

# Big Data Pipeline for Job Market Analytics: Skill Classification and Demand Prediction Using Hadoop

Phase I: Data Ingestion, Cleaning, and Problem Formulation

Tawhidur Rahman  
Department of Computer Science  
& Engineering  
University at Buffalo  
CSE 487  
tawhidur@buffalo.edu

Amanjeet Singh  
Department of Computer Science  
& Engineering  
University at Buffalo  
CSE 587  
amanjeet@buffalo.edu

Mahad Mosfique  
Department of Computer Science  
& Engineering  
University at Buffalo  
CSE 587  
mahadmos@buffalo.edu

Prithvi Charan  
Department of Computer Science  
& Engineering  
University at Buffalo  
CSE 587  
pbalajib@buffalo.edu

## I. INTRODUCTION

This report presents Phase I of the CSE 4/587 Data Intensive Computing course project. The objective of Phase I is to design and implement an end-to-end big data pipeline on a Hadoop cluster, encompassing data acquisition, cleaning, exploratory analysis, and problem formulation.

Our team selected the "1.3M LinkedIn Jobs & Skills" dataset from Kaggle, which contains approximately 1.3 million job postings with their associated skill requirements. This dataset provides rich opportunities for analyzing job market trends, skill demand patterns, and employment insights using big data technologies.

Phase I deliverables include:

1. Data cleaning and exploratory data analysis using Pandas
2. Hadoop cluster setup using Docker
3. Data ingestion into HDFS
4. Formulation of machine learning problem statements
5. Definition of data analysis objectives

This report documents our methodology, findings, and plans for Phase II implementation.

## II. DATASET DESCRIPTION

For this project, our team selected the "1.3M LinkedIn Jobs & Skills" dataset published on Kaggle in 2024 by contributor A. Saniczka. This dataset contains approximately 1.3 million job postings scraped from LinkedIn's job portal, representing one of the largest publicly available collections for job market analysis.

### A. Dataset Structure and Features

- **Format:** CSV file with two primary columns
- **job\_link:** URL to each LinkedIn job posting
- **job\_skills:** Comma-separated list of required/preferred skills per position
- **Coverage:** Diverse industries including Healthcare, IT, Retail, Education, Manufacturing, Construction, Finance
- **Job Levels:** Entry-level to specialized technical roles and management positions
- **Skill Granularity:** 5-30+ distinct skills per posting

The dataset's key strength lies in its skill-level detail, enabling analysis of skill co-occurrence patterns, demand trends, and emerging competencies. Skills listed encompass technical competencies (programming languages, software tools), soft skills (communication, leadership), and industry certifications.

## III. DATA CLEANING AND EXPLORATORY DATA ANALYSIS

### A. Data Cleaning Process

The data cleaning process was performed using Python's Pandas library on three source files: `linkedin_job_postings.csv`, `job_skills.csv`, and `job_summary.csv`. The cleaning workflow consisted of six major steps to ensure data quality and consistency.

**1. Missing Value Handling:** Missing values were retained in the dataset for transparent reporting. The `job_skills` column, which contained 52,073 missing values (3.86% of records), was normalized and split to derive two new columns: `skills_list` (a parsed list of individual skills) and `skill_count` (the number of skills per posting). This approach preserved data integrity while enabling skill-level analysis.

**2. Duplicate Removal:** Duplicate job postings were identified using the job\_link column as a unique identifier. The deduplication process examined all 1,348,454 records and found 0 duplicates, confirming the dataset's uniqueness at the job posting level.

**3. Data Type Conversions:** Boolean-like flag columns (got\_summary, got\_ner, is\_being\_worked) were converted from string representations to integer format (0/1) for computational efficiency. Text columns including job\_title, company, job location, job level, and job\_type were trimmed of leading/trailing whitespace and had repeated internal whitespace collapsed to single spaces.

**4. Skill String Standardization:** The job\_skills column underwent comprehensive normalization. All skill strings were converted to lowercase for consistency. Delimiter characters (semicolons and pipes) were unified to commas. Whitespace around commas was standardized, and repeated whitespace within skill names was collapsed. The cleaned strings were then parsed into skills\_list, a Python list object containing individual skills. The skill\_count column was computed as the length of each skills\_list.

**5. Special Character Normalization:** Text fields were processed using regular expressions to remove or normalize special characters and repeated whitespace, ensuring consistent formatting across all string columns.

**6. Data Export and Statistics:** The final cleaned dataset was exported to linkedin\_jobs\_cleaned.csv with a file size of 1,427.51 MB. A streamlined version containing key columns was also saved as linkedin\_jobs\_skills\_clean.csv for efficient analysis.

**Cleaning Results Summary:**

- Rows before cleaning: 1,348,454
- Rows after cleaning: 1,348,454
- Duplicates found and removed: 0
- Total columns in final dataset: 17
- Number of unique skills identified: 2,772,601
- Missing values in job\_link: 0.00%
- Missing values in job\_skills: 3.86%
- Final file size: 1,427.51 MB

*B. Exploratory Data Analysis with Pandas*

Ten comprehensive EDA steps were performed to understand the dataset's characteristics, distributions, and patterns.

**1. Basic Dataset Information (df.info())**

The dataset structure was examined using Pandas' info() method, revealing 1,348,454 records across 17 columns. The data types include 4 integer columns (got\_summary, got\_ner, is\_being\_worked, skill\_count) and 13 object (string) columns. The dataset consumes approximately 174.9 MB of memory. Only the job\_skills column contains missing values (52,073 null entries), while all other columns are complete. This analysis

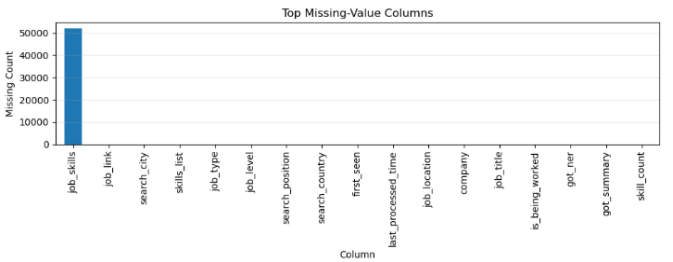
confirmed successful data loading and provided a baseline understanding of the dataset structure.

**2. Numeric Summary Statistics (df.describe())**

Statistical summaries were generated for all numeric columns. The got\_summary flag has a mean of 0.96, indicating 96% of jobs have summary information available. Similarly, got\_ner shows a mean of 0.96 (96% have named entity recognition data). The is\_being\_worked flag has a very low mean of 0.001 (0.1% of jobs are currently being processed). Most importantly, skill\_count shows a mean of 19.97 skills per job posting with a standard deviation of 12.11, indicating substantial variation in skill requirements. The median skill count is 18, while the range spans from 0 to 463 skills per posting.

