

# **Automatyczna Klasyfikacja Tematyczna Polskich Artykulow Newsowych**

Maciej Biegan      Kamil Dziedzic      Jan Bobak

Eksploracja Danych Tekstowych – Grudzien 2025

# Spis treści

<b>1</b>	<b>Wstęp</b>	<b>4</b>
1.1	Opis Problemu . . . . .	4
1.2	Wyzwania dla Języka Polskiego . . . . .	4
1.3	Cel Projektu . . . . .	4
<b>2</b>	<b>Metody Przetwarzania Języka Naturalnego</b>	<b>4</b>
2.1	Pipeline Przetwarzania . . . . .	4
2.2	Normalizacja Tekstu . . . . .	5
2.3	Lematyzacja z spaCy . . . . .	5
2.4	Wektoryzacja TF-IDF . . . . .	6
2.5	Embeddingi Kontekstowe (Transformery) . . . . .	6
<b>3</b>	<b>Zródła Danych</b>	<b>7</b>
3.1	Kanale RSS . . . . .	7
3.2	Architektura Scrapera . . . . .	7
3.3	Ekstrakcja Pełnej Treści . . . . .	8
3.4	Filtracja Jakościowa . . . . .	8
3.5	Charakterystyka Datasetu . . . . .	9
<b>4</b>	<b>Modele</b>	<b>9</b>
4.1	MLP . . . . .	9
4.1.1	Architektura . . . . .	10
4.1.2	Komponenty . . . . .	10
4.1.3	Parametry Konfiguracyjne . . . . .	11
4.1.4	Trening . . . . .	11
4.2	HerBERT . . . . .	12
4.3	BERT Multilingual . . . . .	12
4.4	Konfiguracja Transformerów . . . . .	12
4.5	Fine-tuning . . . . .	13
<b>5</b>	<b>Aplikacja Webowa</b>	<b>13</b>
5.1	Architektura . . . . .	13
5.2	Endpointy . . . . .	14
5.3	Obsługa Wejścia . . . . .	14
5.4	Asynchroniczne Trenowanie . . . . .	15
5.5	Interfejs Użytkownika . . . . .	15
<b>6</b>	<b>Wyniki</b>	<b>17</b>
6.1	Metodologia . . . . .	17
6.2	Porównanie Modeli . . . . .	18
6.3	Analiza Hiperparametrów MLP . . . . .	18
6.3.1	Testy ANOVA . . . . .	19
6.3.2	Macierz Korelacji . . . . .	19
6.4	Najlepszy Model MLP . . . . .	20
6.5	Macierz Pomyłek . . . . .	20
6.6	Metryki Per Kategoria . . . . .	21
6.7	Wpływ Architektury . . . . .	22

<b>7</b>	<b>Podsumowanie</b>	<b>23</b>
7.1	Osiagniecia . . . . .	23
7.2	Wnioski Techniczne . . . . .	23

# 1 Wstep

## 1.1 Opis Problemu

Automatyczne tagowanie artykulow (klasyfikacja tematyczna) to zadanie przypisania dokumentu tekstowego do jednej z predefiniowanych kategorii na podstawie jego tresci. Problem nalezy do domeny klasyfikacji wieloklasowej tekstu w przetwarzaniu jezyka naturalnego (NLP).

**Formalnie:** Dla dokumentu  $d$  i zbioru kategorii  $C = \{c_1, c_2, \dots, c_k\}$  szukamy funkcji  $f : D \rightarrow C$  maksymalizujacej prawdopodobienstwo poprawnej klasyfikacji.

W niniejszym projekcie  $k = 6$  kategorii:

- **Polska** – polityka krajowa, sprawy wewnetrzne
- **Swiat** – wydarzenia miedzynarodowe
- **Biznes** – gospodarka, finanse, rynki
- **Technologie** – IT, innowacje, startupy
- **Moto** – motoryzacja, samochody
- **Sport** – wydarzenia sportowe, wyniki

## 1.2 Wyzwania dla Jezyka Polskiego

Klasyfikacja polskich tekstow jest trudniejsza niz angielskich ze wzgledu na kilka czynnikow. Po pierwsze, skomplikowana budowa jezyka polskiego, z siedmioma przypadkami, co utrudnia proces lematyzacji. Po drugie, polski ma swobodny szyk zdania, gdzie kolejnosc slow nie zawsze determinuje ich funkcje, co utrudnia analize skladniowa. Po trzecie, brak gotowych danych do trenowania modeli.

## 1.3 Cel Projektu

Celem projektu bylo opracowanie systemu pozwalajacego na:

1. Automatyczny scraping artykulow z kanalow RSS
2. Pipeline preprocessingu NLP (normalizacja, lematyzacja, wektoryzacja)
3. Trening i ewaluacja trzech architektur modeli
4. Aplikacja webowa do predykcji i trenowania

# 2 Metody Przetwarzania Jezyka Naturalnego

## 2.1 Pipeline Przetwarzania

Tekst przechodzi sekwencje transformacji:

$$\text{RAW} \xrightarrow{\text{normalize}} \text{CLEAN} \xrightarrow{\text{tokenize}} \text{TOKENS} \xrightarrow{\text{lemmatize}} \text{LEMMAS} \xrightarrow{\text{vectorize}} \mathbb{R}^n$$

## 2.2 Normalizacja Tekstu

Funkcja `normalize_text()` w `Scrapper/scrapper.py` wykonuje:

Listing 1: Implementacja normalizacji tekstu

```
1 import re
2
3 RE_URL = re.compile(r'https?://\S+|www\.\S+')
4 RE_NON_LETTER = re.compile(r'^a-zA-ZacelnoszzACELNOSZZ\s-')
5 RE_MULTI_WS = re.compile(r'\s+')
6
7 def normalize_text(text: str) -> str:
8     if not isinstance(text, str):
9         return ''
10    text = text.strip()
11    text = RE_URL.sub(' ', text) # usuwa URLs
12    text = text.replace('\xa0', ' ') # non-breaking space
13    text = RE_NON_LETTER.sub(' ', text) # tylko litery
14    text = RE_MULTI_WS.sub(' ', text) # normalizacja spacji
15    text = text.lower()
16    return text
```

### Operacje:

- `RE_URL.sub(' ', text)` – usuwa linki HTTP/HTTPS (brak informacji semantycznej)
- `RE_NON_LETTER.sub(' ', text)` – zachowuje tylko litery ASCII i polskie diakrytyki
- `text.lower()` – unifikacja wielkości liter

## 2.3 Lematyzacja z spaCy

Lematyzacja redukuje formy fleksyjne do formy podstawowej. Model `pl_core_news_sm` (15MB) trenowany na NKJP.

