

LinkedIn Job Market Insights

Siddharth Hiraou
UBID: shiraou

Aishvarya Samuel Salvi
UBID: aishvary

I. Introduction:

The job market is constantly evolving, with new roles emerging and skill requirements changing rapidly. Analyzing employment patterns can reveal important information about worker demands as companies adjust to changing economic conditions and technology progress. One such domain where data analysis is particularly impactful is the job market. The need for particular talents and employment responsibilities fluctuates quickly as industries change and new technologies are developed. In this industry, it is crucial for companies, policymakers, and job seekers to comprehend these developments. The 1.3M LinkedIn Jobs & Skills dataset, which includes over 1.3 million job advertisements extracted from LinkedIn, provides a comprehensive picture of the current labor market. For this project, we'll be building an end-to-end big data pipeline using Hadoop and Spark to analyze this dataset. The tasks performed are:

- Perform exploratory data analysis (EDA) using Pandas to understand key patterns and trends.
- Set up a local Hadoop cluster using Docker and verify its components.
- Develop a data ingestion script to store our dataset in HDFS for further processing.

II. Problem Statement:

The following ML problems are addressed using our dataset:

1. **Classification of Job Title:**
The objective of this task is to classify job postings into predefined categories based on the job title, description, and required skills. By combining text data from job descriptions and skills, we can build a robust multi-class classification model to accurately categorize job postings. The expected outcome is a model that automatically classifies job postings into relevant categories, helping job seekers and employers streamline job searches and recruitment processes.

2. **Skill Recommendation system:** The objective of this task is to recommend a list of skills for a given job title or description using the skills data from the dataset. We can create a recommendation system that makes skill suggestions for particular positions by examining the relationship between job titles/descriptions and talents. A system that aids companies in crafting precise and thorough job descriptions and helps job seekers comprehend what skills are required for particular occupations is the expected outcome.
3. **Predicting the Location of Jobs:** The objective of this task is to predict the location of a job posting (city or country) based on the job title, company, and skills required. We can forecast the likely locations of particular employment roles by looking at trends in job names, employers, and skill sets. An approach that helps job searchers concentrate their employment searches geographically and gives employers insight into regional hiring trends is the expected result from this.

III. Data Analysis Objectives:

1. **Identify Top Recruiting Companies:** The objective of this analysis is to identify which companies are posting the most job listings and their geographic distribution. Job seekers can more effectively focus their applications and gain insights into the most active industries or regions by knowing which top recruiting organizations are situated in a given area. Employers looking to comprehend competitive hiring environments and job seekers seeking to narrow down their search can both benefit from these helpful insights.

2. **Evolution of Job Market Trends:**

The objective of this analysis is to explore how job postings and required skills have evolved over time. Examining how job descriptions and necessary competencies have changed over time is the aim of this investigation. Finding these

trends enables stakeholders, including employers, schools, and job seekers, to comprehend how the labor market is evolving and predict demands in the future. We want to create visualizations (such as line charts and area charts) that show seasonal patterns, long-term trends, and new skill demands by monitoring the quantity of job posts and skill requirements over time. These insights can assist stakeholders make well-informed decisions and give a clearer picture of how the labor market has changed.

3. Check for Skills Gaps in the Workforce:

The objective of this analysis is to highlight skills that are in high demand but are not frequently listed in job postings, indicating potential skill gaps. Such insights will guide efforts to bridge skill gaps and align training programs with market needs. The goal of this analysis is to identify skills that are in high demand but are not frequently listed in job postings, indicating potential skill gaps. Finding these gaps is important for educators and training programs, as it helps them design curricula that address industry needs and prepare job seekers for future opportunities.

4. Distribution of Job Opportunities Across various Locations:

The objective of this analysis is to explore the geographic distribution of job postings to identify hotspots (regions with high job activity) and underserved regions (areas with fewer opportunities). Knowing these patterns is essential for employers as it gives them information about possible growth areas and for job seekers as it enables them to select places with more chances. Stakeholders will use these insights to inform data-driven choices on where to concentrate their efforts.

5. Job Duration:

Job seekers can identify positions that are expected to remain open for a longer period of time by knowing the posting duration, which also helps companies streamline their hiring procedures. Each job posting's duration is determined using the `first_seen` and `last_processed_time` fields, and the average duration per job title, business, or industry is examined. In order to give businesses and job searchers clear insights into how long job posts normally remain

active, the findings are displayed using box plots or histograms, which aid in decision-making.

