

### Министерство образования и науки Российской Федерации Федеральное государственное бюджетное образовательное учреждение высшего образования

«Московский государственный технический университет имени Н.Э. Баумана

(национальный исследовательский университет)» (МГТУ им. Н.Э. Баумана)

ФАКУЛЬТЕТ

#### ИНФОРМАТИКА И СИСТЕМЫ УПРАВЛЕНИЯ

КАФЕДРА СИСТЕМЫ ОБРАБОТКИ ИНФОРМАЦИИ И УПРАВЛЕНИЯ

# Отчет по лабораторной работе № 3 «Обработка пропусков в данных, кодирование категориальных признаков, масштабирование данных»

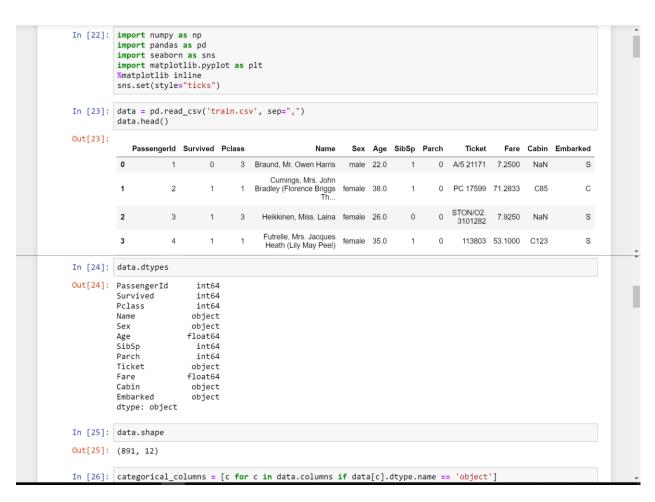
по курсу "Технологии машинного обучения"

Исполнитель: Студент группы ИУ5-63 Желанкина А.С. 03.03.2018

## Задание лабораторной работы

- 1. Выбрать набор данных (датасет), содержащий категориальные 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.)
- 2. Для выбранного датасета (датасетов) на основе материалов лекции решить следующие задачи:
- обработку пропусков в данных;
- кодирование категориальных признаков;
- масштабирование данных.

# Экранные формы с текстом программы и примерами её выполнения



```
numerical_columns = [c for c in data.columns if data[c].dtype.name != 'object']
           print(categorical_columns)
           print(numerical_columns)
           ['Name', 'Sex', 'Ticket', 'Cabin', 'Embarked']
['PassengerId', 'Survived', 'Pclass', 'Age', 'SibSp', 'Parch', 'Fare']
In [27]: data.isnull().sum()
Out[27]: PassengerId
           Survived
                               0
           Pclass
                               0
           Name
                               0
           Sex
                               0
                             177
           Age
           SibSp
           Parch
                               0
           Ticket
                               0
           Fare
           Cabin
                             687
           Embarked
           dtype: int64
In [28]: num_cols = []
    for col in data.columns:
                # Количество пустых значений
                temp_null_count = data[data[col].isnull()].shape[0]
dt = str(data[col].dtype)
                if temp_null_count>0:
                    num_cols.append(col)
                    hom_cols.append(col)
temp_perc = round((temp_null_count / data.shape[0]) * 100.0, 2)
print('Колонка {}. Тип данных {}. Количество пустых значений {}, {}%.'.format(col, dt, temp_nu
           Колонка Age. Тип данных float64. Количество пустых значений 177, 19.87%.
           Колонка Cabin. Тип данных object. Количество пустых значений 687, 77.1%.
Колонка Embarked. Тип данных object. Количество пустых значений 2, 0.22%.
In [29]: data = data.drop(('Cabin'), axis=1)
In [30]: data.isnull().sum()
Out[30]: PassengerId
           Survived
           Pclass
                               0
           Name
                               0
           Sex
           Age
                             177
           SibSp
                               0
           Parch
           Ticket
                               0
           Fare
                               0
           Embarked
           dtype: int64
In [31]: data['Age'].describe()
Out[31]: count
                      714.000000
                       29.699118
           mean
           std
                       14.526497
           min
                        0.420000
           25%
                       20.125000
           50%
                       28.000000
           75%
                       38.000000
           max
                       80.000000
           Name: Age, dtype: float64
In [32]: from sklearn.impute import SimpleImputer
           from sklearn.impute import MissingIndicator
           data['Age'] = SimpleImputer(strategy='mean').fit_transform(data[['Age']])
In [33]: data['Age'].describe()
Out[33]: count
                      891.000000
           mean
                       29.699118
                       13.002015
           std
                        0.420000
           min
           25%
                       22.000000
           50%
                       29.699118
           75%
                       35.000000
                       80.000000
           Name: Age, dtype: float64
In [34]: data['Embarked'].describe()
Out[34]: count
                       889
           unique
```

```
top
                    644
          freq
          Name: Embarked, dtype: object
In [35]: data['Embarked'] = SimpleImputer(missing_values=np.nan, strategy='most_frequent').fit_transform(data[[
In [36]: data['Embarked'].describe()
Out[36]: count
                    891
          unique
                      3
S
          top
          freq
                    646
          Name: Embarked, dtype: object
In [37]: data.isnull().sum()
Out[37]: PassengerId
          Survived
          Pclass
                         0
          Name
                         0
                         0
          Sex
          Age
          SibSp
                         0
                         0
          Parch
          Ticket
          Embarked
                         0
          dtype: int64
In [38]: pd.get_dummies(data).head()
Out[38]:
                                                                               Name_Abbott,
Mr. Rossmore
Edward (Rosa Hunt)
                                                            Fare Name_Abbing,
Mr. Anthony
                                                                                                           Ticket_W./C.
14263
             Passengerld Survived Pclass Age SibSp Parch
                                                                                         0
                                                                                                                    0
           0
                                      3 22.0
                                                       0
                                                          7.2500
                                                                            0
                      2
                                                                            0
                                                                                         0
                                                                                                      0
                                                                                                                    0
           1
                                      1 38.0
                                                       0 71.2833
                                                                                         0
                                                                                                      0 ...
                                                                                                                    0
           2
                      3
                                     3 26.0
                                                 0
                                                       0
                                                          7.9250
                                                                            0
           3
                                                                                         0
                      4
                                      1 35.0
                                                       0 53.1000
                                                                            0
                                                                                                      0
                                                                                                                    0
                                                 0
          4
                      5
                              0
                                     3 35.0
                                                      0 8.0500
                                                                            0
                                                                                         0
                                                                                                      0
                                                                                                                    0
In [46]: from sklearn.preprocessing import MinMaxScaler
          data = MinMaxScaler().fit_transform(data)
          data
Out[46]: array([[0.],
                 [1.],
[1.],
                 [0.],
                 [0.],
```