



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 вычислительная техника»  
Кафедра ИУ5 «Системы обработки информации и управления»

Курс «Технология машинного обучения»

Отчет по лабораторной работе №2  
«Обработка пропусков в данных, кодирование категориальных признаков и  
масштабирование данных»

Выполнил:  
студент группы ИУ5-62Б

Ванина П.В.

Подпись и дата:

Проверил:  
преподаватель каф.  
ИУ5

Гапанюк Ю.Е.

Подпись и дата:

Москва, 2022 г.

# Цель лабораторной работы:

Изучение способов предварительной обработки данных для дальнейшего формирования моделей.

## Описание задания:

- Выбрать набор данных (датасет)
- Для выбранного датасета выполнить следующие задачи:
  1. Обработку пропусков в данных
  2. Кодирование категориальных признаков
  3. Масштабирование данных

## Выполнение лабораторной работы:

Для выполнения данной лабораторной работы будем использовать данные об отзывах клиентов магазина, которые совершили покупки.  
Для начала импортируем Dataset и посмотрим на "шапку". Параметр sep зададим, чтобы использовать разделитель, отличный от стандартного - в нашем случае это запятая.

```
In [2]: import pandas as pd
import numpy as np
```

Отметим, что при таком импорте Dataset необходимо разместить файл с Dataset в том же месте (в той же директории и папке), что и проект Jupyter notebook.

```
In [3]: df = pd.read_csv('olist_order_reviews_dataset.csv', sep=',')

In [4]: df.head()
```

```
Out[4]:
```

	review_id	order_id	review_score	review_comment_title	review_comment_message	review_creation_date
0	7bc2406110b928393aa56f80a40eba40	73fc7af87114b39712e6da79b0a377eb	4	NaN	NaN	2018-01-18 00:00
1	80e541a11e59f04c1ad46dd5545f0fde	a548910a1c0147790b080df73d0eba33	5	NaN	NaN	2018-03-10 00:00
2	228ce500dc1d8e020d8d1322874b8f0	f0e4b658b201a9f2e0decbb34bed034b	5	NaN	NaN	2018-02-17 00:00
3	e84fb393a7b32834bb780f8bb30750e	658977c97b385a9be170737859d3511b	5	NaN	Recebi bem antes do prazo estipulado.	2017-04-21 00:00
4	f7c4243c7fe1938f181bec41a392bdeb	8e8bfb81e283fa7e4f11123a3fb894f1	5	NaN	Parabéns lojas lannister adorei comprar pela l...	2018-03-01 00:00

Для начала удалим дубликаты, так как дублирующие записи искажают показатели Dataset. Но перед удалением дубликатов обязательно узнаем, сколько записей находится в Dataset.

```
In [5]: df.shape

Out[5]: (99224, 7)
```

Посмотрим краткую информацию обо всех параметрах датасета

```
In [6]: df.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 99224 entries, 0 to 99223
Data columns (total 7 columns):
#   column              Non-Null Count  Dtype
---  ---
0   review_id           99224 non-null  object
1   order_id            99224 non-null  object
2   review_score        99224 non-null  int64
3   review_comment_title 11568 non-null  object
4   review_comment_message 48977 non-null  object
5   review_creation_date 99224 non-null  object
6   review_answer_timestamp 99224 non-null  object
dtypes: int64(1), object(6)
memory usage: 5.3+ MB
```

В выбранном датасете находятся преимущественно данные строкового типа, а также целочисленного. Проверим датасет на наличие дубликатов

```
In [7]: df = df.drop_duplicates()
df.shape

Out[7]: (99224, 7)
```

Видно, что Pandas не нашел ни одного дубликата

Обработка пропусков!  
Понимая, что если у признака более 70% пропусков, то такой признак удаляют. Поэтому проверим, насколько наши признаки полны.

```
In [8]: > column_values = df[['order_id', 'review_score', 'review_comment_title', 'review_comment_message', 'review_creation_date']].values
unique_values = pd.unique(column_values)
print(unique_values)

['73fc7af87114b39712e6da79b0a377eb' 4 nan ...
'725825d039fc1f0ceb7635e3f7d9286' '9053136becb1eccc2a1fbb265a0db0508'
'meu produto chegou e ja tenho que devolver, pois está com defeito , não segurar carga']
```

Существует несколько способов обозначить пропуски, и зачастую создатели датасета не описывают данные в достаточной мере, и определять, как обозначены пропуски, приходится вручную.