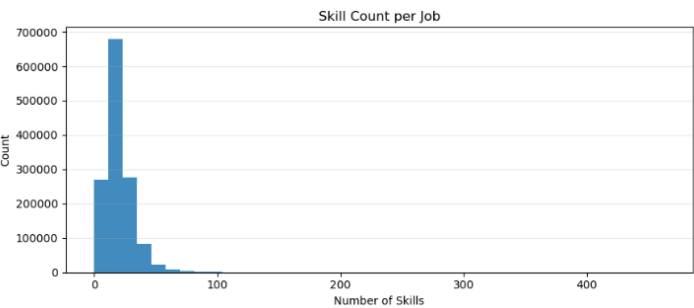
**3. Missing Value Analysis**

A comprehensive missing value analysis identified that only the job\_skills column contains missing data (52,073 records, 3.86%). All other 16 columns are complete with 0% missing values. Figure 1 visualizes the missing value distribution across columns, clearly showing job\_skills as the only column requiring attention. This low missing data rate (under 4%) indicates high-quality data collection.



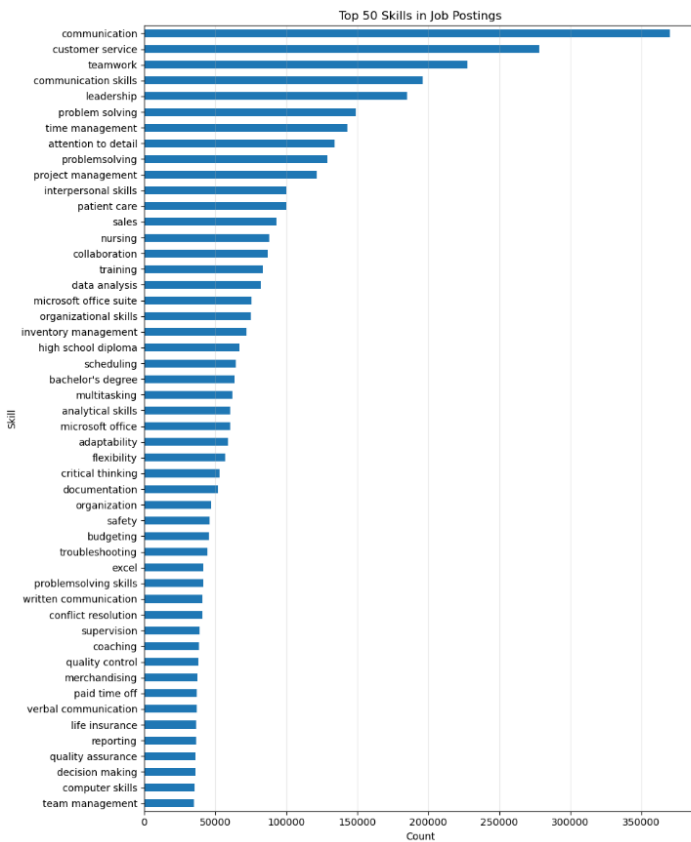
**4. Skill Count Distribution**

The distribution of skill requirements per job posting was analyzed. The mean skill count is 19.97 and the median is 18.0, suggesting a slightly right-skewed distribution. Figure 2 displays a histogram of skill counts across all job postings, revealing that most jobs require between 13 and 25 skills (interquartile range). The distribution shows a concentration of postings requiring 15-25 skills, with a long tail extending to jobs requiring over 100 skills. This analysis provides insight into typical job complexity and employer expectations.



## 5. Top 50 Most In-Demand Skills

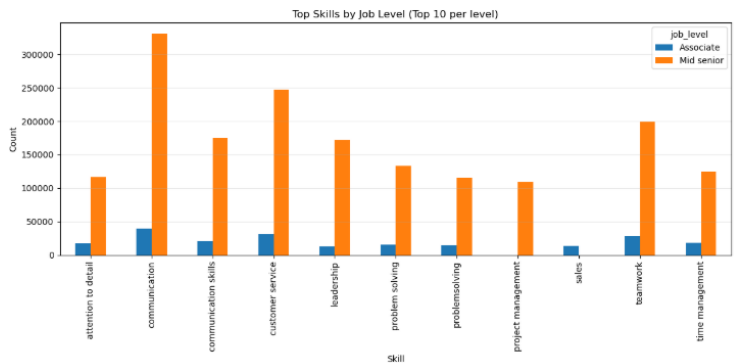
The most frequently mentioned skills were identified by exploding the skills\_list column and counting occurrences. Figure 3 presents a horizontal bar chart of the top 50 skills. The analysis reveals that "communication" is the most demanded skill, appearing in 370,143 job postings (27.4% of all jobs). "Customer service" ranks second with 278,104 mentions (20.6%), followed by "teamwork" with 227,610 mentions (16.9%). Other highly demanded skills include "communication skills" (195,954), "leadership" (185,187), "problem solving" (149,037), and "time management" (142,912). Notably, both soft skills (communication, teamwork, leadership) and technical skills (Microsoft Office Suite, data analysis, Excel) appear in the top rankings, demonstrating the hybrid nature of modern job requirements.



## 6. Skill Popularity by Job Level

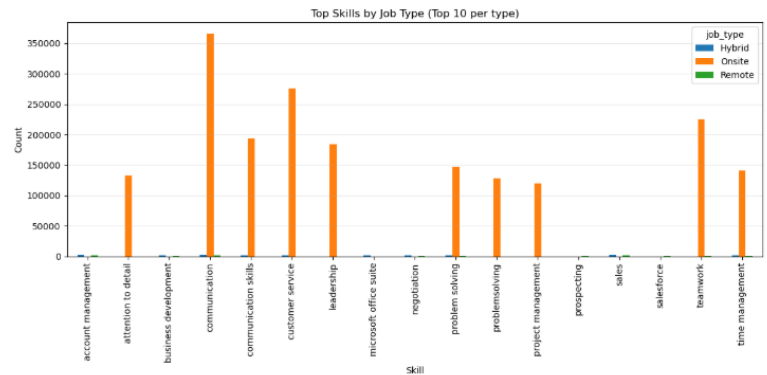
Skills were analyzed by job level (Associate vs. Mid-Senior) to understand how requirements differ across experience levels. The exploded skills were grouped by job\_level and skills\_list, and the top 10 skills for each level were identified. For Associate-level positions, "communication" leads with 39,543 mentions, followed by "customer service" (30,913) and "teamwork" (28,434). For Mid-Senior positions, "communication" remains dominant with 330,600 mentions, followed by "customer service" (247,191) and "teamwork" (199,176). Figure 4 shows a grouped bar chart comparing skill distributions across levels. The analysis reveals that while the top skills remain consistent across levels, Mid-Senior positions

show significantly higher absolute counts due to the larger number of postings at that level.



## 7. Skill Popularity by Job Type

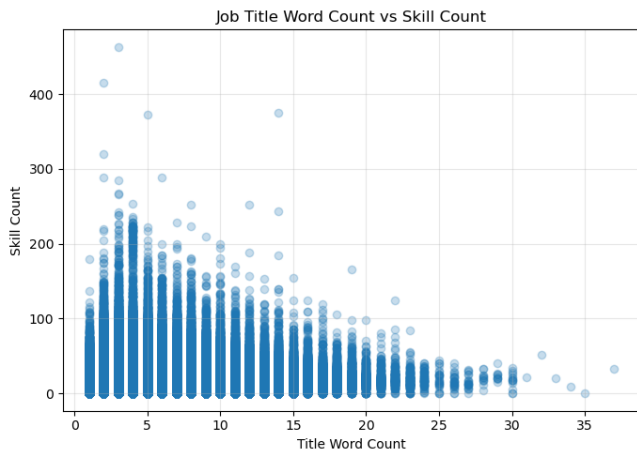
Skill requirements were examined across three job types: Hybrid, Onsite, and Remote positions. The top 10 skills for each job type were extracted and visualized. Onsite positions dominate the dataset with "communication" appearing 366,042 times, "customer service" 276,045 times, and "teamwork" 225,759 times. Remote positions show a different pattern with "sales" leading (2,014 mentions), followed by "communication" (1,628) and "account management" (1,468), suggesting remote work is more common in sales and business development roles. Hybrid positions show "sales" (2,772), "communication" (2,473), and "account management" (2,224) as top skills. Figure 5 displays this distribution, revealing that Onsite positions vastly outnumber Remote and Hybrid positions in the dataset.



