

M.S. Applied Data Science Portfolio Paper

Syracuse University

By: Kenyang Lual

Introduction

The Master of Science in Applied Data Science (MS-ADS) program at Syracuse University has equipped me with the necessary skills to analyze, model, and interpret complex datasets across various domains.

This portfolio paper synthesizes my learning experiences and professional growth by presenting three significant projects undertaken during the program. These projects illustrate my proficiency in data collection, data cleaning, predictive modeling, data visualization, and ethical considerations in data science applications. The projects included in this portfolio are:

1. **Crime Data Analysis (IST 652 Final Project)** – Examining crime trends in Syracuse from 2019-2022 using Python.
2. **Energy Consumption Forecasting (IST 687 Final Project)** – Predicting electricity demand based on weather and building data using R.
3. **Agricultural Yield Prediction (IST 707 Python Final Project)** – Modeling corn production using meteorological data and machine learning in Python.

Each project reflects my ability to apply theoretical knowledge to practical challenges, utilizing data science methodologies to generate actionable insights that support decision-making processes.

Project 1: Crime Data Analysis (IST 652)

Objective

Crime analytics is essential for urban planning, public safety initiatives, and law enforcement resource allocation. This project aimed to analyze crime data from Syracuse, covering the years 2019-2022, to uncover key trends, identify peak crime periods, and assess arrest patterns. The findings can assist policymakers and law enforcement agencies in making informed decisions. Crime trends analysis helps in understanding the evolution of different types of crimes over time. By assessing crime rates across different neighborhoods, law enforcement agencies can allocate resources more effectively. Moreover, identifying correlations between crime occurrences and external factors such as socioeconomic conditions, weather patterns, and time of day provides deeper insights into potential crime prevention strategies. This project also sought to determine the effectiveness of law enforcement interventions by analyzing arrest rates in relation to crime occurrences. By studying which crimes had higher arrest rates and which ones did not, the analysis could provide actionable recommendations for improving policing strategies. Furthermore, the study considered ethical concerns in crime data analytics, ensuring that insights derived from the data would not perpetuate biases in law enforcement practices. By focusing on crime trends rather than individual profiling, the project aimed to contribute to fair and transparent crime prevention policies.

Technologies & Methods Used

- **Data Collection & Cleaning:** Crime dataset from Open Data Syracuse (CSV format)
- **Programming:** Python (Pandas, NumPy, Matplotlib, Seaborn)
- **Analysis Techniques:** Time-series analysis, correlation analysis, categorical variable grouping

- **Visualization:** Line charts, bar graphs, heatmaps
- **Statistical Evaluation:** Trend analysis and correlation analysis

Key Findings

- **Trend Analysis:** Aggravated assault and motor vehicle theft saw a significant increase, while larceny rates showed a steady decline.
- **Temporal Analysis:** Crimes were most prevalent around midnight and surged in August, suggesting seasonal patterns.
- **Arrest Correlation:** Violent crimes such as murder and aggravated assault had higher arrest rates, whereas larceny had lower arrest rates.
- **Policy Implications:** These insights suggest the need for increased night patrols and focused efforts during peak crime seasons.

Learning Outcomes Demonstrated

- Proficiency in data collection and cleaning for large-scale datasets
- Application of visualization techniques to effectively communicate trends
- Use of Python for statistical analysis and data manipulation
- Interpretation of crime data for actionable policy recommendations
- Consideration of ethical concerns in crime reporting and public data use

Project 2: Energy Consumption Forecasting (IST 687)

Objective

Energy consumption forecasting is critical for preventing blackouts, optimizing power distribution, and implementing sustainable energy solutions. This project focused on predicting hourly energy demand based on weather and building data, aiming to support energy providers in better demand planning and infrastructure development. By leveraging historical energy consumption data alongside meteorological factors such as temperature, humidity, and wind speed, this project sought to identify patterns in electricity demand fluctuations. Understanding these variations allows utility companies to implement proactive measures, such as load balancing, demand-side management, and infrastructure improvements. The study also evaluated the impact of extreme weather conditions on energy consumption, providing insights into how climate fluctuations influence residential and commercial electricity usage.

