

International Journal of Research Publication and Reviews

Journal homepage: www.ijrpr.com ISSN 2582-7421

ETL Pipeline using Python and ETL Tools for Data Transformation on Datasets and Visualization

Wala Rahul Jayesh¹, Anil Kumar Kadam²

Department of Engineering, ME AI & DS, AISSMS College of Engineering, Pune, India, <u>rahulwala777@gmail.com</u> Department of Engineering, ME AI & DS, AISSMS College of Engineering, Pune, India, <u>ajkadam@aissmscoe.com</u>

ABSTRACT:

This paper describes the development of an Extract, Transform, Load (ETL) pipeline using Python and a variety of ETL tools to process and transform publicly available datasets. ETL pipelines play a key role in data engineering by streamlining the process of data extraction, cleaning, transformation, and loading, making it ready for analysis. In this study, we utilize Python libraries such as Pandas, Matplotlib, Seaborn, and Pandas Profiling to perform exploratory data analysis (EDA) and data visualization. We demonstrate this pipeline using a customer churn dataset and apply various transformation techniques, including handling missing data, feature engineering, normalization, and categorization. The paper further explores how visualizing transformed data can improve decision-making for businesses. Lastly, we discuss potential improvements and future scalability of ETL pipelines.

Keywords: ETL, Data Processing, Python, Data Transformation, Data Visualization, Exploratory Data Analysis (EDA)

Introduction

As the volume of data continues to grow across industries, the need for efficient, organized methods to process and analyze this data has become critical. ETL (Extract, Transform, Load) pipelines are essential for cleaning, transforming, and structuring raw data into meaningful insights. Businesses and researchers rely on these processes to ensure that their data is accurate and ready for analysis. This paper demonstrates how a Python-based ETL pipeline can automate data preprocessing and visualization for available datasets, showing the impact of such preprocessing on business decision-making.

Literature Review

The concept of ETL (Extract, Transform, Load) pipelines has been extensively examined in data engineering and business intelligence research. Numerous studies highlight the significance of ETL in areas like data warehousing, integration, and analytics. For example, Kimbal and Ross (2013) outline best practices for ETL in the context of data warehouses, addressing challenges such as maintaining data quality, managing transformation complexities, and ensuring scalability. Inmon (2019) also underscores the critical role of structured ETL processes in constructing robust enterprise data architectures.

In more recent research, attention has shifted toward Python-based ETL implementations, praised for their flexibility and powerful ecosystem. McKinney (2017) discusses how Python, particularly with the Pandas library, plays a central role in data preprocessing, offering an efficient way to manage and transform large datasets. Further studies, such as those by Lin et al. (2020), explore the integration of Python-based ETL pipelines with automation tools like Apache NiFi and Airflow, showcasing their application in scalable, cloud-based systems.

This study extends these foundational findings by demonstrating an automated, Python-driven ETL pipeline that incorporates advanced data transformation and visualization techniques. The goal is to enhance workflow efficiency and streamline structured data processing in modern data environments.

Methodology

3.1	ETL.	Pineline	Flow	Fig.	1.	ETL.	Pineline	Flow
2.1	LIL	1 ipenne	1,10.0	rig.			1 ipcinic	1.10 %

Extract Data	\rightarrow	Transform Data	\rightarrow	Load Transformed
(CSV Source)		(Cleaning, Feature Engineering)		Data to CSV/DB

3.2 Extraction Phase

The extraction phase involves retrieving raw data from various sources. In this study, we used the Deutsche Bank Customer Churn Dataset (2024) stored in a CSV file. The extraction process includes:

- Reading Data: The dataset is loaded into a Pandas DataFrame.
- Validating Data Structure: The consistency of the column names and data types is checked.
- Handling Missing or Corrupt Files: The pipeline includes error handling to manage missing or corrupted files during extraction.

3.3 Transformation Phase

Data transformation is crucial in turning raw data into useful insights. Several transformation techniques were applied to clean and structure the data:

- Handling Missing Data: Rows with critical missing values in fields like balance and estimated salary were removed. For other missing values, imputation methods were used.
- Feature Engineering:
 - Age Categorization: The age column was divided into categories like 18-25, 26-35, etc.
 - O Salary-to-Balance Ratio: A new feature was created by calculating the ratio of salary to balance.
 - Customer Tenure Classification: Customers were grouped based on the length of time they had been with the company to examine retention patterns.
- Standardization and Normalization:
 - 0 Numerical data was standardized to maintain consistent scales.
 - O Balance and salary values were normalized to ensure comparability across the dataset.
- Data Encoding and Type Conversion:
 - Categorical data, such as gender and geography, was transformed into numerical values using one-hot encoding.
 - Data types were adjusted to optimize storage and computation.
- Exploratory Data Analysis (EDA):
 - Automated reports using Pandas Profiling provided insights into the dataset, including statistical summaries, missing data counts, and correlation analyses. These insights helped refine the transformation process.

3.4 Loading Phase

The final phase of the ETL process involves loading the transformed data into a structured format for further use. In this case, the cleaned data is saved to a CSV file. Key steps in the loading process include:

- Validation: Ensuring the integrity of the transformed data before loading it.
- Saving to CSV: The final data is written to a CSV file, formatted correctly for analysis.

