# УНИВЕРСИТЕТ ИТМО

# Практика 9. Генетические алгоритмы: отбор признаков и оптимизация гиперпараметров

**Михаил А. Каканов**[1]     Олег А. Евстафьев[1]

[1]Факультет систем управления и робототехники, Университет ИТМО
{makakanov,oaevstafev}@itmo.ru

1. Подготовка среды

2. Разработка модели

3. Валидация

4. Оптимизация модели

5. Анализ результатов

# Содержание

В этом блокноте будет использоваться набор данных 'Santander Customer Satisfaction', поскольку он имеет высокую размерность (300+) и подходит для тестирования процесса оптимизации.

Подготовим среду для построения оптимизационной модели и поиска хорошей прогностической модели (признаки и гиперпараметры) для данного набора данных. В данной работе используются следующие пакеты:

- ▶ Numpy
- ▶ Pandas
- ▶ Scipy
- ▶ Sklearn

Источник: Genetic Algo: Feature Selection Hyperparameters

```python
1  import numpy as np
2  import pandas as pd
3  import numpy.random as rnd
4  from scipy import spatial
5
6  from sklearn.model_selection import train_test_split, cross_val_score, KFold
7  from sklearn.preprocessing import MinMaxScaler, StandardScaler
8  from sklearn.pipeline import Pipeline
9  from sklearn.compose import ColumnTransformer
10 from sklearn.feature_selection import VarianceThreshold
11 from sklearn import metrics
12
13 from sklearn.neighbors import KNeighborsClassifier
14 from sklearn.linear_model import ElasticNet
15 from sklearn.ensemble import RandomForestClassifier
16 from xgboost import XGBClassifier
17
18 from sklearn.utils.testing import ignore_warnings
19 from sklearn.exceptions import ConvergenceWarning
20
21 import matplotlib.pyplot as plt
22 from matplotlib.ticker import MaxNLocator
23 import seaborn as sns
```

```
1  df_train = pd.read_csv("../input/santander-customer-satisfaction/train.csv", index_col="ID")
2  df_test = pd.read_csv("../input/santander-customer-satisfaction/test.csv", index_col="ID")
3
4  print("{} rows and {} columns".format(*df_train.shape))
5
6  df_train.TARGET.value_counts(normalize=True)
```

```
76020 rows and 370 columns

0    0.960431
1    0.039569
Name: TARGET, dtype: float64
```

# Подготовка среды IV

```
1  display(df_train.dtypes.value_counts())
2
3  print("Number of Missing entries: " + str(df_train.isnull().sum().sum()))
4
5  df_train.var3.describe()
```

```
int64     259
float64   111
dtype: int64
Number of Missing entries: 0

count    76020.000000
mean     -1523.199277
std      39033.462364
min     -999999.000000
25%          2.000000
50%          2.000000
75%          2.000000
max        238.000000
Name: var3, dtype: float64
```

# Подготовка среды V

```
1  display(df_train['var3'].value_counts(normalize=True)[0:3])
2
3  df_train['var3'] = df_train['var3'].replace(-999999, 2)
```

```
    2        0.975599
    8        0.001815
-999999    0.001526
Name: var3, dtype: float64
```

В процессе оптимизации к новым обучающим данным будет применена 5-кратная перекрестная валидация, которая поможет получить стабильную оценку того, какие характеристики и гиперпараметры следует включить в окончательную модель.

```
1  df_train , df_validation = train_test_split(df_train,
2                                              test_size=0.2,
3                                              random_state=1989,
4                                              stratify=df_train.TARGET,
5                                              shuffle=True)
6
7  kfold = KFold(n_splits=5, random_state=1989, shuffle=True)
```

```
 1  # Identify columns
 2  fts_num = df_train.drop(axis=1,columns=['TARGET']).select_dtypes(np.number).columns
 3
 4  # Numerical Transformer StandardScaler
 5  trans_num = Pipeline(steps = [('Standarise', StandardScaler()), ('MinMax', MinMaxScaler())])
 6
 7  # Create a single Preprocessing step for predictors
 8  preprocessor_preds = ColumnTransformer(
 9      transformers=[
10          ('num', trans_num, fts_num) # Centre and scale and constrain range
11      ])
12
13  # Apply the transformations to train
14  df_train2 = pd.DataFrame(preprocessor_preds.fit_transform(df_train))
15  df_train2.columns = fts_num
16
17  # Apply the transformations to validation
18  df_validation2 = pd.DataFrame(preprocessor_preds.fit_transform(df_validation))
19  df_validation2.columns = fts_num
20
21  # Apply the transformations to test
22  df_test2 = pd.DataFrame(preprocessor_preds.fit_transform(df_test))
23  df_test2.columns = fts_num
24
25  # Create preprocessed training data
26  df_train = pd.concat([df_train2,
27                        df_train.drop(axis=1,columns=fts_num).reset_index().drop(axis=1,columns=['ID'])],
```

```
28                          axis=1)
29
30 # Create preprocessed validation data
31 df_validation = pd.concat([df_validation2,
32                              df_validation.drop(axis=1,columns=fts_num).reset_index().drop(axis=1,columns=['
      ID'])],
33                              axis=1)
34
35 # Create preprocessed test data
36 df_test = pd.concat([df_test2,
37                        df_test.drop(axis=1,columns=fts_num).reset_index().drop(axis=1,columns=['ID'])],
38                        axis=1)
39
40 # Clear objects
41 del df_train2, df_validation2, df_test2, fts_num, trans_num, preprocessor_preds
```

