

Министерство науки и высшего образования РФ  
федеральное государственное автономное  
образовательное учреждение высшего образования  
«Омский государственный технический университет»

Факультет информационных технологий и компьютерных систем  
Кафедра «Прикладная математика и фундаментальная информатика»

**Расчетно-графическая работа**  
по дисциплине **Машинное обучение и большие данные**  
**Тема: Разработка Web-приложения (дашборда)**  
**для инференса моделей ML и анализа данных**

Студента	Курпенова Куата Ибраимовича
Курс	<u>2</u> Группа <u>ФИТ-222</u>
Направление	<u>02.03.02 Фундаментальная информатика и</u> <u>информационные технологии</u>
Руководитель	<u>доцент</u> <u>Моисеева Наталья Александровна</u>
Выполнил	28.10.2024
Проверил	28.10.2024

<b>Введение .....</b>	<b>3</b>
<b>Основная часть .....</b>	<b>5</b>
<b>Список использованных источников .....</b>	<b>10</b>
<b>Приложение А.....</b>	<b>11</b>
<b>Приложение Б.....</b>	<b>13</b>

## Введение

Тема машинного обучения (ML) является крайне актуальной и востребованной в современном мире. Она находит применение во многих областях, включая медицину, финансы, транспорт, маркетинг и другие. Развитие решений на основе ML позволяет автоматизировать процессы, улучшать прогнозирование, оптимизировать принятие решений и создавать инновационные продукты и услуги.

Тема машинного обучения (ML) является крайне актуальной и востребованной в современном мире. Она находит применение во многих областях, включая медицину, финансы, транспорт, маркетинг и другие. Развитие решений на основе ML позволяет автоматизировать процессы, улучшать прогнозирование, оптимизировать принятие решений и создавать инновационные продукты и услуги.

1. Scikit-learn (Sklearn): Это одна из наиболее популярных библиотек машинного обучения для языка Python. Она предоставляет широкий спектр алгоритмов машинного обучения, обработки данных, предобработки, валидации моделей и метрик для оценки результатов. Sklearn обладает простым и интуитивно понятным API, и его можно использовать для классификации, регрессии, кластеризации и других задач.

2. Matplotlib: Это библиотека визуализации данных для языка Python. Она позволяет создавать простые и сложные графики, включая гистограммы, диаграммы рассеяния, линейные и логарифмические графики и др. Matplotlib играет важную роль в анализе и визуализации результатов ML-моделей, представлении данных и понимании моделирования.

3. Pandas: Это библиотека для обработки и анализа данных в Python. Она предоставляет удобные структуры данных, такие как DataFrames, для эффективной работы с табличными данными. pandas позволяет проводить предобработку данных, выполнение операций фильтрации, агрегации и трансформации, а также интеграцию данных с моделями машинного обучения. Основываясь на полученных знаниях и использованных инструментах, мы разработали web-приложение (дашборд) с использованием библиотеки Streamlit. Главная цель приложения состоит в том, чтобы представить модели машинного обучения и результаты анализа данных с возможностью их визуализации в различных форматах. Это позволяет пользователям легко взаимодействовать с моделями и получать наглядные представления результатов анализа данных.

4. Streamlit: Это библиотека для создания интерактивных веб-приложений для машинного обучения и анализа данных на языке Python. Streamlit позволяет разработчикам быстро и просто создавать красивые дашборды, включая визуализацию данных, интерактивные элементы

управления и интеграцию моделей машинного обучения. Он может быть полезным инструментом для демонстрации результатов и предоставления доступа к моделям в режиме реального времени.

