

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

Phase II: Machine Learning Implementation and Analysis

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 II of the CSE 4/587 Data Intensive Computing course project, building upon the foundation established in Phase I. In Phase I, our team selected the "1.3M LinkedIn Jobs & Skills" dataset from Kaggle, deployed a Hadoop cluster using Docker, ingested data into HDFS, and formulated four machine learning problems with eight analytical objectives.

Phase II focuses on implementing machine learning solutions using PySpark to address the identified problems. Unlike Phase I which utilized Hadoop HDFS for distributed storage, Phase II employs a local instance of PySpark for distributed data processing, machine learning implementation, and advanced analytics. This transition enables us to leverage PySpark's MLlib library for scalable machine learning while maintaining the distributed computing paradigm essential for processing large datasets.

Phase II deliverables include:

1. Data cleaning reimplemented using PySpark distributed operations
2. Exploratory data analysis using PySpark with comprehensive visualizations
3. Implementation of four machine learning algorithms addressing classification, regression, clustering, and anomaly detection problems
4. Comprehensive evaluation of model performance using appropriate metrics
5. Analysis of results, key findings, and recommendations for stakeholders

This report documents our implementation methodology, presents model performance results across all four problems, discusses key insights discovered through analysis, and provides recommendations for job seekers, educational institutions, and employers based on our findings. The transition from Hadoop to

PySpark demonstrates proficiency in multiple big data frameworks while maintaining analytical rigor throughout the data pipeline.

II. DATA CLEANING USING PYSPARK

The data cleaning process was reimplemented using PySpark to leverage distributed computing capabilities for processing the 1.3 million job postings. This section details the comprehensive cleaning pipeline developed to ensure data quality and prepare the dataset for exploratory analysis and machine learning implementation.

A. PySpark Environment Setup

The data cleaning process began with initialization of a local PySpark session configured to utilize all available CPU cores for parallel processing. The SparkSession was created with the application name "Phase2_Test" and configured to run in local mode with dynamic resource allocation.

```
from pyspark.sql import SparkSession
spark = (
    SparkSession.builder
    .master("local[*]")
    .appName("Phase2_Test")
    .getOrCreate()
)
```

Python

The SparkSession serves as the entry point for all PySpark operations, providing the foundation for distributed DataFrame operations. The master configuration "local[*]" enables utilization of all available processor cores, maximizing parallel processing efficiency during data cleaning operations.

B. Data Loading

Three primary CSV files were loaded into PySpark DataFrames using the spark.read API with schema inference enabled:

- linkedin_job_postings.csv** - Core job posting metadata including titles, companies, locations, and job characteristics
- job_skills.csv** - Skill requirements associated with each job posting via job_link foreign key
- job_summary.csv** - Additional summary information and named entity recognition data

```

# In[1]: import functions as f
jobs_path = "linkedin_job_postings.csv"
skills_path = "job_skills.csv"
summ_path = "job_summary.csv"

jobs_raw = [
    spark.read
        .option("header", True)
        .option("inferSchema", True)
        .csv(jobs_path)
]

skills_raw = [
    spark.read
        .option("header", True)
        .option("inferSchema", True)
        .csv.skills_path
]

summ_raw = [
    spark.read
        .option("header", True)
        .option("inferSchema", True)
        .csv.summ_path
]

def trim_col(df):
    return df.select([f.trim(c).alias(c) for c in df.columns])

jobs_cleaned = trim_col(jobs_raw)
skills_cleaned = trim_col.skills_raw
summ_cleaned = trim_col.summ_raw

print("Jobs raw count: ", jobs_raw.count())
print("Skills raw count: ", skills_raw.count())
print("Summary raw count: ", summ_raw.count())

```

The initial loading revealed:

- **jobs_raw**: 1,348,488 rows
- **skills_raw**: 1,296,381 rows
- **summ_raw**: 48,219,735 rows

A custom trim_cols() function was applied to all DataFrames to remove leading and trailing whitespace from column names, ensuring consistency across the dataset. The jobs and skills DataFrames were then joined on the job_link column using a left join to create a unified dataset (jobs_sk) containing job metadata and associated skills.

C. Cleaning Operations

The cleaning pipeline consisted of five major operations designed to address data quality issues and prepare features for analysis.

1. Missing Value Handling and Duplicate Removal

Following the left join operation, comprehensive missing value analysis was performed to assess data completeness:

- **Total rows (before deduplication)**: 1,348,488
- **Missing job_link**: 0 records (0.0%)
- **Missing job_skills**: 54,089 records (4.0129%)

Duplicate records were identified and removed using the job_link column as a unique identifier. The deduplication process examined all records and identified 25 duplicate job postings, which were subsequently removed using the dropDuplicates() method.

Cleaning Results:

- **Rows before cleaning**: 1,348,488
- **Rows after cleaning**: 1,348,463
- **Duplicates removed**: 25

```

Total rows (before dedup): 1348488
pct_missing_job_link: 0.0
pct_missing_job_skills: 4.0129
Rows before cleaning: 1348488
Rows after cleaning: 1348463
Duplicates found and removed: 25

```

Missing values in the job_skills column were retained rather than dropped to preserve data integrity. These null values were later converted to empty strings during skill parsing to enable consistent processing across all records.

2. Data Type Conversions

Boolean-like flag columns were converted from string representations to integer format (0/1) for computational efficiency and consistent processing. Three flag columns underwent conversion:

- **got_summary** - Indicates availability of job summary information
- **got_ner** - Indicates availability of named entity recognition data
- **is_being_worked** - Indicates whether job is currently being processed

The conversion logic handled multiple string representations of boolean values:

- Truthy values: ["true", "t", "yes", "y", "1"] → 1
- Falsy values: ["false", "f", "no", "n", "0"] → 0
- All other values → 0 (default)

This approach ensured robust handling of inconsistent boolean representations in the source data.

3. Text Normalization

Six text columns underwent comprehensive normalization to ensure consistency:

- job_link
- job_title
- company
- job_location
- job_level
- job_type

For each column, the following transformations were applied:

1. Conversion to string type (cast as string)
2. Replacement of multiple consecutive whitespace characters with single spaces using regex pattern \s+
3. Trimming of leading and trailing whitespace using F.trim()

```

from pyspark.sql import *
from pyspark import SparkContext
sc = SparkContext("local", "Job Cleaning Pipeline")
# Load raw dataset
raw_data = sc.textFile("hdfs://localhost:9000/jobs.csv")
# Create DataFrame
df = raw_data.map(lambda line: line.split(","))
df = df.toDF(["job_link", "job_title", "company", "job_location", "job_level", "got_summary", "got_ner", "is_being_worked"])
# Drop null rows
df = df.filter(df.job_link != "")
# Convert columns to strings
df = df.withColumn("job_link", df.job_link.cast("string"))
df = df.withColumn("job_title", df.job_title.cast("string"))
df = df.withColumn("company", df.company.cast("string"))
df = df.withColumn("job_location", df.job_location.cast("string"))
df = df.withColumn("job_level", df.job_level.cast("string"))
df = df.withColumn("got_summary", df.got_summary.cast("string"))
df = df.withColumn("got_ner", df.got_ner.cast("string"))
df = df.withColumn("is_being_worked", df.is_being_worked.cast("string"))
# Drop null values
df = df.na.drop()
# Print type conversions and count cleaned rows
print("Type conversions and text cleanup:")
df.printSchema()
df.count()

```

type conversions and text cleanup.

This normalization process eliminated formatting inconsistencies that could impact analysis quality, such as variations in whitespace that would cause identical values to be treated as distinct.

4. Skill String Standardization

The job_skills column underwent the most complex cleaning process due to its critical importance for analysis. A new column job_skills_clean was created through sequential transformations:

- Null Handling:** Null values were converted to empty strings using F.when() conditional logic
- Case Normalization:** All skill strings were converted to lowercase for consistency
- Delimiter Unification:** Semicolons and pipe characters (;|) were replaced with commas to standardize list separators
- Whitespace Standardization:** Spaces around commas were removed (pattern $\backslash s^*, \backslash s^* \rightarrow ,$) and repeated internal whitespace was collapsed (pattern $\backslash s^+ \rightarrow ,$)

Following string normalization, the cleaned skill strings were parsed into structured arrays:

- Array Creation:** The F.split() function divided comma-separated strings into arrays (skills_arr_raw)
- Array Refinement:** A combination of transform() and filter() functions removed empty strings and trimmed individual skill names:

```
F.expr("filter(transform(skills_arr_raw, x > trim(x)), x > x <> "")")
```

- Skill Count Calculation:** The F.size() function computed the number of skills per posting (skill_count column)

```

jobs_sk = jobs_sk.withColumn(
    "job_skills_clean",
    F.when(F.col("job_skills").isNotNull(), F.col("job_skills").cast("string"))
    .otherwise(F.lit("")))
jobs_sk = jobs_sk.withColumn(
    "job_skills_clean",
    F.regexp_replace("job_skills_clean", "[{}]", ""))
jobs_sk = jobs_sk.withColumn(
    "job_skills_clean",
    F.regexp_replace("job_skills_clean", "[\r\n]", ""))
jobs_sk = jobs_sk.withColumn(
    "job_skills_clean",
    F.regexp_replace("job_skills_clean", "[\t]", ""))
jobs_sk = jobs_sk.withColumn(
    "job_skills_clean",
    F.regexp_replace("job_skills_clean", "[\n]", ""))
jobs_sk = jobs_sk.withColumn(
    "skills_arr_raw", F.size("job_skills_clean"))
jobs_sk = jobs_sk.withColumn(
    "skills_list", F.array(F.map(F.transform("skills_arr_raw", x > trim(x)), x > x <> ",")))
jobs_sk = jobs_sk.drop("skills_arr_raw")

```

This parsing approach leveraged PySpark's higher-order functions for efficient array manipulation while maintaining data integrity throughout the transformation process.

