

Systemy Sztucznej Inteligencji

dokumentacja projektu Drzewo decyzyjne

Mikołaj Molenda, grupa 2/3
Dominik Meisner, grupa 2/3
Bartłomiej Janoszka, grupa 2/3

30 maja 2025

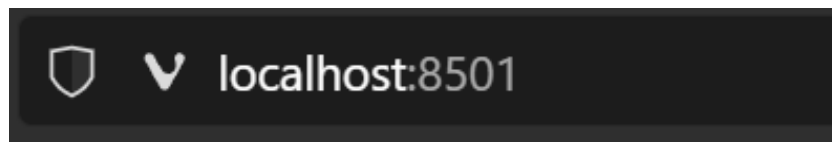
Część I

Opis programu

Stworzyliśmy model **DecisionTree** z warunkiem stopu opartym na **zbiorach miękkich**. Nasz model został zaprojektowany i wytrenowany pod dataset ludzi którzy przeżyli katastrofę Tytanic'a. W naszym projekcie użytkownik może wypełnić formularz odpowiednimi danymi, następnie taki formularz jest przesyłany i analizowany przez nasz wytrenowany model i udziela odpowiedzi, czy taka potencjalna osoba przeżyłaby katastrofę.

Instrukcja obsługi

Program uruchamiany jest przy pomocy streamlit, który uruchamia aplikacje na lokalnej domenie 'localhost'. Aby skorzystać z usługi wystarczy otworzyć przeglądarkę i wyszukać localhost. Następnie wystarczy wypełnić formularz danymi i przesłać, a program wyświetli wynik "Przeżyłeś" bądź "Nie przeżyłeś", bazując na bazie treningowej. Program jest zaopatrzony w proste i czytelne UI.



Rysunek 1: Port na którym uruchomiona jest aplikacja

Projekt SSI

Wypełnij formularz i sprawdź czy byś przeżył!

Płeć:

☒ M
☐ K

Klasa:

1

Liczba rodzeństwa + małżonek na pokładzie:

0

Liczba rodziców/dzieci na pokładzie:

0

Opłata za bilet:

0,00

Port startowy pasażera:

☒ Cherbourg
☐ QuinsTown
☐ Southampton

Sprawdź!

Rysunek 2: Wygląd zewnętrzny aplikacji (frontend)

80,00

Port startowy pasażera:

☐ Cherbourg
☒ QuinsTown
☐ Southampton

Sprawdź!

Nie przeżyłeś!

Rysunek 3: Wyświetlanie wyniku algorytmu

Dodatkowe informacje

Algorytm nie działa w sposób probabilistyczny, co oznacza, że po otrzymaniu konkretnego przykładu (np. kobieta, klasa 1, 0 rodzeństwa, 0 rodziców, opłata za bilet: 70.00, port startowy: Cherbourg), model decyduje jednoznacznie, czy dana osoba przeżyłaby, czy nie. Decyzja podejmowana jest na podstawie porównania cech tej osoby z danymi treningowymi, bez uwzględniania rozkładów prawdopodobieństw.

Jeżeli dana instancja nie wykazuje wystarczającego podobieństwa do przypadków osób, które przeżyły, zostaje automatycznie sklasyfikowana jako osoba, która by zginęła — niezależnie od proporcji takich przypadków w zbiorze danych. Taka deterministyczna natura modelu powoduje, że dla tych samych danych wejściowych wynik klasyfikacji zawsze będzie identyczny.

Dla porównania, w przypadku modelu probabilistycznego, osoba o takich samych cechach mogłaby mieć 30% szans na przeżycie, co skutkowałoby tym, że przy 10-krotnym podaniu tych samych danych moglibyśmy otrzymać około 3 wyniki pozytywne (przeżycie) i 7 negatywnych (śmierć).

Część II

Opis działania

Nasz model do przewidywania, czy dana osoba przeżyje katastrofę Titanica, bazuje na algorytmie drzewa decyzyjnego. Różnica polega jednak na tym, że klasyfikator zbiorów miękkich decyduje, czy drzewo powinno kontynuować podział danych, czy utworzyć liść.

Algorytmy

Dla każdej klasy $c \in \mathbb{Y}$

- oblicza się średnią i odchylenie standardowe każdej cechy:

$$\mu_{cj} = \mathbb{E}[x_j|y = c], \quad \sigma_{cj} = \text{std}[x_j|y = c]$$

- następnie dla próbki \mathbf{x} i danej klasy c , przekształcenie:

$$\tilde{x}_{cj} = \begin{cases} 1 & \text{jeśli } x_j \geq \mu_{cj} - \tau \sigma_{cj} \\ 0 & \text{w przeciwnym razie} \end{cases}$$

- łączna liczba dopasowań dla klasy c :

$$s_c(\mathbf{x}) = \sum_{j=1}^d \tilde{x}_{cj}$$

- predykcja:

$$\hat{y} = \arg \max_c s_c(\mathbf{x})$$

Algorithm 1: Algorytm klasyfikatora zbiorów miękkich użytego jako warunek stopu.

