Московский государственный технический университет им. Н.Э. Баумана

Факультет «Информатика и системы управления» Кафедра «Системы обработки информации и управления»



Лабораторная работа № 2

По курсу «методы машинного обучения в АСОИУ»

«Обработка признаков. Часть 1»

Выполнил:

студент ИУ5-23M Семенов И.А.

Проверил:

Гапанюк Ю.Е.

Подпись:

29.02.2024

- 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.
- Для выбранного датасета (датасетов) на основе материалов лекций решить следующие задачи:
- ✓ устранение пропусков в данных;
- ✓ кодирование категориальных признаков;
- ✓ нормализация числовых признаков.

Текст программы и экранные формы с примерами выполнения программы

У Кодирование категориальных признаков

```
[] import numpy as np
import pandas as pd
import seaborn as sns
import matplotlib.pyplot as plt
%matplotlib inline
sns.set(style="ticks")
```

Загрузка и предобработка данных

Используем данные из соревнования <u>Titanic</u>

```
[ ] # Будем использовать только обучающую выборку data_loaded = pd.read_csv('data/titanic.csv', sep=",")
```

[] # размер набора данных data_loaded.shape

→ (891, 12)

[] data_loaded.head()

<u>-</u>	PassengerId	Survived	Pclass	Name	Sex	Age	SibSp	Parch	Ticket	Fare	Cabin	Embarked
) 1	0	3	Braund, Mr. Owen Harris	male	22.0	1	0	A/5 21171	7.2500	NaN	S
	2	1	1	Cumings, Mrs. John Bradley (Florence Briggs Th	female	38.0	1	0	PC 17599	71.2833	C85	С
	2 3	1	3	Heikkinen, Miss. Laina	female	26.0	0	0	STON/O2. 3101282	7.9250	NaN	S
	3 4	1	1	Futrelle, Mrs. Jacques Heath (Lily May Peel)	female	35.0	1	0	113803	53.1000	C123	S
	5	0	3	Allen, Mr. William Henry	male	35.0	0	0	373450	8.0500	NaN	S

```
[ ] data_features = list(zip(
                                   # признаки
                                   [i for i in data_loaded.columns],
                                   zip(
                                        # типы колонок
                                        [str(i) for i in data_loaded.dtypes],
                                        # проверим есть ли пропущенные значения
                                        [i for i in data_loaded.isnull().sum()]
                                   # Признаки с типом данных и количеством пропусков
                                   data features
                             ('Survived', ('int64', 0)), ('Pclass', ('int64', 0)),
                                    ('Name', ('object', 0)),
                                    ('Sex', ('object', 0)),
('Age', ('float64', 177)),
('SibSp', ('int64', 0)),
                                    ('Parch', ('int64', 0)),
('Ticket', ('object', 0)),
('Fare', ('float64', 0)),
                                    ('Cabin', ('object', 687))
                                    ('Embarked', ('object', 2))]
                             [ ] # Используем только некоторые признаки
                                   data = data_loaded[cols_filter]
                                   data.head()
                             \overline{\Sigma}
                                       Pclass Age SibSp Parch
                                                                                       Sex Cabin Embarked Survived
                                                                            Fare
                                    0
                                             3 22.0
                                                                      0 7.2500
                                                                                    male
                                                                                               NaN
                                                                                                               S
                                                                                                                            0
                                    1
                                              1 38.0
                                                                      0 71.2833 female
                                                                                               C85
                                                                                                               C
                                                                                                                            1
                                    2
                                              3 26.0
                                                                        7.9250 female
                                                                                               NaN
                                    3
                                              1 35.0
                                                             1
                                                                      0 53.1000 female
                                                                                              C123
                                                                                                               S
                                                                                                                            1
[ ] # Заполним пропуски
     data.dropna(subset=['Fare', 'Embarked'], inplace=True)
5: C:\ProgramData\Anaconda3\lib\site-packages\pandas\util\_decorators.py:311: SettingWithCopyWarning:
     A value is trying to be set on a copy of a slice from a DataFrame
     See the caveats in the documentation: <a href="https://pandas.pydata.org/pandas-docs/stable/user_guide/indexing.html#returning-a-view-versus-a-copy">https://pandas.pydata.org/pandas.pydata.org/pandas-docs/stable/user_guide/indexing.html#returning-a-view-versus-a-copy</a>
       return func(*args, **kwargs)
[ ] # От каюты оставляет только первую букву
     # и убираем каюты типа Т так как их мало
     data['Cabin'] = data['Cabin'].astype(str).str[0]
     data = data[data['Cabin'] != 'T']
C:\Users\Paladin\AppData\Local\Temp/ipykernel_9088/3362610652.py:3: SettingWithCopyWarning:
     A value is trying to be set on a copy of a slice from a DataFrame.
    Try using .loc[row_indexer,col_indexer] = value instead
     See the caveats in the documentation: <a href="https://pandas.pydata.org/pandas-docs/stable/user_guide/indexing.html#returning-a-view-versus-a-copy">https://pandas.pydata.org/pandas.pydata.org/pandas-docs/stable/user_guide/indexing.html#returning-a-view-versus-a-copy</a>
       data['Cabin'] = data['Cabin'].astype(str).str[0]
[ ] # Заполним пропуски возраста средними значениями
     def impute_na(df, variable, value):
         df[variable].fillna(value, inplace=True)
     impute_na(data, 'Age', data['Age'].mean())
[ ] # Убедимся что нет пустых значений
    data.isnull().sum()
```

→ Pclass Age

SibSp

Parch

Fare Sex Cabin Embarked Survived 0

0

0

▼ 1. Кодирование категорий целочисленными значениями - <u>label encoding</u>

Label Encoding (LE), который также называют integer encoding, предполагает, что значения категорий заменяются целыми числами в случайном порядке.