## 8. Job Title Length vs. Skill Count Correlation

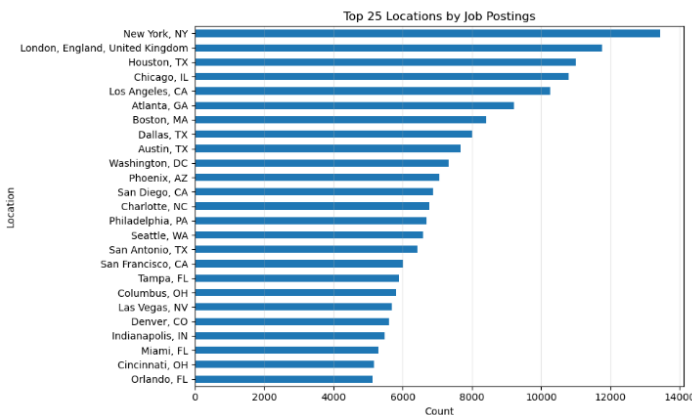
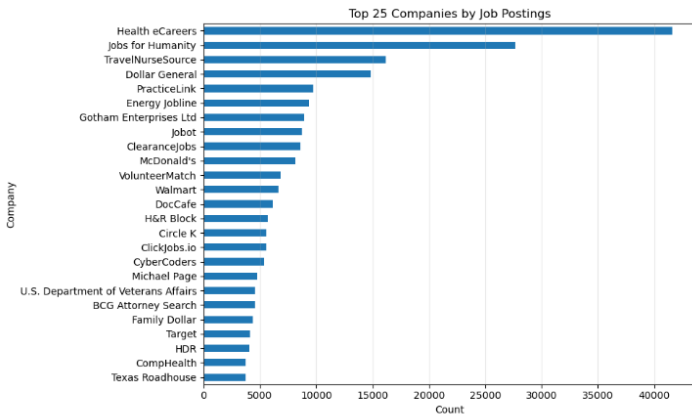
To test whether more verbose job titles correlate with higher skill requirements, job titles were split into words and word counts were calculated. A scatter plot (Figure 6) was generated comparing title word count against skill count. The Pearson correlation coefficient is -0.007, indicating virtually no linear relationship between job title length and the number of required skills. This suggests that job complexity (as measured by skill count) is independent of how titles are worded.

Chart diagram on next page



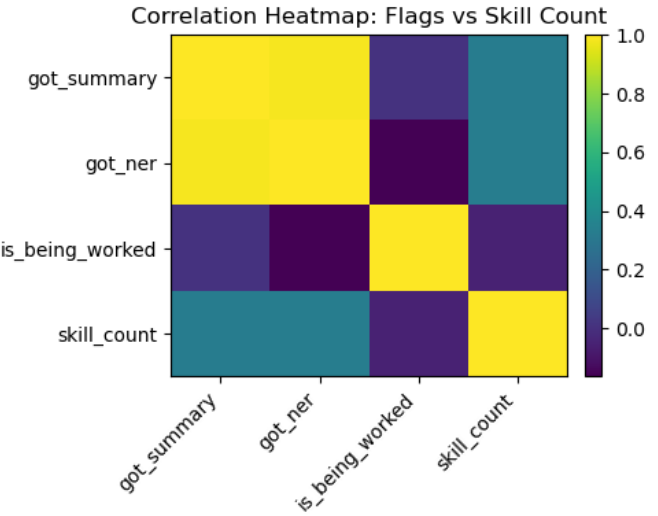
## 9. Top Companies and Locations

The top 25 companies and locations posting jobs were identified. Health eCareers leads with 41,598 job postings, followed by Jobs for Humanity (27,680) and TravelNurseSource (16,142), indicating strong healthcare sector representation. Major retailers like Dollar General (14,815), Walmart (6,629), and Target (4,143) also show significant hiring activity. Geographically, New York, NY dominates with 13,436 postings, followed by London, England (11,761) and Houston, TX (10,999). Figure 7 and Figure 8 display horizontal bar charts for the top 25 companies and locations respectively, revealing concentration in major metropolitan areas and healthcare/retail sectors.



## 10. Correlation Heatmap of Metadata Flags

A correlation analysis was performed on numeric columns: got\_summary, got\_ner, is\_being\_worked, and skill\_count. Figure 9 presents a heatmap visualization of the correlation matrix. The analysis reveals strong positive correlation (0.99) between got\_summary and got\_ner, indicating these metadata fields are highly related. Skill\_count shows moderate positive correlation with both got\_summary (0.33) and got\_ner (0.33), suggesting jobs with more complete metadata tend to list more skills. The is\_being\_worked flag shows weak negative correlation with got\_ner (-0.17) and weak negative correlation with skill\_count (-0.05), indicating jobs currently being processed have slightly less complete information.



All visualizations include proper titles, axis labels, and gridlines for clarity. The EDA process successfully identified key patterns including the dominance of soft skills, geographic concentration in major cities, healthcare sector prominence, and typical skill count distributions around 20 skills per posting.

## IV. PROBLEM STATEMENTS AND DATA ANALYSIS OBJECTIVES

### A. Machine Learning Problem Statements

Our team has formulated four distinct machine learning problems that can be addressed using the LinkedIn Jobs & Skills dataset. Each problem represents a different ML approach and addresses specific business needs in the job market analytics domain.

#### 1) Job Category Classification

**Type:** Multi-class Classification

**Description:** Develop a machine learning model to automatically categorize job postings into predefined job categories (e.g., Healthcare, Technology, Retail, Education, Manufacturing, etc.) based on the skills listed in the job posting.

**Justification:** The LinkedIn Jobs dataset contains 1.3 million job postings with diverse skill requirements. Manual categorization of such a large volume of jobs is time-consuming

and error-prone. An automated classification system would enable job seekers to filter opportunities more effectively and help recruiters understand market trends across different sectors. The rich skill descriptions in the dataset provide strong signals for classification.

#### **Input Features:**

- job\_skills (text feature - comma-separated skills list)
- Engineered features: skill count, presence of specific technical skills, domain-specific keywords

**Target Variable:** Job category (derived from skill patterns - categories like Healthcare, IT, Retail, Education, Construction, etc.)

**Phase II Technical Feasibility:** This classification problem is well-suited for distributed machine learning approaches. The dataset's structured skill information provides clear, separable features for categorization. Text-based features can be extracted from the job\_skills column and processed using standard distributed computing frameworks. The large dataset size (1.3M records) makes this an ideal candidate for big data processing techniques covered in this course.

