

Лабораторная работа №2 по курсу "Методы машинного обучения"

Выполнил: Саврасов П.А. группа ИУ5-24М

Задание

Выбрать набор данных (датасет), содержащий категориальные и числовые признаки и пропуски в данных. Для выполнения следующих пунктов можно использовать несколько различных наборов данных (один для обработки пропусков, другой для категориальных признаков и т.д.) Просьба не использовать датасет, на котором данная задача решалась в лекции.</p>
</div>
<div data-bbox="79 55 587 60" data-label="Text>
<p>Для выбранного датасета (датасетов) на основе материалов лекций решить следующие задачи:</p>
</div>
<div data-bbox="91 64 316 79" data-label="List-Group>

• устранение пропусков в данных;
• кодирование категориальных признаков;
• нормализацию числовых признаков.

</div>
<div data-bbox="26 83 981 121" data-label="Code-Block>
<pre>In [41]: import numpy as np
import pandas as pd
import seaborn as sns
import matplotlib.pyplot as plt
from sklearn.impute import SimpleImputer
from sklearn.preprocessing import LabelEncoder
import scipy.stats as stats
import warnings
warnings.filterwarnings('ignore')</pre>
</div>
<div data-bbox="26 125 981 163" data-label="Code-Block>
<pre>In [42]: def diagnostic_plots(df, variable):
plt.figure(figsize=(15,6))
гистограмма
plt.subplot(1, 2, 1)
df[variable].hist(bins=30)
Q-Q plot
plt.subplot(1, 2, 2)
stats.probplot(df[variable], dist="norm", plot=plt)
plt.show()</pre>
</div>
<div data-bbox="26 167 981 181" data-label="Code-Block>
<pre>In [43]: data = pd.read_csv('BankChurners.csv', sep=",")
data.head()</pre>
</div>
<div data-bbox="26 183 981 233" data-label="Table>
<table>
<tr>
<th></th>
<th>CLIENTNUM</th>
<th>Attrition_Flag</th>
<th>Customer_Age</th>
<th>Gender</th>
<th>Dependent_count</th>
<th>Education_Level</th>
<th>Marital_Status</th>
<th>Income_Category</th>
<th>Card_Category</th>
<th>Months_on_book</th>
<th>...</th>
<th>Credit_Limit</th>
<th>Total</th>
</tr>
<tr>
<td>0</td>
<td>768805383</td>
<td>Existing Customer</td>
<td>45</td>
<td>M</td>
<td>3</td>
<td>High School</td>
<td>Married</td>
<td>60K–80K</td>
<td>Blue</td>
<td>39</td>
<td>...</td>
<td>12691.0</td>
<td></td>
</tr>
<tr>
<td>1</td>
<td>818770008</td>
<td>Existing Customer</td>
<td>49</td>
<td>F</td>
<td>5</td>
<td>Graduate</td>
<td>Single</td>
<td>Less than \$40K</td>
<td>Blue</td>
<td>44</td>
<td>...</td>
<td>8256.0</td>
<td></td>
</tr>
<tr>
<td>2</td>
<td>713982108</td>
<td>Existing Customer</td>
<td>51</td>
<td>M</td>
<td>3</td>
<td>Graduate</td>
<td>Married</td>
<td>80K–120K</td>
<td>Blue</td>
<td>36</td>
<td>...</td>
<td>3418.0</td>
<td></td>
</tr>
<tr>
<td>3</td>
<td>769911858</td>
<td>Existing Customer</td>
<td>40</td>
<td>F</td>
<td>4</td>
<td>High School</td>
<td>NaN</td>
<td>Less than \$40K</td>
<td>Blue</td>
<td>34</td>
<td>...</td>
<td>3313.0</td>
<td></td>
</tr>
<tr>
<td>4</td>
<td>709106358</td>
<td>Existing Customer</td>
<td>40</td>
<td>M</td>
<td>3</td>
<td>Uneducated</td>
<td>Married</td>
<td>60K–80K</td>
<td>Blue</td>
<td>21</td>
<td>...</td>
<td>4716.0</td>
<td></td>
</tr>
</table>
</div>
<div data-bbox="79 235 183 239" data-label="Text>
<p>5 rows × 23 columns</p>
</div>
<div data-bbox="79 241 981 255" data-label="Text>
<p>Заполнение пустых значений</p>
</div>
<div data-bbox="26 257 981 289" data-label="Code-Block>
<pre>In [44]: columnsWithNull = []

