

Московский государственный технический университет им. Н.Э. Баумана
Кафедра «Системы обработки информации и управления»

Лабораторная работа №3
по дисциплине

«Методы машинного обучения»
на тему

«Обработка пропусков в данных, кодирование
категориальных признаков, масштабирование
данных»

Выполнил:
студент группы ИУ5-21М
Исмаил Ахмад

Москва — 2020 г.

1. Цель лабораторной работы

Изучить способы предварительной обработки данных для дальнейшего формирования моделей

2. Задание

Требуется [1]:

1. Выбрать набор данных (датасет), содержащий категориальные признаки и пропуски в данных. Для выполнения следующих пунктов можно использовать несколько

различных наборов данных.

2. Для выбранного датасета (датасетов) на основе материалов лекции решить следующие задачи:

- обработку пропусков в данных;
- кодирование категориальных признаков;
- масштабирование данных.

3. Ход выполнения работы

Подключим все необходимые библиотеки и настроим отображение графиков [2,3]:

```
In [1]: import numpy as np
import pandas as pd
import seaborn as sns
import sklearn.impute
import sklearn.preprocessing
```

```
/usr/local/lib/python3.6/dist-packages/statsmodels/tools/_testing.py:19: FutureWarning: pandas.util.testing is deprecated. Use the functions in the public API at pandas.testing instead.
import pandas.util.testing as tm
```

```
# Enable inline plots
```

```
In [0]: %matplotlib inline
```

```
# Set plot style
```

```
In [0]: sns.set(style="ticks")
```

```
In [0]: from IPython.display import set_matplotlib_formats
set_matplotlib_formats("retina")
```

тобы в дальнейшем текст в отчёте влезал на A4 [4]:

```
In [0]: pd.set_option("display.width", 70)
```

Для выполнения данной лабораторной работы возьмём набор данных по приложениям в US counties COVID 19 dataset

```
In [0]: data = pd.read_csv("/content/sample_data/us-counties.csv")
```

In [7]: data.head()

Out[7]:

	date	county	state	fips	cases	deaths
0	2020-01-21	Snohomish	Washington	53061.0	1	0
1	2020-01-22	Snohomish	Washington	53061.0	1	0
2	2020-01-23	Snohomish	Washington	53061.0	1	0
3	2020-01-24	Cook	Illinois	17031.0	1	0
4	2020-01-24	Snohomish	Washington	53061.0	1	0

In [8]: data.dtypes

Out[8]:

```
date      object
county    object
state     object
fips      float64
cases     int64
deaths    int64
dtype: object
```

In [9]: data.shape

Out[9]: (64707, 6)

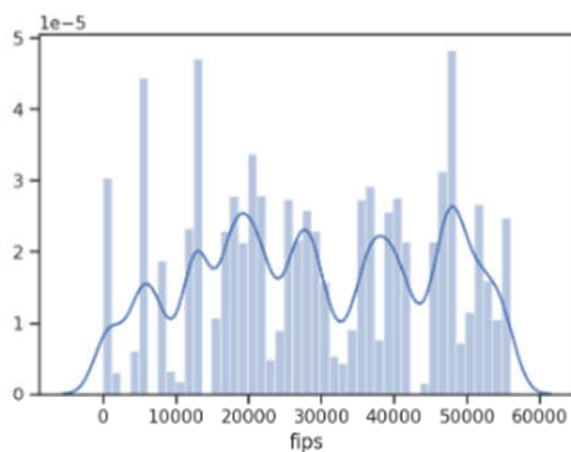
Обработка пропусков в данных

In [10]: data.isnull().sum()

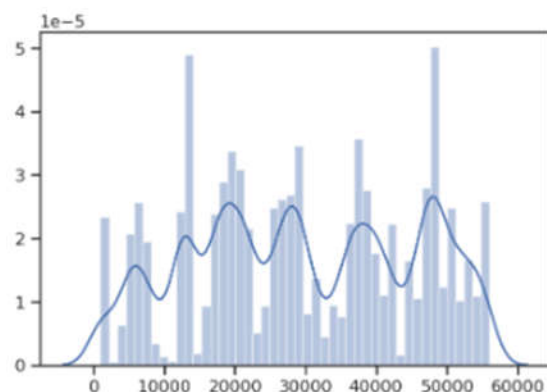
Out[10]:

```
date      0
county    0
state     0
fips      838
cases     0
deaths    0
dtype: int64
```

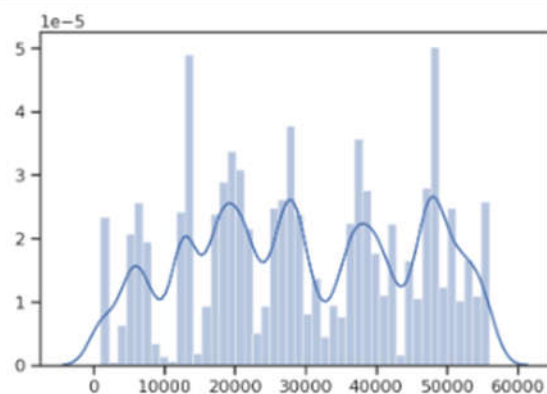
In [11]: sns.distplot(data["fips"].fillna(0));