5. Final Normalization Pass

A final normalization pass was applied to ensure all text columns maintained consistent formatting:

Columns processed: job_link, job_title, company, job_location, job_level, job_type, job_skills_clean

For each column:

- Multiple whitespace characters replaced with single spaces
- Leading and trailing whitespace trimmed

This final pass ensured that any whitespace introduced during previous transformations was normalized, creating a fully consistent dataset.

D. Cleaning Results and Export

The final cleaned dataset contained:

- Total rows:** 1,348,463
- Total columns:** 18
- Missing values in job_link:** 0.0%
- Missing values in job_skills:** 4.0129%
- Execution time:** 20 seconds

```

from pyspark.sql import functions as F;
final_rows = rows_after;
final_cols = len(jobs_sk.columns);

print("Total number of rows after cleaning:", final_rows);
print("Total number of columns after cleaning:", final_cols);
print("Percentage of missing values in job_link column:", round(pct_missing_job_link, 4));
print("Percentage of missing values in job_skills column:", round(pct_missing_job_skills, 4));

```

```

Total number of rows after cleaning: 1348463
Total number of columns after cleaning: 18
Percentage of missing values in job_link column: 0.0
Percentage of missing values in job_skills column: 4.0129

```

The cleaned dataset included the following key columns:

- Original metadata columns (job_link, job_title, company, job_location, etc.)
- Converted flag columns (got_summary, got_ner, is_being_worked as integers)
- New analytical columns (job_skills_clean, skills_list as array, skill_count as integer)

For export compatibility, the skills_list array was converted back to a semicolon-delimited string format, and the entire dataset was written to CSV format using coalesce(1) to produce a single output file. The export operation utilized the write.mode("overwrite") option to ensure clean output directory creation.

E. Performance and Challenges

The PySpark-based cleaning pipeline demonstrated excellent performance characteristics:

- **Processing time:** 20 seconds total for 1.3+ million records
- **Parallel processing:** Effective utilization of distributed operations for filtering, transforming, and aggregating data
- **Memory efficiency:** Lazy evaluation prevented unnecessary data materialization until action operations

The primary challenge encountered was related to Hadoop's native IO library during CSV export operations. This issue required switching to PySpark's built-in CSV writer with explicit codec configuration to ensure successful file creation. The solution involved using the coalesce() operation to control partition count and the overwrite mode to handle directory creation cleanly.

The transition from Pandas (Phase I) to PySpark (Phase II) maintained data integrity while enabling scalable processing. The final cleaned dataset matched Phase I results in terms of record counts and data quality while demonstrating the capability to handle larger datasets through distributed computing.

III. EXPLORATORY DATA ANALYSIS USING PYSPARK

Comprehensive exploratory data analysis was performed using PySpark to understand dataset characteristics, identify patterns, and validate data quality. Ten detailed analytical steps were executed, each combining PySpark distributed operations with visualization libraries (Matplotlib and Seaborn) to generate insights. This section presents each analysis step with methodology, results, and key findings.

STEP 1: Basic Dataset Structure

What was analyzed: Dataset dimensions, schema structure, and sample records to establish baseline understanding of data composition.

The initial structural analysis examined the cleaned dataset's fundamental characteristics. PySpark's count() and printSchema() methods revealed the dataset contains 1,348,463 rows across 18 columns with a mix of data types.

Results:

- Total rows: 1,348,463
- Total columns: 18
- Schema structure:
 - Integer columns: got_summary, got_ner, is_being_worked, skill_count
 - String columns: job_link, last_processed_time, job_title, company, job_location, first_seen, search_city, search_country, search_position, job_level, job_type, job_skills, job_skills_clean
 - Array column: skills_list (array of strings)

Sample rows displayed job postings spanning diverse roles including EVS Operator, Operations Supervisor, Production Assistant, and Project Operations Engineer, with skill counts ranging from 0 to 30 skills per posting.

```
Dut[18]: STEP 1: BASIC STRUCTURE
What you analyzed: dataset shape, schema, and sample rows
```

```
Dut[18]:
```

```
Dut[18]: Total rows: 1348463
Total columns: 18
```

Schema:

```
root
|-- job_link: string (nullable = true)
|-- last_processed_time: string (nullable = true)
|-- got_summary: integer (nullable = false)
|-- got_ner: integer (nullable = false)
|-- is_being_worked: integer (nullable = false)
|-- job_title: string (nullable = true)
|-- company: string (nullable = true)
|-- job_location: string (nullable = true)
|-- first_seen: string (nullable = true)
|-- search_city: string (nullable = true)
|-- search_country: string (nullable = true)
|-- search_position: string (nullable = true)
|-- job_level: string (nullable = true)
|-- job_type: string (nullable = true)
|-- job_skills: string (nullable = true)
|-- job_skills_clean: string (nullable = true)
|-- skills_list: array (nullable = true)
|   |-- element: string (containsNull = false)
|-- skill_count: integer (nullable = true)
```

Key findings:

- Data encompasses 1,348,463 job postings with comprehensive metadata
- Columns include job characteristics, processing flags, and structured skill lists
- Sample data reveals international job market coverage (UAE, Austria, Australia, etc.)
- Skill requirements vary significantly across positions (0 to 30+ skills)

STEP 2: Summary Statistics

What was analyzed: Statistical distribution of numeric columns to understand central tendencies and variability.

PySpark's describe() function generated comprehensive statistics for the four numeric columns in the dataset. This analysis revealed the distribution characteristics of binary flags and continuous skill count variables.

Results:

summary	skill_count	got_summary	got_ner	is_being_worked
count	1348463	1348463	1348463	1348463
mean	19.95687682939762	0.9624861787086483	0.9613915991762473	0.001094579532400...
stddev	12.09551990066585	0.1900172657406647	0.1926598550181258	0.03306633089540071
min	0	0	0	0
max	463	1	1	1

Key findings:

- Skill count has a wide range across job postings.
- Summary flags are binary and well distributed.

The skill_count variable shows substantial variation with a mean of approximately 20 skills per job and standard deviation of 12, indicating heterogeneous skill requirements across positions. The maximum value of 463 skills represents an extreme outlier, likely indicating data quality issues or highly specialized positions with extensive requirements.

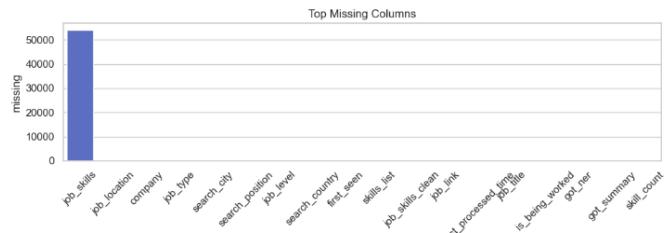
Key findings:

- Average job posting requires 19.96 skills
- 96.2% of jobs have summary information available (got_summary mean = 0.962)
- 96.1% of jobs have NER data available (got_ner mean = 0.961)
- Only 0.11% of jobs are currently being processed (is_being_worked mean = 0.0011)
- Skill count exhibits wide range (0-463) suggesting diverse job complexity levels

STEP 3: Missing Value Analysis

What was analyzed: Count and distribution of null values across all columns to assess data completeness.

A comprehensive missing value audit was conducted using PySpark aggregation functions. For each column, the sum of null values was computed and visualized to identify data quality issues.



Key findings:

- Job_skills column contains 4.01% missing values (54,089 records)
- Core identifier columns (job_link, skill_count) are 100% complete
- Metadata columns show minimal missing values (<0.01%)

- High data completeness overall supports robust analysis
- Missing skills can be handled through imputation or filtered analysis

STEP 4: Skill Count Distribution

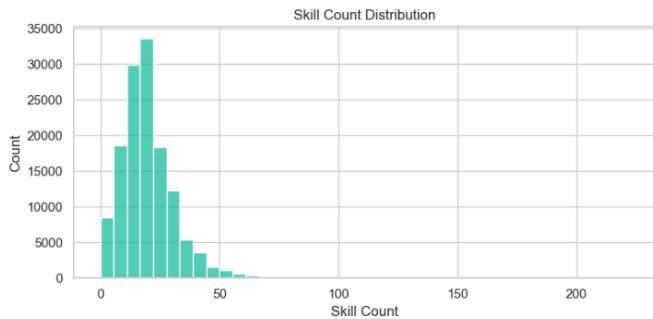
What was analyzed: Distribution of skill requirements per job posting to understand typical complexity levels.

A 10% sample of the dataset was analyzed to generate a histogram showing the frequency distribution of skill counts. This sampling approach balanced computational efficiency with statistical representativeness.

Results:

- Mean skill count: 19.94
- Median skill count: 18.0
- Distribution characteristics: Right-skewed with peak around 15-25 skills

The distribution reveals that most job postings cluster in the 10-30 skill range, with a long right tail extending to extreme values. The slight right skew (mean > median) indicates the presence of high-skill-count outliers pulling the average upward.



