

# Leveraging Data Analysis to Align Education and Employment

Dhanush C R

*Presidency University, Itgalpur, Rajankunte, Bengaluru - 560064*

**Abstract—** The disconnect between academic training and industry requirements is a key driver of graduate unemployability. This Data Analysis project leverages real-time data mining and analytics to bridge that gap. Job and internship postings from various platforms were analyzed alongside Training & Placement Officer (TPO) data from colleges. Using tools like Excel, Power BI, and Microsoft Azure, the project identified high-demand roles, skill gaps, and regional hiring trends. Key outcomes include targeted training recommendations for colleges, corporate training insights, and recruitment alignment, enabling dual revenue opportunities through education and employment services.

**Index Terms—** Data Mining, Data-Driven Training, Microsoft Azure, Power BI, Skill Gap analysis

## I. INTRODUCTION

The education-to-employment ecosystem in India faces a major challenge: a misalignment between academic curricula and evolving industry needs. Despite high educational qualifications, many graduates remain unemployable due to outdated syllabi, unstructured placement processes, and weak campus-industry linkages. Training and Placement Officers (TPOs), who serve as the bridge between institutions and companies, often lack centralized tools and real-time insights to drive effective collaboration.

This project leverages data analytics to address these issues by collecting and analyzing structured information on engineering colleges and corporate hiring patterns. Tools such as Power BI, Excel, PySpark, and Azure platforms were used to uncover trends in regional skill demand, hiring roles, and institutional readiness. Key outcomes include identifying curriculum gaps, mapping high-demand skills by region, and enabling smarter decision-making for colleges and companies.

The project proposes a scalable model that benefits all stakeholders: enhancing placement outcomes for

students, offering training services for colleges, and providing companies with access to pre-trained, job-ready talent—thereby creating sustainable value across the ecosystem.

## II. LITERATURE SURVEY

This chapter reviews existing research on how data-driven technologies can bridge the gap between education and employment. It emphasizes the increasing use of big data, cloud platforms, predictive analytics, and visualization tools to align academic training with industry demands.

The literature highlights the evolution of educational analytics, shifting from traditional metrics to real-time, dynamic insights. Studies show a growing need for curriculum alignment with in-demand skills, enabled by data mining techniques like web scraping and NLP-based preprocessing. These methods help in collecting institutional and job-related data, which is crucial for talent mapping.

Cloud platforms such as Microsoft Azure (ADLS, Databricks) support large-scale data storage and analysis, improving collaboration and reducing infrastructure dependency. However, privacy risks, technical barriers, and limited adoption in smaller institutions are notable challenges.

Predictive analytics is used to forecast employability, suggest upskilling paths, and match candidates to roles. While beneficial, these models face issues like bias and the need for frequent retraining. Visualization tools like Power BI and Tableau help translate insights into actionable strategies for educators, placement officers, and policymakers.

Overall, the literature underlines the transformative potential of integrated, data-centric systems in creating

responsive and effective education-to-employment pathways.

### III. PROPOSED METHODOLOGY

The proposed methodology harnesses data-driven technique using cloud computing, to provide a seamless, actionable solution that bridges the gap between education and employment. The goal is to enhance employability by analyzing educational programs and job market data, identifying skill gaps, and offering personalized recommendations to students, institutions, and employers. The methodology includes the following key components:

#### 3.1 Data Collection from Diverse Sources

Raw data is collected from two major domains:

- Educational Data: Includes academic curricula, placement records, certifications offered, and institutional metadata.
- Employment Data: Covers job descriptions, hiring trends, required skill sets, and regional demand for roles.

These datasets form the foundation for subsequent analysis.

#### 3.2 Data Ingestion into Cloud Storage

All collected data is uploaded to a scalable cloud-based storage platform such as:

- Microsoft Azure Data Lake Storage (ADLS) or
- AWS S3 Buckets

This step ensures centralized, secure, and scalable access to raw and semi-structured data, preparing it for transformation and analysis.

#### 3.3 Data Cleaning and Transformation using PySpark

The raw datasets often contain noise, inconsistencies, or duplicates. Using PySpark (Python + Apache Spark):

- Data is cleaned by removing nulls, correcting formats, and handling duplicates.
- Column transformation, type casting, and merging of datasets is done to create a clean, structured data layer.

