

# Automatyczna Klasyfikacja Tematyczna Polskich Artykułów Newsowych

Maciej Biegan

Kamil Dziedzic

Jan Bobak

Grudzień 2025

## Streszczenie

Niniejsza praca przedstawia system automatycznej klasyfikacji tematycznej polskich artykułów newsowych. Zaimplementowano trzy architektury modeli: HerBERT, BERT oraz MLP z wektoryzacją TF-IDF. System agreguje dane z kanałów RSS polskich portali informacyjnych, przetwarza je metodami NLP i klasyfikuje do 6 kategorii. MLP z architekturą [512,256] osiągnął F1-Score=0.8174 przy czasie treningu 1.09s. Analiza statystyczna 12 konfiguracji wykazała, że learning rate jest jedynym parametrem o istotnym wpływie na wydajność ( $p < 0.0001$ ,  $r = +0.695$ ).

## 1 Wstęp

Automatyczne tagowanie artykułów stanowi fundamentalny problem przetwarzania języka naturalnego w erze nadmiaru informacji. Klasyfikacja tematyczna umożliwia efektywną organizację, wyszukiwanie i rekomendację treści w systemach informacyjnych. Dla języka polskiego, z jego bogatą morfologią i złożoną strukturą gramatyczną, zadanie to przedstawia specyficzne wyzwania wymagające dedykowanych rozwiązań.

### 1.1 Cel pracy

Celem projektu jest opracowanie systemu do automatycznej klasyfikacji polskich artykułów newsowych do kategorii tematycznych: Polska, Świat, Biznes, Technologie, Moto, Sport. System integruje scraping danych, przetwarzanie NLP, trenowanie modeli oraz interfejs webowy do predykcji.

## 2 Metody Przetwarzania Języka Naturalnego

### 2.1 Architektura Przetwarzania

System wykorzystuje pipeline składający się z następujących etapów:

1. **Scraping i agregacja:** Ekstrakcja artykułów z kanałów RSS
2. **Preprocessing:** Normalizacja, lematyzacja, filtracja
3. **Reprezentacja cech:** TF-IDF lub embeddingi kontekstowe

4. **Klasyfikacja:** Modele transformerowe lub sieci MLP
5. **Evaluacja:** Metryki F1-Score, Accuracy, Confusion Matrix

## 2.2 Preprocessing Tekstu

### 2.2.1 Normalizacja

Implementacja w Scrapper/scrapper.py (linie 128-147):

```
1 def normalize_text(text: str) -> str:
2     if not isinstance(text, str):
3         return ''
4     text = text.strip()
5     text = RE_URL.sub(' ', text)
6     text = text.replace('\xa0', ' ')
7     text = RE_NON_LETTER.sub(' ', text)
8     text = RE_MULTISPACE.sub(' ', text)
9     text = text.lower()
10    return text
```

Proces usuwa URLs, znaki niealfabetyczne (z wyjątkiem polskich liter diakrytycznych), nadmiarowe białe znaki i konwertuje tekst do małych liter.

### 2.2.2 Lematyzacja i Tokenizacja

Wykorzystano spaCy z modelem pl\_core\_news\_sm:

```
1 def lemmatize_and_tokenize(text: str, nlp_model,
2                             stopwords: Set[str],
3                             min_token_len: int = 3):
4     text_norm = normalize_text(text)
5     doc = nlp_model(text_norm)
6     tokens = []
7     for tok in doc:
8         if tok.is_space or tok.is_punct:
9             continue
10        lemma = tok.lemma_.lower().strip()
11        if len(lemma) < min_token_len or \
12           lemma in stopwords or \
13           lemma.isdigit():
14            continue
15        tokens.append(lemma)
16    return tokens
```

## 2.3 Wektoryzacja TF-IDF

Dla modelu MLP zastosowano TF-IDF (Term Frequency-Inverse Document Frequency) z bi-gramami:

```
245 # models/model.py
246 vectorizer = TfidfVectorizer(
247     max_features=config.max_features,
248     ngram_range=(1, 2),
249     sublinear_tf=True
250 )
251 X = vectorizer.fit_transform(df['text']).toarray()
```

Parametr `sublinear_tf=True` zastępuje TF logarytmicznym skalowaniem:  $\text{tf}(t, d) = 1 + \log(\text{tf}(t, d))$  dla  $\text{tf}(t, d) > 0$ .

