

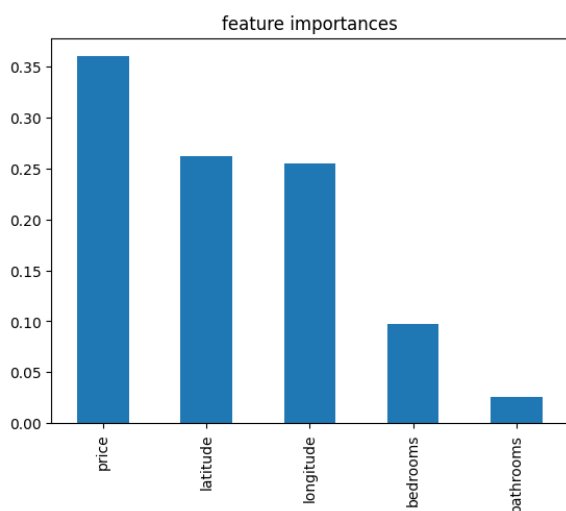
## Отчет.

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.

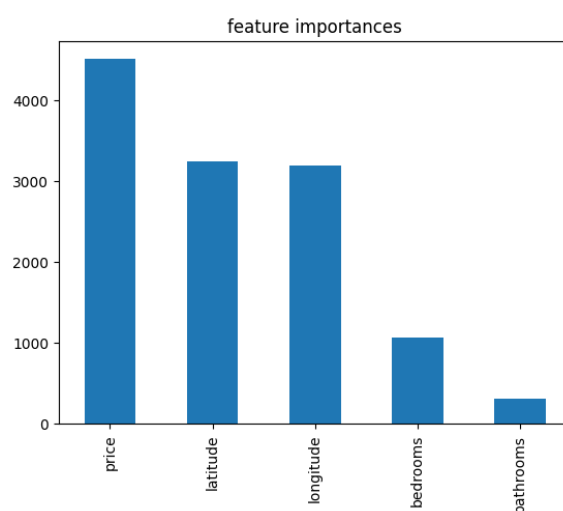
## Эксперименты.

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

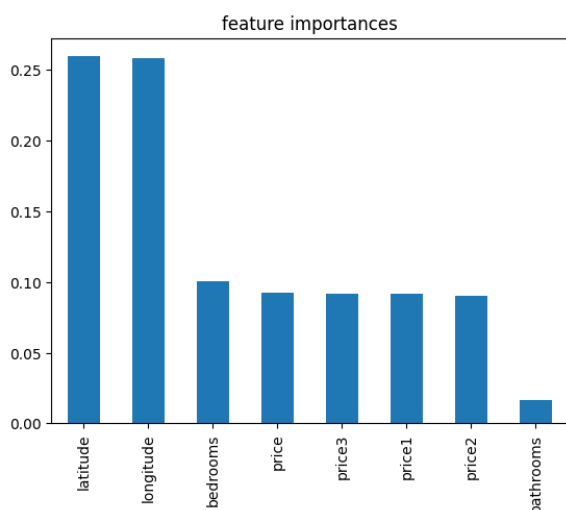
Подбор оптимальных гиперпараметров в реализации sklearn я осуществлял с помощью *Optuna*, и полученные параметры я использовал в реализации на R.