Data: Zbiór danych $D = \{(\mathbf{x}_i, y_i)\}_{i=1}^N$, gdzie $\mathbf{x}_i \in \mathbb{R}^d$ jest wektorem cech, a $y_i \in \mathbb{Y}$ jest etykietą klasy, $\mathbb{Y} = \{0, \dots, C - 1\}$

Result: Struktura drzewa $T : \mathbb{R}^d \rightarrow \mathbb{Y}$

Dla każdego możliwego podziału na cechę j i próg s , dane dzielone są na dwa podzbiory:

$$D_{left} = \{(\mathbf{x}_i, y_i) \in D : x_{ij} < s\}, \quad D_{right} = D \setminus D_{left}.$$

Aby wybrać najlepszy podział, oblicza się zysk informacji:

$$G = I(D) - \left(\frac{|D_{left}|}{|D|} I(D_{left}) + \frac{|D_{right}|}{|D|} I(D_{right}) \right)$$

gdzie $I(\cdot)$ to miara niejednorodności, w klasyfikacji stosuje się m.in.:

- Indeks Giniego:

$$G(D) = 1 - \sum_{k=0}^{C-1} p_k^2$$

- Entropię:

$$H(D) = - \sum_{k=0}^{C-1} p_k \log_2 p_k$$

gdzie p_k to udział klasy k w zbiorze D .

Następnie drzewo tworzone jest rekurencyjnie do momentu, gdy zostanie spełniony z jeden z warunków:

- maksymalna głębokość została osiągnięta: $\text{depth} \geq \text{max_depth}$
- liczba próbek w węźle jest mniejsza niż min_samples_split
- brak poprawy zysku informacji: $G \leq 0$
- klasyfikator miękki osiągnął dokładność powyżej progu:
 $\text{accuracy}_{soft}(D) \geq \text{soft_threshold}$

Po zbudowaniu drzewa, funkcja predykcji $T(\mathbf{x}) = \text{label}(L(\mathbf{x}))$, gdzie $L(\mathbf{x})$ to liść, do którego trafia punkt \mathbf{x} .

Algorithm 2: Algorytm drzewa decyzyjnego, w którym warunkiem stopu jest klasyfikator zbiorów miękkich.

Zbiór danych

Poniżej prezentujemy pierwotną strukturę danych pobranych ze strony kaggle.com. Dane te zostały przez nas poddane edycji i oczyszczeniu w celu dostosowania ich do założeń projektu oraz zwiększenia skuteczności zastosowanych algorytmów predykcyjnych.

Tabela 1: Pierwotna struktura danych

Kolumna	Typ danych	Opis
PassangerId	INT (PK)	Unikalny identyfikator pasażera
Iame	VARCHAR(100)	Imię i nazwisko pasażera
Sex	ENUM('male', 'female')	Płeć
Age	DECIMAL(4,1)	Wiek pasażera (np. 35.0)
Ticket	VARCHAR(50)	Numer biletu
Cabin	VARCHAR(20)	Numer kabiny (jeśli przydzielony)
Fare	DECIMAL(8,2)	Cena biletu
Embarked	ENUM('C', 'Q', 'S')	Miejsce wejścia na pokład (Cherbourg, Queenstown, Southampton)
Pclass	TINYINT	Klasa pasażerska (1, 2, 3)
Survived	BOOLEAN	TRUE – przeżył, FALSE – zginął

W celu przygotowania danych do analizy i budowy modelu predykcyjnego przeprowadziliśmy proces oczyszczania oraz selekcji zmiennych. Usunęliśmy kolumny, które nie miały istotnego wpływu na zmienną docelową "Survived" lub dla których zaobserwowana zależność była znikoma. Dotyczyło to m.in. pól takich jak PassengerId, Name czy Ticket, które były unikalne dla każdego pasażera i nie wносиły wartości predykcyjnej.

Zmieniliśmy także format wybranych zmiennych kategoriowych, stosując kodowanie zero-jedynkowe (one-hot encoding). Zmienną Embarked, która wskazywała miejsce zaokrętownia (C, Q, S), przekształciliśmy na trzy oddzielne kolumny binarne, w których 1 oznaczało obecność danego portu, a 0 jego brak. W podobny sposób zakodowaliśmy zmienną Sex, przypisując wartość 0 dla kobiet oraz 1 dla mężczyzn.