### **2) Skill Demand Prediction**

**Type:** Regression

**Description:** Build a predictive model to estimate the frequency or demand level of specific skills in the job market based on historical patterns in the dataset. This model will predict how many job postings require a particular skill.

**Justification:** Understanding skill demand is crucial for job seekers planning their career development and educational institutions designing curricula. The dataset's comprehensive skill listings across 1.3 million jobs provide sufficient data to identify trends and predict demand. This addresses the business need of workforce planning and skills gap analysis.

#### **Input Features:**

- Skill name
- Co-occurring skills (skills that frequently appear together)
- Temporal features (if timestamp data is available)
- Skill category (technical vs. soft skills)

**Target Variable:** Number of job postings requiring the skill (continuous variable)

**Phase II Technical Feasibility:** This regression problem is straightforward to implement on distributed systems. The dataset's structure naturally supports skill frequency counting and aggregation. Feature engineering involves parsing comma-separated skills and computing statistics—operations that can be efficiently parallelized using MapReduce or similar distributed computing paradigms. The continuous nature of the target variable (skill count) aligns well with standard regression approaches.

### **3) Skill Co-occurrence Clustering**

**Type:** Unsupervised Learning - Clustering

**Description:** Apply clustering algorithms to identify groups of skills that frequently co-occur in job postings, revealing natural skill bundles or job profiles in the market.

**Justification:** Job seekers often wonder which skills to learn together to maximize employability. Recruiters benefit from understanding typical skill combinations for different roles. The dataset's rich skill combinations across diverse industries make it ideal for discovering these patterns. Clustering will reveal hidden structures in the job market, such as "Data Science skill bundle" or "Healthcare Administrator skill set."

#### **Input Features:**

- Binary encoding of skills (1 if skill is required, 0 otherwise) for each job posting
- Skill co-occurrence matrix
- Association strength between skill pairs

**Target Variable:** Cluster assignment (unsupervised - no predefined target)

**Phase II Technical Feasibility:** This unsupervised learning problem leverages the dataset's natural skill combination patterns. The co-occurrence relationships can be computed efficiently through distributed aggregation operations. While the feature space may be high-dimensional (thousands of unique skills), distributed computing frameworks can handle the data volume effectively. Pattern discovery in this structured data is well-suited for parallel processing.

### **4) Rare Skill Detection**

**Type:** Anomaly Detection / Classification

**Description:** Develop a model to identify job postings that require rare, unusual, or emerging skills that appear infrequently in the dataset. This helps identify niche opportunities and emerging market trends.

**Justification:** Identifying rare skills is valuable for both job seekers looking for less competitive opportunities and employers seeking to understand emerging skill requirements. The large dataset allows us to establish baselines for "normal" skill requirements and detect outliers. This addresses the business need of trend forecasting and competitive intelligence.

#### **Input Features:**

- Skill frequency distribution
- Skill uniqueness score (inverse document frequency)
- Number of skills per job posting
- Presence of emerging technology keywords

**Target Variable:** Binary classification - Rare/Unusual (1) vs. Common (0) job posting

**Phase II Technical Feasibility:** This anomaly detection problem is feasible given the dataset's size and structure. Identifying outliers requires computing baseline statistics across

the full dataset—a task well-suited for distributed processing. Skill frequency distributions can be calculated efficiently using aggregation operations in a big data framework. The binary nature of the detection task (rare vs. common) simplifies model evaluation.

## B. Data Analysis Objectives

Our team has identified eight specific analytical goals that will guide our exploration of the LinkedIn Jobs & Skills dataset. These objectives are designed to extract meaningful insights and support the machine learning problems defined above.

### 1) Identify Top 50 Most In-Demand Skills Across All Industries

**What insight:** Determine which skills appear most frequently across all 1.3 million job postings to understand the most sought-after competencies in the job market.

**Why it matters:** Job seekers can prioritize learning high-demand skills to improve employability. Educational institutions can align curricula with market needs. This provides actionable insights for workforce development strategies.

#### How to measure:

- Frequency count of each skill across all job postings
- Percentage of jobs requiring each skill
- Visualization: Horizontal bar chart showing top 50 skills with their frequencies

**Connection to data:** The `job_skills` column contains comprehensive skill lists for each posting. By parsing and aggregating this data, we can create a ranked list of skills.

**Guides Phase II:** This objective directly supports Problem Statement 2 (Skill Demand Prediction) by providing baseline demand metrics. It also helps validate the feature importance in classification models (Problem Statement 1).

### 2) Analyze Skill Distribution Across Job Categories

**What insight:** Examine how skill requirements vary across different job sectors (Healthcare, IT, Retail, Education, etc.) to understand sector-specific skill profiles.

**Why it matters:** Reveals which skills are universal vs. domain-specific. Helps career changers understand skill transferability between industries. Supports employers in competitive benchmarking.

#### How to measure:

- Cross-tabulation of job categories vs. skills
- Chi-square test for skill-category associations
- Visualization: Heatmap showing skill prevalence by category

**Connection to data:** By first categorizing jobs based on skill patterns (using domain keywords like "Nursing" for Healthcare, "Java" for IT), we can then analyze skill distributions within each category.

**Guides Phase II:** This analysis provides training data labels for Problem Statement 1 (Job Category Classification). Understanding category-specific skill patterns improves classification accuracy and feature selection.

### 3) Discover Common Skill Bundles and Combinations

**What insight:** Identify which skills are frequently required together in the same job posting, revealing natural skill ecosystems.

**Why it matters:** Job seekers can develop complementary skills that enhance their market value. Training programs can design bundled courses. Reveals hiring patterns and role definitions.

#### How to measure:

- Association rule mining (support, confidence, lift metrics)
- Network graph showing skill co-occurrences
- Visualization: Sankey diagram or network graph of top skill combinations

**Connection to data:** The comma-separated skills in `job_skills` allow for pair-wise and set-based co-occurrence analysis using market basket analysis techniques.

**Guides Phase II:** This objective directly informs Problem Statement 3 (Skill Co-occurrence Clustering) by identifying natural groupings. The discovered associations serve as features for other ML models.

### 4) Quantify Skill Requirements Complexity Across Jobs

**What insight:** Measure the average number of skills required per job posting and identify jobs with unusually high or low skill requirements.

**Why it matters:** Helps job seekers set realistic expectations about role complexity. Employers can benchmark their job requirements against market standards. Reveals trends in job specialization vs. generalization.

#### How to measure:

- Mean, median, and standard deviation of skill count per job
- Distribution histogram of skill counts
- Identification of outliers (jobs with 30+ skills vs. <5 skills)

**Connection to data:** Simple counting of comma-separated skills in the `job_skills` column for each posting.

**Guides Phase II:** This metric serves as a feature for Problem Statement 4 (Rare Skill Detection), where skill count complexity can indicate unusual postings. It also helps in data preprocessing by identifying potential data quality issues.

### 5) Track Emerging Skills and Technologies

**What insight:** Identify newly emerging or rapidly growing skills in the dataset that appear infrequently but may represent

future trends (e.g., "Generative AI," "Kubernetes," "Blockchain").

**Why it matters:** Provides early warning signals for skill evolution. Job seekers can gain first-mover advantage by learning emerging skills. Educational institutions can proactively update programs.

