

---

PROFESSIONAL SUMMARY

---

Data Engineer with nearly 4 years of experience who turns messy source feeds into fast, reliable analytics. At DRG, I built metadata-driven Azure Data Factory pipelines ingesting ~1M rows/week and cut daily refresh ~50% with monitored execution and alerting. I also improved PostgreSQL performance ~35% through partitioning and indexing, boosting report responsiveness. Outside work, my Databricks projects (Autoloader, Delta Live Tables, Unity Catalog) showcase production-ready Delta Lake patterns.

---

SKILLS

---

- **Cloud/Platform:** Azure (ADLS Gen2, Azure SQL Database, Key Vault), Databricks (DLT, Autoloader, Structured Streaming, Unity Catalog, SQL Warehouse)
- **Data Engineering:** ETL/ELT, CDC, SCD Types, Dimensional Modeling (Star/Snowflake), Delta Lake, Medallion Architecture, Data Quality (DLT Expectations), Lineage/Observability
- **Orchestration/Automation:** Azure Data Factory (pipelines, Mapping Data Flows, SHIR), Azure Logic Apps, GitHub Actions, Databricks Asset Bundles (CI/CD)
- **Programming/Query:** Python (Pandas, NumPy, PySpark, SQLAlchemy, Flask), SQL/T-SQL (CTEs, window functions, performance tuning), Bash, DAX, PostgreSQL, MySQL, REST/HTTP APIs, JSON/XML
- **Analytics/BI:** Power BI (DAX, RLS, Incremental Refresh, semantic models), Excel (Pivot Tables, VBA Automation), Grafana, Matplotlib, Seaborn
- **Data Acquisition:** REST APIs, BeautifulSoup, Selenium, Web Scraping, Managed Identity Authentication, Python Automation Scripts
- **Other:** Docker, Git/GitHub, Documentation & Runbooks

---

WORK HISTORY

---

**Data Analyst and Engineer****May 2024 – Present****The DRG, United Kingdom**

- Designed and optimized robust, scalable database and **schema** models in collaboration with Finance and key stakeholders, ensuring **data integrity**, accessibility, and alignment with business needs.
- **Delivered** parameterized ADF pipelines processing **~1M+ rows/week** from EPOS, loyalty, and REST APIs; standardized schemas with **idempotent loads** and drift-tolerant ingestion.
- **Cut** daily refresh time **~50%** by **parallelizing** safe paths and removing redundant stages; reduced end-to-end runtime and increased on-time report readiness for Finance/Ops.
- **Migrated** MySQL to **Azure PostgreSQL** and **re-modeled** facts/dims with targeted **BTREE indexes** and date/location **partitioning**, reducing key report query times **~35%**.
- Automated daily data extraction from **RESTful APIs** and web scraping (**BeautifulSoup**, **Selenium**) and converting JSON, XML data into structured formats, improving database performance and usability.
- **Modeled** governed gold datasets for **Power BI** with **RLS** and **Incremental Refresh**, increased dashboard adoption and reduced refresh failures per week.
- **Partnered** with Marketing/Ops to analyze seasonality and run **A/B tests**; **translated** findings into campaign targeting and service planning.
- Implemented automated **Power BI** dashboards for real-time reporting and strategic decision-making, increasing quarterly revenue by 15% through enhanced sales insights.

**Data Analyst****April 2021 - March 2022****NSEIT, Nashik**

- **Automated** SQL/Python ingestion and preprocessing with parameterized validations (types, duplicates, referential checks), **freeing ~10 hrs/week** and shortening release cycles.
- **Co-designed** a governed preprocessing framework with a **6-person** team (standardized schemas, data dictionary, access controls), **reducing rework** and handoffs across monthly releases.
- Automated recurring reconciliations and data-processing workflows, freeing **~10 hours/week** and shortening analysis turnaround for stakeholders.
- Generated comprehensive 20+ reports and dashboards to support informed decision-making by stakeholders.
- Monitored data integrity and quality, promptly addressing discrepancies and improving reporting reliability.

