

MIN (MIA) SHI

469-403-7557 ◇ minmiashi@gmail.com ◇ [Personal Portfolio](#) ◇ <https://www.linkedin.com/in/min-mia-shi/>

EDUCATION

The University of Texas at Dallas	Dec. 2024 (<i>Expected</i>)
Ph.D. Candidate in Political Science – Quantitative Statistical Modeling Focused	GPA: 3.95/4.0
The University of Texas at Dallas	May 2024
Master of Science in Business Analytics (STEM) – Data Science & Data Engineer Track	GPA: 4.0/4.0

WORK EXPERIENCES

Data Engineer Intern	Jun. 2024 - Present
<i>The Sunwater Institute</i>	<i>North Bethesda, MD / Remote</i>
Developed scripts to collect data, created and managed data pipelines, and ensured data quality.	
<ul style="list-style-type: none">Implemented web scraping solutions to extract data from websites, storing over 1 million records in databases.Created ETL process for ingesting data using AWS S3 and Glue, boosting data processing efficiency by 40%.Automated speech-to-text and speaker identification using AWS Transcribe, achieving over 99% accuracy.	
Data Analyst & Research Assistant	May 2020 - May 2024
<i>The University of Texas at Dallas</i>	<i>Richardson, TX</i>
Took responsibility for data analysis for 10+ global health/policy projects using advanced statistical models.	
<ul style="list-style-type: none">Managed data collection in diverse methods including Qualtrics surveys and web scraping using R and Python.Developed 20+ robust statistical models (multi-variable and fixed-effect regression, difference-in-difference, time-series) combined ML models and NLP skills to support correlation and causal inference in research.Led a team of five junior assistants, ensuring collaboration and timely project completion and publication.	
Data Scientist Student Consultant	Aug. 2023 - Dec. 2023
<i>Working for Onyx CenterSource through The University of Texas at Dallas</i>	<i>Dallas, TX</i>
Led the creation of an AI-driven chatbot, enhancing customer engagement through advanced NLP techniques.	
<ul style="list-style-type: none">Employed NLP and MySQL for analyzing and querying an extensive database containing over 10 million entries.Achieved 25% improvement in response efficiency and provided 99% accurate predictions using XGBoost model.Contributed to a 15% rise in user engagement, increasing customer satisfaction and bolstering company's image.	

PROJECTS

Twitter Clone: High-throughput Social Media Backend — <i>Python, Django</i>	May 2024 - Present
Working on a six-month solo project developing a social media platform's backend using HBase, MySQL, and Redis with Django framework in Python.	
<ul style="list-style-type: none">Maximizing query efficiency by storing objects with HBase, MySQL & Amazon S3 based on query complexity.Addressing N+1 slow query issues by implementing Redis caching and denormalization.Integrating Celery and RabbitMQ to establish asynchronous workers with varying priority levels.Implementing a push model for distributing news feeds to followers efficiently.Optimizing memory and resource allocation using recursive small batches of asynchronous tasks.	
Kaggle Plant Pathology Competition: Leveraging Deep Learning CNNs	Nov. 2023 - Dec. 2023
Implemented deep learning models using Python and PyTorch to enhance disease identification accuracy in crops.	
<ul style="list-style-type: none">Utilized transfer learning on CNNs with 13042 images in 12 categories, enhancing disease identification accuracy.Conducted image transformation, including rotation, flipping, zooming, and noise injections to augment data.Fine-tuned ConvNext DL CNN models and achieve 86.8% accuracy, securing a Top 3 ranking in the competition.	
Big Data Risk Analysis and Data Visualization for a Trucking Company	Aug. 2022 - Dec. 2022
Engineered data visualization dashboards using Tableau, linked to Hadoop, for business risk analysis.	
<ul style="list-style-type: none">Processed and analyzed geospatial data with Hadoop, Hive, and Spark, reducing processing time by 40%.Developed Tableau visualizations linked to Hadoop and built interactive dashboards for business analysis.Conducted linear regression and multivariate analysis, contributing to predictive accuracy by 15%.	

SKILLS

Programming & Tools: Python, SQL, R, SAS, Stata, Tableau, Power BI, Alteryx
Database & Big Data: MySQL, PostgreSQL, AWS S3, AWS Glue, Hadoop, Sqoop, Hive, Impala, Pig, Spark
Certificates: Certificate in Applied Machine Learning, AWS Certified Cloud Practitioner, Google Analytics