- Преимущества:
 - Простота реализации.
 - Не расширяется признаковое пространство (не появляется дополнительных колонок).

Недостатки:

- Не использует информацию о распределении значений категорий.
- Не подходит для линейных моделей, так как создает фиктивное отношение порядка между значениями.

```
[ ] from sklearn.preprocessing import LabelEncoder

[ ] le = LabelEncoder()
    cat_enc_le = le.fit_transform(data['Cabin'])

[ ] data['Cabin'].unique()

array(['n', 'C', 'E', 'G', 'D', 'A', 'B', 'F'], dtype=object)

[ ] np.unique(cat_enc_le)

array([0, 1, 2, 3, 4, 5, 6, 7])

[ ] le.inverse_transform([0, 1, 2, 3])

array(['A', 'B', 'C', 'D'], dtype=object)
```

✓ 2. Кодирование категорий наборами бинарных значений - one-hot encoding

One-hot encoding предполагает, что значение категории заменяется на отдельную колонку, которая содержит бинарные значения.

- Преимущества:
 - Простота реализации.
 - Подходит для любых моделей, так как НЕ создает фиктивное отношение порядка между значениями.

Недостатки:

• Расширяется признаковое пространство.

[0., 0., 0., 0., 0., 0., 1.]])

```
[ ] from sklearn.preprocessing import OneHotEncoder
[ ] ohe = OneHotEncoder()
    cat_enc_ohe = ohe.fit_transform(data[['Cabin']])
    cat_enc_ohe
< <888x8 sparse matrix of type '<class 'numpy.float64'>'
            with 888 stored elements in Compressed Sparse Row format>
cat_enc_ohe.todense()[0:10]
→ matrix([[0., 0., 0., 0., 0., 0., 0., 1.],
            [0., 0., 1., 0., 0., 0., 0., 0.],
            [0., 0., 0., 0., 0., 0., 0., 1.],
            [0., 0., 1., 0., 0., 0., 0., 0.]
            [0., 0., 0., 0., 0., 0., 0., 1.],
            [0., 0., 0., 0., 0., 0., 1.],
            [0., 0., 0., 0., 1., 0., 0., 0.],
            [0., 0., 0., 0., 0., 0., 1.],
            [0., 0., 0., 0., 0., 0., 0., 1.],
```

∨ Использование библиотеки Category Encoders



Count encoding предполагает что значение категории заменяется на количество раз, которое оно встречается в категории.

В случае frequency encoding вместо количества используется доля (процент) от количества записей.

Преимущества:

• Простота реализации.

886 27.000000

687

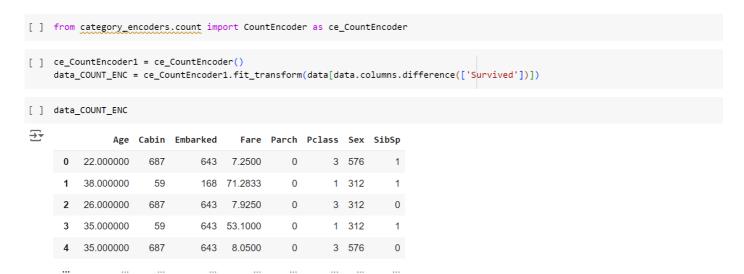
• Не расширяется признаковое пространство.

Недостатки:

• Если два значения встречаются одинаковое количество раз, то они будут заменены на одно и то же количество и становятся неразличимы.

У Использование библиотеки <u>Category Encoders</u>

643 13.0000



2 576

0

```
[ ] data['Cabin'].unique()
    array(['n', 'C', 'E', 'G', 'D', 'A', 'B', 'F'], dtype=object)
[ ] data_COUNT_ENC['Cabin'].unique()
    array([687, 59, 32,
                            4, 33, 15, 45, 13], dtype=int64)
[ ] ce_CountEncoder2 = ce_CountEncoder(normalize=True)
    data_FREQ_ENC = ce_CountEncoder2.fit_transform(data[data.columns.difference(['Survived'])])
[ ] data_FREQ_ENC
₹
                Age
                       Cabin Embarked
                                          Fare Parch Pclass
                                                                   Sex SibSp
          22.000000 0.773649 0.724099
                                        7.2500
                                                            3 0.648649
          38.000000 0.066441
                                                            1 0.351351
      1
                             0.189189 71.2833
                                                    0
                                                                            1
         26.000000 0.773649
                             0.724099
                                                            3 0.351351
                                        7.9250
      3
          35.000000 0.066441
                              0.724099
                                       53.1000
                                                    0
                                                              0.351351
                                                                            1
          35.000000 0.773649
                              0.724099
                                        8.0500
                                                            3 0.648649
     886
         27.000000 0.773649
                              0.724099 13.0000
                                                            2 0.648649
                                       30.0000
                                                            1 0.351351
     887
          19.000000 0.050676
                              0.724099
                                                    0
                                                                            0
          29.620492 0.773649
                              0.724099
                                       23.4500
                                                            3 0.351351
     889
          26.000000 0.066441
                             0.189189 30.0000
                                                    0
                                                              0.648649
                                                                            0
         32.000000 0.773649 0.086712
                                        7.7500
                                                            3 0.648649
    888 rows × 8 columns
[ ] data_FREQ_ENC['Cabin'].unique()
→ array([0.77364865, 0.06644144, 0.03603604, 0.0045045 , 0.03716216,
```

На основе полученных результатов можем сделать вывод о том, что качество яблока напрямую зависит от его зрелости. Также имеется прямая зависимость между зрелостью яблока и его хрупкостью. Сладость же и сочность почти не зависят друг от друга.

[✓] Подключено к среде выполнения "Серверный ускоритель Python 3 на базе Google Compute Engine ().".