- PySpark enables handling of large-scale datasets with distributed computing, improving speed and reliability.

#### 3.4 Structured Data Storage and Processing

Post-cleaning, the structured data is stored in cloud-based formats like Parquet, CSV, or SQL tables, enabling seamless integration with BI tools. This layer forms the analytics-ready base for dashboards.

#### 3.5 Data Visualization and Insight Generation

The transformed data is fed into BI tools such as:

- Power BI
- Tableau

These tools are used to:

- Build interactive dashboards
- Display regional hiring trends, skill gaps, and qualification requirements
- Support decision-making for colleges (curriculum alignment), students (career direction), and companies (targeted hiring)

#### 3.6 Insight-Driven Stakeholder Recommendations

Based on the visualized data:

- Students are guided on in-demand skills and certifications.
- Colleges receive insights on updating programs and improving placement efforts.
- Employers identify talent pools and collaborate with institutions accordingly.

#### 3.7 Cloud-Based Scalability and Real-Time Updates

Using cloud infrastructure allows:

- Real-time updates to the dataset
- Easy onboarding of more institutions or job roles
- Future integration with advanced analytics or AI models if needed

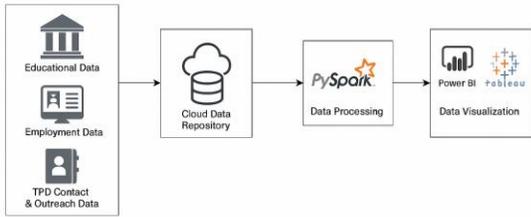


Fig. 1. System Architecture describing the end-to-end process

The system architecture outlines the complete data pipeline—from cloud-based data storage in Azure Data Lake, to data cleaning with PySpark in Azure Databricks, and finally, visualization using Power BI.

#### IV. IMPLEMENTATION

The implementation followed a modular and cloud-integrated approach. Below are the key steps and their sub-tasks:

##### 4.1 Cloud-Based Data Storage

- Azure Data Lake Storage Gen2 was used for organizing and storing data.
- Data was divided into containers:
  - raw/ → unprocessed data from sources
  - cleaned/ → data ready for analysis and export post PySpark processing
- Storage ensured scalability, access control, and integration with analytics tools.

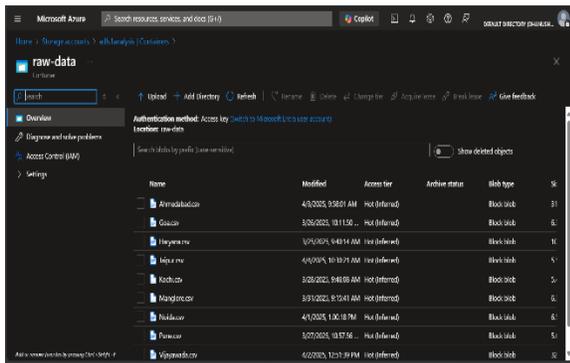


Fig. 2. Raw Data stored in Azure Data Lake Storage Gen 2

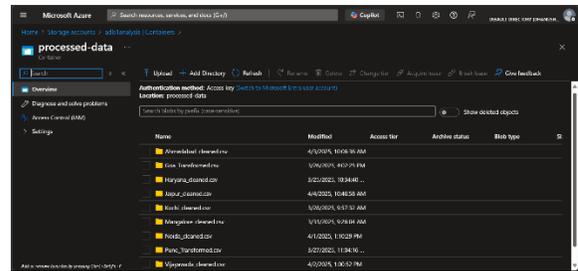


Fig. 3. Processed data in Azure is ready for visualization

##### 4.2 Data Ingestion

- Data was collected from various public and institutional sources.
- Ingested through Python scripts or manual uploads, depending on source format.
- Metadata (e.g., source name, ingest date) was added for traceability.
- Formats handled: .csv, .xlsx, .json.

##### 4.3 Data Cleaning with PySpark

- Executed within Azure Databricks using Python-based PySpark notebooks.
- Operations performed:
  - Dropping null/missing values and rows
  - Removing duplicates
  - Standardizing column names
  - Parsing and formatting date/time fields
  - Converting categorical values to lowercase for consistency

##### 4.4 Data Transformation & Integration

- Combined datasets (e.g., job data + college/training info + contact lists).
- Created calculated fields like:
  - Region-wise job counts
  - Top hiring domains by location
- Final data stored in .csv format for use in BI tools.