## Основная часть

### 1. Описание выбранных моделей.

a. Naïve Bayes: простой вероятностный классификатор, основанный на применении теоремы Байеса со строгими (наивными) предположениями о независимости. В зависимости от точной природы вероятностной модели, наивные байесовские классификаторы могут обучаться очень эффективно. Во многих практических приложениях для оценки параметров для наивных байесовых моделей используют метод максимального правдоподобия; другими словами, можно работать с наивной байесовской моделью, не веря в байесовскую вероятность и не используя байесовские методы. Несмотря на наивный вид и, несомненно, очень упрощенные условия, наивные байесовские классификаторы часто работают намного лучше нейронных сетей во многих сложных жизненных ситуациях.

b. Бэггинг (Bagging): Бэггинг - это метод ансамблевого обучения, в котором создается несколько независимых моделей на основе подвыборки обучающих данных. Затем, для получения окончательного прогноза, результаты этих моделей усредняются (в случае регрессии) или используется голосование (в случае классификации). Бэггинг помогает уменьшить дисперсию модели и повысить её стабильность.

c. Градиентный бустинг (Gradient Boosting): Градиентный бустинг - это метод ансамблевого обучения, который обучает последовательность слабых моделей (например, деревьев решений) в итеративном режиме. Каждая новая модель приближает остатки предыдущей модели. Затем результаты этих моделей суммируются для получения окончательного прогноза. Градиентный бустинг является мощным методом, позволяющим получить лучшую модель, чем только одиночная модель.

d. Стекинг (Stacking): Стекинг - это метод ансамблевого обучения, который комбинирует результаты нескольких базовых моделей, обученных на одних и тех же данных. Базовые модели предоставляют свои предсказания, а затем метамодель используется для объединения этих предсказаний и получения окончательного прогноза.

e. Глубокая полносвязная нейронная сеть: Глубокая полносвязная нейронная сеть - это модель машинного обучения, состоящая из нескольких слоев нейронов, где каждый нейрон из предыдущего слоя соединен с каждым нейроном в следующем слое. Эта модель может иметь множество слоев и обучается на больших объемах данных. Глубокая полносвязная нейронная сеть широко используется для выполнения сложных задач, таких как распознавание образов и обработка естественного языка.

### 2. Процесс сериализации моделей.

a. Обученные модели были сериализованы с использованием библиотеки joblib.

### 3. Процесс разработки дашборда.

а. Web-страница 1: Информация о разработчике. На рисунке 1 приложения А показано как выглядит разработанная web-страничка. Ниже представлен код реализации данной web-страницы.

```
1 import streamlit as st
2
3 st.set_page_config(page_title="Author")
4
5 st.title("ML model inferences demo")
6
7 full_name = "Kuat Kurpenov"
8 st.write(f"Full name: {full_name}")
9
10 study_group = "FIT-222"
11 st.write(f"Study group: {study_group}")
```

б. Web-страница 2: Информация о наборе данных. На рисунке 2 приложения А показано как выглядит разработанная web-страничка. Ниже представлена часть кода реализации данной web-странички.

```
1 import pandas as pd
2 import streamlit as st
3
4 st.set_page_config(page_title="Data")
5
6 st.title("Information about data")
7
8 st.write(
9     """
10     The dataset is dedicated to monitoring car prices in Moldova.
11     Using this dataset, you can solve the regression problem -
12     predict car prices depending on the brand, mileage, and so on.
13     """
14 )
15
16 st.subheader("Example of preprocessed data")
17 df = pd.read_csv("../data/preprocessed_moldova_cars_dataset.csv")
18 df.drop("Unnamed: 0", axis=1, inplace=True)
19 st.dataframe(df.head())
20
21 st.subheader("Description of columns")
22 st.markdown(
23     """
24     - `make` &mdash; brand of car
25     - `model` &mdash; model of car
```

