

Hala Arar

705-977-6821 | hala.arar.02@gmail.com | linkedin.com/in/hala-arar | github.com/halaarar | <https://hala-arar.github.io/hala-arar-portfolio/>

SUMMARY

Data Scientist with a background in biomedical science and hands-on experience building NLP and machine learning systems. Developed end-to-end pipelines using PubMedBERT and XGBoost, deployed cloud-based solutions on AWS, and built interactive dashboards for clinical and public-sector use cases. Passionate about translating complex scientific and operational data into clear, actionable insights.

SKILLS

- **Languages:** Python, R, SQL, Bash
- **ML & NLP:** scikit-learn, XGBoost, LightGBM, TensorFlow, PyTorch, spaCy, Hugging Face, transformer models
- **Data Engineering:** ETL, data wrangling, REST APIs, SQL optimization
- **Tools:** Git, GitHub Actions, Jupyter, VS Code, Docker, AWS (S3, EC2, SageMaker)
- **Visualization:** Shiny, Plotly, Dash, geospatial dashboards

EXPERIENCE

Data Scientist – Capstone Project <i>University of Helsinki, Finland (Remote)</i>	Apr 2025 – Present
<ul style="list-style-type: none">Built a two-stage PubMedBERT pipeline for automated detection and classification of protein autoregulatory mechanisms across biomedical literature, achieving 96.0 percent accuracy in Stage 1 and 96.2 percent macro-F1 in Stage 2.Deployed the model across 252,880 unseen PubMed abstracts to generate the largest resource to date for predicted autoregulatory mechanisms and built an interactive Shiny app supporting search, filtering, and ontology-linked exploration.	
Data Scientist – Research Collaboration <i>University of British Columbia, Vancouver, BC (Remote)</i>	Jul 2025 – Oct 2025
<ul style="list-style-type: none">Developed XGBoost and adversarial LLM models to classify clinical gene variants using structured patient features, improving consistency between ACMG categories and reducing borderline classifications.Engineered domain-specific features aligned with cardiogenetics guidelines and collaborated with geneticists to design interpretable ML workflows for future clinical decision support.	

EDUCATION

Master of Data Science <i>University of British Columbia, Vancouver, BC</i>	Sep 2024 – Jun 2025
<ul style="list-style-type: none">Core Coursework: Advanced ML, Causal Inference, Data Viz, Cloud & Database Systems	
Bachelor of Science, Biomedical Sciences <i>Trent University, Peterborough, ON</i>	Sep 2020 – Apr 2024
<ul style="list-style-type: none">Minor in Mathematics Renewable Entrance Scholarship Dean's Honour Roll	

TECHNICAL PROJECTS

SOORENA: Autoregulatory Mechanism Detection (GitHub)	
<ul style="list-style-type: none">Built a two-stage PubMedBERT NLP system achieving 96.0 percent accuracy in detecting autoregulation and 96.2 percent macro-F1 across seven biological mechanism classesRan inference on 252,880 PubMed abstracts and developed an interactive Shiny dashboard with ontology-linked search and confidence scores	
Bank Marketing Predictions (GitHub)	
<ul style="list-style-type: none">Built a logistic regression model using 45,000+ customer records to predict term deposit subscriptions, improving targeting precision by 20 percentDelivered insights that reduced outreach waste and improved campaign ROI	
NYPD Arrest Tracker App (GitHub)	
<ul style="list-style-type: none">Developed a Plotly- and Pandas-based dashboard to visualize arrest patterns by borough, precinct, and offense categorySupported NYC government analysis with dynamic, filterable geospatial visualizations	

PUBLICATIONS

Arar, H., Aldahdooh, J., Nickchi, P., & Jafari, M. (2025). SOORENA: Self-loop containing or autoRegulatory Nodes in biological network Analysis. bioRxiv. <https://doi.org/10.1101/2025.11.03.685842> ([preprint](#)).