

# Kexin (Summer) Shang

4044097577 | [ks4254@drexel.edu](mailto:ks4254@drexel.edu) | [linkedin.com/in/kexin-shang5301](https://www.linkedin.com/in/kexin-shang5301)

## EDUCATION

---

<b>Drexel University, PA</b> Doctor of Philosophy in Information Science (Focus on LLM in Healthcare)	Sep 2023 - Present <b>GPA: 4.00/4.00</b>
<b>Washington University in St. Louis, MO</b> Master of Science in Biostatistics and Data Science	Sep 2021 - Dec 2022 <b>GPA: 3.94/4.00</b>
<b>Georgia State University, GA</b> Bachelor of Science in Mathematics (Statistics) Bachelor of Science in Biology (Double Major, Co-diploma)	Jun 2019 - May 2023 <b>GPA: 3.85/4.30</b>
<b>Southwest Jiaotong University, Chengdu, China</b> Bachelor of Engineering in Bioengineering (Co-diploma)	Sep 2017 - Jun 2019 <b>GPA: 3.63/4.00</b>
Main courses: <i>Recommender System, Nature Language Processing (Pytorch), Network Analysis, Applied Machine Learning in Data Science (Sklearn), Biostatistics (SAS), Analysis, Optimization, Survival Analysis, Bioinformatics (Linux), etc.</i>	

## RESEARCH EXPERIENCE

---

<b>Healthcare Informatics Research Lab, CCI, Drexel</b> Topic: Multi-LLM Collaboration on Medical QA 1 <sup>st</sup> author manuscript ready for International Conference on Healthcare Informatics (ICHI) 2025 submission <ul style="list-style-type: none"><li>Developed an agent-based LLM collaboration framework that improves both consistency and accuracy across all three LLM participants on medical QA task</li><li>Measured the “confidence” of a LLM and investigate how it affects LLM’s teamworking strategy</li><li>Serve local LLMs through Text Generation Inference (TGI) method to significantly reduce inference time cost</li><li>Elicit LLMs’ generation via a combined prompting of Zero-shot Chain-of-Thought and Self-Consistency</li></ul>	Research Assistant	Sep 2023 - Present
<b>Center for Healthy Weight and Wellness, Psychiatry, WUSTL</b> Topic: Harnessing Mobile Technology to Reduce Mental Health Disorders in College Population Poster presented on International Conference on Eating Disorders (ICED) 2023 <ul style="list-style-type: none"><li>Constructed composite variables from over 200 features via PCA regression, which determines nearly 40% of the variation of response rate to follow-up surveys</li><li>Designed a factorial design on 4 treatment components and identified moderator variables with each component using logistic regression models and simple slope analysis</li><li>Conducted a cross-sectional survey the prevalence of 11 types of clinical and subclinical eating disorders in rural areas, suburban areas, and urban areas in U.S. applying pairwise T-tests with Holm’s corrections</li></ul>	Intern	May - Dec 2022
<b>Department of Developmental Biology, WUSTL</b> Topic: Role of Transposable Element in Transcript-level Expression Regulation <ul style="list-style-type: none"><li>Developed a Shell-based pipeline to obtain TE-derived transcripts’ expression contribution from GTEx database</li><li>Located age-sensitive TE-derived transcripts in skin tissue by plotting time-series Z-scores across age intervals</li><li>Removed unwanted variation using residuals (RUVr with k=4) from RNAseq data and plot a 3D PCA which successfully showed clear separations between sun-exposed skin genes and sun-unexposed skin genes</li></ul>	Research Assistant	Sep 2021- Jun 2022

## IN-CLASS PROJECTS

---

<b>DSCI 691 Nature Language Processing</b>	Drexel	2024 Summer
--------------------------------------------	--------	-------------

Topic: Text-based Emotion Recognition across BERT Family and LSTM

Source of data: a dataset of English Twitter messages with six basic emotions: anger, fear, joy, love, sadness, and surprise. <https://huggingface.co/datasets/dair-ai/emotion>

- Benchmarked four BERT family models (BERT, RoBERTa, DistilBERT, and XLNet) on ‘emotion’ dataset with bidirectional LSTM, the previous generation state-of-the-art model as the baseline
- Analyzed the potential reasons of the significant difference between our result and a [similar research](#)

GitHub page: <https://github.com/summer5301/Text-based-Emotion-Recognition-across-BERT-Family-and-LSTM>

## **DSCI 511 Data Acquisition and Preprocessing**

Drexel

2023 Fall

Topic: Analysis of the Canadian Wildfires Effect Posed on the Air Quality in US Cities.

Source of data: Scraped Wikipedia for Top 20 most populous US Cities and the “Open-Meteo” API for weather data and air quality data.

Individual contribution:

- Used Pandas Python package to restructured time-series weather data of each city scraped from Wikipedia and API into 1-year span by date and month.
- Represented continuous variables such as pm 2.5 index by mean and categorical variables such as “air quality level” by major vote and store cleaned data in Json files.

GitHub page: [https://github.com/summer5301/4\\_smokwatchers\\_project/tree/main](https://github.com/summer5301/4_smokwatchers_project/tree/main)

## **MSB 660 01 Biomedical Data Mining**

WUSTL

2022 Spring

- Leveraged the Medical Expenditure Panel Survey (MEPS) database to predict medical cost across 3376 patients by fitting models of multiple linear regression, bagged random forest regression, and logistic regression w/t lasso penalty
- Adopted LDA and Naïve Bayes classifier to classify patients with high medical cost, achieving 96.9% and 94.1% specificity respectively
- Used Inverse normal transformation (INT) to normalize highly skewed data (change skewness from 4.9 to 0.074)

## **CONFERENCE**

2022 ASA Women in Statistics and Data Science Conference, St. Louis, MO

Audience

International Conference on Eating Disorders (ICED) 2023, Washington, DC

Poster Presenter

## **HONOURS & AWARDS**

Valedictorian of the Recognition Ceremony, WUSTL

2022

Merit Scholarship (\$11,886), WUSTL

2021

Wiley M. Suttles Math Award (\$750), GSU

2023

In-state Scholarship; Presidential List; Member of the Honors College, GSU

2020 - 2021

Second-class Scholarship (¥3000); National Scholarship Nominated, SWJTU

2018 - 2019

## **EXTRACURRICULAR ACTIVITY**

Publicity Department of the Chinese Student Union, GSU

Minister

- 2020 Atlanta Chinese Students and Scholars Spring Festival Gala
- Social media account management and operation

## **SKILLS**

Analytics: Machine Learning Models, Deep Learning Architectures (implementation on CNN, RNN, transformers) on various data types (text, image, and tabular), OpenAI API, LangChain, vLLM and Text Generation Inference (TGI) inference

Programming: Python (Pandas, Numpy, Tensorflow, Pytorch), Shell, R, SAS, MySQL, Latex

Language: English (fluent); Chinese (native)