

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

Курс «Технологии машинного обучения»

Отчёт по рубежному контролю №2

«Технологии разведочного анализа и обработки данных.»

Вариант № 18

Выполнил:

Табахов Е.В.

группа ИУ5-62Б

Дата: 22.05.25

Подпись:

Проверил:

Гапанюк Ю.Е.

Дата:

Подпись:

2025 г.

Задание(вариант 18):

Для заданного набора данных (по Вашему варианту) постройте модели классификации или регрессии (в зависимости от конкретной задачи, рассматриваемой в наборе данных). Для построения моделей используйте методы 1 и 2 (по варианту для Вашей группы). Оцените качество моделей на основе подходящих метрик качества (не менее двух метрик). Какие метрики качества Вы использовали и почему? Какие выводы Вы можете сделать о качестве построенных моделей? Для построения моделей необходимо выполнить требуемую предобработку данных: заполнение пропусков, кодирование категориальных признаков, и т.д.

ИУ5-62Б, ИУ5Ц-82Б: Метод опорных векторов, Случайный лес.

Датасет: <https://www.kaggle.com/datasets/agrafintech/world-happiness-index-and-inflation-dataset>

Загрузка библиотек и необходимых модулей

```
[5] 1 import pandas as pd
    2 import numpy as np
    3 import matplotlib.pyplot as plt
    4 import seaborn as sns
    5 from sklearn.model_selection import train_test_split, GridSearchCV
    6 from sklearn.preprocessing import StandardScaler, OneHotEncoder
    7 from sklearn.compose import ColumnTransformer
    8 from sklearn.pipeline import Pipeline
    9 from sklearn.impute import SimpleImputer
   10 from sklearn.svm import SVR
   11 from sklearn.ensemble import RandomForestRegressor
   12 from sklearn.metrics import mean_squared_error, r2_score, mean_absolute_error
```

Загрузка данных

```
1 df = pd.read_csv('WHI_Inflation.csv')
2
3 # Просмотр первых строк и базовой информации
4 print(df.head())
5 print(df.info())
6 print(df.isnull().sum())
```

```
Country Year  Headline Consumer Price Inflation \
0  Afghanistan  2015      -0.660
1  Afghanistan  2016       4.380
2  Afghanistan  2017       4.976
3  Afghanistan  2018       0.630
4  Afghanistan  2019       2.302

Energy Consumer Price Inflation  Food Consumer Price Inflation \
0      -4.250000      -0.840000
1       2.070000       5.670000
2       4.440000       6.940000
3       1.474185      -1.045952
4      -2.494359       3.794770

Official Core Consumer Price Inflation  Producer Price Inflation \
0              0.219999              NaN
1             5.192760              NaN
2             5.423228              NaN
```

Предобработка данных

```
[7] 1 # Удаляем строки с пропущенными значениями в целевой переменной
    2 df = df.dropna(subset=['Score'])
```

Определяем признаки и целевую переменную

```
[8] 1 X = df.drop(['Score', 'Country'], axis=1) # Удаляем целевую переменную и идентификатор страны
    2 y = df['Score']
```

Разделение на категориальные и числовые признаки

```
[9] 1 categorical_features = ['Continent/Region']
    2 numeric_features = [col for col in X.columns if col not in categorical_features]
```

Разделение на обучающую и тестовую выборки

```
[10] 1 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)
```

Создание препроцессоров

```
[11] 1 numeric_transformer = Pipeline(steps=[
    2     ('imputer', SimpleImputer(strategy='median')), # Заполнение пропусков медианой
    3     ('scaler', StandardScaler()) # Стандартизация числовых признаков
    4 ])
    5
    6 categorical_transformer = Pipeline(steps=[
    7     ('imputer', SimpleImputer(strategy='most_frequent')), # Заполнение пропусков наиболее частым значением
    8     ('onehot', OneHotEncoder(handle_unknown='ignore')) # One-hot кодирование категориальных признаков
    9 ])
```

Объединение препроцессоров

```
1 preprocessor = ColumnTransformer(
2     transformers=[
3         ('num', numeric_transformer, numeric_features),
4         ('cat', categorical_transformer, categorical_features)
5     ]
6 )
```