с. Web-страница 3: Визуализация зависимостей в наборе данных. На рисунках 3-7 приложения А показано как выглядит разработанная web-страничка. Ниже представлена часть кода реализации данной web-страницы.

```
1 from collections import Counter
2
3 import matplotlib.pyplot as plt
4 import pandas as pd
5 import seaborn as sns
6 import streamlit as st
7
8 st.set_page_config(page_title="Data visualization")
9
10 st.title("Data visualization")
11
12 df = pd.read_csv("../data/preprocessed_moldova_cars_dataset.csv")
13 df.drop("Unnamed: 0", axis=1, inplace=True)
14
15 fig, axs = plt.subplots(2, 2, figsize=(20, 16))
16
17 unique_makes = list(df["make"].unique())
18 counts = df["make"].value_counts() / len(df)
19 axs[0, 0].bar(unique_makes, counts)
20 axs[0, 0].set_title("Distribution of car brands")
21 axs[0, 0].set_xlabel("Car brand")
22 axs[0, 0].set_ylabel("Fraction")
23
24 sns.heatmap(df.corr(), ax=axs[0, 1])
25 axs[0, 1].set_title("Correlations heatmap")
26
27 prices_distributions = Counter(df["price"].values)
28 axs[1, 0].bar(list(prices_distributions.keys()), list(prices_distributions.values()))
29 axs[1, 0].set_title("Distribution of car prices")
30 axs[1, 0].set_xlabel("Price")
31 axs[1, 0].set_ylabel("Count")
32
33 st.pyplot(fig)
```

d. Web-страница 4: Предсказания моделей на введенных пользователем данных. На рисунках 9-10 приложения А показано как выглядит разработанная web-страничка. Ниже представлена часть кода реализации данной web-страницы.

```

90 st.set_page_config(page_title="Models")
91
92 st.title("Machine learning models")
93
94 make = st.text_input("Make:")
95 model = st.text_input("Model:")
96 year = st.text_input("Year:")
97 style = st.text_input("Style:")
98 distance = st.text_input("Distance:")
99 engine_capacity = st.text_input("Engine capacity (cm3):")
100 fuel_type = st.text_input("Fuel type:")
101 transmission = st.text_input("Transmission:")
102
103 if st.button("Get price"):
104     data = [
105         float(make),
106         float(model),
107         float(year),
108         float(style),
109         float(distance),
110         float(engine_capacity),
111         float(fuel_type),
112         float(transmission),
113     ]
114
115     predicted_prices = predict([data])
116
117     st.success(f"\
118 [+] Predicted prices (in euro):\n\
119 [+] Linear regression: {int(predicted_prices[0])}\n\
120 [+] Catboost: {int(predicted_prices[1])}\n\
121 [+] Bagging: {int(predicted_prices[2])}\n\
122 [+] Stacking: {int(predicted_prices[3])}\n\
123 [+] Neural network: {int(predicted_prices[4])}")

```

f. GitHub-репозиторий и размещение исходных файлов. Все файлы размещены в открытом репозитории GitHub. Подробное описание проекта находится в файле README. [Ссылка на репозиторий](#).

е. Веб-сервис, разработанный с использованием библиотеки streamlit был размещён на streamlit cloud и доступен по следующей [ссылке](#).



## Заключение

В рамках данной работы было создано web-приложение (dashboard) с использованием библиотеки Streamlit. Основная цель приложения заключалась в выводе результатов инференса моделей машинного обучения и их визуализации.

В рамках данной работы было создано web-приложение (dashboard) с использованием библиотеки Streamlit. Основная цель приложения заключалась в выводе результатов инференса моделей машинного обучения и их визуализации.

Выбранные модели были сериализованы и сохранены на жесткий диск, используя соответствующие инструменты (joblib).

С помощью библиотеки streamlit было реализовано web-приложение (dashboard), содержащее несколько страниц:

1. Страница с информацией о разработчике моделей машинного обучения.
2. Страница с информацией о выбранном наборе данных.
3. Страница с визуализациями зависимостей в наборе данных, используя библиотеки matplotlib и seaborn.
4. Страница, позволяющая получить предсказание соответствующей модели машинного обучения для данных вводимых пользователем.

