

V L Sivasai Mani Harshith Bhattaram | NY, USA | +1(716)232-0462

harshithbhattaram@gmail.com | <https://www.linkedin.com/in/harshith68/> | <https://github.com/maniharshith68/>

SUMMARY

Data Scientist with 2.5 years building end-to-end Data pipelines | XGBoost · ETL · NLP · SHAP · LIME Explainability · Hypothesis Testing | Predictive modeling across fleet telematics, finance and healthcare domains | Python · XAI · Scikit-Learn · AWS · Render

EDUCATION

University at Buffalo, The State University of New York, USA
Master of Science, Computer Science

Feb 2026

WORK EXPERIENCE

LocoNav Fleet Management Solutions Software Engineer (Data Scientist)

Gurugram, India
Jun 2022 - Jun 2024

- Built a driver drowsiness and intoxication classifier using OpenCV, Tensorflow on 10k+ video frames, achieving 87% accuracy.
- Reduced false negative drowsiness detections by 25% by incorporating facial landmark tracking and yawn angle-based alertness.
- Cut model inference response latency by 20% via optimized production pipeline that triggered real-time low-latency in-app alerts.
- Developed GBM, LSTM ETA prediction models on GPS, telematics data across 50k+ trips, improved forecast accuracy by 15%.
- Reduced average trip duration by 12% via graph-based route optimization using real-time traffic features, historical travel data.
- Collaborated with cross-functional teams in Agile workflows, supporting planning, stand-ups, and retrospectives.

Software Engineer (Data Science Internship)

Jan 2022 - May 2022

- Built a predictive maintenance pipeline from time-series telemetry data, handled class imbalance in failure labels via SMOTE.
- Achieved 85% precision in component failure prediction using ensemble classifiers with engineered degradation features.
- Reduced unplanned vehicle downtime by 18%, deployed automated ML alert pipeline integrated with fleet operations dashboard.

PROJECTS

Credit Card Customer Churn Prediction System with Explainable AI Dashboard | [Link](#)

- Built end-to-end ML pipeline using XGBoost and Logistic Regression on 10k+ customers on kaggle credit card dataset.
- Achieved ROC-AUC of 0.99 and 90.5% recall after engineering on 21 features and validating feature design with TreeExplainer.
- Deployed an interactive Streamlit application which includes EDA, model performance, and churn probability per customer.
- Integrated SHAP explainability for global feature importance & waterfall plots for model interpretability to show key drivers.

Healthcare Patient Readmission Risk Modeling with Explainability Report | [Link](#)

- Analyzed the UCI diabetes readmission dataset (100k+ rows) with SQL profiling queries to find readmission rates across cohorts.
- Engineered 44 features, resolved class imbalance using XGBoost instead of SMOTE to avoid interpolation errors on mixed data.
- Trained a gradient boosted classifier, achieved statistical improvements, validated 11 risk factors via chi-square & Welch's t-tests.
- Autogenerated clinical decision report with SHAP waterfall charts, confusion matrix, feature importance and deployed on cloud.

NLP Sentiment Analysis & Topic Modeling Platform for Product Reviews | [Link](#)

- Designed a full NLP pipeline processing 24k+ Amazon reviews via VADER and DistilBERT sentiment analysis, LDA modeling.
- Engineered modular text preprocessing via spaCy lemmatization & custom cleaning logic, parquet schema with 20 columns.
- Reduced 27+ raw records to 24k analysis ready rows, tracked sentiment scores, topic labels & model agreement via NLP model.
- Deployed a 4-page interactive Plotly dashboard on streamlit with time-series trend analysis by sentiment, visualizations per topic.

Clinical Radiology Report Generation via Multimodal Retrieval-Augmented Generation (RAG) | [Link](#)

- Architected a multimodal RAG pipeline on the MIMIC-CXR dataset (250k+ chest X-ray image-report pairs).
- Used BioMedCLIP (ViT-B/16) for cross-modal embeddings, retrieval via FAISS indexing, dual layer generation module.
- Orchestrated a memory-constrained subprocess architecture to get <1second query latency on knowledge base.
- Evaluated with CheXbert, BLEU, ROUGE-L, BERTScore to quantify accuracy against ground-truth reports.

SKILLS

- **Languages and Tools:** Python, SQL, PostgreSQL, R, Git, Tableau, Plotly, Streamlit, Render
- **Data Handling:** Data Cleaning, Data Preprocessing, Structured data, Unstructured data, EDA, Feature engineering
- **Machine Learning & Statistics:** Scikit-Learn, XGBoost, Gradient Boost, Random Forest, LSTM, ARIMA, Transformers, Numpy, Keras, Cross-Validation, Hypothesis Testing, SMOTE, PCA, Pandas, Time-Series Analysis, Matplotlib
- **Deep Learning & NLP:** TensorFlow, PyTorch, DistilBERT, BioMedCLIP, OpenCV, BERTopic, TF-IDF, LDA, Named Entity Recognition (NER), SHAP, LIME, Hugging Face
- **Generative AI:** RAG, FAISS, Prompt Engineering, QLoRA, LangChain, LangGraph
- **LLMs:** LLM APIs (Gemini, OpenAI, Groq, Llama), Vector Databases, VLMs, Fine-tuning
- **Evaluation & Cloud:** ROC-AUC, ROUGE, BLEU, BERTScore, FastAPI, Docker, AWS (S3, Lambda, RDS), Azure, GCP