Future Scalability: The data can easily be loaded into SQL databases or cloud-based storage solutions for more extensive analysis.

Results and Visualization

After transforming the data, we created visualizations to uncover patterns and insights. Some of the visualizations included:

- Customer Age Distribution: A histogram depicting how customers are distributed across different age groups.
- Salary vs. Balance Relationship: A scatter plot showing the correlation between salary and balance.
- Churn Rate Comparison: A bar chart illustrating the number of active versus churned customers
- Geographical Distribution: A pie chart showing how customers are spread across different regions.
- Tenure vs. Churn: A box plot analyzing how customer tenure affects the likelihood of churn.

These visualizations provided valuable insights into factors that influence customer retention and financial stability, which can guide business decisions.

Few images of post_transformation_eda report for visualization

Exploratory Data Analy	ysis			0	verview Variables Interactions Correlations Missing values Sample
Ove	rview				
					Brought to you by YData
Overv	iew Alerts 3	Reproduction			
Datas	et statistics			Variable type	es
Numb	er of variables		13	Numeric	6
Numb	er of observations		10000	Text	1
Missin	g cells		D	Categorical	6
Missin	g cells (%)		0.0%		
Duplic	ate rows		0		
Total s	ize in memory		2.6 MiB		
Averag	ge record size in men	iory	268.4 B		
Exploratory Data Anal	vcic			c	Overview Variables Interactions Correlations Missing values Sample
	,				
Vari	مامامه				
Varia	ables				
Select Co	olumns 🗸				
Custo	omerid				
Real nu	mber (R)				
Unique					
Disting	et 1	10000	Minimum	15565701	adua Infationalia aduation
Disting	rt (%)	100.0%	Maximum	15815690	
Missin	g (2	Zeros	0	
Missin	g(%) (2.0%	Zeros (%)	0.0%	# # L & P
lanat	- 403	0	Negative	0	
Mean	e (%)	15690941	Memory size	78.2 KiB	
	Menery d CreditS Rel softe Delicet	2019 c (8) 4/2	cité de cité		
Emloratory Data Analysis	Distinct (1) Missing Missing (1) Indiates (6) Misson	0 405 0 105 0 105 0 105 105 105	Nasimen Zeres (N) Higgstus Higgstus (S) Nerrory siz	651 0 0.0% 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0	Normal Content forces and the states the
Diploting even may					
	Geography Categorical				
	Distinct		1	france Germany	5014
	Distinct (%) Missing		< 0.1% 0	Spain	2477
	Missing (%)		00%		
	Memory size		615.4108		Note details
	Gender Categorical				
	Distinct		2	Male	5457
	Distinct (%)		< 0.1%	Female	4543
	Minning (%)		0.0%		
	Memory size		664.7 835		
					More details
	Age				
pioratory Data Analysis					Overview Variables Interactions Correlations Missing values Sam
	Age Bui number (II)				
	Distinct Distinct (%)	79	Minimum	18	.lli
	Missing	٥	Zeres	0	
	Minning (%)	0.0%	Zares (%) Negative	0.0%	
	infinite (%)	0.0%	Negative (%)	0.0%	
	Mean	38.5216	Melhory size	782.KE	More details
	Tenure Real number (II)				
	Distinct	11	Minimum	0	
	Distinct (%)	0.1%	Maximum	10	
	Missing Missing (%)	0	Zeros Zeros (X)	413 4.1X	
	Infinite	0	Negative	0	5 1 2 4 4 3
	Infinite (%)	0.0%	Negative (%)	outre.	
	Mean	540328	Memory size	182.84	



fig Fig. 2 - post_transformation_eda report for visualization

Conclusion and Future Directions

This paper successfully demonstrates how a Python-based ETL pipeline can automate the extraction, transformation, and visualization of data. It highlights the power of Python in preprocessing and analyzing data, offering a robust solution for handling various datasets. The pipeline could be extended with cloud-based tools for scalability.

Future directions include:

- Real-time Data Processing: Implementing ETL pipelines that process data in real time using technologies like Apache Kafka.
- Anomaly Detection: Enhancing the transformation phase by incorporating machine learning models to identify anomalies.

Cloud Integration: Deploying the ETL pipeline on cloud platforms, such as Google BigQuery or AWS Redshift, for large-scale analytics.

Acknowledgements

The authors would like to express their sincere gratitude to their mentors and faculty members for their valuable guidance and support throughout this project. Special thanks to the institution for providing the necessary infrastructure and resources. We also acknowledge the developers and maintainers of the open-source tools and libraries such as Python, Pandas, Matplotlib, Apache NiFi, and Talend, which were instrumental in implementing this work.

REFERENCES

[1] A. Katal, M. Wazid, and R. Goudar, "Big Data: Issues, Challenges, Tools and Good Practices," in 2013 Sixth International Conference on

Contemporary Computing (IC3), 2013.

[2] Apache NiFi Documentation. https://nifi.apache.org/docs.html

[3] Talend Open Studio Documentation. https://help.talend.com

[4] Wes McKinney, 'Python for Data Analysis', O'Reilly Media, 2017.

[5] Hunter, J. D. (2007). Matplotlib: A 2D Graphics Environment. Computing in Science & Engineering.