Listing 2: Funkcja lematyzacji i tokenizacji

```
1 import spacy
2 from stop_words import get_stop_words
3
4 def get_polish_stopwords() -> Set[str]:
5     nlp = spacy.load('pl_core_news_sm')
6     spacy_stopwords = set(nlp.Defaults.stop_words)
7     base_stopwords = set(get_stop_words('polish'))
8     extra_stops = {'z', 'na', 'i', 'w', 'o', 'ze', 'sie',
9                   'roku', 'r', 'godz', 'fot', 'foto'}
10    return base_stopwords | spacy_stopwords | extra_stops
11
12 def lemmatize_and_tokenize(text, nlp_model, stopwords,
13                             min_token_len=3):
14     text_norm = normalize_text(text)
15     doc = nlp_model(text_norm)
16     tokens = []
17     for tok in doc:
18         if tok.is_space or tok.is_punct:
19             continue
```

```

20     lemma = tok.lemma_.lower().strip()
21     if len(lemma) < min_token_len:
22         continue
23     if lemma in stopwords or lemma.isdigit():
24         continue
25     tokens.append(lemma)
26     return tokens

```

### Przykład transformacji:

Input: "Prezydent podpisał ustawę o reformie szpitali"

Output: ["prezydent", "podpisac", "ustawa", "reforma", "szpital"]

Redukcja: 7 słów → 5 lematów (28% kompresji).

## 2.4 Wektoryzacja TF-IDF

Term Frequency-Inverse Document Frequency przekształca tekst w wektor numeryczny.

**Definicja:**

$$\text{TF-IDF}(t, d, D) = \text{TF}(t, d) \cdot \text{IDF}(t, D) \quad (1)$$

Listing 3: Konfiguracja TfidfVectorizer

```

1 from sklearn.feature_extraction.text import TfidfVectorizer
2
3 vectorizer = TfidfVectorizer(
4     max_features=5000,      # limit wymiarowosci
5     ngram_range=(1, 2),    # unigramy + bigramy
6     sublinear_tf=True,     # logarytmiczne TF
7     min_df=2,              # min 2 dokumenty
8     max_df=0.95            # max 95% dokumentow
9 )
10 X = vectorizer.fit_transform(df['text']).toarray()

```

**N-gramy:** Bigramy chwytają frazy: „liga\_mistrzow”, „premier\_polski”, „gielda\_papierow”.

## 2.5 Embeddingi Kontekstowe (Transformery)

W przeciwieństwie do TF-IDF (bag-of-words), transformery generują embeddingi zależne od kontekstu.

**Self-Attention:**

$$\text{Attention}(Q, K, V) = \text{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right)V \quad (2)$$

gdzie  $Q$ ,  $K$ ,  $V$  to projekcje liniowe wejścia,  $d_k$  – wymiar klucza.

**Tokenizacja WordPiece:**

Listing 4: Tokenizacja dla HerBERT

```

1 from transformers import AutoTokenizer
2
3 tokenizer = AutoTokenizer.from_pretrained(
4     'allegro/herbert-base-cased'
5 )
6

```

```

7 def tokenize_fn(examples):
8     return tokenizer(
9         examples["text"],
10        padding="max_length",
11        truncation=True,
12        max_length=256
13    )

```

WordPiece dzieli nieznane słowa na subwordy:

„niesamowity” → [„nie”, „##sam”, „##owity”]

## 3 Zrodla Danych

### 3.1 Kanaly RSS

Dane pochodza z 26 kanalow RSS polskich portali informacyjnych:

Tabela 1: Zrodla RSS per kategoria

Kategoria	Portal	Domena
Polska	Polsat News, TVN24, WP, Fakt, Do Rzeczy, Wprost	6 zrodel
Swiat	Polsat News, TVN24, RMF24, Newsweek	4 zrodla
Biznes	Polsat News, TVN24 Biznes, Puls Biznesu	3 zrodla
Technologie	Polsat News, TVN24, Computerworld, Spider’s Web	4 zrodla
Moto	Polsat News, TVN24 Moto	2 zrodla
Sport	Polsat News, TVN24 Sport, Onet Sport	3 zrodla

### 3.2 Architektura Scrapera

Implementacja w Scraper/scrapper.py:

Listing 5: Slownik zrodel RSS

```

1 FEEDS = {
2     'Polska': [
3         'https://www.polsatnews.pl/rss/polska.xml',
4         'https://tvn24.pl/polska.xml',
5         'https://wiadomosci.wp.pl/rss.xml',
6         'https://www.fakt.pl/rss',
7         'https://dorzeczy.pl/rss',
8         'https://www.wprost.pl/rss'
9     ],
10    'Swiat': [
11        'https://www.polsatnews.pl/rss/swiat.xml',
12        'https://tvn24.pl/swiat.xml',
13        'https://www.rmf24.pl/feed/',
14        'https://www.newsweek.pl/rss'
15    ],
16    # ... pozostale kategorie
17 }

```

### 3.3 Ekstrakcja Pełnej Treści

RSS zawiera tylko tytuł i streszczenie. Pełna treść wymaga pobrania HTML:

Listing 6: Ekstrakcja tekstu z URL

```
1 from newspaper import Article
2 from bs4 import BeautifulSoup
3 import requests
4
5 def extract_full_text(url: str, timeout: int = 10) -> str:
6     headers = {'User-Agent': 'Mozilla/5.0'}
7
8     # Metoda 1: newspaper3k
9     try:
10         article = Article(url)
11         article.download()
12         article.parse()
13         if article.text and len(article.text) > 100:
14             return article.text
15     except Exception:
16         pass
17
18     # Metoda 2: BeautifulSoup fallback
19     try:
20         resp = requests.get(url, headers=headers, timeout=timeout)
21         soup = BeautifulSoup(resp.content, 'lxml')
22
23         selectors = ["article", "div[class*='article']",
24                     "div[class*='content']"]
25         for sel in selectors:
26             el = soup.select_one(sel)
27             if el:
28                 ps = el.find_all('p')
29                 text = '\n'.join([p.get_text(strip=True)
30                                 for p in ps])
31                 if len(text) > 100:
32                     return text
33     except Exception:
34         pass
35
36     return None
```

### 3.4 Filtracja Jakościowa

Listing 7: Kryteria filtracji w funkcji fetch\_articles\_from\_feeds

```
1 def fetch_articles_from_feeds(feeds, min_length=200):
2     # ...
3     for entry in category_entries:
4         text = extract_full_text(link)
5         text = clean_text(text)
6
7         # Filtr 1: minimalna długość
8         if len(text) < min_length:
9             continue
10
11         # Filtr 2: wykrywanie języka
```



```

12         if detect(text) != 'pl':
13             continue
14
15         records.append({...})
16
17     # Filtrowanie 3: deduplikacja
18     df.drop_duplicates(subset=['url', 'text'], inplace=True)
19     return df

```

## 3.5 Charakterystyka Datasetu

Wynikowy plik: Scrapper/polsatnews\_articles\_clean.csv

Tabela 2: Rozkład artykułow per kategoria

Kategoria	Liczba	Udział [%]	Zrodla
Polska	934	26.8	6
Sport	762	21.9	3
Swiat	708	20.3	4
Biznes	463	13.3	3
Technologie	389	11.2	4
Moto	231	6.6	2
<b>Razem</b>	<b>3487</b>	<b>100</b>	<b>26</b>

**Niezbalansowanie:** Polska (26.8%) vs Moto (6.6%) – stosunek 4:1. Zastosowano `class_weight='balanced'`.

## 4 Modele

W projekcie zaimplementowano i porównano trzy różne architektury modeli uczenia maszynowego do klasyfikacji tekstu. Każdy z modeli wykorzystuje inną reprezentację tekstu oraz inny mechanizm uczenia. Model MLP operuje na wektorach TF-IDF i stanowi klasyczne podejście oparte na płytkim uczeniu. Modele transformerowe HerBERT i BERT wykorzystują mechanizm self-attention i pretrenowane reprezentacje językowe, co pozwala na lepsze zrozumienie kontekstu i semantyki tekstu. Wszystkie modele zostały zaimplementowane w PyTorch z wykorzystaniem biblioteki Hugging Face Transformers dla modeli BERT.

