

**Московский государственный технический
университет им. Н.Э. Баумана**

Факультет «Информатика и системы управления»
Кафедра ИУ5 «Системы обработки информации и
управления»

Курс «Парадигмы и конструкции языков
программирования» Отчет по ДЗ

Выполнил:

студент группы ИУ5-33Б
каф.ИУ5
Захаров Ф.А.

Проверил:

преподаватель
Гапанюк Ю.Е.

Подпись и дата:

Подпись и дата:

Москва, 2023г.

Задача:

Задача взята с сайта <https://www.kaggle.com/competitions/titanic>,

Цель:

создание модели бинарной классификации, предсказывающей выжил ли пассажир Титаника или нет

Метрика:

точность(accuracy)

Решение:

- 1) Заполнить пропуски
- 2) Закодировать категориальные признаки
- 3) Визуализировать данные
- 4) Приведение данных к нормальному распределению
- 5) Тренировка моделей
- 6) Подбор гиперпараметров

1) и 2)

```
1 dict_ = {'male':True, 'female':False}
2 train['is_male'] = train['Sex'].replace(dict_)
3 test['is_male'] = test['Sex'].replace(dict_)
4
5 train.Embarked.fillna('S', inplace=True)
6
7 df = pd.concat([train,test])
8 encoder = OneHotEncoder()
9 encoder.fit(df.Embarked.to_numpy().reshape(-1,1))
10 temp = encoder.transform(train.Embarked.to_numpy().reshape(-1,1))
11 temp = pd.DataFrame(temp.astype(int).toarray())
12 for i, u in enumerate(list(encoder.categories_[0])):
13     train[u] = temp[i]
14     train[u] = train[u].astype(bool)
15 temp = encoder.transform(test.Embarked.to_numpy().reshape(-1,1))
16 temp = pd.DataFrame(temp.astype(int).toarray())
17 for i, u in enumerate(list(encoder.categories_[0])):
18     test[u] = temp[i]
19     test[u] = test[u].astype(bool)
20
21 temp = encoder.transform(df.Embarked.to_numpy().reshape(-1,1))
22 temp = pd.DataFrame(temp.astype(int).toarray())
23 for i, u in enumerate(list(encoder.categories_[0])):
24     df[u] = temp[i]
25     df[u] = df[u].astype(bool)
```

```

1 train['stat'] = train.Name.str.split(',').apply(lambda x:x[1].split('.')[0])
2 test['stat'] = test.Name.str.split(',').apply(lambda x:x[1].split('.')[0])
3
4 fullnames = ['Mister', 'Missis', 'Miss', 'Master', 'Don', 'Reverend', 'Doctor', 'Madame',
5             'Miss', 'Major', 'Lady', 'Sir', 'Madmuaselle', 'Colonel', 'Captain',
6             'Countess', 'Jonkheer']
7 dict_ = dict(zip(train.stat.unique(), fullnames))
8
9 train['stat'] = train.stat.replace(dict_)
10 test['stat'] = test.stat.replace(dict_)
11
12 train['family_size'] = train.SibSp + train.Parch
13 test['family_size'] = test.SibSp + test.Parch
14
15 dict_ = {'Mister': 'Sir',
16         'Missis': 'Lady',
17         'Miss': 'Kid',
18         'Master': 'Kid',
19         'Don': 'Sir',
20         'Reverend': 'Rare',
21         'Doctor': 'Rare',
22         'Madame': 'Lady',
23         'Major': 'Rare',
24         'Madmuaselle': 'Kid',
25         'Colonel': 'Rare',
26         'Captain': 'Rare',
27         'Countess': 'Lady',
28         'Jonkheer': 'Rare',
29         'Dona': 'Lady'
30        }
31
32 train['is_civil'] = train.stat.replace(dict_)
33 test['is_civil'] = test.stat.replace(dict_)
34
35
36 train = pd.concat([train.drop('is_civil', axis=1), pd.get_dummies(train['is_civil'])], axis=1)
37 test = pd.concat([test.drop('is_civil', axis=1), pd.get_dummies(test['is_civil'])], axis=1)

```

```

1 X['deck'] = X.Cabin.str[0]
2
3 temp = X.Ticket.value_counts().loc[X.Ticket.value_counts()>1]
4 for i in temp.index:
5     t = X.loc[(X.Ticket==i)]
6     if (t.Cabin.count()!=t.shape[0])&(t.Cabin.count()!=0):
7         id = t.loc[t.Cabin.isna()].index
8         X.loc[id, 'Cabin'] = t.Cabin.mode().iloc[0]
9         X.loc[id, 'deck'] = t.deck.mode().iloc[0]
10
11
12
13 X['deck'] = X['Cabin'].apply(lambda s: s[0] if pd.notnull(s) else 'M')
14 X['deck'] = np.where((X.deck == 'T'),'A',X.deck)
15
16 X['deck'].replace(['A', 'B', 'C'], 'ABC', inplace=True)
17 X['deck'].replace(['D', 'E'], 'DE', inplace=True)
18 X['deck'].replace(['F', 'G'], 'FG', inplace=True)
19
20 train = X.iloc[:train.shape[0]]
21 test = X.drop('Survived',axis=1).iloc[train.shape[0]:]

