

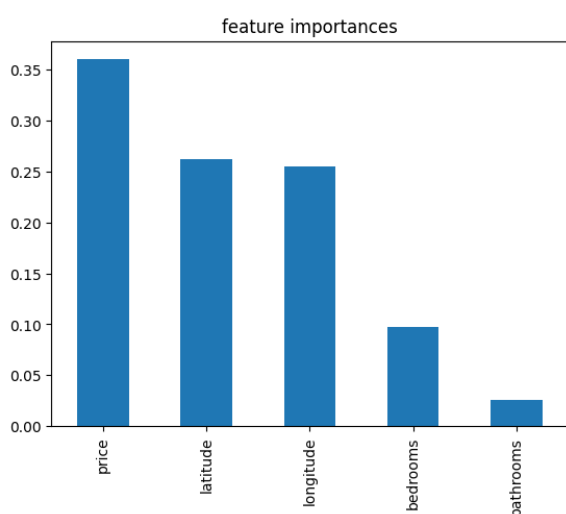
Отчет.

В ходе работы я провел эксперименты с feature importance и добавил все в библиотеку на python. Сравнил feature importance в реализации sklearn и R, а также сравнил sklearn с R-ranger.

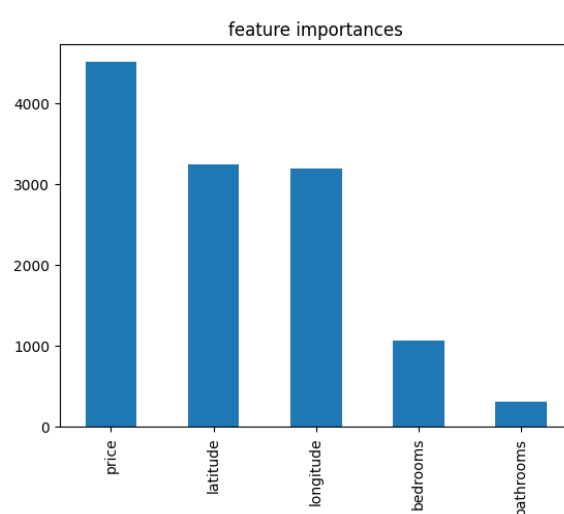
Эксперименты.

Использовал я датасет *rent.csv* из статьи, рассмотрел случаи, когда признаки: 1) разнородны, 2) присутствуют коррелированные признаки, 3) присутствуют нерелевантные признаки.

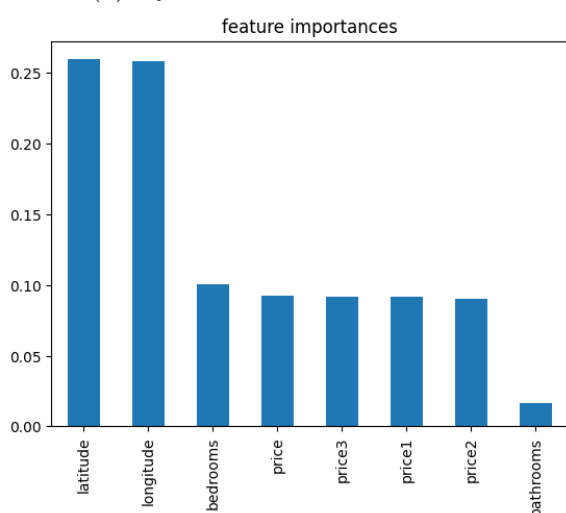
Подбор оптимальных гиперпараметров в реализации sklearn я осуществлял с помощью *Optuna*, и полученные параметры я использовал в реализации на R.