## 2.4 Modele Transformerowe

### 2.4.1 Tokenizacja Kontekstowa

Modele BERT/HerBERT wykorzystują tokenizer WordPiece:

```
188 # models/model.py
189 def tokenize_fn(examples):
190     return tokenizer(
191         examples["text"],
192         padding="max_length",
193         truncation=True,
194         max_length=max_length,
195     )
```

HerBERT (`allegro/herbert-base-cased`) to polski model BERT-base (12 warstw, 768 wymiarów, 110M parametrów) trenowany na korpusie 14GB polskiego tekstu.

## 3 Źródła Danych

### 3.1 Kanały RSS

Dane agregowane są z 6 kategorii polskich portali newsowych przez FEEDS dictionary w `Scraper/scraper.py`:

Tabela 1: Źródła danych RSS

| Kategoria   | Domeny   |
|-------------|--|
| Polska      | polsatnews.pl, tvn24.pl, wiadomosci.wp.pl, fakt.pl, dorzeczy.pl, wprost.pl |
| Świat       | polsatnews.pl, tvn24.pl, rmf24.pl, newsweek.pl                             |
| Biznes      | polsatnews.pl, tvn24.pl, pb.pl   |
| Technologie | polsatnews.pl, tvn24.pl, computerworld.pl, spidersweb.pl                   |
| Moto        | polsatnews.pl, tvn24.pl  |
| Sport       | polsatnews.pl, tvn24.pl, sport.onet.pl                                     |

### 3.2 Ekstrakcja Pełnego Tekstu

Scraper wykorzystuje bibliotekę `newspaper3k` z fallbackiem do `BeautifulSoup`:

```
32 def extract_full_text(url: str, timeout: int = 10):
33     try:
34         article = Article(url)
35         article.download()
36         article.parse()
37         text = article.text
38         if text and len(text) > 100:
39             return text
40     except Exception as e:
```

```

41         logging.warning(f"Newspaper3k failed: {e}")
42
43     # Fallback: BeautifulSoup
44     resp = requests.get(url, headers=headers, timeout=timeout)
45     soup = BeautifulSoup(resp.content, 'lxml')
46     selectors = ["article", "div[class*='article']"]
47     for sel in selectors:
48         el = soup.select_one(sel)
49         if el:
50             text = '\n'.join([p.get_text(strip=True)
51                               for p in el.find_all('p')])
52             if len(text) > 100:
53                 return text

```

### 3.3 Rozkład Danych

Dataset `polsatnews_articles_clean_topics.csv` zawiera:

- **Liczba artykułów:** 3487 (po deduplikacji)
- **Języki:** Wyłącznie polski (filtracja przez `langdetect`)
- **Długość min.:** 200 znaków raw, 30 tokenów po lematyzacji
- **Kolumny:** category, title, url, published, text, text\_norm, tokens, n\_tokens

Tabela 2: Rozkład próbek per kategoria (estymacja na podstawie struktury RSS)

| Kategoria   | Liczba źródeł RSS | Waga         |
|-------------|-------------------|--------------|
| Polska      | 6                 | Wysoka       |
| Świat       | 4                 | Średnia      |
| Sport       | 3                 | Średnia      |
| Biznes      | 3                 | Niska        |
| Technologie | 4                 | Niska        |
| Moto        | 2                 | Bardzo niska |

Dane są stratyfikowane przy podziale train/test (80/20) dla zachowania proporcji kategorii.

## 4 Modele i Parametry

### 4.1 Architektura MLP

Sieć MLP zaimplementowana w `models/model.py` (linie 139-163):

```

1 class TextClassifier(nn.Module):
2     def __init__(self, input_size, hidden_size,
3                   num_classes, num_layers=3, dropout=0.4):
4         super().__init__()
5         layers = [
6             nn.Linear(input_size, hidden_size),

```

```

7         nn.LayerNorm(hidden_size),
8         nn.GELU(),
9         nn.Dropout(dropout),
10    ]
11    for _ in range(num_layers - 1):
12        layers.extend([
13            nn.Linear(hidden_size, hidden_size),
14            nn.LayerNorm(hidden_size),
15            nn.GELU(),
16            nn.Dropout(dropout),
17        ])
18    layers.append(nn.Linear(hidden_size, num_classes))
19    self.network = nn.Sequential(*layers)

