

Yu (Zoey) Zhu

Linkedin: <https://www.linkedin.com/in/yu-zoey-zhu-8a330113a/>
Personal Web: <https://yuzoeyzhu.github.io/>

Email : yzhu201@ucsc.edu
Mobile : +1-(530)-304-6852

SUMMARY

Data Scientist and Machine Learning enthusiast with a strong academic background and more than 7 years of practical experience in analyzing structured and unstructured data. Passionate about uncovering valuable insights to drive business growth and product improvement through statistical analysis and machine learning techniques. Seeking a Data Scientist/Machine Learning Engineer Intern role to leverage my skills and contribute to a dynamic and innovative team.

EDUCATION

- **University of California, Santa Cruz** Santa Cruz, CA
PhD in Statistics; GPA: 4.0 *Sep 2020 - Jun 2025*
- **University of California, Davis** Davis, CA
MS in Statistics; GPA: 3.88 *Sep 2017 - Jun 2019*

SKILLS SUMMARY

- **Languages:** R, Python, SQL, Java, C++, JavaScript, Machine Learning Frameworks: Keras, PyTorch, Tensorflow
- **Tools:** GIT, Docker, Springboot, JIRA, Matlab
- **Data Manipulation:** Data Cleaning, Exploratory Data Analysis, Data Visualization
- **Hypothesis Testing:** A/B Testing, Test Design and Interpretation
- **Statistical Modeling:** Bayesian Parametric/Non-parametric Modeling, Time Series Analysis, Machine Learning (GLMs, Tree-based Methods, SVMs, Supervised/Semi-supervised/Unsupervised Learning, etc.), Model Selection, Feature Engineering, Causal Inference, Optimization

WORK EXPERIENCE

- **ThoughtWorks** Shanghai, China
Software Development Engineer (Data Track) *Aug 2019 - Aug 2020*
 - **Coca-Cola China Consumer Engagement Platform:** Designed and developed data processing pipelines using Python to extract, transform, and load data from various sources. Developed algorithms to analyze consumer engagement data and generate insights for business decision-making.
 - **IKEA PAX Cabinet AI Design System:** Developed machine learning algorithms to recommend personalized design options for customers based on their preferences and past behavior. Optimized the system's performance by analyzing user data and implementing feature engineering techniques.
 - **Starbucks APP:** Analyzed user engagement data using SQL and Python to identify patterns and trends in customer behavior. Developed machine learning models to predict user churn and optimize the customer reward program.
- **Bank of Suning** Nanjing, China
Data Scientist Intern - Risk Compliance Department, Big-data Team *Jun 2018 - Sep 2018*
 - Extracted, integrated, and cleaned 3 million rows of customer transaction data using Python, and conducted exploratory data analysis to identify patterns in missing data and feature distributions.
 - Developed a machine learning model to predict target customers for a financial loan product. Conducted data pre-processing by up-sampling with SMOTE and reducing data dimensionality using Principle Component Analysis (PCA). Compared the performance of several models, including K-Nearest Neighbors, SVM with RBF Kernels, and Ensemble Methods. Applied Accuracy, Precision, Recall and AUROC comparisons after hyper-parameter tuning with Grid Search and 5-fold cross-validation.
 - Achieved a **83% Accuracy** for telemarketing, representing a significant improvement from the previous 10% rate, with controlled variables.

ACADEMIC RESEARCH

- **Bayesian Latent Variable Models for Mortality Surveillance** Santa Cruz, CA
Research Assistant - Prof. Zehang (Richard) Li *Sep 2020 - Current*
 - Developed **Bayesian hierarchical models with latent structure** to infer the stratum-specific prevalence (SSP) of COVID-19 related death stratified with time and age with **distribution shift** across domains.
 - Demonstrated the formal framework to analyze **partially verified** Verbal Autopsy data under the **non-ignorable training data selection mechanism**.

- Introduced the novel use of **structured priors** to improve prevalence estimation for small sub-populations by more efficiently borrowing information from different sub-populations.
- Improved flexibility in modeling symptom distributions with respect to cause of death by incorporating advanced **tensor decomposition** techniques to capture the clusters of symptoms and the corresponding dependence.