**Data Analyst****March 2020 – March 2021****Swami Vivekanand Education and Research Centre, Mumbai**

- Integrated and standardized multi-source data with Python (Pandas/NumPy), cutting processing time by ~20% and improving data reliability for downstream reporting.
- Built interactive **Power BI** dashboards for KPI tracking, enabling self-service analysis and faster decision cycles.
- Authored and optimized complex **SQL** to profile and segment customers, delivering a ~25% uplift in decision-ready insights for stakeholders.
- Increased stakeholder engagement by ~30% through targeted walkthroughs, clear documentation, and iterative dashboard improvements.

---

## EDUCATION

---

- MSc in **Advance Computer Science with Data science (2:1)** University of Strathclyde, School of Computing Science January 2023 – January 2024
- Bachelor of **Computer Applications (2:1)** Savitribai Phule Pune university May 2016 – March 2020

---

## PROJECTS

---

### **End-to-End Azure Data Engineering Lakehouse ([GitHub Link](#))**

- Medallion (Bronze/Silver/Gold) on **ADLS Gen2**; parameterized, **metadata-driven ADF** with incremental loads, backfill controls, and JSON CDC watermarking.
- **Databricks Autoloader + Structured Streaming** refine Bronze to Silver **Delta**; **Unity Catalog** governance (metastore, external locations, access connector).
- **Delta Live Tables** for Gold: **SCD 2 dims** (sequence\_by), SCD1 fact upserts, **Expectations + visual lineage DAG**; analytics-ready for **SQL Warehouse**/Power BI.
- Ops & DevEx: Python utilities (cleaning, de-dupe), checkpointing & idempotent runs, empty-file cleanup, **Logic Apps** alerts, secrets via **Managed Identity/Key Vault**, CI/CD with **Databricks Asset Bundles** + GitHub

### **DataBricks Declarative Pipelines ([GitHub Link](#))**

- Built End-to-end **DLT Lakehouse** with **Medallion**; **unified batch + streaming** and automatic **DAG/lineage**.
- **Bronze**: streaming ingests with **Expectations** (warn/drop/fail) and Append Flow to merge multi-region sources; runtime audit metrics.
- Silver: **Auto-CDC** Type-1 upserts using reusable Python **utilities** (casting, de-dup, common transforms) for **incremental** accuracy.
- **Gold & Ops**: **SCD2** dimensions, Type-1 fact, and a **materialized revenue view**; **Unity Catalog governance**, parameterized dev/prod configs, **SQL seed scripts**, and alerting.

### **Azure Data Factory End-to-End (On-Prem + API + SQL): ([GitHub Link](#))**

- ADF hybrid ingestion: **SHIR (on-prem shares)**, **HTTP/REST**, and Azure SQL **incremental loads** without watermark tables to **idempotent re-runs** and drift-tolerant landings.
- Orchestrated Bronze, Silver, Gold: **ADF Mapping Data Flows** standardize/clean/type-cast to **Silver Delta** with Alter Row upserts on business keys; **Gold views (joins, aggregations, dense-rank, Top-N)** optimized for **BI**.
- Centralized control via a **parent pipeline** chaining on-prem/API/SQL paths, **parameterizing file lists & batch sizes**, and parallelizing safe steps; less manual coordination and **predictable backfills/replays**.
- Shipped **ARM publish artifacts** and runbooks for **repeatable CI/CD**, with guardrails on **schema-drift vs validation, partitioning** options, and Gold snapshot overwrite strategy.

### **DBT Databricks Scd2-Data Quality: ([GitHub Link](#))**

- Built a **Medallion** model with **dbt Core** on Databricks using env-driven **profiles.yml** and per-layer materializations.
- Enforced data quality with **dbt tests**: generic, **singular**, and **custom generic**.
- Implemented **SCD2** via **dbt snapshots** (timestamp strategy) for **time-travel analysis**.
- Drove reuse & speed with **Jinja macros**, **seeds**, **incremental** models, node selection, and **dbt build + docs/lineage** for CI.