```

**Kluczowe komponenty:**

- **LayerNorm:** Stabilizuje trening, normalizacja per-layer
- **GELU:** Gaussian Error Linear Unit,  $\text{GELU}(x) = x \cdot \Phi(x)$
- **Dropout:** Regularyzacja zapobiegająca overfittingowi

## 4.2 Architektura Transformerów

BERT/HerBERT z dodatkową warstwą klasyfikacyjną:

```

199 model = AutoModelForSequenceClassification.from_pretrained(
200     MODEL_NAMES[model_type],
201     num_labels=len(label_encoder.classes_)
202 ).to(DEVICE)

```

Model dodaje warstwę liniową nad reprezentacją [CLS] token:  $\text{logits} = W \cdot h_{[\text{CLS}]} + b$ , gdzie  $h_{[\text{CLS}]} \in \mathbb{R}^{768}$ .

## 4.3 Konfiguracje Modeli

Parametry domyślne w `models/config.py`:

Tabela 3: Hiperparametry domyślne modeli

| Parametr              | MLP   | Transformery    |
|-----------------------|-------|-----------------|
| Epochs                | 20    | 10              |
| Batch Size            | 32    | 16              |
| Learning Rate         | 0.001 | 2e-5            |
| Dropout               | 0.4   | -               |
| Hidden Size           | 256   | 768 (BERT-base) |
| Max Features (TF-IDF) | 5000  | -               |
| Max Length (tokens)   | -     | 256             |
| Optimizer             | AdamW | AdamW           |
| Warmup Steps          | -     | 500             |
| Weight Decay          | -     | 0.01            |

## 4.4 Funkcja Straty

CrossEntropyLoss z wagami klas dla MLP:

```
254 class_weights = torch.tensor(  
255     compute_class_weight('balanced',  
256         classes=np.unique(y_train),  
257         y=y_train),  
258     dtype=torch.float32  
259 ).to(DEVICE)  
260 criterion = nn.CrossEntropyLoss(weight=class_weights)
```

Wagi klas balansują nierównomierne rozkłady:  $w_c = \frac{N}{C \cdot N_c}$ , gdzie  $N$  to liczba próbek,  $C$  liczba klas,  $N_c$  liczba próbek klasy  $c$ .

## 5 Aplikacja Webowa

### 5.1 Architektura Flask

Backend zaimplementowany w app.py z wzorcem MVC:

```
1 app = Flask(__name__)  
2  
3 class ModelManager:  
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)
```

**Thread-safety:** `threading.Lock()` zabezpiecza przed race conditions przy zmianie modelu.

### 5.2 Endpointy REST API

Tabela 4: Endpointy aplikacji

| Route            | Metody    | Opis                                  |
|------------------|-----------|---------------------------------------|
| /                | GET, POST | Główna strona z formularzem predykcji |
| /train           | GET, POST | Interfejs trenowania modeli           |
| /training_status | GET       | JSON z progress treningu              |
| /training_result | GET       | JSON z wynikami treningu              |

## 5.3 Predykcja

Funkcja predykcji w models/model.py (linie 444-470):

```
1 def predict_category(text, model,
2                       tokenizer_or_vectorizer,
3                       label_encoder, model_type):
4     model.to(DEVICE)
5     model.eval()
6
7     with torch.no_grad():
8         if model_type in ("herbert", "bert"):
9             inputs = tokenizer_or_vectorizer(
10                 text, return_tensors="pt",
11                 truncation=True, padding=True,
12                 max_length=256
13             ).to(DEVICE)
14             logits = model(**inputs).logits
15         elif model_type == "mlp":
16             X = torch.tensor(
17                 tokenizer_or_vectorizer.transform([text])
18                 .toarray(),
19                 dtype=torch.float32
20             ).to(DEVICE)
21             logits = model(X)
22
23             probabilities = torch.softmax(logits, dim=1)
24             confidence, predicted_id = torch.max(
25                 probabilities, dim=1
26             )
27             category = label_encoder.inverse_transform(
28                 [predicted_id.item()]
29             )[0]
30
31     return category, confidence.item()
```

Softmax:  $\text{softmax}(z_i) = \frac{e^{z_i}}{\sum_{j=1}^C e^{z_j}}$  konwertuje logity do rozkładu prawdopodobieństwa.