```

```

1 def size(x):
2     if x<1: return 'alone'
3     elif x>4: return '1-3'
4     else: return '4+'
5 train['family_cat'] = train.family_size.apply(size)
6 test['family_cat'] = test.family_size.apply(size)
7
8
9 X = pd.concat([train,test]).drop('Survived', axis=1)
10 temp = pd.get_dummies(X.family_cat)
11
12 train = pd.concat([train, temp.iloc[:train.shape[0]]], axis=1)
13 test = pd.concat([test, temp.iloc[train.shape[0]:]], axis=1)
14

```

```

1 X = pd.concat([train,test]).drop('Survived', axis=1)
2 temp = pd.get_dummies(X.Pclass)
3
4 train = pd.concat([train, temp.iloc[:train.shape[0]]], axis=1)
5 test = pd.concat([test, temp.iloc[train.shape[0]:]], axis=1)
6
7 train = train.rename({1:'1st class', 2:'2nd class', 3:'3rd class'}, axis=1)
8 test = test.rename({1:'1st class', 2:'2nd class', 3:'3rd class'}, axis=1)

```

```

1 X = pd.concat([train,test]).drop('Survived', axis=1)[['Ticket', 'Fare']]
2 tickets_count = X.Ticket.value_counts().to_dict()
3
4 X['tickcount'] = X.Ticket.replace(tickets_count)
5 X['Fare'] = X['Fare']/X['tickcount']
6 train['Fare'] = X['Fare'].iloc[:train.shape[0]]
7 test['Fare'] = X['Fare'].iloc[train.shape[0]:]
8 # test.loc[test.loc[test.Fare.isna()].index, 'Age'] = df.loc[(df.Pclass==3)&(df.Sex=='male')].Age.median()
9 test['Fare'].fillna(df.loc[(df.Pclass==3)&(df.Sex=='male')].Age.median(), inplace = True)

```

```

1 df = pd.concat([train,test])
2 for k in df.is_male.unique():
3     for i in df.Pclass.unique():
4         id = df.loc[(df.Pclass==i)&(df.Age.isna())&(df.is_male==k)].index
5         temp = df.loc[(df.Pclass == i)&(df.is_male==k)].Age
6         std = temp.std()
7         mean = temp.mean()
8         repl_arr = np.random.normal(mean,len(id))
9         df.loc[id, 'Age'] = repl_arr
10
11 df['Age'] = df.Age.round()
12 train['Age'] = df.Age.iloc[:train.shape[0]].astype(int)
13 test['Age'] = df.Age.iloc[train.shape[0]:].astype(int)

```

```

1 X = pd.concat([X, pd.get_dummies(X.deck)], axis=1)
2 train['Age'] = X.Age.iloc[:train.shape[0]].astype(int)
3 test['Age'] = X.Age.iloc[train.shape[0]:].astype(int)

```

```

1 X = pd.concat([X, pd.get_dummies(X.deck)], axis=1)
2 train = X.iloc[:train.shape[0]]
3 test = X.iloc[train.shape[0]:]
4 lst = ['Sex', 'Embarked', 'Name', 'Ticket', 'Cabin', 'stat', 'family_cat', 'Pclass',
5        'SibSp', 'Parch', 'Sir', 'deck', 'family_size']
6 train.drop(lst,axis=1, inplace=True)
7 test.drop(lst,axis=1, inplace=True)

```

4)Привидение к нормальному виду

scaling

```

1 lst = ['Age', 'Fare']#, 'family_size']
2 scaler = StandardScaler()
3 temp = scaler.fit_transform(train[lst])
4 temp = pd.DataFrame(temp, columns=lst).set_index(train.index)
5 train = pd.concat([train.drop(lst, axis=1), temp], axis=1)
6
7 temp = scaler.transform(test[lst])
8 temp = pd.DataFrame(temp, columns=lst).set_index(test.index)
9 test = pd.concat([test.drop(lst, axis=1), temp], axis=1)

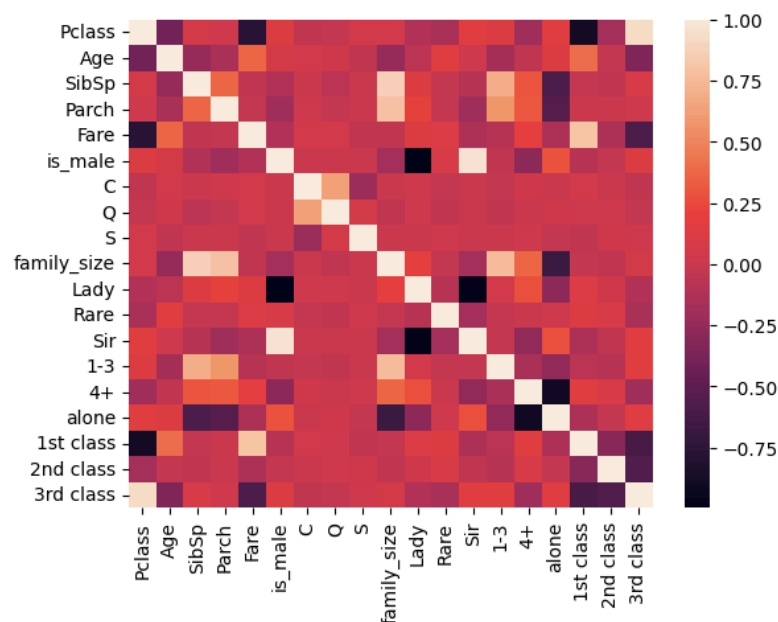
```

```
1 df = pd.concat([train.drop('Survived', axis=1), test])
2 sns.heatmap(df.corr())
```