##### 4.5 Visualization with Power BI

- Connected processed data to Power BI using Azure connectors.
- Built interactive dashboards showing:
  - Top in-demand roles across regions
  - Required skills vs available training
  - College-wise placement insights
  - Filters for role, location, and qualification
- Charts used: bar graphs, heatmaps, slicers, maps.

4.6 Scalability & Performance Optimization

- Modular pipeline design allowed repeated use with updated data.
- Jobs can be scheduled or automated in Azure for near-real-time data refresh.
- Cloud storage and compute ensured performance on large datasets.

Technology / Tool	Description / Usage	Importance in Project
Azure Data Lake Gen2	Cloud-based storage used to store raw and processed data.	Centralized and scalable data storage solution for efficient access and management.
Azure Databricks	Cloud-based Apache Spark environment for big data processing and analytics.	Enabled distributed data processing and transformation using PySpark.
PySpark	Python API for Apache Spark used for cleaning, transforming, and preparing large data.	Performed scalable and efficient data preprocessing for analysis and visualization.
Power BI	Business intelligence tool used to create dashboards and visual reports.	Delivered interactive visualizations to communicate insights and hiring trends.
DAX (Data Analysis Expressions)	Formula language in Power BI used to create calculated columns and custom metrics.	Enabled creation of meaningful KPIs and enhanced dashboard interactivity.
Data Mining (Manual/Web)	Collected data on TPOs and hiring companies across cities (name, email, platform, role).	Built the foundation dataset required for further analysis and visualization.
Data Cleaning	Standardization, deduplication, null handling, and format unification.	Ensured data accuracy and consistency before further processing and storage.

Fig. 4. Technical Components Overview

The figure illustrates the tools and technologies used across the project pipeline—from data collection to visualization. Each component contributed to building a scalable and efficient system focused on connecting education with industry needs. The table highlights the purpose and role of each tool in the process.

V. RESULTS

The project delivered meaningful insights and functional outcomes that align with the objective of bridging the gap between education and employment.

5.1 Data Insights Generated

- Identified top in-demand job roles across different regions.
- Highlighted key skills frequently required by employers.
- Mapped job roles with related college training availability.

5.2 Skill Gap Identification

- Revealed mismatches between college curricula and industry needs.
- Showed regions where students lacked training for high-demand skills.

5.3 Visualization Outputs

- Developed interactive dashboards in Power BI to display:
  - Role-wise job demand
  - Location-based hiring trends
  - Qualification and eligibility filters

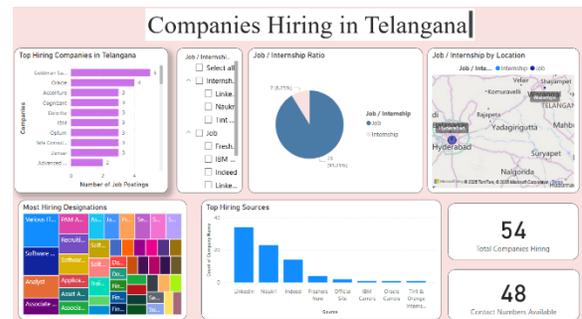


Fig. 5. Hiring Insights Dashboard

This Power BI dashboard visualizes hiring patterns in Telangana, showcasing top recruiting companies, job vs internship ratios, popular hiring sources, and most sought-after job designations. It provides regional insights to help align student training with current market demands.

5.4 Outreach Enhancement

- TPO (Training and Placement Officer) data enabled region-wise college targeting.

- Helped build a contact pipeline for training and placement collaboration.

#### 5.5 Platform Scalability

- The cloud-based setup handled large volumes of raw and processed data.
- Modular design allowed easy updates and integration of new datasets.

### VI. CONCLUSION

The project successfully demonstrates how a data-This project presents a practical and scalable approach to tackling one of the most pressing issues in the academic-to-industry transition: graduate unemployment due to misalignment between educational outcomes and market needs. By leveraging real-time data from various sources and analysing it through a structured cloud-based architecture, the system effectively identifies regional hiring patterns, skill demand trends, and gaps in college-level training programs.