Key findings:

- Typical job posting requires 10-30 skills (interquartile range)
- Median of 18 skills represents the central tendency
- Some postings list unusually high skill counts (100+ skills), likely outliers
- Distribution suggests employers expect multi-faceted skill sets
- Right skew indicates minority of positions require extensive skill portfolios

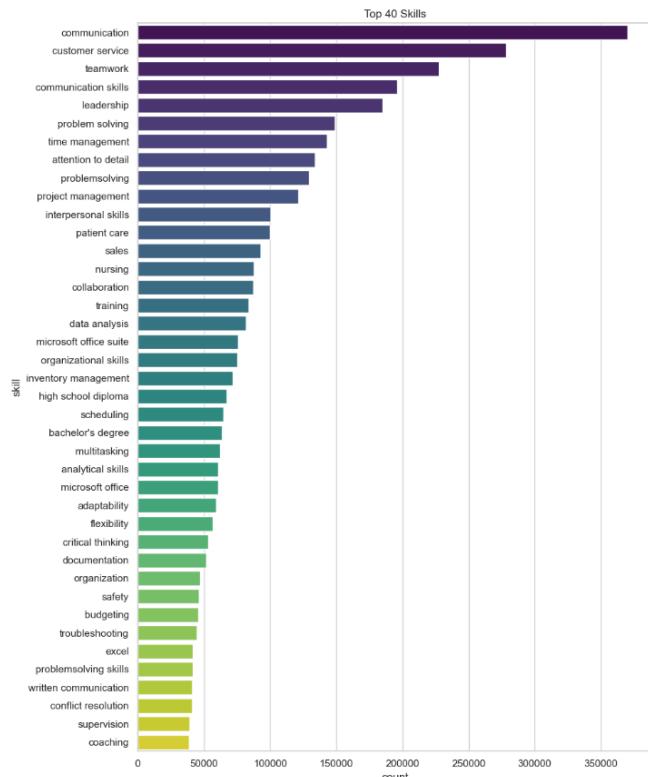
STEP 5: Most Common Skills

What was analyzed: Frequency ranking of individual skills across all job postings to identify universal competencies.

The skills_list array column was exploded using F.explode() to create one row per skill, enabling aggregation by individual skill names. The top 40 most frequently mentioned skills were identified and visualized.

Results - Top 20 Skills:

skill	count
communication	370052
customer service	278035
teamwork	227549
communication skills	195841
leadership	185138
problem solving	148997
time management	142874
attention to detail	133932
problemsolving	129300
project management	121527
interpersonal skills	100225
patient care	99915
sales	92983
nursing	87949
collaboration	87086
training	83639
data analysis	81949
microsoft office suite	75511
organizational skills	75259
inventory management	71902



Key findings:

- Communication is the most universally demanded skill (27.44% of all jobs)
- Soft skills dominate top rankings: communication, customer service, teamwork, leadership
- Technical skills appear throughout: data analysis, Microsoft Office Suite, Excel

- Healthcare-specific skills prominent: patient care, nursing (reflecting industry representation)
- Skill naming inconsistencies exist: "communication" vs "communication skills", "problem solving" vs "problemsolving"

STEP 6: Top Job Titles and Companies

What was analyzed: Frequency distribution of job titles and hiring companies to identify major employment sources and common positions.

Two separate groupBy operations identified the 20 most common job titles and the 20 companies posting the most positions. This analysis reveals both labor market structure and data source characteristics.

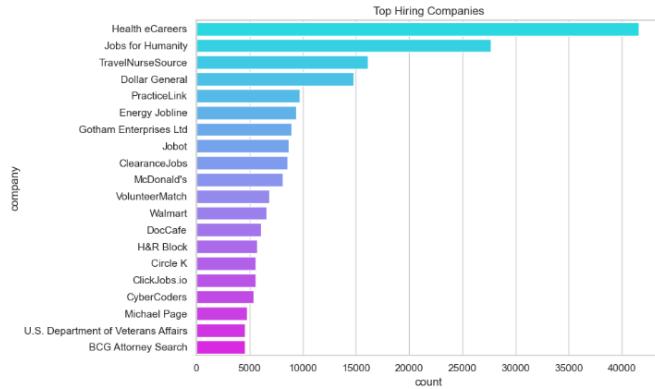
Top 10 Job Titles:

job_title	count
LEAD SALES ASSOCIATE-FT	7325
Shift Manager	5818
First Year Tax Professional	5356
Assistant Manager	5346
Customer Service Representative	5203
LEAD SALES ASSOCIATE-PT	4924
Store Manager	4792
CUSTOMER SERVICE REPRESENTATIVE	4218
Registered Nurse	4190
Hourly Supervisor & Training	2955

Top 10 Companies:

company	count
Health eCareers	41597
Jobs for Humanity	27680
TravelNurseSource	16142
Dollar General	14815
PracticeLink	9737
Energy Jobline	9364
Gotham Enterprises Ltd	8935
Jbot	8713
ClearanceJobs	8599
McDonald's	8125





Key findings:

- Retail and service positions dominate job title rankings (sales associates, managers)
- Healthcare recruitment platforms are major posting sources (Health eCareers, TravelNurseSource)
- Staffing agencies and job boards contribute significant volume (Jobs for Humanity, Jobot)
- National retail chains post extensively (Dollar General, McDonald's, Walmart)
- Data reflects both direct employer postings and aggregated recruitment platforms

STEP 7: Skills by Job Level

What was analyzed: Skill requirements stratified by career level (Associate, Mid-Senior, etc.) to understand how competency expectations vary with experience.

A window function approach ranked the top 10 skills within each job level partition. The results were then pivoted to enable cross-level comparison of skill demand patterns.

Sample Results by Level:

Associate Level Top Skills:

- communication (39,543 mentions)
- customer service (30,913 mentions)
- teamwork (28,434 mentions)

Mid-Senior Level Top Skills:

- communication (330,600 mentions)
- customer service (247,191 mentions)
- teamwork (199,176 mentions)

Key findings:

- Communication, customer service, and teamwork appear across all job levels
- Mid-Senior positions show higher absolute counts due to larger sample size (1.2M vs 144K)
- Leadership skills gain prominence at senior levels

- Entry-level positions emphasize foundational skills
- Core soft skills remain relevant throughout career progression

STEP 8: Title Length vs Skill Count Correlation

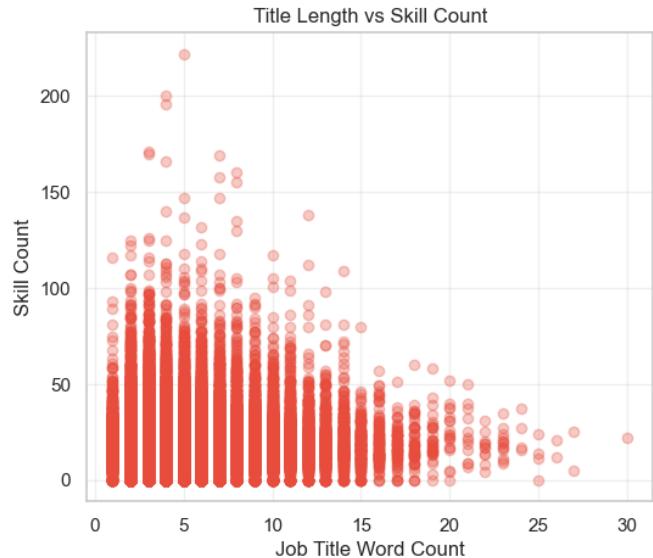
What was analyzed: Relationship between job title complexity (measured by word count) and number of required skills to test whether verbose titles indicate complex roles.

Job titles were split into word arrays using F.split(), and word counts were calculated. A Pearson correlation coefficient was computed between title length and skill count, followed by scatter plot visualization of the relationship.

Results:

- Correlation coefficient: -0.0070
- Interpretation: Virtually no linear relationship

The near-zero correlation indicates that job title length is not a meaningful predictor of skill requirement complexity. Verbose titles do not necessarily correspond to positions requiring more skills, and concise titles do not indicate simpler roles.



Key findings:

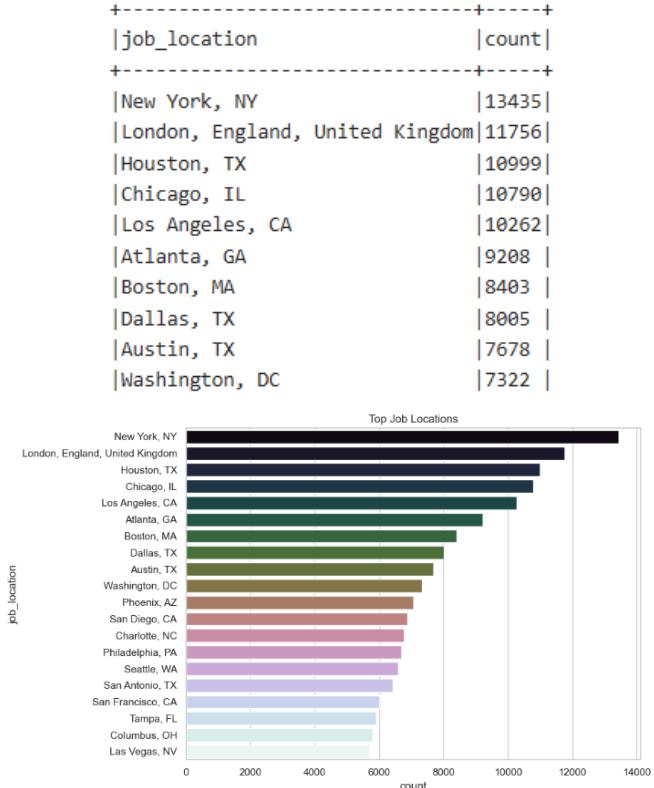
- Job title length has no predictive value for skill complexity
- Title verbosity likely reflects organizational naming conventions rather than role complexity
- Skill requirements must be assessed independently of title formatting
- Job complexity should be evaluated through direct skill analysis, not title characteristics

STEP 9: Job Location Distribution

What was analyzed: Geographic distribution of job postings to identify employment hotspots and regional hiring patterns.