<ipython-input-39-cc734fa1c52f>:2: FutureWarning: The default value of numeric_only in DataFrame.corr is deprecated. In a future version, it will default to False. Select only valid columns or specify the value of numeric_only to silence this warning.

```
sns.heatmap(df.corr())
```

<Axes: >

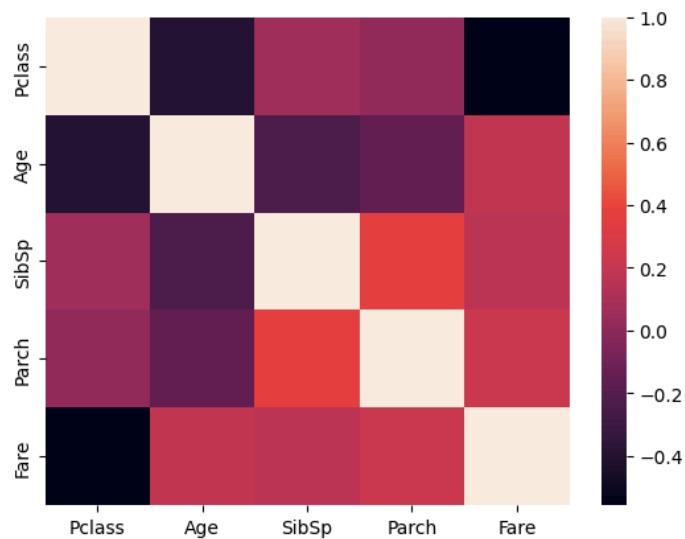


```
1 sns.heatmap(df.corr())
```

<ipython-input-155-aa4f4450a243>:1: FutureWarning: The default value of numeric_only in DataFrame.corr is deprecated. In a future version, it will default to False. Select only valid columns or specify the value of numeric_only to silence this warning.

```
sns.heatmap(df.corr())
```

<Axes: >



```

1: 1 train['stat'] = train.Name.str.split(',').apply(lambda x:x[1].split('.')[0])
2 test['stat'] = test.Name.str.split(',').apply(lambda x:x[1].split('.')[0])
3
4 fullnames = ['Mister', 'Missis', 'Miss', 'Master', 'Don', 'Reverend', 'Doctor', 'Madame',
5             'Miss', 'Major', 'Lady', 'Sir', 'Madmuaselle', 'Colonel', 'Captain',
6             'Countess', 'Jonkheer']
7 dict_ = dict(zip(train.stat.unique(), fullnames))
8
9 train['stat'] = train.stat.replace(dict_)
10 test['stat'] = test.stat.replace(dict_)
11
12 df = pd.concat([train, test])
13
14 for i in df.loc[df.Age.isna()].stat.unique():
15     id = df.loc[(df.stat==i)&(df.Age.isna())].index
16     if i == 'Doctor':
17         df.loc[id, 'Age'] = df.loc[df.stat == i].Age.median()
18     else:
19         temp = df.loc[df.stat == i].Age
20         if temp.isna().sum() != temp.shape[0]:
21             std = temp.std()
22             mean = temp.mean()
23         else:
24             std = df.Age.std()
25             mean = df.Age.mean()
26         # id = df.loc[(df.stat==i)&(df.Age.isna())].index
27         repl_arr = np.random.normal(mean, std, len(id))
28         while (repl_arr < 2).sum() != 0:
29             repl_arr = np.random.normal(mean, std, len(id))
30         # repl_arr = [k if k > 1 else mean for k in repl_arr]
31         df.loc[id, 'Age'] = repl_arr
32 df['Age'] = df.Age.round()
33
34 train['Age'] = df.Age.iloc[:train.shape[0]].astype(int)
35 test['Age'] = df.Age.iloc[train.shape[0]:].astype(int)