Создание и обучение моделей

1. Метод опорных векторов (SVR)

```
[13] 1 svr_pipeline = Pipeline(steps=[
    2     ('preprocessor', preprocessor),
    3     ('regressor', SVR(kernel='rbf'))
    4 ])
    5
    6 # Подбор гиперпараметров для SVR
    7 param_grid_svr = {
    8     'regressor__C': [0.1, 1, 10, 100],
    9     'regressor__gamma': ['scale', 'auto', 0.1, 0.01]
   10 }
   11
   12 grid_search_svr = GridSearchCV(svr_pipeline, param_grid_svr, cv=5, scoring='neg_mean_squared_error', n_jobs=
   13 grid_search_svr.fit(X_train, y_train)
   14
   15 # Выбор лучшей модели SVR
   16 best_svr = grid_search_svr.best_estimator_
   17 print(f"Лучшие параметры SVR: {grid_search_svr.best_params_}")
```

Лучшие параметры SVR: {'regressor__C': 10, 'regressor__gamma': 'scale'}

2. Случайный лес (Random Forest)

```
[14] 1 rf_pipeline = Pipeline(steps=[
2     ('preprocessor', preprocessor),
3     ('regressor', RandomForestRegressor(random_state=42))
4 ])
5
6 # Подбор гиперпараметров для Random Forest
7 param_grid_rf = {
8     'regressor__n_estimators': [50, 100, 200],
9     'regressor__max_depth': [None, 10, 20, 30],
10    'regressor__min_samples_split': [2, 5, 10]
11 }
12
13 grid_search_rf = GridSearchCV(rf_pipeline, param_grid_rf, cv=5, scoring='neg_mean_squared_error', n_jobs=-1)
14 grid_search_rf.fit(X_train, y_train)
15
16 # Выбор лучшей модели Random Forest
17 best_rf = grid_search_rf.best_estimator_
18 print(f"Лучшие параметры Random Forest: {grid_search_rf.best_params_}")
```

Лучшие параметры Random Forest: {'regressor__max_depth': 20, 'regressor__min_samples_split': 2, 'regressor__n_estimators': 200}

Предсказания на тестовой выборке

```
[15] 1 svr_pred = best_svr.predict(X_test)
2 rf_pred = best_rf.predict(X_test)
```

Оценка качества моделей

Метрика 1: Средняя квадратичная ошибка (MSE)

```
[16] 1 svr_mse = mean_squared_error(y_test, svr_pred)
2 rf_mse = mean_squared_error(y_test, rf_pred)
```

Метрика 2: Коэффициент детерминации (R^2)

```
[17] 1 svr_r2 = r2_score(y_test, svr_pred)
2 rf_r2 = r2_score(y_test, rf_pred)
```

Метрика 3: Средняя абсолютная ошибка (MAE)

```
[18] 1 svr_mae = mean_absolute_error(y_test, svr_pred)
2 rf_mae = mean_absolute_error(y_test, rf_pred)
```

Вывод результатов

```
[19] 1 print("\nРезультаты оценки моделей:")
2 print("Метод опорных векторов (SVR):")
3 print(f"MSE: {svr_mse:.4f}")
4 print(f"R²: {svr_r2:.4f}")
5 print(f"MAE: {svr_mae:.4f}")
6
7 print("\nСлучайный лес (Random Forest):")
8 print(f"MSE: {rf_mse:.4f}")
9 print(f"R²: {rf_r2:.4f}")
10 print(f"MAE: {rf_mae:.4f}")
```