### **Список использованных источников**

1. user guide / [Электронный ресурс] // scikit-learn : [сайт]. — URL: [https://scikit-learn.ru/user\\_guide/](https://scikit-learn.ru/user_guide/) (дата обращения: 20.10.2024).
2. API Documentation / [Электронный ресурс] // TensorFlow : [сайт]. — URL: [https://www.tensorflow.org/api\\_docs](https://www.tensorflow.org/api_docs) (дата обращения: 20.10.2024).
3. Streamlit documentation / [Электронный ресурс] // Streamlit : [сайт]. — URL: <https://docs.streamlit.io> (дата обращения: 25.10.2024).

## Приложение А

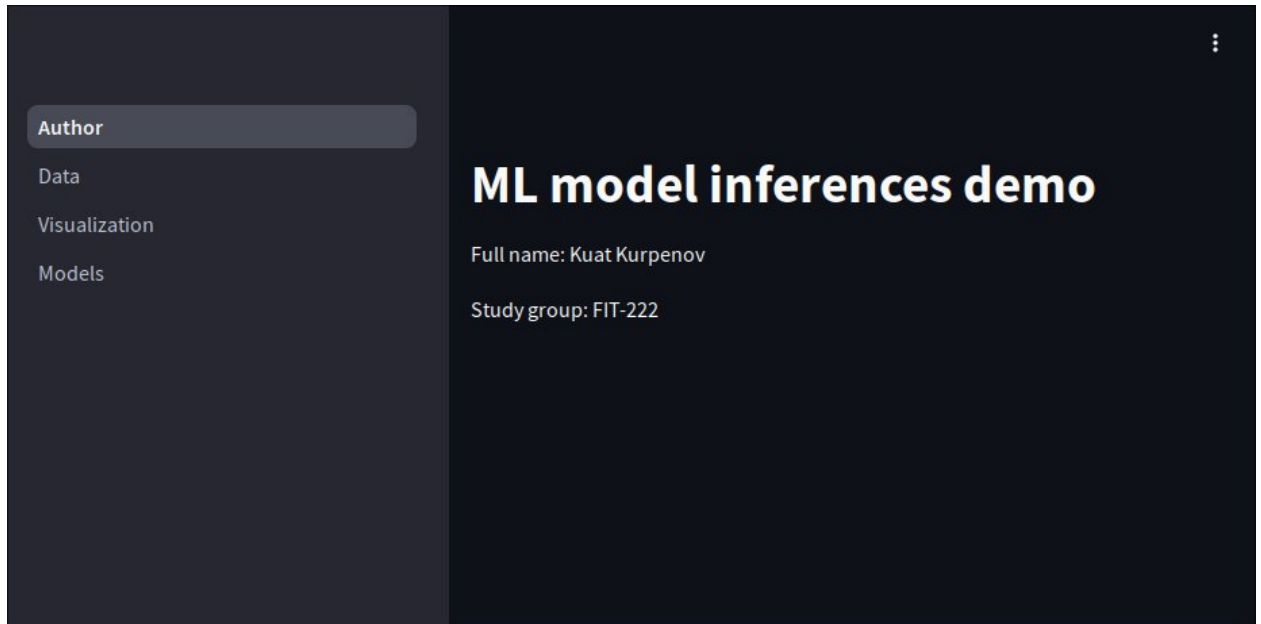


Рисунок 1 — Информация о разработчике моделей ML

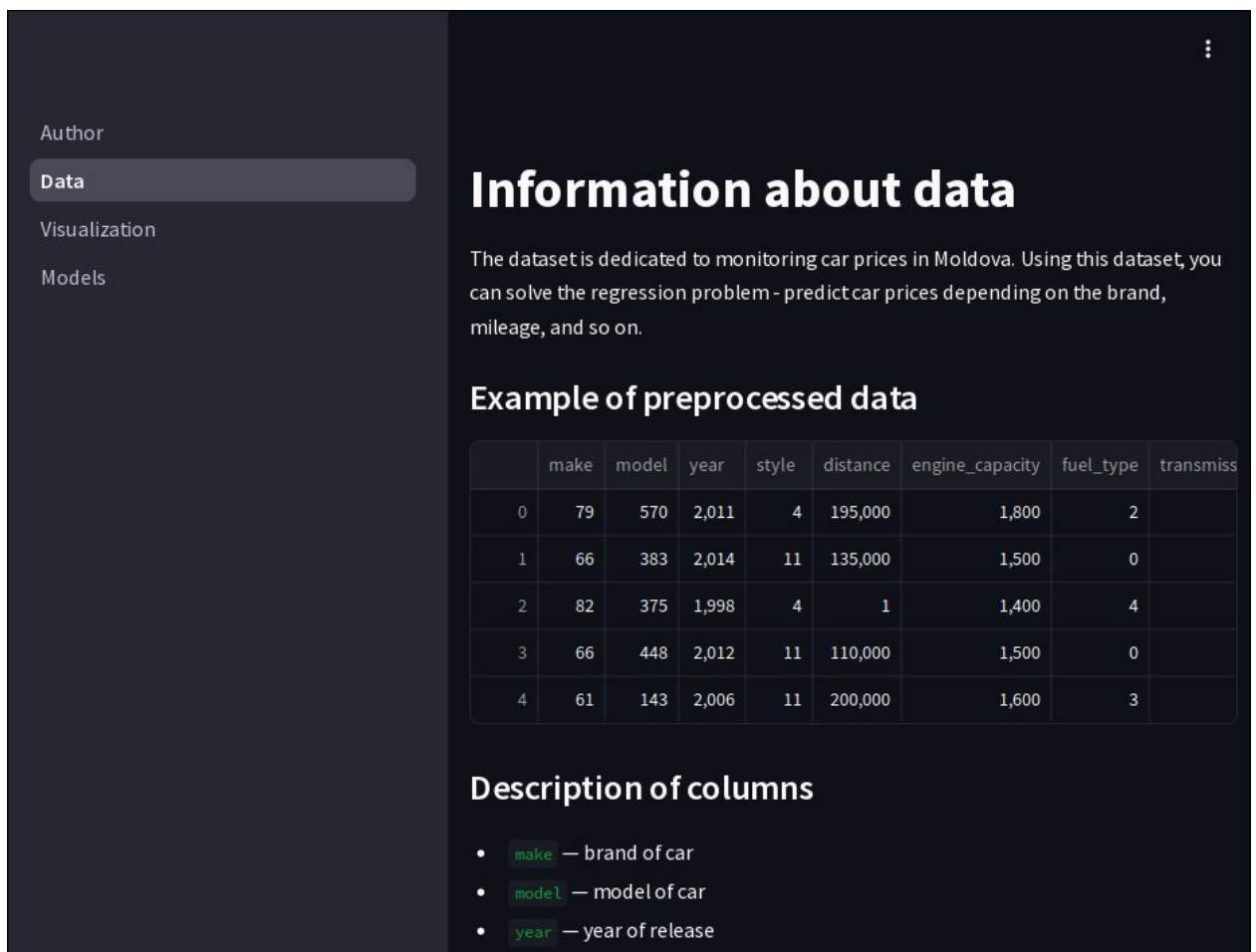


Рисунок 2 — Информация о наборе данных

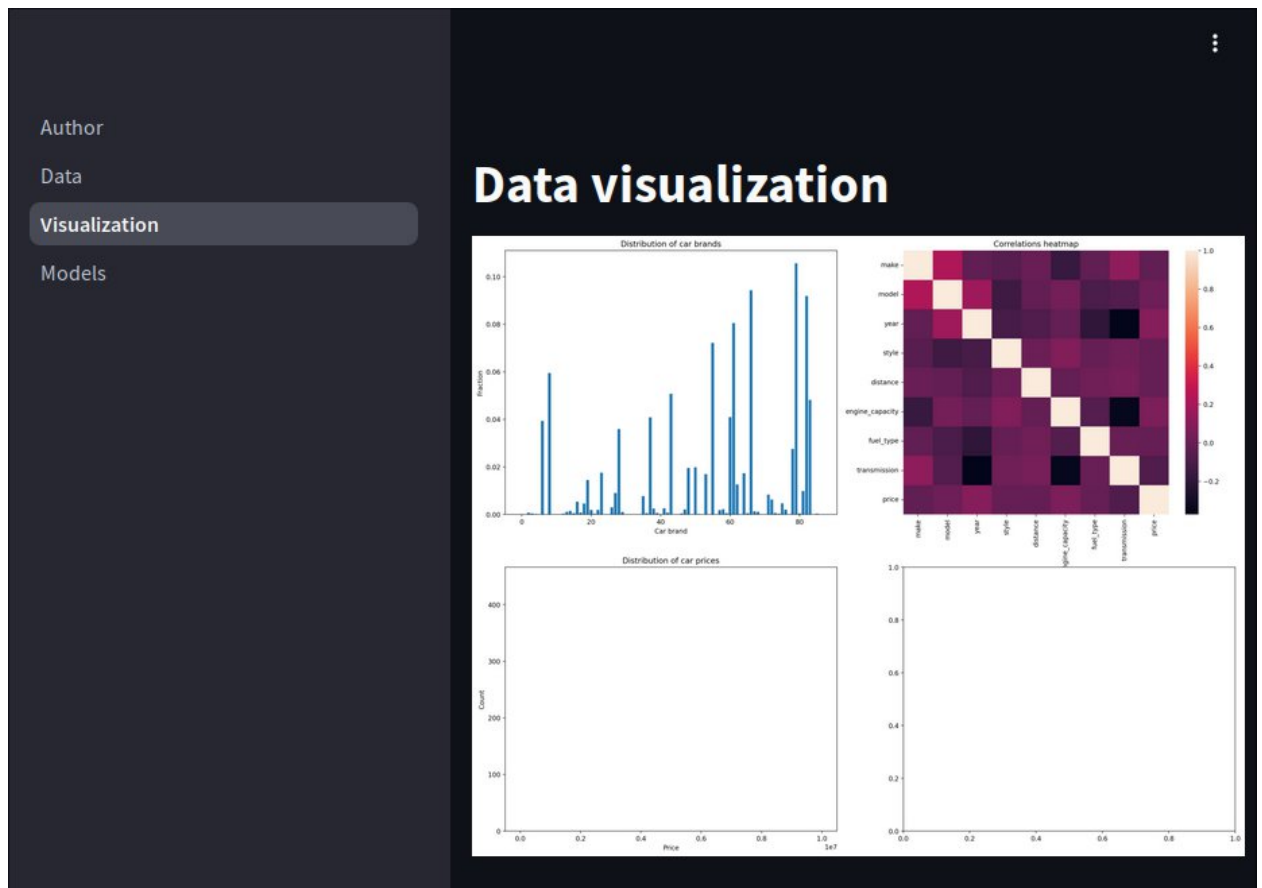


Рисунок 3 — Визуализация зависимостей в наборе данных

The dashboard features a sidebar with navigation links: Author, Data, Visualization, and Models (selected). The main content area is titled "Machine learning models" and contains four input fields for user data:

- Make:** A text input field.
- Model:** A text input field.
- Year:** A text input field.
- Style:** A text input field.