```

```

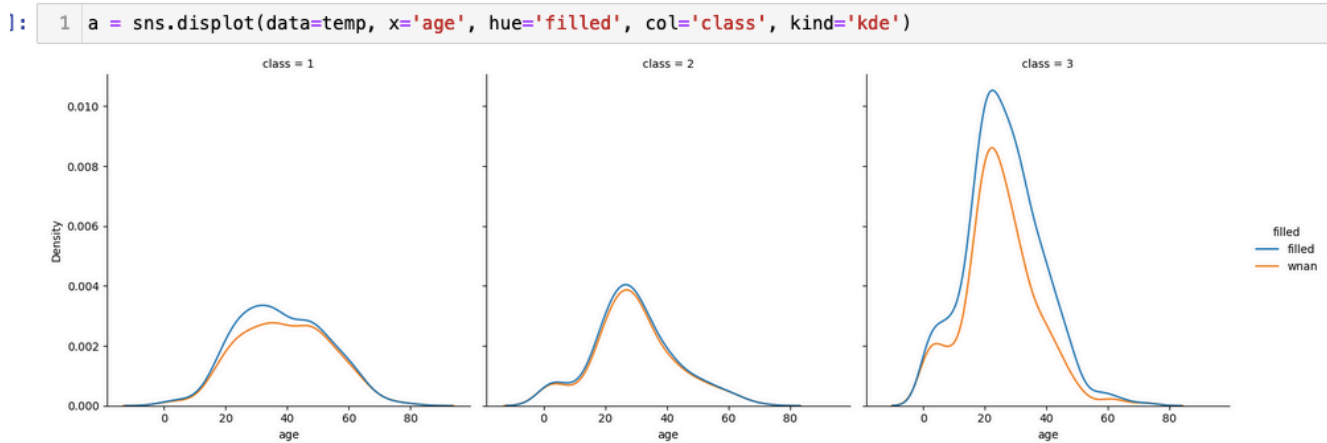
1: 1 df = pd.concat([train, test])
2
3 for i in df.Pclass:
4     id = df.loc[(df.Pclass==i)&(df.Age.isna())].index
5     temp = df.loc[df.Pclass == i].Age
6     std = temp.std()
7     mean = temp.mean()
8     repl_arr = np.random.normal(mean, std, len(id))
9     df.loc[id, 'Age'] = repl_arr
10 df['Age'] = df.Age.round()
11
12 train['Age'] = df.Age.iloc[:train.shape[0]].astype(int)
13 test['Age'] = df.Age.iloc[train.shape[0]:].astype(int)

```

```

1: 1 train_age = pd.concat([train['Age'], train_bu.Age, ], axis=1)
2 test_age = pd.concat([test['Age'], test_bu.Age], axis=1)
3
4 train_age.columns=['filled_train', 'wnan_train']# 'Sex',
5 test_age.columns=['filled_test', 'wnan_test']
6
7 sex_col = pd.concat([train.Sex, train.Sex, test.Sex, test.Sex]).reset_index(drop=True)
8 class_col = pd.concat([train.Pclass, train.Pclass, test.Pclass, test.Pclass]).reset_index(drop=True)
9
10
11 d1 = {'filled_train': 'train',
12       'filled_test': 'test',
13       'wnan_train': 'train' }
14 d2 = {'filled_train': 'filled',
15       'filled_test': 'filled',
16       'wnan_test': 'wnan',
17       'wnan_train': 'wnan'}
18
19 temp = pd.concat([train_age.melt(), test_age.melt()]).reset_index(drop=True)
20 temp = pd.concat([temp['value'], temp['variable'].replace(d1), temp['variable'].replace(d2), sex_col, class_col],
21                 )
22 temp.columns = ['age', 'subset', 'filled', 'sex', 'class']

