# Natural Language Processing Analysis of the Tata Steel Annual Report (FY 2023-24)

Rachit Jain | 70572200036 | L040

## Abstract

This report presents a comprehensive textual analysis of the Tata Steel Annual Report for the fiscal year 2023-24 using Natural Language Processing (NLP) techniques. The primary objective was to extract, process, and analyze the unstructured text from the 581-page PDF document to uncover key themes, sentiments, and underlying topics. The methodology involved data extraction using PyPDF2, preprocessing with NLTK, sentiment analysis via TextBlob, and topic modeling using Latent Dirichlet Allocation (LDA) with Gensim. Key findings reveal a predominantly neutral-to-positive sentiment throughout the report. The most frequent terms, such as "steel," "tata," "company," "financial," and "crore," highlight the report's core focus on corporate and financial performance. Finally, the LDA model successfully identified 10 distinct topics, including financial statements, corporate governance, production & operations, and sustainability efforts (emissions, safety). This project demonstrates the efficacy of NLP in transforming large corporate documents into structured, actionable insights.

## 1. Introduction

Corporate annual reports are dense, text-rich documents that provide a comprehensive overview of a company's performance, strategy, and outlook. However, their length and unstructured nature make manual analysis a challenging task. Natural Language Processing (NLP) offers a suite of powerful tools to automate the analysis of such large textual datasets, enabling the extraction of meaningful patterns and insights that might otherwise remain hidden.

This project aims to apply a standard NLP pipeline to the **Tata Steel Annual Report for fiscal 2023-24**. The goal is to dissect the report's content to understand its overall sentiment, identify the most significant keywords and themes, and discover the latent topics discussed throughout the document.

## 2. Methodology and Tools

The project was executed in a Jupyter Notebook environment using Python 3. The analysis followed a sequential pipeline, with each step building upon the previous one. The core tasks and the primary libraries used are outlined below:

- **PDF Text Extraction:** PyPDF2 was used to read the 581-page PDF file and extract the raw text from each page.
- **Data Structuring:** Pandas was employed to organize the extracted text into a structured DataFrame, mapping each page's content to its corresponding page number.
- **Text Preprocessing:** The NLTK (Natural Language Toolkit) library was used for fundamental cleaning tasks, including converting text to lowercase, tokenization (splitting text into sentences and words), and removing stopwords, punctuation, and digits.
- **Sentiment Analysis:** TextBlob was used to perform sentence-level sentiment analysis, calculating polarity (positive/negative tone) and subjectivity for each sentence.
- **Feature Engineering:** Scikit-learn's TfidfVectorizer converted the cleaned text into a TF-IDF (Term Frequency-Inverse Document Frequency) matrix, a numerical representation of word importance.
- **Topic Modeling:** Gensim was utilized to build a Latent Dirichlet Allocation (LDA) model to identify 10 underlying topics from the text corpus.
- **Data Visualization:** Matplotlib and WordCloud were used to create visual representations of the data, such as histograms and a word cloud, to aid in interpretation.

## 3. Implementation and Results

### 3.1. Data Extraction and Preprocessing

The text from the 581-page 'Tata Steel Annual Report 2023-24.pdf' was successfully extracted and loaded into a Pandas DataFrame. A custom preprocessing function was applied to clean the text. This involved:

1. Conversion to lowercase.
2. Removal of all punctuation and digits.
3. Elimination of common English stopwords (e.g., 'the', 'a', 'is').

This step resulted in a clean, normalized text corpus ready for analysis.

### 3.2. Sentiment Analysis

Sentiment analysis was performed on each sentence of the original text to gauge the overall tone of the report. Polarity scores, ranging from -1.0 (negative) to +1.0 (positive), were calculated.

The distribution of these scores is shown in the histogram below.

**Analysis:** The histogram is heavily concentrated around the 0.0 to 0.25 range, indicating a

predominantly neutral-to-slightly-positive tone. This is expected for a formal document like an annual report, which is typically objective and fact-based. The absence of significant negative polarity suggests the report focuses on achievements and positive framing.