Additionally, the project explored potential strategies for reducing peak energy demand through data-driven recommendations. Ensuring ethical considerations in energy forecasting was also a key focus. The project acknowledged the risks of biased model predictions and sought to mitigate potential inaccuracies by implementing robust validation techniques and assessing data fairness. This ensures that the model's recommendations support equitable energy distribution without disproportionately affecting specific consumer groups.

Technologies & Methods Used

- **Data Sources:** House static data, hourly electricity usage, weather data (Amazon S3, CSV, Parquet)
- **Programming:** R (Tidyverse, ggplot2, Shiny App)
- **Modeling Techniques:** Linear regression, panel data regression, time-series forecasting

- **Visualization:** Interactive dashboards, histograms, scatter plots
- **Model Evaluation:** RMSE, adjusted R^2 , and cross-validation

Key Findings

- **Temperature Impact:** Higher temperatures led to a notable increase in electricity consumption, primarily due to cooling demands.
- **Peak Demand Periods:** Energy usage peaked at 6 PM, aligning with residential electricity usage patterns, and was lowest at 10 AM.
- **Model Performance:** The linear regression model achieved 91% accuracy (Adjusted $R^2 = 0.91$), demonstrating its effectiveness in predicting consumption trends.
- **Future Implications:** Utility companies could use this model to implement demand-side management strategies, shifting peak loads to off-peak hours.

Learning Outcomes Demonstrated

- Proficiency in data collection and cleaning for large-scale datasets
- Application of visualization techniques to effectively communicate trends
- Use of Python for statistical analysis and data manipulation
- Interpretation of crime data for actionable policy recommendations
- Consideration of ethical concerns in crime reporting and public data use

Project 3: Agricultural Yield Prediction (IST 707)

Objective

With climate change impacting agricultural production, predictive modeling is crucial for farmers and policymakers. This project leveraged meteorological and biodiesel market data to predict corn yields in Iowa, supporting decision-making in the agricultural sector. Agricultural production is highly dependent on external factors such as temperature, precipitation, soil conditions, and market demand. This project aimed to develop a data-driven model to forecast crop yields based on historical climate patterns and economic indicators. By integrating meteorological data from 90 weather stations and USDA biodiesel market data, the study identified key variables influencing corn yield variations. The project not only focused on predicting crop output but also assessed how fluctuations in climate conditions affect long-term agricultural sustainability. By leveraging machine learning techniques, such as Random Forest Regression, the model provided actionable insights into how farmers can optimize planting schedules and anticipate environmental risks. Additionally, ethical considerations in agricultural forecasting were addressed. The project acknowledged potential biases in historical data, ensuring that model predictions were interpretable and transparent. This approach prevents misleading recommendations that could disproportionately impact small-scale farmers or regions with climate instability.

Technologies & Methods Used

- **Data Collection & Processing:** Meteorological data (90 stations), USDA biodiesel market data, corn production data
- **Programming:** Python (Pandas, NumPy, Scikit-learn, Matplotlib, Seaborn)
- **Machine Learning Models:** Polynomial regression, Random Forest Regressor, cross-validation

- **Hyperparameter Tuning:** GridSearchCV for optimizing Random Forest parameters
- **Feature Engineering:** Standardization of variables, categorical encoding, time-based aggregations

Key Findings

- **Weather & Yield Correlation:** Rainfall and temperature were the most significant predictors of corn yield variations.
- **Market Influence:** Biodiesel price fluctuations had a minor effect on corn production trends.
- **Model Performance:** The Random Forest model achieved high predictive accuracy, indicating its robustness for agricultural forecasting.
- **Decision-Making Applications:** Farmers and policymakers could use these models to optimize planting schedules and anticipate market shifts.

