

# 数据处理报告

- 1. 数据提取
- 2. 数据特征提取
- 3. 数据预处理
  - 3.1. 缺失值
  - 3.2. 异常值处理
  - 3.3. 数据转换
  - 3.4. 数据标准化
- 4. 特征工程（计算相关性）
  - 4.1. 计算相关性
  - 4.2. 计算多重共线性VIF
- 5. 模型评估
  - 5.1. 客户星级
  - 5.2. 信用等级

## 1. 数据提取

extract.sqlSQL复制代码

```
1 SELECT uid,
2       tran_flag,
3       tran_amt,
4       cac_intc_pr,
5       dr_cr_code,
6       pay_term
7 FROM dm.dm_v_tr_huanx_mx INTO OUTFILE '/home/vboxuser/dm_v_tr_huanx_mx.csv'
8      FORMAT CSV
9 .....
```

使用了clickhouse自带的csv导出功能  
导出了24张表中需要的数据，导出的数据汇总在”数据库导出的csv数据（已筛选）”文件夹中。

## 2. 数据特征提取

使用pandas库read\_csv函数读取提取好的csv文件并创建对应列表。

```
feature_credit.ipynb Python 复制代码
1 es = ft.EntitySet(id = "my_set")
2 es.entity_from_dataframe(entity_id='pri_credit_info',dataframe=pri_credit_i
  nfo,index='index_column')
3 .....
```

使用featuretool库创建实体集EntitySet以便操作数据，然后使用entity\_from\_dataframe函数添加数据，其中entity\_id表示要从哪个表获取数据，dataframe表示要从哪个DataFrame中获取数据

```
credit_info.ipynb Python 复制代码
1 relationships = [
2     ft.Relationship(es["pri_star_info"]["uid"], es["dm_v_as_djk_info"]["ui
  d"]),
3     ft.Relationship(es["pri_star_info"]["uid"], es["dm_v_as_djkfq_info"]
  ["uid"]),
4     ft.Relationship(es["pri_star_info"]["uid"], es["dm_v_tr_contract_mx"]
  ["uid"]),
5     ft.Relationship(es["pri_star_info"]["uid"], es["dm_v_tr_djk_mx"]
  ["uid"]),
6     ft.Relationship(es["pri_star_info"]["uid"], es["dm_v_tr_dsf_mx"]
  ["uid"]),
7     ft.Relationship(es["pri_star_info"]["uid"], es["dm_v_tr_duebill_mx"]
  ["uid"]),
8     ft.Relationship(es["pri_star_info"]["uid"], es["dm_v_tr_huanb_mx"]
  ["uid"]),
9     ft.Relationship(es["pri_star_info"]["uid"], es["dm_v_tr_huanx_mx"]
  ["uid"]),
10    ft.Relationship(es["pri_star_info"]["uid"], es["dm_v_tr_sa_mx"]
  ["uid"]),
11    .....]
12 es.add_relationships(relationships)
```

数据添加完成后，再使用relationship函数为不同表中相关数据（uid）添加关系；

▼ credit\_info.ipynb

Python

📄 复制代码

```
1 feature_matrix, features_defs = ft.dfs(  
2     entityset=es,  
3     target_entity='pri_star_info',  
4 )  
5 feature_matrix  
6 print(type(feature_matrix))  
7 feature_matrix.to_csv('./result.csv', sep=',', index=True, header=True)
```

最后使用 `ft.dfs` 函数（代表深度特征合成）提取数据特征，返回`feature_matrix`和`features_defs`，其中`features_matrix`表示特征矩阵，`features_defs`表示每个维度上的注释信息。打印`feature_matrix`及其类型信息，并将`feature_matrix`转换成csv文件保存。

### 3. 数据预处理

#### 3.1. 缺失值

▼ deal\_data\_credit.ipynb

Python

📄 复制代码

```
1 for column in df.columns.tolist():  
2     if df[column].dtype == 'object':  
3         df[column].fillna(df[column].mode()[0], inplace=True)  
4     else:  
5         df[column].fillna(df[column].median(), inplace=True)  
6 df
```

首先读取到文件中的所有列，并将列名转换为列表。

其次对于数据类型进行判断，如果是object类别型变量则使用众数填充，使用`df[column].mode()[0]`计算当前数值出现的次数，并且读取到第一个数据为众数；如果不是类别型变量，则使用中位数进行填充，使用`df[column].median()`计算当前列的中位数

#### 3.2. 异常值处理

```
1 for column in df.columns.tolist():
2     if df[column].dtype != 'object' and column != 'credit_level':
3         print(column)
4     if(column == 'credit_level'):
5         print("error")
6     if(column == 'uid'):
7         print("error")
8     low = df[column].quantile(0.05)
9     high = df[column].quantile(0.95)
10    # 使用clip方法替换小于5%和大于95%的值
11    df[column] = df[column].clip(lower=low, upper=high)
```

首先，如果当前列的数据类型不是 object 且列名不为 'credit\_level'，则执行以下代码：计算当前列的 5% 分位数，即数据中排在第 5% 位置的数值，将其赋值给变量 low。计算当前列的 95% 分位数，即数据中排在第 95% 位置的数值，将其赋值给变量 high。

接着，使用 clip 方法对当前列的数据进行剪裁，将小于 low 的值替换为 low，将大于 high 的值替换为 high。

### 3.3. 数据转换

```
1 df['MODE(dm_v_tr_huanb_mx.tran_type)'] = pd.to_numeric(df['MODE(dm_v_tr_huanb_mx.tran_type)'], errors="coerce")
2 df['MODE(dm_v_tr_huanb_mx.tran_type)'] = df['MODE(dm_v_tr_huanb_mx.tran_type)'].astype(float)
3 print(df['MODE(dm_v_tr_huanb_mx.tran_type)'].unique())
4 df['MODE(dm_v_tr_huanb_mx.tran_type)'] = df['MODE(dm_v_tr_huanb_mx.tran_type)'].fillna(0.0)
5 print(df['MODE(dm_v_tr_huanb_mx.tran_type)'].unique())
```

