

EDUCATION

University of Pennsylvania (School of Engineering and Applied Science)

Philadelphia, PA

Master of Science in Engineering (MSE) in Data Science; GPA 3.9/4.0

Aug 2023 – May 2025

Coursework: Databases/Big Data Analytics, Machine Learning, NLP, AI, Computer Vision, Deep Learning, Statistical Modeling

Kirori Mal College (KMC – University of Delhi)

Delhi, India

Bachelor of Science (BS) in Statistics; GPA: 8.11/10

July 2017 – Aug 2020

EXPERIENCE

ML Researcher, [Computational Social Listening Lab \(UPenn\)](#)

May 2024 – Present

Project 1: (Misinformation)

- Worked on identifying misinformation on social media and whether it relates to health outcomes segmented by race.
- Extracted linguistic features from posts using [DLATK](#) and applied LDA for topic modeling and performed correlation analysis.
- Built NLP pipelines to detect health-related misinformation from ~20K social media posts using RoBERTa and LDA topic modeling; improved **classification precision by 20%** through alignment-based entailment.
- Unified data from various online survey platforms (**10K+ responses**) in a secured server via MySQL and performed feature engineering in Pandas to prepare data for further downstream tasks.

Project 2: (IH Risk Model)

- Developed a pipeline to predict Incisional Hernia (IH) risk from unstructured operative notes and structured EHR data in a surgical cohort of 10k+ patients.
- Engineered a few-shot GPT-based extraction pipeline to label operative features (e.g., incision, ostomy) from redacted notes **and reduced noisy extractions by 30% vs. BERT embeddings**.

Data Science Intern (Full-time), [Universal Media \(PA, USA\)](#)

May 2024 – Aug 2024

- Architected Azure SQL Database solutions, encompassing DDL scripts to enhance data management and reporting solutions.
- Led the development of **3+ data pipelines** using Azure Data Factory (ADF), facilitating the seamless ingestion and transformation of diverse data sources into the Azure environment.
- Developed python scripts for data transformation, stored them in Blob storage and executed them via batch activity in ADF.
- Drove product marketing insights by building **Mixed Media models (MMM) in Azure Synapse**, analyzing marketing channel impacts on media diversity metrics. Built **Power BI** dashboards to deliver actionable insights for optimizing client strategies.
- Authored **5 stored procedures** in SQL, automating repetitive tasks and improving query performance by over 30%.

Assistant Manager (Full-time), [IIFL Finance Ltd](#)

Apr 2022 – July 2023

- Analyzed ETL process failures and created **10+** paginated reports in SSIS to help the management track 1000+ branches.
- Optimized & migrated complex SQL queries from an obsolete database server that improved the **reporting services by ~40%**.
- Digital Adoption:** Led a product-focused initiative to identify digitally savvy customers by engineering features and building ADF pipelines to track campaign behavior. Trained and deployed a Random Forest model (**with a 90% accuracy**) in Azure ML Studio; **exposed it as a REST endpoint consumed by marketing campaigns**, driving digital **disbursal adoption by 50%**.

SELECTED PROJECTS

- Azure ETL (2025):** Built a scalable ETL pipeline in Azure Data Factory that was capable of ingesting and transforming 1M+ rows daily. Ingested raw data into Data Lake Storage Gen2, achieved **~50% reduction in query latency** via Azure Databricks optimizations, and analyzed it in Synapse Analytics - delivering a seamless end-to-end solution. [\[Link\]](#)
- FitBit(2024):** Engineered a Django health chatbot leveraging PostgreSQL for robust patient data management, featuring an LLM-agnostic architecture with seamless model switching via **Langchain that reduced overhead by 40%**. Optimized memory usage for handling long conversations and implemented entity extraction to dynamically enhance medical context. [\[Link\]](#)
- Diffusion Transformer (2024):** Implemented PatchVAE with convolutional encoders and patch-based decoding for fine-grained feature extraction. Trained a Diffusion Transformer to sample from the latent space of PatchVAE, achieving a **30% reduction** in FID score and **2x greater feature diversity** compared to VAE-generated samples. [\[Link\]](#)
- Statistical Segmentation (2024):** Modeled customer purchase behavior using Pareto II and Weibull distributions; optimized parameters using Excel Solver and back tested on historical unit sales, improving segment alignment with **actual revenue by 21%** and informing future campaign targeting. (Solver, Marketing). [\[Link\]](#)

TECHNICAL SKILLS

Programming Languages: Python, SQL, C/C++, R programming, JavaScript

Databases/Frameworks: MySQL, PostgreSQL, SSMS, MongoDB, Neo4j, PyTorch, React, NodeJS, A/B Testing, PySpark, Django

Cloud/Big Data Orchestration: AWS (S3, Glue), Azure (Data Factory, Synapse, DevOps), Docker, DataBricks, Kafka, Airflow, DBT