## 3.3. Frequent Words and Word Cloud

After preprocessing, the frequency of each word in the entire document was calculated. The top 20 most frequent words were:

[('steel', 2857), ('tata', 2547), ('company', 2080), ('year', 1824), ('financial', 1534), ('crore', 1431), ('limited', 1345), ('march', 1266), ('fy', 978), ('report', 975), ('h', 947), ('value', 860), ('annual', 809), ('integrated', 784), ('accounts', 712), ('statements', 709), ('th', 683), ('india', 670), ('f', 669), ('assets', 635)]

A word cloud was generated to visualize these frequencies, with more frequent words appearing larger.

**Analysis:** The most frequent words clearly reflect the document's nature. Terms like "steel," "tata," "company," "limited," and "india" establish the subject, while "financial," "crore," "march," "fy" (Fiscal Year), and "year" confirm the focus on financial reporting within a specific timeframe.

## 3.4. TF-IDF Matrix

The processed text was converted into a TF-IDF matrix of shape (581, 1000). This means the content of the 581 pages was represented by the 1000 most important terms across the document. This matrix serves as the numerical input for the topic modeling algorithm.

## 3.5. Topic Modeling (Latent Dirichlet Allocation)

An LDA model was trained on the corpus to discover 10 latent topics. Each topic is a collection of related words. The results are as follows:
- **Topic #1:** limited, company, tata, steel, board, ltd, private, shares, mr, year
  - **Interpretation:** Corporate Governance and Board Structure.
- **Topic #2:** tata, company, steel, report, management, board, business, also, policy, integrated
  - **Interpretation:** Business Strategy and Management Reporting.
- **Topic #3:** crore, tata, steel, company, march, safety, h, limited, loan, increase
  - **Interpretation:** Financials and Safety Metrics.
- **Topic #4:** year, crore, march, h, f, financial, net, ended, assets, tax
  - **Interpretation:** Annual Financial Statements (likely Profit & Loss).
- **Topic #5:** crore, march, court, ü, year, company, amount, scheme, tata, financial
  - **Interpretation:** Legal Proceedings and Financial Schemes.
- **Topic #6:** steel, tata, fy, year, production, company, india, growth, also, products
  - **Interpretation:** Production, Sales, and Market Growth.

- **Topic #7:** steel, tata, emissions, fy, uk, scope, energy, water, year, steelmaking
  - **Interpretation:** Sustainability and Environmental Impact (especially UK operations).
- **Topic #8:** steel, tata, fy, employees, year, report, business, women, annual, human
  - **Interpretation:** Human Resources and Employee Welfare.
- **Topic #9:** fy, year, shares, steel, h, tata, total, employees, ordinary, due
  - **Interpretation:** Shareholding and Employee Data.
- **Topic #10:** financial, statements, value, march, assets, year, company, consolidated, report, crore
  - **Interpretation:** Consolidated Financial Reporting (likely Balance Sheet).

## 4. Discussion

The NLP analysis provided several key insights into the Tata Steel Annual Report:

1. **Objective Tone:** The sentiment analysis confirms the report's formal and objective tone, with a slight positive skew, likely from management discussions on performance and strategy.
2. **Financial Dominance:** The word frequency and topic modeling results overwhelmingly point to a strong emphasis on financial reporting. Topics 3, 4, 5, and 10 are almost entirely composed of financial and accounting terms ("crore," "march," "assets," "financial statements"). This is the primary function of an annual report.
3. **Key Strategic Areas:** Beyond finance, the topics reveal other key areas of focus for Tata Steel. Topic 7 clearly points to a discussion on **sustainability**, with terms like "emissions," "energy," and "water," and a specific mention of "UK" operations, likely related to the Port Talbot steelworks transition. Topic 8 highlights a focus on **Human Resources**, including "employees" and "women," suggesting discussions on diversity and workforce management.
4. **Operational Focus:** Topic 6 relates directly to the core business of **production and sales**, with words like "production," "products," and "growth."

## 5. Conclusion

This project successfully applied an NLP pipeline to analyze the Tata Steel Annual Report 2023-24. By transforming the large, unstructured PDF into processed data, it was possible to quantify the report's sentiment, identify its most important terms, and uncover its main thematic pillars through topic modeling.

The analysis confirms that while the report is heavily centered on mandatory financial disclosures, it also contains significant discussions on corporate governance, sustainability, and human resources, reflecting the company's key strategic priorities for the fiscal year. This demonstrates the value of NLP as a tool for efficiently extracting high-level business intelligence from complex corporate documents.