只是工具，思维方式才是核心竞争力。

环境搭建

- Mac 用户：安装 Miniconda（比 Anaconda 轻量），创建独立的虚拟环境
- 编辑器：Jupyter Notebook（探索分析）+ VS Code（写脚本）
- 版本管理：安装 Git，注册 GitHub 账号（用于存放学习项目，也是简历加分项）

学习资源原则

- 官方文档是最权威的参考，学会查阅（Pandas、NumPy 文档写得很好）
- 遇到报错先读错误信息，再搜索 Stack Overflow
- 不要囤课，选定一套教程跟到底

阶段一：Python 核心语法（10-14 天）

目标

能读懂基础代码，能独立写出解决简单问题的脚本。

核心内容

知识点	重点掌握	优先级
变量与数据类型	字符串、数字、列表、字典、元组	★★★
流程控制	if-else、for 循环、while 循环	★★★
函数	def 定义函数、参数传递、返回值	★★★
列表推导式	[<code>x</code> for <code>x</code> in <code>list</code> if <code>condition</code>]	★★★
文件操作	读写 txt、csv 文件	★★
异常处理	try-except 基本用法	★★
模块导入	import 语法、pip 安装第三方库	★★

阶段性练习

- 写一个猜数字小游戏
- 读取一个 txt 文件，统计词频
- 用字典实现一个简单的通讯录

推荐资源

- 《Python Crash Course》前 10 章
- 廖雪峰 Python 教程（中文免费）

阶段二：统计学基础思维（5-7 天）

为什么需要这个阶段？

很多人跳过统计直接学工具，结果会用 `groupby` 却不知道什么时候该用、算出来的数字代表什么。

核心内容

描述性统计

- 集中趋势：均值（mean）、中位数（median）、众数（mode）
- 离散程度：标准差（std）、方差（var）、四分位距（IQR）
- 分布形态：偏度、峰度、正态分布概念

关系与比较

- 相关系数（Pearson、Spearman）及其局限性
- 相关性 ≠ 因果性（这个认知极其重要）

业务常用指标

- 同比、环比怎么算
- 转化率、留存率、复购率的定义
- 什么是 A/B 测试、为什么需要对照组

阶段性练习

- 拿到一组数据，能口述其分布特征
- 解释“某功能上线后转化率提升 10%”这个结论是否可靠

推荐资源

- 《深入浅出统计学》 (Head First Statistics)
- 可汗学院统计学课程 (免费)

阶段三：SQL 数据提取（7-10 天）

为什么必须学？

真实工作中，数据存在数据库里，你需要先用 SQL 把数据“捞”出来，才能用 Python 分析。面试时 SQL 也是必考项。

核心内容

知识点	说明	优先级
SELECT 基础	选择列、WHERE 条件筛选、ORDER BY 排序	★★★
聚合函数	COUNT、SUM、AVG、MAX、MIN	★★★
GROUP BY	分组聚合（对应 Pandas 的 groupby）	★★★
JOIN	INNER JOIN、LEFT JOIN（核心难点）	★★★
子查询	嵌套查询、WITH 语句	★★
窗口函数	ROW_NUMBER、RANK、LAG/LEAD	★★

阶段性练习

- 在 LeetCode 或 HackerRank 刷 30 道 SQL 题
- 用 SQL 写出“每个用户的首单时间”、“连续登录 3 天的用户”

推荐资源

- Mode Analytics SQL 教程（免费在线练习环境）
- 《SQL 必知必会》（薄且实用）

阶段四：数据处理核心库（14-18 天）

目标

这是数据分析的主战场，投入最多时间是值得的。

1. NumPy（3-4 天）

快速理解即可，不用深钻：

- ndarray 多维数组概念
- 向量化运算（避免 for 循环）
- 基本索引与切片
- 常用函数：`np.mean()`、`np.sum()`、`np.where()`

2. Pandas（10-14 天） ★核心中的核心

数据结构

- Series（一维）与 DataFrame（二维）
- 索引（Index）的概念与重要性

数据读写

- `pd.read_csv()`、`pd.read_excel()`、`pd.read_sql()`
- `df.to_csv()`、`df.to_excel()`