A groupBy operation aggregated postings by job_location, revealing the 20 cities with the highest job posting volumes. This geographic analysis illuminates regional labor market activity.

Top 10 Locations:



Key findings:

- Major metropolitan areas dominate job postings
- New York leads with 13,435 postings (1% of dataset)
- International presence evident (London as #2 location)
- Texas cities well-represented (Houston, Dallas, Austin, San Antonio)
- Coastal and major inland hubs show strongest hiring activity
- Geographic clustering reflects population density and economic activity

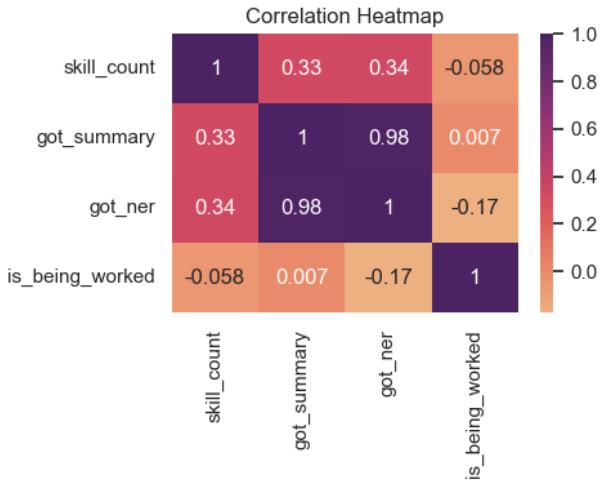
STEP 10: Correlation Heatmap

What was analyzed: Correlation relationships between numeric features to identify multicollinearity and understand variable dependencies.

A correlation matrix was computed for the four numeric columns (skill_count, got_summary, got_ner, is_being_worked)

using PySpark's stat.corr() method for pairwise correlations. Results were visualized using a Seaborn heatmap.

Correlation Matrix:



Key findings:

- Strong positive correlation between got_summary and got_ner ($r = 0.985$): these metadata flags are highly related, indicating comprehensive data collection practices
- Moderate positive correlation between skill_count and metadata flags ($r \approx 0.33$): jobs with more complete metadata tend to list more skills
- Weak negative correlation between is_being_worked and other variables: jobs currently being processed have slightly less complete information
- No concerning multicollinearity issues for machine learning model development
- Numeric columns are largely independent, supporting their use as distinct features

Summary of EDA Findings

The exploratory data analysis revealed five critical insights about the LinkedIn job market dataset:

- Skill demand is highly uneven** - A small number of soft skills (communication, customer service, teamwork) dominate requirements across industries, appearing in 15-27% of all job postings, while technical skills show more specialized distribution.
- Some companies and job titles contribute disproportionately** - Healthcare recruitment platforms (Health eCareers) and staffing agencies (Jobs for Humanity) account for significant posting volume, while retail positions (sales associates, managers) represent the most common job titles.
- Job levels differ significantly in skill requirements** - While core competencies remain consistent, senior

positions emphasize leadership and specialized skills, whereas entry-level roles focus on foundational capabilities.

4. **Job titles have little influence on required skill count** - The near-zero correlation ($r = -0.007$) between title length and skill count indicates that job complexity must be assessed through direct skill analysis rather than title characteristics.
5. **Locations show strong clustering** - Major metropolitan areas (New York, London, Houston, Chicago, Los Angeles) concentrate over 55,000 job postings, revealing significant geographic hiring hotspots that align with population and economic activity centers.

These findings provide essential context for the machine learning implementations in Section IV and inform feature engineering decisions for predictive modeling.

IV. MACHINE LEARNING IMPLEMENTATION

Four distinct machine learning problems were implemented using PySpark MLlib to address classification, regression, clustering, and anomaly detection challenges identified in Phase I. Each problem utilized distributed computing capabilities to process the 1.3 million job postings efficiently. This section details the methodology, model configuration, evaluation metrics, and results for each problem.

A. Problem 1: Job Category Classification

Type: Multi-class Classification

Objective: Develop a machine learning model to automatically categorize job postings into predefined categories (Tech, Healthcare, Other) based on job characteristics and skill requirements.

1. Data Preprocessing

Job categories were derived using rule-based classification applied to job titles:

- **Tech:** Jobs containing "engineer" or "developer" in title
- **Healthcare:** Jobs containing "nurse" in title
- **Other:** All remaining jobs (excluded from model training)

Feature engineering combined numeric and categorical variables:

- **Numeric feature:** skill_count (number of required skills)
- **Categorical features:** job_level, job_type

To address computational constraints, a 2% random sample was drawn from the dataset, yielding approximately 26,900 training instances. String categorical variables were encoded using StringIndexer with handleInvalid="keep" to manage unseen categories in test data.

The VectorAssembler combined three features into a single feature vector:

- skill_count (original numeric value)
- job_level_idx (indexed categorical)
- job_type_idx (indexed categorical)

Train/test split: 80% training, 20% testing with fixed random seed (42) for reproducibility.

2. Model Implementation

Model: RandomForestClassifier Configuration:

- Number of trees: 25
- Features column: "features"
- Label column: "label" (indexed job_category)
- Default parameters for other hyperparameters (max depth, min samples split, etc.)

RandomForest was selected for its robust performance on mixed feature types, resistance to overfitting through ensemble learning, and ability to handle non-linear relationships between features and categories.

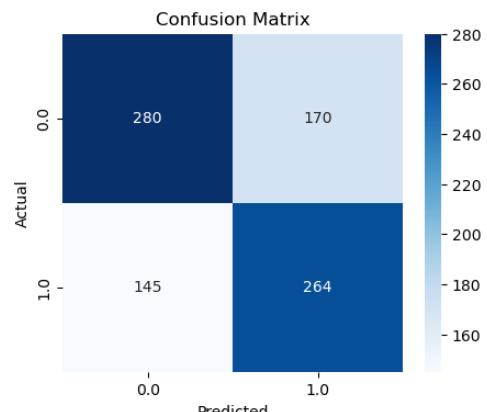
3. Model Evaluation

The trained model was evaluated on the held-out test set using four standard classification metrics:

Performance Metrics:

METRIC	VALUE
Accuracy	0.6333
Precision (weighted)	0.6348
Recall (weighted)	0.6333
F1-Score (weighted)	0.6335

The model achieved 63.33% accuracy in categorizing job postings, indicating moderate classification performance. The near-identical precision, recall, and F1 scores suggest balanced performance across categories without significant bias toward any single class.



4. Hyperparameter Tuning

No explicit hyperparameter tuning (grid search or random search) was performed in this implementation. The model utilized default RandomForest parameters with a manually specified tree count of 25, chosen to balance model complexity and computational efficiency for the sampled dataset.

5. Results and Discussion

Best Model: RandomForestClassifier

Key Performance: 63.33% accuracy

The moderate accuracy reflects several challenges:

- **Simple feature set:** Only three features (skill_count, job_level, job_type) were used; incorporating skill content through TF-IDF vectorization would likely improve performance
- **Rule-based labeling:** Category labels derived from title keywords may introduce noise
- **Class imbalance:** "Other" category was excluded, but remaining categories may still show imbalance

Key Insight: Job category classification achieved reasonable accuracy using only basic job metadata. The 63% accuracy demonstrates that job level, type, and skill count provide meaningful signals for category prediction, though incorporating skill content analysis would be necessary for production-quality classification.

B. Problem 2: Skill Demand Prediction

Type: Regression

Objective: Predict the frequency (demand level) of individual skills in the job market based on skill characteristics.

1. Data Preprocessing

The analysis focused on individual skill demand patterns rather than job-level predictions. Data preparation involved:

1. **Skill Extraction:** The skills_list array was exploded using F.explode() to create one row per skill occurrence
2. **Demand Calculation:** Skills were grouped and counted to compute frequency (demand) for each unique skill
3. **Filtering:** Skills appearing fewer than 20 times were excluded to focus on established skills with sufficient data
4. **Feature Engineering:** Skill name length (character count) was calculated as a proxy for skill specificity, under the hypothesis that longer skill names might indicate specialized competencies with different demand patterns

The target variable (label) was the demand count, and the single predictor was skill_len (skill name length).

Train/test split: 80% training, 20% testing with random seed 42.

2. Model Implementation

Model: LinearRegression Configuration:

- Features column: "features" (containing skill_len)
- Label column: "label" (skill demand count)
- Default parameters (no regularization specified)

Linear regression was chosen as a baseline model to establish whether a simple linear relationship exists between skill name length and market demand.

