# Job Description Intelligence and Market Insights (2023–2024)

## Essentials of Data Science Project Report

*Team 19*

# Contents

# 1. Project Introduction

# 2. Introduction

### 2.1 Project Overview

The *Job Description Intelligence and Market Insights (2023–2024)* project transforms unstructured job description text into structured, actionable data using Natural Language Processing (NLP) and data analytics. By automatically extracting key attributes—job titles, skills, locations, salaries, and company names—from the LinkedIn Job Dataset (2023–2024), we aim to uncover critical job market trends and provide data-driven insights for career planning and recruitment strategies.

### 2.2 Objectives and Scope

The project delivers four core objectives:

- Develop an NLP model using spaCy and BERT to extract structured information from job descriptions
- Build analytical dashboards revealing salary trends, skill demands, and geographic patterns
- Create a functional prototype for real-time job description analysis
- Generate comprehensive market insights including top technologies, hiring hotspots, and compensation benchmarks

**In Scope:** Data cleaning, NLP model development, exploratory analysis, visualization, and system integration.
**Out of Scope:** Real-time web scraping and production deployment.

### 2.3 Team Structure and Technical Approach

Our four-member team brings complementary expertise:

- **Mohammed (Data Engineer):** Dataset preprocessing and quality assurance
- **Sanjeet (ML Engineer(NLP)):** Text extraction model development
- **Dheeraj (Data Analyst):** Market insights and visualization
- **Rachit (Software Engineer):** System integration and API development

The technical stack leverages Python 3.10+ with key libraries including Pandas, NumPy, spaCy, Transformers (BERT), Matplotlib, and Seaborn. Development occurs in Jupyter/Colab environments with GitHub version control and LaTeX documentation via Overleaf.

### 2.4 Project Timeline and Deliverables

The 8-week project follows this schedule:

- **Weeks 1–2:** Dataset acquisition and cleaning
- **Weeks 3–4:** NLP model development
- **Weeks 5–6:** Exploratory analysis and visualization
- **Weeks 7–8:** System integration and documentation

**Key Deliverables:** Cleaned dataset, NLP extraction model, visual market reports, functional prototype, and comprehensive documentation with GitHub repository.

### 2.5 Expected Impact

This project addresses the critical need for automated job market intelligence by providing:

- Real-time structured extraction from unstructured job postings
- Data-driven insights on emerging skills and compensation trends
- Geographic analysis of hiring patterns and salary distributions
- A practical tool for job seekers and recruiters to analyze opportunities

By combining advanced NLP techniques with robust data analytics, we deliver actionable intelligence for navigating the evolving 2023–2024 job market landscape.

## 3. Team Members and Contributions

| Member | Role | Key Contributions |
| --- | --- | --- |
| Mohammed | Data Engineer | Data cleaning, preprocessing, handling missing values. |
| Sanjeet | ML Engineer | Model development for text extraction using spaCy/BERT. |
| Dheeraj | Data Analyst | Exploratory data analysis, insight generation, and visualization. |
| Rachit | SW Engineer | Create API to integrate model, System integration of NLP and analysis modules |

# 4. Detailed Contributions

| Member | Role | Key Contributions |
|--------|------|-------------------|
| Mohammed | Data Engr | <ul><li>Data cleaning and preprocessing using Pandas and NumPy</li><li>Handle missing values and inconsistent formats in salary/location fields</li><li>Create ETL pipelines for data ingestion and transformation</li><li>Implement data validation and quality assurance checks</li><li>Standardize multi-format data into canonical representations</li><li>Document data schema and maintain data dictionary</li></ul> |
| Sanjeet | ML Engr(NLP) | <ul><li>Develop hybrid NLP model using spaCy and BERT</li><li>Extract job-specific entities (skills, experience, requirements)</li><li>Fine-tune transformer models for skill extraction tasks</li><li>Create skill taxonomy and classification system</li><li>Build text preprocessing pipeline for unstructured data</li><li>Implement model evaluation metrics and confidence scoring</li></ul> |

| Member | Role | Key Contributions |
|---|---|---|
| Dheeraj | Data Analyst | <ul><li>Perform exploratory data analysis on job market trends</li><li>Generate insights on salary, location, and skill patterns</li><li>Create visualizations using Matplotlib, Seaborn, and Plotly</li><li>Build interactive dashboards for trend analysis</li><li>Conduct statistical analysis on hiring patterns</li><li>Identify emerging technologies and in-demand skills</li></ul> |
| Rachit | Software Engr | <ul><li>Create RESTful API for model integration using FastAPI</li><li>Develop system architecture for NLP and analytics modules</li><li>Build Streamlit prototype for job description analysis</li><li>Implement batch processing and asynchronous operations</li><li>Set up Docker containerization for deployment</li><li>Establish testing framework and error handling systems</li></ul> |

# 5. Tools, Technologies, and Languages

| Category | Tools / Technologies |
|---|---|
| Programming Language | Python 3.10+ |
| Data Cleaning and Analysis | Pandas, NumPy |
| NLP Processing | spaCy, NLTK, Transformers (BERT) |
| Visualization | Matplotlib, Seaborn, Plotly |
| IDE / Platform | Jupyter Notebook / Google Colab |
| Version Control | Git, GitHub |
| Documentation | Overleaf (LaTeX) |

# 6. Dataset Information

**Dataset Used:** LinkedIn Job Dataset (2023–2024)

The dataset includes job titles, companies, locations, descriptions, salary estimates, and required skills.

**License:** Publicly available for educational and research use on Kaggle.

# 7. Progress Tracking and Verification

Project progress has been tracked weekly using an Excel tracker, documenting milestones such as:

- Dataset collection and cleaning
- NLP model development
- Visualization and insight generation
- Final integration and documentation

All tasks meet the project requirements and are ready for stakeholder review.

**GitHub Repository:** https://github.com/job-description-intelligence

The repository contains:

1. Updated project files (data cleaning, model, and visualization scripts)
2. This PDF report (exported from Overleaf)
3. The Excel progress tracker

# 8. Expected Outcomes

By the end of the project, the team expects to deliver:

- A system that extracts structured information from job descriptions
- Data-driven insights about job markets in 2023–2024
- Visual dashboards summarizing salary, skill, and location trends
- A complete data science workflow from data preprocessing to insight generation