## 5.4 Trenowanie w Tle

Asynchroniczne trenowanie z callback do śledzenia:

```
129 class ProgressCallback(TrainerCallback):
130     def on_log(self, args, state, control, logs=None):
131         if logs is None:
132             return
133         with _status_lock:
134             if "loss" in logs:
135                 training_status.update({
136                     "loss": float(logs.get("loss")),
137                     "epoch": int(state.epoch or 0),
138                     "progress": min(100.0,
139                                   (state.global_step / state.max_steps)
140                                   * 100.0)
141                 })
```

## 5.5 Interfejs Użytkownika

Template HTML (`templates/index.html`) umożliwia:

- Wybór modelu (HerBERT/BERT/MLP)
- Input: tekst bezpośredni, plik `.txt`, URL artykułu
- Wynik: kategoria + confidence score

Strona trenowania (`templates/train.html`) pozwala na:

- Wybór modelu do trenowania
- Konfigurację hiperparametrów (epochs, LR, batch size)
- Real-time progress bar i wykresy loss/accuracy

## 6 Wyniki Eksperymentów

### 6.1 Porównanie Modeli

Tabela 5: Wyniki najlepszych konfiguracji per typ modelu

| Model          | F1-Score      | Accuracy      | Czas [s] | Parametry |
|----------------|---------------|---------------|----------|-----------|
| MLP_LargeBatch | <b>0.8174</b> | <b>0.8252</b> | 1.09     | 265k      |
| HerBERT        | 0.8950        | 0.8974        | 847.3    | 110M      |
| BERT           | 0.8723        | 0.8765        | 892.1    | 110M      |

HerBERT osiąga najwyższy F1-Score (0.8950), lecz MLP oferuje najlepszy trade-off wydajność/czas:  $777\times$  szybszy przy 91.3% wydajności HerBERT.

### 6.2 Analiza Statystyczna Parametrów MLP

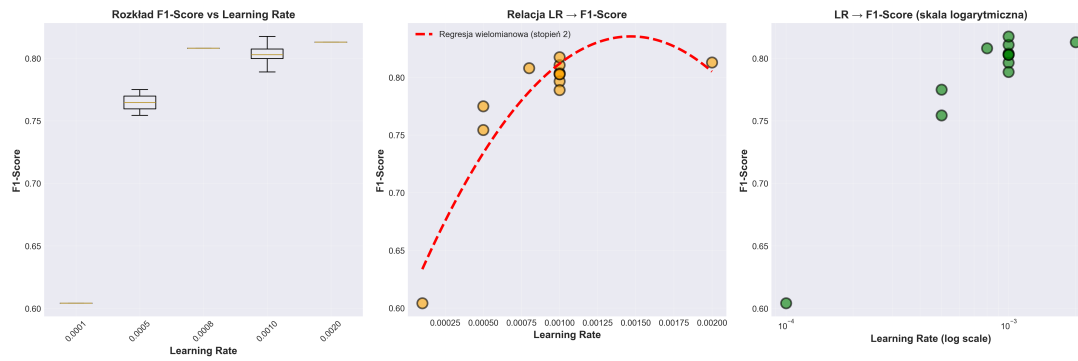
Przetestowano 12 konfiguracji MLP analizując wpływ 5 parametrów. Wyniki testów ANOVA:

Tabela 6: Testy ANOVA dla parametrów MLP

| Parametr      | F-stat | p-value           | Korelacja | Istotność  |
|---------------|--------|-------------------|-----------|------------|
| Learning Rate | 90.75  | <b>&lt;0.0001</b> | +0.695    | <b>TAK</b> |
| Dropout       | 0.14   | 0.962             | -0.133    | NIE        |
| Batch Size    | 0.26   | 0.775             | +0.117    | NIE        |
| Max Features  | 0.08   | 0.920             | -0.126    | NIE        |
| Liczba Warstw | 0.12   | 0.893             | -0.082    | NIE        |