print("Столбцы с пустыми значениями (название (число уникальных значений): число пустых):\n")
for column in data.columns:
 if data[column].isnull().sum() > 0:
 columnsWithNull.append(column)
 print("\t", column,"(", len(data[column].unique()), "):\t", data[column].isnull().sum())
data[columnsWithNull].head(10)</pre>
</div>
<div data-bbox="83 291 598 295" data-label="Text>
<p>Столбцы с пустыми значениями (название (число уникальных значений): число пустых):</p>
</div>
<div data-bbox="26 297 981 375" data-label="Table>
<table>
<tr>
<th></th>
<th>Education_Level</th>
<th>Marital_Status</th>
<th>Income_Category</th>
</tr>
<tr>
<td>0</td>
<td>High School</td>
<td>Married</td>
<td>60K–80K</td>
</tr>
<tr>
<td>1</td>
<td>Graduate</td>
<td>Single</td>
<td>Less than \$40K</td>
</tr>
<tr>
<td>2</td>
<td>Graduate</td>
<td>Married</td>
<td>80K–120K</td>
</tr>
<tr>
<td>3</td>
<td>High School</td>
<td>NaN</td>
<td>Less than \$40K</td>
</tr>
<tr>
<td>4</td>
<td>Uneducated</td>
<td>Married</td>
<td>60K–80K</td>
</tr>
<tr>
<td>5</td>
<td>Graduate</td>
<td>Married</td>
<td>40K–60K</td>
</tr>
<tr>
<td>6</td>
<td>NaN</td>
<td>Married</td>
<td>\$120K +</td>
</tr>
<tr>
<td>7</td>
<td>High School</td>
<td>NaN</td>
<td>60K–80K</td>
</tr>
<tr>
<td>8</td>
<td>Uneducated</td>
<td>Single</td>
<td>60K–80K</td>
</tr>
<tr>
<td>9</td>
<td>Graduate</td>
<td>Single</td>
<td>80K–120K</td>
</tr>
</table>
</div>
<div data-bbox="26 383 981 397" data-label="Code-Block>
<pre>In [45]: imputer = SimpleImputer(missing_values=np.nan, strategy="constant", fill_value = "Unknown")
nullFixedData = pd.DataFrame(data = imputer.fit_transform(data), columns=data.columns)
nullFixedData[columnsWithNull].head(10)</pre>
</div>
<div data-bbox="26 400 981 465" data-label="Table>
<table>
<tr>
<th></th>
<th>Education_Level</th>
<th>Marital_Status</th>
<th>Income_Category</th>
</tr>
<tr>
<td>0</td>
<td>High School</td>
<td>Married</td>
<td>60K–80K</td>
</tr>
<tr>
<td>1</td>
<td>Graduate</td>
<td>Single</td>
<td>Less than \$40K</td>
</tr>
<tr>
<td>2</td>
<td>Graduate</td>
<td>Married</td>
<td>80K–120K</td>
</tr>
<tr>
<td>3</td>
<td>High School</td>
<td>Unknown</td>
<td>Less than \$40K</td>
</tr>
<tr>
<td>4</td>
<td>Uneducated</td>
<td>Married</td>
<td>60K–80K</td>
</tr>
<tr>
<td>5</td>
<td>Graduate</td>
<td>Married</td>
<td>40K–60K</td>
</tr>
<tr>
<td>6</td>
<td>Unknown</td>
<td>Married</td>
<td>\$120K +</td>
</tr>
<tr>
<td>7</td>
<td>High School</td>
<td>Unknown</td>
<td>60K–80K</td>
</tr>
<tr>
<td>8</td>
<td>Uneducated</td>
<td>Single</td>
<td>60K–80K</td>
</tr>
<tr>
<td>9</td>
<td>Graduate</td>
<td>Single</td>
<td>80K–120K</td>
</tr>
</table>
</div>
<div data-bbox="79 470 292 474" data-label="Text>
<p>Кодирование категориальных признаков</p>
</div>
<div data-bbox="26 478 981 500" data-label="Code-Block>
<pre>In [46]: labelEnc = LabelEncoder()
for column in nullFixedData.columns:
 if nullFixedData[column].dtype not in ['float', 'int']:
 nullFixedData[column] = pd.DataFrame(labelEnc.fit_transform(nullFixedData[column].astype(str)), columns=[column])