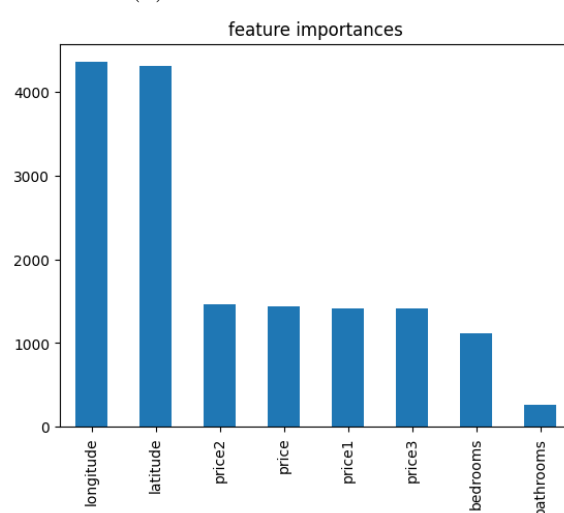
(a) Python: Исходный датасет



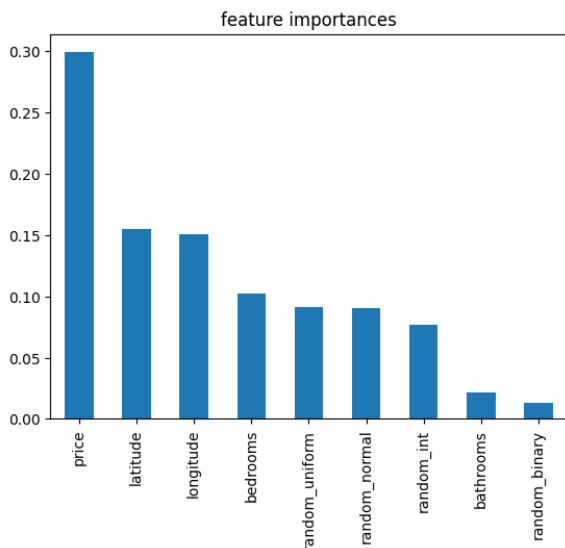
(a) R: Исходный датасет



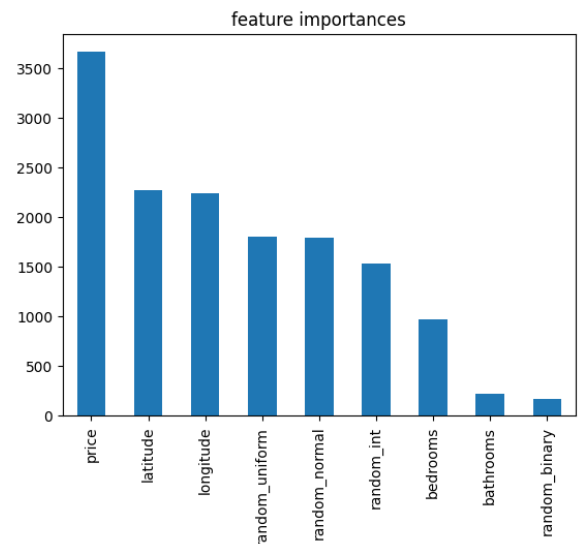
(b) Python: добавлены коррелированные признаки



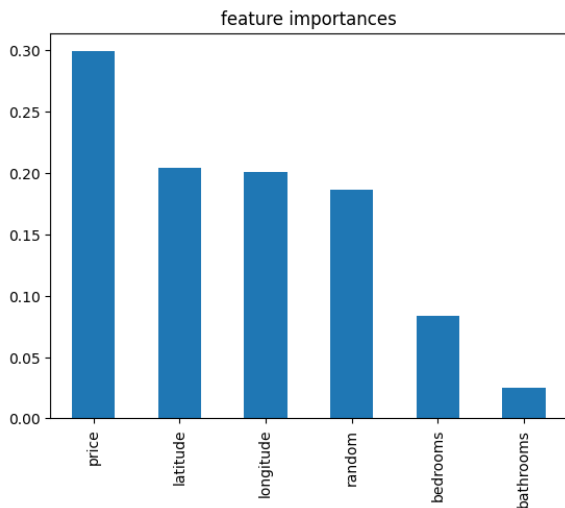
(b) R: добавлены коррелированные признаки