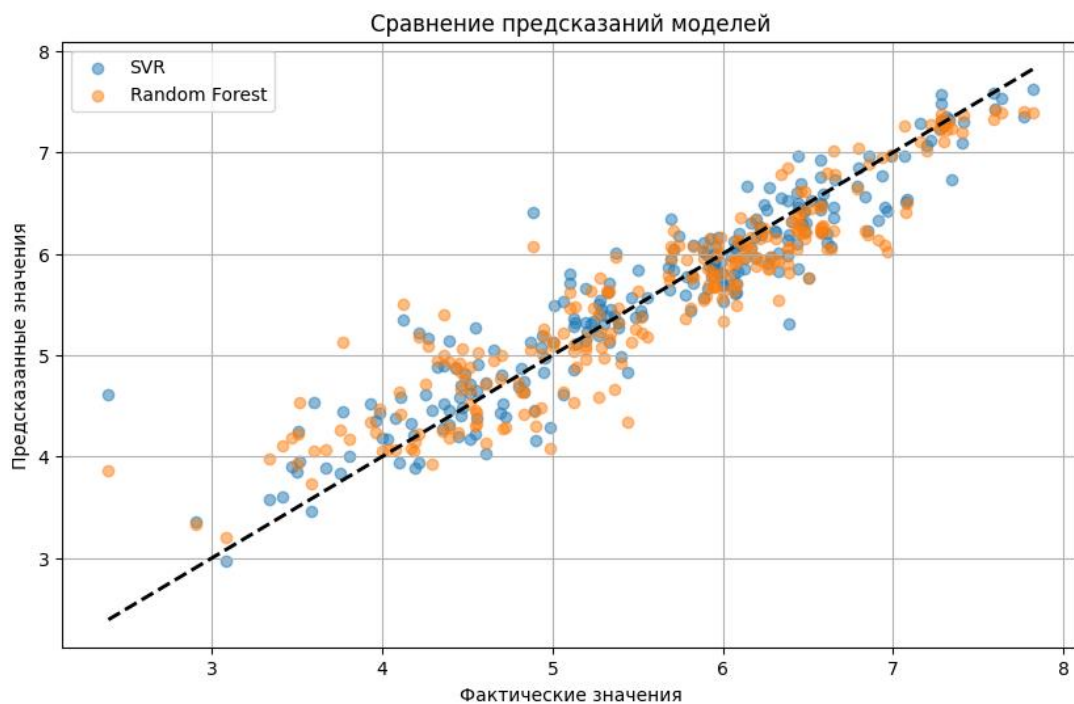


Результаты оценки моделей:
Метод опорных векторов (SVR):
MSE: 0.1393
R²: 0.8790
MAE: 0.2656

Случайный лес (Random Forest):
MSE: 0.1627
R²: 0.8587
MAE: 0.3088

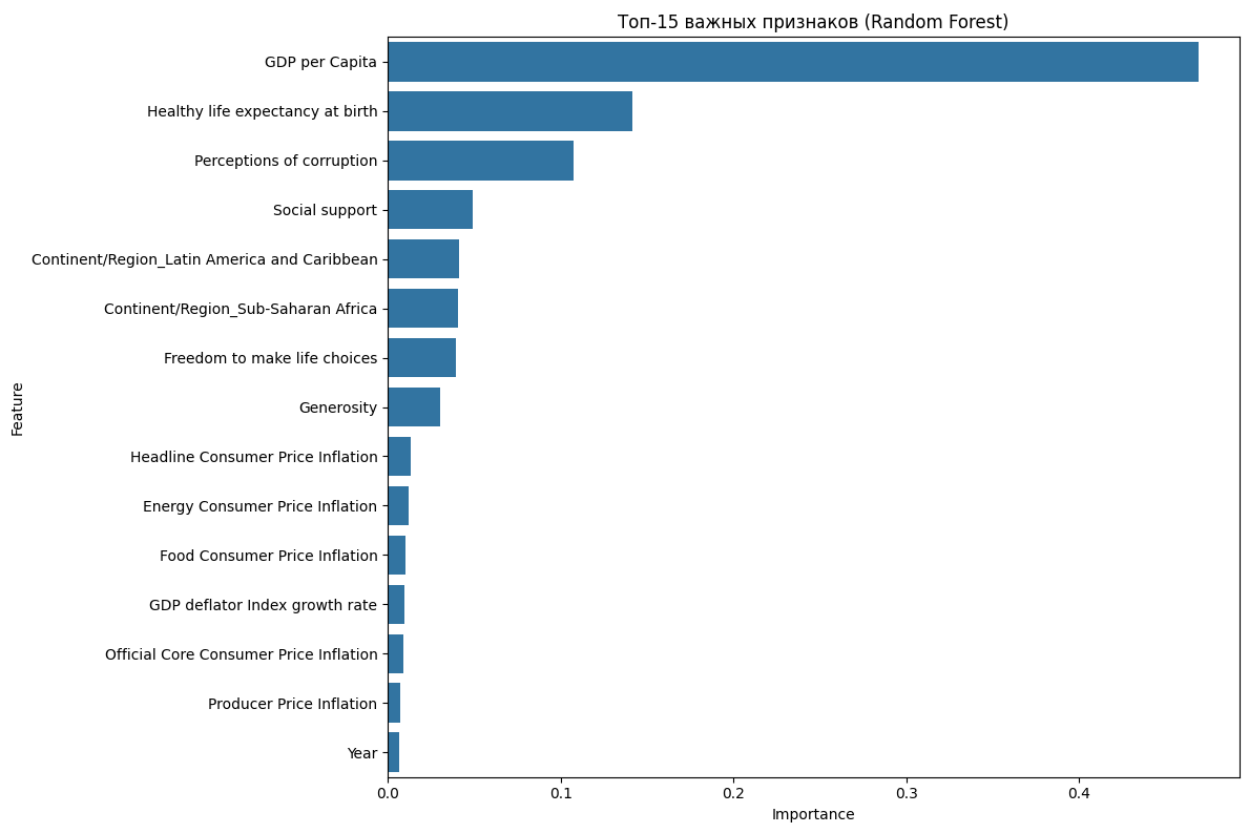
Визуализация результатов

```
1 plt.figure(figsize=(10, 6))
2 plt.scatter(y_test, svr_pred, alpha=0.5, label='SVR')
3 plt.scatter(y_test, rf_pred, alpha=0.5, label='Random Forest')
4 plt.plot([y_test.min(), y_test.max()], [y_test.min(), y_test.max()], 'k--', lw=2)
5 plt.xlabel('Фактические значения')
6 plt.ylabel('Предсказанные значения')
7 plt.title('Сравнение предсказаний моделей')
8 plt.legend()
9 plt.grid(True)
10 plt.show()
```