```



```
1 t1 = pd.Series(np.random.choice(t, a1.shape))
2 a1.reset_index(drop=True, inplace=True)

1 def get_best_distribution(data):
2     dist_names = ["norm", "exponweib", "weibull_max", "weibull_min", "pareto", "genextreme"]
3     dist_results = []
4     params = {}
5     for dist_name in dist_names:
6         dist = getattr(st, dist_name)
7         param = dist.fit(data)
8
9         params[dist_name] = param
10        # Applying the Kolmogorov-Smirnov test
11        D, p = st.kstest(data, dist_name, args=param)
12        print("p value for "+dist_name+" = "+str(p))
13        dist_results.append((dist_name, p))
14
15    # select the best fitted distribution
16    best_dist, best_p = (max(dist_results, key=lambda item: item[1]))
17    # store the name of the best fit and its p value
18
19    print("Best fitting distribution: "+str(best_dist))
20    print("Best p value: "+ str(best_p))
21    print("Parameters for the best fit: "+ str(params[best_dist]))
22
23    return best_dist, best_p, params[best_dist]
24 get_best_distribution(data.loc[~data.isna()])
```

```
p value for norm = 4.090734030520837e-06
p value for exponweib = 0.0010394734025497118
p value for weibull_max = 0.00036793342881385316
p value for weibull_min = 9.937416623563527e-06
p value for pareto = 1.749205034758647e-86
p value for genextreme = 0.0003678772194694297
Best fitting distribution: exponweib
Best p value: 0.0010394734025497118
Parameters for the best fit: (15.26537334292596, 2.860196772430548, -82.61306733997591, 74.93359349480347)

('exponweib',
 0.0010394734025497118,
 (15.26537334292596, 2.860196772430548, -82.61306733997591, 74.93359349480347))
```

```

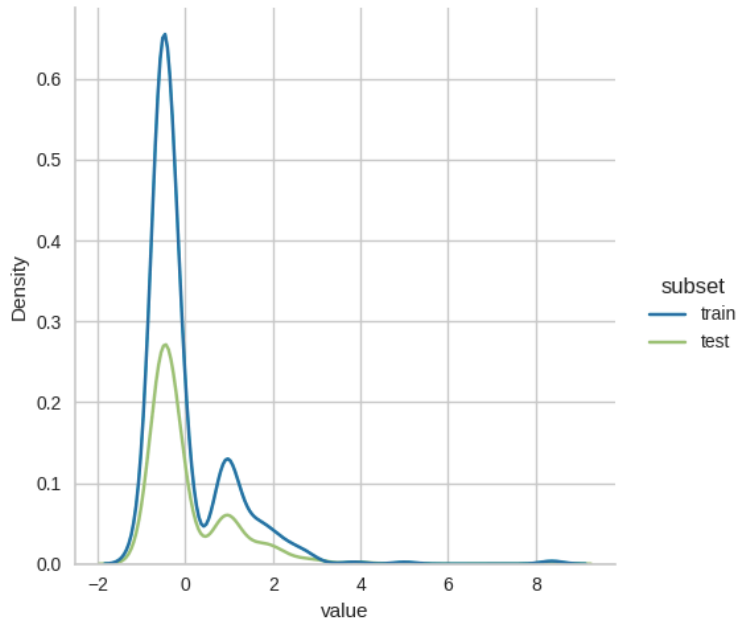
]: 1 a = pd.Series(['train' if i<test.shape[0] else 'test' for i in range(temp.shape[0])])
2    temp = pd.concat([pd.Series(np.random.choice(train.Fare, test.Fare.shape[0])), test.Fare])
3
4    temp = pd.concat([train.Fare, test.Fare]).reset_index(drop=True)
5    a = pd.Series(['train' if i<train.shape[0] else 'test' for i in range(temp.shape[0])])
6    temp = pd.concat([temp, a], axis=1)
7    temp.columns = ['value', 'subset']
8    sns.displot(data=temp, x='value', kind='kde', hue='subset', )

```

```

]: <seaborn.axisgrid.FacetGrid at 0x79780921d870>

```

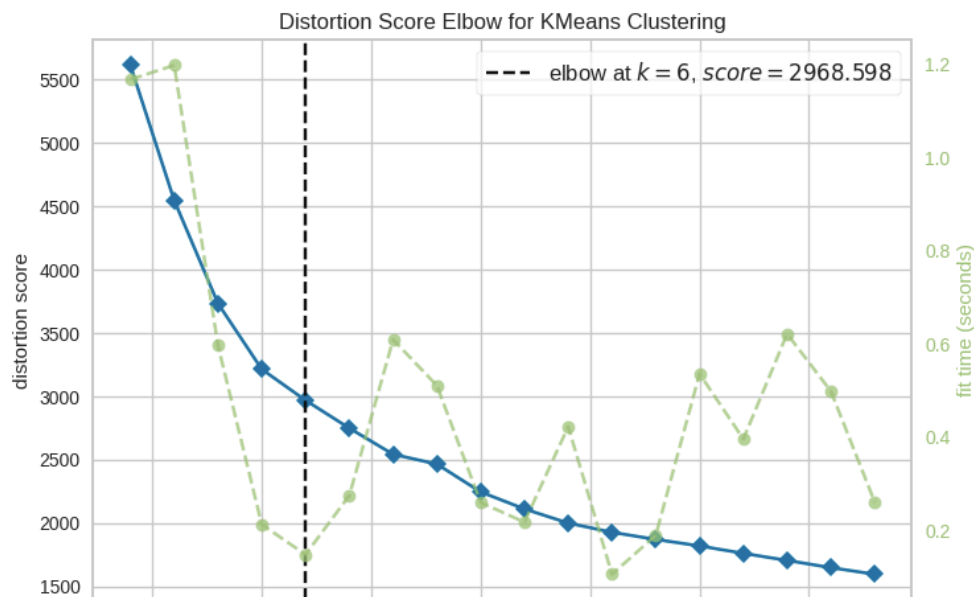


clustering vis

```

]: 1 from yellowbrick.cluster import KElbowVisualizer
2
3    X = pd.concat([train.drop('Survived', axis=1), test])
4
5    km = KMeans(random_state=42)
6    visualizer = KElbowVisualizer(km, k=(2,20))
7
8    visualizer.fit(X)          # Fit the data to the visualizer
9    visualizer.show()