Zrezygnowaliśmy ze zmiennej Age, mimo jej potencjalnej przydatności, ze względu na dużą liczbę brakujących wartości. Spośród 891 rekordów aż ponad 100 nie zawierało informacji o wieku. Uzyskana w ten sposób luka w danych ograniczała jakość modelu oraz skutkowała zmniejszeniem zbioru treningowego, dlatego uznaliśmy, że bezpieczniej będzie całkowicie wykluczyć tę zmienną z dalszej analizy.

Poniżej przedstawiamy ostateczną wersję przekształconej struktury danych, która została wykorzystana w dalszych etapach projektu.

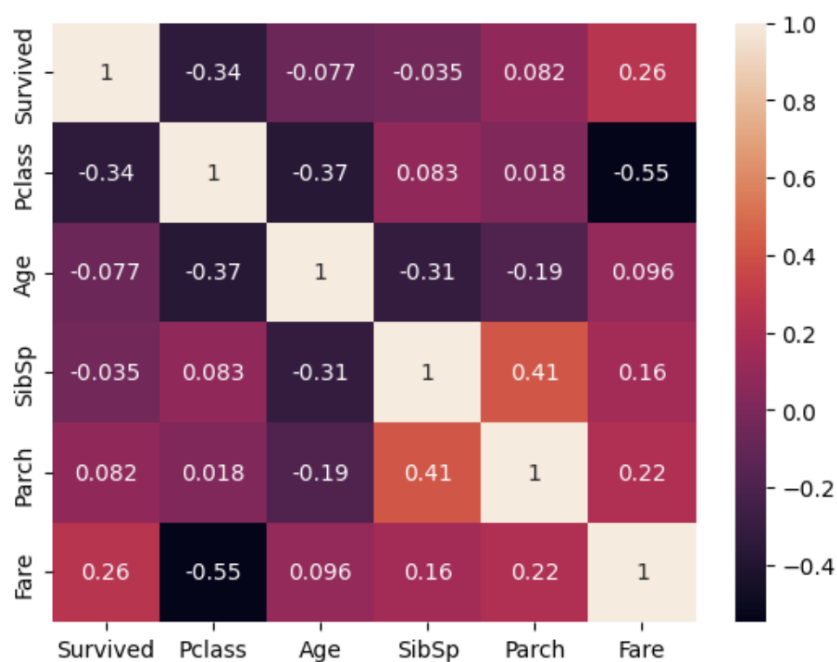
Tabela 2: użytkowa struktura danych

Kolumna	Typ danych	Opis
id	INT (PK)	Unikalny identyfikator pasażera
sex	ENUM('male','female')	Płeć
fare	DECIMAL(8,2)	Cena biletu
embarked	ENUM('C','Q','S')	Miejsce wejścia na pokład (Cherbourg, Queenstown, Southampton)
class	TINYINT	Klasa pasażerska (1, 2, 3)
survived	BOOLEAN	TRUE – przeżył, FALSE – zginął

Zależności pomiędzy zmiennymi a ich znaczenie względem zmiennej docelowej Survived zostały przeanalizowane z wykorzystaniem funkcji `seaborn.heatmap`, co przedstawiono na poniższej ilustracji.

```
In [5]: sns.heatmap(
df_train.drop(["PassengerId", "Name", "Ticket", "Cabin", "Embarked", "Sex"], axis=1).corr(),
annot=True,
)
```