# Содержание

# Разработка оптимизационной модели I

Функция ниже генерирует начальные решения как для характеристик, так и для гиперпараметров.

Выбор начальных характеристик:

► Для каждого решения в популяции случайным образом выбирается процент в диапазоне 10-91%, который будет применен к общему количеству признаков.

► Для каждого решения случайным образом выбираются признаки до тех пор, пока не будет получен процент столбцов от общего числа столбцов.

Выбор начальных гиперпараметров:

► Для каждого решения генерируйте случайное значение для каждого гиперпараметра между заданными пользователем min и max

```
1   # Randomly generate candidates
2   def f_random_candidates(features_name, population, hyperparams, output_type, df_pop=False):
3       '''create an initial population'''
4
5       # Create solution for features
6       if output_type == 'feature':
7
8           # Initial population will have between 10-91% of features
9           feature_size = rnd.choice(a=range(10,91),size=population, replace=True)
10          feature_size = [np.round(pct / 100 * len(features_name)) for pct in feature_size]
11
12          # Create a list of feature positions for each candidate
13          selection = [rnd.choice(a=range(0,len(features_name)-1), replace=False, size=cols.astype('int')) \
14                       for cols in feature_size]
15
16          selection = [list(selection[i]) for i in range(len(selection))]
17
18          return selection
19
20      # Create solution for hyperparameters
21      elif (output_type == 'hyperparams') & (hyperparams != False):
22
23          # Generate random numbers in range for each hyperparameter
24          random_hyperparams = []
25          for j in range(len(hyperparams['names'])):
26              temp = (np.random.uniform(hyperparams['min_value'][j],
```

```
27                                              hyperparams['max_value'][j],
28                                              population))
29          random_hyperparams.append(temp)
30
31      # Get length of features
32      n_features = df_pop['features'].apply(len).tolist()
33
34      # Store hyperparameters in diction
35      hyperparam_vals = []
36      for i in range(population):
37          val = {'name':[],'value':[]}
38          for j in range(len(hyperparams['names'])):
39              val['name'].append(hyperparams['names'][j])
40              temp = random_hyperparams[j][i]
41              if hyperparams['type'][j] == 'int':
42                  temp = np.int64(round(temp))
43              if hyperparams['names'][j] == 'max_features':
44                  temp = min(temp, n_features[i])
45              val['value'].append(temp)
46
47          hyperparam_vals.append(val)
48          del val
49
50      return hyperparam_vals
```

Скрещивание признаков:

- ► Случайным образом генерируется целое число, которое представляет собой точку пересечения между первой и последней характеристикой (в конечном счете, индекс столбца).
- ► Взвешенная выборка решений предыдущего поколения и выбор двух родителей
- ► Создание дочернего решения, которое имеет все признаки до точки пересечения (индекс столбца) от первого родителя и все признаки от второго родителя после точки пересечения.

Скрещивание гиперпараметров:

- ► Взвешенная выборка решений предыдущего поколения размером в число гиперпараметров
- ► Случайный выбор гиперпараметров из родительского решения