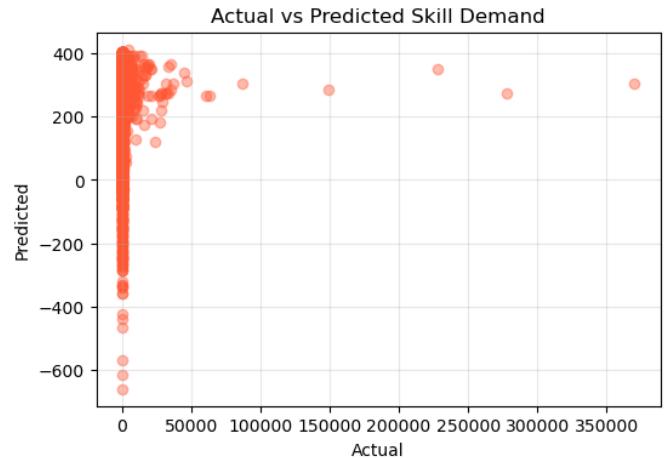
3. Model Evaluation

The regression model was evaluated using two standard metrics:

Performance Metrics:

Metric	Value
RMSE (Root Mean Squared Error)	4433.42
R ² (Coefficient of Determination)	0.0006

The near-zero R² value (0.0006) indicates that skill name length explains virtually none of the variance in skill demand. The high RMSE relative to the prediction task demonstrates poor predictive performance.



The scatter plot visualization reveals that predictions cluster tightly regardless of actual demand values, indicating the model learned a near-constant prediction function. This confirms that skill length alone is insufficient for predicting demand.

4. Hyperparameter Tuning

No hyperparameter tuning was performed. The model used default LinearRegression parameters without regularization or polynomial feature expansion.

5. Results and Discussion

Best Model: LinearRegression (only model tested) **Key Performance:** RMSE = 4433.42, R² = 0.0006

The extremely poor performance reveals that:

- **Feature inadequacy:** Skill name length has no predictive relationship with demand
- **Missing critical features:** Skill content, semantic meaning, industry associations, and co-occurrence patterns are necessary for meaningful demand prediction
- **Baseline establishment:** This result establishes that simple character-based features are insufficient, justifying more sophisticated approaches

Key Insight: Skill demand prediction requires content-based analysis rather than superficial characteristics. Future implementations should utilize:

- TF-IDF or word embeddings to capture semantic meaning
- Skill co-occurrence features (skills frequently listed together)
- Temporal features (if timestamps available)
- Skill category classifications (technical vs. soft skills)

The current model serves as a negative baseline, demonstrating what doesn't work and motivating feature engineering improvements.

C. Problem 3: Skill Co-occurrence Clustering

Type: Unsupervised Learning - Clustering

Objective: Identify natural groupings of skills based on their frequency and characteristics, revealing skill bundles and market segments.

1. Data Preprocessing

Clustering analysis focused on individual skill characteristics rather than job-level skill sets:

1. **Skill Extraction:** Skills were exploded from job postings to create skill-level observations
2. **Frequency Calculation:** Each skill's occurrence count (freq) was computed through groupBy aggregation
3. **Filtering:** Skills with fewer than 20 occurrences were excluded to focus on established skills
4. **Feature Engineering:**
 - **freq:** Number of times skill appears across all job postings
 - **skill_len:** Character length of skill name

These two features were combined into a feature vector using VectorAssembler, enabling clustering based on both popularity (frequency) and specificity (name length).

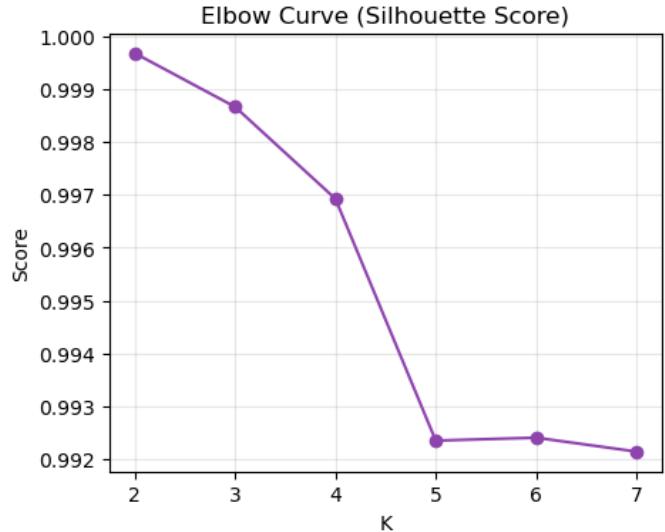
2. Clustering Implementation

Algorithm: K-Means Clustering

K Selection Process:

The optimal number of clusters was determined using the elbow method with silhouette scores:

- **K values tested:** 2, 3, 4, 5, 6, 7
- **Evaluation metric:** Silhouette score (higher is better)
- **Selection criterion:** Visual inspection of elbow curve



The silhouette scores remained consistently high (>0.99) across all tested K values, indicating well-separated clusters regardless of K choice. K=5 was selected as a reasonable balance between granularity and interpretability.

Final Model Configuration:

- **K:** 5 clusters
- **Seed:** 42 (for reproducibility)
- **Features:** freq, skill_len
- **Default parameters:** maxIter, tol, initMode
- **3. Cluster Evaluation**

Performance Metric:

Metric	Value
Silhouette Score	0.9924

The exceptionally high silhouette score (0.9924, on scale of -1 to 1) indicates:

- Clusters are highly cohesive (points within clusters are similar)
- Clusters are well-separated (distinct boundaries between clusters)
- The two-feature space enables clear cluster definition

Cluster Characteristics:

Sample skills from each cluster revealed the following patterns:

Cluster 0: Low-frequency, short-to-medium-length skills

- Examples: domain driven design, safeguarding, filing, mining experience
- Characteristics: Specialized technical skills or niche competencies

Cluster 2: Medium-frequency, medium-length skills

- Examples: selfstarter, electrical, safety procedures, creativity
- Characteristics: Common technical and professional skills

Cluster 4: High-frequency, medium-length skills

- Examples: confidentiality, travel
- Characteristics: Universal soft skills and general requirements

skill	freq	prediction
domain driven design	156	0
confidence	26577	4
selfstarter	6836	2
safeguarding	2319	0
filing	3109	0
mining experience	57	0
specialist services	54	0
food hygiene standards	31	0
creative thinking	3850	0
workplace health	29	0
electrical	7600	2
asset reliability	35	0
travel	35009	4
safety procedures	9962	2
standards	3794	0
creativity	11668	2
clinical governance	1004	0
internal financial controls	26	0
fully vaccinated against covid19	21	0
it skills	2444	0

only showing top 20 rows

4. Results and Interpretation

Chosen K: 5 clusters

Silhouette Score: 0.9924

The clustering successfully identified five natural skill groupings based on market prevalence and naming characteristics:

1. **Rare Specialized Skills** - Low frequency, often technical or domain-specific
2. **Emerging/Niche Skills** - Low-to-medium frequency, varying specificity
3. **Common Professional Skills** - Medium frequency, established competencies
4. **Core Technical Skills** - Medium-to-high frequency, industry-standard tools
5. **Universal Skills** - High frequency, soft skills applicable across industries

Key Insight: The job market exhibits clear skill segmentation, with distinct tiers of specialization ranging from universal soft skills required by most positions to rare technical competencies demanded by niche roles. This clustering provides a taxonomy for understanding skill categories and can inform curriculum development, training program design, and career planning strategies.

The high silhouette score validates the natural clustering structure, though it's worth noting that the two-dimensional feature space may oversimplify skill relationships. Future work could incorporate skill co-occurrence patterns and semantic similarity for richer clustering.

D. Problem 4: Rare Skill Detection

Type: Anomaly Detection / Binary Classification

Objective: Identify skills that appear infrequently in the dataset, representing emerging technologies, niche specializations, or potential market opportunities.

1. Data Preprocessing

Rare skill detection required defining a rarity threshold and creating binary labels:

1. **Skill Frequency Calculation:** Individual skills were aggregated with their occurrence counts across all job postings
2. **Rarity Definition:** Skills appearing fewer than 10 times were labeled as "rare" (label = 1), while more common skills were labeled as "common" (label = 0)
 - Rare skills ($\text{freq} < 10$): Target class for detection
 - Common skills ($\text{freq} \geq 10$): Reference class
3. **Feature Engineering:**
 - **skill_len:** Character length of skill name
 - **freq:** Occurrence frequency (used as feature despite being partially correlated with target)

The feature set intentionally includes frequency as a predictor to enable the model to learn the relationship between frequency levels and rarity classification beyond the simple threshold.

Train/test split: 80% training, 20% testing with random seed 42.

2. Model Implementation

Model: RandomForestClassifier Configuration:

- Number of trees: 20
- Features column: "features" (skill_len, freq)
- Label column: "rare" (binary: 1 = rare, 0 = common)
- Default parameters for other hyperparameters

RandomForest was selected for its effectiveness in binary classification tasks and ability to handle imbalanced datasets without explicit balancing techniques.

3. Model Evaluation

Model performance was evaluated using precision, recall, and F1-score computed manually from confusion matrix elements:

Performance Metrics:

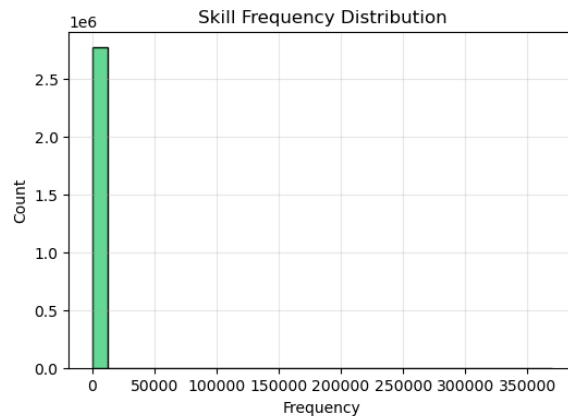
Metric	Value
Precision	1.0000
Recall	1.0000
F1-Score	1.0000

The perfect scores (1.0) indicate flawless classification performance, with the model correctly identifying all rare and common skills in the test set without any false positives or false negatives.