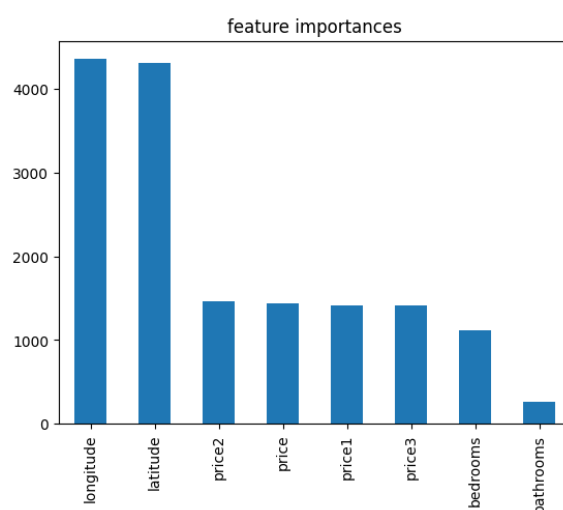
(a) Python: Исходный датасет



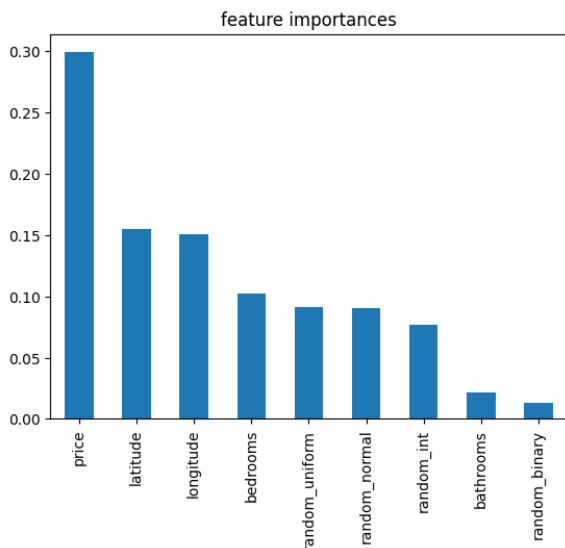
(a) R: Исходный датасет



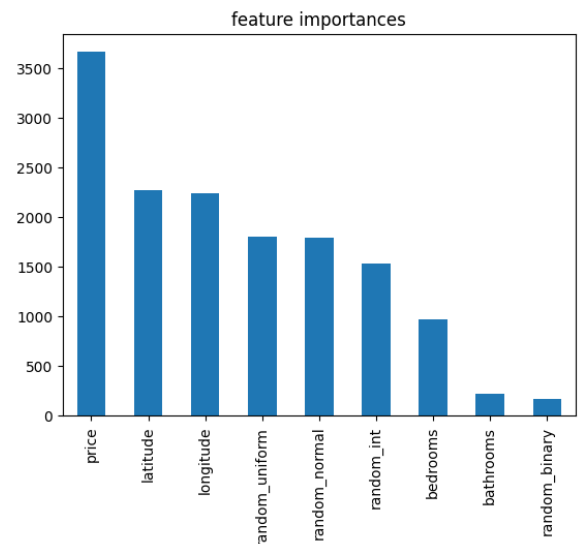
(b) Python: добавлены коррелированные признаки



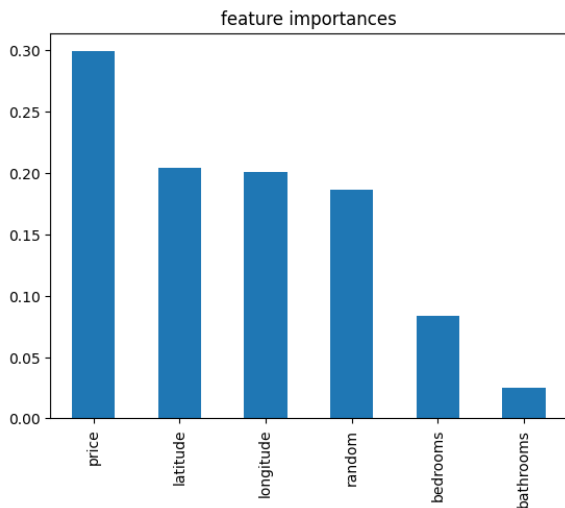
(b) R: добавлены коррелированные признаки