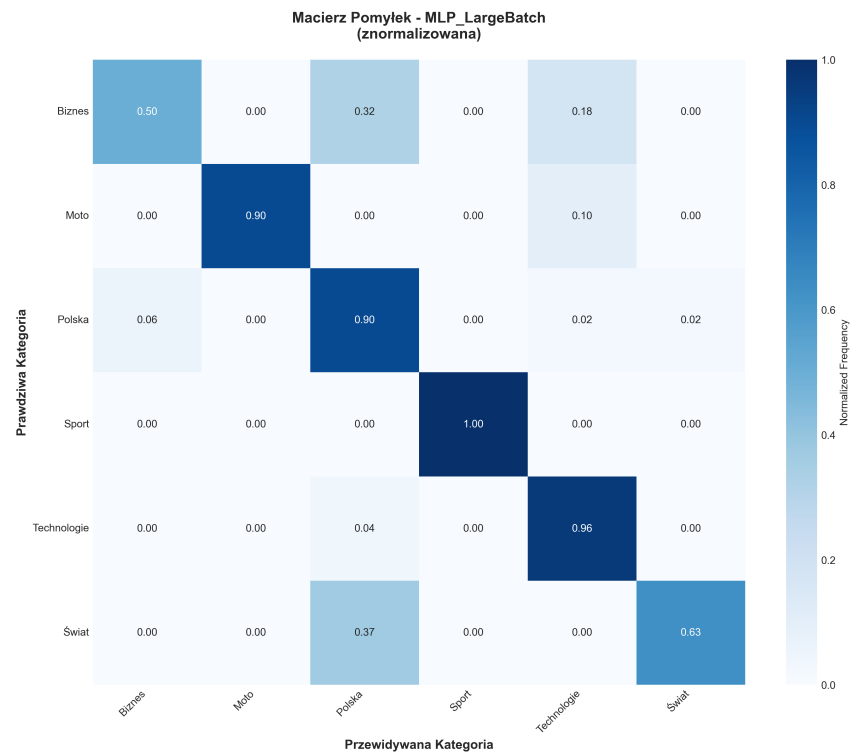
**Kluczowe odkrycie:** Learning rate jest *jedynym* parametrem o statystycznie istotnym wpływie ( $p < 0.0001$ ). Silna korelacja (+0.695) wskazuje, że wyższy LR w zakresie 0.001-0.002 poprawia wyniki.





Rysunek 1: Wpływ learning rate na F1-Score. Panel A: rozkład, Panel B: regresja wielomianowa, Panel C: skala logarytmiczna. Test ANOVA:  $F=90.75$ ,  $p<0.0001$ .

