

# MIN (MIA) SHI

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

## EDUCATION

---

### The University of Texas at Dallas

Ph.D. Candidate in Political Science – Quantitative Statistical Modeling Focused

Master of Science in Business Analytics (STEM) – Data Science & Data Engineer Track

GPA: 3.95/4.0

(*Expected*) Dec. 2024

May 2024

## WORK EXPERIENCES

---

### Data Scientist

*The Sunwater Institute*

Jun. 2024 - Present

*North Bethesda, MD / Remote*

Developed scripts to collect data, created and managed data pipelines, and ensured data quality.

- 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

*The University of Texas at Dallas*

May 2020 - May 2024

*Richardson, TX*

Took responsibility for data analysis for 10+ global health/policy projects using advanced statistical models.

- 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

*Working for Onyx CenterSource through The University of Texas at Dallas*

Aug. 2023 - Dec. 2023

*Dallas, TX*

Led the creation of an AI-driven chatbot, enhancing customer engagement through advanced NLP techniques.

- 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

---

### US Top 4 Airlines Financial Performance Analytics

Jan. 2024 - May 2024

- Analyzed over 10,000 records spanning 20 years to identify financial trends and shifts in the US airline industry.
- Pinpointed key strategic turning points affected by major events and changes in alliances and partnerships.
- Provided specific business model recommendations for enhancing the competitive stance of each top airline.

### 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.

- 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.

### Forecasting Stock Prices Through NLP Examination of Newspaper Articles

May 2023 - Dec. 2023

Developed automated web scraping tools and machine learning models in Python to predict stock market trends.

- Developed automated web scraping for 7,000+ WSJ articles, increasing data acquisition efficiency by 30%.
- Employed various vectorizers for WSJ article analysis, such as Tfidf Vectorizer, n-grams Count Vectorizer, etc.
- Utilized Naïve Bayes and Random Forest models, enhancing S&P 500 prediction accuracy by 12%.

### 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.

- 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, Git, SQL, R, SAS, Stata, Tableau, Power BI, Alteryx

**Database & Big Data:** MySQL, PostgreSQL, AWS S3, AWS Glue, Hadoop, Sqoop, Hive, Impala, Pig, Spark

**Data Analysis:** Machine Learning, Statistical Modeling, Data Visualization, Experimentation