The architecture employed in the project—utilizing Azure Data Lake for storage, Azure Databricks and PySpark for data processing, and Power BI for visualization—demonstrated high efficiency, scalability, and modularity. This framework enabled the team to clean, transform, and analyse large volumes of data with minimal latency and high accuracy. The integration of visual dashboards added immense value by making the insights accessible, interactive, and actionable for stakeholders such as Training and Placement Officers, college administrators, and hiring partners.

One of the major contributions of this system is its ability to extract meaningful insights from raw, unstructured job and institutional data. Through the analysis, the project highlighted the most in-demand roles and skillsets across regions, thereby helping colleges re-evaluate their curriculum and training focus. It also allowed companies to discover clusters of qualified candidates and institutions that could potentially meet their hiring needs.

Overall, the project successfully bridges the gap between academic offerings and industry expectations by building a data-backed ecosystem. It not only

supports students in gaining better placement opportunities but also opens dual revenue possibilities—via targeted training programs for colleges and recruitment services for companies. With continued data updates and system enhancements, this model holds strong potential to evolve into a sustainable platform for employment readiness and strategic academic reform.

### VII. FUTURE SCOPE

This project lays a strong foundation for data-driven decision-making in bridging the education-to-employment gap, and there are several opportunities to expand and enhance its impact in future iterations. One immediate direction is to scale the system geographically to include more states and regions across India. This would allow for a comprehensive national-level understanding of hiring trends and academic readiness, offering even more targeted insights for institutions and policymakers.

Another promising area for development is the integration of institutional data such as curriculum outlines, placement histories, and student skill assessments. By correlating this internal academic data with external hiring trends, the system can provide deeper recommendations—for example, suggesting specific curriculum enhancements or identifying training modules that align with high-demand job roles in a particular region.

Additionally, the platform can evolve into a dynamic dashboard for real-time monitoring of the job market. With automated data pipelines and scheduled refreshes, stakeholders could access up-to-date analytics on job postings, new hiring patterns, and emerging technologies.

This real-time capability would empower institutions to adapt quickly and remain aligned with industry developments.

Finally, incorporating feedback loops from recruiters and academic stakeholders can further refine the system's recommendations. Over time, machine learning models may be introduced to forecast job market shifts or suggest institution-specific strategies for improving placement performance. By continuously learning from new data and feedback, the system can become more intelligent, personalized,

and impactful in shaping the future of higher education and workforce development.

#### REFERENCE

- [1] Barton, D., Farrell, D., & Mourshed, M. (2012). "Education to Employment: Designing a System that Works." McKinsey & Company
- [2] Zhu, G., Kopalle, N. A., Wang, Y., Liu, X., Jona, K., & Börner, K. (2020). "Community Based Data Integration of Course and Job Data in Support of Personalized Career-Education Recommendations."
- [3] Fernandes, J. (2023). "The Role of Data-Driven Decision-Making in Effective Educational Leadership." *Academy of Educational Leadership Journal*, 27(S2), 1-3
- [4] Bhatia, M. K., Nigam, S., Adwani, S., & Kothari, D. (2024). "Applying Data-Driven Decision-Making to Academic Hiring Procedures: An HR Analytics in Higher Education." *Library Progress International*, 44(3).
- [5] Sharma, R. (2024). "Boosting Student Employability with the Effective Use of Data." Hurix Digital.
- [6] Martinez, K. (2021). "Data-Driven Decision-Making in Education: Using Data to Improve Instruction." *Journal of Arts, Society, and Education Studies*.
- [7] Liu, Y. (2024). "Data-Driven Decision Making in Higher Education Institutions: State-of play."
- [8] BibliU. (2023). "Data-Driven Decision-Making in Education: A Guide for Higher Learning Leaders."
- [9] McKinsey & Company. (2018). "Creating an Effective Workforce System for the New Economy."
- [10] McKinsey & Company. (2020). "Closing the Skills Gap: Creating Workforce Development Programs That Work for Everyone."
- [11] McKinsey & Company. (2021). "Creating Opportunities for Meaningful Employment."
- [12] National Association of Colleges and Employers. (2023). "Career Services Benchmark Survey."
- [13] World Economic Forum. (2023). "The Future of Jobs Report 2023."
- [14] OECD. (2022). "Education at a Glance 2022: OECD Indicators." OECD Publishing.
- [15] UNESCO. (2021). "Global Education Monitoring Report 2021: Non-State Actors in Education."