```



Подбор моделей, обучение

testing models

```
яд [ ]: 1 x_train, x_test, y_train, y_test = train_test_split(train.drop('Survived', axis=1), train.Survived,
2                                                    test_size=.25, random_state=42, shuffle=True,
3                                                    stratify=train.Survived)
```

```
яд [ ]: 1 model_1 = LogisticRegression(max_iter=1000)
2 model_1.fit(x_train, y_train)
3 y_preds = model_1.predict(x_test)
4 # result_list[0] = (f'f1: {round(f1_score(y_test, y_preds), 4)}\naccuracy: {round(accuracy_score(y_test, y_preds),
5 print(f'f1: {round(f1_score(y_test, y_preds), 4)}\naccuracy: {round(accuracy_score(y_test, y_preds), 4)}'))
```

f1: 0.75
accuracy: 0.8117

```
яд [ ]: 1 model_2 = RandomForestClassifier(random_state=42)
2 model_2.fit(x_train, y_train)
3 y_preds = model_2.predict(x_test)
4 print(f'f1: {round(f1_score(y_test, y_preds), 4)}\naccuracy: {round(accuracy_score(y_test, y_preds), 4)}'))
```

f1: 0.6795
accuracy: 0.7758

```
яд [ ]: 1 model_3 = DecisionTreeClassifier(random_state=42)
2 model_3.fit(x_train, y_train)
3 y_preds = model_3.predict(x_test)
4 print(f'f1: {round(f1_score(y_test, y_preds), 4)}\naccuracy: {round(accuracy_score(y_test, y_preds), 4)}'))
```

f1: 0.7232
accuracy: 0.7803

```
яд [ ]: 1 model_4 = LinearSVC(max_iter=2000, random_state=42)
2 model_4.fit(x_train, y_train)
3 y_preds = model_4.predict(x_test)
4 print(f'f1: {round(f1_score(y_test, y_preds), 5)}\naccuracy: {round(accuracy_score(y_test, y_preds), 5)}'))
5 # print(result_list[3])
```

```
яд [ ]: 1 from xgboost import XGBClassifier, XGBRFClassifier
2 model_8 = XGBClassifier(random_state=42)
3 model_8.fit(x_train, y_train)
4 y_preds = model_8.predict(x_test)
5 print(f'XGBC /nf1: {round(f1_score(y_test, y_preds), 5)}\naccuracy: {round(accuracy_score(y_test, y_preds), 5)}'))
6
7 model_9 = XGBRFClassifier(random_state=42)
8 model_9.fit(x_train, y_train)
9 y_preds = model_9.predict(x_test)
10 print(f'XGBRFC /nf1: {round(f1_score(y_test, y_preds), 5)}\naccuracy: {round(accuracy_score(y_test, y_preds), 5)}'))
```

XGBC /nf1: 0.75449
accuracy: 0.81614
XGBRFC /nf1: 0.725
accuracy: 0.80269

```
яд [ ]: 1 from sklearn.linear_model import SGDClassifier
2 model_10 = SGDClassifier(random_state=42)
3 model_10.fit(x_train, y_train)
4 y_preds = model_10.predict(x_test)
5 print(f'SGDC /nf1: {round(f1_score(y_test, y_preds), 5)}\naccuracy: {round(accuracy_score(y_test, y_preds), 5)}'))
```

SGDC /nf1: 0.74713
accuracy: 0.80269

```
яд [ ]: 1 from sklearn.linear_model import PassiveAggressiveClassifier
2 model_11 = PassiveAggressiveClassifier(random_state=42)
3 model_11.fit(x_train, y_train)
4 y_preds = model_11.predict(x_test)
5 print(f'SGDC /nf1: {round(f1_score(y_test, y_preds), 5)}\naccuracy: {round(accuracy_score(y_test, y_preds), 5)}'))
```

SGDC /nf1: 0.71628
accuracy: 0.72646

```
1 [ ]: 1 from sklearn.linear_model import RidgeClassifier
2       model_12 = RidgeClassifier(random_state=42)
3       model_12.fit(x_train, y_train)
4       y_preds = model_12.predict(x_test)
5       print(f'SGDC /nf1: {round(f1_score(y_test,y_preds),5)}\naccuracy: {round(accuracy_score(y_test,y_preds), 5)}')
```

SGDC /nf1: 0.75862
accuracy: 0.81166

```
1 [ ]: 1 from sklearn.neural_network import MLPClassifier
2       model_13 = MLPClassifier(random_state=42)
3       model_13.fit(x_train, y_train)
4       y_preds = model_13.predict(x_test)
5       print(f'SGDC /nf1: {round(f1_score(y_test,y_preds),5)}\naccuracy: {round(accuracy_score(y_test,y_preds), 5)}')
```

