

# Reuben Chatterjee

San Diego, CA | reuben.a.chatterjee@gmail.com | LinkedIn - Reuben Chatterjee | github.com/ReubenChatterjee

## Professional Summary

Data Scientist with a Master's in Data Science from UC San Diego, specializing in building scalable ML and analytics pipelines that drive measurable business outcomes. Expertise in data warehousing (Snowflake, Fivetran), supervised and unsupervised learning, customer segmentation modeling, churn prediction, and fraud detection. Skilled in feature engineering, statistical analysis, and time-series forecasting using Python, R, SQL, and cloud platforms (AWS, GCP, Azure). Proven track record of delivering \$2M+ business impact and 60% cost reductions in both academia and industry settings

## Experience

<b>Alcamo Marketing</b> <i>Data Scientist</i>	<i>Jul 2025 - Present</i>
<ul style="list-style-type: none"><li>Reduced operational costs by 60% (<b>\$20K annually</b>) by migrating data infrastructure from Adverity to Snowflake, building Snowflake tables and views with 100+ field mappings and transformations for 4 clients.</li><li>Enabled real-time reporting with 99% uptime for 4 clients by engineering 15 automated ETL pipelines (Fivetran to Snowflake to Looker Studio) processing 50GB+ daily data from Google Ads, LinkedIn Ads, Meta, Salesforce, and HubSpot</li><li>Engineered ML-based customer segmentation with clustering algorithms achieving 90% accuracy scores, analyzing behavioral patterns across 100K+ customer records to optimize engagement strategies and prioritize high-value opportunities</li></ul>	
<b>Ellis Lab, UC San Diego</b> <i>Graduate Research Assistant - Data Science</i>	<i>Sep 2024 - Jun 2025</i>
<ul style="list-style-type: none"><li>Developed automated data cleaning and validation pipeline in R using Regex and NLP to standardize 5,000+ free-text responses, reducing manual processing time by 90%</li><li>Created longitudinal panel dataset tracking 1,000+ students over time, and conducted demographic analysis across gender, major, and experience and visualized equity gaps in course outcomes using ggplot2 to inform course structure redesigns</li><li>Applied ANOVA, Tukey HSD tests and correlation analysis to examine gender based differences male vs female dominant groups across team contributions, leadership roles, and project outcomes</li></ul>	
<b>Datamatics Global Services</b> <i>Data Scientist Intern</i>	<i>Jun 2024 - Sep 2024</i>
<ul style="list-style-type: none"><li>Built Python and SQL based ETL pipelines to process 50GB+ of economic zone data for RAKEZ, enabling demand forecasting of industrial land lease activity</li><li>Trained a Random Forest model to optimize zone planning decisions, achieving 92% accuracy and validating results via A/B testing</li></ul>	
<b>Halicioglu Data Science Institute</b> <i>Data Analyst</i>	<i>Dec 2023 - Jun 2024</i>
<ul style="list-style-type: none"><li>Increased online engagement by 30% by building interactive Tableau dashboards to visualize social media metrics</li><li>Analyzed platform-specific metrics across Twitter, Instagram, and Facebook, generating recommendations that improved user retention and conversion rates by 15%</li></ul>	

## Projects

<b>Afterpay Customer Retention Prediction Model</b> <i>Python, XGBoost, LightGBM, Scikit-learn, Pandas, SHAP, Snowflake, SQL</i>	
<ul style="list-style-type: none"><li>Built gradient boosting churn prediction model achieving 89% accuracy and 0.92 ROC-AUC by engineering 50+ behavioral features including RFM metrics, payment patterns, and temporal trends from 350K+ transaction records.</li><li>Identified \$2.1M in at-risk revenue by scoring 12,500 high-risk customers through 4-tier risk segmentation system, enabling targeted retention campaigns that improved customer lifetime value by 18%.</li></ul>	
<b>Credit Card Anomaly Detection using Gradient Boosting</b> <i>Python, SciKit-learn, Hugging Face, CNN, XGBoost, LightGBM, Random Forest, Seaborn</i>	
<ul style="list-style-type: none"><li>Engineered 3,200+ behavioral features from 97,852 credit card transactions using domain-specific encodings and behavioral signals.</li><li>Tuned LightGBM via multi-model comparison (RF, XGBoost, CatBoost), achieving 92% accuracy and 0.59 OOT AUC</li><li>Reduced false positives by 10% via threshold tuning and SMOTE, contributing to \$2M+ projected annual savings</li></ul>	

## Education

<b>University of California, San Diego</b> <i>Master of Science in Data Science (CGPA:3.82)</i>	<i>Sep 2023 - Jun 2025</i>
<i>Courses: Statistical Models, Scalable Data Systems, Fraud &amp; Pricing Analytics, Deep Learning</i>	
<b>University of Mumbai</b> <i>Bachelor of Engineering in Computer Engineering (CGPA:3.93)</i>	<i>Aug 2019 - May 2023</i>
<i>Courses: Data Structures, Big Data Analytics, DBMS, AI &amp; Machine Learning, NLP</i>	

## Technical Skills

<b>Programming:</b> Python (Pandas, NumPy, Scikit-learn, PySpark), SQL, R, C++, JavaScript	
<b>Data Science &amp; Machine Learning:</b> Supervised & Unsupervised Learning (Classification, Regression, Clustering), Feature Engineering, Statistical Modeling, Time Series Forecasting, Model Deployment, A/B Testing, Customer Segmentation	
<b>Big Data &amp; Cloud Platforms:</b> Snowflake, Databricks, AWS, Azure, Hadoop, dbt, Ray, Docker, Kubernetes	
<b>Visualization &amp; Communication:</b> Tableau, Looker Studio, Power BI, Excel, ggplot2, Shiny, D3.js	