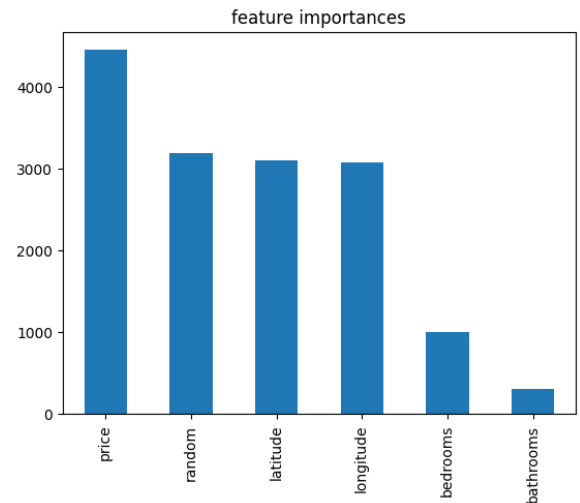
(c) Python: много нерелевантных признаков



(c) R: много нерелевантных признаков



(d) Python: 1 нерелевантный признак



(d) R: 1 нерелевантный признак

Как можно заметить, результаты в некоторых экспериментах неожиданные.

В примере (b) я скопировал признак *price* 3 раза, и получилось, что важность этого признака уменьшилась. Вообще, при добавлении коррелированных признаков, важности этих признаков будут уменьшаться.

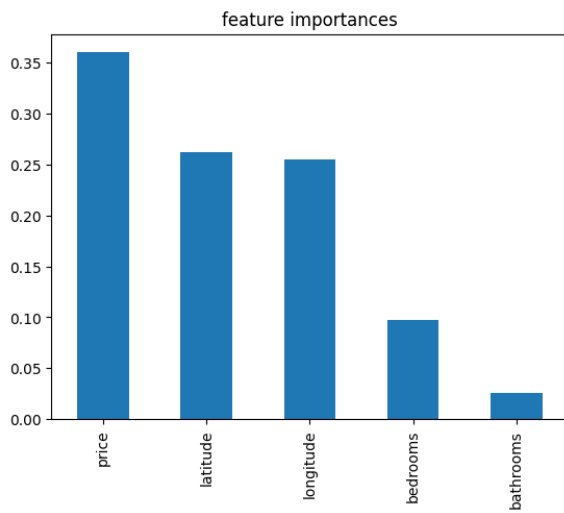
В примерах (c) и (d) я создал несколько признаков, которые никак не связаны с *target*.

Возникают вопросы, глядя на рисунки с нерелевантными признаками разного типа: 1) Почему эти признаки имеют важность? 2) Почему признаки *random_uniform*, *random_normal*, *random_int* имеют важность больше, чем признак *random_binary*?

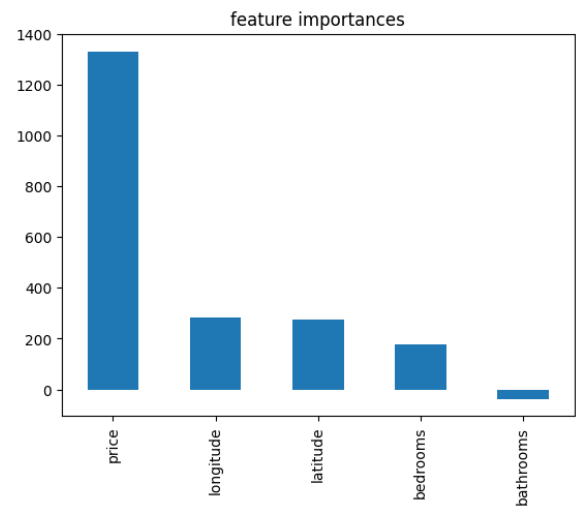
Ниже приведена таблица с оптимальными параметрами для каждого из примеров (a)-(d):