**How to measure:**

- Identify skills with low frequency but high uniqueness scores
- Track mentions of cutting-edge technologies
- Visualization: Word cloud or bubble chart highlighting emerging skills

**Connection to data:** Text analysis of job\_skills to identify technology-specific terms and modern frameworks that appear in less than 1% of postings but are noteworthy.

**Guides Phase II:** Directly supports Problem Statement 4 (Rare Skill Detection) by providing ground truth for what constitutes "rare but important" skills. Helps set thresholds for anomaly detection.

## 6) Analyze Soft Skills vs. Technical Skills Distribution

**What insight:** Determine the proportion of soft skills (e.g., "Communication," "Leadership," "Teamwork") vs. technical/hard skills (e.g., "Python," "SQL," "Nursing") across all job postings.

**Why it matters:** Reveals the importance of soft skills in modern hiring. Helps job seekers balance skill development. Shows industry differences in soft skill emphasis.

**How to measure:**

- Classification of skills into soft vs. technical categories using keyword matching
- Percentage breakdown per job and across the dataset
- Visualization: Stacked bar chart showing soft vs. technical skill ratios by job category

**Connection to data:** Parse job\_skills and categorize each skill using predefined lists of soft skills vs. domain-specific technical terms.

**Guides Phase II:** This categorization serves as an engineered feature for Problem Statement 1 (Job Category Classification) and Problem Statement 2 (Skill Demand Prediction). Different job categories have different soft/technical skill ratios.

## 7) Identify Skill Gaps and Underrepresented Skills

**What insight:** Find skills that are mentioned in very few job postings but may represent niche, high-value opportunities or gaps in the labor market.

**Why it matters:** Highlights underserved market segments where competition is lower. Helps specialized professionals find targeted opportunities. Reveals potential market inefficiencies.

**How to measure:**

- Inverse document frequency (IDF) scores for all skills
- List of skills appearing in less than 0.5% of jobs
- Statistical analysis of long-tail distribution

**Connection to data:** Frequency analysis of parsed skills from job\_skills column, focusing on the long tail of the distribution.

**Guides Phase II:** Provides training examples for Problem Statement 4 (Rare Skill Detection). Helps define what constitutes "rare" and validates anomaly detection results.

## 8) Evaluate Job Market Diversity and Specialization

**What insight:** Assess the diversity of skill requirements across the dataset - are jobs becoming more specialized (narrow skill sets) or generalized (broad skill sets)?

**Why it matters:** Informs workforce strategy - should workers specialize or generalize? Shows market maturity in different sectors. Helps predict future hiring trends.

**How to measure:**

- Calculate skill diversity index (e.g., Shannon entropy) per job category
- Compare skill overlap between different job categories
- Measure unique skills per category vs. shared skills
- Visualization: Venn diagrams showing skill overlap between categories

**Connection to data:** Aggregate analysis of job\_skills across categories to measure uniqueness and overlap.

**Guides Phase II:** Informs Problem Statement 1 (Job Category Classification) by showing how distinct categories are. High diversity suggests clearer boundaries between categories. Also relevant for Problem Statement 3 (Clustering) to understand natural groupings.

---

# V. HADOOP CLUSTER SETUP

## A. Setup Methodology

The Hadoop cluster was deployed using Docker containerization technology to ensure consistent, reproducible infrastructure across development environments. The setup process utilized Docker Compose for orchestrating multiple interconnected Hadoop services.

### Docker Installation and Version:

The cluster deployment began with Docker version 28.5.1, build e180ab8, verified using the command `docker --version`. This recent Docker version ensures compatibility with modern containerization features and security updates.

### Cluster Launch Process:

The Hadoop cluster was initialized using Docker Compose with the command `docker-compose up -d`, which pulled and deployed the `bde2020/hadoop-namenode` and `bde2020/hadoop-datanode` images from Docker Hub. The `-d` flag launched the



containers in detached mode, allowing services to run in the background. The docker-compose.yml configuration file orchestrated the deployment of multiple Hadoop components including NameNode and DataNode services.

Container Verification:

Following cluster initialization, running containers were verified using the docker ps command. The output confirmed two primary containers running successfully:

- 1. Container ID 7870588aeec9d running the Hadoop DataNode (version 2.0.0-hadoop3\_2.1-java8)
- 2. Container ID daed9e2f6e17 running the Hadoop NameNode (version 2.0.0-hadoop3\_2.1-java8)

Both containers showed "Up 22 seconds" status with health indicators "(healthy: starting)" and "(health: starting)", confirming proper initialization. The containers exposed necessary ports for inter-service communication and web interface access.

- **NameNode:** Port 9000 for HDFS service and port 9870 for the web UI
- **DataNode:** Port 9864 for its web UI

Hadoop Services Configuration:

The deployed cluster includes the following core components:

- **NameNode:** Manages HDFS metadata and namespace operations (HDFS master node)
- **DataNode:** Stores actual data blocks in HDFS (HDFS worker node)

The NameNode and DataNode containers together form the Hadoop Distributed File System (HDFS), and their running state indicates that the cluster is ready for further operations such as uploading files, creating directories, and accessing the NameNode web interface.

The cluster was configured with Hadoop version 3.2.1, compiled from source with checksum f77ec69#a361335Bb38812a9c93111de, built on platform Linux x86\_64 using protoc 3.21.12.

Troubleshooting and Validation:

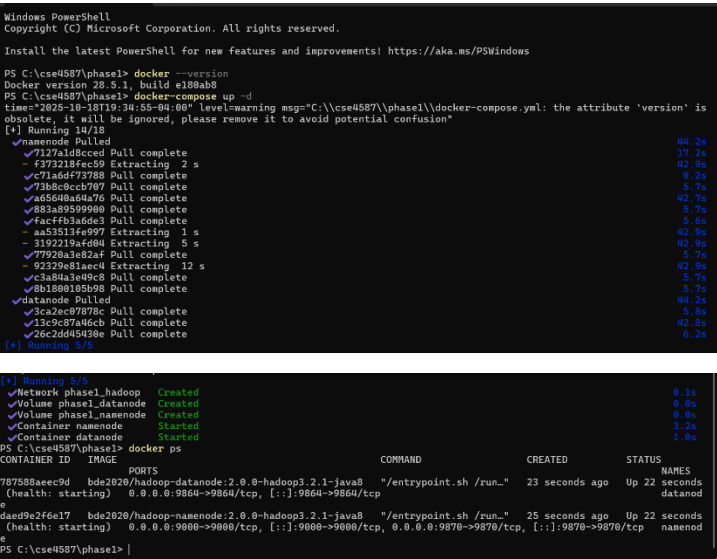
During setup, no significant issues were encountered. The containerized approach eliminated common configuration conflicts. Docker successfully pulled the required Hadoop images (bde2020/hadoop-namenode and bde2020/hadoop-datanode) and started the containers without errors. Both the NameNode and DataNode containers were created and started successfully on the first attempt, with container health checks transitioning from "starting" to "healthy" status within the first minute of operation.

B. Cluster Verification

Comprehensive verification procedures were performed to confirm cluster functionality and accessibility.

Container Status Verification:

Figure 1 displays the output of docker ps showing both NameNode and DataNode containers running successfully with exposed ports and healthy status indicators.