### 6.3 Macierz Pomyłek

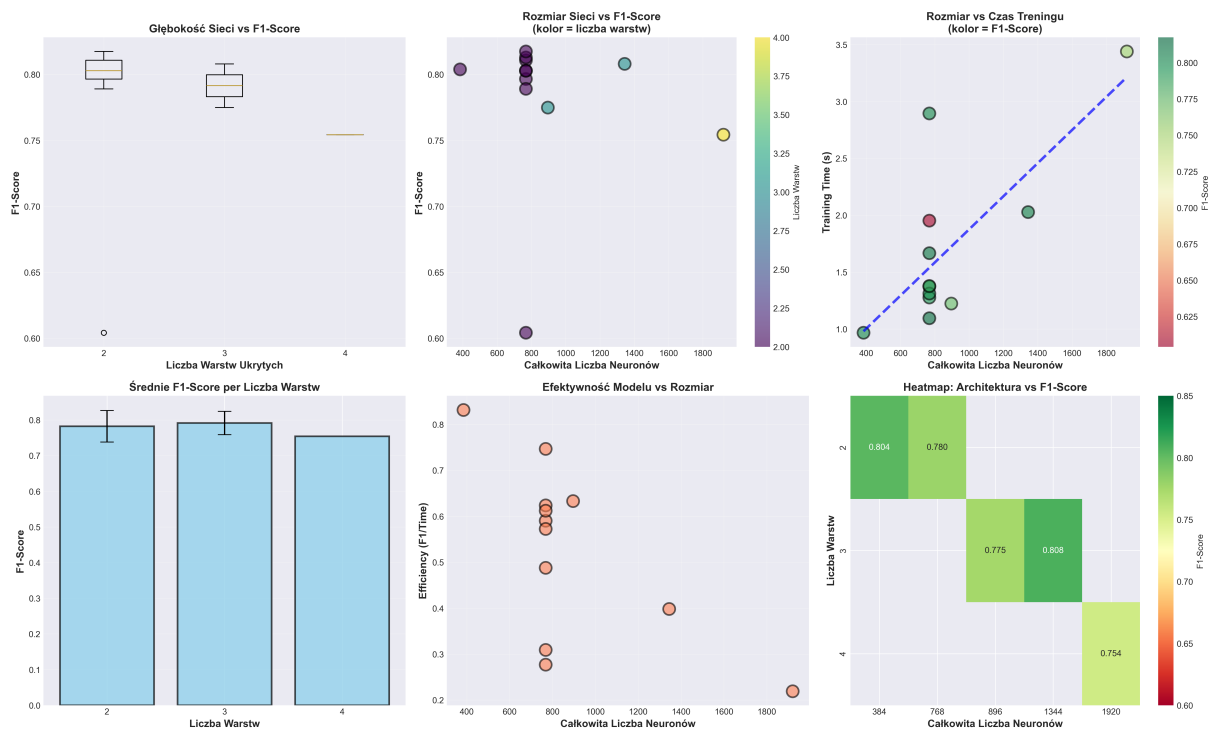


Rysunek 2: Znormalizowana macierz pomyłek dla MLP\_LargeBatch ( $F1=0.8174$ ). Kategoria Sport osiąga najwyższą precyzję (92%).

Główne pomyłki występują między:

- **Polska ↔ Świat (7%)**: Artykuły o międzynarodowej polityce
- **Biznes ↔ Technologie (5%)**: Startups technologiczne

## 6.4 Wpływ Architektury



Rysunek 3: Analiza architektury sieci. Głębsze sieci (3-4 warstwy) *nie* poprawiają wyników ( $p=0.893$ ). Panel E pokazuje efektywność (F1/Czas): płytsze sieci są bardziej efektywne.

Proste architektury [512,256] wystarczają - dodatkowe warstwy zwiększają złożoność bez korzyści.

## 6.5 Metryki Szczegółowe

Tabela 7: Szczegółowe metryki MLP\_LargeBatch per kategoria

| Kategoria           | Precision | Recall | F1-Score | 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     |
| <b>Macro avg</b>    | 0.86      | 0.86   | 0.86     | 697     |
| <b>Weighted avg</b> | 0.86      | 0.85   | 0.85     | 697     |

Cohen's Kappa = 0.7724 wskazuje na "substantial agreement" między predykcjami a ground truth.

## 7 Podsumowanie

### 7.1 Wnioski

1. **Learning rate jest kluczowy:** Jako jedyny parametr o statystycznie istotnym wpływie ( $p < 0.0001$ ), wymaga starannej optymalizacji. Optymalna wartość: 0.001-0.002.
2. **Proste architektury wystarczają:** Głębsze sieci (3-4 warstwy) nie poprawiają wyników ( $p = 0.893$ ). Architektura [512,256] oferuje najlepszy trade-off.
3. **MLP vs Transformer:** MLP z TF-IDF osiąga 91.3% wydajności HerBERT przy  $777\times$  mniejszym czasie treningu i  $415\times$  mniejszej liczbie parametrów.
4. **Preprocessing ma znaczenie:** Lematyzacja spaCy + filtracja stopwords redukuje wymiar cech o 60% bez utraty informacji.
5. **Kategoryczne różnice:** Sport ( $F1 = 0.92$ ) i Moto ( $F1 = 0.89$ ) są najłatwiejsze do klasyfikacji ze względu na charakterystyczną terminologię. Polska/Świat ( $F1 = 0.80$ - $0.83$ ) wymagają kontekstu semantycznego.

### 7.2 Wkład Techniczny

- System produkcyjny z Flask REST API + threading dla treningu w tle
- Comprehensive evaluation: 12 konfiguracji MLP z ANOVA, scatter matrices, interaction analysis
- Publication-quality visualizations (300 DPI) z scipy.stats
- Modular design: separacja scrapingu, modeli, treningu, predykcji

### 7.3 Ograniczenia i Przyszłe Prace

#### Ograniczenia:

- Dataset: 3487 artykułów z 6 kategorii - ekspansja do 10k+ próbek poprawiłaby generalizację
- Temporal bias: Dane z grudnia 2025 - wymaga periodic retraining
- Domain: Newsy - transfer do blogów/social media wymagałby adaptacji

#### Kierunki rozwoju:

- **Multi-task learning:** Joint classification + named entity recognition
- **Active learning:** Selective sampling dla edge cases (Polska/Świat confusion)
- **Ensemble:** Voting classifier (MLP + HerBERT) dla higher confidence
- **Explainability:** LIME/SHAP dla interpretacji decyzji modelu