Пример	<i>n_estimators</i>	<i>max_depth</i>	<i>min_samples_leaf</i>	<i>max_features</i>
(a)	1760	26	5	sqrt
(b)	1289	14	3	1.0
(c)	1064	24	13	0.5
(d)	1096	19	3	0.3333333333333333

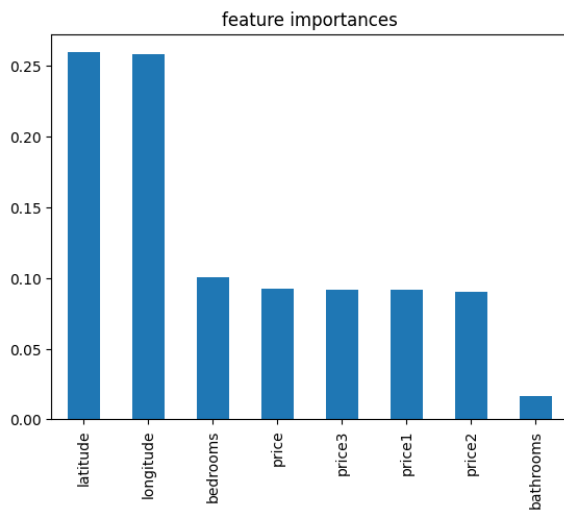
Sklearn vs R-ranger



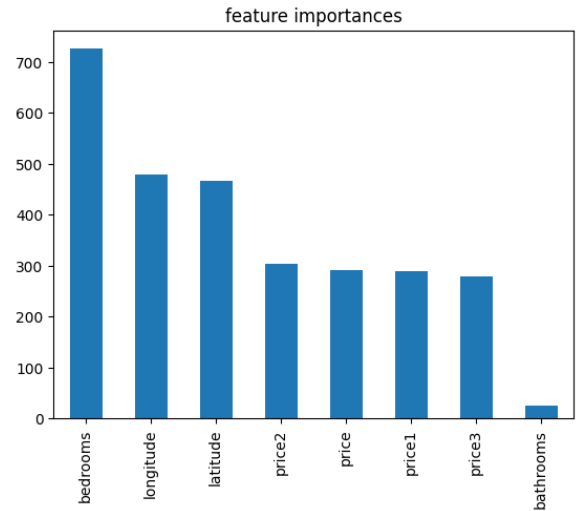
(a) Sklearn: Исходный датасет



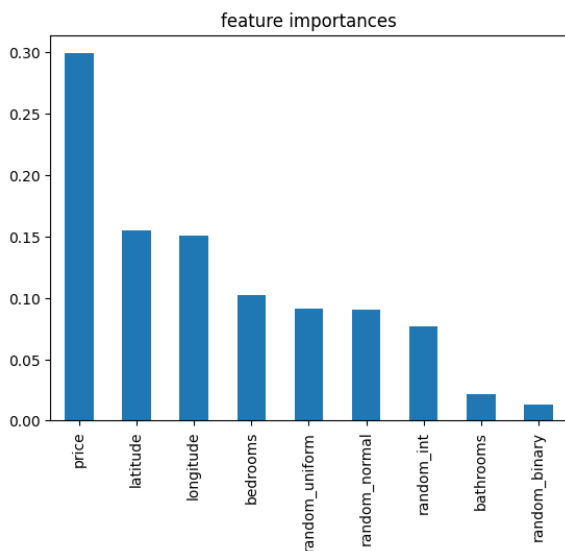
(a) R-ranger: Исходный датасет



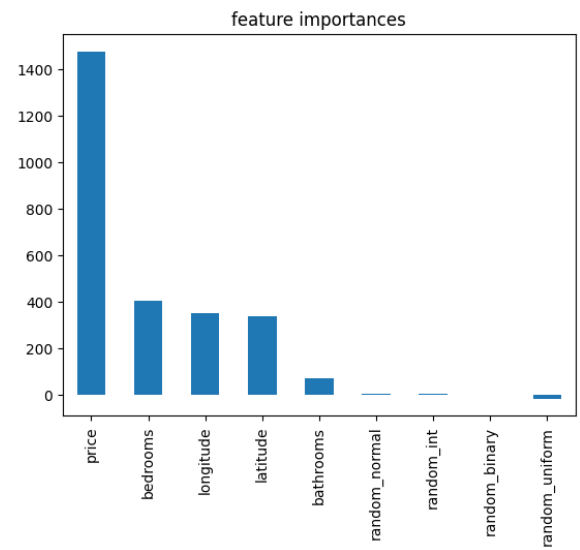
(b) Sklearn: добавлены коррелированные признаки



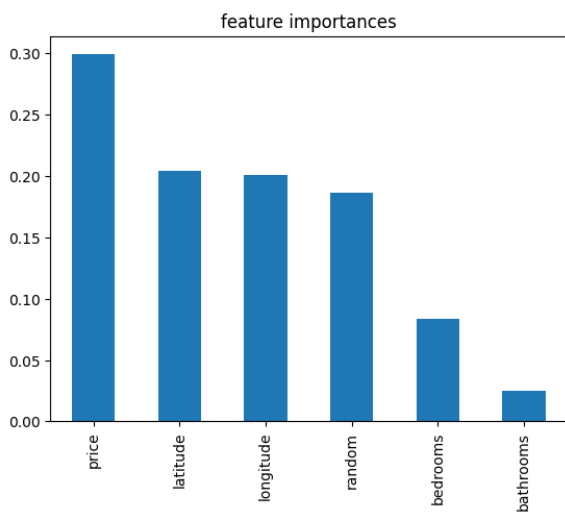
(b) R-ranger: добавлены коррелированные признаки



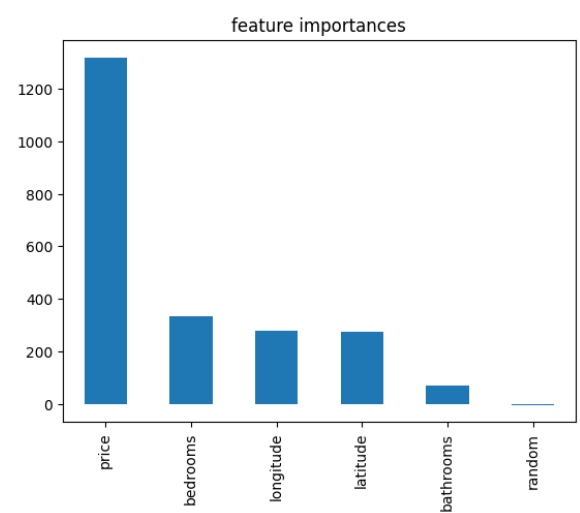
(c) Sklearn: много нерелевантных признаков



(c) R-ranger: много нерелевантных признаков



(d) Sklearn: 1 нерелевантный признак



(d) R-ranger: 1 нерелевантный признак

Код

```
1 from .python_implementations import (
2     sklearn_importance,
3     objective_classifier
4 )
5
6 from .r_implementations import (
7     r_randomforest_importance,
8     clean_feature_names,
9     r_ranger_importance_air
10 )
11
12
13 from .pic import (
14     picture
15 )
16
17 __version__ = "1.0.0"
18
19 __all__ = ['sklearn_importance', 'picture',
20            'r_randomforest_importance', 'r_ranger_importance_air']
```

Листинг 1: Файл `__init__.py`

```
1 import matplotlib.pyplot as plt
2
3
4 def picture(importance_list, titles):
5     n_plots = len(importance_list)
6     ncols = 3
7     nrows = int(np.ceil(n_plots / ncols))
8
9     fig, axes = plt.subplots(
10         nrows=nrows,
11         ncols=ncols,
12         figsize=(5 * ncols, 6 * nrows)
13     )
14
15     axes = axes.flatten()
16
17     for i, (importance_series, title) in enumerate(zip(importance_list,
18                                                         titles)):
19         ax = axes[i]
20
21         importance_series.plot(kind='bar', ax=ax)
22         ax.set_title(title, fontsize=14)
23         ax.set_ylabel('')
24         ax.tick_params(axis='x', rotation=90)
25         ax.grid(axis='y', linestyle='--')
26
27     for j in range(n_plots, nrows * ncols):
28         fig.delaxes(axes[j])
29
30     plt.tight_layout()
31     plt.show()
32     return
```