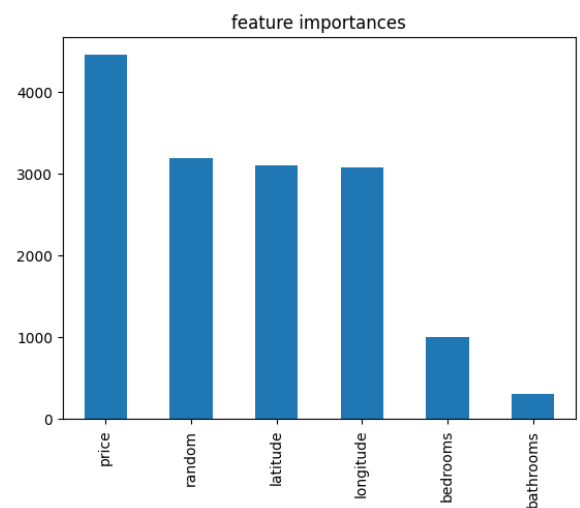
(c) Python: много нерелевантных признаков



(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.

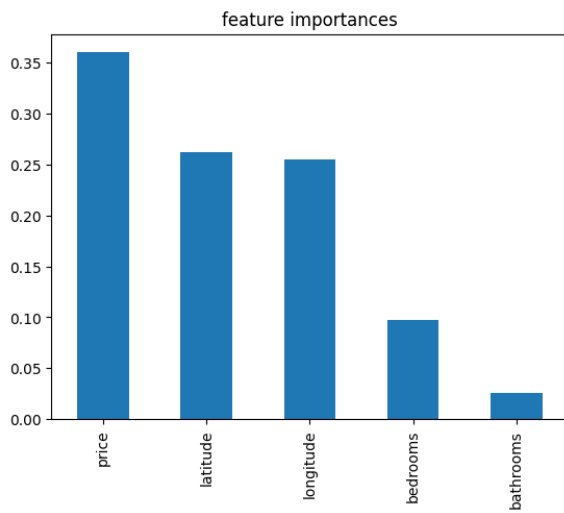
В примерах (c) и (d) я создал несколько признаков, которые никак не связаны с *target*.

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*?

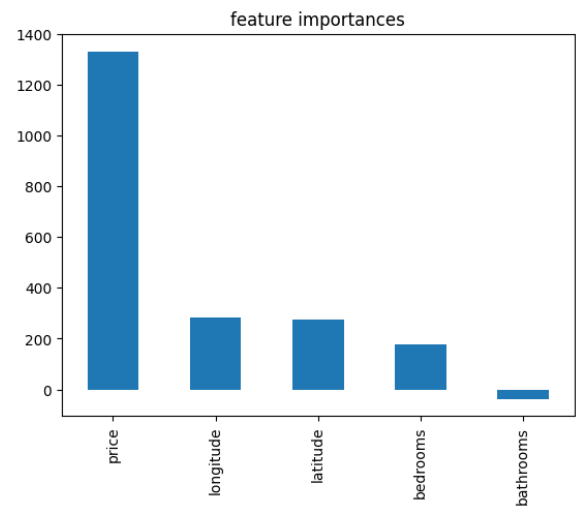
Ниже приведена таблица с оптимальными параметрами для каждого из примеров (a)-(d):