▶ Создать дочерние гиперпараметры из выбранного родительского решения

```
1   # Crossover function
2   def f_gen_child_crossover(df, features_name, hyperparams, output_type):
3       '''Mutate 2 parents to create a child'''
4
5       # Crossover features
6       if output_type == 'feature':
7
8           # Create an integer list of features
9           l_features = list(range(0,len(features_name)))
10
11          # Identify a random cross over point
12          cross_point = np.int(rnd.randint(low=0, high=len(features_name), size=1))
13
14          # Extract Two Parents
15          selection = np.random.choice(df.features,
16                                       size=2,
17                                       replace=False,
18                                       p=df.probability)
19          par1 = list(selection[0])
20          par2 = list(selection[1])
21
22          # Convert to Boolean
```

```
23          par1 = [item in par1 for item in l_features]
24          par2 = [item in par2 for item in l_features]
25
26          # Single point cross over and convert to indices
27          child = par1[0:cross_point] + par2[cross_point:]
28          child = [i for i,x in enumerate(child) if x == True]
29
30          # Return
31          return child
32
33      # Crossover hyperparameters
34      elif (output_type == 'hyperparams') & (hyperparams != False):
35
36          # Identify the number of parameters
37          n_hyperparameters = len(hyperparams['min_value'])
38
39          # Extract n Parents
40          selection = np.random.choice(df.hyperparameters,
41                                        size=n_hyperparameters,
42                                        replace=False,
43                                        p=df.probability)
44
45          # Randomly choose which parent to select each parameter from
46          parent_choice = list(np.random.choice(range(n_hyperparameters),
47                                        size = n_hyperparameters,
48                                        replace=False))
49
```

```
50        # Copy the parent as the child
51        child = selection[0]
52
53        # Update child vector with choosen parent
54        for i in range(n_hyperparameters):
55            child['value'][i] = selection[parent_choice[i]]['value'][i]
56
57        # Return
58        return child
```

Мутация признаков:

- ► Для каждого признака в кадре данных сгенерируйте случайное число от 0 до 1.
- ► Если сгенерированная вероятность ниже указанной пользователем скорости мутации, то поменяйте местами переключатели для этого столбца (т.е. если функция включена, то удалите ее и наоборот).

Мутация гиперпараметров:

- ► Для каждого гиперпараметра в выбранной модели сгенерируйте случайное число от 0 до 1.
- ► Если случайное число ниже указанной пользователем скорости мутации, сгенерируйте случайное число в указанном диапазоне.