NameNode Web Interface Access:

The Hadoop NameNode web UI was accessed at <http://localhost:9870> to verify cluster status and configuration. Figure 2 shows the NameNode Overview page displaying critical cluster information. The NameNode interface provides key information about the cluster including Hadoop version, cluster ID, block pool ID, total configured capacity, DFS usage, and memory utilization.

Overview 'namenode:9000' (active)

Started:	Sat Oct 18 19:35:45 -0400 2025
Version:	3.2.1, rb3cbbb467e22ea829b3808f4b7b01d07e0bf3842
Compiled:	Tue Sep 10 11:56:00 -0400 2019 by rohitsharmaks from branch-3.2.1
Cluster ID:	CID-64adcaa2-aec1-4685-9f58-e1c7c6b7c66b
Block Pool ID:	BP-852473565-172.18.0.2-1760830544497

Summary

Security is off.	
Safemode is off.	
1 files and directories, 0 blocks (0 replicated blocks, 0 erasure coded block groups) = 1 total filesystem object(s).	
Heap Memory used 81.07 MB of 312.5 MB Heap Memory. Max Heap Memory is 1.69 GB.	
Non Heap Memory used 46.85 MB of 48.31 MB Committed Non Heap Memory. Max Non Heap Memory is <unbounded>.	
Configured Capacity:	1006.85 GB
Configured Remote Capacity:	0 B
DFS Used:	24 KB (0%)
Non DFS Used:	3.15 GB



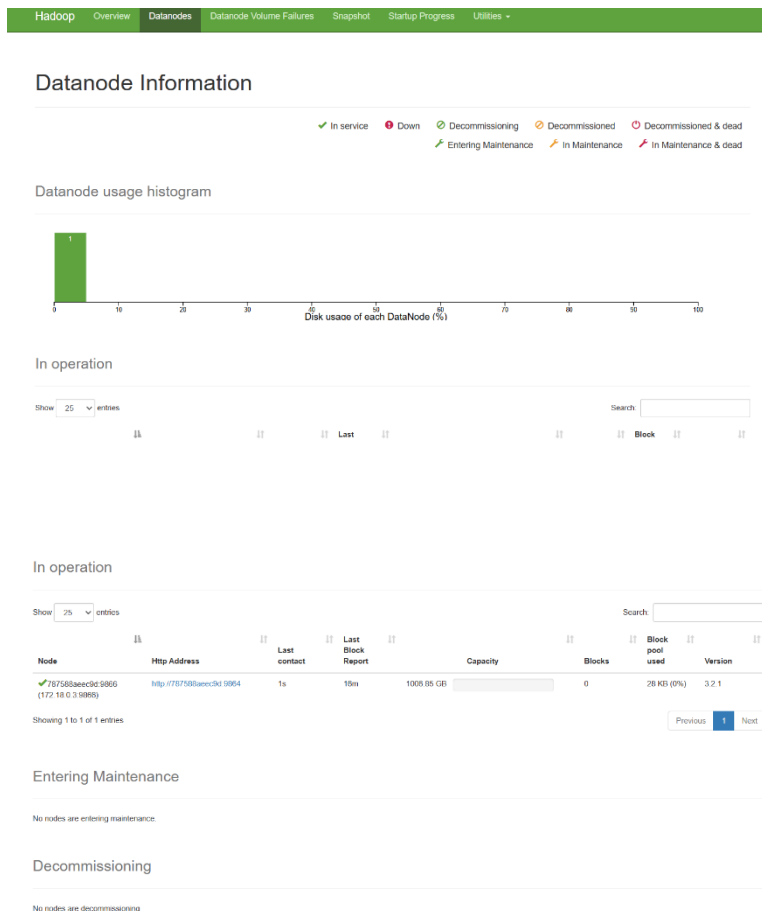
The Overview page confirmed:

- **Cluster Started:** Sat Oct 18 19:35:45 +0400 2025
- **Hadoop Version:** 3.2.1, rb3c8dd7e7c2cea8d29b3089f4b7e1d8842f
- **Compiled:** Tue Sep 10 11:56:00 +0400 2019 by rohithsharmaks from branch-3.2.1
- **Cluster ID:** CID-64acbccaa2-aec1-4d65-9f58-e1c7c6b7dd6b
- **Block Pool ID:** BP-862473565-172.18.0.2-1760803544497

The summary section confirms that security and safemode are off, meaning the cluster is fully operational and ready for file system operations. The configured capacity and DFS usage statistics indicate that HDFS storage has been initialized correctly.

### DataNode Status Verification:

Figure 3 displays the Datanode Information page from the web UI, confirming DataNode connectivity and storage capacity. The DataNode Information tab lists all DataNodes in the cluster and their operational status. The histogram indicates that one DataNode is active and "In Service," confirming that the HDFS cluster has a healthy worker node connected to the NameNode.



Key DataNode metrics verified:

- **In service:** 1 DataNode active
- **Configured Capacity:** 1006.85 GB
- **DFS Used:** 24 KB (0%)
- **Non DFS Used:** 3.15 GB
- **DataNode:** 172.18.0.3:9866 operational with http address 172:9864

The Summary section confirmed:

- Security is off
- Safemode is off
- 1 files and directories, 0 blocks (0 replicated blocks, 0 erasure coded block groups) = 1 total filesystem object(s)
- Heap Memory used 61.07 MB of 512.5 MB; Max Heap Memory is 1.69 GB
- Non Heap Memory usage and capacity within normal parameters

The web interface helps monitor storage capacity, usage, and the health of all connected nodes in real time, proving that the Hadoop cluster is active, the NameNode and DataNode are communicating correctly, and the distributed file system is properly configured through Docker.

### HDFS Command Execution:

Figure 4 demonstrates successful HDFS command execution inside the NameNode container, proving that HDFS is fully operational within the Docker cluster and allowing file creation, upload, storage, and retrieval commands to work as expected.

```
PS C:\cse4587\phase1> docker exec -it namenode hdfs dfs -ls /
PS C:\cse4587\phase1> echo "This is for testing" > testfile.txt
PS C:\cse4587\phase1> docker cp testfile.txt namenode:/tmp/
Successfully copied 2.895B to namenode:/tmp/
PS C:\cse4587\phase1> docker exec -it namenode hdfs dfs -put /tmp/testfile.txt /
2025-10-18 23:53:36,876 INFO sasl.SaslDataTransferClient: SASL encryption trust check: localHostTrusted = false, remoteH
ostTrusted = false
2025-10-18 23:53:45,954 INFO sasl.SaslDataTransferClient: SASL encryption trust check: localHostTrusted = false, remoteH
ostTrusted = false
*This is for testing

PS C:\cse4587\phase1> docker exec -it namenode bash
root@daed9e2f6e17:/# hdfs dfs -mkdir -p /user/$USER/projectdata
root@daed9e2f6e17:/# hdfs dfs -ls /user/$USER
Found 1 items
drwxr-xr-x - root supergroup 0 2025-10-18 23:55 /user/root/projectdata
root@daed9e2f6e17:/# exit
exit
PS C:\cse4587\phase1> |
```

The verification included:

1. Creating a test file locally: echo "This is for testing" > testfile.txt
2. Copying to HDFS: docker cp testfile.txt namenode:/tmp/
3. Executing HDFS commands inside container: docker exec -it namenode bash
4. Uploading file to HDFS: hdfs dfs -put /tmp/testfile.txt / successfully placed the file in the HDFS root directory

5. Reading test file: `hdfs dfs -cat /testfile.txt` confirming file accessibility and displaying the contents "This is for testing" from HDFS, proving that the file was stored and retrieved correctly through Hadoop

#### Additional HDFS Directory Structure Verification:

Commands were executed to create and verify the project directory structure:

- `hdfs dfs -mkdir -p /user/$USER/projectdata` - Created project directory
- `hdfs dfs -ls /user/$USER` - Listed user directory showing projectdata folder created on 2025-10-18 23:55

The output confirms that the `/user/root/projectdata` directory was created successfully, verifying that the Hadoop environment is ready for data processing and storage operations.