Learning Outcomes Demonstrated

- Proficiency in data collection and cleaning for large-scale datasets
- Application of visualization techniques to effectively communicate trends
- Use of Python for statistical analysis and data manipulation
- Interpretation of crime data for actionable policy recommendations
- Consideration of ethical concerns in crime reporting and public data use

Conceptual Works and Their Influence

Throughout my coursework and projects, I have drawn upon several key conceptual works that shaped my understanding of statistical modeling, machine learning, and data ethics. *An Introduction to Statistical Learning* by James et al. [1] provided a structured approach to regression analysis and predictive modeling, influencing my work on energy consumption forecasting, where I applied linear regression to predict power demand. Similarly, *The Elements of Statistical Learning* by Friedman et al. [2] guided my use of machine learning algorithms, particularly Random Forest Regression, which I used to predict corn yields in my agricultural yield prediction project. For time-series analysis in both crime data trends and energy demand forecasting, I relied on the foundational work of Box and Jenkins [3], which helped me understand seasonality and autocorrelation in large datasets. Additionally, *The Visual Display of Quantitative Information* by Tufte [4] significantly influenced how I structured my data visualizations, ensuring clarity and effectiveness in conveying insights. Finally, ethical considerations were a critical component of all projects, and *Fairness in Machine Learning* by Barocas et al. [5] informed my approach to bias detection and mitigation, ensuring that models produced fair and unbiased predictions across different demographic groups. These works provided the theoretical foundation that allowed me to successfully design and implement data-driven solutions throughout the MS-ADS program.

Conclusion

Through these projects, I have demonstrated the ability to apply data science techniques to a range of real-world challenges, from crime analysis and energy forecasting to agricultural yield prediction. Each project required a combination of data collection, cleaning, modeling, and visualization, as well as a strong understanding of ethical considerations in data science. By integrating multiple datasets and employing diverse analytical techniques, I have developed solutions that provide valuable insights into public safety, energy consumption, and agricultural planning.

Key Skills Developed:

- **Data Processing & Management:** Collecting, cleaning, and integrating large datasets from multiple sources.
- **Predictive Modeling & Machine Learning:** Developing and fine-tuning statistical and machine learning models for forecasting and decision support.
- **Visualization & Communication:** Creating interactive reports, dashboards, and presentations that effectively communicate data-driven insights to technical and non-technical audiences.
- **Programming Expertise:** Utilizing Python and R to analyze complex datasets, automate workflows, and build scalable models.
- **Ethical Responsibility:** Ensuring fairness, transparency, and accountability in predictive modeling, while mitigating bias and maintaining data privacy.

Beyond technical skills, these projects have reinforced my ability to think critically, identify patterns in complex datasets, and develop solutions that can inform policy, optimize resource allocation, and support strategic planning. My work in crime analytics has provided insights into trends that can enhance public safety measures, while my energy forecasting project demonstrated how predictive modeling can help manage power grids more efficiently. Additionally, my research in agricultural yield prediction has highlighted the importance of data-driven decision-making in food production and sustainability.

Future Directions:

Moving forward, I aim to deepen my expertise in big data analytics, cloud computing, and deep learning, with a focus on building scalable, ethical, and impactful data science solutions. I am particularly

interested in applying advanced data science methodologies to domains such as smart cities, environmental monitoring, and business intelligence. By continuing to refine my skills and staying at the forefront of emerging technologies, I am committed to making meaningful contributions to the field of applied data science and helping organizations leverage data for improved decision-making and innovation.

References (IEEE Format):

1. G. James, D. Witten, T. Hastie, and R. Tibshirani, *An Introduction to Statistical Learning: With Applications in R*, 2nd ed. New York, NY, USA: Springer, 2021.
2. J. Friedman, T. Hastie, and R. Tibshirani, *The Elements of Statistical Learning: Data Mining, Inference, and Prediction*, 2nd ed. New York, NY, USA: Springer, 2009.
3. G. E. P. Box, G. M. Jenkins, G. C. Reinsel, and G. M. Ljung, *Time Series Analysis: Forecasting and Control*, 5th ed. Hoboken, NJ, USA: Wiley, 2015.
4. E. R. Tufte, *The Visual Display of Quantitative Information*, 2nd ed. Cheshire, CT, USA: Graphics Press, 2001.
5. S. Barocas, M. Hardt, and A. Narayanan, *Fairness in Machine Learning: Limitations and Opportunities*, 2019. [Online]. Available: <https://fairmlbook.org/>