Рисунок 4 — Предсказания моделей на введённых пользователем данных

## Приложение Б

```
15 def save_model(model, model_name: str) -> None:
16     joblib.dump(model, model_name)
17
18
19 def load_model(model_dump_path: str):
20     return joblib.load(model_dump_path)
21
22
23 def fit_models(X_train, y_train, X_test, y_test) -> list[float]:
24     scaler = StandardScaler().fit(X_train)
25     X_train = scaler.transform(X_train)
26     X_test = scaler.transform(X_test)
27     save_model(scaler, "dumps/scaler.joblib")
28
29     lr = LinearRegression().fit(X_train, y_train)
30     lr_pred = lr.predict(X_test)
31     lr_mape = mean_absolute_percentage_error(y_test, lr_pred)
32     save_model(lr, "dumps/lr.joblib")
33
34     dbscan = DBSCAN(eps=0.5, min_samples=5).fit(X_train)
35     dbscan_ss = silhouette_score(X_train, dbscan.labels_)
36     save_model(dbscan, "dumps/dbscan.joblib")
37
38     cb = CatBoostRegressor(verbose=False).fit(X_train, y_train)
39     cb_pred = cb.predict(X_test)
40     cb_mape = mean_absolute_percentage_error(y_test, cb_pred)
41     save_model(cb, "dumps/catboost.joblib")
42
43     br = BaggingRegressor().fit(X_train, y_train)
44     br_pred = br.predict(X_test)
45     br_mape = mean_absolute_percentage_error(y_test, br_pred)
46     save_model(br, "dumps/bagging.joblib")
47
48     estimators = [
49         ("lr", RidgeCV()),
50         ("svr", LinearSVR()),
51     ]
52     sr = StackingRegressor(estimators).fit(X_train, y_train)
53     sr_pred = sr.predict(X_test)
54     sr_mape = mean_absolute_percentage_error(y_test, sr_pred)
55     save_model(sr, "dumps/stacking.joblib")
56
57     nn = MLPRegressor().fit(X_train, y_train)
58     nn_pred = nn.predict(X_test)
59     nn_mape = mean_absolute_percentage_error(y_test, nn_pred)
60     save_model(nn, "dumps/nn.joblib")
61
62     return [lr_mape, dbscan_ss, cb_mape, br_mape, sr_mape, nn_mape]
63
```

```

65 def predict(data: list[list[float]]) -> list[float]:
66     scaler = joblib.load("dumps/scaler.joblib")
67     data = scaler.transform(data)
68
69     lr = joblib.load("dumps/lr.joblib")
70     lr_pred = lr.predict(data)
71
72     # dbscan = joblib.load("dumps/dbscan.joblib")
73     # dbscan_pred = dbscan.predict(data)
74
75     cb = joblib.load("dumps/catboost.joblib")
76     cb_pred = cb.predict(data)
77
78     br = joblib.load("dumps/bagging.joblib")
79     br_pred = br.predict(data)
80
81     sr = joblib.load("dumps/stacking.joblib")
82     sr_pred = sr.predict(data)
83
84     nn = joblib.load("dumps/nn.joblib")
85     nn_pred = nn.predict(data)
86
87     return [lr_pred, cb_pred, br_pred, sr_pred, nn_pred]
88
90 st.set_page_config(page_title="Models")
91
92 st.title("Machine learning models")
93
94 make = st.text_input("Make:")
95 model = st.text_input("Model:")
96 year = st.text_input("Year:")
97 style = st.text_input("Style:")
98 distance = st.text_input("Distance:")
99 engine_capacity = st.text_input("Engine capacity (cm3):")
100 fuel_type = st.text_input("Fuel type:")
101 transmission = st.text_input("Transmission:")
102
103 if st.button("Get price"):
104     data = [
105         float(make),
106         float(model),
107         float(year),
108         float(style),
109         float(distance),
110         float(engine_capacity),
111         float(fuel_type),
112         float(transmission),
113     ]
114
115     predicted_prices = predict([data])

```

```
116
117     st.success(f"\n\n
118     [+] Predicted prices (in euro):\n\n
119     [+] Linear regression: {int(predicted_prices[0])}\n\n
120     [+] Catboost: {int(predicted_prices[1])}\n\n
121     [+] Bagging: {int(predicted_prices[2])}\n\n
122     [+] Stacking: {int(predicted_prices[3])}\n\n
123     [+] Neural network: {int(predicted_prices[4])}")
124
```