6. Is Using `hdfs dfs -put` a Good Way to Write Files to an HDFS Cluster?

The `hdfs dfs -put` command is a fairly simple, and straightforward way to write files out to an HDFS cluster and should only be used, for simple adhoc or manual uploads. However when dealing with the ingestion of significant amounts of data, this may not be the best option, due to network constraints, or NameNode performance. In our project where we are processing a dataset with a size of 5GB, we are more inclined to set up some automatic ingestion through the use of the Hadoop FileSystem API in Java. This allows us to get more elegant error handling, means of supporting integration into other workflows we want to do processing, and scalability. Automated ingestion is also able to be tuned for parallel writes, compression, and conversion to various formats, and can also be more suitable in production.

7. Are You Satisfied with the Speed of Writing Using `hdfs dfs -put`? How Can You Improve It?

Although the command `'hdfs dfs -put'` provides a simple way to put, it may not be fast enough for large files (like the 5GB dataset we will work with in this project). There are many factors creating the slowness, such as network latency (the time it takes to send data over the network), disk I/O (which can slow things down), and overloading the NameNode with many files in a single call. To improve performance, you work with `'distcp'`, a command designated to copy files in parallel; use a more efficient file format such as Parquet or ORC; use different HDFS block sizes to read/write simultaneously and gain better performance. At scale, integrating data ingestion with Apache Kafka or Flume may improve scalability and speed for data where the pipeline is continuously submitting updating records.

8. What Format Should You Use to Store Data in HDFS for Easier Analysis?

To facilitate storage and further analysis in HDFS, a structured, columnar format like Parquet or ORC is preferable than storing raw CSV. During this remainder of the project, as we need to handle a relatively small 5GB dataset, a Parquet format will have considerable benefits in terms of query performance and storage compression, while also lowering storage costs by bringing faster access to relevant data. A CSV file is row-based, meaning it is inefficient for large-scale analytics. By saving a

file in Parquet, it optimizes for columnar reads. With columnar formats, there is an additional performance boost when using Apache Spark or Hive, as many distributed processing frameworks are optimized for analytical workloads, which increases overall efficiency.

IV. Visualization:

We find out which job titles are in high demand, which firms are hiring, what talents employers are seeking, and how the labor market is evolving over time by using exploratory data analysis (EDA) and visualizations. We've created clear and intriguing visualizations using tools like Seaborn, Matplotlib, and Folium to make these insights simple to comprehend. These results give comprehensive view of the current employment market so that job seeker, employer, or educator, may make better decisions.

