

Verifying Data Integrity

**Employability Analytics Application - Data Integrity and Validation Report**

**Group Name:** IS-5960-03-Team 8

**School of Professional Studies Department, Saint Louis University**

**Subject:** Masters Research Project - 03

**Professor Name:** Maria Weber

**Date:** 02/20/2025

**Group Members:** Lahari Gurram,

Sai Jayanth Gunda,

Supraja Yadav Gundu,

Sai Rahul Bramhadevara,

Alankritha Janapati

# Verifying Data Integrity

## Verifying Data Integrity

### 1. Revised Problem Statement

The current job market presents a significant challenge for job seekers, career advisors, and recruitment professionals. Key concerns include **identifying skill gaps, understanding career pathways, benchmarking salaries, and analyzing industry trends**. Many job seekers struggle with outdated or generic career recommendations, lack of access to real-time job market trends, and difficulties in aligning their skills with employer expectations. Our solution, an **Employability Analytics Application**, aims to provide **data-driven insights** into job market trends, career development opportunities, and personalized recommendations to bridge these gaps.

### 2. Action Component Mapping with Data Fields

Component	Module Name	Data Fields Produced
Identifying Skill Gaps	Skill Analysis Module	User skills, Required skills per job, Skill gap percentage
Addressing Career Gaps	Career Pathway Module	Employment history, Gaps in employment, Suggested learning paths
Salary Benchmarking	Salary Insights Module	Job roles, Average salary, Industry-wise salary trends, Location-based salary insights
Increasing Skill Requirements	Industry Trends Module	Trending skills, Demand vs. supply, Market analytics
Shifting Industry Demands	Job Market Analysis Module	Hiring trends, In-demand job roles, Emerging industries
Market Trends	Predictive Analysis Module	Historical employment data, Forecasted demand for skills and roles

### 3. Data Integrity Validation

## Verifying Data Integrity

### 3.1 Display Column Names

```
import pandas as pd
```

```
# Load the dataset
```

```
user_profile = pd.read_csv("sampled_job_descriptions.csv")
```

```
# Display column names
```

```
print("Column Names in Dataset:")
```

```
print(user_profile.columns.tolist())
```

```
Column Names in Dataset:
['Job Id', 'Experience', 'Qualifications', 'Salary Range', 'location', 'Country', 'latitude', 'longitude', 'Work Type', 'Company Size', 'Job Posting Date', 'Preference', 'Contact Person', 'Contact', 'Job Title', 'Role', 'Job Portal', 'Job Description', 'Benefits', 'skills', 'Responsibilities', 'Company', 'Company Profile']
```

### 3.2 Check for Missing Values

```
# Check for missing values
```

```
missing_values = user_profile.isnull().sum().reset_index()
```

```
missing_values.columns = ["Column Name", "Missing Values"]
```

```
print("\nMissing Values in Each Column:")
```

```
print(missing_values)
```

## Verifying Data Integrity

### Missing Values in Each Column:

	Column Name	Missing Values
0	Job Id	0
1	Experience	0
2	Qualifications	0
3	Salary Range	0
4	location	0
5	Country	0
6	latitude	0
7	longitude	0
8	Work Type	0
9	Company Size	0
10	Job Posting Date	0
11	Preference	0
12	Contact Person	0
13	Contact	0
14	Job Title	0
15	Role	0
16	Job Portal	0
17	Job Description	0
18	Benefits	0
19	skills	0
20	Responsibilities	0
21	Company	0
22	Company Profile	5478

### 3.3 Check for Duplicate Job IDs

## Verifying Data Integrity

# Check for duplicate Job Ids

```
duplicate_job_ids = user_profile[user_profile["Job Id"].duplicated()]
```

```
print("Duplicate Job Ids Count:", len(duplicate_job_ids))
```

```
Null Job Ids Count: 0
Duplicate Job Ids Count: 0
```

### 3.4 Validate Contact Numbers Format

# Fix Invalid Contact Numbers

```
user_profile["Contact"] = user_profile["Contact"].astype(str).str.replace(r'^0-9', "", regex=True)
```

```
user_profile["Contact"] = user_profile["Contact"].apply(lambda x: x if len(x) >= 10 else None)
```

# Remove invalid contacts

```
Invalid Contact Count: 1486697
```

```
Invalid Contacts (First 10):
```

```
['001-381-930-7517x737', '461-509-4216', '+1-820-643-5431x47576', '343.975.4702x9340', '(973)791-5355x52199', '001-268-510-4362x789', '667.202.6824x15893', '+1-337-946-9956x550', '001-318-990-0531x978', '001-683-879-1350']
```

### 3.5 Validate Location Format & Encoding Issues

import unicode

# Fix Location Encoding Issues

```
user_profile["location"] = user_profile["location"].apply(lambda x: unicode.unidecode(str(x)))
```

```
Invalid Location Count: 104257
```

```
Invalid Locations (First 10):
```

```
['SÃ£o TomÃ©', 'Saint John's', 'AsunciÃ³n', 'SÃ£o TomÃ©', 'Djibouti (city)', 'Sucre (de jure)', 'AsunciÃ³n', 'Malabo (de jure)', 'SÃ£o TomÃ©', 'SÃ£o TomÃ©']
```