▶ Наконец, проверьте, не находится ли гиперпараметр за пределами диапазона min-max, и при необходимости уменьшите его до этого диапазона.

```
1  # Mutate function
2  def f_gen_child_mutate(candidate, features_name, p_mutate,
3                         hyperparams, output_type,
4                         hyperparams_increment):
5      '''Mutate 2 parents to create a child'''
6
7      # Mutate Features
8      if output_type == 'feature':
9
10         # Create an integer list of features
11         l_features = list(range(0,len(features_name)))
12
13         # Convert feature into boolean vector
14         candidate = [item in candidate for item in l_features]
15
16         # Conditionally mutate features in chromosome (reverse binary flag)
17         candidate_new = []
18         for item in candidate:
19             if rnd.rand() <= p_mutate:
20                 candidate_new.append(not item)
```

```
21              else:
22                  candidate_new.append(item)
23
24          # Convert to indicies
25          candidate_new = [i for i,x in enumerate(candidate_new) if x == True]
26
27          # Return
28          return candidate_new
29
30      # Mutate hyperparameters
31      elif (output_type == 'hyperparams') & (hyperparams != False):
32
33          # Identify size of mutation
34          v_mutate = (np.random.uniform((1-hyperparams_increment),
35                                          (1+hyperparams_increment), 1)).item()
36
37          # Identify Min and Max for parameters
38          l_min = hyperparams['min_value']
39          l_max = hyperparams['max_value']
40
41          # Identify the number of parameters
42          n_hyperparameters = len(l_min)
43
44          # Probabilistically mutate certain parameters
45          candidate_new = []
46          for i in range(n_hyperparameters):
47              if rnd.rand() <= p_mutate:
```

```
48                    temp = candidate['value'][i] * v_mutate
49                    if hyperparams['type'][i] == 'int':
50                        temp = np.int64(round(temp))
51                    candidate_new.append(temp)
52                else:
53                    candidate_new.append(candidate['value'][i])
54
55            # Ensure that value is between ranges
56            for i in range(n_hyperparameters):
57                if (candidate_new[i] < l_min[i]):
58                    candidate_new[i] = l_min[i]
59                elif (candidate_new[i] > l_max[i]):
60                    candidate_new[i] = l_max[i]
61
62            # Update values
63            candidate['value'] = candidate_new
64
65            # return
66            return candidate
```

```python
1   # Function to generate a population of candidates
2   def f_generate_population(inital_flag, population, features_name,
3                             p_crossover, p_mutate,
4                             hyperparams, hyperparams_increment,
5                             hyperparams_multiple,
6                             df=False, generation=0, initalise=False):
7       '''Generates all candidates in population'''
8
9       # Create initial population
10      if inital_flag == True:
11
12          # Check if there is an initial solution & reduce
13          # population by one if there is
14          if initalise != False:
15              population = population - 1
16
17          # generate random features
18          df_pop = pd.DataFrame({'generation':generation,
19                                 'candidate':range(0,population),
20                                 'features':f_random_candidates(features_name,
21                                                                population,
22                                                                hyperparams,
23                                                                output_type = 'feature')})
24
25          # Duplicate rows for population range
26          df_pop = df_pop.loc[df_pop.index.repeat(hyperparams_multiple)]
```

```python
27
28        # Generate population
29        df_pop['hyperparameters'] = \
30            f_random_candidates(features_name=features_name,
31                                population = population * hyperparams_multiple,
32                                hyperparams=hyperparams,
33                                output_type = 'hyperparams',
34                                df_pop=df_pop)
35
36        # If Initial solution then add in
37        if initalise != False:
38            df = pd.DataFrame({'generation':generation,
39                               'candidate':range(population, population + 1),
40                               'features':[initalise['features']],
41                               'hyperparameters':[initalise['hyperparameters']]},
42                              index=[population])
43
44            df_pop = df_pop.append(df)
45
46
47        # Reset Index
48        df_pop.index = range(0, population * hyperparams_multiple)
49
50        # Return
51        return df_pop
52    else:
53        # Distribute the population
```

```python
54          population_crossover = round(population * p_crossover)
55          population_remainder = population-population_crossover
56
57          # ----- Create crossover candidates -----
58
59          # Create crossover populate for feature selection
60          df_pop = pd.DataFrame({'generation':generation,
61                                 'candidate':range(0,population_crossover)})
62          df_pop['features'] = [f_gen_child_crossover(df=df,
63                                                      features_name=features_name,
64                                                      hyperparams=hyperparams,
65                                                      output_type = 'feature') \
66                                for _ in range(population_crossover)]
67
68          # Duplicate rows for population range
69          df_pop = df_pop.loc[df_pop.index.repeat(hyperparams_multiple)]
70
71          # Create crossover population for hyperparameters
72          df_pop['hyperparameters'] = \
73              [f_gen_child_crossover(df=df,
74                                     features_name=features_name,
75                                     hyperparams=hyperparams,
76                                     output_type = 'hyperparams') \
77              for _ in range(population_crossover * hyperparams_multiple)]
78
79          # Reset Index
80          df_pop.index = range(0, population_crossover * hyperparams_multiple)
```

```python
81
82      # ————— Create Randomly Selected candidates —————
83
84      # Initialise population
85      df_temp = pd.DataFrame({'generation':generation,
86                              'candidate':range(population_crossover,
87                                                population)})
88      # Randomly select candidates
89      selected_index = \
90          df.sample(n=population_remainder,
91                    replace = False,
92                    weights=df.probability).candidate.tolist()
93
94      # Extract hyperparameters
95      selected_features = df.iloc[selected_index,:].features.tolist()
96      selected_params = df.iloc[selected_index,:].hyperparameters.tolist()
97
98      # Update temp dataframe
99      df_temp['features'] = [selected_features[i]
100                            for i in range(len(selected_features))]
101     df_temp['hyperparameters'] = [selected_params[i]
102                            for i in range(len(selected_params))]
103
104     # Duplicate rows for population range
105     df_temp = df_temp.loc[df_temp.index.repeat(population_remainder)]
106
107     # Append to population dataframe
```

```python
108            df_pop = df_pop.append(df_temp,ignore_index=True)
109
110            # Clear up
111            del selected_features, selected_params, df_temp
112
113            # ----- Mutate Population -----
114
115            # Mutate existing candidate features
116            df_pop['features'] = \
117                df_pop.features.apply(f_gen_child_mutate,
118                                      features_name=features_name,
119                                      p_mutate=p_mutate,
120                                      hyperparams=hyperparams,
121                                      output_type = 'feature',
122                                      hyperparams_increment=hyperparams_increment)
123
124            # Mutate existing candidate hyperparameters
125            df_pop['hyperparameters'] = \
126                df_pop.hyperparameters.apply(f_gen_child_mutate,
127                                             features_name=features_name,
128                                             p_mutate=p_mutate,
129                                             hyperparams=hyperparams,
130                                             output_type = 'hyperparams',
131                                             hyperparams_increment=
132                                                 hyperparams_increment)
133
134            # ----- Hyperparameter fix -----
```

```
135         if hyperparams != False:
136
137             # Get length of features
138             n_features = df_pop['features'].apply(len).tolist()
139
140             # Hyperparameter fix
141             for i in range(population):
142                 for j in range(len(hyperparams['names'])):
143                     if hyperparams['names'][j] == 'max_features':
144                         if df_pop.hyperparameters[i]['value'][j] > n_features[i]:
145                             df_pop.hyperparameters[i]['value'][j] = \
146                                 n_features[i]
147
148         # Return
149         return df_pop
```

