jxx151730@utdallas.edu   Cell: (469) 835-4885

# Jing Xia

Address: Dallas, Texas 75082

Data Scientist/Bioinformatics Scientist

## Summary

PhD candidate in bioinformatics, AI enthusiast who are interested in applying novel Machine Learning methods to solve biological questions. Strong statistical background and solid programming skills in Python, SQL and R.

## Education

| | | |
|---|---|---|
| **University of Texas at Dallas** *Richardson, Texas* | | Expected December 2022 |
| Ph.D., Bioinformatics | | GPA: 3.74/4.0 |
| **University of Texas at Dallas** *Richardson, Texas* | | Expected December 2022 |
| Master of Science in Data Science | | GPA: 3.84/4.0 |

Coursework: Machine Learning, Advanced Statistical Methods (I&II), Multivariate Analysis, Statistical Inference, Big Data, Graph Theory, Data Structure & Algorithm, Quantitative Biology

## SKILLS

**Programming:** Python, R, SAS, Bash, SQL, Apache Spark, PyTorch, TensorFlow, Scikit-learn, Advanced Excel

**Machine Learning Techniques:**

o   Linear Regression
o   Logistic Regression, Linear/Quadratic Discriminant Analysis, KNN, Random Forest, SVM
o   Neural Network, RNN(LSTM), CNN
o   K-Means, Hierarchical, Latent Dirichlet Allocation
o   Principal Component Analysis

## Career Experience

**University of Texas at Dallas**   *Richardson, Texas*                    January 2018 – December 2022

### *PhD Projects - Bioinformatics*

**RNA Variants Calling by Rigorous Statistical Framework**

o   Built a dictionary-based annotation system in Python to map billions of RNA sequencing sites to the genome.
o   Designed and implemented pipeline on clustering-based computing system to analyze 426 cancer cell RNA sequencing data by comparing three RNA identification packages.
o   Applied Binomial Test and Chi-squared Test on ~4 millions of RNA variants site to filter out systematic errors and obtained ~2 millions of RNA editing sites

**Causality Inference between RNA Editing and RNA Binding Protein**

o   Identified ~1 million common RNA editing sites from multiple datasets.
o   Conducted Fisher's Exact Test to find biological significant editing sites. Followed by several multiple test correction methods (Bonferroni and FDR).

**Protein Binding Prediction with RNA Sequence Data**

o   Preprocessed sequencing data and encoded sequence in one-hot encoding in PyTorch.
o   Trained a Convolution Neural Network to predict the binding of a specific RNA binding protein.

### *Master Projects – Data Science*

**Comprehensive linear Regression Analysis of 1000 Cardiovascular Patients' Data in R**

o   Performed linear assumption check by Shapiro-Wilk Test, Brown-Forsythe Test and Lack of Fit Test. Outliers were check by Studentized Residuals.
o   Carried out Box-Cox transformation to get normality of data.
o   Detected multicollinearity by Variance Inflation Factor and Condition Index.
o   Implemented three forms of linear models: Ordinary Least Square, Lasso Regression and Ridge Regression. Hyper parameters were optimized by Leave-One-Out Cross Validation.
o   Estimated test MSE for each model by k-fold cross-validation.

### In-depth Classification Analysis of 569 Breast Cancer Patients in Python
o   Preprocessed data by data cleaning, categorical feature transformation and standardization, etc.
o   Trained supervised machine learning models including Logistic Regression, KNN, LDA, QDA, SVM, Random Forest and Neural network. Regularization was applied to avoid overfitting with parameters determined by k-fold cross validation. Optimization methods such as Adam and mini-batch gradient descent were applied for Neural Network to speed up training. Various kernel functions were utilized for SVM.
o   Coded numerical searching algorithm (Newton-Raphson) to solve the likelihood function of logistic regression.
o   Evaluated performance (Accuracy or AUC) of each classifier via k-fold cross-validation and obtained the feature importance to find top factors for the prediction.

### Stock Price Prediction based on Deep Learning in Python
o   Exploratory data analysis of 15000 S&P 500 records followed by data normalization and splitting.
o   Built a simple RNN model from scratch to predict the open price. Truncated Backpropagation Through Time algorithm was coded to save computational power and regularization terms were implemented to avoid gradient explosion.  Fine-tuned the model by trying different combinations of hyperparameters.
o   Trained LSTM and GRU models via TensorFlow on GPU to validate the performance of hand built RNN. Various activation functions were tested.
o   GRU outperformed naïve RNN and LSTM on the testing data with respect to Mean Square Error.

### Natural Language Processing and Topic Modeling of 5000 Cell Phone Reviews in Python
o   Pretreated customer reviews by tokenization, stemming, removing stop words. Extracted features by Term Frequency – Inverse Document Frequency (TF-IDF).
o   Utilized three unsupervised learning models of K-means Clustering, Hierarchical Clustering and Latent Dirichlet Allocation.
o   Recognized latent topics and key words of each review for clustering.
o   Visualized training results by Principal Component Analysis (PCA).
o   Three algorithms provided meaningful clustering of customer reviews with latent semantic structures retrieved.

**UT Southwestern Medical Center**   *Dallas, Texas*                                    January 2018 – May 2019

*Research Assistant – Biochemistry & Bioinformatics*
**Participated in research project on biomedical data, applying molecular biology techniques and t-test for the biological difference.**
o   Publication: Ebrahim H.GHazvini Zadeh, ZhiJiang Huang, Jing Xia, Dalliang Li, Howard W. Davidson, and Wenhong Li, Functional study of ZnT8 in pancreatic cells by immunofluorescence and electron microscopy and quantified the expression level of three different pancreatic hormones in hundreds of pancreatic cells. (2020). https://doi.org/10.1016/j.celrep.2020.107904