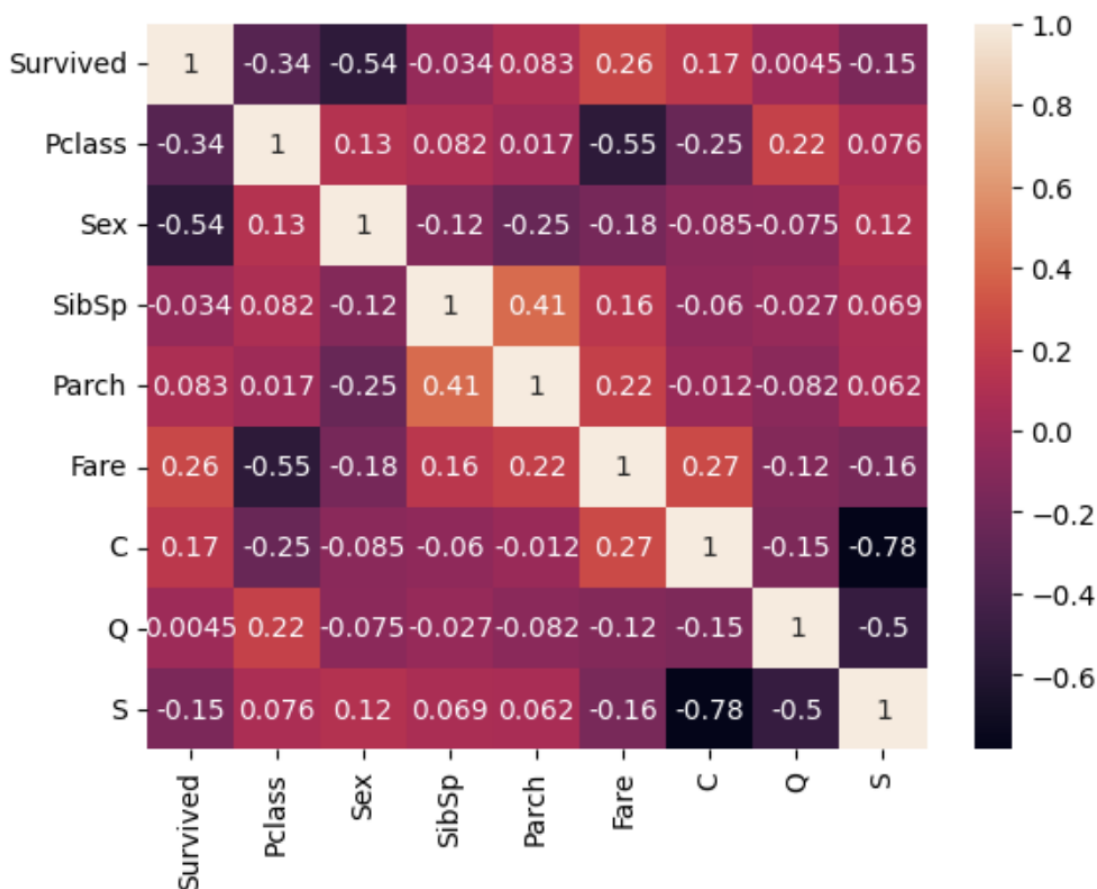
Out[5]: <Axes: >



Rysunek 4: Heatmap 1


```
In [24]: sns.heatmap(df_train.corr(), annot=True)
```

```
Out[24]: <Axes: >
```



Rysunek 5: Heatmap 2

Na podstawie powyższej mapy (heatmap) można zauważyć, że zmiennymi najsilniej skorelowanymi z prawdopodobieństwem przeżycia pasażera (Survived) są płeć (Sex, współczynnik korelacji ≈ -0.54) oraz klasa (Pclass, ≈ -0.34). Wartości ujemne oznaczają korelację odwrotną, co w tym przypadku oznacza, że niższe wartości zmiennej Sex (gdzie 0 oznacza kobiety, a 1 mężczyzn) są związane z wyższym prawdopodobieństwem przeżycia — innymi słowy, kobiety miały większe szanse na przetrwanie katastrofy.

Podobna zależność występuje w przypadku zmiennej Pclass, gdzie pasażerowie z klasy 1. (najwyższej) przeżywali częściej niż pasażerowie klasy 3. (najniższej).

Pozostałe zmienne wykazują znacznie słabszą korelację z przeżyciem. Warto wspomnieć o zmiennej Fare (cena biletu), która jest częściowo skorelowana z klasą (Pclass) – współczynnik korelacji wynosi około -0.55 – dlatego jej wpływ na Survived może być pośredni.

Zmienne takie jak PassengerId, Name, Ticket, Cabin oraz Embarked okazały się nieistotne w kontekście modelu predykcyjnego, głównie z powodu ich unikalności (np. Name, Ticket) lub dużej liczby braków (Cabin). Zmienna Age również nie wykazała silnej korelacji z przeżyciem, a dodatkowo zawierała wiele brakujących danych (891 rekordów ogółem, z czego 714 zawie-