Визуализация важности признаков для Random Forest

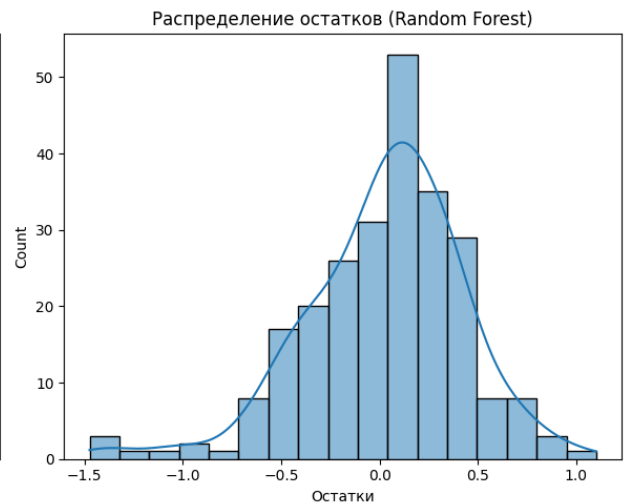
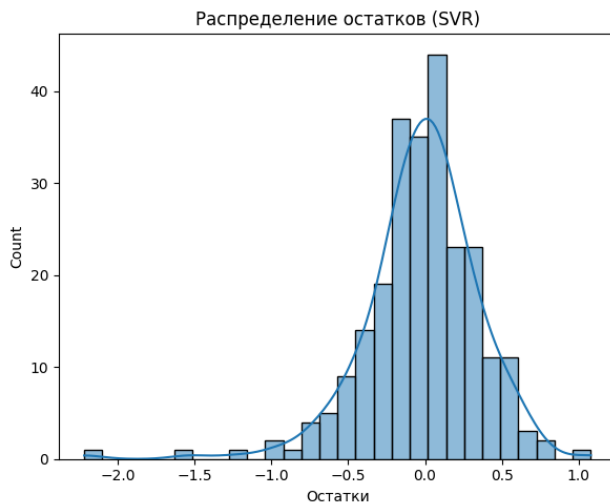
```
[21] 1 if hasattr(best_rf[-1], 'feature_importances_'):
2     # Получение имен признаков после преобразования
3     preprocessor = best_rf.named_steps['preprocessor']
4     feature_names = []
5
6     for name, trans, cols in preprocessor.transformers_:
7         if name == 'num':
8             feature_names.extend(cols)
9         elif name == 'cat':
10            # Получение названий закодированных категориальных признаков
11            ohe = trans.named_steps['onehot']
12            cat_features = ohe.get_feature_names_out(cols)
13            feature_names.extend(cat_features)
14
15    importances = best_rf[-1].feature_importances_
16
17    # Создание DataFrame для удобства сортировки
18    feature_importance_df = pd.DataFrame({
19        'Feature': feature_names,
20        'Importance': importances
21    }).sort_values(by='Importance', ascending=False)
22
23    # Визуализация топ-15 важных признаков
24    plt.figure(figsize=(12, 8))
25    sns.barplot(x='Importance', y='Feature', data=feature_importance_df.head(15))
26    plt.title('Топ-15 важных признаков (Random Forest)')
27    plt.tight_layout()
28    plt.show()
```



