

Bhanu Prakash Reddy

bhanuprakash2511@gmail.com (214) 218-6957 Dallas, Texas LinkedIn

SUMMARY: Dynamic and results-oriented Data Scientist with over 7+ years of experience in leveraging data-driven insights to drive business improvement and innovation. Proficient in a range of advanced analytical techniques and tools – Python, SQL, Pyspark, including Predictive Analytics, NLP, Time Series Forecasting, Gen AI, AB Testing, Statistical Analysis, Big Data and Business Intelligence tools. Demonstrated expertise in developing and implementing end to end predictive models and algorithms that have successfully informed business strategies and solutions.

PROFESSIONAL WORK EXPERIENCE:

Citigroup Inc., Dallas, USA

Sr. Generative AI Engineer (Deep Learning, Gen AI, Python, SQL, Fast API, OpenShift/Kubernetes) Apr 2024 – Present

- AI/ML & Platform Engineering: Engineered and deployed a Hybrid AI Engine (Traditional ML and Gemini 2.5 Pro GenAI) to automate critical regulatory report tagging and rationale generation. Scaled the solution to support 26 high-value financial reports (18 Risk reports via ML, 8 Finance reports via GenAI), processing complex instruction sets up to 750 pages (10 MB).
- Optimized Traditional ML performance for 18 risk reports, conducting rigorous experimentation with Random Forest, XGBoost, SVM, and Logistic Regression. Deployed tailored models with performance-enhancing techniques like SMOTE, achieving a recall of about 0.78.
- Advanced GenAI solution design for 8 Finance reports, implementing advanced patterns including Retrieval-Augmented Generation (RAG) and similarity-based few-shot learning. Final deployment on Gemini 2.5 Pro resulted in a 70% improvement in automated rationale quality and interpretability.
- Drove end-to-end MLOps and Deployment on Openshift, configuring deployment.yaml, cronjob.yaml, and values.yaml for automated scheduling and scaling. Architected and exposed the model via multiple high-performance FastAPI endpoints.
- Maximized throughput and efficiency by integrating Async and Threadpool parallel processing with LangChain for GenAI workflow orchestration, resulting in reduction in end-to-end latency for report processing.
- COBOL Modernization & LLM POC: Led a rapid COBOL Code Generation Proof-of-Concept (POC), leveraging LLMs to transform data across multiple source/target copybooks based on complex mapping rules. Successfully generated runnable COBOL code covering 5,000+ data mapping rules (80% simple, 20% complex) in under one week, demonstrating a potential 50x acceleration in migration development cycles.

Globe Life Insurance, Dallas, USA

Sr. Data Scientist (Gen AI, Claude, Python, SQL, Tableau, QuickSight, AWS Glue, Redshift) Apr 2024 – Mar 2025

- Developed a RAG chatbot for Globe Life Insurance – Sales Team using fine-tuned Claude models on AWS Cloud enabling the content creation, Question Answering, text to text generation, Classification and optimizing language understanding and generation for enhanced user interactions.
- Applied chunking strategies to effectively process and analyze large blocks of text for more coherent and contextually relevant responses.
- Utilized vector representations of text and semantic search techniques to enhance the accuracy and speed of query processing, leading to more satisfying user experiences.

Metamorphix Inc, Remote, USA

Data Scientist (NLP, Python, SQL, Pyspark, ADF, Databricks, A/B Testing) Jul 2023 – Dec 2023

- Developed ensemble models using NLP techniques, achieving a 20% performance improvement through hybrid feature engineering module, sentiment aware embedding layer (BERT) and error analysis module.
- Designed scalable data pipelines using Apache Kafka and Spark Streaming processing 500GB+ daily for low latency real-time sentiment analysis at scale.
- Validated context-aware recommendation engines that drove a 20% increase in user engagement on Cascades Platform by incorporating real-time user data such as location, time of day, and device type to deliver personalized recommendations.

- Spearheaded A/B testing initiatives for new UI/UX features, achieving a 10% uplift in conversion rates. Automated data preprocessing workflows, reducing manual effort by 25% and accelerating model deployment cycles. Additionally, developed a Sequential Probability Ratio Testing (SPRT) module for continuous experiment monitoring, enabling early detection of significant results.
- Collaborated across departments to implement ML-driven product features benefiting 70,000+ active users at Metamorphix - Cascades Platform.

Yoamigos Webservices, India

Jun 2017 – Jul 2022

Senior Data Scientist (NLP, Timeseries, Classification, Regression)

- Detecting Industry Trends on social media using NLP: Processed and analyzed extensive data, including 2+ million unique F&B products in 125+ categories. Utilized techniques like parsing, lemmatizing, and POS tagging, to process tweets. Identified precise topics and detected trends accurately using Bayesian classifiers, fuzzy matching algorithms, and Latent Dirichlet Allocation models.
- Forecasted sales for more than 300 products across 5 retailers in multiple countries resulting in trade promotion optimization (TPO). The built Auto Regressive Distributed Lag (ARDL) model achieved 13% WMAPE that far exceeded the accuracy of manual forecasts.
- Deployed and productionized LightGBM and Ordinary Least Squares models on Azure platform, trained on data normalized data using ALS factorization, to recommend product placement for retail stores in various European markets on a weekly basis.

Data Scientist

- Designed a payment integrity solution to classify claims into different levels of risk using Support Vector Classifier (SVC) with a recall of 91% at AUC of 0.95. Integrated the SVC model into the existing software development framework. Reduction in false positives increased the efficiency of auditors by 40%.
- Customer segmentation using RFMT model: Applied the Recency, Frequency, Monetary, and Time (RFMT) model to analyze patient data. Employed K-Means, Gaussian, and DBSCAN algorithms to classify patients into distinct groups. Conducted cluster factor analysis using methods such as elbow, dendrogram, silhouette, Calinsky–Harabasz, Davies–Bouldin, and Dunn index to ensure the meaningfulness of patient groupings. Implemented the majority voting (mode version) technique to select the most relevant patient clusters.
- Increased customer lead utilization by 20% by using Random Forest to classify untouched leads as potential customers or not. Achieved accuracy of 98% and F2 score of 35%.

ACADEMIC PROJECT EXPERIENCE

- Utilized **BERT model** to identify SEO keywords and built a model to rate their importance based on search volume and rank using cosine similarity. Prompt-engineering experiments conducted in GPT3/GPT4 based models from OpenAI to minimize hallucinations and improve text generation.
- **Text Summarization:** Engineered a deep learning-based model utilizing advanced NLP techniques, including LSTMs, to generate concise and coherent summaries of long articles with a ROUGE-2 score of approximately 0.73.
- **Keyword Extraction:** Created an automated system that extracts the most important keywords from large bodies of text using advanced NLP techniques and machine learning algorithms, resulting in a high F1 score of approximately 0.85.

CERTIFICATIONS

Google Cloud Certified Professional Machine Learning Engineer (Series ID: 4158)

May 2023 – May 2025

VOLUNTEER EXPERIENCE

National Service Scheme, President

Aug 2015 – Jul 2016

EDUCATION

- **Master of Science, Business Analytics (Data Science Track)**
- **Bachelor of Engineering, BITS Pilani, India**

Dec 2023

Aug 2017