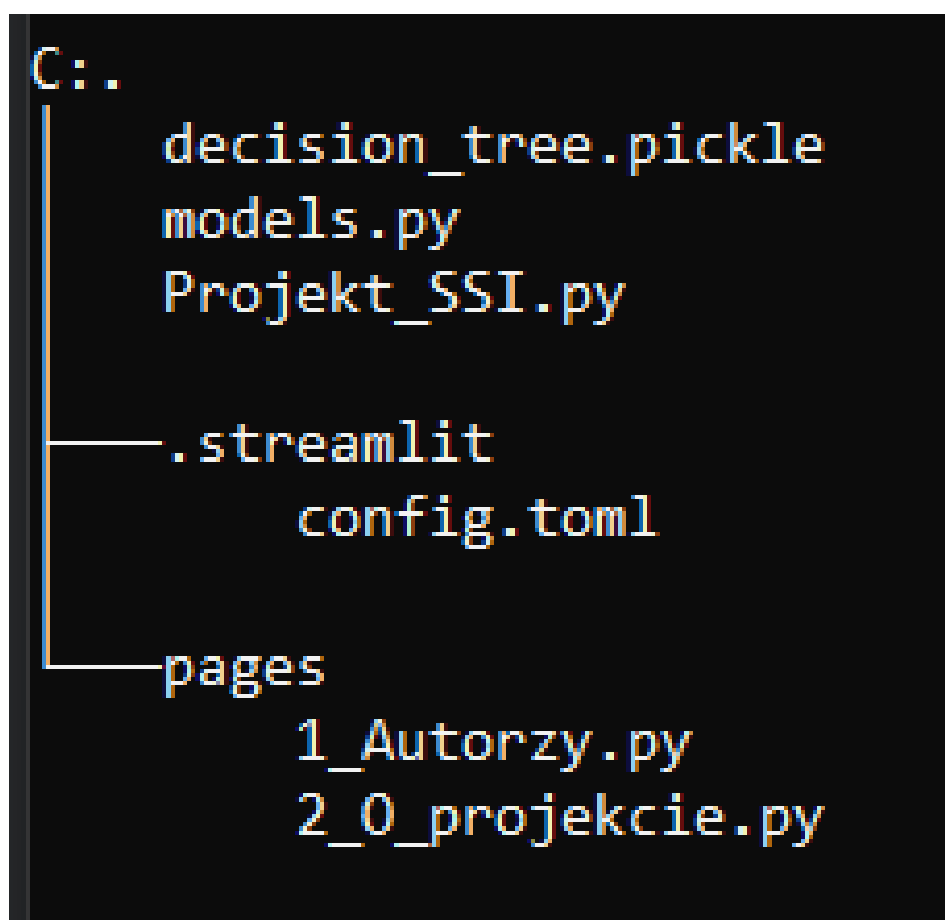
rało wartość Age, co oznacza brak danych w 177 przypadkach), co ograniczyło jej wartość analityczną.

```
: df_train.info()

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 891 entries, 0 to 890
Data columns (total 12 columns):
#   Column          Non-Null Count  Dtype
---  -
0   PassengerId     891 non-null    int64
1   Survived        891 non-null    int64
2   Pclass          891 non-null    int64
3   Name            891 non-null    object
4   Sex             891 non-null    object
5   Age             714 non-null    float64
6   SibSp           891 non-null    int64
7   Parch           891 non-null    int64
8   Ticket          891 non-null    object
9   Fare            891 non-null    float64
10  Cabin           204 non-null    object
11  Embarked        889 non-null    object
dtypes: float64(2), int64(5), object(5)
memory usage: 83.7+ KB
```