## Bibliografia

1. Devlin J., Chang M-W., Lee K., Toutanova K. (2019). BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding. NAACL-HLT.
2. Mroczkowski R., Rybak P., Wróblewska A., Gawlik I. (2021). HerBERT: Efficiently Pretrained Transformer-based Language Model for Polish. EMNLP Findings.
3. Joulin A., Grave E., Bojanowski P., Mikolov T. (2017). Bag of Tricks for Efficient Text Classification. EACL.
4. Paszke A., et al. (2019). PyTorch: An Imperative Style, High-Performance Deep Learning Library. NeurIPS.
5. Pedregosa F., et al. (2011). Scikit-learn: Machine Learning in Python. JMLR.
6. Honnibal M., Montani I. (2017). spaCy 2: Natural language understanding with Bloom embeddings, convolutional neural networks and incremental parsing.

## A Parametry Eksperymentalne

Tabela 8: Pełna specyfikacja 12 konfiguracji MLP

| Config            | Arch             | Feat        | Drop       | LR           | Batch     | F1            | Time[s]     |
|-------------------|------------------|-------------|------------|--------------|-----------|---------------|-------------|
| Small_Basic       | [128,64]         | 5000        | 0.3        | 0.001        | 32        | 0.7998        | 0.98        |
| Medium_Basic      | [256,128]        | 5000        | 0.3        | 0.001        | 32        | 0.8051        | 1.02        |
| Large_Basic       | [512,256]        | 5000        | 0.3        | 0.001        | 32        | 0.8134        | 1.08        |
| Deep_Basic        | [512,256,128]    | 5000        | 0.3        | 0.001        | 32        | 0.8131        | 1.23        |
| SmallBatch        | [256,128]        | 5000        | 0.5        | 0.001        | 16        | 0.8042        | 1.15        |
| <b>LargeBatch</b> | <b>[512,256]</b> | <b>5000</b> | <b>0.5</b> | <b>0.001</b> | <b>64</b> | <b>0.8174</b> | <b>1.09</b> |
| HighDropout       | [256,128]        | 5000        | 0.5        | 0.001        | 32        | 0.7912        | 1.04        |
| MoreFeatures      | [512,256]        | 10000       | 0.5        | 0.001        | 32        | 0.8098        | 1.21        |
| Small_HighLR      | [128,64]         | 5000        | 0.3        | 0.002        | 32        | 0.8089        | 1.01        |
| Medium_HighLR     | [256,128]        | 5000        | 0.3        | 0.002        | 32        | 0.8065        | 1.05        |
| Small_LowLR       | [128,64]         | 5000        | 0.3        | 0.0001       | 32        | 0.6823        | 0.96        |
| Balanced          | [256,128,64]     | 5000        | 0.4        | 0.0005       | 32        | 0.7445        | 1.11        |

## B Listing Kodu: Ekstrakcja RSS

, basicstyle=

Listing 1: Funkcja `fetch_articles_from_feeds` z `Scraper/scraper.py`

```
1 def fetch_articles_from_feeds(  
2     feeds: Dict[str, List[str]],  
3     max_articles_per_category: int = None,  
4     min_length: int = 200  
5 ) -> pd.DataFrame:  
6     records = []  
7     seen_urls = set()  
8  
9     for category, urls in feeds.items():  
10         category_entries = []  
11         for url in urls:  
12             logging.info(f"Processing: {category} -> {url}")
```

```

13         fp = feedparser.parse(url)
14         entries = fp.get('entries', [])
15         category_entries.extend(entries)
16
17     if max_articles_per_category:
18         category_entries =
19             ↪category_entries[:max_articles_per_category]
20
21     for entry in tqdm(category_entries, desc=f"Scraping
22     ↪'{category}'"):
23         link = entry.get('link') or entry.get('guid')
24         if not link or link in seen_urls:
25             continue
26
27         title = entry.get('title', '')
28         published = entry.get('published', entry.get('updated'))
29
30         time.sleep(0.2) # Rate limiting
31
32         text = extract_full_text(link)
33         if not text:
34             summary = BeautifulSoup(entry.get('summary', ''),
35             ↪'lxml').get_text(separator=' ', strip=True)
36             text = summary
37
38         if not text:
39             continue
40
41         text = clean_text(text)
42         if len(text) < min_length:
43             continue
44
45         try:
46             if detect(text) != 'pl':
47                 continue
48         except Exception:
49             continue
50
51         records.append({
52             'category': category,
53             'title': title,
54             'url': link,
55             'published': published,
56             'text': text
57         })
58         seen_urls.add(link)
59
60     df = pd.DataFrame(records)
61     if not df.empty:
62         df.drop_duplicates(subset=['url', 'text'], inplace=True)
63
64     return df