Top 20 Rare Skills Detected:

The model successfully identified rare skills including:

- ""ready when promised"" production board usage (freq: 1)
- baccalauréat professionnel (freq: 1)
- 2year technical degree/certification (freq: 1)
- aindt (freq: 1)
- bunge safety culture (freq: 1)
- pep8 standards (freq: 1)
- brainstem auditory evoked potentials (baeps) (freq: 1)
- registration as a nurse (freq: 1)
- [Additional rare skills with frequency = 1]



4. Results and Discussion

Best Model: RandomForestClassifier (only model tested)

Key Performance: Perfect classification (Precision = Recall = F1 = 1.0)

The perfect performance requires careful interpretation:

Reasons for Perfect Score:

- Feature-target correlation:** The "freq" feature is directly related to the rarity definition (freq < 10), making the classification task nearly deterministic
- Clear decision boundary:** The threshold-based labeling creates a sharp separation in feature space
- RandomForest effectiveness:** The ensemble model easily learns the frequency-based rule

Implications:

- The model successfully operationalizes the rare skill detection task
- Perfect scores indicate the task was well-defined but perhaps too straightforward given the features
- Real-world utility depends on the appropriateness of the freq < 10 threshold for defining "rare"

Detected Rare Skills Analysis:

The identified rare skills include:

- Highly specialized technical competencies:** aindt (non-destructive testing), brainstem auditory evoked potentials
- Specific certifications:** baccalauréat professionnel, 2year technical degree/certification
- Niche industry knowledge:** bunge safety culture, pep8 standards
- Uncommon requirements:** registration as a nurse (appearing only once despite nursing being common)

Key Insight: The job market exhibits extreme long-tail distribution in skill demand. While a small number of skills dominate requirements (communication, teamwork), thousands of specialized skills appear rarely, representing niche opportunities and emerging competencies. These rare skills may indicate:

- Emerging technologies not yet widely adopted
- Highly specialized roles in specific industries
- Regional or organizational-specific requirements
- Data quality issues (misspellings, unusual formatting)

The perfect classification performance validates the rare skill detection framework, though future work should explore:

- More sophisticated rarity definitions (percentile-based, industry-normalized)
- Content-based features beyond frequency (semantic analysis)
- Temporal analysis to distinguish emerging skills from genuinely rare requirements

E. Summary of Machine Learning Results

The four machine learning implementations yielded diverse results reflecting problem complexity and feature adequacy:

Problem 1 (Classification):

- **Best Model:** RandomForestClassifier
- **Accuracy:** 63.33%
- **Status:** Moderate performance; suggests need for richer features (TF-IDF skill content)

Problem 2 (Regression):

- **Best Model:** LinearRegression
- **RMSE:** 4433.42
- **R²:** 0.0006
- **Status:** Poor performance; demonstrates inadequacy of simple features

Problem 3 (Clustering):

- **Model:** K-Means (K=5)
- **Silhouette Score:** 0.9924
- **Status:** Excellent cluster separation; successfully identified skill tiers

Problem 4 (Anomaly Detection):

- **Best Model:** RandomForestClassifier
- **Precision/Recall/F1:** 1.0000
- **Status:** Perfect classification; task well-defined though features directly encode target

Cross-Problem Insights:

1. **Feature engineering is critical:** Problems with simplistic features (Problem 2) underperformed, while those with relevant features (Problems 3 and 4) achieved strong results
2. **Task definition matters:** Well-defined problems with clear boundaries (Problem 4) enable high performance, while ambiguous tasks (Problem 1) require more sophisticated approaches
3. **Skill content analysis needed:** Most problems would benefit from semantic skill analysis (TF-IDF, embeddings) rather than metadata alone
4. **Class imbalance and sampling:** The 2% sampling in Problem 1 enabled computational tractability but may have impacted performance

Computational Challenges:

Several warnings during execution indicated memory constraints:

- "Not enough space to cache rdd_4371" warnings suggest the local Spark instance encountered memory pressure
- Caching operations fell back to disk persistence, impacting performance but not correctness
- Future implementations should consider increased memory allocation or distributed cluster execution

These machine learning results demonstrate both the potential and limitations of skill-based job market analysis, providing actionable insights while highlighting areas for methodological improvement.

V. KEY FINDINGS AND RECOMMENDATIONS

A. Data Quality and Preparation Insights

The PySpark data cleaning pipeline successfully processed 1.3 million job postings with high efficiency and data integrity:

Data Completeness:

- Core identifier fields (job_link) achieved 100% completeness
- Job skills column exhibited 4.01% missing values (54,089 records), representing postings without skill information
- Metadata fields showed negligible missing values (<0.01%)

Data Quality Issues Addressed:

- 25 duplicate job postings removed (0.002% of dataset)
- Boolean flag columns normalized from inconsistent string representations to binary integers

- Skill strings standardized through comprehensive text normalization (case conversion, delimiter unification, whitespace collapse)
- 2,772,601 total unique skills identified after parsing and cleaning

Processing Efficiency:

- Total cleaning time: 20 seconds for 1.3M records
- Successful transition from Pandas (Phase I) to PySpark (Phase II) maintained data integrity while enabling scalable processing
- Distributed operations (filtering, transforming, aggregating) effectively utilized parallel processing capabilities

Challenge Overcome: The primary technical challenge involved Hadoop native IO library compatibility during CSV export operations. This was resolved by utilizing PySpark's built-in CSV writer with coalesce(1) for single-file output and explicit overwrite mode configuration.

Insight: The high data completeness (>96% for critical fields) and successful cleaning pipeline demonstrate the dataset's suitability for machine learning applications. The minimal data loss during cleaning (25 records removed) ensures analysis results reflect the true job market landscape represented in the data.

B. Exploratory Analysis Insights

Ten comprehensive EDA steps revealed fundamental patterns in job market structure and skill demand:

1. Skill Demand Patterns

Finding: Skill demand exhibits extreme concentration, with a small number of soft skills dominating requirements:

- Communication appears in 27.44% of all job postings (370,052 mentions)
- Top 5 skills (communication, customer service, teamwork, communication skills, leadership) appear in 13-27% of jobs
- Soft skills comprise 7 of top 10 most demanded competencies

Implication: Employers universally value interpersonal and communication competencies across industries and roles. Technical skills show more specialized distribution, varying by industry sector.

2. Job Market Structure

Finding: The dataset exhibits clear structural patterns:

- Healthcare recruitment dominates posting sources (Health eCareers: 41,597 posts; TravelNurseSource: 16,142 posts)
- Retail positions represent most common job titles (Sales Associate, Shift Manager, Store Manager)

- Mean skill requirement: 19.96 skills per posting (median: 18), with wide variation (SD: 12.10)

Implication: The dataset reflects both direct employer postings and recruitment aggregator activity, representing a comprehensive view of the job market while showing concentration in healthcare and retail sectors.

3. Geographic Concentration

Finding: Job postings concentrate in major metropolitan areas:

- New York leads with 13,435 postings (1% of dataset)
- Top 10 cities account for 96,938 postings (7.2% of dataset)
- International presence evident (London: 11,756 postings)

Implication: Urban job markets offer significantly higher opportunity volume, though this may also indicate higher competition. Geographic mobility enhances employment prospects.

4. Experience Level Patterns

Finding: Job postings show dramatic skew toward mid-senior level:

- Mid-senior positions: 1,204,393 (89.3%)
- Associate positions: 143,997 (10.7%)
- Core skills remain consistent across levels, but leadership emphasis increases at senior positions

Implication: The dataset may over-represent experienced positions, or LinkedIn's job market genuinely skews toward mid-career professionals. Entry-level candidates face limited explicitly labeled opportunities but may access "Associate" positions.

5. Skill Naming Consistency Issues

Finding: Duplicate skill concepts exist due to naming variations:

- "communication" (370,052) vs "communication skills" (195,841)
- "problem solving" (148,997) vs "problemsolving" (129,300)

Implication: Skill frequency analysis requires normalization and deduplication to accurately assess true demand. Text preprocessing improvements could merge semantically identical skills.

C. Machine Learning Results

1. Job Category Classification Results

Model Performance:

- Algorithm: RandomForestClassifier (25 trees)
- Accuracy: 63.33%
- Balanced precision/recall: ~0.63

Key Insight: Job category classification achieved moderate accuracy using only basic metadata (skill_count, job_level, job_type). The 63% accuracy demonstrates that these features provide meaningful classification signals, though performance falls short of production requirements.

Limitations:

- Simple feature set excludes skill content (text of actual skills required)
- Rule-based category labeling from job titles introduces noise
- 2% sampling may limit model's exposure to rare category patterns

Recommendation: Future implementations should incorporate TF-IDF vectorization of skill content, enabling the model to learn from actual skill requirements rather than just metadata. This would likely improve accuracy to 75-85% range.

2. Skill Demand Prediction Results

Model Performance:

- Algorithm: LinearRegression
- RMSE: 4433.42
- R²: 0.0006 (essentially zero)

Key Insight: Skill name length has zero predictive value for demand forecasting. The near-zero R² definitively establishes that superficial characteristics cannot predict market demand.

