

Finanziato dall'Unione europea NextGenerationEU







NeuroStarMap Stefano Tortora, ALTEN

Spoke 3 III Technical Workshop, Perugia 26-29 Maggio, 2025

ICSC Italian Research Center on High-Performance Computing, Big Data and Quantum Computing









Project overview

- Main goal: development of a neural network-based algorithm to estimate astronomical distances, trained and validated on the Gaia DR3 stellar catalog.
- Impact: enabling astronomers to obtain more reliable estimation of stellar and cosmological distances.

Scientific coordinations

- Deborah Busonero
- Mario Gai

Team composition

- Alessio Petrone (PM)
- Luca Maina (TL)
- Stefano Tortora (TL)
- Simone Zimotti (DEV)
- Carmine Fruncillo (DEV)











Project steps

- **Dataset preparation:** collection and processing of data from GAIA satellite on three datasets concerning cepheids, RRLyrae and eclipse binaries.
- Neural network architecture study: analysis of potential neural network architectures.
- **Network development and implementation:** design and coding of the neural network models selected from the previous study.
- Training of networks: training of the networks on the datasets, monitoring performance metrics and
 optimizing them by tuning parameters and adapting the architecture to the specific characteristics of the
 datasets.
- Validation of the networks: preliminary testing of the neural networks on the datasets to ensure that they meet the established accuracy and efficiency criteria.
- **Final pipeline development:** creating decision mechanisms based on simple averages, weighted averages or other techniques to be defined.
- **Final system validation:** testing of the completed system for overall effectiveness, robustness and ease of use, ensuring that the final product meets design requirements.











- Classification and regression algorithms
- Data preprocessing (scaling, encoding, normalization etc.)
- Feature selection
- Model evaluation (accuracy, confusion matrix etc.)

Ideal for problems involving structured data



- Full support for building and training deep neural networks
- Management of layers, activation functions, and optimizers
- Training and validation on GPU
- Model saving, exporting, and reuse

Ideal for working with deep models









Data Preparation – Merge and clean datasets

The Cepheid and RRLyrae datasets individually contain a relatively low number of records. It was decided to merge them to maximize the amount of data available for a single training and validation session.

 The neural models can achieve greater generalization by learning from a wider variety of examples.

Fields common to both Cepheids and RR Lyrae stars were selected to ensure data consistency during the training and validation phases.

Name	Description				
phot_g_mean_mag	G-band photometric magnitude				
bp_rp	BP/RP Spectrum Photometry				
peak_to_peak_g	Peak-to-peak amplitude of the G- band				
metallicity	Metallicity of the star				
phi21_g	Phase difference in 1 and 2 harmonic				
r21_g	ratio of amplitude of the 1 and 2 harmonics				
zp_mag_g	Zero point of the final G-band light curve				
parallax_error	Parallax measurement error				
type_best_classification	Type classification between: "DCEP", "T2CEP", "ACEP".				
parallax	Parallax measurement				
Llood data act fields					

Used dataset fields











ICSC Italian Research Center on High-Performance Computing, Big Data and Quantum Computing









Deep Neural Network (DNN) architecture

- Very deep network (7 dense layers)
- One neuron for each feature used from the dataset
 - in the input layer
- Dropout = 5×10^{-2}
- Weight initialization = "He Normal"
- Batch Normalization = present
- Activation function = "LeakyReLU" (alpha = 10⁻⁴)
- L2 Regularization = absent



ICSC Italian Research Center on High-Performance Computing, Big Data and Quantum Computing









Multilayer Perceptron (MLP) architecture

- Shallow network (4 dense layers)
- One neuron for each feature used from the dataset
 - in the input layer
- Dropout = absent
- Weight initialization = "glorot uniform"
- Batch Normalization = present
- Activation function = "LeakyReLU" (alpha = 10⁻¹)
- L2 Regularization = present

11 neuroni 256 neuroni 128 neuroni 64 neuroni Neurone

ICSC Italian Research Center on High-Performance Computing, Big Data and Quantum Computing

Missione 4 • Istruzione e Ricerca

Hidden Layer

Input layer

Output Layer



Both the DNN and MLP achieve very similar and progressively better results as the parallax values increase.

For values very close to zero the predictions become less accurate. This can be attributed to the great distance of the celestial bodies involved and to the numerical difficulties encountered by the neural networks.

ICSC Italian Research Center on High-Performance Computing, Big Data and Quantum Computing

values

σ

dicte









Data Preparation – Clean datasets

Parallax error related fields were removed from the dataset: they may not be available in datasets that will be used for parallax prediction.





100

0.1 values

0.001

100µ

10µ

Predicted 0.01







Real vs Predicted values (Log Scale)

DNN



Mean Squared Error (MSE): 3.176 x 10⁻²

10

Mean Absolute Error (MAE): 9.104 x 10⁻²

Coefficient of determination (R²): 2.254 x 10⁻¹

100

Mean Squared Error (MSE): 3.075 x 10⁻² Mean Absolute Error (MAE): 9.233 x 10⁻² Coefficient of determination (R²): 2.5 x 10⁻¹

Removing the parallax error related fields caused a significant drop in performance: the previous results were heavily influenced by these fields.