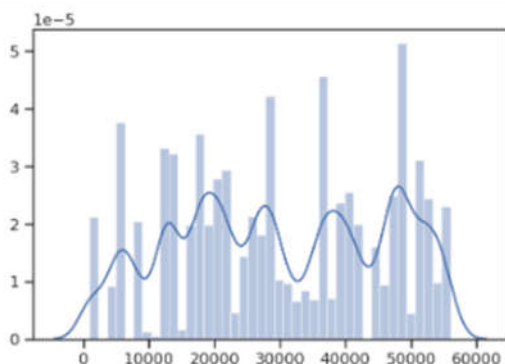
```
In [12]: mean_imp = sklearn.impute.SimpleImputer(strategy="mean")
mean_rat = mean_imp.fit_transform(data[["fips"]])
sns.distplot(mean_rat);
```



```
In [13]: med_imp = sklearn.impute.SimpleImputer(strategy="median")
med_rat = med_imp.fit_transform(data[["fips"]])
sns.distplot(med_rat);
```



```
In [14]: freq_imp = sklearn.impute.SimpleImputer(strategy="most_frequent")
freq_rat = freq_imp.fit_transform(data[["fips"]])
sns.distplot(freq_rat);
```



Кодирование категориальных признаков

```
In [15]: types = data["state"].dropna().astype(str)
types.value_counts()
```

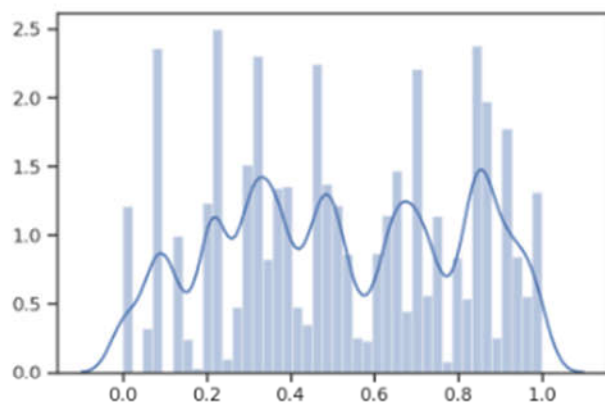
```

Out[15]: Texas 3877
         Georgia 3814
         Virginia 2706
         Indiana 2236
         North Carolina 2226
         Tennessee 2175
         Mississippi 2080
         Kentucky 2058
         Ohio 2052
         California 1990
         Florida 1896
         Missouri 1895
         Illinois 1869
         Michigan 1860
         New York 1780
         Pennsylvania 1724
         Iowa 1711
         Louisiana 1696
         Arkansas 1626
         Alabama 1597
         Minnesota 1590
         Colorado 1531
         Wisconsin 1516
         Oklahoma 1361
         Washington 1287
         South Carolina 1266
         Kansas 1258
         Oregon 850
         West Virginia 840
         Nebraska 839
         South Dakota 814
         New Jersey 763
         Idaho 723
         Maryland 718
         Montana 627
         North Dakota 613
         Utah 593
         New Mexico 591
         Massachusetts 559
         Arizona 490
         Wyoming 474
         Maine 420
         Vermont 403
         New Hampshire 316
         Nevada 283
         Connecticut 278
         Alaska 242
         Hawaii 165
         Rhode Island 162
         Delaware 104
         District of Columbia 41
         Puerto Rico 35
         Virgin Islands 34
         Guam 33
         Northern Mariana Islands 20
         Name: state, dtype: int64

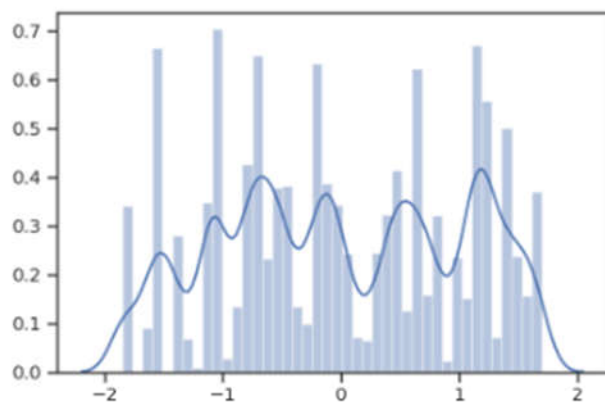
```


Масштабирование данных

```
In [19]: mm = sklearn.preprocessing.MinMaxScaler()  
sns.distplot(mm.fit_transform(data[["fips"]]]);
```



```
In [20]: ss = sklearn.preprocessing.StandardScaler()  
sns.distplot(ss.fit_transform(data[["fips"]]]);
```



Список литературы

- [1] Гапанюк Ю. Е. Лабораторная работа «Обработка пропусков в данных, кодирование категориальных признаков, масштабирование данных» [Электронный ресурс] // GitHub. — 2019. — Режим доступа: https://github.com/ugapanyuk/ml_course/wiki/LAB_MISSING (дата обращения: 05.04.2019).
- [2] Team The IPython Development. IPython 7.3.0 Documentation [Electronic resource] // Read the Docs. — 2019. — Access mode: <https://ipython.readthedocs.io/en/stable/> (online; accessed: 20.02.2019).
- [3] Waskom M. seaborn 0.9.0 documentation [Electronic resource] // PyData. — 2018. — Access mode: <https://seaborn.pydata.org/> (online; accessed: 20.02.2019).
- [4] pandas 0.24.1 documentation [Electronic resource] // PyData. — 2019. — Access mode: <http://pandas.pydata.org/pandas-docs/stable/> (online; accessed: 20.02.2019).
- [5] Gupta L. Google Play Store Apps [Electronic resource] // Kaggle. — 2019. — Access mode: <https://www.kaggle.com/lava18/google-play-store-apps> (online; accessed: 05.04.2019)