# Содержание

```python
1   # Evaluate solution fitness
2   @ignore_warnings(category=ConvergenceWarning)
3   def f_fitness(model, eval_metric, features, target,
4                 feature_idx, kfold, hyperparams):
5       '''Evaluates fitness of proposed solution'''
6
7       # Extract the hyperparameters
8       n_hyperparams = len(hyperparams['name'])
9       hyperparameters = {hyperparams['name'][0]: hyperparams['value'][0]}
10      if n_hyperparams > 1:
11          for i in range(n_hyperparams):
12              tempparameters = {hyperparams['name'][i]: hyperparams['value'][i]}
13              hyperparameters = {**hyperparameters, **tempparameters}
14
15      # Determine CV strategy
16      if kfold == False:
17          kfold = 5
18      else:
19          kfold = kfold
20
21      # Apply cross validation to the modells
22      results = cross_val_score(model.set_params(**hyperparameters),
23                                features.iloc[:, feature_idx],
24                                target,
25                                cv=kfold,
26                                scoring=eval_metric)
27
```

```
28      # Replace NA's with 0
29      results[np.isnan(results)] = 0
30
31      return results
```

```
1   # Apply evaluation score to current population
2   def f_evaluation_score(df, features, target, eval_metric, model,
3                          kfold, hyperparams):
4       '''Apply f_fitness to each candidate'''
5
6       # Calculate the evaluation metric
7       evaluation_score = []
8       for val in range(0, len(df)):
9           eval_score = f_fitness(model=model,
10                                  eval_metric=eval_metric,
11                                  features = features,
12                                  target=target,
13                                  feature_idx=df['features'][val],
14                                  kfold=kfold,
15                                  hyperparams=df['hyperparameters'][val])
16
17          # Average evaluation metric across folds
18          evaluation_score.append(eval_score.mean())
```

```
19
20          # Clear object
21          del eval_score
22
23      # Clear object
24      del val
25
26      # return evaluation score
27      return evaluation_score
```

```
1  # Calculate jaccard similarity
2  def f_j_sim(list1, list2):
3      s1 = set(list1)
4      s2 = set(list2)
5      return float(len(s1.intersection(s2)) / len(s1.union(s2)))
6
7  # Calculate cosine similarity
8  def f_c_sim(l_other, l_best_score):
9
10      # Extract hyperparameter values
11      l_other = l_other['value']
12
13      # calculate similarity
```

```python
14      sim = 1 - spatial.distance.cosine(l_best_score, l_other)
15
16      # return
17      return sim
18
19  # Calculate similarity between candidates and probability for next gen selection
20  def f_sim_n_prob(df):
21
22      # Calculate similarity of solutions with best solutions - Features
23      l_best_score = df.features[df['fitness_score'].idxmax()]
24      df['similarity_features'] = df['features'].apply(f_j_sim, list2=l_best_score)
25      del l_best_score
26
27      # Calculate similarity of solutions with best solutions
28      l_best_score = df.hyperparameters[df['fitness_score'].idxmax()]['value']
29      df['similarity_hyperparameters'] = df.hyperparameters.apply(f_c_sim, l_best_score=l_best_score)
30      del l_best_score
31
32      # Calculate cumulative probability for future stages
33      df['probability'] = (df['fitness_score'] / sum(df['fitness_score']))
34
35      # return
36      return df
```

```python
# Function to populate attributes of candidates
def f_population_features(df, features, target, desiriability,
                          eval_metric, model, kfold, hyperparams):
    '''Get features of all candidates in population'''

    # Calculate feature size for candidates
    df['feature_size'] = df['features'].apply(len)

    # Calculate evaluation score for candidates
    df['evaluation_score'] = f_evaluation_score(df,
                                                features,
                                                target,
                                                eval_metric,
                                                model,
                                                kfold,
                                                hyperparams)

    # Conditionally create desirability fitness score
    if desiriability != False:

        # Create scalars – Features
        v_lb_features = desiriability['lb'][1]
        v_ub_features = desiriability['ub'][1]
        v_s_features = desiriability['s'][1]

        # Create scalars – Evaluation Metric
```

```
27            v_lb_eval = desiriability['lb'][0]
28            v_ub_eval = desiriability['ub'][0]
29            v_s_eval = desiriability['s'][0]

31            # Calculate desirability for features
32            df['desire_features']  = [0 if x > v_ub_features else 1
33                                if x < v_lb_features else
34                                    ((x-v_ub_features)/
35                                     (v_lb_features-v_ub_features))**
36                                    v_s_features
37                                for x in df['feature_size']]

39            # Calculate desirability for evaluation metric
40            df['desire_eval']  = [0 if x < v_lb_eval else 1
41                             if x > v_ub_eval else
42                                 ((x-v_lb_eval)/
43                                  (v_ub_eval-v_lb_eval))**
44                                  v_s_eval
45                             for x in df['evaluation_score']]

47            # calculate fitness score
48            df['fitness_score'] = (df['desire_features'] * df['desire_eval'])**0.5

50            # Drop fields
51            df = df.drop(columns=['desire_features', 'desire_eval'])

53        else:
```

```
54          # calculate fitness score
55          df['fitness_score'] = df['evaluation_score']
56
57      # Return
58      return df
```

```
1   # Main Optimisation Function
2   def f_model_optimisation(df,
3                               target_var,
4                               generations,
5                               population,
6                               eval_metric,
7                               model,
8                               kfold=False,
9                               hyperparams_multiple = 3,
10                              hyperparams = False,
11                              desiriability=False,
12                              p_crossover=0.8,
13                              p_mutate=0.01,
14                              hyperparams_increment = 0.1,
15                              elitism=False,
