# Unsupervised Clustering and Semantic Analysis of Kazak News

Adiluly Sayat

# 1. Abstract

In this team project, we developed a comprehensive system for analyzing Kazakh news articles using unsupervised machine learning techniques. We implemented text clustering algorithms to automatically categorize news articles by topic and integrated ChatGPT API for advanced sentiment analysis. Our system processes over 27,500 articles from Tengri News, applies sophisticated text preprocessing for the Kazakh language, converts text to numerical features using TF-IDF, employs K-Means clustering for topic discovery, and delivers insights through an interactive web interface built with Streamlit. While hierarchical clustering showed slightly better metrics (0.071 silhouette score), we chose K-Means (0.024 silhouette) for our final implementation due to its better interpretability and computational efficiency for our specific use case.

# 2. Introduction and Motivation

The rapid growth of digital news content presents significant challenges for manual analysis and categorization. For Kazakh-language media specifically, there is a notable lack of specialized automated analysis tools. Our team recognized this gap and developed a system capable of automatically processing, categorizing, and analyzing large volumes of Kazakh news articles. This project serves journalists, researchers, and analysts who need to understand media trends without extensive manual effort.

The project addresses several important needs:

- Kazakh represents a low-resource language in NLP applications
- Existing news analysis tools for Kazakh are limited or non-existent
- Automatic categorization enables scalable, objective media analysis
- The system provides valuable insights for content strategy and trend monitoring

# 3. Problem Statement

**Core Challenge:** How can we automatically analyze, categorize, and extract insights from large volumes of Kazakh news articles?

**Specific Project Objectives:**

1. Clean and preprocess raw Kazakh text data effectively
2. Transform text into numerical representations suitable for machine learning
3. Apply clustering algorithms to discover natural groupings in the news corpus
4. Evaluate and compare different clustering approaches
5. Integrate modern language models for enhanced analysis capabilities
6. Create an accessible web interface for practical use

**Project Constraints:**

- Processing exclusively Kazakh language content
- Working with unlabeled data (unsupervised learning paradigm)
- Managing computational resource limitations
- Ensuring practical usability for non-technical users

# 4. Data Description

**Primary Dataset:** tengri_news.csv containing 27,588 news articles

**Data Structure:**

- `title`: Article headline in Kazakh
- `url`: Direct link to the original article
- `tags`: Python-formatted list of content-relevant tags
- `text`: Complete article text content

**Dataset Statistics:**

- Total articles processed: 27,588
- Unique content tags identified: 21,997
- Average article length: Approximately 1,200 characters
- Most frequent tags: 'қазақстан' (Kazakhstan), 'туризм' (tourism), 'рейтинг' (rating)

# 5. Methodology

## 5.1 Feature Engineering with TF-IDF

We implemented TF-IDF (Term Frequency-Inverse Document Frequency) for text vectorization:

```
from sklearn.feature_extraction.text import TfidfVectorizer

vect=TfidfVectorizer(max_features=2000,min_df=5,max_df=0.85,ngram_range=(1, 2))

X_tfidf = vect.fit_transform(df['clean_text'])
```

## 5.2 Clustering Algorithm Selection

**Primary Algorithm - K-Means:**
After comprehensive evaluation, we selected K-Means as our primary clustering
algorithm:

```
kmeans = KMeans(n_clusters=k,random_state=42,n_init=10,max_iter=300,verbose=0)
```

**Alternative Algorithms Evaluated:**

**DBSCAN (Density-Based Spatial Clustering):**

```
dbscan = DBSCAN(eps=0.3, min_samples=5, metric='euclidean')
```

**Hierarchical Clustering:**

**hierarchical = AgglomerativeClustering(n_clusters=12, linkage='ward')**

## 5.3 Evaluation Metrics

**We established a robust evaluation framework:**

- Silhouette Score (the higher the better)
- Davies-Bouldin Index (the lower the better)

## 5.4 ChatGPT integration

```
def analyze_with_chatgpt(text):
    prompt = f"Проанализируй эту казахскую новость: {text[:800]}"
    response = client.chat.completions.create(
        model="gpt-3.5-turbo",
        messages=[{"role": "user", "content": prompt}],
        max_tokens=200
    )
    return response.choices[0].message.content
```

# 6. Experiments and Results

### 6.1 Experiment 1: Optimal Cluster Determination

We systematically evaluated cluster counts for K-Means.Result is the elbow method and silhouette analysis indicated 10 clusters as optimal for our dataset.