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

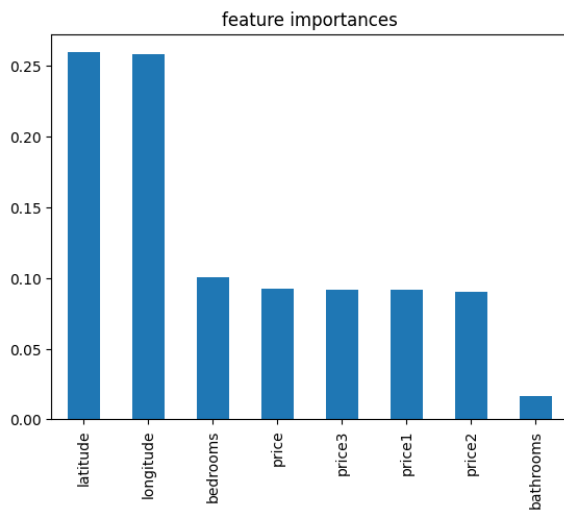
# Sklearn vs R-ranger



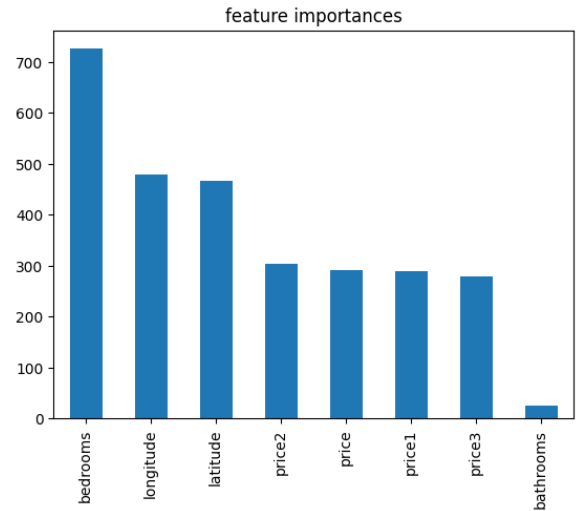
(a) Sklearn: Исходный датасет



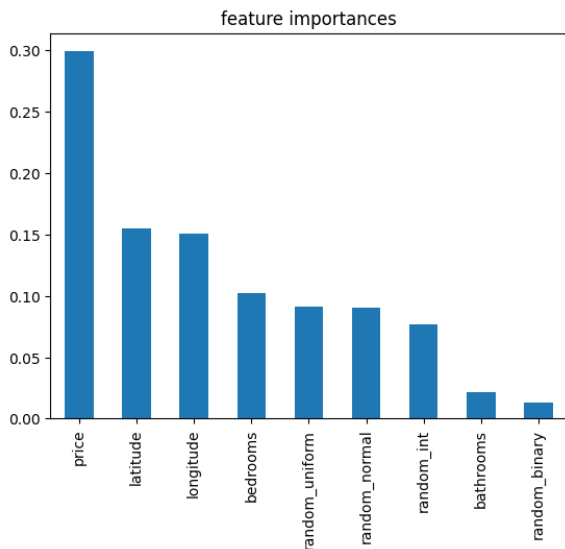
(a) R-ranger: Исходный датасет



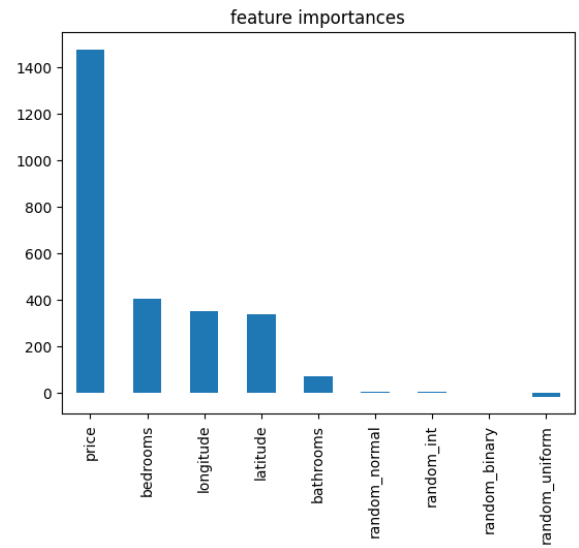
(b) Sklearn: добавлены коррелированные признаки