16                              gens_no_improve = False,
17                              initalise = False):
```

```
18      '''Function uses GA's to choose features and tune hyperparameters'''
19
20      # Print Model Stats
21      print('Model Initialisation')
22
23      # --------- Split features and target ---------
24      features = df.drop(target_var, axis=1)
25      features_name = features.columns
26      target = df[target_var]
27
28      # --------- First Generation ---------
29
30      # Generate inital candidate features solutions
31      df_pop_cur = f_generate_population(inital_flag=True,
32                                         population=population,
33                                         features_name=features_name,
34                                         p_crossover=p_crossover,
35                                         p_mutate=p_mutate,
36                                         hyperparams=hyperparams,
37                                         hyperparams_increment=hyperparams_increment,
38                                         hyperparams_multiple=hyperparams_multiple,
39                                         initalise=initalise)
40
41      # Enrich candidate solutions with features
42      df_pop_cur = f_population_features(df=df_pop_cur,
43                                         features=features,
44                                         target=target,
```

```
45                                          desiriability=desiriability,
46                                          eval_metric=eval_metric,
47                                          model=model,
48                                          kfold=kfold,
49                                          hyperparams=hyperparams
50                                          )
51
52      # Extract best score for each candidate
53      df_pop_cur = df_pop_cur.loc[df_pop_cur.reset_index().\
54                                  groupby(['candidate'])['fitness_score'].\
55                                  idxmax()]
56
57      # Enrich candidate solutions with similarity & probability
58      df_pop_cur = f_sim_n_prob(df_pop_cur)
59
60      # --------- Create search storage ---------
61      df_output = df_pop_cur.copy()
62
63      # Print Model Stats
64      print('Gen: 00' +
65            ' - Generation Mean:' + str(round(df_output.fitness_score.mean(), 4)).zfill(4) +
66            ' - Generation Best:' + str(round(df_output.fitness_score.max(), 4)).zfill(4) +
67            ' - Global Best:' + str(round(df_output.fitness_score.max(), 4)).zfill(4)
68            )
69
70      # Track best solution
71      if gens_no_improve != False:
```

```python
72          count = 0
73          v_best = df_output.fitness_score.max()
74
75      # --------- Run additional generations ---------
76
77      # Loop for additional generations
78      for gen in range(1, generations):
79
80          # --------- Elitism ---------
81          if elitism > 0:
82
83              # Create a dataframe with elite candidates
84              df_elite = df_output.nlargest(columns='fitness_score', n=elitism)
85              df_elite['candidate'] = population - 1
86              df_elite['generation'] = gen
87              df_elite = df_elite.drop(columns=['similarity_features',
88                                                'similarity_hyperparameters', 'probability'])
89          # --------- New Population ---------
90
91          # Generate next candidate solutions
92          df_pop_cur = f_generate_population(inital_flag=False,
93                                             generation = gen,
94                                             population=(population-elitism),
95                                             features_name=features_name,
96                                             df=df_pop_cur,
97                                             p_crossover=p_crossover,
98                                             p_mutate=p_mutate,
```

```
 99                                                    hyperparams=hyperparams ,
100                                                    hyperparams_increment=hyperparams_increment ,
101                                                    hyperparams_multiple=hyperparams_multiple
102                                                    )
103
104            # Enrich candidate solutions with features
105            df_pop_cur = f_population_features (df=df_pop_cur ,
106                                                features=features ,
107                                                target=target ,
108                                                desiriability=desiriability ,
109                                                eval_metric=eval_metric ,
110                                                model=model ,
111                                                kfold=kfold ,
112                                                hyperparams=hyperparams )
113
114            # Add elite
115            if elitism > 0:
116                df_pop_cur = pd.concat ([ df_pop_cur , df_elite ]) .reset_index () .drop (columns=['index '])
117                del df_elite
118
119            # Extract best score for each candidate
120            df_pop_cur = df_pop_cur .loc [ df_pop_cur .reset_index () .\
121                                            groupby ([ 'candidate ']) [ 'fitness_score '].\
122                                            idxmax () ]
123
124            # Enrich candidate solutions with similarity & probability
125            df_pop_cur = f_sim_n_prob (df=df_pop_cur )
```

```
126
127            # Update Output
128            df_output = df_output.append(df_pop_cur,ignore_index=True)
129
130            # Print Model Stats
131            print('Gen: ' + str(gen).zfill(2) +
132               ' - Generation Mean:' + str(round(df_output[df_output.generation == gen].fitness_score.mean()
          , 4)).zfill(4) +
133               ' - Generation Best:' + str(round(df_output[df_output.generation == gen].fitness_score.max(),
           4)).zfill(4) +
134               ' - Global Best:' + str(round(df_output.fitness_score.max(), 4)).zfill(4)
135                  )
136
137            # Track number of generations with no improvement
138            if gens_no_improve != False:
139                if df_output.fitness_score.max() > v_best:
140                    count = 0
141                    v_best = df_output.fitness_score.max()
142                else:
143                    count += 1
144
145                # Conditionally break loop
146                if count == gens_no_improve:
147                    break
148
149     # Return df
150     return df_output
```