nullFixedData.head(10)</pre>
</div>
<div data-bbox="26 503 981 568" data-label="Table>
<table>
<tr>
<th></th>
<th>CLIENTNUM</th>
<th>Attrition_Flag</th>
<th>Customer_Age</th>
<th>Gender</th>
<th>Dependent_count</th>
<th>Education_Level</th>
<th>Marital_Status</th>
<th>Income_Category</th>
<th>Card_Category</th>
<th>Months_on_book</th>
<th>...</th>
<th>Credit_Limit</th>
<th>Total</th>
</tr>
<tr>
<td>0</td>
<td>7152</td>
<td>1</td>
<td>19</td>
<td>1</td>
<td>3</td>
<td>3</td>
<td>1</td>
<td>2</td>
<td>0</td>
<td>26</td>
<td>...</td>
<td>518</td>
<td></td>
</tr>
<tr>
<td>1</td>
<td>9812</td>
<td>1</td>
<td>23</td>
<td>0</td>
<td>5</td>
<td>2</td>
<td>2</td>
<td>4</td>
<td>0</td>
<td>31</td>
<td>...</td>
<td>5723</td>
<td></td>
</tr>
<tr>
<td>2</td>
<td>3053</td>
<td>1</td>
<td>25</td>
<td>1</td>
<td>3</td>
<td>2</td>
<td>1</td>
<td>3</td>
<td>0</td>
<td>23</td>
<td>...</td>
<td>3718</td>
<td></td>
</tr>
<tr>
<td>3</td>
<td>7204</td>
<td>1</td>
<td>14</td>
<td>0</td>
<td>4</td>
<td>3</td>
<td>3</td>
<td>4</td>
<td>0</td>
<td>21</td>
<td>...</td>
<td>3612</td>
<td></td>
</tr>
<tr>
<td>4</td>
<td>501</td>
<td>1</td>
<td>14</td>
<td>1</td>
<td>3</td>
<td>5</td>
<td>1</td>
<td>2</td>
<td>0</td>
<td>8</td>
<td>...</td>
<td>4463</td>
<td></td>
</tr>
<tr>
<td>5</td>
<td>2544</td>
<td>1</td>
<td>18</td>
<td>1</td>
<td>2</td>
<td>2</td>
<td>1</td>
<td>1</td>
<td>0</td>
<td>23</td>
<td>...</td>
<td>4080</td>
<td></td>
</tr>
<tr>
<td>6</td>
<td>9493</td>
<td>1</td>
<td>25</td>
<td>1</td>
<td>4</td>
<td>6</td>
<td>1</td>
<td>0</td>
<td>1</td>
<td>33</td>
<td>...</td>
<td>3747</td>
<td></td>
</tr>
<tr>
<td>7</td>
<td>9818</td>
<td>1</td>
<td>6</td>
<td>1</td>
<td>0</td>
<td>3</td>
<td>3</td>
<td>2</td>
<td>3</td>
<td>14</td>
<td>...</td>
<td>3135</td>
<td></td>
</tr>
<tr>
<td>8</td>
<td>1401</td>
<td>1</td>
<td>11</td>
<td>1</td>
<td>3</td>
<td>5</td>
<td>2</td>
<td>2</td>
<td>0</td>
<td>23</td>
<td>...</td>
<td>2206</td>
<td></td>
</tr>
<tr>
<td>9</td>
<td>5893</td>
<td>1</td>
<td>22</td>
<td>1</td>
<td>2</td>
<td>2</td>
<td>2</td>
<td>3</td>
<td>0</td>
<td>23</td>
<td>...</td>
<td>332</td>
<td></td>
</tr>
</table>
</div>
<div data-bbox="79 571 189 575" data-label="Text>
<p>10 rows × 23 columns</p>
</div>
<div data-bbox="79 584 268 588" data-label="Text>
<p>Нормализация числовых признаков</p>
</div>
<div data-bbox="26 592 981 606" data-label="Code-Block>
<pre>In [47]: data.hist(figsize=(20,20))
plt.show()</pre>
</div>
<div data-bbox="79 608 981 769" data-label="Figure>

</div>
<div data-bbox="26 773 981 787" data-label="Code-Block>
<pre>In [57]: diagnostic_plots(data, 'Total_Trans_Ct')</pre>
</div>
<div data-bbox="79 789 981 871" data-label="Figure>

</div>
<div data-bbox="26 879 981 893" data-label="Code-Block>
<pre>In [68]: data['Total_Trans_Ct_sqrt'] = data['Total_Trans_Ct']**(1/2)
diagnostic_plots(data, 'Total_Trans_Ct_sqrt')</pre>
</div>
<div data-bbox="79 895 981 977" data-label="Figure>

</div>
<div data-bbox="26 985 981 999" data-label="Code-Block>
<pre>In []:</pre>
</div>
</div>