Rysunek 6: info

Implementacja



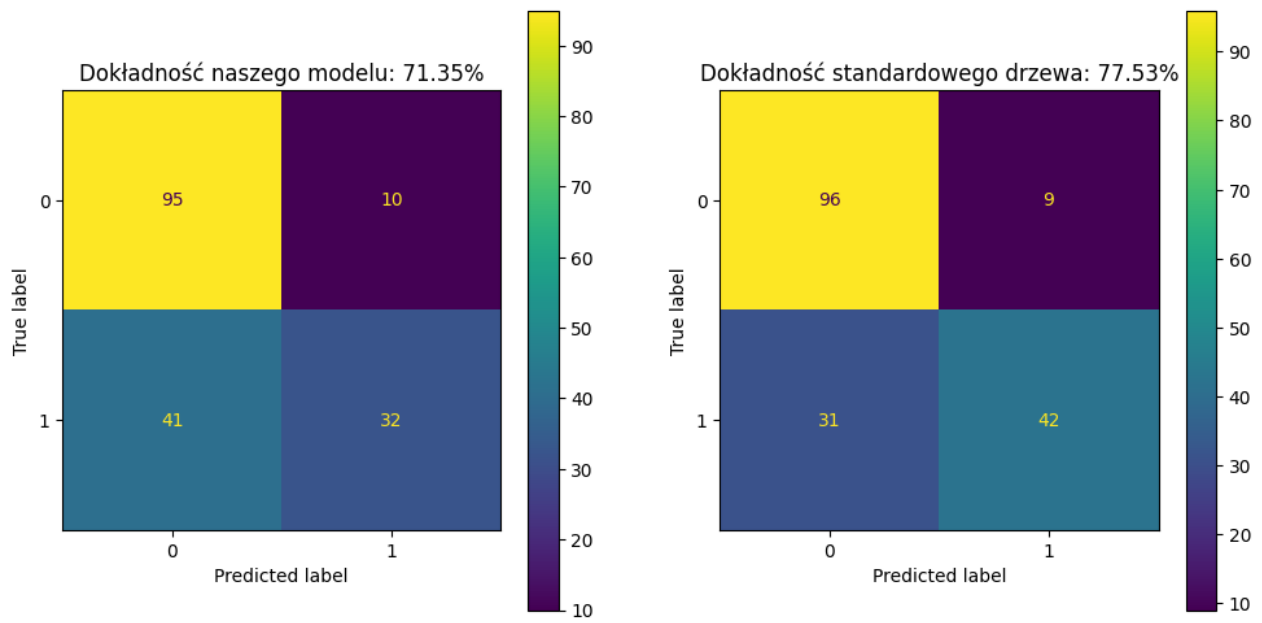
Rysunek 7: Struktura plików

Testy

Nasz model porównaliśmy ze standardowym drzewem decyzyjnym z biblioteki *scikit-learn* z tymi samymi parametrami tj. `max_dept` ustawiliśmy na 100, za miarę niejednorodności posłużył nam indeks Gini’ego, a wartość `min_samples_split` ustawiliśmy na 2.

Dokładność naszego modelu okazała się być niższa aż o 6 pkt. procentowych od dokładności standardowego drzewa decyzyjnego. Porównując macierze pomyłek obu modeli można zauważyć, że standardowe drzewo lepiej poradziło sobie z klasyfikacją osób, które katastrofę przeżyły. Osoby, które katastrofę przeżyły, nasz model częściej klasyfikował jako osoby, które owej katastrofy nie przeżyły. Jeżeli chodzi o klasyfikację osób, które katastrofy nie przeżyły - oba modele poradziły sobie równie dobrze.

Porównanie naszego modelu ze standardowym drzewem decyzyjnym



Rysunek 8: Porównanie macierzy pomyłek obu modeli

Pełen kod aplikacji

Pełen kod wraz z plikiem notatnika Jupyter, w którym trenowaliśmy model znajduje się również tutaj.

Frontend

```
1 import streamlit as st
2 import pickle as pk
3 import numpy as np
4
5 st.title("Projekt SSI")
6 with st.form(key='data_form'):
7     st.text('Wypełnij formularz i sprawdź czy by przeżył!')
8     sex = st.radio('Płeć:', ['M', 'K'])
9     cls = st.selectbox('Klasa:', ['1', '2', '3'])
10    sibsp = st.number_input('Liczba rodzeństwa + matki na pokładzie:', 0, 10)
11    parch = st.number_input('Liczba rodziców/dzieci na pokładzie:', 0, 10)
12    fare = st.number_input('Opłata za bilet:', 0.0, 512.0)
13    embarked = st.radio('Port startowy pasażera:', ['Cherbourg', 'QuinsTown', 'Southampton'])
14    submit = st.form_submit_button('Sprawdź!')
15
16 if submit:
17     data = np.array([[int(cls), sex=='M', sibsp, parch, fare, embarked=='Cherbourg', embarked=='QuinsTown', embarked=='Southampton']])
18     with open('decision_tree.pickle', 'rb') as f:
19         dt = pk.load(f)
20     surv = dt.predict(data)
21     if surv[0] == 1:
22         st.success('Przeżył!')
23         st.balloons()
24     else:
25         st.error('Nie przeżył!')
```

Plik "models.py", zawierający modele DecisionTree oraz SoftClassifier

```
1 import numpy as np
2 from typing import Optional
3 from dataclasses import dataclass
4
5
6 @dataclass
7 class Split:
8     feature_index: int
9     threshold: float
10    data_left: np.ndarray
11    data_right: np.ndarray
```

```

12     gain: float
13
14
15 @dataclass
16 class TreeNode:
17     feature_index: Optional[int] = None
18     threshold: Optional[float] = None
19     left: Optional["TreeNode"] = None
20     right: Optional["TreeNode"] = None
21     label: Optional[int] = None
22
23
24 class SoftClassifier:
25     """
26     A soft classifier that transforms features based on statistical
27     thresholds
28     and assigns labels by matching samples to the closest class profile.
29     """
30     def __init__(self, threshold=0.5):
31         self.threshold = threshold
32         self.feature_stats_ = None
33         self.classes_ = None
34
35     def fit(self, X, y):
36         """
37         Calculates mean and standard deviation of each feature per class
38         .
39         """
40         self.classes_ = np.unique(y)
41         self.feature_stats_ = {}
42
43         for class_ in self.classes_:
44             X_class = X[y == class_]
45             means = np.mean(X_class, axis=0)
46             stds = np.std(X_class, axis=0)
47             self.feature_stats_[class_] = {"mean": means, "std": stds}
48
49         return self
50
51     def _transform_features(self, X, class_):
52         """
53         Transforms features to binary indicators based on proximity to
54         class means.
55         """
56         means = self.feature_stats_[class_]["mean"]
57         stds = self.feature_stats_[class_]["std"]
58         transformed = (X >= (means - self.threshold * stds)).astype(int)
59         return transformed
60
61     def predict(self, X):
62         """
63         Predicts class labels based on soft matching of features to
64         class profiles.
65         """

