
API 参考

发行版本 v1.0

lihetong

2025 年 12 月 26 日

内容目录:

1 介绍	3
2 模块参考	5
2.1 data_loader	5
2.2 data_clean	6
2.3 data_explore	6
2.4 data_visualize	7
2.5 feature_engineer	7
2.6 model_train	7
2.7 model_evaluate	8
2.8 config	8
2.9 main	8
3 Indices and tables	9
4 项目介绍	11
4.1 功能特性	11
4.2 技术栈	11
Python 模块索引	13
索引	15

欢迎来到 EcommerceSalesAnalysis 项目文档!

这是一个电商平台销售数据分析项目，包含数据加载、清洗、探索、可视化、特征工程、模型训练和评估等模块。

CHAPTER 1

介绍

这是一个电商平台销售数据分析系统，主要功能包括：

- 数据加载与预处理
- 数据探索与分析
- 数据可视化
- 特征工程
- 模型训练与评估

CHAPTER 2

模块参考

```
data_loader  
data_clean  
data_explore  
data_visualize  
feature_engineer  
model_train  
model_evaluate  
config  
main
```

2.1 data_loader

Functions

<code>data_loader(data, information)</code>	
<code>get_basic_info(df)</code>	获取数据基本信息
<code>load_raw_data()</code>	加载原始用户数据 - 读取 CSV / 数据库 - 不做任何 清洗与修改 - 保证“原始性” :return: 返回加载好的 原始数据

`data_loader.data_loader(data, information) → DataFrame`

`data_loader.get_basic_info(df: DataFrame) → dict`

获取数据基本信息

:param df: 传入原始数据 :return: 返回部分数据

```
data_loader.load_raw_data() → DataFrame
```

加载原始用户数据 - 读取 CSV / 数据库 - 不做任何清洗与修改 - 保证“原始性” :return: 返回加载好的原始数据

2.2 data_clean

Functions

<code>clean_data(df, numeric_cols, categorical_cols)</code>	统一保存: return:
<code>handle_categorical_missing(df, categorical_cols)</code>	类别型特征缺失值处理: 填充为'Unknown'
<code>handle_missing_values(df, numeric_cols)</code>	数值型特征缺失值处理: - 使用 -1 进行占位填充 - 同时构造缺失指示变量: param df: 原始数据: param numeric_cols: 数值列: return:
<code>mark_abnormal_values(df)</code>	异常值标记, 不修改原始取值
<code>remove_duplicates(df_new)</code>	查询是否有重复项, 去除重复项: return:

```
data_clean.clean_data(df: DataFrame, numeric_cols: list, categorical_cols: list) → DataFrame
```

统一保存: return:

```
data_clean.handle_categorical_missing(df: DataFrame, categorical_cols: list) → DataFrame
```

类别型特征缺失值处理: 填充为'Unknown'

```
data_clean.handle_missing_values(df: DataFrame, numeric_cols: list) → DataFrame
```

数值型特征缺失值处理: - 使用 -1 进行占位填充 - 同时构造缺失指示变量: param df: 原始数据: param numeric_cols: 数值列: return:

```
data_clean.mark_abnormal_values(df: DataFrame) → DataFrame
```

异常值标记, 不修改原始取值

```
data_clean.remove_duplicates(df_new: DataFrame) → DataFrame
```

查询是否有重复项, 去除重复项: return:

2.3 data_explore

Functions

<code>categorical_by_lifecycle(df, col)</code>	分析不同用户中, 某类别特征的分布情况: param df: :param col: :return:
<code>explore_categorical_features(df)</code>	类别型特征分布: return: 返回用户各类别数量
<code>explore_missing_values(df)</code>	统计各字段缺失率: param df: :return:
<code>explore_numeric_features(df)</code>	数值型特征描述性统计: param df: 原始的数据: return: 返回数值型特征的描述性统计信息, 并转置结果 (T) 使特征作为行显示

```
data_explore.categorical_by_lifecycle(df, col)
```

分析不同用户中, 某类别特征的分布情况: param df: :param col: :return:

```
data_explore.explore_categorical_features(df: DataFrame) → dict
```

类别型特征分布: return: 返回用户各类别数量

`data_explore.explore_missing_values(df: DataFrame) → DataFrame`
 统计各字段缺失率: param df: :return:

`data_explore.explore_numeric_features(df: DataFrame) → DataFrame`
 数值型特征描述性统计: param df: 原始的数据: return: 返回数值型特征的描述性统计信息, 并转置结果
 (T) 使特征作为行显示

2.4 data_visualize

Functions

<code>plot_box_by_category(df, cat_col, num_col)</code>	类别-数值关系箱线图
<code>plot_categorical_by_lifecycle(df, col)</code>	
<code>plot_categorical_distribution(df, col)</code>	类别型特征分布
<code>plot_correlation_heatmap(df)</code>	数值特征相关性热力图
<code>plot_numeric_distribution(df, col)</code>	数值型特征分布