### 4.1 MLP

Model MLP jest wielowarstwowa sieć neuronowa typu feed-forward, która przyjmuje na wejściu wektor TF-IDF o wymiarze 5000 i klasyfikuje go do jednej z 6 kategorii. Architektura sieci jest w pełni konfigurowalna i pozwala na zmianę liczby warstw ukrytych, rozmiaru warstw oraz współczynnika dropout.

### 4.1.1 Architektura

Klasa `TextClassifier` dziedziczy po `nn.Module` i buduje siec sekwencyjnie. Pierwsza warstwa liniowa redukuje wymiarowosc z 5000 cech TF-IDF do rozmiaru warstwy ukrytej (domyslnie 256). Nastepnie stosowana jest normalizacja `LayerNorm`, funkcja aktywacji `GELU` oraz regularyzacja `Dropout`. Ten blok powtarzany jest dla kazdej warstwy ukrytej. Ostatnia warstwa liniowa mapuje reprezentacje na 6 klas wyjsciowych.

Implementacja w `models/model.py`:

Listing 8: Klasa `TextClassifier`

```
1 import torch
2 import torch.nn as nn
3
4 class TextClassifier(nn.Module):
5     def __init__(self, input_size: int, hidden_size: int,
6                   num_classes: int, num_layers: int = 3,
7                   dropout: float = 0.4):
8         super().__init__()
9
10        layers = [
11            nn.Linear(input_size, hidden_size),
12            nn.LayerNorm(hidden_size),
13            nn.GELU(),
14            nn.Dropout(dropout),
15        ]
16
17        for _ in range(num_layers - 1):
18            layers.extend([
19                nn.Linear(hidden_size, hidden_size),
20                nn.LayerNorm(hidden_size),
21                nn.GELU(),
22                nn.Dropout(dropout),
23            ])
24
25        layers.append(nn.Linear(hidden_size, num_classes))
26        self.network = nn.Sequential(*layers)
27
28        def forward(self, x: torch.Tensor) -> torch.Tensor:
29            return self.network(x)
```

### 4.1.2 Komponenty

Architektura MLP sklada sie z kilku kluczowych komponentow, ktore wspolpracuja w celu efektywnego przetwarzania i klasyfikacji tekstu.

**LayerNorm** normalizuje aktywacje po wymiarze cech, a nie po wymiarze batcha jak `BatchNorm`. Dla kazdego przykladu batchu obliczana jest srednia i odchylenie standardowe po wszystkich cechach, a nastepnie aktywacje sa standaryzowane. Parametry `gamma` i `beta` sa uczone podczas treningu i pozwalaja sieci na skalowanie i przesuwanie znormalizowanych wartosci. `LayerNorm` jest szczegolnie przydatna w NLP, poniewaz dziala stabilnie niezaleznie od rozmiaru batcha.

$$\text{LayerNorm}(x) = \gamma \cdot \frac{x - \mu}{\sigma + \epsilon} + \beta \quad (3)$$

**GELU** jest funkcja aktywacji, ktora mnozy wejscie przez wartosc dystrybuanty rozkladu normalnego. W przeciwienstwie do `ReLU`, ktora ostro obcina wartosci ujemne do

zera, GELU zapewnia gładkie przejście w okolicy zera. Dzięki temu gradienty są lepiej propagowane przez sieć, co jest szczególnie korzystne w zadaniach NLP. GELU jest standardowa funkcja aktywacji w modelach BERT i GPT.

$$\text{GELU}(x) = x \cdot \Phi(x) = x \cdot \frac{1}{2} \left[ 1 + \text{erf} \left( \frac{x}{\sqrt{2}} \right) \right] \quad (4)$$

**Dropout** jest technika regularyzacji, która podczas treningu losowo zeruje część neuronów z prawdopodobieństwem  $p$ . Zapobiega to nadmiernemu dopasowaniu modelu do danych treningowych (overfitting), ponieważ sieć nie może polegać na żadnym pojedynczym neuronie. Podczas inferencji dropout jest wyłączany, a wagi są skalowane przez  $(1-p)$  aby zachować oczekiwaną wartość aktywacji.

$$\text{Dropout}(x) = \begin{cases} 0 & \text{z prawdop. } p \\ \frac{x}{1-p} & \text{wpp.} \end{cases} \quad (5)$$

### 4.1.3 Parametry Konfiguracyjne

Definicja w `models/config.py`:

Listing 9: MLPConfig dataclass

```
1 @dataclass
2 class MLPConfig:
3     epochs: int = 20
4     batch_size: int = 32
5     learning_rate: float = 0.001
6     max_features: int = 5000
7     hidden_size: int = 256
8     num_layers: int = 3
9     dropout: float = 0.4
10    early_stopping: bool = False
11    early_stopping_patience: int = 3
```

### 4.1.4 Trening

Listing 10: Pętla treningowa MLP (fragment `_train_mlp`)

```
1 from sklearn.utils.class_weight import compute_class_weight
2
3 # Wagi klas dla niezbalansowanego datasetu
4 class_weights = torch.tensor(
5     compute_class_weight('balanced',
6                           classes=np.unique(y_train),
7                           y=y_train),
8     dtype=torch.float32
9 ).to(DEVICE)
10
11 criterion = nn.CrossEntropyLoss(weight=class_weights)
12 optimizer = torch.optim.AdamW(model.parameters(),
13                                lr=config.learning_rate)
14
15 for epoch in range(config.epochs):
16     model.train()
17     for X_batch, y_batch in train_loader:
```

```

18     optimizer.zero_grad()
19     outputs = model(X_batch.to(DEVICE))
20     loss = criterion(outputs, y_batch.to(DEVICE))
21     loss.backward()
22     optimizer.step()

```

## 4.2 HerBERT

Polski transformer BERT pretrenowany przez Allegro.pl.