SGDC /nf1: 0.70064
accuracy: 0.78924

```
1 [ ]: 1 from sklearn.svm import NuSVC
2       model_14 = NuSVC(random_state=42, probability=True)
3       model_14.fit(x_train, y_train)
4       y_preds = model_14.predict(x_test)
5       print(f'SGDC /nf1: {round(f1_score(y_test,y_preds),5)}\naccuracy: {round(accuracy_score(y_test,y_preds), 5)}')
```

SGDC /nf1: 0.76471
accuracy: 0.82063

```
1 [ ]: 1 from catboost import CatBoostClassifier
2       model_15 = CatBoostClassifier(random_state=42, verbose=False)
3       model_15.fit(x_train, y_train)
4       y_preds = model_15.predict(x_test)
5       print(f'catboost /nf1: {round(f1_score(y_test,y_preds),5)}\naccuracy: {round(accuracy_score(y_test,y_preds), 5)}')
```

Collecting catboost
catboost /nf1: 0.73418
accuracy: 0.81166

```
1 [ ]: 1 from lightgbm import LGBMClassifier
2       model_16 = LGBMClassifier(num_leaves=25, verbose=-1, random_state=42)
3       model_16.fit(x_train, y_train)
4       y_preds = model_16.predict(x_test)
5       print(f'catboost /nf1: {round(f1_score(y_test,y_preds),5)}\naccuracy: {round(accuracy_score(y_test,y_preds), 5)}')
```

catboost /nf1: 0.73054
accuracy: 0.79821

```
1 [ ]: 1 ests = [('logreg', model_1), ('tree', model_3), ('linsvc', model_4), ('gbc', model_5), ('knn', model_7),
2           ('rf', model_2), ('abc', model_6), ('XGBC', model_8), ('XGBRFC', model_9), ('SGDC', model_10), ('PAC', model_11),
3           ('RC', model_12), ('MLPC', model_13), ('NuSVC', model_14), ('cbc', model_15), ('lightgbm', model_16)]
4       ensemble_1 = VotingClassifier(estimators=ests)#, voting='soft')
5       ensemble_1.fit(x_train, y_train)
6       y_preds = ensemble_1.predict(x_test)
7       print(f'f1: {round(f1_score(y_test,y_preds),5)}\naccuracy: {round(accuracy_score(y_test,y_preds), 5)}')
```

f1: 0.77576
accuracy: 0.83408

```
1 [ ]: 1 # ests = [('logreg', model_1), ('tree', model_3), ('linsvc', model_4), ('gbc', model_5), ('knn', model_7)]#, ('abc', model_6), ('XGBC', model_8), ('XGBRFC', model_9), ('MLPC', model_13), ('NuSVC', model_14), ('cbc', model_15), ('lightgbm', model_16)]
2       ensemble_2 = StackingClassifier(estimators=ests)#, final_estimator=model_4)
3       ensemble_2.fit(x_train, y_train)
4       y_preds = ensemble_2.predict(x_test)
5       print(f'f1: {round(f1_score(y_test,y_preds),5)}\naccuracy: {round(accuracy_score(y_test,y_preds), 5)}')
```

f1: 0.7619
accuracy: 0.82063

```
1 [ ]: 1 from sklearn.naive_bayes import CategoricalNB
2       ensemble_3 = StackingClassifier(estimators=ests, final_estimator=model_1)
3       ensemble_3.fit(x_train, y_train)
4       y_preds = ensemble_3.predict(x_test)
5       print(f'f1: {round(f1_score(y_test,y_preds),5)}\naccuracy: {round(accuracy_score(y_test,y_preds), 5)}')
```

f1: 0.76836
accuracy: 0.81614

```
1 [ ]: 1 ests = [('logreg', model_1), ('tree', model_3), ('gbc', model_5), ('knn', model_7),
2           ('rf', model_2), ('abc', model_6), ('XGBC', model_8), ('XGBRFC', model_9), ('MLPC', model_13), ('NuSVC', model_14), ('cbc', model_15), ('lightgbm', model_16)]
3       ensemble_4 = VotingClassifier(estimators=ests, voting='soft')
4       ensemble_4.fit(x_train, y_train)
5       y_preds = ensemble_4.predict(x_test)
6       print(f'f1: {round(f1_score(y_test,y_preds),5)}\naccuracy: {round(accuracy_score(y_test,y_preds), 5)}')
```

f1: 0.76829
accuracy: 0.8296

```
1 [ ]: 1 from sklearn.ensemble import BaggingClassifier
2       ensemble_5 = BaggingClassifier(estimator=model_2)
3       ensemble_5.fit(x_train, y_train)
4       y_preds = ensemble_5.predict(x_test)
5       print(f'f1: {round(f1_score(y_test,y_preds),5)}\naccuracy: {round(accuracy_score(y_test,y_preds), 5)}')
```

f1: 0.74699
accuracy: 0.81166

6) Подбор гиперпараметров

optimizers

```
1 [ ]: 1 model_5 = GradientBoostingClassifier(cv.best_params_)
2      2 model_5.fit(x_train, y_train)
3      3 y_preds = model_5.predict(x_test)
4      4 print(f'f1: {round(f1_score(y_test, y_preds), 5)} \n accuracy: {round(accuracy_score(y_test, y_preds), 5)}')
```

