

# Unlocking the Power of Data: Enhancing Public Policy through Advanced Data Infrastructure and Language Model Analysis

Zahid Asghar, School of Economics, Quaid-i-Azam University, Islamabad, Pakistan

## Abstract

Data is the fundamental building block for advancements in artificial intelligence (AI), general AI (GAI), machine learning (ML), and large language models (LLMs). This study emphasizes the critical need for robust data infrastructure, arguing that without it, countries cannot fully benefit from technological advancements across various economic sectors. Governments possess vast repositories of both structured and unstructured data across domains such as the judiciary, parliaments, and civil bureaucracy. However, these potential goldmines remain largely untapped due to inadequate data management capabilities and a lack of appreciation for the necessity of high-quality data.

The research identifies key issues in public data management, including the non-uniform representation of key datasets and the prevalence of non-machine-readable formats, which complicate data utilization. By analyzing inconsistencies in standard data conventions within public datasets, this study underscores the challenges posed by messy data that require specialized skills to transform into tidy, consistently formatted datasets.

The objectives of this research are twofold: to explore the effective utilization of public policy data and to harness natural language processing (NLP) and LLMs to analyze critical policy documents, such as monetary policy statements issued by the State Bank of Pakistan. This study aims to demonstrate how large amounts of unstructured policy document data can be leveraged to analyze policy objectives, enhance public policy formulation and implementation, and realize the potential of data as a strategic asset in governance.

## Introduction

Data is the backbone of artificial intelligence (AI), general AI (GAI), machine learning (ML), and large language models (LLMs), playing a critical role in the advancement of these technologies. The advent of the 4th Industrial Revolution has positioned data as a pivotal resource, often likened to the “new currency” of the modern economy. Data fuels the development of AI and ML applications, which are reshaping industries and redefining competitive advantage on a global scale.

Governments worldwide hold vast amounts of data, particularly in the public sector, encompassing records from the judiciary, legislative bodies, and civil services. Despite this abundance, it is difficult to capitalize on these assets due to inadequate data infrastructure, lack of standardization, and insufficient appreciation of data’s strategic value. In countries like Pakistan, the public

sector's data infrastructure is often inadequate, hindering the realization of the full potential of data-driven technologies. The challenges include non-uniform data representation, data in non-machine-readable formats, and a general lack of data management practices.

Tools like natural language processing (NLP) and LLMs can be leveraged to analyze policy documents and enhance public policy formulation and implementation. However, the challenges posed by non-uniform data representation and non-machine-readable formats need to be addressed to unlock the power of data in governance. Moreover, 80% of data are unstructured, and the public sector holds a significant portion of this data. By effectively utilizing this data, governments can make informed decisions and improve public services. Nevertheless, the lack of appreciation for the importance of high-quality data and the absence of robust data management practices pose significant challenges to leveraging this data effectively.

Historically, the sources of national competitive advantage have evolved across different eras, reflecting shifts in global power dynamics. Civilizations gained dominance based on unique strengths, ranging from cultural influence and military prowess to technological advancements. For example, Ancient India was known for its profound knowledge and cultural richness, establishing it as a center of learning and philosophy. Similarly, the Roman Empire's expansion was propelled by its organized legions and technological innovations like catapults, demonstrating how strategic military capabilities can establish a dominant position. As history progressed, the Mongol Empire leveraged its mobility and trade networks, while the Ottoman Empire capitalized on heavy artillery and cannons. The British Empire expanded through colonization backed by naval superiority and gunpowder. Moving into the 20th century, the United States emerged as a global leader through a combination of economic strength and military power.

Today, the trend indicates that the next global superpower will be defined by its command over data. Nations capable of harnessing, analyzing, and leveraging data effectively will gain a competitive edge, signifying a shift from traditional military and economic power to digital and information dominance. This shift underscores the importance of data as a strategic asset, essential for innovation, economic growth, and national security.

- **Explore effective utilization of public policy data.**
- **Harness NLP and LLMs to analyze critical policy documents, specifically monetary policy statements issued by the State Bank of Pakistan.**
- **Demonstrate how text policy document data can be leveraged to analyze policy objectives.**
- **Enhance public policy formulation and implementation through data-driven insights.**
- **Realize the potential of data as a strategic asset in governance.**

The 4th Industrial Revolution is characterized by the integration of digital technologies into all areas of business and society. Data plays a central role in this transformation, enabling the development of AI and ML applications that can process vast amounts of information to generate insights and automate processes. Previous studies have highlighted the importance of data quality and accessibility in realizing the benefits of these technologies (Schwab, 2016).

**Challenges in Public Sector Data Management** The public sector often lags behind in data management practices compared to the private sector. Issues such as data silos, legacy systems, and lack of standardization hinder effective data utilization (Janssen et al., 2017). Furthermore, public datasets are frequently stored in non-machine-readable formats, complicating efforts to apply advanced analytical techniques (Ubaldi, 2013).