Tabela 3: Specyfikacja HerBERT

Parametr	Wartosc
Model	allegro/herbert-base-cased
Warstwy Transformer	12
Wymiar ukryty	768
Attention heads	12
Parametry	110M
Korpus treningowy	14GB polskiego tekstu
Tokenizer	SentencePiece (BPE)
Vocab size	50,000

## 4.3 BERT Multilingual

Model bert-base-multilingual-cased trenowany na 104 jezykach:

- Vocab size: 119,547 tokenow
- Polski: ok. 3% korpusu
- Gorsze wyniki niz HerBERT (brak specjalizacji)

## 4.4 Konfiguracja Transformerow

Listing 11: TransformerConfig dataclass

```

1  @dataclass
2  class TransformerConfig:
3      epochs: int = 10
4      batch_size: int = 16
5      eval_batch_size: int = 64
6      warmup_steps: int = 500
7      weight_decay: float = 0.01
8      learning_rate: float = 2e-5
9      max_length: int = 256
10     early_stopping: bool = False
11     early_stopping_patience: int = 3

```

## 4.5 Fine-tuning

Listing 12: Fine-tuning transformera (fragment `_train_transformer`)

```
1 from transformers import (AutoModelForSequenceClassification,
2                             Trainer, TrainingArguments)
3
4 model = AutoModelForSequenceClassification.from_pretrained(
5     MODEL_NAMES[model_type],
6     num_labels=len(label_encoder.classes_)
7 ).to(DEVICE)
8
9 training_args = TrainingArguments(
10     output_dir=str(PATHS.results_dir),
11     num_train_epochs=config.epochs,
12     per_device_train_batch_size=config.batch_size,
13     learning_rate=config.learning_rate,
14     warmup_steps=config.warmup_steps,
15     weight_decay=config.weight_decay,
16     eval_strategy='epoch',
17     save_strategy='epoch',
18     load_best_model_at_end=True,
19     metric_for_best_model='accuracy',
20     fp16=torch.cuda.is_available(),
21 )
22
23 trainer = Trainer(
24     model=model,
25     args=training_args,
26     train_dataset=train_dataset,
27     eval_dataset=test_dataset,
28     compute_metrics=compute_metrics,
29 )
30 trainer.train()
```

## 5 Aplikacja Webowa

System klasyfikacji został wdrożony jako aplikacja webowa, która umożliwia użytkownikom klasyfikowanie artykułów newsowych w czasie rzeczywistym oraz trenowanie własnych modeli. Aplikacja została zbudowana z wykorzystaniem nowoczesnych technologii webowych i zaprojektowana z myślą o skalowalności i łatwości użytkowania.

### 5.1 Architektura

Aplikacja została zaimplementowana w języku Python z wykorzystaniem frameworka Flask, który jest lekkim frameworkiem webowym opartym na bibliotece Werkzeug i silniku szablonów Jinja2. Architektura aplikacji opiera się na wzorcu MVC, gdzie modele uczenia maszynowego stanowią warstwę Model, szablony HTML stanowią warstwę View, a funkcje widoków Flask stanowią warstwę Controller.

Główny plik aplikacji `app.py` definiuje endpointy HTTP i zarządza logiką biznesową. Modele ML są enkapsulowane w klasie `ModelManager`, która zapewnia thread-safe dostęp do modeli z wielu wątków. Szablony HTML wykorzystują Bootstrap do stylowania i JavaScript do dynamicznej aktualizacji interfejsu.

Listing 13: Klasa ModelManager

```

1 class ModelManager:
2     """Thread-safe manager dla modeli ML"""
3
4     def __init__(self, default_model_type="herbert"):
5         self.model_type = default_model_type
6         self.model = None
7         self.tokenizer_or_vectorizer = None
8         self.label_encoder = None
9         self._lock = threading.Lock()
10        self.load()
11
12    def load(self, model_type=None):
13        with self._lock:
14            if model_type:
15                self.model_type = model_type
16            (self.model,
17             self.tokenizer_or_vectorizer,
18             self.label_encoder) = load_model(self.model_type)
19
20    def predict(self, text: str):
21        if self.model is None:
22            raise RuntimeError("Model not loaded")
23        with self._lock:
24            return predict_category(
25                text, self.model,
26                self.tokenizer_or_vectorizer,
27                self.label_encoder,
28                self.model_type
29            )

```

## 5.2 Endpointy

Tabela 4: REST API

Endpoint	Metody	Opis
/	GET, POST	Predykcja kategorii
/train	GET, POST	Formularz/start treningu
/training_status	GET	Status treningu (JSON)
/training_result	GET	Wyniki treningu (JSON)

## 5.3 Obsługa Wejscia

Trzy źródła tekstu: bezpośredni input, upload pliku, URL artykułu.

Listing 14: Ekstrakcja tekstu z requestu

```

1 def extract_text_from_request(form, files):
2     # Opcja 1: tekst bezpośredni
3     if form.get("text", "").strip():
4         return form["text"].strip()
5
6     # Opcja 2: upload pliku .txt

```

```

7     if "file" in files and files["file"].filename:
8         return files["file"].read().decode("utf-8")
9
10    # Opcja 3: URL artykułu
11    if form.get("url", "").strip():
12        return extract_text_from_url(form["url"].strip())
13
14    return None

```

## 5.4 Asynchroniczne Trenowanie

Trening w osobnym watku, aby nie blokować UI:

Listing 15: Trenowanie w tle

```

1  def _background_train(params):
2      try:
3          train_model_with_params(params)
4          model_manager.load(params.get("model_type"))
5
6          with model_module._status_lock:
7              model_module.training_status["message"] = "completed"
8      except Exception as e:
9          with model_module._status_lock:
10             model_module.training_status.update(
11                 {"running": False, "error": str(e)}
12             )
13      finally:
14          with model_module._status_lock:
15             model_module.training_status["running"] = False
16
17  @app.route('/train', methods=['POST'])
18  def train():
19      # ... parsowanie parametrow ...
20      thread = threading.Thread(
21          target=_background_train,
22          args=(params,),
23          daemon=True
24      )
25      thread.start()
26      return render_template('train.html', training=True)

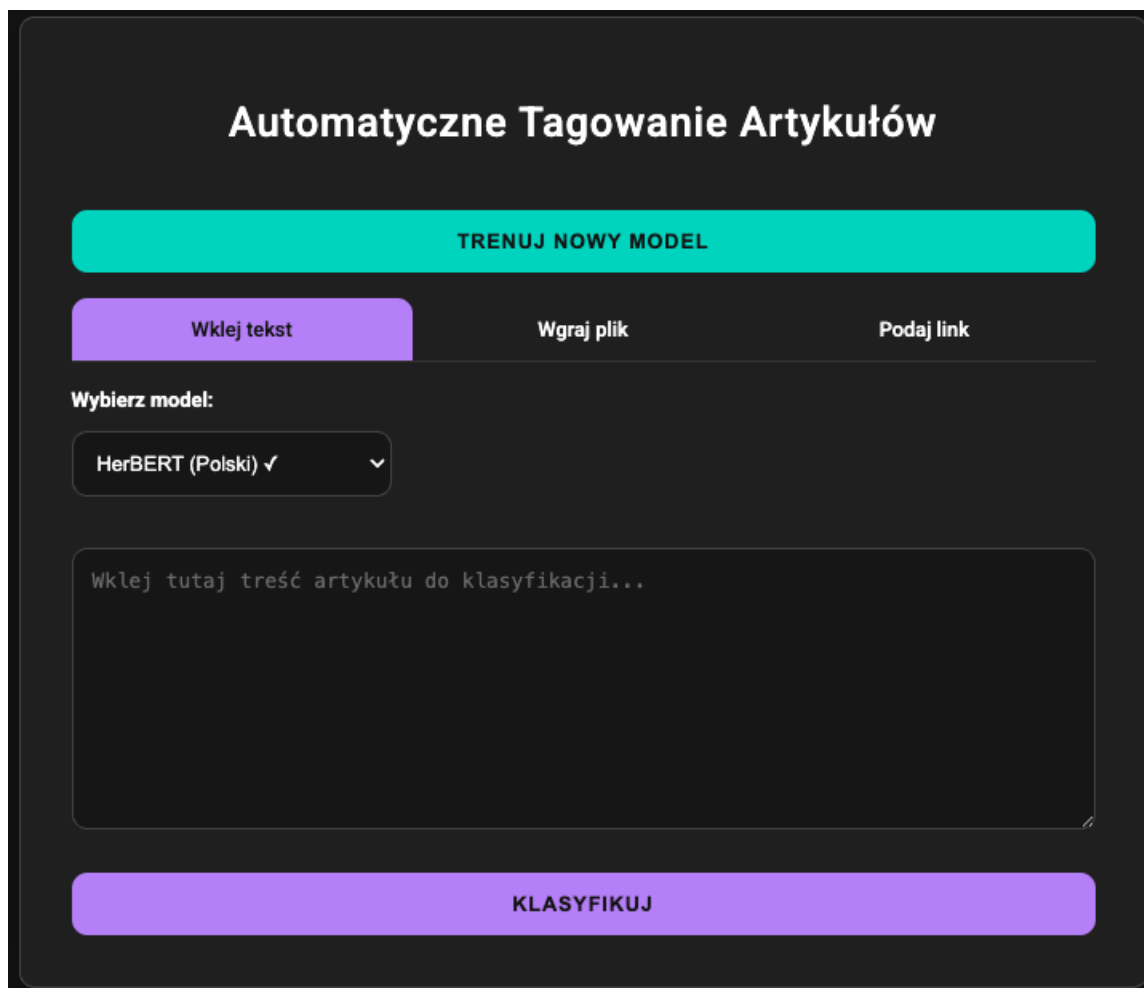
```

## 5.5 Interfejs Użytkownika

Interfejs użytkownika został zaprojektowany z myślą o prostocie i intuicyjności. Aplikacja składa się z dwóch głównych widoków: strony głównej do klasyfikacji tekstu oraz panelu trenowania modeli.

Strona główna umożliwia użytkownikowi wprowadzenie tekstu do klasyfikacji na trzy sposoby. Pierwszy sposób to bezpośrednie wpisanie lub wklejenie tekstu artykułu w pole tekstowe. Drugi sposób to upload pliku tekstowego w formacie UTF-8. Trzeci sposób to podanie URL artykułu, z którego tekst zostanie automatycznie wyodrębniony za pomocą biblioteki newspaper3k. Użytkownik może również wybrać model do klasyfikacji spośród dostępnych: HerBERT, BERT multilingual lub MLP. Po wysłaniu formularza wyświetlana jest przewidziana kategoria wraz z poziomem pewności modelu wyrażonym w procentach.

Panel trenowania pozwala na konfigurację i uruchomienie treningu nowego modelu. Użytkownik może wybrać typ modelu oraz dostosować hiperparametry takie jak liczba epok, learning rate, batch size, dropout oraz maksymalna liczba cech TF-IDF dla modelu MLP. Podczas treningu wyświetlany jest pasek postępu oraz bieżące metryki: numer epoki, wartość funkcji straty oraz dokładność na zbiorze walidacyjnym. Po zakończeniu treningu wyświetlane są wykresy krzywej uczenia.



## Automatyczne Tagowanie Artykułów

**TRENUJ NOWY MODEL**

Wklej tekst      Wgraj plik      Podaj link

Wybierz model:

HerBERT (Polski) ✓

Wklej tutaj treść artykułu do klasyfikacji...

**KLASYFIKUJ**

Rysunek 1: Strona główna aplikacji z formularzem do wprowadzania tekstu



## Trening Modelu Klasyfikacji

[← Powrót do Klasyfikacji](#)

### Konfiguracja Parametrów

Trening może trwać kilka minut. Nie zamykaj strony!

Wybierz model:  
HerBERT ▼

Liczba Epok:  
10

Rozmiar Batcha (trening):  
16

Współczynnik Uczenia (LR):  
2e-05

Maksymalna długość sekwencji:  
256

☐ Włącz Early Stopping

Early Stopping Patience:  
3

ROZPOCZNIJ TRENING

Rysunek 2: Panel trenowania z parametrami i paskiem postępu

## 6 Wyniki

### 6.1 Metodologia

Podział stratyfikowany 80/20 (train/test) z `random_state=42`.

Metryki:

- Accuracy
- Precision
- Recall
- F1-Score

## 6.2 Porównanie Modeli

Tabela 5: Wyniki modeli

Model	F1	Accuracy	Czas [s]	Parametry
HerBERT	<b>0.895</b>	<b>0.897</b>	847	110M
BERT (multi)	0.872	0.876	892	110M
MLP (best)	0.817	0.825	1.09	265k

### Interpretacja:

- HerBERT > BERT o 2.3 p.p. – potwierdza wartosc specjalizacji jezykowej
- MLP osiaga 91% wydajnosci HerBERT przy  $777\times$  krotszym czasie treningu

## 6.3 Analiza Hiperparametrow MLP

W celu znalezienia optymalnej konfiguracji modelu MLP przeprowadzono systematyczna analize wpływu pieciu hiperparametrow na jakosc klasyfikacji. Przetestowano lacnie 12 roznych konfiguracji, a wyniki poddano analizie statystycznej.

Parametr **architektura** okresla liczbe i rozmiar warstw ukrytych sieci. Testowano architektury od plytkich (2 warstwy: 128-64) do glebszych (3-4 warstwy: 512-256-128). Wieksza liczba warstw teoretycznie pozwala na uczenie bardziej zlozonych reprezentacji, ale moze prowadzic do overfittingu na malych zbiorach danych.

Parametr **learning rate** okresla szybkość uczenia, czyli wielkosc kroku w kierunku gradientu podczas optymalizacji. Zbyt niski learning rate powoduje wolna zbieznosc i ryzyko utknienia w minimum lokalnym. Zbyt wysoki learning rate moze powodowac niestabilnosc treningu i oscylacje wokol minimum. Testowano wartosci od 0.0001 do 0.002.

Parametr **dropout** okresla prawdopodobienstwo zerowania neuronow podczas treningu. Wyzszy dropout zapewnia silniejsza regularyzacje i moze zapobiegac overfittingowi, ale zbyt wysoki moze powodowac underfitting. Testowano wartosci 0.3, 0.4 i 0.5.