Листинг 2: Файл `pic.py`

```

1 from sklearn.ensemble import RandomForestClassifier
2 import pandas as pd
3 import optuna
4 from sklearn.model_selection import train_test_split
5 from sklearn.metrics import accuracy_score
6
7
8 def objective_classifier(trial, X_train, y_train):
9     n_estimators = trial.suggest_int('n_estimators', 50, 2000)
10    max_depth = trial.suggest_int('max_depth', 1, 30, log=True)
11    min_samples_leaf = trial.suggest_int('min_samples_leaf', 1, 40)
12    max_features = trial.suggest_categorical(
13        'max_features', ["sqrt", 0.25, 1/3, 0.5, 0.7, 1.0])
14
15    model = RandomForestClassifier(
16        n_estimators=n_estimators,
17        max_depth=max_depth,
18        min_samples_leaf=min_samples_leaf,
19        max_features=max_features,
20        random_state=42,
21        oob_score=True,
22        bootstrap=True,
23        n_jobs=-1
24    )
25
26    model.fit(X_train, y_train)
27
28    return model.oob_score_
29
30
31 def sklearn_importance(X_train, y_train):
32     study = optuna.create_study(direction='maximize')
33     study.optimize(lambda trial: objective_classifier(
34         trial, X_train, y_train), n_trials=100)
35     best_params = study.best_params
36     rf = RandomForestClassifier(
37         **best_params,
38         random_state=42,
39         oob_score=True,
40         bootstrap=True,
41         n_jobs=-1
42     )
43     print(best_params)
44     rf.fit(X_train, y_train)
45     fi = pd.Series(rf.feature_importances_, index=X_train.columns)
46     return fi.sort_values(ascending=False), rf, study

```

Листинг 3: Файл python_implementations.py

```

1 from rpy2.robj.ctors import StrVector
2 import pandas as pd
3 import numpy as np
4 import re
5 import rpy2.robj as robj
6 from rpy2.robj import pandas2ri
7 from rpy2.robj.packages import importr
8 import warnings
9 import traceback
10 from rpy2 import robj
11

```

```

12
13 def clean_feature_names(names):
14     cleaned_names = [name.replace('.', '_').replace(
15         '-', '_').replace(' ', '_') for name in names]
16     return cleaned_names
17
18
19 def r_randomforest_importance(X, y, data_path=None, n_estimators=100,
20     max_depth=None, min_samples_leaf=1, max_features='sqrt'):
21     try:
22         base = importr('base')
23         utils = importr('utils')
24         randomForest = importr('randomForest')
25
26         original_features = list(X.columns)
27         cleaned_features = clean_feature_names(original_features)
28
29         df_for_r = X.copy()
30         df_for_r.columns = cleaned_features
31         df_for_r['target'] = y.values
32
33         with objects.conversion.localconverter(objects.default_converter +
34             pandas2ri.converter) as cv:
35             r_df = objects.conversion.py2rpy(df_for_r)
36
37             formula_str = "target ~ " + " + ".join(cleaned_features)
38             formula = objects.Formula(formula_str)
39
40             print(f"Training R Random Forest with formula: {formula_str}")
41
42             n_samples = X.shape[0]
43             n_features = len(cleaned_features)
44             if max_features == 'sqrt':
45                 max_features = round(np.sqrt(n_features))
46             else:
47                 max_features = round(n_features * max_features)
48
49             if max_depth is None:
50                 maxnodes = objects.NULL
51             else:
52                 maxnodes=2**max_depth
53
54             rf_result = randomForest.randomForest(
55                 formula,
56                 data=r_df,
57                 ntree=n_estimators,
58                 nodesize=min_samples_leaf,
59                 maxnodes=maxnodes,
60                 mtry=max_features,
61                 importance=True
62             )
63
64             importance_r = randomForest.importance(rf_result, type=2, scale=
65                 False)
66
67             importance_matrix = np.array(importance_r)
68
69             feature_names_r = list(objects.r['rownames'](importance_r))
70
71             print(f"Importance matrix shape: {importance_matrix.shape}")