**2.3 Natural Language Processing in Policy Analysis** NLP has emerged as a powerful tool for analyzing large volumes of text data, such as policy documents and legal texts. Studies have demonstrated the potential of NLP in extracting insights from unstructured data, informing decision-making processes, and enhancing transparency (Young et al., 2018). The application of LLMs, such as GPT models, has further advanced the capabilities in understanding and generating human-like text, offering new avenues for policy analysis (Brown et al., 2020).

## **Methodology**

### **Data Collection**

The study focuses on monetary policy statements issued by the State Bank of Pakistan. These documents are rich in information regarding the country's economic outlook, policy objectives, and regulatory changes. The data collected includes all available monetary policy statements over the past decade, sourced directly from the State Bank's official website.

### **Data Preprocessing**

Given that the policy documents are unstructured text data, preprocessing is essential to prepare the data for analysis. The preprocessing steps include: Converting documents into machine-readable text formats. Removing irrelevant information such as headers, footers, and disclaimers. Normalizing text by converting to lowercase and removing punctuation. Tokenization to break down text into individual words or terms. Stop-word removal to eliminate common words that do not contribute to meaningful analysis.

**3.3 Analytical Techniques** The study employs NLP techniques and LLMs to analyze the monetary policy statements. The analytical approach includes: Sentiment Analysis: To gauge the tone of the policy statements and assess shifts in optimism or caution over time. Topic Modeling: Using algorithms like Latent Dirichlet Allocation (LDA) to identify underlying themes and topics discussed in the statements. Keyword Extraction: Identifying frequently used terms and phrases that signify policy focus areas. Temporal Analysis: Examining changes in language and emphasis over time to understand policy evolution.

**3.4 Tools and Software** The analysis is conducted using Python programming language, leveraging libraries such as: NLTK (Natural Language Toolkit): For basic NLP tasks. spaCy: For advanced NLP processing. Gensim: For topic modeling. Transformers (by Hugging Face): To utilize pre-trained LLMs for text analysis.

**4. Results**

**4.1 Sentiment Analysis** The sentiment analysis reveals fluctuations in the tone of monetary policy statements correlating with economic events. Periods of economic instability show a more cautious tone, while stable periods reflect a more optimistic outlook.

**4.2 Topic Modeling** Topic modeling identifies key themes such as inflation control, exchange rate policies, and financial sector reforms. Over the decade, there is a noticeable shift from focusing primarily on inflation to incorporating growth and stability objectives.

**4.3 Keyword Trends Analysis** Analysis of keyword frequency highlights terms like "inflation," "GDP growth," "interest rates," and "foreign

exchange reserves.” The prominence of these terms varies with economic conditions, indicating shifts in policy priorities. 4.4 Temporal Analysis Temporal analysis shows an evolution in policy emphasis, with recent statements placing greater importance on digital banking and financial inclusion, reflecting global trends in financial technology adoption. 5. Discussion 5.1 Implications for Policy Formulation The findings demonstrate how textual analysis of policy documents can provide valuable insights into the central bank’s priorities and responses to economic challenges. This approach can assist policymakers in identifying areas requiring attention and assessing the effectiveness of past policies. 5.2 Challenges in Data Utilization The study underscores the difficulties posed by unstructured and non-standardized data formats in the public sector. These challenges highlight the need for improving data management practices, adopting uniform data standards, and ensuring that public data is machine-readable. 5.3 Leveraging Data as a Strategic Asset By effectively utilizing existing data, governments can enhance transparency, improve decision-making, and foster trust among stakeholders. The application of advanced analytical tools can unlock the potential of data, transforming it into actionable intelligence that drives economic and social development.

## Conclusion

Data holds immense potential as a strategic asset in governance and economic development. This study illustrates the benefits of applying NLP and LLMs to analyze unstructured policy documents, providing insights that can enhance public policy formulation and implementation. To fully realize these benefits, it is imperative to address challenges related to data quality, standardization, and accessibility. Investing in data infrastructure, promoting data literacy, and fostering collaboration between government, academia, and industry are crucial steps toward harnessing the power of data. As the global economy increasingly relies on data-driven technologies, countries that prioritize data management and utilization will gain a competitive advantage in the 4th Industrial Revolution.

## References

Brown, T. B., Mann, B., Ryder, N., et al. (2020). Language Models are Few-Shot Learners. *Advances in Neural Information Processing Systems*, 33, 1877-1901. Janssen, M., Charalabidis, Y., & Zuiderwijk, A. (2017). Benefits, Adoption Barriers and Myths of Open Data and Open Government. *Information Systems Management*, 29(4), 258-268. Schwab, K. (2016). *The Fourth Industrial Revolution*. World Economic Forum. Ubaldi, B. (2013). *Open Government Data: Towards Empirical Analysis of Open Government Data Initiatives*. OECD Working Papers on Public Governance, No. 22. Young, A., Schenk, T., & Verhulst, S. (2018). The Potential of Signal-Driven and Collaborative Approaches to Address Wicked Problems: Lessons Learned from Data Collaboratives. *Data & Policy*, 1, e1.

Appendix A. Monetary Policy Simulator Portal The Monetary Policy Simulator (MPS) portal, available at <https://zahidasghar.com/mps/mpspk>, provides a practical example of how data-driven tools can support economic planning and analysis.