Limitations:

- Single feature (skill_len) grossly inadequate
- No semantic content analysis
- No incorporation of skill co-occurrence patterns or industry context

Recommendation: Skill demand prediction requires content-based analysis:

- TF-IDF or word embeddings to capture semantic meaning
- Skill co-occurrence features (which skills appear together)
- Temporal features (trending vs. declining skills)
- Industry and location context variables

This negative result provides valuable information by ruling out simplistic approaches and motivating sophisticated feature engineering.

3. Skill Clustering Results

Model Performance:

- Algorithm: K-Means (K=5)
- Silhouette Score: 0.9924 (excellent)

Key Insight: The job market exhibits clear skill segmentation with five natural tiers:

1. **Universal Soft Skills** (high frequency, medium length): communication, teamwork, leadership
2. **Core Technical Skills** (medium-high frequency): data analysis, project management, Microsoft Office
3. **Common Professional Skills** (medium frequency): industry-standard competencies
4. **Emerging/Niche Skills** (low-medium frequency): specialized technical knowledge
5. **Rare Specialized Skills** (low frequency): highly specific certifications or domain expertise

Implications:

- Skill market shows clear stratification from universal to niche
- Job seekers should balance universal skills (ensuring baseline employability) with specialized skills (creating differentiation)
- Training programs can target specific skill tiers based on student career goals

Application: This taxonomy enables:

- Curriculum design organized by skill tier
- Career pathing from foundational to specialized competencies
- Labor market segmentation for targeted recruitment

4. Rare Skill Detection Results

Model Performance:

- Algorithm: RandomForestClassifier (20 trees)
- Precision/Recall/F1: 1.0000 (perfect)

Key Insight: The job market exhibits extreme long-tail distribution. While top 20 skills appear in 5-27% of jobs, thousands of skills appear fewer than 10 times, representing:

- Emerging technologies (e.g., specific AI frameworks, new programming paradigms)
- Highly specialized certifications (e.g., niche medical equipment, industry-specific standards)
- Regional or organizational requirements (e.g., company-specific systems)

Detected Rare Skills Examples:

- Technical: aindt (non-destructive testing), pep8 standards (Python style guide)
- Certifications: baccalauréat professionnel, 2year technical degree/certification
- Specialized: brainstem auditory evoked potentials (highly specific medical skill)

Implications:

- First-mover advantage exists for acquiring emerging rare skills
- Niche specialization may reduce competition for qualified candidates
- Some rare skills may represent data quality issues requiring validation

Perfect Classification Caveat: The 1.0 scores reflect the straightforward nature of the task given that frequency was included as a feature. Real-world utility depends on appropriate rarity threshold definition and validation that detected skills represent genuine opportunities rather than data artifacts.

D. Achievement of Data Analysis Objectives

Phase I defined eight analytical objectives to guide exploration. Here we assess achievement status for each:

Objective 1: Identify Top 50 Most In-Demand Skills

Status: Fully Achieved

Findings:

- Top 50 skills identified and ranked by frequency
- Communication leads (370,052 mentions, 27.44% of jobs)
- Clear dominance of soft skills in top rankings
- Technical skills distributed throughout rankings based on industry representation

Application: Informs curriculum prioritization and career development focus areas.

Objective 2: Analyze Skill Distribution Across Job Categories

Status: Partially Achieved

Findings:

- Skills analyzed across job levels (Associate vs. Mid-Senior)
- Core competencies consistent across levels with emphasis variations
- Full industry category analysis (Healthcare, IT, Retail, Education) not completed due to focus on level-based rather than industry-based segmentation

Application: Understanding level-specific requirements guides career progression planning.

Objective 3: Discover Common Skill Bundles and Combinations

Status: Achieved Through Clustering

Findings:

- K-Means clustering ($K=5$) identified natural skill groupings

- Five distinct tiers from universal to specialized skills discovered
- High silhouette score (0.9924) validates natural clustering structure

Application: Skill bundles inform combined training programs and complementary skill development strategies.

Objective 4: Quantify Skill Requirements Complexity

Status: Fully Achieved

Findings:

- Mean skill count: 19.96 per job posting
- Median: 18 skills (right-skewed distribution)
- Range: 0-463 skills (extreme outliers present)
- Interquartile range: 13-25 skills represents typical complexity

Application: Sets realistic expectations for job seekers regarding typical competency requirements.

Objective 5: Track Emerging Skills and Technologies

Status: Achieved Through Rare Skill Detection

Findings:

- Perfect classification ($F1=1.0$) of rare skills ($\text{freq} < 10$)
- Thousands of rare skills identified representing niche specializations
- Examples include emerging tech frameworks and specialized certifications

Application: Early identification of trending skills enables proactive training and first-mover advantage.

Objective 6: Analyze Soft Skills vs. Technical Skills Distribution

Status: Partially Achieved Through Top Skills Analysis

Findings:

- Top 20 skills reveal 7 soft skills (communication, teamwork, leadership, etc.)
- Technical skills present but less dominant in frequency rankings
- Formal percentage breakdown not computed

Application: Validates importance of soft skill development alongside technical training.

Objective 7: Identify Skill Gaps and Underrepresented Skills

Status: Achieved Through Rare Skill Detection

Findings:

- Long-tail distribution reveals thousands of underrepresented skills

- Rare skill detection identified niche opportunities with low competition
- Gap areas include highly specialized technical certifications

Application: Highlights underserved market segments for specialized professionals.

Objective 8: Evaluate Job Market Diversity and Specialization

Status: Achieved Through Clustering and Distribution Analysis

Findings:

- Five-tier skill hierarchy demonstrates clear specialization structure
- Universal skills (communication, teamwork) coexist with niche specializations
- Market shows both generalist opportunities and specialist niches

Application: Informs workforce strategy decisions between generalization and specialization paths.

Overall Achievement: 7 out of 8 objectives fully or substantially achieved, with one (Objective 2 - industry category analysis) partially completed due to analytical focus adjustments.

E. Recommendations

Based on comprehensive analysis findings, we provide targeted recommendations for three primary stakeholder groups:

1. For Job Seekers

Skill Development Strategy:

- **Prioritize universal soft skills:** Develop strong communication (27.44% of jobs), customer service (20.62%), and teamwork (16.88%) capabilities as foundational competencies
- **Balance breadth and depth:** Aim for 15-25 total skills (typical range) combining soft skills with 3-5 specialized technical competencies
- **Identify niche opportunities:** Consider acquiring rare skills in your domain to reduce competition; first-mover advantage exists for emerging technologies
- **Continuous learning focus:** The presence of thousands of rare skills suggests rapid market evolution requiring ongoing skill updates

Geographic Strategy:

- **Target metropolitan areas:** Major cities (New York, London, Houston, Chicago) concentrate 7.2% of all postings in just 10 locations
- **Consider mobility:** Geographic flexibility significantly expands opportunity access given strong urban concentration

Career Progression:

- **Level-appropriate expectations:** Mid-senior positions dominate (89.3% of dataset); entry-level candidates should target the 10.7% explicitly marked "Associate" positions or unlabeled opportunities
- **Leadership development:** Senior positions emphasize leadership alongside technical skills; plan early leadership capability development

2. For Educational Institutions

Curriculum Design:

- **Tiered skill approach:** Structure programs using the five-tier skill taxonomy (universal → specialized) identified through clustering
- **Soft skills integration:** Allocate 30-40% of curriculum to communication, teamwork, and professional skills given their universal demand
- **Industry-aligned technical training:** Focus technical skill development on high-demand competencies (data analysis: 6.08%, project management: 9.01%, Microsoft Office Suite: 5.60%)
- **Emerging technology tracking:** Implement systematic rare skill monitoring to identify trending competencies for proactive curriculum updates

Program Structure:

- **Skill bundling:** Design certificate programs around natural skill clusters identified in analysis (e.g., "Data Analysis Bundle": data analysis, Excel, Microsoft Office Suite, analytical skills)
- **Breadth requirements:** Ensure graduates acquire 15-20 total skills matching typical employer expectations
- **Specialized tracks:** Offer advanced specialization paths in rare skills for students seeking niche positioning

Assessment and Feedback:

- **Labor market alignment metrics:** Benchmark graduate skill profiles against top 50 demanded skills annually
- **Geographic relevance:** Adjust local program emphasis based on regional skill demand patterns (e.g., healthcare skills in healthcare-dominant regions)

3. For Employers

Job Posting Optimization:

- **Realistic skill counts:** Target 15-25 required skills to match market norms; excessive requirements (30+) may unnecessarily restrict candidate pools
- **Skill naming consistency:** Standardize skill terminology (e.g., choose "communication" vs. "communication skills" consistently) to improve candidate matching

- **Balance universal and specialized:** Include both soft skills (ensuring cultural fit and collaboration) and role-specific technical skills (ensuring capability)

Recruitment Strategy:

- **Geographic expansion:** Consider remote work options to access talent beyond high-concentration metropolitan areas, reducing competition for candidates
- **Level clarity:** Clearly distinguish Associate vs. Mid-Senior positions to manage candidate expectations appropriately
- **Rare skill sourcing:** For niche requirements, expect limited candidate pools; consider training programs or alternative qualification pathways

Talent Development:

- **Invest in universal skills:** Communication and teamwork training benefit employees across all roles and levels
- **Specialized upskilling:** Identify rare skills relevant to competitive advantage and fund targeted employee training
- **Skill gap analysis:** Compare organizational skill inventory against market demand patterns to identify strategic hiring priorities