Bayesian Non-parametric Bernstein Polynomial Model for ROC Curve

Research Assistant - Prof. Zehang (Richard) Li, Claudia Wehrhahn

Santa Cruz, CA

Apr 2022 - Current

- Proposed to model the Receiver Operating Characteristic (ROC) curve to validate a newly designed system for performing portable molecular diagnostic testing which denominated solid state PCR.
- Developed a flexible covariate dependent Bayesian non-parametric Bernstein polynomial model using stick-breaking process to accommodate to the bounded outcomes of the SS-PCR test.

Stochastic Nearest Neighbor Multiple Imputation of the TAST Database

Research Assistant - Prof. James Sharpnack

Davis, CA

Feb 2019 - Dec 2019

- Developed a new python package 'SDataFrame' to realize similar functions such as 'groupby' for data frame in Pandas, based on the imputation method of Stochastic Nearest Neighbors (SNN) with Euclidean distance.
- Simulated data with missing at random (MAR) missingness under Beta distribution with Guassian Mixture Model; Proposed SNN multiple imputation methods to compare with Multivariate Imputation by Chained Equations (MICE) as well as MissForest, and presented the advantages of SNN.

COURSE PROJECTS

Semi-Supervised Learning (SSL) on Causal and Anti-Causal Structure [\[Report\]](#)

Advisor: Prof. Zehang (Richard) Li

Santa Cruz, CA

Sep 2022 - Current

- Implemented the semi-generative model and conditional self-learning algorithm under the real-world Verbal Autopsy data to validate the assumption that SSL works for Anti-Causal other than Causal relationship.

Bayesian Non-parametric Approaches for Stochastic Order in ROC Analysis [\[Report\]](#)

Advisor: Prof. Anathasis Kottas

Santa Cruz, CA

March 2022 - June 2022

- Applied the Bayesian non-parametric approaches Dirichlet process mixtures (DPM) and Mixtures of finite Polya tree (MPT) with stochastic order constraint to model the ROC curve with a meaningful value of AUROC that strictly larger than 0.5.

Image Recognition with Bayesian CNN for Simpsons Characters [\[Report\]](#)

Advisor: Prof. Juhee Lee

Santa Cruz, CA

March 2022 - June 2022

- Proposed and compared the Non-Bayesian Convolutional Neural Network (regularCNN) with Bayesian Convolutional Neural Network (BayesCNN) with Variational Inference based on the predictive performance for the image recognition task under Simpsons data set.
- Measure the uncertainty estimation in BayesCNN and interpreted the uncertainty based on the 95% credible intervals of posterior predicted class assignment probabilities for some of the test images.

Robust PCA and Extreme Classification

Advisor: Prof. Cho-Jui Hsieh

Davis, CA

Nov 2017 - Dec 2017

- Applied the ADMM algorithm to solve a robust PCA problem under the non-convex condition and tested in MNIST
- Conducted the Conjugate Gradient Descent algorithm to solve a multi-label classification problem with an extremely large number of labels in MATLAB and R

CONFERENCE

WNAR

A Bayesian Hierarchical Model for Mortality Surveillance using Partially Verified Verbal Autopsy Data [\[Talk\]](#)

Anchorage, AK

Jun 2023

Objective Bayes

A Bayesian Hierarchical Model for COVID-19 Related Cause-of-death Assignment Using Verbal Autopsies [\[Poster\]](#)

Santa Cruz, CA

Sep 2022

TEACHING ASSISTANT

Fall 2021, Spring 2022 **[STAT 07] Statistical Methods for the Biological, Environmental and Health Science**

Fall 2020, WInter 2021, Winter 2022 **[STAT 05] Statistics**

HONORS AND REWARDS

- 2023 UCSC Statistics Summer Research Fellowship
- 2023 UCSC Graduate Dean's Travel Fellowship
- 2023 WNAR Student Paper Competition Travel Fellowship
- 2022 UCSC Statistics Summer Research Fellowship
- 2022 UW Biostatistics Summer Institutes Scholarship
- 2016 National Scholarship for Students with Excellent Academic Performance