Анализ остатков

✓
0
сек.

```
[22] 1 plt.figure(figsize=(12, 5))
      2
      3 # Остатки для SVR
      4 plt.subplot(1, 2, 1)
      5 residuals_svr = y_test - svr_pred
      6 sns.histplot(residuals_svr, kde=True)
      7 plt.title('Распределение остатков (SVR)')
      8 plt.xlabel('Остатки')
      9
     10 # Остатки для Random Forest
     11 plt.subplot(1, 2, 2)
     12 residuals_rf = y_test - rf_pred
     13 sns.histplot(residuals_rf, kde=True)
     14 plt.title('Распределение остатков (Random Forest)')
     15 plt.xlabel('Остатки')
     16
     17 plt.tight_layout()
     18 plt.show()
```



Сравнение фактических и предсказанных значений

✓
0
сек.

```
[23] 1 results_df = pd.DataFrame({
2     'Actual': y_test,
3     'SVR_Predicted': svr_pred,
4     'RF_Predicted': rf_pred,
5     'SVR_Residuals': residuals_svr,
6     'RF_Residuals': residuals_rf
7 })
8
9 print("\nПримеры предсказаний и остатков:")
10 print(results_df.head(10))
```



Примеры предсказаний и остатков:

	Actual	SVR_Predicted	RF_Predicted	SVR_Residuals	RF_Residuals
549	4.1200	5.350135	5.502786	-1.230135	-1.382786
1050	7.2840	7.485482	7.328939	-0.201482	-0.044939
244	5.3390	5.446765	5.225187	-0.107765	0.113813
553	5.7900	5.771703	5.941561	0.018297	-0.151561
1165	6.7908	6.662474	6.637771	0.128326	0.153029
907	5.1020	5.712792	5.614744	-0.610792	-0.512744
732	5.2850	5.451921	5.074410	-0.166921	0.210590
526	5.8900	6.015968	5.956664	-0.125968	-0.066664
101	6.9370	6.767037	6.957135	0.169963	-0.020135
128	5.8220	5.713713	5.775080	0.108287	0.046920

Вывод итоговых выводов

✓
0
сек.

```
1 # Вывод итоговых выводов
2 print("\nИтоговые выводы:")
3 if rf_r2 > svr_r2:
4     print(f"Модель Random Forest показала лучшие результаты с R² = {rf_r2:.4f} против R² = {svr_r2:.4f} у SVR")
5     best_model = "Random Forest"
6 else:
7     print(f"Модель SVR показала лучшие результаты с R² = {svr_r2:.4f} против R² = {rf_r2:.4f} у Random Forest")
8     best_model = "SVR"
9
10 print(f"\nВыбранные метрики качества:")
11 print("1. Коэффициент детерминации (R²) - показывает долю дисперсии зависимой переменной, объясняемую моделью.")
12 print("2. Средняя абсолютная ошибка (MAE) - среднее абсолютное отклонение предсказанных значений от фактических.")
13 print("3. Средняя квадратичная ошибка (MSE) - среднее квадратов разностей между предсказанными и фактическими значениями.")
```



Итоговые выводы:

Модель SVR показала лучшие результаты с $R^2 = 0.8790$ против $R^2 = 0.8587$ у Random Forest.

Выбранные метрики качества:

1. Коэффициент детерминации (R^2) - показывает долю дисперсии зависимой переменной, объясняемую моделью.
2. Средняя абсолютная ошибка (MAE) - среднее абсолютное отклонение предсказанных значений от фактических.
3. Средняя квадратичная ошибка (MSE) - среднее квадратов разностей между предсказанными и фактическими значениями.