

Сорокин А.Э, ИУ5-24м, Вариант 11

Задача 11 Для набора данных проведите устранение пропусков для одного (произвольного) категориального признака с использованием метода заполнения наиболее распространенным значением.

Задача 31 Для набора данных проведите удаление повторяющихся признаков.

Решение

```
# оценим информацию в столбцах датасета
df.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 1309 entries, 0 to 1308
Data columns (total 14 columns):
#   Column      Non-Null Count  Dtype
---  -
0   pclass      1309 non-null   int64
1   survived    1309 non-null   int64
2   name        1309 non-null   object
3   sex         1309 non-null   object
4   age         1046 non-null   float64
5   sibsp       1309 non-null   int64
6   parch       1309 non-null   int64
7   ticket      1309 non-null   object
8   fare        1308 non-null   float64
9   cabin       295 non-null    object
10  embarked    1307 non-null   object
11  boat        486 non-null    object
12  body        121 non-null    float64
13  home.dest    745 non-null    object
dtypes: float64(3), int64(4), object(7)
memory usage: 143.3+ KB
```

Задача 1: Устранение пропусков для категориального признака

```
# получаем самое частоповторяемое значение в столбце
most_common = df['home.dest'].mode()[0]
most_common
```

[17]

... 'New York, NY'

```
# заполняем все пустые значения самым частым
df['home.dest'].fillna(most_common, inplace=True)
```

[18]

```
df.info()
```

[19]

```
... <class 'pandas.core.frame.DataFrame'>
RangeIndex: 1309 entries, 0 to 1308
Data columns (total 14 columns):
 #   Column      Non-Null Count  Dtype
---  -
 0   pclass      1309 non-null   int64
 1   survived    1309 non-null   int64
 2   name        1309 non-null   object
 3   sex         1309 non-null   object
 4   age         1046 non-null   float64
 5   sibsp       1309 non-null   int64
 6   parch       1309 non-null   int64
 7   ticket      1309 non-null   object
 8   fare        1308 non-null   float64
 9   cabin       295 non-null    object
10   embarked    1307 non-null   object
11   boat        486 non-null    object
12   body        121 non-null    float64
13   home.dest    1309 non-null   object
dtypes: float64(3), int64(4), object(7)
memory usage: 143.3+ KB
```

Задача 2: Удаление повторяющихся признаков

```
df.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 1309 entries, 0 to 1308
Data columns (total 14 columns):
#   Column      Non-Null Count  Dtype
---  -
0   pclass      1309 non-null   int64
1   survived    1309 non-null   int64
2   name        1309 non-null   object
3   sex         1309 non-null   object
4   age         1046 non-null   float64
5   sibsp       1309 non-null   int64
6   parch       1309 non-null   int64
7   ticket      1309 non-null   object
8   fare        1308 non-null   float64
9   cabin       295 non-null    object
10  embarked    1307 non-null   object
11  boat        486 non-null    object
12  body        121 non-null    float64
13  home.dest    1309 non-null   object
dtypes: float64(3), int64(4), object(7)
memory usage: 143.3+ KB
```

```
# Транспонирование DataFrame для превращения признаков в строки
df_transposed = df.T
df_transposed.head(3)
```

```
# Удаление дублирующих строк (признаков)
df_transposed.drop_duplicates(inplace=True)
```

```
# Обратное транспонирование для возвращения к исходному виду
df = df_transposed.T
```

```
df.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 1309 entries, 0 to 1308
Data columns (total 14 columns):
#   Column      Non-Null Count  Dtype
---  -
0   pclass      1309 non-null   object
1   survived    1309 non-null   object
2   name        1309 non-null   object
3   sex         1309 non-null   object
4   age         1046 non-null   object
5   sibsp       1309 non-null   object
6   parch       1309 non-null   object
7   ticket      1309 non-null   object
8   fare        1308 non-null   object
9   cabin       295 non-null    object
10  embarked    1307 non-null   object
11  boat        486 non-null    object
12  body        121 non-null    object
13  home.dest    1309 non-null   object
dtypes: object(14)
memory usage: 143.3+ KB
```