Например:

- 1) NaN / NaT (упрощенно: "не число" / "не время")
- 2) Пустая ячейка
- 3) Для числовых признаков – радикальный выброс. К примеру, для столбца "День" это число 999.
- 4) Маркер или нестандартный символ

Встроенные методы Pandas позволяют с легкостью справиться с первыми двумя разновидностями таких пробелов. Разберемся для начала с категориальными переменными, объединив их в один список.

```
In [9]: > df = df.replace({"День": 999,
                        "order_id": "Неизвестно",
                        "review_score": "Неизвестно",
                        "review_comment_title": "Неизвестно",
                        "review_comment_message": "Неизвестно",
                        "review_creation_date": "Неизвестно",
                        }, np.nan)

df.head()
```

```
Out[9]:
```

	review_id	order_id	review_score	review_comment_title	review_comment_message	review_creation_date
0	7bc2406110b920393aa50f0a40eba40	73fc7af87114b39712e6da79b0a377eb	4	NaN	NaN	2018-01-18 00:00
1	80e041a11e50f04c1ad40d05645f0fde	a548910a1c0147709b080f73d0eba33	5	NaN	NaN	2018-03-10 00:00
2	2280e5000dc1d8e020d8d1322874b0f0	f9e4b058b201a9f2ecdecbb34bed034b	5	NaN	NaN	2018-02-17 00:00
3	e04fb303e7b32834bb789f0bb30750e	058077c97b385a0be170737850d3511b	5	NaN	Recebi bem antes do prazo estipulado.	2017-04-21 00:00

Среди всех признаков 88% пропусков находится у признака "review\_comment\_title" и почти 59% - у "review\_comment\_message", первый признак поделжит удалению.

```
In [10]: > df.isnull().mean() * 100

Out[10]:
```

	review_id	order_id	review_score	review_comment_title	review_comment_message	review_creation_date	review_answer_timestamp
	0.000000	0.000000	0.000000	88.341530	58.702532	0.000000	0.000000
dtype:	float64						

```
In [11]: > df = df.drop(columns=['review_comment_title'])
df.head()

Out[11]:
```

	review_id	order_id	review_score	review_comment_message	review_creation_date	review_answer_timestamp
0	7bc2406110b920393aa50f0a40eba40	73fc7af87114b39712e6da79b0a377eb	4	NaN	2018-01-18 00:00:00	2018-01-18
1	80e041a11e50f04c1ad40d05645f0fde	a548910a1c0147709b080f73d0eba33	5	NaN	2018-03-10 00:00:00	2018-03-11
2	2280e5000dc1d8e020d8d1322874b0f0	f9e4b058b201a9f2ecdecbb34bed034b	5	NaN	2018-02-17 00:00:00	2018-02-18
3	e04fb303e7b32834bb789f0bb30750e	058077c97b385a0be170737850d3511b	5	Recebi bem antes do prazo estipulado.	2017-04-21 00:00:00	2017-04-21
4	f7c4243c7fe1938f181bec41a302b0eb	8e0fb081e283fa7e4f1123a3fb094f1	5	Parabéns lojas lannister adorei comprar pela l...	2018-03-01 00:00:00	2018-03-02

Процесс обработки пропусков, к счастью, можно сократить с помощью sklearn.impute.SimpleImputer. Мы выбираем все категориальные переменные и применяем стратегию "вставить вместо пропуска" самое распространенное значение":

```
In [12]: > from sklearn.impute import SimpleImputer

imputer = SimpleImputer(missing_values = np.nan, strategy = 'most_frequent')

df["order_id"] = imputer.fit_transform(df["order_id"].values.reshape(-1,1))[:,0]

df["review_comment_message"] = imputer.fit_transform(df["review_comment_message"].values.reshape(-1,1))[:,0]

df["review_creation_date"] = imputer.fit_transform(df["review_creation_date"].values.reshape(-1,1))[:,0]

df["review_answer_timestamp"] = imputer.fit_transform(df["review_answer_timestamp"].values.reshape(-1,1))[:,0]
```

Подобным образом заполняются пустоты в числовых переменных, только стратегия теперь – "вставить среднее значение".

```
In [13]: > imputer = SimpleImputer(missing_values = np.nan, strategy = 'mean')

df["review_score"] = imputer.fit_transform(df["review_score"].values.reshape(-1,1))[:,0]
```