# Содержание

Применим следующие модели к нашему процессу оптимизации:

► ElasticNet

► XGBoost

```
 1  # Run Optimisation – Optimise for AUC
 2  df_ENet_AUC = f_model_optimisation(df=df_train,
 3                                     target_var='TARGET',
 4                                     generations=7,
 5                                     population=20,
 6                                     p_crossover=0.8,
 7                                     p_mutate=0.02,
 8                                     hyperparams_increment=0.01,
 9                                     hyperparams_multiple = 5,
10                                     eval_metric='roc_auc',
11                                     kfold=False,
12                                     model=ElasticNet(),
13                                     hyperparams = {'names':['alpha', 'l1_ratio'],
14                                                    'min_value': [0, 0],
15                                                    'max_value': [0.01, 1],
16                                                    'type':['float', 'float']}
17                                     )
```

# ElasticNet II

```
Model Initialisation
Gen: 00 - Generation Mean:0.7282 - Generation Best:0.7849 - Global Best:0.7849
Gen: 01 - Generation Mean:0.7397 - Generation Best:0.7897 - Global Best:0.7897
Gen: 02 - Generation Mean:0.7486 - Generation Best:0.7875 - Global Best:0.7897
Gen: 03 - Generation Mean:0.7369 - Generation Best:0.7801 - Global Best:0.7897
Gen: 04 - Generation Mean:0.7428 - Generation Best:0.7799 - Global Best:0.7897
Gen: 05 - Generation Mean:0.7521 - Generation Best:0.7802 - Global Best:0.7897
Gen: 06 - Generation Mean:0.7563 - Generation Best:0.7787 - Global Best:0.7897
```

```
1  # Run Optimisation — Optimise for AUC
2  df_xgb_AUC = f_model_optimisation(df=df_train,
3                                    target_var='TARGET',
4                                    generations=5,