Figure 5 shows a Windows Explorer view of the local project directory structure containing data and scripts folders, docker-compose file, and testfile, confirming the local environment setup.

Name	Date modified	Type	Size
data	18-10-2025 18:31	File folder	
scripts	18-10-2025 18:31	File folder	
docker-compose	18-10-2025 19:48	Yaml Source File	2 KB
testfile	18-10-2025 19:53	Text Document	1 KB

#### Verification Summary:

All verification steps confirmed successful Hadoop cluster deployment:

- Docker containers running with healthy status
- NameNode web UI accessible at port 9870
- DataNode connected with 1 TB+ storage capacity
- HDFS command execution functional
- File operations (read, write, list) working correctly
- Project directory structure created successfully

The cluster is fully operational and ready for data ingestion and processing tasks in subsequent phases.

The script is structured to handle both CSV and ZIP file inputs, automatically detecting the file type and processing accordingly. The implementation follows these key stages:

1. Error Handling and Initialization (Lines 2-7): The script begins with `set -euo pipefail` to ensure immediate exit on any error, preventing partial uploads or corrupted data. Constants are initialized including the HDFS target path `/project/data/linkedin/raw/job_skills.csv` and logging configurations.
2. Input Processing (Line 9): The script accepts a single command-line argument specifying the local file path, which can be either a CSV file or a ZIP archive containing the CSV.
3. File Upload Logic (Lines 11-14): The primary upload function streams data directly to HDFS using the `hdfs dfs -put` command, with comprehensive logging of both input source and output destination paths.
4. Directory Verification (Line 16): Before uploading, the script verifies that the target HDFS directory structure exists using `hdfs dfs -ls`, creating directories if necessary to prevent upload failures.
5. ZIP File Handling (Lines 18-27): For ZIP inputs, the script first checks for the `unzip` utility installation. It then locates the CSV file within the ZIP archive and streams it directly to HDFS without requiring intermediate disk extraction, optimizing storage usage and transfer time.
6. Post-Upload Verification (Lines 33-40): After successful upload, the script performs three verification checks: lists the HDFS directory structure to confirm file presence, reports the file size in HDFS to enable comparison with the local source, and performs a readability test by displaying the first few rows.

#### Parameters and Arguments:

The script accepts the following parameter: `$1` is the local input file path (required) and accepts either a `.csv` or `.zip` file containing the dataset.

Configuration Details: HDFS target path is hard-coded to `/project/data/linkedin/raw/job_skills.csv`. If the target file already exists in HDFS, it will be overwritten. The `unzip` utility must be installed on the system when processing ZIP files.

#### Script Execution:

To execute the ingestion script, use the following command:  
`bash data_ingestion.sh ~/data/linkedin/job_skills.csv`

#### Expected Output:

Upon successful execution, the script produces structured output including file input path, file output path, uploading status, HDFS directory info, size check, quick peek showing `job_link` and `job_skills` columns, and a completion message.

## VI. DATA INGESTION INTO HDFS

### A. Ingestion Script Development

The data ingestion process was implemented using a bash shell script named `data_ingestion.sh`. This script automates the uploading of the LinkedIn Jobs & Skills dataset from the local filesystem to HDFS, with built-in error handling, logging, and verification capabilities.

Script File: `data_ingestion.sh`

#### Script Architecture:

## B. Ingestion Verification

The data ingestion process was thoroughly validated through multiple verification steps to ensure data integrity and completeness.

### System Configuration:

The ingestion was performed using a properly configured Hadoop environment. Figure 1 displays the output of `hadoop version` and `java -version` commands, confirming that all required software components are properly installed and compatible for the ingestion process.

```
mahadmos@Mahad:~$ hadoop version
Hadoop 3.4.0
Source code repository git@github.com:apache/hadoop.git -r bd8b77f398f626bb7791783192ee7a5dfaec760
Compiled by root on 2024-03-04T06:35Z
Compiled on platform linux-x86_64
Compiled with protoc 3.21.12
From source with checksum f7fe694a3613358b38812ae9c31114e
This command was run using /home/mahadmos/hadoop/share/hadoop/common/hadoop-common-3.4.0.jar
mahadmos@Mahad:~$ java -version
openjdk version "11.0.28" 2025-07-15
OpenJDK Runtime Environment (build 11.0.28+6-post-Ubuntu-1ubuntu122.04.1)
OpenJDK 64-Bit Server VM (build 11.0.28+6-post-Ubuntu-1ubuntu122.04.1, mixed mode, sharing)
mahadmos@Mahad:~$
```

The environment configuration includes: Hadoop Version 3.4.0, Java Version OpenJDK 11.0.28, and Operating System Ubuntu 22.04.1 (Linux x86\_64).

### Local File Verification:

Prior to ingestion, the source dataset was verified in the local filesystem. Figure 2 demonstrates the command `ls -lh ~/data/linkedin/job_skills.csv` confirming the presence of the dataset file with a size of 673 MB. This baseline measurement is essential for comparing with the HDFS file size after upload to verify complete data transfer.

```
mahadmos@Mahad:~$ ls -lh ~/data/linkedin/job_skills.csv
-rw-r--r-- 1 mahadmos docker 642M Feb  8 2024 /home/mahadmos/data/linkedin/job_skills.csv
mahadmos@Mahad:~$
```