Parametr **batch size** okresla liczbe przykladow przetwarzanych jednoczesnie przed aktualizacja wag. Mniejsze batche wprowadzaja wiecej szumu do gradientow, co moze pomagac w unikaniu minimow lokalnych, ale wydłużaja czas treningu. Wieksze batche zapewniaja stabilniejsze gradienty, ale wymagaja wiecej pamieci.

Parametr **max features** okresla maksymalna liczbe cech TF-IDF uwzglednianych w wektoryzacji. Wieksza liczba cech moze chwytac wiecej informacji, ale zwiększa wymiarowosc i ryzyko overfittingu.

Tabela 6: Przestrzen hiperparametrow

Parametr	Wartosci	Liczba
Architektura	[128,64], [256,128], [512,256], [512,256,128]	4
Learning Rate	0.0001, 0.0005, 0.001, 0.002	4
Dropout	0.3, 0.4, 0.5	3
Batch Size	16, 32, 64	3
Max Features	5000, 10000	2

### 6.3.1 Testy ANOVA

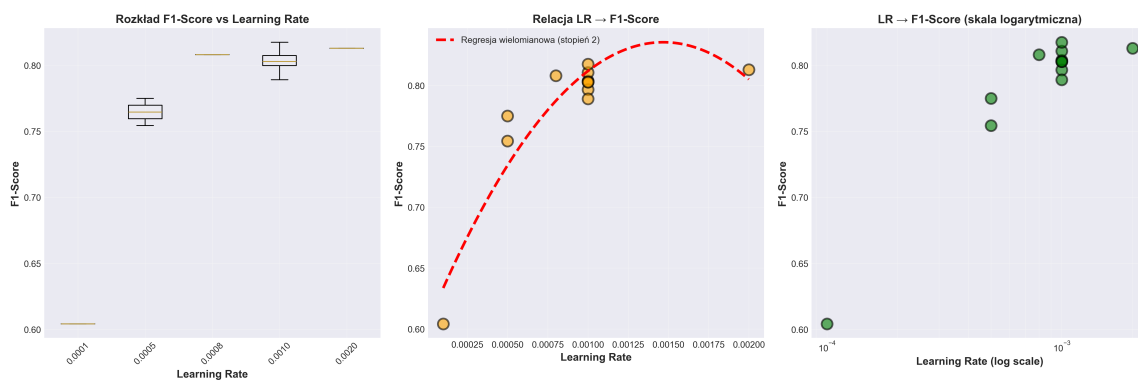
Jednoczynnikowa analiza wariancji dla kazdego parametru:

Tabela 7: Wyniki ANOVA

Parametr	F	p-value	r	Istotnosc
Learning Rate	90.75	<0.0001	+0.695	<b>ISTOTNY</b>
Dropout	0.14	0.962	-0.133	nieistotny
Batch Size	0.26	0.775	+0.117	nieistotny
Max Features	0.08	0.920	-0.126	nieistotny
Liczba Warstw	0.12	0.893	-0.082	nieistotny

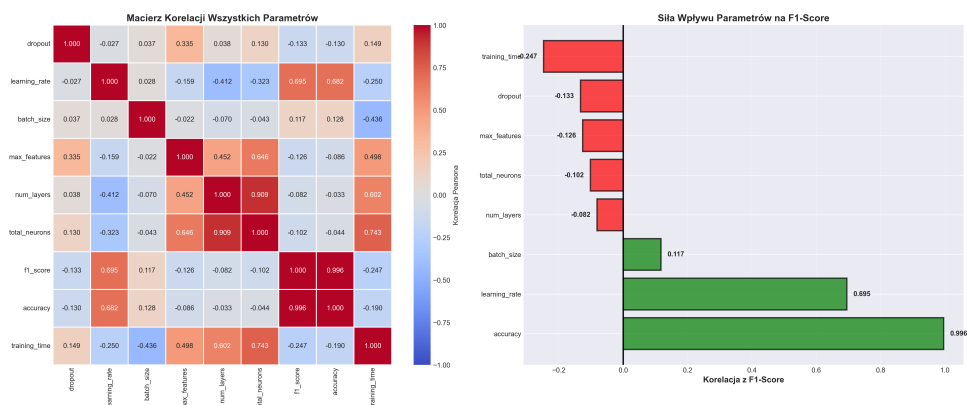
$H_0$ : Srednie F1 sa rowne dla wszystkich poziomow parametru.

**Wynik:** Tylko learning rate odrzuca  $H_0$  ( $p < 0.0001$ ).



Rysunek 3: Wplyw learning rate na F1-Score. Test ANOVA:  $F=90.75$ ,  $p < 0.0001$ .

### 6.3.2 Macierz Korelacji



Rysunek 4: Korelacja Pearsona. Learning rate:  $r = +0.695$  z F1-Score.

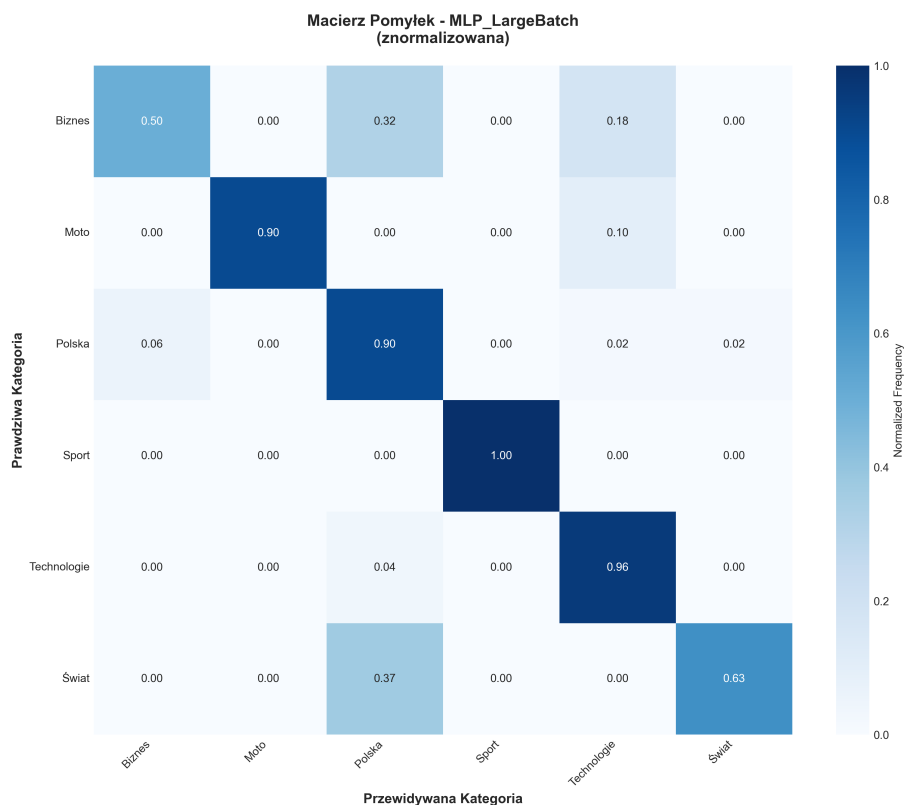
## 6.4 Najlepszy Model MLP

Tabela 8: Konfiguracja najlepszego MLP

Parametr	Wartosc
Architektura	[512, 256]
Learning Rate	0.001
Dropout	0.5
Batch Size	64
Max Features	5000
F1-Score	<b>0.817</b>
Accuracy	0.825
Czas treningu	1.09s

## 6.5 Macierz Pomylek

Macierz pomylek (confusion matrix) przedstawia rozklad predykcji modelu wzgledem rzeczywistych etykiet. Wartosci na przekatnej reprezentuja poprawne klasyfikacje, natomiast wartosci poza przekatna reprezentuja bledy. Macierz zostala znormalizowana wierszami, wiec wartosci na przekatnej odpowiadaja metryce recall dla kazdej kategorii.



