

# YUXUAN FU

(646) 462-0849 | [yf2592@columbia.edu](mailto:yf2592@columbia.edu) | [Linkedin](#)

## EDUCATION

<b>Columbia University</b>	New York, NY
<b>M.A. in Statistics</b>	Sep 2021 – Dec 2022
<ul style="list-style-type: none"><li>Courses: Statistical Computing &amp; Data Science, database design, deep learning, Algorithm Analysis, Cloud Computing &amp; Big Data, NLP</li><li>Leadership: Columbia University Chinese Students and Scholars Association – Media and Public Relations member</li></ul>	
<b>East China Normal University</b>	Shanghai, China
<b>B.S. in Statistics</b>	Sep 2017 – Jul 2021
<ul style="list-style-type: none"><li>University Outstanding Scholar-2<sup>nd</sup> prize; Supervisor of department's Student Office, Practice &amp; Innovation Center</li><li>Courses: Sampling Survey, Machine Learning, Bayesian, Stochastic, Experimental Design, Nonparametric, C++, Business Analytics</li></ul>	
<b>University of California, Berkeley - Exchange Student</b>	Berkeley, CA
Courses: Statistical Computing (graduate level), Time Series, Linear Analytics and Regression.	Sep 2019 - Jan 2020

## PROFESSIONAL EXPERIENCE

<b>Children and Screens: Institute of Digital Media and Child Development</b>	New York, NY
<i>Research Assistant Volunteering</i>	Jan 2022 - Present
<ul style="list-style-type: none"><li>Review latest academic literatures for scholar database and weekly digest presentation, minute webinar and write YouTube summaries.</li><li>Analyze and visualize attendance/outreach/feedback webinar data with Tableau and Demo to offer outreach improvement advice.</li></ul>	
<b>Roche Diagnostic</b>	Shanghai, China
<i>Data Science Intern - Product Development Global Clinical Execution Team (PDG)</i>	Jun 2020 – Mar 2021
<ul style="list-style-type: none"><li>Devised a data visualization platform with clinical management team to promote over 10 corporate-scale R&amp;D tasks.</li><li>Transmitted data with R and python, including reports crawling from Almac and Rave database, R preprocessing, Chrome GAuth tokens' authentication and Google Data Studio interactive template design.</li><li>Specified 7 - 10 paged dashboard templates respectively for six medical experiments by daily communicating with CRAs.</li><li>Probed advanced NLP tasks with CDM to mine the similarity of drug description text, using LDA and deep learning algorithm.</li></ul>	
<b>Accenture</b>	Shanghai, China
<i>Project Assistant – IT Consulting Department</i>	Jan 2020 – Mar 2020
<ul style="list-style-type: none"><li>Provided customer with IT support, cloud service and cloud office for an energy development project.</li><li>Completed over 40 documents' content organization, formatting, partial proofreading, translation and data sorting.</li><li>Integrated a 2GB team scheme based on customer inquiries and fulfilled the customer delivery task with manager and partners.</li></ul>	

## RESEARCH PROJECTS

<b>Multi-model human face expression video recognition (Python)</b>	Oct 2021- Dec 2021
<ul style="list-style-type: none"><li>Operated human facial expression recognition system for video data, applying retina net for face detection, face alignments.</li><li>Literature review for VGG-n in combination with LSTM architecture to reduce hyperparameter numbers and raise accuracy.</li><li>Gained fusion score with C3D, VGG16_LSTM and achieved prediction results with C3D network on Google Cloud Platform.</li></ul>	
<b>Anomaly Detection of Civil Aviation Flights based on QAR data (R)</b>	Jan 2021 – Jun 2021
<ul style="list-style-type: none"><li>Preprocessed 49 flights QAR dataset within 7944 time-varying variables through parallel computation, considering time constraints, strong interference, operational meanings and rough-set's encoding patterns.</li><li>Visualized 644 variables' groupings to respectively 35 and 45 groups with multi-time series method: autocorrelation-function-based Fuzzy C-means Clustering combining SBD distance and K-means Clustering combining DTW distance.</li><li>Implemented KPCA for representative feature extraction and detected 7-10 potential abnormal flights with single variable's functional boxplots and feature-vector based DBSCAN clustering.</li></ul>	
<b>Infants' Sleep Quality, Neurobehavioral Development, Language &amp; Social Interaction Research (R)</b>	Jan 2020 – Apr 2021
<ul style="list-style-type: none"><li>Collaborated with hospital statistics of more than 200,000 standardized questionnaires across nation, measured the effects of gestational age at birth on infants' sleep quality, neurobehavioral development, and language and social interaction.</li><li>Experimental design and linear regression, logistic regression, nonparametric and multilevel statistical analysis modeling with continuous and discrete dependent varies, independent varies and covariates. Paper in submission for publication.</li></ul>	
<b>Prediction of Bitcoin Price using Machine Learning and Sentiment Analysis (Python)</b>	Sep 2020 – Nov 2020
<ul style="list-style-type: none"><li>Extracted sentiment from 4.9M Twitter &amp; Reddit comments with VADER, Text-Blob and LM dict. as prediction's proxy.</li><li>Integrated time by weighted-moving averaging sentiment polarity &amp; subjectivity, bitcoin close price, and platforms' popularity index.</li><li>Extracted 88% explanation of variability feature with LightGBM to simplified DL network's Input.</li><li>Forecasted data trend with LSTM, LSTM_CNN, GRU and identified GRU as better with 62% lower test RMSE than other two.</li></ul>	
<b>User-Generated Content (UGC) Information Accuracy Evaluation and Influencing Factor Analysis (Python)</b>	Mar 2019 – Apr 2020
<ul style="list-style-type: none"><li>Directed the independent subject on UGC accuracy evaluation based on Baidu, Insight-China Pharma database.</li><li>Procured over 1000 opensource medical prescriptions' text quickly by utilizing <i>urllib</i>, <i>BeautifulSoup</i> for Baidu webpage and <i>ajax</i>, <i>selenium</i>, <i>requests</i> for database, overcoming asynchronous login and JavaScript obfuscator issues.</li><li>Implemented LDA combining perplexity, KL divergence model in self-written functions and parallel preprocessed 6 GB unstructured text data, including text segmentation, word tagging and keyword extraction.</li><li>Visualized text comparison of html output with <i>difflib</i> through python and terminal commands.</li></ul>	

## TECHNIQUE SKILLS

- Languages:** R, Python, SQL, PySpark, SAS, C++, Ubuntu, Bash Shell, Git, JavaScript, Latex
- Big Data Toolkits:** MapReduce, Snowflake, GCP, AWS, Azure | **ML Toolkits:** Scikit-Learn, TensorFlow, MXNet
- Data Visualization:** Tableau, R Shiny, Google Data Studio, Power BI | **Methodology:** A/B testing, experiment design, data mining