ICSC Italian Research Center on High-Performance Computing, Big Data and Quantum Computing

Missione 4 • Istruzione e Ricerca



MLP

 Predicted vs Re ---- Regression Line - - y = x









Data Preparation – Separate datasets

Analyze the Cepheid and RRLyrae datasets separately and extract more fields with better correlations.

Because of the sparsity of the datasets more attention was given to preprocessing the data, removing outliers and rows with few information fields.

Name	Description	Name	Description				
zp_mag_g	Zero point (mag) of the final model of the G-band light curve for Cepheids and RR Lyrae stars.		Goodness-of-fit statistic of the astrometric solution for the source in the along-scan direction. This is the 'gaussianized chi-square', which for good fits should approximately follow a normal distribution with zero mean value and unit standard deviation. Values exceeding, say, +3 thus indicate a bad fit to				
zp_mag_bp	Zero point (mag) of the final model of the GBP band light curve for Cepheids and RR Lyrae stars.	astrometric_gof_al					
peak_to_peak_g	beak_g This parameter is filled with the peak-to-peak amplitude value		the data.				
	of the G band light curve. The peak-to-peak amplitude is calculated as the (maximum) - (minimum) of the modelled folded light curve in the G band.	astrometric_chi2_al	Astrometric goodness-of-fit (χ 2) in the AL direction. χ 2 values were computed for the 'good' AL observations of the source, without taking into account the astrometric excess				
peak_to_peak_bp	Peak-to-peak amplitude of the BP band light curve (float, Magnitude[mag])		the attitude excess noise (if any) of each observation.				
fund_freq1	First frequency of the non-linear Fourier modelling. It applies to all three G, GBP, and GRP bands and the radial velocity curve.	astrometric_excess _noise	This is the excess noise of the source. It measures the disagreement, expressed as an angle, between the observations of a source and the best-fitting standard astrometric model (using five astrometric parameters).				
•••							

RRLyrae used dataset fields









Data Preparation – variables correlation (in progress)

Two variables highly correlated contain redundant information. Keeping both can:

- Increase noise;
- Make the model more complex than necessary;
- Slow down training;
- Confuse the optimization process (weights "share" the effect);
- Worsen generalization.
- Variables highly correlated with the target: potentially useful;
- Variables weakly correlated with everything: may not carry useful information;
- Variables too strongly correlated with each other: useless, select only one.

parallax -	1.00	-0.02	-0.02	0.04	-0.02	0.02	0.02	-0.02		- 1.0	
ruwe -	-0.02	1.00	-0.03	0.08	0.01	0.07	-0.03	-0.03		- 0.8	
phot_g_mean_mag -	-0.02	-0.03	1.00	0.27	-0.18	0.31	-0.06	1.00		- 0.6	
bp_rp -	0.04	0.08	0.27	1.00	-0.21	0.21	0.09	0.27		- 0.4	orrelation
peak_to_peak_g -	-0.02	0.01	-0.18	-0.21	1.00	0.16	-0.45	-0.19		- 0.2	Index of co
r21_g -	0.02	0.07	0.31	0.21	0.16	1.00	-0.23	0.31		- 0.0	
phi21_g -	0.02	-0.03	-0.06	0.09	-0.45	-0.23	1.00	-0.06		- –0.2	2
zp_mag_g -	-0.02	-0.03	1.00	0.27	-0.19	0.31	-0.06	1.00		- - 0.4	ı
	parallax -	- Inwe	phot_g_mean_mag -	bp_rp -	peak_to_peak_g -	r21_9 -	phi21_g -	zp_mag_g -			







Real vs Predicted values (Log Scale



Actual results (in progress)

Cepheid

Real vs Predicted values (Log Scale)



Mean Squared Error (MSE): 3.16505 x 10 ⁻³ Mean Absolute Error (MAE): 3.86555 x 10 ⁻² Coefficient of determination (R²): 8.46768 x 10 ⁻¹



RRLyrae

Mean Squared Error (MSE): 1.77831 x 10⁻¹ Mean Absolute Error (MAE): 3.68798 x 10⁻¹ Coefficient of determination (R²): 5.22587 x 10⁻¹

After dataset separation and variables correlation study, results are slightly better than the previous ones.

ICSC Italian Research Center on High-Performance Computing, Big Data and Quantum Computing

Missione 4 • Istruzione e Ricerca

Predicted vs Real
 Regression Line
 Y = Y









Future development phases

Once the phases in progress will be completed:



ICSC Italian Research Center on High-Performance Computing, Big Data and Quantum Computing



Finanziato dall'Unione europea NextGenerationEU







Thank you

ICSC Italian Research Center on High-Performance Computing, Big Data and Quantum Computing