5                                    population=20,
6                                    p_crossover=0.8,
7                                    p_mutate=0.02,
8                                    hyperparams_increment=0.01,
9                                    hyperparams_multiple = 5,
10                                   eval_metric='roc_auc',
11                                   kfold=False,
12                                   model=XGBClassifier(objective="binary:logistic", scale_pos_weight = 25),
13                                   hyperparams = {'names':['learning_rate', 'max_depth',
14                                                           'min_child_weight', 'gamma', 'colsample_bytree'],
15                                                  'min_value': [0.03, 2,   1, 0,   0.3],
16                                                  'max_value': [0.3,  15,  7, 0.5, 0.7],
17                                                  'type':['float', 'int', 'int', 'float', 'float']}
18                                   )
```
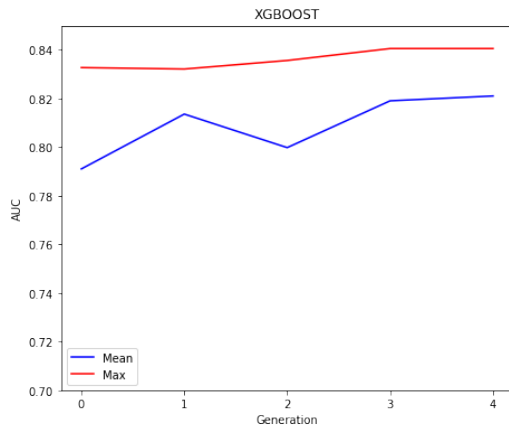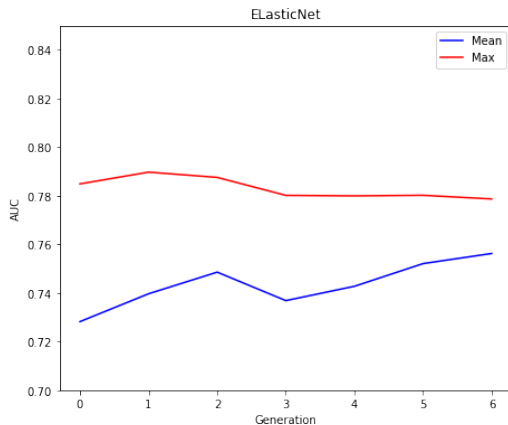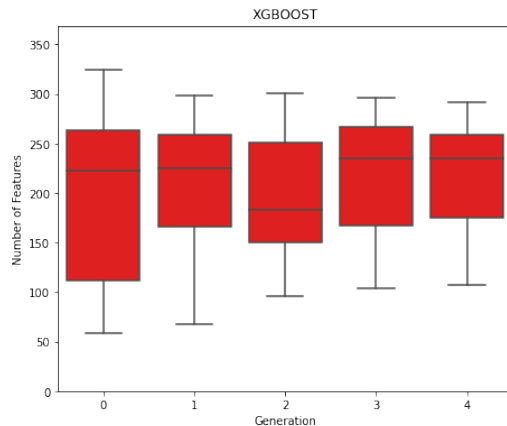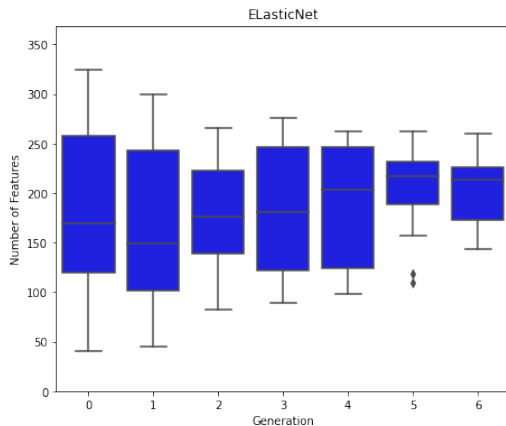
```
Model Initialisation
Gen: 00 - Generation Mean:0.791 - Generation Best:0.8327 - Global Best:0.8327
Gen: 01 - Generation Mean:0.8136 - Generation Best:0.8321 - Global Best:0.8327
Gen: 02 - Generation Mean:0.7998 - Generation Best:0.8356 - Global Best:0.8356
Gen: 03 - Generation Mean:0.819 - Generation Best:0.8405 - Global Best:0.8405
Gen: 04 - Generation Mean:0.821 - Generation Best:0.8405 - Global Best:0.8405
```
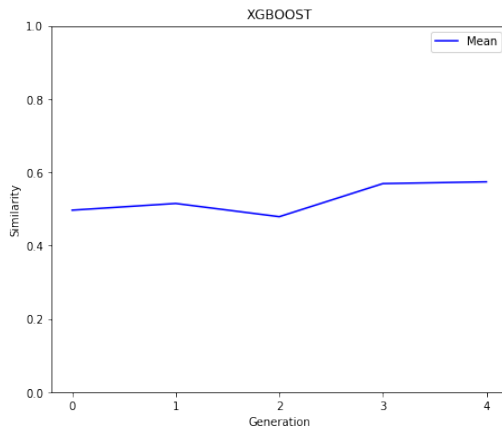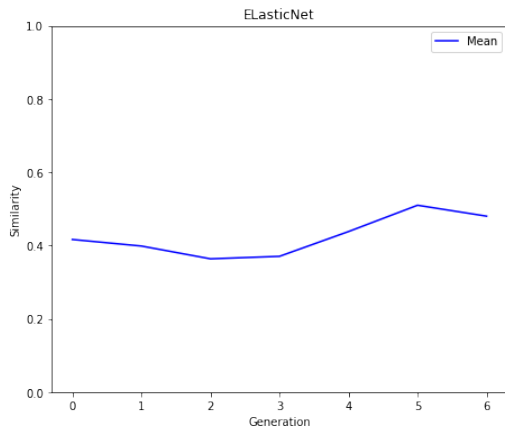
# Содержание

Average and Best Evaluation Scores Per Generation

Distribution of Feature Size Per Generation

Average Similarity Between Candidate Features and Best Solution's Per Generation

Average Similarity Between Candidate Hyperparameters and Best Candidate Per Generation