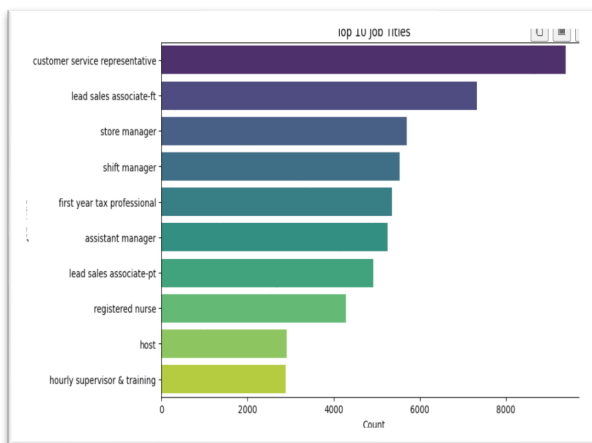


Figure: Top 10 Job Titles

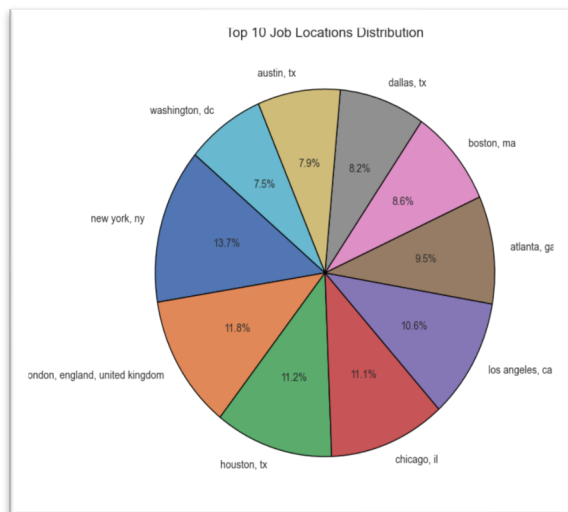


Figure: Top 10 Job locations

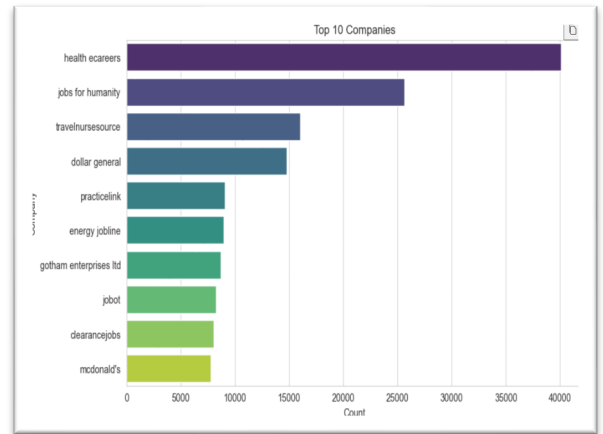
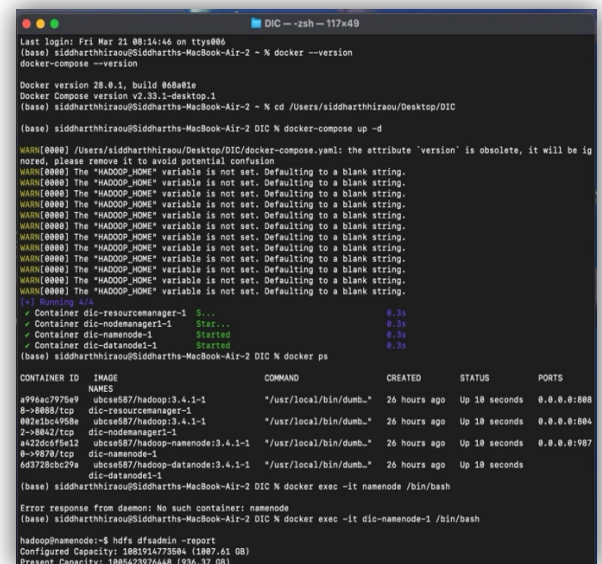