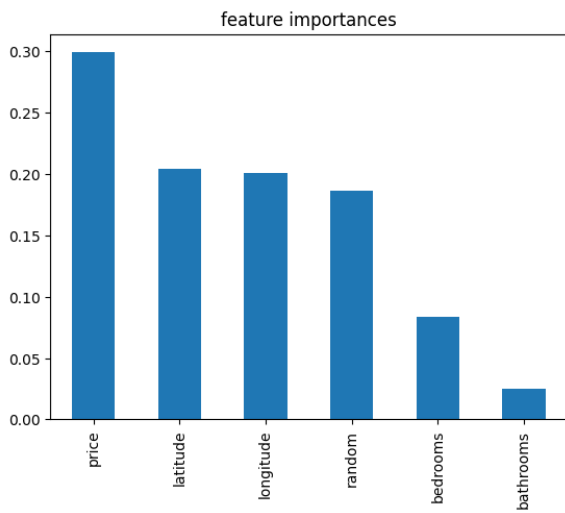
(b) R-ranger: добавлены коррелированные признаки



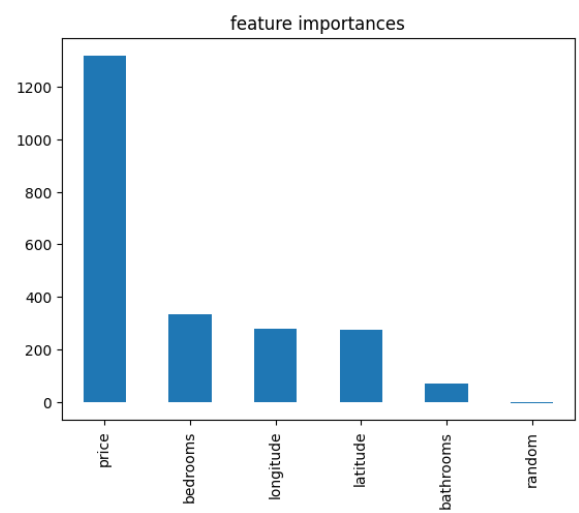
(c) Sklearn: много нерелевантных признаков



(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 )  
10  
11 from .pic import (  
12     picture  
13 )  
14  
15  
16 __version__ = "1.0.0"  
17 __all__ = ['sklearn_importance', 'picture',  
18            'r_randomforest_importance']
```

Листинг 1: Файл `__init__.py`

```
1 import matplotlib.pyplot as plt  
2  
3  
4 def picture(fi):  
5     fi.plot(kind='bar')  
6     plt.title('feature importances')  
7     plt.show()  
8     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', 5, 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)

```

ЛИСТИНГ 3: Файл python\_implementations.py

```

1 from rpy2.robjjects.vectors import StrVector
2 import pandas as pd
3 import numpy as np
4 import re
5 import rpy2.robjjects as robjjects
6 from rpy2.robjjects import pandas2ri
7 from rpy2.robjjects.packages import importr
8 import warnings
9 import traceback
10 from rpy2 import robjjects
11
12
13 def clean_feature_names(names):
14     """
15                                     R."""
16     cleaned_names = [name.replace('.', '_').replace(
17         '-', '_').replace(' ', '_') for name in names]
18     return cleaned_names
19
20 def r_randomforest_importance(X, y, data_path=None, n_estimators=1760,
21     max_depth=26, min_samples_leaf=5, max_features='sqrt'):
22     try:
23         base = importr('base')
24         utils = importr('utils')
25         randomForest = importr('randomForest')
26
27         original_features = list(X.columns)
28         cleaned_features = clean_feature_names(original_features)
29
30         df_for_r = X.copy()
31         df_for_r.columns = cleaned_features
32         df_for_r['target'] = y.values
33
34         with robjjects.conversion.localconverter(robjjects.default_converter +
35             pandas2ri.converter) as cv:
36             r_df = robjjects.conversion.py2rpy(df_for_r)

```