...

```
1 [ ]: 1 gbc_params = {
2      2     'loss': ['exponential', 'log_loss'],
3      3     'learning_rate': [.07, .06, .05, .04, .03 ],
4      4     'n_estimators': [i for i in range(83, 89, 2)],
5      5     'criterion': ['friedman_mse', 'squared_error'],
6      6     'tol': [1e-4, 9e-5, 2e-4]
7      7 }
8
9      8 cv = GridSearchCV(model_5, param_grid= gbc_params, scoring='accuracy')
10
11     9 cv.fit(x_train, y_train)
12     10 cv.best_score_
```

[389]: 0.8278083267871171

```
1 [ ]: 1 model_1 = LogisticRegression()
2      2 model_1.fit(x_train, y_train)
3      3 y_preds = model_1.predict(x_test)
4      4 # result_list[0] = (f'f1: {round(f1_score(y_test, y_preds), 4)} \n accuracy: {round(accuracy_score(y_test, y_preds),
5      5     print(f'f1: {round(f1_score(y_test, y_preds), 4)} \n accuracy: {round(accuracy_score(y_test, y_preds), 4)}')
```

```
1 [ ]: 1 cv.best_params_
```

[12]: {'max_iter': 100, 'penalty': 'l1', 'solver': 'saga', 'tol': 0.0001}

```
1 [ ]: 1 lr_params = {
2      2     'penalty': ['l1', 'l2', 'elasticnet'],
3      3     'tol': [1e-4, 5e-5, 1e-5, 1e-3, 5e-4],
4      4     'solver': ['lbfgs', 'liblinear', 'newton-cg', 'newton-cholesky', 'sag', 'saga'],# ['liblinear', 'saga'],
5      5     'max_iter': [10, 15, 20, 25, 50, 60, 70, 80, 90, 100, 150, 200, 250, 300, 350, 400]
6      6 }
7
8      7 model = LogisticRegression()
9      8 cv = HalvingGridSearchCV(model, param_grid= lr_params, scoring='accuracy', verbose=2)
10     9 cv.fit(train.drop(['Survived'], axis=1), train.Survived)
```

```
1 [ ]: 1 cv.best_score_
```

[97]: 0.7984942886812045

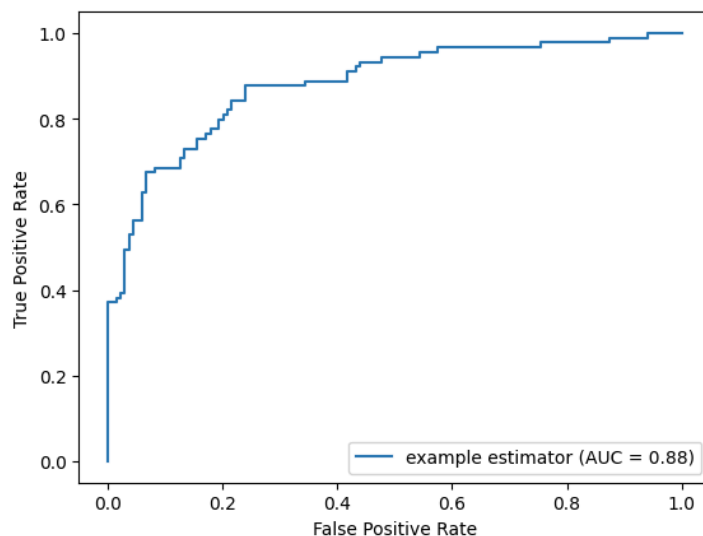
```
1 [ ]: 1 cv.best_params_
```

[52]: {'n_neighbors': 12, 'weights': 'uniform'}

```

1 from sklearn import metrics
2
3 y_score = model_1.predict_proba(x_test)
4
5 fpr, tpr, _ = metrics.roc_curve(y_test, y_score[:,1])
6
7 roc_auc = metrics.auc(fpr, tpr)
8 display = metrics.RocCurveDisplay(fpr=fpr, tpr=tpr, roc_auc=roc_auc,
9                                   estimator_name='example estimator')
10 display.plot()
11
12 plt.show()
13

```



Результат:

Несмотря на хорошие результаты большого количества моделей на тренировочной выборке. На тестовой выборке лучше всего себя проявила Логистическая регрессия с минимумом обработка данных, 0.78468

Модель	Результат
Логистическая регрессия	0.78468
KNN	0.77272
Градиентный бусины	0.76076
Стэкинг	0.75598