

# SUIYUAN ZHANG

✉ a\_dayuanzi@163.com · ☎ (+86) 136-5115-6109

## EDUCATION BACKGROUND

<b>Institute of Computing Technology, Chinese Wuhan University of Technology</b>	Key Laboratory of Network Data Computer Science	<i>Master Scholar(2/110)</i>	2015 – Present 2011 – 2015
--	--	----------------------------------	-------------------------------

## AWARDS

<b>Intelligent Traffic Prediction Challenge</b>	1 / 1716	2017.07 – 2016.08
<b>Kaggle Quora Pair Question Match Contest</b>	31 / 3304	2016.05 – 2016.07
<b>Ali Music Trends Forecasting Competition</b>	16 / 5475	2016.05 – 2016.07
<b>Churn Rate Prediction</b>	1 / 200+	2016.07 – 2016.09
<b>Hotel Production Forecast for the Next 30 Days</b>	1 / 200+	2016.07 – 2016.09
<b>National Collegiate Mathematical Modeling Contest</b>	First Prize	top 1% 2013.09

## INTERNSHIP

<b>Didi Chuxing</b>	Map Department	<i>Algorithm Engineer Intern</i>	2017.06 – 2017.08
<ul style="list-style-type: none"><li>• Be familiar with the framework of line recommendation, developed by pyspark in experimental data</li><li>• The link topological structure of road is established, and the traffic condition is excavated by GCN (graph convolutional neural network model)</li></ul>			
<b>Ping An Group of China</b>	Commercial Real Estate	<i>Indoor Positioning Engineer</i>	2014.09 – 2015.01

## PROJECT EXPERIENCE

<b>Distributed Information Retrieval System</b>	<b>Core Member</b>	2016.09 – Present
---	--------------------	-------------------

Written by C++, the underlying system of multiple systems

- Invert index data structure, supporting for real-time indexing, a new document can be retrieved for no more than 1 minutes
- Support node hot plug, load balancing, construction index exception handling, support document segmentation
- Update improving, delete efficiency, use Thrift framework for communication for transmission higher efficiency
- Use mmseg segmentation algorithm, based on cosine similarity as a basis to sort, use DBSCAN to cluster results

<b>Ctrip Cloud - Division Tournament Big Data Contest</b>	<b>Captain</b>	2016.07 – 2016.09
---	----------------	-------------------

To predict the daily demand of 500 hotels in February of 2016 according to the historical hotel night volume in one year.

- Abstract as two sub-problems, select the hotel attributes, do the of time window statistics features
- The double-split point based random forest algorithm is better than the normal random forest and XGBoost

<b>Trec 2017    OpenSearch</b>	<b>Captain</b>	2017.04 – Present
--------------------------------	----------------	-------------------

A paper sorting system according to the paper search site CiteSeerX feedback information. Obtained the score of 0.133 (logloss), 31/3304 in Quora Match Contest of the same features and model in the just-concluded Kaggle platform

- Build TF-IDF matrix. Do the TSVD, NMF decomposition. Train LDA topic model.
- Train WordToVector, DocToVector and LDAToVector to convert sentences to vector.
- Use wordnet to extract the feature of synonyms, evaluating the similarity by editing distance, cosine distance, etc.
- Total of more than 2000 multi-dimensional features, we use LightGBM、FM、FFM to train.
- Use LSTM and CDSSM. To combine the models, use stacking. DBGD to meet the real-time demand

## IT COMPETENCE

- Programming language: Python、C++、Java

## OTHERS

- blog: <http://blog.csdn.net/chedan541300521> (XGBoost source reading notes)
- GitHub: <https://github.com/dayuanzi> (The double-split point based random forest algorithm)