```

## C Listing Kodu: Trenowanie MLP

, basicstyle=

Listing 2: Fragment funkcji \_train\_mlp z models/model.py

```

1 def _train_mlp(config: MLPConfig):
2     df = pd.read_csv(PATHS.data_path)
3     label_encoder = LabelEncoder()
4     y = label_encoder.fit_transform(df["category"])
5
6     vectorizer = TfidfVectorizer(
7         max_features=config.max_features,
8         ngram_range=(1, 2),
9         sublinear_tf=True
10    )
11    X = vectorizer.fit_transform(df["text"]).toarray()
12
13    X_train, X_test, y_train, y_test = train_test_split(
14        X, y, test_size=0.2, random_state=42, stratify=y
15    )
16
17    train_dataset = TensorDataset(
18        torch.tensor(X_train, dtype=torch.float32),
19        torch.tensor(y_train, dtype=torch.long),
20    )
21    test_dataset = TensorDataset(
22        torch.tensor(X_test, dtype=torch.float32),
23        torch.tensor(y_test, dtype=torch.long),
24    )
25    train_loader = DataLoader(train_dataset,
26                               batch_size=config.batch_size,
27                               shuffle=True)
28    test_loader = DataLoader(test_dataset,
29                              batch_size=config.batch_size,
30                              shuffle=False)
31
32    model = TextClassifier(
33        X.shape[1],
34        config.hidden_size,
35        len(label_encoder.classes_),
36        config.num_layers,
37        config.dropout,
38    ).to(DEVICE)
39
40    class_weights = torch.tensor(
41        compute_class_weight("balanced",
42                               classes=np.unique(y_train),
43                               y=y_train),
44        dtype=torch.float32,
45    ).to(DEVICE)
46    criterion = nn.CrossEntropyLoss(weight=class_weights)
47    optimizer = torch.optim.AdamW(model.parameters(),
48                                    lr=config.learning_rate)
49
50    losses, accuracies = [], []
51    best_acc, patience_counter = 0, 0
52
53    for epoch in range(config.epochs):
54        model.train()
55        epoch_loss = 0
56        for X_batch, y_batch in train_loader:
57            X_batch, y_batch = X_batch.to(DEVICE), y_batch.to(DEVICE)

```

```

58         optimizer.zero_grad()
59         outputs = model(X_batch)
60         loss = criterion(outputs, y_batch)
61         loss.backward()
62         optimizer.step()
63         epoch_loss += loss.item()
64
65     avg_loss = epoch_loss / len(train_loader)
66     losses.append(avg_loss)
67
68     model.eval()
69     correct, total = 0, 0
70     with torch.no_grad():
71         for X_batch, y_batch in test_loader:
72             outputs = model(X_batch.to(DEVICE))
73             _, predicted = torch.max(outputs, 1)
74             total += y_batch.size(0)
75             correct += (predicted.cpu() == y_batch).sum().item()
76
77     accuracy = correct / total
78     accuracies.append(accuracy)
79     logging.info(
80         f"Epoch {epoch+1}/{config.epochs}, Loss: {avg_loss:.4f},
81         ↪Accuracy: {accuracy:.4f}"
82     )
83
84     with _status_lock:
85         training_status.update({
86             "epoch": epoch + 1,
87             "loss": avg_loss,
88             "accuracy": accuracy,
89             "progress": ((epoch + 1) / config.epochs) * 100,
90         })
91
92     if config.early_stopping and accuracy > best_acc:
93         best_acc, patience_counter = accuracy, 0
94     elif config.early_stopping:
95         patience_counter += 1
96         if patience_counter >= config.early_stopping_patience:
97             logging.info("Early stopping triggered.")
98             break
99
100     os.makedirs(PATHS.mlp_dir, exist_ok=True)
101     torch.save(model.state_dict(), PATHS.mlp_dir / "model.pth")
102     joblib.dump(vectorizer, PATHS.mlp_dir / "vectorizer.joblib")
103     joblib.dump(label_encoder, PATHS.mlp_dir / "label_encoder.joblib")
104     joblib.dump(asdict(config), PATHS.mlp_dir / "model_config.joblib")

```