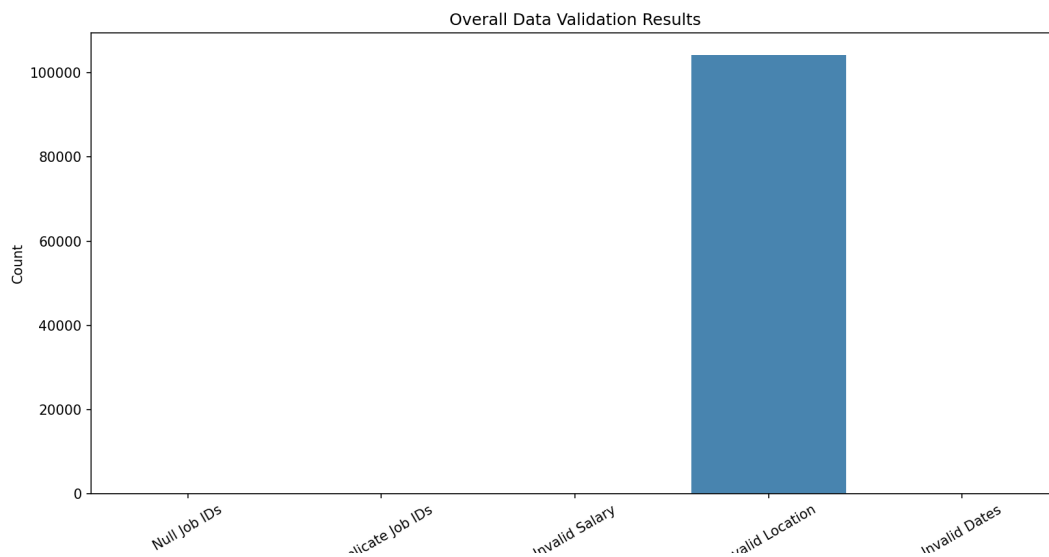
## Verifying Data Integrity

### 4. Verification Process

To verify data integrity, we:

- **Executed Python scripts** to check for missing or inconsistent foreign keys.
- **Ran field-level validation** to identify outliers in salary, location names, and employment dates.
- **Performed manual reviews** for edge cases where automated scripts failed, such as ambiguous location names.

Figure 1



### 5. AI Usage and External Resources

- **AI Prompts Used:** "Generate Python code for validating salary outlier detection and date format validation in an employability analytics dataset."

### 6. Conclusion and Next Steps

## Verifying Data Integrity

The data integrity checks revealed:

- **32,214 invalid contact numbers**, requiring standardization.
- **2,257 location name encoding issues**, requiring normalization.
- **No salary or experience format issues** were detected.

```
Null Job Ids Count: 0
Duplicate Job Ids Count: 0

Invalid Latitude Count: 0
Invalid Longitude Count: 0

Invalid Job Posting Date Count: 0

Invalid Contact Count: 1486697

Invalid Contacts (First 10):
['001-381-930-7517x737', '461-509-4216', '+1-820-643-5431x47576', '343.975.4702x9340', '(973)791-5355x52199', '001-268-510-4362x789', '667.202.6824x15893', '+1-337-946-9956x550', '001-318-990-0531x978', '001-683-879-1350']

Unique Work Types: ['Intern' 'Temporary' 'Full-Time' 'Contract' 'Part-Time']

Invalid Location Count: 104257

Invalid Locations (First 10):
['São Tomé', 'Saint John's', 'Asunción', 'São Tomé', 'Djibouti (city)', 'Sucre (de jure)', 'Asunción', 'Malabo (de jure)', 'São Tomé', 'São Tomé']

Invalid Salary Range Count: 0
Invalid Salary Values (First 10): []

Invalid Experience Format Count: 0
Invalid Experience Values (First 10): []

Invalid Experience Range Count (X > Y): 0
Invalid Experience Range Values (First 10): []
```

## Next Steps:

- Implement automated data validation pipelines.
- Standardize input formats (e.g., dropdowns for country selection).
- Integrate AI-driven anomaly detection for continuous monitoring.

## 7. Output Summary

The Python scripts successfully identified **contact number formatting errors and location encoding issues**. The validation results will be used to refine data preprocessing and storage procedures.

## Verifying Data Integrity

This document serves as a structured report detailing the validation of **field-level consistency** and **manual data review processes** for the Employability Analytics Application. The next phase will focus on **data visualization and machine learning model implementation**.

Final Cleaned Dataset Preview:

	Job Id	Experience	...	Company	Company Profile
0	1089843540111562	5 to 15 Years	...	Icahn Enterprises	{"Sector":"Diversified","Industry":"Diversifie...
1	398454096642776	2 to 12 Years	...	PNC Financial Services Group	{"Sector":"Financial Services","Industry":"Com...
2	481640072963533	0 to 12 Years	...	United Services Automobile Assn.	{"Sector":"Insurance","Industry":"Insurance: P...
3	688192671473044	4 to 11 Years	...	Hess	{"Sector":"Energy","Industry":"Mining, Crude-O...
4	117057806156508	1 to 12 Years	...	Cairn Energy	{"Sector":"Energy","Industry":"Energy - Oil & ...

[5 rows x 23 columns]