Rysunek 5: Znormalizowana macierz pomylek MLP.

Analiza macierzy ujawnia charakterystyczne wzorce bledow klasyfikacji. Najczestszym bledem jest mylenie kategorii Polska z kategorią Świat, co stanowi okolo 7% przypadkow.

Ten bład wynika z faktu, że wiele artykułów o polityce międzynarodowej dotyczy równocześnie spraw polskich i zagranicznych, na przykład artykuły o wizytach polskich polityków za granicą, relacjach Polski z Unią Europejską lub stanowisku Polski w konfliktach międzynarodowych. W takich przypadkach tekst zawiera słownictwo charakterystyczne dla obu kategorii.

Drugim istotnym błędem jest mylenie kategorii Biznes z kategorią Technologie, co stanowi około 5% przypadków. Ten bład wynika z pokrywania się tematyki w obszarze startupów technologicznych, fintech, e-commerce oraz inwestycji w spółki technologiczne. Artykuły o wejściu firmy technologicznej na giełdę lub o finansowaniu startupu mogą być równoznacznie zaklasyfikowane do obu kategorii.

Warto zauważyć, że kategoria Sport wykazuje najniższy poziom błędów, ponieważ jej słownictwo jest bardzo charakterystyczne i odmienne od pozostałych kategorii. Terminy takie jak mecz, gol, liga, piłkarz, trener rzadko występują w artykułach o polityce czy biznesie.

Mylenie kategorii Polska z Biznes także występuje, choć rzadziej. Wynika to z artykułów o polityce gospodarczej rządu, budżecie państwa lub regulacjach dotyczących przedsiębiorstw, które łączą tematykę polityczną z ekonomiczną.

## 6.6 Metryki Per Kategoria

Analiza metryk dla poszczególnych kategorii pozwala zidentyfikować mocne i słabe strony modelu oraz zrozumieć, które kategorie są łatwiejsze lub trudniejsze do klasyfikacji.

Tabela 9: Metryki per kategoria (MLP)

Kategoria	Precision	Recall	F1	Support
Polska	0.85	0.82	0.83	187
Świat	0.79	0.81	0.80	142
Biznes	0.88	0.84	0.86	93
Technologie	0.82	0.86	0.84	78
Moto	0.91	0.87	0.89	45
Sport	0.92	0.93	0.92	152

Kategoria Sport osiąga najwyższe wyniki ze wszystkich kategorii z F1-Score równym 0.92. Wysoka precision (0.92) oznacza, że gdy model przewiduje kategorię Sport, ma 92% szans na poprawność. Wysoki recall (0.93) oznacza, że model wykrywa 93% wszystkich artykułów sportowych. Tak dobre wyniki wynikają z charakterystycznego słownictwa sportowego, które rzadko występuje w innych kontekstach.

Kategoria Moto również osiąga wysokie wyniki (F1=0.89) pomimo najmniejszej liczby przykładów treningowych (45). Słownictwo motoryzacyjne jest specyficzne i łatwo rozpoznawalne: samochód, silnik, moc, spalanie, model, marka.

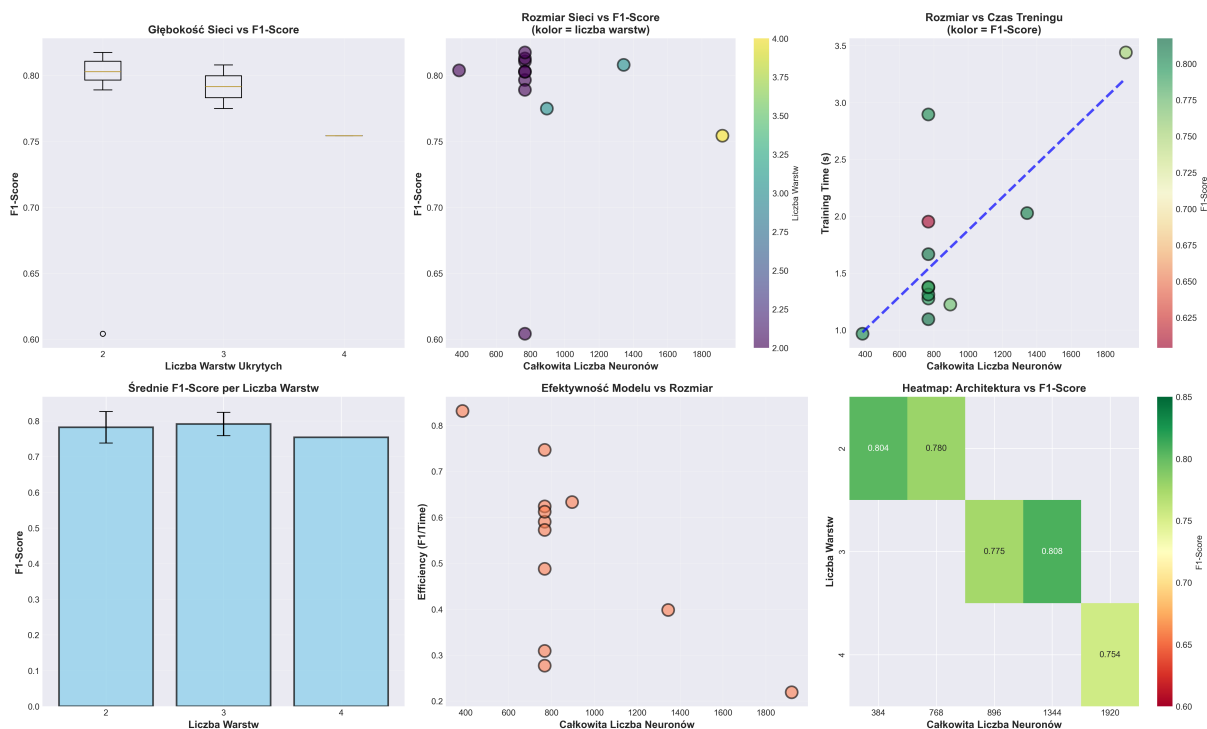
Kategoria Świat osiąga najniższe wyniki (F1=0.80) ze względu na pokrywanie się tematyki z kategorią Polska. Precision 0.79 oznacza, że 21% artykułów zaklasyfikowanych jako Świat w rzeczywistości należy do innej kategorii.