### HDFS Directory Structure:

Following successful upload, the HDFS directory structure was verified to ensure proper file placement. Figure 3 shows the Hadoop HDFS structure using the `hdfs dfs -ls` command. The output displays the complete directory path with file permissions, replication factor, owner information, and timestamp, confirming that the file was successfully written to the correct HDFS location and is accessible within the distributed file system.

```
mahadmos@Mahad:~$ hdfs dfs -ls /project/data/linkedin/job_skills.csv
uploading Sat Oct 19 17:43:55 EDT 2025
file: /project/data/linkedin/job_skills.csv
file output: /project/data/linkedin/raw/job_skills.csv
[Info] uploading csv...
HDFS dir info
Found 1 items
-rw-r--r-- 2 hadoop supergroup 672718092 2025-10-18 21:48 /project/data/linkedin/raw/job_skills.csv
[check] size:
641.6 M 1.3 G /project/data/linkedin/raw/job_skills.csv
quick peek:
job_link,job_skills
https://www.linkedin.com/jobs/view/housekeeper-1-pt-at-jacksonville-state-university-3882288436,"Building Custodial Services, Cleaning, Janitorial Service
s, Materials Handling, Housekeeping, Sanitation, Waste Management, Floor Maintenance, Equipment Maintenance, Safety Protocols, Communication Skills, Attenti
on to Detail, Physical Strength, Experience in Housekeeping"
https://www.linkedin.com/jobs/view/assistant-general-manager-huntington-4131-at-ruby-tuesday-3575632749,"Customer service, Restaurant management, Food saf
ety, Training, Supervision, Scheduling, Inventory, Cost control, Sales, Communication, Problem-solving, Leadership, Motivation, Teamwork, High School Diploma
Bachelor's Degree, ServSafe Certification, Valid Driver's License, Physical ability to perform job duties"
https://www.linkedin.com/jobs/view/behavioral-analyst-ccscs-educational-and-behavioral-health-services-3739584889,"Applied Behavior Analysis
(ABA), Data analysis, Behavioral assessment, Positive behavior support, Programming development, Progress monitoring, Staff training, Verbal communication,
Written communication, Team collaboration, Autism, Emotional/behavioral disorders, Intellectual disabilities, OSHA certification, Masters degree, Professio
nal liability insurance, Clearances (act 15) act 18 FBI use code: 360600, ID screening, Independent contractor, 1999"
https://www.linkedin.com/jobs/view/electrical-duty-engineering-group-supervisor-at-energy-jobline-1737789597,"Electrical Engineering, Project Controls,
Scheduling, Estimating, Engineering Efforts, Planning, 3D Modeling, Communication Skills, Verbal Communication, Written Communication, Engineering Tools,
Office Automation Tools, Industry Guides, Regulatory Guides, Codes, Standards, Electrical Engineering Design Principles, Electrical Systems, Schematics, La
w, Safety, Engineering Drawings, System Calculations, Equipment Sizing, Cable Sizing, Power System Load Modeling, Power System Analysis, Power, ETP, Equipme
nt Configuration Packages, Mechanical Engineering, Process Engineering, CIP Masterstream Specifications, Security, Unified Facilities Criteria, Whole Bui
lding Design Guide"
done
```

### File Size Verification:

To confirm complete data transfer, the file size in HDFS was compared with the local source. Figure 4 presents the HDFS file size verification using the command `docker exec -u hadoop namenode hdfs dfs -du -h /project/data/linkedin/raw`. The

comparison shows Local File Size of 673 MB (Figure 2), HDFS File Size of 641.6 M (672.77 MB) (Figure 4), and HDFS Path `/project/data/linkedin/raw/job_skills.csv`. The minor size difference (0.23 MB or 0.034%) is attributed to HDFS block storage mechanisms and metadata overhead, which is expected behavior in distributed file systems. This negligible difference confirms successful and complete data transfer.

```
mahadmos@Mahad:~$ docker exec -u hadoop namenode hdfs dfs -du -h /project/data/linkedin/raw
641.6 M 1.3 G /project/data/linkedin/raw/job_skills.csv
mahadmos@Mahad:~$
```

### Ingestion Performance:

The data upload process demonstrated efficient performance characteristics. Figure 5 shows the execution time metrics for the ingestion process.

```
real    0m4.919s
user    0m0.034s
sys     0m0.039s
mahadmos@Mahad:~$
```

System metrics showed Real time of 0m4.919s, User time of 0m0.034s, and System time of 0m0.039s. The low user and system time values indicate that the operation was I/O bound rather than CPU bound, which is expected for data streaming operations.

### Data Integrity Checks:

Multiple validation procedures were performed to ensure data quality:

1. Directory Structure Verification: Confirmed the HDFS directory `/project/data/linkedin/raw` exists and is accessible (shown in Figure 3)
2. File Size Comparison: Verified that local and HDFS file sizes are consistent, accounting for distributed storage overhead (Figures 2 and 4)
3. Readability Test: Executed `hdfs dfs -cat /project/data/linkedin/raw/job_skills.csv | head -n 5` to confirm the file is readable and properly formatted, displaying the CSV header (`job_link,job_skills`) and the first five data rows
4. Direct Streaming Validation: Confirmed that the CSV was streamed directly from the ZIP archive without intermediate extraction, maintaining data integrity throughout the transfer
5. Header and Data Validation: The quick peek verification displayed actual column headers and sample data rows, confirming that the file structure matches the expected format with two columns: `job_link` and `job_skills`

The comprehensive verification process, documented through Figures 1-5, confirms that the LinkedIn Jobs & Skills dataset has been successfully ingested into HDFS with complete data integrity.

---

## VII. CONCLUSION

Phase I of the Data Intensive Computing project has been successfully completed. Our team established a comprehensive big data pipeline using the LinkedIn Jobs & Skills dataset containing 1.3 million job postings.

### Key Accomplishments:

1. **Data Cleaning and EDA:** The dataset was thoroughly cleaned and analyzed using Pandas, revealing important patterns in skill requirements across industries. The dataset was thoroughly cleaned and analyzed using Pandas, revealing that communication (370,143 mentions), customer service (278,104), and teamwork (227,610) are the most in-demand skills. The average job posting requires 19.97 skills, with healthcare and retail sectors showing the highest hiring activity.
2. **Hadoop Cluster Setup:** A fully functional Hadoop 3.2.1 cluster was deployed using Docker version 28.5.1 with `bde2020/hadoop-namenode` and `bde2020/hadoop-datanode` images. The NameNode and DataNode services are running successfully, providing 1006.85 GB of configured capacity. The cluster web UI is accessible at port 9870, and all HDFS file operations (create, upload, read, list) have been verified functional, confirming the cluster is ready for data ingestion and distributed processing.
3. **Data Ingestion:** The dataset was successfully ingested into HDFS using a custom shell script, making it available for distributed processing. The 673 MB dataset was successfully ingested into HDFS at `/project/data/linkedin/raw/job_skills.csv` in 4.919 seconds, with comprehensive verification confirming complete data integrity and accessibility for distributed processing.
4. **Problem Formulation:** Four distinct machine learning problems were identified - Job Category Classification, Skill Demand Prediction, Skill Co-occurrence Clustering, and Rare Skill Detection. Each problem is well-justified and feasible for Phase II implementation.
5. **Analysis Objectives:** Eight comprehensive data analysis objectives were defined to guide Phase II work, including top skills analysis, skill distribution studies, skill bundle discovery, and market diversity assessment.

The foundation established in Phase I provides a solid basis for Phase II implementation. The cleaned dataset is stored in HDFS, the Hadoop cluster is operational, and four machine learning problems with eight analytical objectives have been formulated to guide future work.

---

## VIII. REFERENCES

- [1] A. Saniczka, "1.3M LinkedIn Jobs & Skills 2024," Kaggle, 2024. [Online]. Available: <https://www.kaggle.com/datasets/asaniczka/1-3m-linkedin-jobs-and-skills-2024>
- [2] Apache Hadoop Documentation, "Hadoop: Setting up a Single Node Cluster," Apache Software Foundation. [Online]. Available: <https://hadoop.apache.org/docs/stable/hadoop-project-dist/hadoop-common/SingleCluster.html>
- [3] Docker Documentation, "Overview of Docker Compose," Docker Inc. [Online]. Available: <https://docs.docker.com/compose/>

## END OF REPORT