```

36 formula_str = "target ~ " + " + ".join(cleaned_features)
37 formula = robjects.Formula(formula_str)
38
39 print(f"Training R Random Forest with formula: {formula_str}")
40
41 n_samples = X.shape[0]
42 n_features = len(cleaned_features)
43 if max_features == 'sqrt':
44     max_features = round(np.sqrt(n_features))
45 else:
46     max_features = round(n_features * max_features)
47
48 rf_result = randomForest.randomForest(
49     formula,
50     data=r_df,
51     ntree=n_estimators,
52     nodesize=min_samples_leaf,
53     maxnodes=2**max_depth,
54     mtry=max_features,
55     importance=True
56 )
57
58 importance_r = randomForest.importance(rf_result, type=2, scale=
False)
59
60 importance_matrix = np.array(importance_r)
61
62 feature_names_r = list(robjects.r['rownames'](importance_r))
63
64 print(f"Importance matrix shape: {importance_matrix.shape}")
65 print(f"Feature names from R: {feature_names_r}")
66
67 if importance_matrix.ndim == 2:
68     importance_values = importance_matrix[:, 0]
69 else:
70     importance_values = importance_matrix
71
72 importance_series = pd.Series(
73     importance_values, index=original_features)
74
75 print("Successfully computed variable importance")
76 print(importance_series)
77
78 return importance_series.sort_values(ascending=False)
79
80 except Exception as e:
81     warnings.warn(f"R implementation failed: {e}")
82     import traceback
83     traceback.print_exc()
84     return None
85
86
87 def r_ranger_importance_air(X, y, n_estimators=1760, max_depth=26,
min_samples_leaf=5, max_features='sqrt'):
88     try:
89         base = importr('base')
90         ranger = importr('ranger')
91
92         original_features = list(X.columns)
93

```

```

94     cleaned_features = clean_feature_names(original_features)
95
96     df_for_r = X.copy()
97     df_for_r.columns = cleaned_features
98
99     df_for_r['target'] = y.astype(str).values
100
101     name_map = dict(zip(cleaned_features, original_features))
102
103     with objects.conversion.localconverter(objects.default_converter +
104 pandas2ri.converter):
105         r_df = objects.conversion.py2rpy(df_for_r)
106
107         r_target = base.factor(r_df.rx2('target'))
108
109         r_cleaned_features = StrVector(cleaned_features)
110
111         r_df_features_only = r_df.rx(True, r_cleaned_features)
112
113         r_df = base.cbind(r_df_features_only, target=r_target)
114
115         formula_str = "target ~ " + " + ".join(cleaned_features)
116         formula = objects.Formula(formula_str)
117
118         print(f"Training R Ranger Forest with formula: {formula_str}")
119
120         target_type = objects.r['class'](r_df.rx2('target'))
121         target_levels = objects.r['levels'](r_df.rx2('target'))
122         print(f"Target type in R: {list(target_type)}")
123         print(f"Target levels: {list(target_levels)}")
124
125         n_features = len(cleaned_features)
126         if max_features == 'sqrt':
127             mtry_val = round(np.sqrt(n_features))
128         elif isinstance(max_features, float):
129             mtry_val = round(n_features * max_features)
130         else:
131             mtry_val = max_features
132
133         rf_result = ranger.ranger(
134             formula,
135             data=r_df,
136             num_trees=n_estimators,
137             min_node_size=min_samples_leaf,
138             mtry=mtry_val,
139             importance="impurity_corrected",
140             max_depth=max_depth,
141             classification=True
142         )
143
144         task_type = rf_result.rx2('treetype')[0]
145         print(f"Ranger task type: {task_type}")
146
147         importance_r = rf_result.rx2('variable.importance')
148         importance_values = np.array(importance_r)
149         feature_names_r = list(objects.r['names'](importance_r))
150
151         importance_series = pd.Series(importance_values, index=
feature_names_r)
152         importance_series.index = importance_series.index.map(name_map)

```



```

152
153     print("Successfully computed variable importance for CLASSIFICATION"
154
155     )
156
157     return importance_series.sort_values(ascending=False)
158
159 except Exception as e:
160     warnings.warn(f"Ranger implementation failed: {e}")
161     traceback.print_exc()
162     return None

```

Листинг 4: Файл r\_implementations.py