Figure: Top 10 Companies



Figure: Word cloud



CONTAINER ID	IMAGE NAMES	COMMAND	CREATED	STATUS	PORTS
a8e96797ee7	ucsbcs87/hadoop:3.1.1-1	"usr/local/bin/dumb."	26 hours ago	Up 10 seconds	0.0.0.0:888
8-e968d7f1	ucsbcs87/hadoop:3.1.1-1	"usr/local/bin/dumb."	26 hours ago	Up 10 seconds	0.0.0.0:888
e02be1da98ba	ucsbcs87/hadoop:3.1.1-1	"usr/local/bin/dumb."	26 hours ago	Up 10 seconds	0.0.0.0:888
> x9ac84c0	dio-namenode-1	"usr/local/bin/dumb."	26 hours ago	Up 10 seconds	0.0.0.0:888
x262efaf12	ucsbcs87/hadoop-namenode:4.1.1-1	"usr/local/bin/dumb."	26 hours ago	Up 10 seconds	0.0.0.0:888
-e967bf7c	dio-namenode-1	"usr/local/bin/dumb."	26 hours ago	Up 10 seconds	0.0.0.0:888
b3572bcb2f9	ucsbcs87/hadoop-datanode:3.1.1-1	"usr/local/bin/dumb."	26 hours ago	Up 10 seconds	0.0.0.0:888
(base) sidhatthirao@sidhatthirs-MacBook-Air:~\$	dio-datanode-1	DIO N docker exec -it namenode /bin/bash			
Error response from daemon: No such container: namenode					
(base) sidhatthirao@sidhatthirs-MacBook-Air:~\$	dio-datanode-1	DIO N docker exec -it dio-namenode-1 /bin/bash			
hadoopnamenode>	hdfs fsadmin -report				
Configured Capacity:	188154277584 (1807.61 GB)				
Report Capacity:	945939744 (9.17 GB)				

```
Cache Used: 0 (0 B)
Cache Remaining: 0 (0 B)
Cache Used: 100.00%
Cache Remaining: 0.00%
Xceiver: 0
Last contact: Sat Mar 22 14:10:45 UTC 2025
Last Block Report: Sat Mar 22 14:08:36 UTC 2025
Num of Blocks: 45

hadoop@namenode:~$ hdfs dfs -mkdir /test_directory
hdfs dfs -ls /
mkdir: /test_directory: File exists
Found 2 items
drwxr-xr-x - hadoop supergroup 0 2025-03-21 12:31 /test_directory
drwxr-xr-x - hadoop supergroup 0 2025-03-21 12:40 /user
hadoop@namenode:~$ exit
(base) siddharthhiraou@siddharths-MacBook-Air-2: ~ % docker cp /Users/siddharthhiraou/Desktop/DIC/processed_dataset.csv dic-namenode-1:/processed_dataset.csv
Successfully copied 5.9308 to dic-namenode-1:/processed_dataset.csv
(base) siddharthhiraou@siddharths-MacBook-Air-2: ~ % docker exec -it dic-namenode-1 /bin/bash
hadoop@namenode:~$ ls -l /processed_dataset.csv
-rw-r--r-- 1 501 dialout 5932914417 Mar 21 12:38 /processed_dataset.csv
hadoop@namenode:~$ hdfs dfs -mkdir -p /user/hadoop/datasets
hdfs dfs -put /processed_dataset.csv /user/hadoop/datasets/
put: /user/hadoop/datasets/processed_dataset.csv: File exists
hadoop@namenode:~$ hdfs dfs -ls /user/hadoop/datasets
Found 1 item
-rw-r--r-- 1 hadoop supergroup 5932914417 2025-03-21 12:58 /user/hadoop/datasets/processed_dataset.csv
hadoop@namenode:~$ python3 ingest_data.py
python3: can't open file '/opt/hadoop/ingest_data.py': [Errno 2] No such file or directory
hadoop@namenode:~$ exit
(base) siddharthhiraou@siddharths-MacBook-Air-2: ~ % python3 ingest_data.py
Copying file to NameNode container...
Successfully copied 5.9308 to dic-namenode-1:/processed_dataset.csv
Creating HDFS directory (if not exists)...
Removing existing file from HDFS (if exists)...
Deleted /user/hadoop/datasets/processed_dataset.csv
Uploading file to HDFS...
Verifying file in HDFS...
Found 1 item
-rw-r--r-- 1 hadoop supergroup 5932914417 2025-03-22 14:10 /user/hadoop/datasets/processed_dataset.csv
Data ingestion completed successfully
(base) siddharthhiraou@siddharths-MacBook-Air-2: ~ %
```

