

Samhitha Nuka

✉ samhithanuka51@gmail.com ☎ +1 812-361-5059 📄 SamhithaNuka 🌐 github.com/snuka75 📁 Portfolio

PROFESSIONAL SUMMARY

Data Science professional with a Master's degree from Indiana University and over two years of industry experience in machine learning, data engineering, and AI-powered solutions. Skilled in designing scalable ETL pipelines, developing serverless applications, and implementing real-time analytics across AWS, Azure Databricks, and Snowflake platforms. Certified in AWS Machine Learning, Databricks Data Engineering, and Google Data Analytics, with hands-on expertise in building LLM applications, cloud-based deployments, and data-driven research initiatives.

EDUCATION

MS in Data Science, Indiana University Bloomington Aug 2023 – May 2025 | Bloomington, United States
GPA: 3.6/4

SKILLS

Data Engineering: ETL/ELT (Databricks, Spark, Informatica), Data Warehousing (Snowflake, Redshift, BigQuery), Stream & Batch Processing (Kafka, Hadoop), NoSQL (MongoDB, Cassandra), Orchestration (Airflow), **Machine Learning:** Supervised and Unsupervised Learning Algorithms, Long Short term Memory (LSTM), NLP, CNN, LLM (Gemini Pro, Gemma, Groq), Transformers, Random forest, K-mean Clustering, SVM., **Generative AI & Agentic Systems:** LLMs (Gemini, Gemma, Groq), Prompt Engineering, RAG, FAISS, LangChain, Agentic AI Workflows, **Data Analysis & Visualization:** Excel (Pivot Tables, Power Query, VLOOKUP), Tableau, Power BI, Looker Studio, Pandas, NumPy, Matplotlib, Seaborn, **Programming:** Python, SQL, R, Java, JavaScript, PHP, HTML, CSS, Bash

WORK EXPERIENCE

Indiana University Bloomington, Graduate Research Assistant Jan 2025 – present | Bloomington, United States

- Built a Python pipeline to process and label over **10,000 time-series events**, preparing data for classification based on temporal sequence patterns.
- Engineered 15+ statistical features (e.g., inter-event intervals, rolling variance, delta rates) and trained Random Forest, SVM, and Logistic Regression models, achieving **92% accuracy** with **5-fold cross-validation** and **GridSearchCV**.
- Modeled temporal trends using linear regression, yielding an **average R² of 0.87** and identifying timing anomalies through residual analysis.
- Created 10+ visualizations with Matplotlib and Seaborn to explain model performance, feature contributions, and temporal behavior to non-technical stakeholders.

Cognizant Technology Solutions, Junior Machine Learning Engineer Aug 2021 – Jul 2023 | Chennai, India

- Catalyzed the integration of OpenText ECM with AWS SageMaker, creating **42+ RESTful APIs** to automate data extraction and analysis, establishing a reusable integration standard.
- Enhanced legacy systems to support ML workflows, resulting in a **30% improvement in processing efficiency** and a **40% boost in search index speed**.
- Orchestrated **serverless inference endpoints** using API Gateway and Lambda, sustaining **<100ms latency**, **99.9% uptime**, and handling **500+ concurrent requests daily** for reliable ML pipeline execution.
- Implemented real-time monitoring with Amazon CloudWatch and automated deployment with AWS CodePipeline, accelerating model release cycles and ensuring **production-grade MLOps**.

ACADEMIC PROJECTS

YOLOv7-Based Object Detection for Face Mask Compliance, YOLOv7 | Computer Vision | OpenCV | CNN | Tensorflow 🔗

- Fine-tuned YOLOv7 on a custom Kaggle face mask dataset annotated using CVAT, **achieving ~84% mAP@0.5** and strong real-world detection accuracy.
- Optimized training pipeline with **data augmentation**, **hyperparameter tuning**, and performance tracking using TensorBoard.

End-to-End AI Application Suite for Multimodal Interaction, Streamlit | Gemini Pro/Flash | FAISS | Groq | LLM 🔗

- Developed a suite of 4 AI-powered web applications using Streamlit, enabling natural language Q&A, PDF analysis, image interpretation, and no-code SQL query generation with Gemini Pro/Flash and SQLite.
- Integrated Groq and Gemma for high-speed document processing with response times as low as **30–50 milliseconds**, and leveraged **FAISS** to enable semantic search across 30+ documents with sub-second retrieval times.

NYC Crime Analytics Dashboard, Cloud-Based Public Data Pipeline, BigQuery | SQL | Looker Studio | ETL | Feature Engineering | GCP 🔗

- Built an end-to-end ETL pipeline to ingest and transform **7M+ NYC crime records** using 12+ SQL-based feature engineering steps.
- Designed a 3-page interactive **Looker Studio** dashboard with 20+ visualizations to analyze crime trends by time, location, and demographics.

NIFTY - 50 Stock Market Prediction, LSTM | Financial Market Analysis | Time-Series Forecasting 🔗

- Refined a predictive model using Long Short-Term Memory (LSTM) networks to forecast Nifty 50 stock movements within **Automotive and Finance**, with hyperparameter tuning achieving **95% accuracy**, showcasing LSTM capabilities.

Credit Card Fraud Detection with PCA, SMOTE, and Ensemble Learning, PCA | SMOTE | XGBoost | Ensemble Learning 🔗

- Built fraud detection pipeline achieving **95.6% F1-score**, using PCA and SMOTE on imbalanced datasets.
- Compared Random Forest, XGBoost, and Logistic Regression, deploying a stacked model that **Increased recall by 12%**.

Scalable Data Pipeline Development for Real-Time Traffic Analytics, Azure Databricks | SparkSQL | CI/CD | ETL 🔗

- Engineered a real-time ETL pipeline handling **100,000+ traffic records/day**, using Spark SQL and Medallion Architecture on Azure.
- Designed and deployed 3 Power BI dashboards that improved decision-making speed by **30%**.
- Automated deployments via Azure DevOps, reducing manual deployment time by **80%**.

CERTIFICATIONS

• AWS Machine Learning Engineer – Associate 🔗	• Databricks Data Engineer Associate 🔗	• Google Data Analytics Professional Certificate 🔗	• IBM Python for Data Science and AI 🔗
• Machine Learning offered by Stanford University 🔗	• Rest API (Intermediate) assessment in HackerRank 🔗	• OpenText Content Server Developer v20.4	