

# KAI LIU

kai.liu@utexas.edu • 512.917.6781

GitHub: [github.com/KaiUT](https://github.com/KaiUT) • LinkedIn: [www.linkedin.com/in/kai-liu-utaustin](https://www.linkedin.com/in/kai-liu-utaustin)

## PROFILE

Ph.D. in Computational Biology (expected December 2018). Solid skills in machine learning, statistical inference. Comprehensive computing skills including Python, R, SQL, Hadoop, Spark, TensorFlow. Excellent problem solving skills in both independent and team environments. Skilled presenter of technical materials to both technical and non-technical audiences. Quick, thorough and effective learner.

## EDUCATION

<b>Ph.D. in Computational Biology</b> ◇ <i>The University of Texas at Austin, Austin, TX</i>	Expected December 2018
<b>M.S. in Microbiology</b> ◇ <i>Huazhong Agricultural University, Wuhan, China</i>	Grad. June 2013
<b>B.S. in Biotechnology</b> ◇ <i>Huazhong Agricultural University, Wuhan, China</i>	Grad. June 2011

## WORK EXPERIENCE

<b>Machine Learning Intern</b> ◇ <i>QuintilesIMS, Plymouth Meeting, PA</i>	June 2017 - August 2017
<b>Predicted Quality of Investigators in Future Clinical Trials   Python, Spark</b>	
· Predicted outliers of investigators per Key Risk Indicator using <u>distribution based approach</u> ;	
· Built multiple machine learning models ( <u>Lasso Regression, Neural Network, Random Forests</u> ) to predict the quality of investigators in a future study, which is one of the core projects in the investigator recommender system.	
<b>Graduate Research Assistant</b> ◇ <i>The University of Texas at Austin</i>	December 2014 - Present
<b>Developed Infectious Diseases Surveillance App   Python</b>	
· Developed a <u>regression model</u> and a <u>Multivariate Exponentially Weighted Moving Average (MEWMA)</u> model to detect emerging outbreaks with an accuracy of 0.9;	
· Combined above models with <u>stepwise variable selection algorithms</u> to select best data sources for infectious diseases surveillance (more than 400 data sources in total);	
· Built up data pipeline to automate the process of retrieving and cleaning data from CDC, BSVE, Google Trends, Wikipedia, Twitter etc; integrated the App into Cloud Ecosystem.	
<b>Assessed Real-time Zika Risk in the State of Texas   R</b>	
· Collaborated with other researchers in developing a <u>branching process model framework</u> that captures variation and uncertainty in Zika case reporting, importations, and transmission;	
· Applied the framework to assess county-level epidemic risk throughout Texas.	

## PROJECTS

<b>Developing a R Package for Big Data Analysis   R &amp; Rcpp</b>	
· Implementing following algorithms in the package: Stochastic gradient descent using line search and quasi-Newton methods to determine step size · The lasso · The proximal gradient method · Laplacian smoothing solved by sparse Cholesky/LU, the Gauss-Seidel method, the Jacobi iterative method, and conjugate gradient method · Graph fused lasso solved by Alternating Direction Method of Multipliers (ADMM) · Sparse matrix factorization.	
<b>Predicted the Direction of Exchange-Traded Fund (ETF) movement   Python</b>	
· Retrieved nine historical ETF sector datasets from Yahoo Finance;	
· Implemented <u>Logistic regression</u> , <u>Ridge &amp; Lasso regression</u> , and <u>Artificial Neural Network</u> to predict the direction of ETF movement;	
· Achieved an <u>accuracy of 55% ~ 60%</u> for predicting nine ETF sectors movement; and the trading strategy based on my prediction <u>outperforms</u> baseline strategies.	
<b>Predicted Yelp Rating Based on User Review Enhanced Collaborative Filtering   R</b>	
· Extracted user opinions from restaurants dataset from Yelp (~10GB) using <u>Stanford coreNLP tool</u> ;	
· Developed a <u>new Collaborative Filtering-based method</u> to improve the accuracy of user's rating prediction and solve the sparseness of dataset by <u>combining item's features and user opinions from all reviews</u> ;	
· Improved the prediction accuracy by <u>4.23%</u> compared to the traditional KNN method, and the <u>coverage is 100%</u> .	

## SKILLS

<b>Programming</b>	Fluency in Python(NumPy, SciPy, pandas, scikit-learn), R, Git · Familiar with MATLAB, Linux, LaTeX · Experience in SQL, Hadoop, Spark, TensorFlow, C++
<b>Machine Learning</b>	Deep Neural Network · Regression with regularization · Support Vector Machine · Random Forests · Hidden Markov Model · Clustering · Time series and dynamic models · Frequent Pattern Mining