

数据获取 → 数据整理 → 对比分析
关联关系
分析建模

数据聚合 `df.groupby('star').sum()`

创建新列 `df['new_star'] = df['star'].map(star_to_number)`

多列排序 `df.sort_values(by=['col1', 'col2'], ascending=[T, F])` 深度学习数据集 [sklearn]
国家统计局

表的横向连接 `df.merge(data1, data2, on='group', how='inner')`

连接键类型, 解决没有公共列问题

`pd.merge(data1, data2, left_index='a', right_index='b', on=)`
指定连接键:

连接成: 默认 inner 取两表公共部分
left 以左表为基础, 右表往左表拼接
right 以右表为基础, 左表往右表拼接
outer 取两表并集

纵向拼接: `pd.concat([data1, data2])` 相同列
字典

字符串
--str-- | --getitem-- | --setitem-- | --delitem--
--iter-- | --call-- | --eq-- | --gt-- | --ge--
--add-- | --sub-- | --hash--
数据预处理 pandas
Python 与 pandas 交互可用魔法

聚合: `df2.groupby('type').aggaggregate` 多种多样的聚合
(['type': 'count', 'Feb': 'sum'])
`data.groupby('group').transform('mean')`
取完之后写入到每一个值上面

数据透视表: `pd.pivot-table`
(csv) / to-pickle 速度快, 只有 pandas 兼容

数据导出: `df.to_excel(excel_writer='rfile.xlsx', sheet_name='sheet1', index=False, columns=['col1', 'col2'])`

缺失值处理: `na_rep=0 inf_rep=0`

数据集

数据结构

预处理

数据调整

数值运算

多表操作

数据导出

Series (tuple) Pd.series
底层: numpy
Dataframe

缺失值 Series: x.hasnans
Dataframe: isnull()

重复值 x.drop_duplicates()

索引

行列选择 `df[1:3]` 前3行

数值替换

数值删除 `df.drop(df.columns, axis=1)`

新行列插入

索引重塑

基本运算

汇总

数据分组

多表拼接

连接键

连接方式

导出类型

可视化 matplotlib, pyplot

连接数据库 pymysql

`import seaborn as sns`
`sns.set_style('darkgrid')`
`plt.scatter(df.index, df['A'])`
`plt.show()`