```

```

69     print(f"Feature names from R: {feature_names_r}")
70
71     if importance_matrix.ndim == 2:
72         importance_values = importance_matrix[:, 0]
73     else:
74         importance_values = importance_matrix
75
76     importance_series = pd.Series(
77         importance_values, index=original_features)
78
79     print("Successfully computed variable importance")
80     print(importance_series)
81
82     return importance_series.sort_values(ascending=False)
83
84 except Exception as e:
85     warnings.warn(f"R implementation failed: {e}")
86     import traceback
87     traceback.print_exc()
88     return None
89
90
91 def r_ranger_importance_air(X, y, n_estimators=100, max_depth=None,
92 min_samples_leaf=1, max_features='sqrt'):
93     try:
94         base = importr('base')
95         ranger = importr('ranger')
96
97         original_features = list(X.columns)
98
99         cleaned_features = clean_feature_names(original_features)
100
101         df_for_r = X.copy()
102         df_for_r.columns = cleaned_features
103
104         df_for_r['target'] = y.astype(str).values
105
106         name_map = dict(zip(cleaned_features, original_features))
107
108         with robjects.conversion.localconverter(robjects.default_converter +
109 pandas2ri.converter):
110             r_df = robjects.conversion.py2rpy(df_for_r)
111
112             r_target = base.factor(r_df.rx2('target'))
113
114             r_cleaned_features = StrVector(cleaned_features)
115
116             r_df_features_only = r_df.rx(True, r_cleaned_features)
117
118             r_df = base.cbind(r_df_features_only, target=r_target)
119
120             formula_str = "target ~ " + " + ".join(cleaned_features)
121             formula = robjects.Formula(formula_str)
122
123             print(f"Training R Ranger Forest with formula: {formula_str}")
124
125             target_type = robjects.r['class'](r_df.rx2('target'))
126             target_levels = robjects.r['levels'](r_df.rx2('target'))
127             print(f"Target type in R: {list(target_type)}")
128             print(f"Target levels: {list(target_levels)}")

```



```

127     n_features = len(cleaned_features)
128     if max_features == 'sqrt':
129         mtry_val = round(np.sqrt(n_features))
130     elif isinstance(max_features, float):
131         mtry_val = round(n_features * max_features)
132     else:
133         mtry_val = max_features
134
135     if max_depth is None:
136         max_depth = robjects.NULL
137
138     rf_result = ranger.ranger(
139         formula,
140         data=r_df,
141         num_trees=n_estimators,
142         min_node_size=min_samples_leaf,
143         mtry=mtry_val,
144         importance="impurity_corrected",
145         max_depth=max_depth,
146         classification=True
147     )
148
149     task_type = rf_result.rx2('treetype')[0]
150     print(f"Ranger task type: {task_type}")
151
152     importance_r = rf_result.rx2('variable.importance')
153     importance_values = np.array(importance_r)
154     feature_names_r = list(robots.r['names'](importance_r))
155
156     importance_series = pd.Series(importance_values, index=
feature_names_r)
157     importance_series.index = importance_series.index.map(name_map)
158
159     print("Successfully computed variable importance for CLASSIFICATION")
160 )
161
162     return importance_series.sort_values(ascending=False)
163
164 except Exception as e:
165     warnings.warn(f"Ranger implementation failed: {e}")
166     traceback.print_exc()
167     return None

```

Листинг 4: Файл r_implementations.py