`data_visualize.plot_box_by_category(df: DataFrame, cat_col: str, num_col: str)`
 类别-数值关系箱线图

`data_visualize.plot_categorical_by_lifecycle(df, col)`

`data_visualize.plot_categorical_distribution(df: DataFrame, col: str)`
 类别型特征分布

`data_visualize.plot_correlation_heatmap(df: DataFrame)`
 数值特征相关性热力图

`data_visualize.plot_numeric_distribution(df: DataFrame, col: str)`
 数值型特征分布

2.5 feature_engineer

Functions

<code>build_features(df, categorical_cols)</code>	特征工程主入口
<code>one_hot_encode(df, categorical_cols)</code>	对类别型特征进行 One-Hot 编码只需要对类别列进行翻译

`feature_engineer.build_features(df: DataFrame, categorical_cols: list) → DataFrame`
 特征工程主入口

`feature_engineer.one_hot_encode(df: DataFrame, categorical_cols: list) → DataFrame`
 对类别型特征进行 One-Hot 编码只需要对类别列进行翻译

2.6 model_train

Functions

`train_model(df, target_col[, test_size, ...])`

模型训练模块:param df: 特征工程后的数据:param target_col: 目标列:return: 训练好的模型 + 测试集

model_train.**train_model**(*df*: DataFrame, *target_col*: str, *test_size*: float = 0.2, *random_state*: int = 123, *model_type*: str = 'rf')

模型训练模块:param df: 特征工程后的数据:param target_col: 目标列:return: 训练好的模型 + 测试集

2.7 model_evaluate

Functions

`evaluate_model(model, x_test, y_test)`

模型评估函数

model_evaluate.**evaluate_model**(*model*, *x_test*, *y_test*)

模型评估函数

2.8 config

2.9 main

Functions

`main()`

main.**main**()

CHAPTER 3

Indices and tables

- genindex
- modindex
- search

CHAPTER 4

项目介绍

EcommerceSalesAnalysis 是一个电商平台销售数据分析项目。该项目旨在通过数据科学的方法对电商平台的销售数据进行深入分析，提供有价值的业务洞察。

4.1 功能特性

- **数据加载:** 支持多种格式的销售数据加载
- **数据清洗:** 自动识别和处理数据中的异常值、缺失值
- **数据探索:** 提供丰富的统计分析功能
- **数据可视化:** 生成直观的图表展示数据特征
- **特征工程:** 提取有用的特征用于模型训练
- **模型训练:** 使用机器学习算法预测销售趋势
- **模型评估:** 评估模型性能并提供可视化结果

4.2 技术栈

- Python 3.x
- Pandas - 数据处理
- NumPy - 数值计算
- Matplotlib/Seaborn - 数据可视化
- Scikit-learn - 机器学习
- Sphinx - 文档生成

Python 模块索引

c

config, 8

d

data_clean, 6
data_explore, 6
data_loader, 5
data_visualize, 7

f

feature_engineer, 7

m

main, 8
model_evaluate, 8
model_train, 7

索引

B

build_features() (在 feature_engineer 模块中), 7

C

categorical_by_lifecycle() (在 data_explore 模块中), 6

clean_data() (在 data_clean 模块中), 6

config

module, 8

D

data_clean
 module, 6

data_explore
 module, 6

data_loader
 module, 5

data_loader() (在 data_loader 模块中), 5

data_visualize
 module, 7

E

evaluate_model() (在 model_evaluate 模块中), 8

explore_categorical_features() (在 data_explore 模块中), 6

explore_missing_values() (在 data_explore 模块中), 6

explore_numeric_features() (在 data_explore 模块中), 7

F

feature_engineer
 module, 7

G

get_basic_info() (在 data_loader 模块中), 5

H

handle_categorical_missing() (在 data_clean 模块中), 6

handle_missing_values() (在 data_clean 模块中), 6

L

load_raw_data() (在 data_loader 模块中), 5

M

main

 module, 8
main() (在 main 模块中), 8
mark_abnormal_values() (在 data_clean 模块中), 6

model_evaluate
 module, 8

model_train
 module, 7

module
 config, 8
 data_clean, 6
 data_explore, 6
 data_loader, 5
 data_visualize, 7
 feature_engineer, 7
 main, 8
 model_evaluate, 8
 model_train, 7

O

one_hot_encode() (在 feature_engineer 模块中), 7

P

plot_box_by_category() (在 data_visualize 模块中), 7

plot_categorical_by_lifecycle() (在 data_visualize 模块中), 7

plot_categorical_distribution() (在
data_visualize 模块中) , [7](#)
plot_correlation_heatmap() (在
data_visualize 模块中) , [7](#)
plot_numeric_distribution() (在
data_visualize 模块中) , [7](#)

R

remove_duplicates() (在 data_clean 模块中) ,
[6](#)

T

train_model() (在 model_train 模块中) , [8](#)