对于数据进行转换处理，过程如下：

- 1 将该列的数据类型转换为数值型（numeric），如果无法转换则将其转为缺失值（coerce）。
- 2 将该列的数据类型转换为浮点型（float）。
- 3 输出该列不同的唯一值（unique）。
- 4 将该列中的缺失值填充为 0。
- 5 再次输出该列不同的唯一值，观察是否存在缺失值。

最终将非数值的编码转换为数值类型，同时对于缺失值进行处理

### 3.4. 数据标准化

deal\_data\_credit.ipynb

Python

复制代码

```
1 scaler = StandardScaler()
2 df = pd.read_csv("./credit_encoded.csv", dtype=None)
3 df_columns = df.columns.to_list()
4 df_columns.remove('uid')
5 df_columns.remove('credit_level')
6 print(df_columns)
7 df[df_columns] = scaler.fit_transform(df[df_columns])
8 df.head()
```

代码通过 `df_columns = df.columns.to_list()` 将数据框中的列名保存到名为 `df_columns` 的列表中，并使用 `df_columns.remove('uid')` 和 `df_columns.remove('credit_level')` 从列表中删除了不需要处理的 'uid' 和 'credit\_level' 两列。接下来，代码使用 `StandardScaler` 对数据框中除 'uid' 和 'credit\_level' 以外的所有列进行标准化处理，并将处理后的结果保存回数据框中。

## 4. 特征工程（计算相关性）

### 4.1. 计算相关性

feature\_credit

Python

复制代码

```
1 df = pd.read_csv("./data/credit_final.csv")
2 corr_matrix = df.corr()
```

计算相关性主要使用pandas库中的corr函数，从之前生成的csv文件中生成相关系数矩阵

```
1 to_drop = corr_series[abs(corr_series) >= 0.7].index.tolist()
2
3 drop = []
4 index = 0
5 for t in to_drop:
6     if (t[0][0:4] == 'MEAN' and t[1][0:4] == 'MAX(') or (t[1][0:4] == 'MEAN' and t[0][0:4] == 'MAX('):
7         if t[1][0:4] == 'MAX(':
8             to_drop.pop(index)
9             drop.append(t[1])
10        else:
11            to_drop.pop(index)
12            drop.append(t[0])
13        if (t[0][0:4] == 'MEAN' and t[1][0:4] == 'MIN(') or (t[1][0:4] == 'MEAN' and t[0][0:4] == 'MIN('):
14            if t[1][0:4] == 'MIN(':
15                drop.append(t[1])
16                to_drop.pop(index)
17            else:
18                drop.append(t[0])
19                to_drop.pop(index)
20        index = index + 1
21
22 df.drop(columns=drop,axis=1,inplace=True)
23 df.to_csv("./data/credit_after_corr.csv")
```

（由于删去数据之后导致训练结果不好，最后选择没有删去）

筛选相关系数高于指定值的，并记录其列名），并且从dataframe中删除这些相关系数高于指定值的，得到处理后的数据，并保存到对应csv文件中。

## 4.2. 计算多重共线性VIF

```
1 import pandas as pd
2 import statsmodels.api as sm
3
4 df = pd.read_csv("./credit_after_corr.csv")
5
6 df.drop(columns='Unnamed: 0',axis=1,inplace=True)
7 x_list = df.columns.to_list()
8 print(x_list)
9
10 from statsmodels.stats.outliers_influence import variance_inflation_factor
11
12 x_list.remove('Unnamed: 0.1')
13 x_list.remove('uid')
14 x_list.remove('credit_level')
15 X = df[x_list]
16 Y = df['credit_level']
17 model = sm.OLS(Y, X).fit()
18 vif = pd.DataFrame()
19 vif["VIF Factor"] = [variance_inflation_factor(X.values, i) for i in range
    (X.shape[1])]
20 vif["features"] = X.columns
21
22 vif.to_csv("./credit_vif_1.csv")
```

计算VIF方差膨胀因子时，从statsmodels库中引入了variance\_inflation\_factor函数以计算VIF值，在计算前提前删除了无用的数据，计算时VIF结果存储在一个dataframe中，其中VIF Factor是计算出的值，features是列名。

```
1 vif_threshold = 10
2 high_vif_features = list(vif[vif['VIF Factor'] > vif_threshold]["features"]
3 )
4 df.drop(columns=high_vif_features,axis=1,inplace=True)
5 df.to_csv('./credit_after_vif.csv')
```

（由于删去数据之后导致训练结果不好，最后选择没有删去）

并筛选VIF值高于指定值的，并记录其features（即列名），并且从dataframe中删除这些VIF值高于指定值的得到处理后的数据，并保存到对应csv文件中。

## 5. 模型评估

### 5.1. 客户星级

模型	准确率	精确率	召回率	F1分数	Cohen's Kappa系数
随机森林模型	0.9738573071604875	0.9740772053204735	0.9738573071604875	0.973903174880576	0.9582329939228935
XGBoost模型	0.7379495272495462	0.9202710328280544	0.5432862868389917	0.5654842522505414	0.8731477277261196

### 5.2. 信用等级

模型	准确率	精确率	召回率	F1分数	Cohen's Kappa系数
随机森林模型	0.9633147603033224	0.957026256269053	0.8455545355692532	0.8913416340143352	0.9203203653612596
XGBoost模型	0.8522306355918314	0.6773776797965199	0.5835956243077896	0.6156770921619316	0.7289625204956423