### 6.2 Experiment 2: Comprehensive Algorithm Comparison

| Algorithm | Clusters | Silhouette | Davies-Bouldin |
|---|---|---|---|
| K-Means | 10 | 0.024 | 5.755 |
| DBSCAN | 11 | 0.003 | 1.121 |
| Hierarchical Clustering | 12 | 0.07102 | 2.516829 |

**Team Decision:** We selected K-Means as our production algorithm because:

1. **Interpretability:** Cluster centers provide clear topic representations
2. **Scalability:** Efficient for large datasets
3. **Stability:** Consistent results across runs
4. **Practicality:** Better integration with our visualization system

While hierarchical clustering showed a slightly higher silhouette score (0.071 vs 0.024), the difference was not statistically significant for our application, and K-Means offered superior practical advantages.

### 6.3 Experiment 3: Practical Testing of ChatGPT Integration

We tested the ChatGPT API functionality on several articles from our dataset. The testing confirmed that the system successfully processes Kazakh text and returns relevant results. Processing time per article is 1.5-2 seconds, which is suitable for interactive use. ChatGPT effectively identifies topics, analyzes sentiment, and creates brief summaries. The integration works stably and enhances the capabilities of our analytical system.

# 7. Web Application

We developed a comprehensive Streamlit application.

**Key Application Features:**

### Interactive Text Analysis:

- Real-time clustering predictions
- ChatGPT-enhanced insights
- Sentiment visualization

### Data Exploration Tools:

- Cluster distribution charts
- Topic word clouds
- Article similarity networks

# 8. Discussion

## 8.1 Technical Achievements

Our team successfully:

1. **Developed specialized preprocessing** for Kazakh text handling
2. **Implemented a hybrid analysis system** combining traditional ML with modern LLMs
3. **Created an accessible interface** for users with varying technical backgrounds
4. **Established a reproducible pipeline** from raw data to actionable insights

## 8.2 Algorithm Selection Rationale

While our evaluation showed hierarchical clustering with 12 clusters achieved a slightly higher silhouette score (0.071) compared to K-Means with 10 clusters (0.024), we deliberately selected K-Means for our final implementation. This decision was based on several practical considerations:

1. **Computational Efficiency:** K-Means scales better to our dataset size

2. **Interpretability:** Cluster centroids in K-Means provide clearer topic representations
3. **Visualization Simplicity:** K-Means results are easier to visualize and explain
4. **Resource Constraints:** Lower memory and processing requirements

The marginal improvement in silhouette score from hierarchical clustering did not justify the additional complexity and resource requirements for our specific application.

### 8.3 Practical Application of the System

Our system is ready for the following applications:

- **Automatic news categorization:** Grouping articles by topic for quick overview
- **Sentiment analysis:** Determining positive, negative, or neutral tone in texts
- **Finding similar articles:** Locating relevant content based on a given sample
- **Data visualization:** Graphical representation of topic distribution and trends
- **Demonstration of ML methods:** A clear example of clustering algorithms for educational purposes

The system represents a working prototype that can be used for automated processing of Kazakh news and obtaining basic analytics on media content.

# 9. Conclusion and Future Work

### 9.1 Project Summary

We have successfully developed a working prototype of a Kazakh news analysis system based on machine learning methods. Our approach using TF-IDF and K-Means clustering has proven effective for automatic grouping of articles by topic. Integration with the ChatGPT API added sentiment analysis and text summarization capabilities. The system demonstrates practical applicability for basic media content analysis.

### 9.2 Development Plans

To improve the system, we plan to:

1. Test more powerful models (BERT, GPT-4) to enhance analysis quality
2. Integrate with news APIs for automatic data updates
3. Add real-time trend monitoring
4. Optimize performance for handling larger data volumes

5. Create an API for integration with other systems

## 10. References

https://scikit-learn.org/stable/modules/generated/sklearn.cluster.KMeans.html

https://scikit-learn.org/stable/modules/generated/sklearn.cluster.AgglomerativeClustering.html

https://scikit-learn.org/stable/modules/generated/sklearn.cluster.DBSCAN.html

https://platform.openai.com/docs/api-reference/backward-compatibility

https://scikit-learn.org/stable/modules/generated/sklearn.feature_extraction.text.TfidfVectorizer.html

https://scikit-learn.org/stable/modules/generated/sklearn.decomposition.PCA.html

https://www.kaggle.com/datasets/kuantaiulysalamat/kazakh-news-articles-dataset/data?select=tengri_news.csv