```

```

63         scores = []
64         for class_ in self.classes_:
65             transformed = self._transform_features(X, class_)
66             score = np.sum(transformed, axis=1)
67             scores.append(score)
68
69         scores = np.array(scores).T
70         y_pred = self.classes_[np.argmax(scores, axis=1)]
71         return y_pred
72
73     def score(self, X, y):
74         """
75         Computes the accuracy of the soft classifier.
76         """
77         return np.sum(self.predict(X) == y) / len(X)
78
79
80 class DecisionTree:
81     """
82     A decision tree classifier with an optional stopping criterion based
83     on soft classification.
84     """
85     def __init__(
86         self,
87         criterion: str = "gini",
88         max_depth: int = 4,
89         min_samples_split: int = 2,
90         soft_classifier: Optional[SoftClassifier] = None,
91         soft_threshold: float = 0.9,
92     ) -> None:
93         """
94         Initializes the decision tree classifier.
95
96         Args:
97             criterion (str): The function to measure the quality of a
98                             split ('gini' or 'entropy').
99             max_depth (int): The maximum depth of the tree.
100             min_samples_split (int): The minimum number of samples
101                                    required to split a node.
102             soft_classifier (SoftClassifier): Optional soft classifier
103                                              to evaluate stopping condition.
104             soft_threshold (float): Minimum soft accuracy to stop
105                                    splitting.
106         """
107         assert criterion in [
108             "gini",
109             "entropy",
110         ], f"Invalid criterion: {criterion}. Choose 'gini' or 'entropy'."
111
112         self.criterion = criterion
113         self.max_depth = max_depth
114         self.min_samples_split = min_samples_split
115         self.soft_classifier = soft_classifier

```

```

112         self.soft_threshold = soft_threshold
113
114     def _entropy(self, probabilities: np.ndarray) -> float:
115         """
116         Computes the entropy of a set of probabilities.
117
118         Args:
119             probabilities (np.ndarray): Array of probabilities for each
120                                     class.
121
122         Returns:
123             float: The entropy value.
124         """
125         probabilities = probabilities[probabilities > 0]
126         return np.sum([-p * np.log2(p) for p in probabilities])
127
128     def _gini(self, probabilities: np.ndarray) -> float:
129         """
130         Computes the Gini impurity of a set of probabilities.
131
132         Args:
133             probabilities (np.ndarray): Array of probabilities for each
134                                     class.
135
136         Returns:
137             float: The Gini impurity value.
138         """
139         return 1 - np.sum([p**2 for p in probabilities])
140
141     def _get_probabilities(self, y: np.ndarray) -> np.ndarray:
142         """
143         Computes the probabilities for each class label in the dataset.
144
145         Args:
146             y (np.ndarray): Array of class labels.
147
148         Returns:
149             np.ndarray: Array of probabilities for each class.
150         """
151         labels = np.unique(y)
152         probabilities = np.empty(labels.shape)
153
154         for i, label in enumerate(labels):
155             probabilities[i] = len(y[y == label]) / len(y)
156
157         return probabilities
158
159     def _calculate_gain(
160         self, parent: np.ndarray, left_child: np.ndarray, right_child:
161         np.ndarray
162     ) -> float:
163         """
164         Calculates the information gain from a split.
165
166         Args:

```



```

164         parent (np.ndarray): The parent node's data.
165         left_child (np.ndarray): The left child node's data.
166         right_child (np.ndarray): The right child node's data.
167
168     Returns:
169         float: The calculated information gain.
170     """
171     weight_left = len(left_child) / len(parent)
172     weight_right = len(right_child) / len(parent)
173
174     parent_probabilities = self._get_probabilities(parent)
175     left_probabilities = self._get_probabilities(left_child)
176     right_probabilities = self._get_probabilities(right_child)
177
178     if self.criterion == "gini":
179         return self._gini(parent_probabilities) - (
180             weight_left * self._gini(left_probabilities)
181             + weight_right * self._gini(right_probabilities)
182         )
183     else:
184         return self._entropy(parent_probabilities) - (
185             weight_left * self._entropy(left_probabilities)
186             + weight_right * self._entropy(right_probabilities)
187         )
188
189     def _split(
190         self, data: np.ndarray, feature_index: int, threshold: float
191     ) -> tuple[np.ndarray, np.ndarray]:
192         """
193         Splits the dataset into two subsets based on the threshold of a
194         feature.
195
196         Args:
197             data (np.ndarray): The dataset to split.
198             feature_index (int): The index of the feature to split on.
199             threshold (float): The value used to split the data.
200
201         Returns:
202             Tuple[np.ndarray, np.ndarray]: The left and right splits of
203             the dataset.
204         """
205         mask = data[:, feature_index] < threshold
206         return data[mask, :], data[~mask, :]
207
208     def _get_best_split(self, data: np.ndarray) -> Split:
209         """
210         Finds the best split for the dataset based on the criterion.
211
212         Args:
213             data (np.ndarray): The dataset to find the best split for.
214
215         Returns:
216             Split: The best split for the dataset.
217         """
218         best_gain = float("-inf")