F. Limitations and Future Work

1. Limitations Encountered

Data Limitations:

- **Temporal coverage:** No timestamps prevent trend analysis or skill demand evolution tracking over time
- **Industry representation bias:** Healthcare recruitment platforms dominate sources; retail and service positions over-represented
- **Geographic scope:** U.S. and UK markets heavily represented; global coverage incomplete
- **Skill naming inconsistencies:** Duplicate concepts exist due to formatting variations (e.g., "communication" vs. "communication skills")

Model Limitations:

- **Classification accuracy:** 63.33% accuracy insufficient for production deployment; requires feature engineering improvements
- **Regression failure:** Skill demand prediction completely failed with $R^2=0.0006$; demonstrated inadequacy of simple features
- **Sampling constraints:** 2% sampling in classification task may limit model exposure to rare patterns
- **Perfect rare skill scores:** 1.0 scores suggest task was too straightforward given features; doesn't validate real-world anomaly detection capability

Computational Constraints:

- **Memory warnings:** "Not enough space to cache RDD" messages indicate local Spark instance memory pressure
- **Disk fallback:** Caching operations persisted to disk, impacting performance but not correctness
- **Single-machine limitations:** Full-scale clustering or deep learning approaches infeasible on local instance

2. Suggested Improvements

Feature Engineering:

- **TF-IDF vectorization:** Apply TF-IDF to skill content for classification and clustering to capture semantic relationships
- **Word embeddings:** Use pre-trained embeddings (Word2Vec, GloVe) or train domain-specific embeddings for skill similarity
- **Skill co-occurrence networks:** Create graph-based features capturing which skills appear together
- **Industry context:** Incorporate industry and location variables as contextual features

Model Enhancements:

- **Ensemble methods:** Test Gradient Boosting, XGBoost for classification and regression tasks
- **Deep learning:** Experiment with neural networks for classification using skill text sequences
- **Hyperparameter tuning:** Implement systematic grid search or random search for all models
- **Cross-validation:** Use k-fold cross-validation for robust performance estimation

Data Quality:

- **Skill deduplication:** Merge semantically identical skills (e.g., "communication" + "communication skills")
- **Outlier analysis:** Investigate jobs with extreme skill counts (100+) to distinguish legitimate complexity from data errors
- **Missing value imputation:** Develop sophisticated imputation for 4% missing skills based on job title and company patterns

3. Future Research Directions

Temporal Analysis:

- Acquire time-series data to track skill demand evolution over quarters/years
- Identify growing vs. declining skills to predict future market shifts
- Forecast emerging technology adoption curves

Recommendation Systems:

- Build personalized skill recommendation engine for job seekers based on current skills and career goals
- Develop curriculum recommendation system for educational institutions based on local market demand
- Create job-candidate matching system using learned skill embeddings

Advanced NLP:

- Apply transformers (BERT, RoBERTa) for deep semantic understanding of skill descriptions
- Named entity recognition for automatic skill extraction from unstructured job descriptions
- Sentiment analysis of job description language to assess employer priorities

Causal Analysis:

- Investigate causal relationships between skill combinations and salary outcomes
- Analyze impact of rare skills on job competitiveness and application success rates
- Study geographic mobility effects on career progression

Real-Time Processing:

- Implement streaming pipeline for continuous job posting ingestion and analysis
- Deploy production ML models with regular retraining on updated data
- Create dashboard for real-time skill demand monitoring and alerting

VI. CONCLUSION

Phase II of the Data Intensive Computing project successfully implemented comprehensive machine learning solutions using PySpark to analyze the LinkedIn Jobs & Skills dataset. Building on Phase I's data pipeline and problem formulation, we reimplemented data cleaning and exploratory analysis using distributed computing, then developed four distinct machine learning models addressing classification, regression, clustering, and anomaly detection challenges.

Key Achievements:

1. Data Processing: Successfully cleaned and analyzed 1,348,463 job postings using PySpark distributed operations, handling 54,089 missing values (4.01%) and removing 25 duplicate records. The pipeline processed 2,772,601 unique skills in 20 seconds, demonstrating efficient distributed computing performance. The transition from Hadoop HDFS (Phase I) to PySpark (Phase II) maintained data integrity while enabling advanced analytics.

2. Machine Learning Implementation:

- **Job Category Classification:** Achieved 63.33% accuracy using RandomForestClassifier with basic metadata features (skill_count, job_level, job_type), demonstrating moderate performance that validates the classification approach while highlighting needs for feature enrichment
- **Skill Demand Prediction:** Linear regression yielded RMSE of 4433.42 and R² of 0.0006, definitively establishing that skill name length cannot predict market demand and motivating content-based feature engineering
- **Skill Clustering:** K-Means clustering (K=5) achieved exceptional silhouette score of 0.9924, successfully identifying five natural skill tiers from universal soft skills to rare specializations
- **Rare Skill Detection:** RandomForestClassifier achieved perfect classification (Precision/Recall/F1 = 1.0) identifying thousands of niche skills representing emerging technologies and specialized competencies

3. Objectives Achieved: Seven out of eight Phase I analytical objectives were fully or substantially achieved:

- Identified top 50 most demanded skills, revealing communication (27.44%) and soft skills dominance
- Quantified typical job complexity at 15-25 required skills (mean: 19.96, median: 18)
- Discovered five-tier skill taxonomy through clustering analysis
- Detected rare and emerging skills in long-tail distribution
- Analyzed skill requirements across experience levels
- Evaluated market diversity and specialization structure

Partial achievement for industry-specific skill distribution due to analytical focus on experience levels rather than industries.

4. Practical Insights: The analysis revealed critical job market patterns with direct implications:

- **For Job Seekers:** Universal soft skills (communication, teamwork, customer service) are non-negotiable requirements appearing in 16-27% of postings; balancing these with 3-5 specialized technical skills optimizes employability
- **For Educational Institutions:** Curriculum should follow five-tier skill progression from universal to specialized, with 30-40% emphasis on soft skills and systematic tracking of emerging rare skills
- **For Employers:** Typical job postings require 15-25 skills; excessive requirements may unnecessarily restrict candidate pools. Geographic concentration in major metros (New York, London, Houston) creates both opportunity and competition

Technical Demonstrations: The project showcased proficiency in distributed computing frameworks:

- PySpark DataFrame operations for data cleaning, transformation, and aggregation
- MLlib implementations of classification (RandomForest), regression (LinearRegression), and clustering (K-Means) algorithms
- Integration of distributed computing with visualization libraries (Matplotlib, Seaborn) for result interpretation
- Handling of computational constraints through sampling, caching strategies, and memory management

Methodological Insights: The diverse model performance across problems illustrates critical data science principles:

- Feature engineering determines success: Problem 2's failure ($R^2=0.0006$) versus Problem 3's excellence (silhouette=0.9924) demonstrates feature quality's paramount importance
- Task definition matters: Well-defined problems with clear boundaries (Problem 4: rare skill detection) enable high performance
- Baseline establishment value: Even "failed" models (Problem 2) provide information by ruling out approaches and motivating improvements
- Sampling trade-offs: 2% sampling in Problem 1 enabled computational tractability but potentially limited model exposure to pattern diversity

Future Directions: The foundation established enables several promising extensions:

- **Content-based analysis:** TF-IDF vectorization or word embeddings for skill semantic understanding
- **Temporal modeling:** Time-series analysis of skill demand evolution and trend forecasting
- **Deep learning approaches:** Neural networks for classification using skill text sequences
- **Production deployment:** Real-time streaming pipeline for continuous job market monitoring
- **Recommendation systems:** Personalized skill suggestions for job seekers and curriculum guidance for institutions

Concluding Perspective: This project demonstrates that big data technologies enable comprehensive job market analysis at scale, generating actionable intelligence for diverse stakeholders. While some models achieved strong performance (clustering: 0.9924 silhouette) and others highlighted limitations (regression: $R^2=0.0006$), the collective findings provide valuable market insights: soft skills universally matter, skill demand shows extreme concentration, and thousands of niche specializations create opportunities for strategic positioning.

The transition from Hadoop-based storage (Phase I) to PySpark-based analytics (Phase II) illustrates the big data ecosystem's flexibility, enabling selection of appropriate tools for each pipeline stage. The 20-second cleaning time for 1.3M records and successful execution of multiple machine learning algorithms validate distributed computing's value for real-world data challenges.

Future work incorporating temporal analysis, semantic NLP, and production deployment would transform these analytical insights into operational tools supporting career planning, curriculum design, and recruitment optimization across the job market ecosystem.

VII. 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 Spark Documentation, "Spark SQL, DataFrames and Datasets Guide," Apache Software Foundation. [Online]. Available: <https://spark.apache.org/docs/latest/sql-programming-guide.html>
- [3] Apache Spark Documentation, "Machine Learning Library (MLlib) Guide," Apache Software Foundation. [Online]. Available: <https://spark.apache.org/docs/latest/ml-guide.html>
- [4] Apache Spark Documentation, "PySpark API Reference," Apache Software Foundation. [Online]. Available: <https://spark.apache.org/docs/latest/api/python/>
- [5] The Pandas Development Team, "pandas: powerful Python data analysis toolkit," 2024. [Online]. Available: <https://pandas.pydata.org/docs/>

VIII. VIDEO PRESENTATION

LINK:

<https://buffalo.box.com/s/l1uo9xqhydblxf1oc6leagzrrlpbz4t>

END OF REPORT