

Work Experience

- 2022.01 - Now **Machine Learning Engineer (offer holder)** @ Garena | Sea Group

Just join the company and work remotely in a week. I find the daily routine of this role is a little bit inconsistent with its job description. It builds the backend part using Golang for Booyah! recommendation system, and it should be referred as "backend engineer" :)

- 2020.09 - 2021.12 **Data Scientist** @ Barclays

Tech stack: Spark / PySpark (on Elastic Data Platform), Amazon Deep Java Library (DJL), Tensorflow / Keras, BitBucket, Pandas, Jupyter, Pretrained Transformers / Transformers / Likelihood Ratio

- Company address matching and entity matching without internal GPU and labeled data available. Solved using an active learning method. Start from constructing some small datasets only with external data and training an XGBoost tree, then label samples in the boundary and fine-tune BERT models in an iterative way. Finish the inference on 6 million internal pair-wise samples with this model on a CPU cluster, using a DJL based pipeline built from scratch on my own. It achieved a very satisfying result of 94% F1 score on a noisy testing dataset. The model does inference offline on our Spark cluster in a distributed way. For 6 million pair-wise samples, its running time is under 1 hour (on a cluster with 80 CPUs).
- Predict the aggregated user's transaction activity (volume and value), using historical mean and Informer model, a variant of Transformer for timeseries modeling; and then construct counterfactual, to provide an evaluation to how much finance loss of this bank suffers from a system downtime, and to find out the critical time period for the system reliability.
- Maintain the Spark cluster for the team, and build up pipelines for distributed inference by combining DJL / PySpark UDF with models. Cooperated with one of my colleagues, we created a team-wise package to start a Spark session within 4 lines of codes, which significantly reduces the overhead of using Spark for colleagues who are not with a distributed computing background.
- Participated in a fully immersed 6-weeks cloud DevOps training, which involves the deployment of a working pipeline including GitHub, DockerHub, Jenkins, and AWS EKS (Kubernetes) cluster, using Terraform and Ansible.

- 2019.08 - 2020.07 **Java Backend Developer** @ Barclays

Tech stack: Jenkins, Jira, Confluence, BitBucket, OpenShift (Kubernetes), Docker, GridGain, Maven, Gradle, Wiremock, Mockito, Spring Boot, SonarQube, Karate, AppDynamics

- End-to-end function development, testing (unit, functional, performance), deployment (CD)
- Add cache layer to the existing APIs to reduce the latency for repetitive data access
- Migrate legacy codes to internal Spring Boot templates, with refactors to enhance code readability and performance
- Build up handy internal tools (e.g. git hooks) and scripts (python / bash) from scratch to automate software development processes

- 2020.08 - 2020.09 **Natural Language Processing Engineer** @ Kwai Inc (Kuaishou)

- Model Migrations for an newly built internal model management system
- Retaining a transformer-based user language model using newly incoming data (~ millions, fetched by Hive SQL) to mitigate data shift issue

The work happens while my time-off in Barclays - I was in a career choice at that time point where Kwai gave me this campus offer, and Barclays offers me a role of change at the same time. It was a great experience. My leave was for a better work-life balance.

Education

- 2018 - 2019 **University College London** , MSc Web Science and Big Data Analytics, Distinction

Probability Graphical Models; Deep Learning; Complex Network; Affective Computing; Statistical NLP; Information Retrieval

- 2016 - 2018 **University of Liverpool** , BSc Internet Computing, First class

- 2014 - 2016 **Xi'an Jiaotong-Liverpool University** , BSc Information and Computing Science

2+2 pathway routine (first 2 years in Suzhou, China and final 2 years in Liverpool, UK), dual degree.

Selected Projects

- 2021.09 - 2021.10 **Wechat chat history analysis**

It was a gift for one of my important friends for a friendship anniversary. I collected all of our chat history (in Chinese). Take an analysis for the following aspects: the emotion appeared in our single chat sentences (80% accuracy for fine-tuning a Chinese version Roberta model for a 6-classes dataset including happy, natural, angry, fearful, anxious, exciting), word cloud generation, Wechat emoji counting, and hourly chat statistics. The final delivery is a mobile html5 page constructed using a library wechat-h5-boilerplate. This was an end-to-end project from data collection (WeChat does not provide any public way of exporting chat history).

- 2021.08 - 2021.09 **Implementation of Conflict-Free Replicated Data Type (CRDT) - set and graph**

A project from one of my previous take-home interview exercises, finished in a week. It is a self-contained, fully-functional and properly-tested implementation of Last Write Wins (LWW) Graph and Set in Python, which is one type of CRDT. CRDTs can be replicated across systems, they can be updated independently and concurrently without coordination between the replicas, and it is always mathematically possible to resolve inconsistencies that might result. Github: <https://github.com/billweasley/Last-write-win-CRDT-graph>

- 2019.06 - 2019.09 **Project Internship (Master Degree Thesis) @ Astroscreen**

Social media posting language source identification (tweets and gabs) project. Finished a crawler for collecting language (posts) data from Gab.com, pre-processed data using Regular Expression, built models for classifying the source of these data by fine-tuning BERT and XLNet, visualized results using t-SNE, did "leave-one-hashtag-out" cross-validation, and evaluated models using some common metrics (Accuracy, F1 score, Confusion Matrix, Matthews Correlation Coefficient).

- 2019.02 - 2019.03 **Integrating BERT and Embeddings into CommonsenseQA Challenge**

We fine-tuned Google BERT to CommonsenseQA challenge 1.0 (with 3 options of each question) and then integrated Conceptnet Numberbatch and ELMo embeddings attempting to improve the model performance. The challenge involves a set of MCQ questions requiring human commonsense knowledge. We achieved 68.79% of accuracy on validation set using BERT + ELMo (soly BERT : 67.47%; BERT + Numberbatch: 67.68%).

Technical Article

- "Accelerating Deep Learning on the JVM with Apache Spark and NVIDIA GPUs"**

Author: Haoxuan Wang, Qin Lan [AWS], Carol McDonald [Nvidia]; Link: https://www.infoq.com/articles/deep-learning-apache-spark-nvidia-gpu/?itm_source=articles_about_ai-ml-data-eng&itm_medium=link&itm_campaign=ai-ml-data-eng