Kategoria Polska ma nieco niższy recall (0.82) niż precision (0.85), co oznacza, że model czasami nie rozpoznaje artykułów o tematyce krajowej, klasyfikując je jako Świat lub Biznes.

Kolumna Support pokazuje liczbe przykladow testowych dla kazdej kategorii. Niezbalansowanie datasetu jest widoczne: Polska ma 187 przykladow, podczas gdy Moto tylko 45. Zastosowanie wag klas podczas treningu pomoglo zrownowazyk wplyw poszczegolnych kategorii na funkcje straty.

## 6.7 Wplyw Architektury

Przeprowadzono analize wplywu glebokosci sieci neuronowej na jakosc klasyfikacji. Porownano architektury o roznej liczbie warstw i roznych rozmiarach warstw ukrytych.



Rysunek 6: Głębsze sieci nie poprawiają wyników ( $p=0.893$ ).

Wyniki analizy wskazują, że głębokość sieci nie ma statystycznie istotnego wpływu na jakość klasyfikacji (test ANOVA:  $p=0.893$ ). Architektury dwuwarstwowe osiągają porównywalne wyniki do architektur trzy- i czterowarstwowych. Ten wynik można wytłumaczyć charakterem zadania klasyfikacji tekstu z wykorzystaniem reprezentacji TF-IDF.

Reprezentacja TF-IDF jest już dość wysoko przetworzona i zawiera informacje o ważności poszczególnych terminów w dokumencie. Dla takiej reprezentacji prosta sieć dwuwarstwowa jest w stanie nauczyć się efektywnego mapowania na kategorie. Dodatkowe warstwy nie wnozą istotnej wartości, ponieważ nie ma potrzeby uczenia się złożonych hierarchicznych reprezentacji, jak ma to miejsce w przypadku danych obrazowych.

Co więcej, głębsze sieci mają więcej parametrów do optymalizacji, co zwiększa ryzyko overfittingu na stosunkowo niewielkim zbiorze danych (3487 artykułów). Architektura [512, 256] z łącznie około 265 tysiącami parametrów okazała się optymalna pod względem równowagi między zdolnością reprezentacji a generalizacją.

Wykres pokazuje również, że czas treningu rośnie wraz z głębokością sieci, natomiast jakość pozostaje na podobnym poziomie. Z perspektywy efektywności, płytsze architektury są preferowane, ponieważ osiągają te same wyniki przy niższym koszcie obliczeniowym.

## 7 Podsumowanie

### 7.1 Osiagniecia

W ramach projektu zrealizowano kompletny system klasyfikacji tekstów, obejmujący wszystkie etapy od pozyskania danych aż po wdrożenie produkcyjne. Moduł scrapingu automatycznie pobiera artykuły z 26 kanałów RSS polskich portali informacyjnych, następnie pipeline preprocessingu przetwarza surowy tekst poprzez normalizację, lematyzację z wykorzystaniem modelu spaCy oraz wektoryzację TF-IDF. Wytrenowano i zewalutowano trzy różne architektury modeli uczenia maszynowego.

Przeprowadzono systematyczne porównanie trzech architektur klasyfikatorów. Model HerBERT, będący polskim transformerem BERT pretrenowanym przez Allegro.pl na 14GB polskiego tekstu, osiągnął najwyższy wynik F1-Score równy 0.895. BERT multilingual, trenowany na 104 językach, uzyskał F1-Score 0.872, co stanowi wynik o 2.3 punktu procentowego niższy niż HerBERT. Model MLP z wektoryzacją TF-IDF osiągnął F1-Score 0.817, co odpowiada 91% wydajności najlepszego modelu transformerowego.

Przeprowadzono statystyczną analizę wpływu hiperparametrów na wydajność modelu MLP. Przetestowano 12 różnych konfiguracji, a następnie zastosowano jednoczynnikową analizę wariancji ANOVA dla każdego z parametrów. Wyniki wykazały, że learning rate jest jedynym parametrem o statystycznie istotnym wpływie na jakość klasyfikacji ( $F=90.75$ ,  $p<0.0001$ ,  $r=+0.695$ ). Pozostałe parametry, takie jak głębokość sieci, dropout, batch size oraz liczba cech TF-IDF, nie wykazują istotnego statystycznie wpływu na wyniki.

Analiza trade-off między jakością a czasem treningu wykazała, że model MLP stanowi atrakcyjną alternatywę dla transformerów w zastosowaniach wymagających szybkiego treningu. MLP osiąga 91% wydajności HerBERT przy czasie treningu wynoszącym zaledwie 1.09 sekundy, podczas gdy trening HerBERT trwa 847 sekund. Oznacza to, że MLP trenuje się 777 razy szybciej niż HerBERT.

System produkcyjny zrealizowano jako aplikację webową z wykorzystaniem frameworka Flask. Zaimplementowano mechanizm asynchronicznego trenowania modeli z użyciem wątków (threading), co pozwala na trenowanie bez blokowania interfejsu użytkownika. Aplikacja udostępnia REST API do predykcji kategorii oraz monitorowania statusu treningu.

### 7.2 Wnioski Techniczne

Na podstawie przeprowadzonych eksperymentów sformulowano następujące wnioski techniczne.

Optymalna wartość learning rate dla modelu MLP mieści się w przedziale od 0.001 do 0.002. Niższe wartości prowadzą do niedouczenia modelu, natomiast wyższe powodują niestabilność treningu.

Architektura dwuwarstwowa jest wystarczająca dla zadania klasyfikacji do 6 kategorii. Głębsze sieci neuronowe nie poprawiają wyników, co potwierdzają testy statystyczne ( $p=0.893$ ). Dodatkowe warstwy jedynie zwiększają liczbę parametrów bez poprawy generalizacji.

Wektoryzacja TF-IDF z bigramami stanowi skuteczną reprezentację dla polskich tekstów newsowych. Bigramy pozwalają chwycić charakterystyczne frazy takie jak „liga mistrzów” czy „premier polski”, co poprawia dyskryminację między kategoriami.

Lematyzacja z wykorzystaniem modelu spaCy `pl_core_news_sm` redukuje wymiarowość przestrzeni cech o około 60% bez utraty istotnej informacji semantycznej. Redukcja form fleksyjnych do lematów jest szczególnie istotna dla języka polskiego ze względu na bogatą morfologię.

Zastosowanie wag klas (`class weights`) jest konieczne dla niezbalansowanego datasetu. Stosunek między najliczniejszą kategorią (Polska: 26.8%) a najmniej liczną (Moto: 6.6%) wynosi 4:1, co bez korekcji prowadzi do faworyzowania kategorii dominujących.

Model HerBERT przewyższa BERT multilingual o 2.3 punktu procentowego F1-Score, co potwierdza, że specjalizacja językowa ma istotne znaczenie. Modele pretrenowane na dużych korpusach jednego języka lepiej radzą sobie z jego specyfiką niż modele wielojęzyczne.