数据清洗

操作	方法	使用场景
查看数据概况	<code>.info()</code> 、 <code>.describe()</code> 、 <code>.head()</code>	拿到数据第一步
缺失值处理	<code>.isna()</code> 、 <code>.dropna()</code> 、 <code>.fillna()</code>	几乎每个项目都用
重复值处理	<code>.duplicated()</code> 、 <code>.drop_duplicates()</code>	数据去重
类型转换	<code>.astype()</code> 、 <code>pd.to_datetime()</code>	日期、数值格式问题
字符串处理	<code>.str.contains()</code> 、 <code>.str.split()</code>	文本字段清洗

数据操作

操作	方法	重要性
筛选	<code>.loc[]</code> 、 <code>.iloc[]</code> 、布尔索引	★★★
排序	<code>.sort_values()</code> 、 <code>.sort_index()</code>	★★
分组聚合	<code>.groupby()</code> + 聚合函数	★★★
透视表	<code>.pivot_table()</code>	★★★
合并	<code>.merge()</code> (类似 SQL JOIN)、 <code>.concat()</code>	★★★
新增列	<code>.apply()</code> 、 <code>.assign()</code>	★★★
时间序列	日期索引、resample、时间差计算	★★

阶段性练习

- 拿到任意 CSV 文件，能在 10 分钟内完成数据概况了解
- 用 Pandas 复现你之前写的 SQL 查询
- 完成一个完整的数据清洗流程

阶段五：数据可视化（5-7 天）

目标

让数据说话，掌握从“能出图”到“出好图”的能力。

1. Matplotlib (2 天)

- 理解底层逻辑：Figure → Axes → Plot
- 能画出基础图表：折线图、柱状图、散点图
- 会调整标题、标签、图例、颜色

2. Seaborn (2-3 天)

- 统计图表：`histplot` (分布)、`boxplot` (箱线图)、`heatmap` (热力图)
- 关系图表：`scatterplot`、`pairplot`
- 学会用 `hue` 参数做分组对比

3. 交互式图表 (1-2 天)

- Plotly Express：快速生成可交互图表
- 了解 PyEcharts (国内业务场景常用)

可视化原则

- 先想清楚要表达什么，再选图表类型
- 比较用柱状图、趋势用折线图、分布用直方图、关系用散点图
- 少即是多，避免 3D 图表和过多颜色

阶段六：实战项目（10-14 天）

目标

把前面学的串起来，形成可展示的作品集。

项目一：电商用户行为分析（初级）

数据源：阿里天池、Kaggle 公开数据集

分析内容：

- 用户活跃时段分布
- 转化漏斗分析（浏览→加购→下单→支付）
- RFM 用户分层
- 复购率计算

交付物：Jupyter Notebook + 分析结论

项目二：销售数据仪表盘（中级）

使用工具：Streamlit 或 Dash

功能：

- 上传 Excel 文件
- 自动生成销售趋势图、TOP10 商品、区域分布
- 支持筛选时间范围

交付物：可运行的 Web 应用，部署到 Streamlit Cloud

项目三：自选主题深度分析（进阶）

选择你感兴趣的领域：

- MBTI 性格分析
- 电影/音乐/游戏数据
- 公开的社会经济数据

要求：

- 自己定义分析问题
- 完成数据获取、清洗、分析、可视化全流程
- 撰写一份完整的分析报告（背景、方法、发现、建议）

阶段七：进阶补充（持续学习）

视时间和需求选修

方向	内容	适用场景
数据获取	requests + BeautifulSoup 爬虫基础	需要自己采集数据时
自动化	Python 脚本定时执行、自动发邮件	重复性工作自动化
机器学习入门	scikit-learn 基础（分类、回归、聚类）	进阶分析师/转算法
BI 工具	Tableau / Power BI 基础	部分公司要求

学习节奏建议

周一至周五：学习新知识（每天 1-2 小时）

周末：动手练习 + 复盘（每天 3-4 小时）

每完成一个阶段：做一个小项目巩固

遇到卡壳超过 30 分钟：先跳过，标记后回来

祝学习顺利！记住：数据分析师的核心竞争力是思维方式，Python 只是工具。