

# HANCHAO ZHANG

Chair of CSSA at NYU Grossman School of Medicine  
Website: hanchaozhang.xyz  
Address: 180 Madison Ave, 5-31C, New York, NY, 10016

Mobile: (+1) - 646 - 206 - 6662  
Email: hanchao.zhang@nyu.edu

## EDUCATION

- COLUMBIA UNIVERSITY, SCHOOL OF ART AND SCIENCE** New York, New York  
Visiting Ph.D. Student - Statistics Department through IUDC;  
Jun. 2023 - Present  
*Research* : Stochastic Process
- NEW YORK UNIVERSITY, GROSSMAN SCHOOL OF MEDICINE & SCHOOL OF ART AND SCIENCE** New York, New York  
Doctor of Philosophy - Biostatistics;  
Aug. 2019 - Present  
*Thesis Advisor*: Professor Thaddeus Tarpey (Chair, and Program Director)  
*Research Interest*: Manifold Learning, Self-Consistency, Functional Data, Supervised Learning and Clustering Methods in High-Dimensional Data  
*Awards*: Special MacCracken Awards 2023 for Research Outstanding
- JOHNS HOPKINS UNIVERSITY, SCHOOL OF PUBLIC HEALTH** Baltimore, Maryland  
Visiting Student - Biostatistics, Applied Mathematics and Statistics;  
Sept. 2018 - Jun. 2019
- CORNELL UNIVERSITY, SCHOOL OF INFORMATION SCIENCES & SCHOOL OF MEDICINE** New York, New York  
Master of Science - Biostatistics and Data Science; GPA: 4.1/4.3  
Aug. 2017 - Sept. 2018
- CAPTIAL UNIVERSITY OF ECONOMICS AND BUSINESS, SCHOOL OF FINANCE** Beijing, China  
Bachelor of Art - International Finance; GPA: 3.8/4.0 (WES Certified)  
Aug. 2013 - Jun. 2017

## EXPERIENCE

- GOOGLE** Mountain View, California  
Ph.D. Research Scientist Intern  
Jun. 2022 - Sep. 2022
  - Developed statistical model that quantifies cross-questions inter-rater reliability and detects unreliable raters with cross-questions information based on linear mixed effect model with repetitive measurements
  - Developed a semi-supervised statistical model to estimate latent variables representing unreliable raters. The model improves the prediction accuracy of rating by 10.25% compared to the linear mixed effect model
- HEDGEHOG LAB INC. ( [HTTPS://HLAB.APP/](https://hlab.app/) )** Remote  
Co-founder, Chief Scientist  
Jun. 2021 - Jun. 2022
  - Co-founder and Chief Scientist of Hedgehog Lab, a web-based computing language that runs machine learning algorithm
  - Contributed machine learning libraries based on javascript and hedgehogscripts.
- JOHNS HOPKINS UNIVERSITY, BLOOMBERG SCHOOL OF PUBLIC HEALTH** Baltimore, Maryland  
Research Assistant  
Oct. 2018 - Jun. 2019
  - Performed data cleaning & manipulation over text data; built automatic data pipeline with python to extract features based on NLP (nltk, gensim)
  - Customized machine learning model to make predictions based on multiple data sources (NHANES, Universities & Google Search dataset)
  - Conducted statistical analysis to interpret & optimize model performance based on multiple metrics (ROC, AUC, R-square, confusion matrix)
- TECHNICAL CONSULTING & RESEARCH, INC.** New York, New York  
Data Analyst Intern  
Oct. 2018 - Jun. 2019
  - Performed web scraping with Python (Selenium, BeautifulSoup, Multiprocess) to get 5k + medical data from National Library
  - Utilized NLP techniques to extract features (location, age, history, condition) from medical script data with Python (nltk, genism, regr)
  - Conducted statistical analysis to interpret & optimize model performance based on multiple metrics (ROC, AUC, R-square, confusion matrix)

## PUBLICATIONS & TALKS

- METHODOLOGY PAPER: K-TENSORS: CLUSTERING POSITIVE SEMI-DEFINITE MATRICES:** Hanchao Zhang, Thaddeus Tarpey., arXiv
- METHODOLOGY PAPER: PRINCIPAL AND SELF-CONSISTENCY FOR POSITIVE SEMI-DEFINITE MATRICES ON RIEMANNIAN MANIFOLD:** Hanchao Zhang, Thaddeus Tarpey., Working Paper
- COLLABERATIVE PAPER: LOW MUSCLE MASS IS ASSOCIATED WITH A HIGHER RISK OF ALL-CAUSE AND CARDIOVASCULAR DISEASE-SPECIFIC MORTALITY IN CANCER SURVIVORS:** Dongyu Zhang, Hanchao Zhang, etc., Nutrition
- COLLABERATIVE PAPER: VALIDATION OF EHR MEDICATION FILL DATA OBTAINED THROUGH ELECTRONIC LINKAGE WITH PHARMACIES:** Saul Blecker, Samrachana Adhikari, Hanchao Zhang, etc., Journal of Managed Care and Specialty Pharmacy