```

```

217         best_split = None
218
219     for feature_index in range(self.n_features_in_):
220         thresholds = np.unique(data[:, feature_index])
221
222         for threshold in thresholds:
223             data_left, data_right = self._split(data, feature_index,
224                                                 threshold)
225
226             if data_left.size > 0 and data_right.size > 0:
227                 y = data[:, -1]
228                 y_left = data_left[:, -1]
229                 y_right = data_right[:, -1]
230
231                 gain = self._calculate_gain(y, y_left, y_right)
232                 if gain > best_gain:
233                     best_split = Split(
234                         feature_index=feature_index,
235                         threshold=threshold,
236                         data_left=data_left,
237                         data_right=data_right,
238                         gain=gain,
239                     )
240
241     if best_split is None:
242         return Split(
243             feature_index=0,
244             threshold=0,
245             data_left=data,
246             data_right=data,
247             gain=0.0,
248         )
249     return best_split
250
251 def _get_most_frequent_label(self, y: np.ndarray) -> int:
252     """
253     Returns the most frequent label in the dataset.
254
255     Args:
256         y (np.ndarray): The array of labels.
257
258     Returns:
259         int: The most frequent label.
260     """
261     values, counts = np.unique(y, return_counts=True)
262     return values[np.argmax(counts)]
263
264 def _init_tree(self, data: np.ndarray, current_depth: int = 0) ->
265     TreeNode:
266     """
267     Recursively builds the decision tree.
268
269     Args:
270         data (np.ndarray): The dataset to build the tree from.
271         current_depth (int): The current depth of the tree.

```

```

270
271 Returns:
272     TreeNode: The root node of the tree.
273 """
274 n_samples = len(data)
275 y = data[:, -1]
276 X = data[:, :-1]
277
278 if n_samples >= self.min_samples_split and current_depth <= self
    .max_depth:
279     best_split = self._get_best_split(data)
280
281     if best_split.gain > 0:
282         if self.soft_classifier:
283             self.soft_classifier.fit(X, y)
284             if self.soft_classifier.score(X, y) >= self.
                soft_threshold:
285                 return TreeNode(label=self.
                    _get_most_frequent_label(y))
286
287         left = self._init_tree(best_split.data_left,
            current_depth + 1)
288         right = self._init_tree(best_split.data_right,
            current_depth + 1)
289
290         return TreeNode(
291             feature_index=best_split.feature_index,
292             threshold=best_split.threshold,
293             left=left,
294             right=right,
295         )
296
297     return TreeNode(label=self._get_most_frequent_label(y))
298
299 def fit(self, X: np.ndarray, y: np.ndarray) -> None:
300     """
301     Trains the decision tree on the given data.
302
303     Args:
304         X (np.ndarray): Feature matrix (n_samples, n_features).
305         y (np.ndarray): Target labels (n_samples,).
306     """
307     self.n_features_in_ = X.shape[1]
308     data = np.concatenate((X, y.reshape((-1, 1))), axis=1)
309     self.tree_ = self._init_tree(data)
310
311 def _predict(self, x: np.ndarray, node: TreeNode) -> int:
312     """
313     Recursively predicts the label for a given sample.
314
315     Args:
316         x (np.ndarray): The input data point.
317         node (TreeNode): The current node in the tree.
318
319     Returns:

```

```

320         int: The predicted label.
321     """
322     if node.label is not None:
323         return node.label
324
325     if x[node.feature_index] < node.threshold:
326         return self._predict(x, node.left)
327     else:
328         return self._predict(x, node.right)
329
330 def predict(self, X: np.ndarray) -> np.ndarray:
331     """
332     Predicts class labels for the input data.
333
334     Args:
335         X (np.ndarray): Feature matrix (n_samples, n_features).
336
337     Returns:
338         np.ndarray: Predicted class labels (n_samples,).
339     """
340     return np.array([self._predict(x, self.tree_) for x in X])
341
342 def score(self, X: np.ndarray, y: np.ndarray) -> float:
343     """
344     Computes the accuracy of the classifier.
345
346     Args:
347         X (np.ndarray): Feature matrix (n_samples, n_features).
348         y (np.ndarray): True class labels (n_samples,).
349
350     Returns:
351         float: Accuracy of the classifier.
352     """
353     return np.sum(self.predict(X) == y) / len(X)

```
