

Министерство науки и высшего образования 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»

TMO PK №1

Горкунов Н.М. ИУ5-63Б 15 апреля 2024 г.

- 1 ТМО РК1 ИУ5-63Б Горкунов Николай
- 2 Задача №1.
- 2.1 Для заданного набора данных проведите корреляционный анализ. В случае наличия пропусков в данных удалите строки или колонки, содержащие пропуски. Сделайте выводы о возможности построения моделей машинного обучения и о возможном вкладе признаков в модель.
- 3 Набор данных №3.
- 3.1 Toy Dataset. A dataset to play around with!

```
[66]: import warnings
warnings.filterwarnings("ignore")
import pandas as pd
import numpy as np
import seaborn as sns
import time
from sklearn.datasets import make_classification
import matplotlib.pyplot as plt
from kaggle.api.kaggle_api_extended import KaggleApi
```

- [67]: Dataset URL: https://www.kaggle.com/datasets/carlolepelaars/toy-dataset
  - 3.2 Смотрю, что в данных

1 Dallas Male 41 40367.0

```
[68]: df = pd.read_csv('toy_dataset.csv')
print(df.shape)
    df.head()

(150000, 6)
[68]: Number City Gender Age Income Illness
```

```
54 45084.0
       2 Dallas
                 Male
                                       No
2
       3 Dallas
                 Male
                       42 52483.0
                                       No
3
       4 Dallas
                 Male
                        40 40941.0
                                       No
       5 Dallas
4
                 Male
                        46 50289.0
                                       No
```

#### 3.3 Сношу лишнее

```
[69]: del df["Number"]
     df.head()
[69]:
          City Gender
                      Age
                            Income Illness
     0 Dallas
                Male
                       41 40367.0
     1 Dallas
                Male
                       54 45084.0
                                       Νo
     2 Dallas
                Male
                       42 52483.0
                                       No
     3 Dallas
                Male
                       40 40941.0
                                       No
     4 Dallas
                Male
                       46 50289.0
                                       No
```

## 3.4 Проверяю типы данных, как и ожидалось, категориальные признаки – строки

#### 3.5 Проверяю значения категориальных признаков

```
[73]: df.Illness.unique()
[73]: array(['No', 'Yes'], dtype=object)
         Проверяю пропуски, их нет
[74]: df.isna().sum()
[74]: City
     Gender
                0
     Age
     Income
     Illness
     dtype: int64
          Выполняю требование: Если отсутствуют пропуски, замените на
          пропуски часть значений в одном или нескольких признаках
[75]: df.Income[df.Income < 10000] = np.nan
     df.isna().sum()
[75]: City
                 0
     Gender
                 0
     Age
                0
     Income
     Illness
     dtype: int64
          Обнаружил пропуски в численном признаке "Income", удаляю за-
     3.8
          писи с пропусками
[76]: df = df.dropna(axis=0, how="any")
     df.isna().sum()
[76]: City
                0
     Gender
                0
     Age
     Income
                0
     Illness
```

dtype: int64

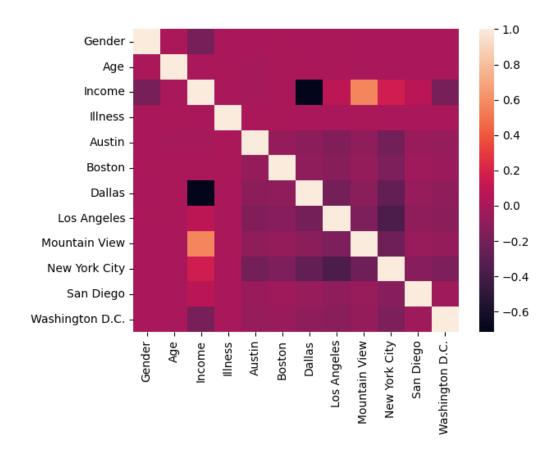
#### 3.9 Преобразую категориальные признаки (one hot encoding)

```
[77]: for to_enc in ["City"]:
         one_hot = pd.get_dummies(df[to_enc]).astype(int)
         del df[to_enc]
         df = df.join(one_hot)
      df.Illness = df.Illness.replace({'No': 0, 'Yes': 1})
      df.Gender = df.Gender.replace({'Male': 0, 'Female': 1})
      df.head()
[77]:
        Gender Age
                      Income Illness Austin Boston Dallas Los Angeles
                 41 40367.0
                                           0
                                                   0
     0
             0
                                   0
     1
             0 54 45084.0
                                                   0
                                                           1
                                                                        0
                                    0
                                           0
     2
             0 42 52483.0
                                   0
                                           0
                                                   0
                                                           1
                                                                        0
     3
             0 40 40941.0
                                  0
                                           0
                                                   0
                                                           1
                                                                        0
             0
                 46 50289.0
                                  0
                                            0
                                                   0
                                                           1
        Mountain View New York City San Diego Washington D.C.
     0
                    0
                                   0
                                             0
                                                              0
     1
                    0
                                   0
                                             0
                                                              0
     2
                    0
                                   0
                                             0
                                                              0
     3
                    0
                                   0
                                             0
                                                              0
                    0
                                   0
                                                              0
     4
```

# 3.10 Провожу корреляционный анализ, с натяжкой можно утверждать о зависимости дохода от проживания в г. Даллас

[78]:	Dallas Mountain View New York City  Los Angeles Income Washington D.C.	Income Income Los Angeles Dallas Mountain View Austin Dallas Gender Income	0.715624 0.567808 0.371263 0.276162 0.229916 0.212268 0.203141 0.198565 0.194212
	New York City	Boston	0.171963

dtype: float64



### 3.11 Выводы о возможности построения моделей машинного обучения и о возможном вкладе признаков в модель

Имеет смысл исследовать зависимость дохода от места проживания и пола, все прочие, как можно было убедиться выше, совсем слабо выражены. Таким образом можно попытаться построить модель, предсказывающую доход по месту проживания и полу. Однако, мало смысла строить такую модель, так как для неё нет различий между жителями одного города и пола.

# 3.12 Выполняю дополнительное требование: Для студентов групп ИУ5-63Б, ИУ5Ц-83Б - для произвольной колонки данных построить график "Ящик с усами (boxplot)"

[79]: sns.boxplot(y=df["Income"])

[79]: <Axes: ylabel='Income'>