- **COLLABERATIVE PAPER: QUANTITATIVE EVALUATION OF REJUVENATION TREATMENT OF NASOLABIAL FOLD WRINKLES BY REGRESSION MODEL AND 3D PHOTOGRAPHY:** Rou-Yu Fang, Hanchao Zhang, etc., Journal of Cosmetic Dermatologyh
- **INVITED TALK: CLUSTERING POSITIVE SEMI-DEFINITE MATRICES: A METRIC LEARNING APPROACH TO DISEASES SUBTYPING:** International Biometric Society, Invited Paper Session (ENAR 2023)
- **INVITED TALK: OPTIMAL TRANSFORMATIONS OF HIGH-DIMENSIONAL FUNCTIONAL DATA FOR CLUSTERING METHODS:** Joint Statistical Meeting 2022, Invited Paper Session (JSM 2022)
- **INVITED TALK: FUNCTIONAL DATA CLUSTERING AND REGRESSION METHODS:** Columbia University Functional Data Working Group Invited Talk (FDAWG 2022)
- **INVITED TALK: OPTIMAL LINEAR TRANSFORMATIONS OF FUNCTIONAL DATA FOR CLUSTERING METHODS:** Joint Statistical Meeting 2021, Invited Paper Session (JSM 2021)

## PROJECTS

---

- **PRINCIPAL AND SELF-CONSISTENT POSITIVE SEMI-DEFINITE MANIFOLD:**
  - Developed and defined a novel self-consistent submanifold framework tailored for positive semi-definite manifold data
  - Investigated and established the statistical properties governing this submanifold, showcasing its principal and self-consistency properties
  - Defined a novel distance metric for positive semi-definite matrices, contributing to the advancement of distance-based analysis on manifolds
- **K-TENSORS: A SELF-CONSISTENT ALGORITHM FOR CLUSTERING POSITIVE SEMI-DEFINITE MATRICES:**
  - Developed a self-consistent clustering algorithm tailored specifically for positive semi-definite matrix-valued data
  - Extended the clustering algorithm to seamlessly accommodate manifold data, enhancing its applicability and robustness on manifold learning
  - Improved the clustering algorithm's capabilities in detecting the true data generation process
- **SELF-CONSISTENT CONVEXITY-BASED CLUSTERING ALGORITHM FOR BREGMAN DISTANCE:**
  - Developed a self-consistent support hyperplane approach for convexity-based clustering as a weighted maximum volume problem
  - Conducted analysis of the intricate relationship between convexity-based clustering and clustering algorithms based on the Bregman distance
- **FAIRELASTICNET: A FRAMEWORK FOR FAIRNESS VARIABLE SELECTION:**
  - Developed the FairElasticnet model, a novel and impactful approach for fairness variable selection in prediction models
  - Applied augmentation methods and ADMM optimizer to simplify and optimize  $\ell_1$  and  $\ell_2$  problems within the FairElasticnet framework
- **AN OUTLIERS DETECTION ALGORITHM FOR FUNCTIONAL DATA:**
  - Developed an outliers detection algorithm for Electroencephalography (EEG) data using a random consecutive window method
  - Applied the developed outliers detection algorithm on the real EEG dataset and improved the AUC from 88% to 96%
  - Built a functional regression model on the data preprocessed by the outliers detection algorithm, and improved prediction accuracy by 12%
- **ASSOCIATION BETWEEN MATERNAL PRENATAL STRESS AND HUMAN FETAL BRAIN DEVELOPMENT:**
  - Acquired the fMRI and questionnaire data, preprocessed and winsorized the survey data and fMRI data
  - Applied EDA for previewing the data and obtained the association between cortisol value and multiple stress scores
  - Utilized unsupervised clustering methods (PCA and T-SNE) to cluster the patients and identify subgroups of the patients for further analysis
- **DEEP LEARNING U-NET MODEL IN SEGMENTATION OF BRAIN MR IMAGES:**
  - Utilized the MRI data released for the MICCAI 2012 Grand Challenge on Multi-Atlas Segmentation
  - Established a baseline network for two basic segmentation tasks: brain/non-brain and grey matter/white matter/cerebrospinal fluid
  - Provided the network full 2D axial slices of the MR volumes for training and testing
  - Interpreted the changes that are being made to the network design by exploring intermediate feature maps created in the alternate networks, looking for differences in the organization of features
- **RISK PREDICTION AND EVALUATION OF THE EFFECT OF MYOCARDIAL INFARCTION ON STROKE:**
  - Performed data & feature engineering over 1.7 M healthcare data(Strokes) including extraction, manipulation, feature generation & selection
  - Applied Survival Analysis and machine learning methods in computer clusters, including Cox Regression, Accelerated Failure model Survival Random Forest. Selected the Cox Regression without interaction as the best prediction model by Cross- Validation using R (survival, randomForestSRC, and pec) (Dr. Diaz's Lab)
- **ROBUST REGRESSION FOR OUTLIERS IN LABORATORY DATA:**
  - Applied three robust regression methods, M-estimation, S-estimation, and MM-estimation, to evaluate the association between urinary PGE-M level and urinary PGD-M level in obese and lean mice, with and without celecoxib by R (MASS, ggplot2); Concluded that M-estimation is slightly better than the other two robust regressions
  - Improved certainty of variance analysis by leveraging each data point with weight and new loss function
- **TRANSPORTATION AND ELECTRONIC HEALTH RECORD (EHR) DATABASE CONSTRUCTION:**
  - Established database schema and populated tables based on electronic health records and New York City transportation data from the U.S. Department of Transportation on the MySQL server
  - Applied Logistic Regression to search association of transportation preference, concluding that number of patients having Benign Essential Hypertension was associated with the patient's location in New York City, and Hyperlipidemia and Atrial Fibrillation were associated with the patient's transportation preference using R (stat, ggplot2, randomForest)

## SKILLS SUMMARY

---

- **PROGRAMMING:** Python, R, Linux, SQL, SAS (Advanced Certified)
- **LANGUAGES:** English, Mandarin