Detecting M Giants in Space Using XGBoost

Dr. Zesheng Chen Department of Computer Science Purdue University Fort Wayne



An Efficient Spectral Selection of M Giants Using XGBoost

Zhenping Yi,¹ Zesheng Chen,² Jingchang Pan,¹ Lili Yue,¹ Yuxiang Lu,¹ Jia Li,¹ and A-Li Luo³

¹School of Mechanical, Electrical and Information Engineering, Shandong University, Weihai, 264209, Shandong, China ²Department of Computer Science, Purdue University Fort Wayne, Fort Wayne, IN 46805, USA ³Key Laboratory of Optical Astronomy, National Astronomical Observatories, Chinese Academy of Sciences, Beijing 100012, China



PURDUE FORT WAYNE THE Nobel Prize in Physics 2019









- Red giants with spectral type M
- Lower surface temperature (≤ 4000K)
- Extremely bright with typical luminosities of 10³ L_O
- M giants provide a way for researchers to explore the substructures of the halo of the Milky Way



AGE MAP The ages of tens of thousands of red giant stars are charted atop a map of the Milky Way. The oldest stars are in red, near the galactic center. The youngest stars are blue.







Data
XGBoost
Results

PURDUE VNIVERSITY FORT WAYNE Data

LAMOST DR4 data

- LAMOST is a new type of wide-field telescopes with a large aperture and a large field of view
- Currently, LAMOST DR4 has released 7.68 million spectra
- We used 6,311 M giant spectra and 5,883 M dwarf spectra, with labels
- We randomly selected about 70% as the training data



- Extreme Gradient Boosting
- A scalable machine learning system for tree boosting
- An open source package
- Widely recognized in many machine learning and data mining challenges (e.g., Kaggle)
- Use slides from "Introduction to Boosted Trees" by Tianqi Chen

PURDUE FORT WAYNE Regression Tree (CART)

- Decision rules same as in decision tree
- Contains one score in each leaf value

Input: age, gender, occupation, ...

Like the computer game X





Prediction score is the sum of scores predicted by each of the tree.



- Why do we want to contain two components in the objective?
- Optimizing training loss encourages predictive models
 - Fitting well in training data at least get you close to training data which is hopefully close to the underlying distribution
- Optimizing regularization encourages simple models
 - Simpler models tends to have smaller variance in future predictions, making prediction stable





PURDUE Shallow Learning vs. FORT WAYNE Deep Learning

- Shallow learning algorithms learn the parameters of a model directly from the features of training samples and build a structurally understandable model
- We focus on shallow learning to identify most important features to separate M giants from M dwarfs

PURDUE Performance Comparison of FORT WAYNE Four Machine Learning Methods

Algorithm	Accuracy	Precision	Recall
XGBoost	99.79	96.87	98.93
SVM	99.53	92.08	98.94
Random forests	99.29	90.05	96.23
ELM	98.75	80.71	97.85

FORT WAYNE Important Features

- We found that 287 features among 3,951 pixels of input data are used in XGBoost
- The more times a feature is used in XGBoost tree, the more important it is

 $IPS = N_{times}/SUM_{times}$

FORT WAYNE IMPORTANT Features









PURDUE FORT WAYNE CONClusions

- XGBoost is used to discern M giants from M dwarfs for spectroscopic surveys
- The important feature bands for distinguishing between M giants and M dwarfs are accurately identified by the XGBoost method
- We think that our XGBoost classifier will perform effectively for other spectral surveys as well if the corresponding features wavelength bands are covered

PURDUE FORT WAYNE Thanks For Your Attention