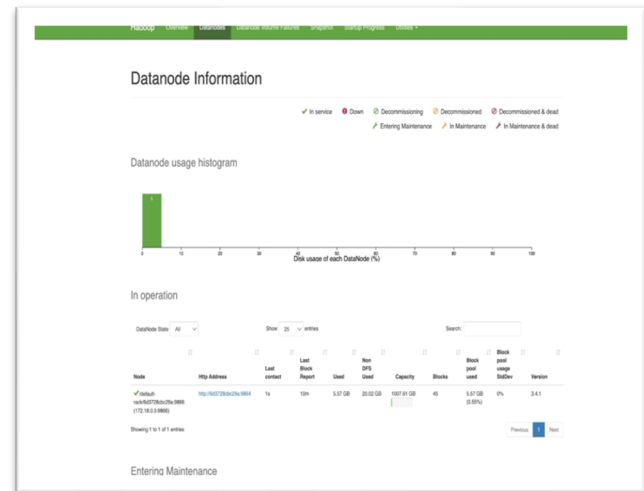


Figure: Datanode 1

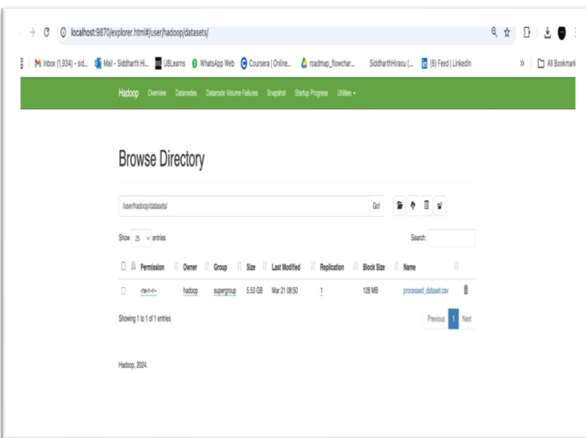


Figure: Browse Directory

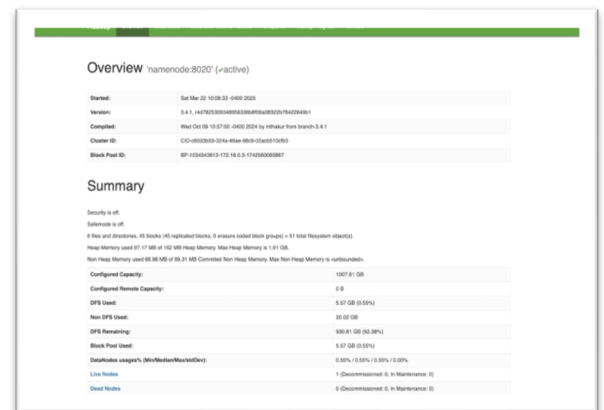


Figure: Datanode 2

```
hadoop@namenode:~$ hdfs dfadmin -report
Configured Capacity: 100.00 GB (100.00 GB)
Present Capacity: 100.00 GB (100.00 GB)
DFS Remaining: 94.03 GB (94.03 GB)
DFS Used: 5.97 GB (5.97 GB)
DFS Used: 0.59%
Replicated Blocks:
Under replicated blocks: 0
Blocks with corrupt replicas: 0
Missing blocks: 0
Missing blocks (with replication factor 1): 0
Low redundancy blocks with highest priority to recover: 0
Pending deletion blocks: 0
Erasure Coded Block Groups:
Low redundancy block groups: 0
Block groups with corrupt internal blocks: 0
Missing block groups: 0
Low redundancy blocks with highest priority to recover: 0
Pending deletion blocks: 0

Live datanodes (1):
Name: 172.18.0.5:8020 (dic-datanode-1.dic.default)
Hostname: 6d9728b2c29a
Decommission Status: Normal
Configured Capacity: 100.00 GB (100.00 GB)
DFS Used: 5.97 GB (5.97 GB)
Non DFS Used: 20.00 GB (20.00 GB)
DFS Remaining: 94.03 GB (94.03 GB)
DFS Used: 0.59%
DFS Remaining: 94.03 GB (94.03%)
Configured Cache Capacity: 0 (0 B)
Cache Used: 0 (0 B)
Cache Remaining: 0.00%
Xceiver: 0
Last contact: Sat Mar 22 14:10:45 UTC 2025
Last Block Report: Sat Mar 22 14:08:36 UTC 2025
Num of Blocks: 45

hadoop@namenode:~$ hdfs dfs -mkdir /test_directory
hdfs dfs -ls /
mkdir: /test_directory: File exists
Found 2 items
drwxr-xr-x - hadoop supergroup 0 2025-03-21 12:31 /test_directory
```

V. Hadoop Cluster Running:

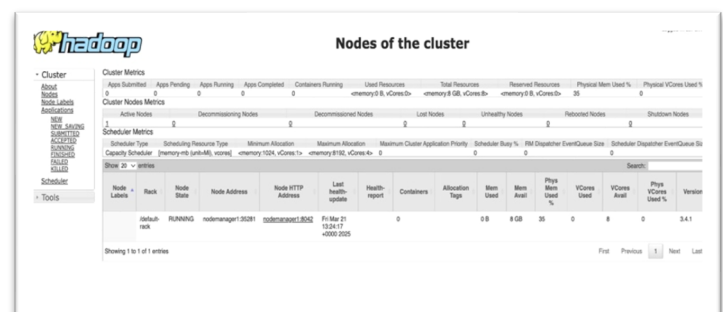


Figure: Cluster Running

VI. References:

1. <https://github.com/UBCSE587/2025Spring-projectphase1>
2. <https://hadoop.apache.org/docs/current/hadoop-project-dist/hadoop-common/FileSystemShell.html>
3. <https://github.com/spinlud/py-linkedin-jobs-scraper>
4. <https://github.com/spinlud/linkedin-jobs-scraper>