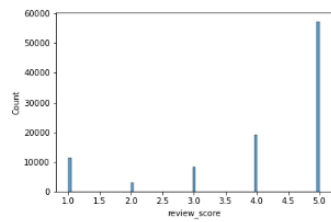
```
In [14]: > df.head()

Out[14]:
```

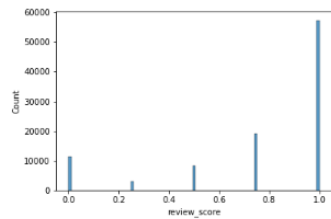
	review_id	order_id	review_score	review_comment_message	review_creation_date	review_answer_timestamp
0	7bc2406110b920393aa50f0a40eba40	73fc7af87114b39712e6da79b0a377eb	4.0	Muito bom	2018-01-18 00:00:00	2018-01-18
1	80e041a11e50f04c1ad40d05645f0fde	a548910a1c0147709b080f73d0eba33	5.0	Muito bom	2018-03-10 00:00:00	2018-03-11
2	2280e5000dc1d8e020d8d1322874b0f0	f9e4b058b201a9f2ecdecbb34bed034b	5.0	Muito bom	2018-02-17 00:00:00	2018-02-18
3	e04fb303e7b32834bb789f0bb30750e	058077c97b385a0be170737850d3511b	5.0	Recebi bem antes do prazo estipulado.	2017-04-21 00:00:00	2017-04-21
4	f7c4243c7fe1938f181bec41a302b0eb	8e0fb081e283fa7e4f1123a3fb094f1	5.0	Parabéns lojas lannister adorei comprar pela l...	2018-03-01 00:00:00	2018-03-02

# Масштабирование данных

```
In [15]: import seaborn as sns
sns.histplot(df['review_score'])
Out[15]: <AxesSubplot: xlabel='review_score', ylabel='Count'>
```



```
In [16]: from sklearn.preprocessing import MinMaxScaler, StandardScaler, Normalizer
df['review_score'] = MinMaxScaler().fit_transform(df['review_score'].values.reshape(-1, 1))
sns.histplot(df['review_score']);
```



# Кодирование категориальных признаков

```
In [24]: from sklearn.preprocessing import LabelEncoder
le = LabelEncoder()
le.fit(df.order_id)
df['order_id_le'] = le.transform(df.order_id)
df
```

```
Out[24]:
```

review_id	order_id	review_score	review_comment_message	review_creation_date	review_answer_timestamp	Hours	order_id_le
3a40eba40	73fc7af87114b39712e8da79b0a377eb	0.75	Muito bom	2018-01-18 00:00:00	2018-01-18 21:48:59	20	44992
35645fd8de	a548910a1c6147799b98df73dbaba33	1.00	Muito bom	2018-03-10 00:00:00	2018-03-11 03:05:13	20	63412
22874b8f0	f0e4b558b201a9f2e0ecbb34bed034b	1.00	Muito bom	2018-02-17 00:00:00	2018-02-18 14:38:24	20	96304
3bb30750e	659877c67b385a9be17073789d3511b	1.00	Recebi bem antes do prazo estipulado.	2017-04-21 00:00:00	2017-04-21 22:02:08	20	39089
1a392bdeb	8e8bfb81e283fa7e4f11123a3fb894f1	1.00	Parabéns lojas lannister adorei comprar pela i...	2018-03-01 00:00:00	2018-03-02 10:28:53	20	54800
...	...	...	...	...	...	...	...
bb7f9d7c9	2a9c23fee101d4d5992fa870396eb8da	1.00	Muito bom	2018-07-07 00:00:00	2018-07-14 17:18:30	20	16272
3d0f4ed7a	22ec9f069f784db00fa8b035cf8802	1.00	Muito bom	2017-12-09 00:00:00	2017-12-11 20:08:42	20	13435
6f3c302472	55d40047443889571d1f90031933e4	1.00	Excelente mochila, entrega super rápida. Super...	2018-03-22 00:00:00	2018-03-23 09:10:43	20	33099
33eb85149	7725825d039fc1f0ceb7935e3f7d9208	0.75	Muito bom	2018-07-01 00:00:00	2018-07-02 12:59:13	20	45831
3cd4cc548f	90531360eb1ee02a1fb2b285a0db0508	0.00	meu produto chegou e ja tenho que devolver, po...	2017-07-03 00:00:00	2017-07-03 21:01:49	20	55335

```
In [ ]:
```