



# **Anomaly Detection in Stellar Data: Unveiling the Unusual in the Cosmos**

**Exploring Celestial Mysteries through Data Analysis**



# Executive Summary

**Project Aim:** An in-depth analysis of stellar data to identify anomalies through unsupervised learning models, aiming to uncover stars with unusual properties that might indicate rare or unique celestial phenomena.

**Objectives:**

- Implement Isolation Forest and Local Outlier Factor (LOF) models to detect anomalies within stellar characteristics.
- Evaluate the effectiveness of these models in distinguishing potential celestial objects of interest.

**Key Findings:**

- Identified anomalies exhibit significantly different properties, such as higher temperatures and luminosities, suggesting the potential discovery of rare stars.
- Isolation Forest model demonstrated a higher silhouette score, indicating a clearer distinction of anomalies compared to LOF.



## Data Overview:

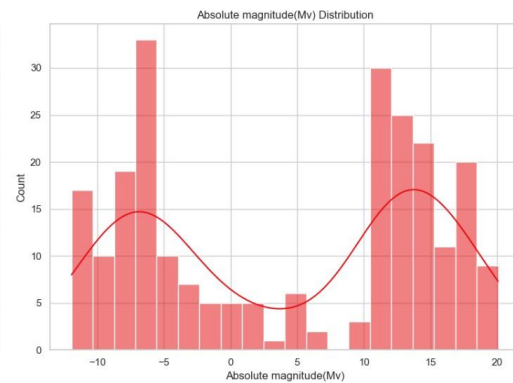
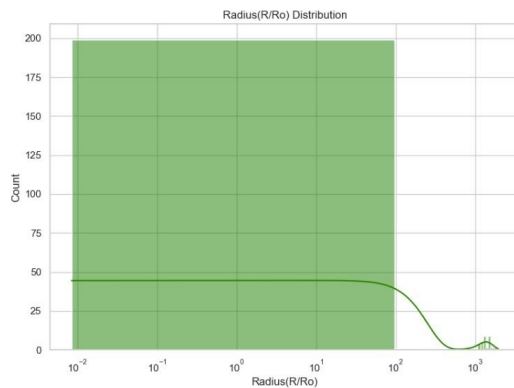
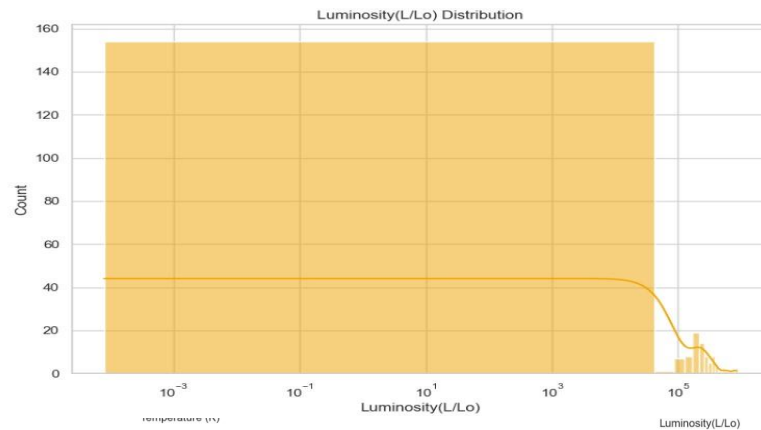
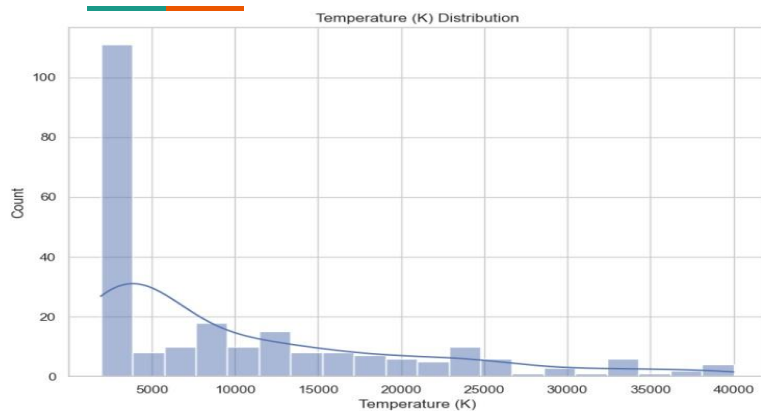
**Dataset:** This analysis focused on a comprehensive dataset of stellar attributes, including Temperature (K), Luminosity (L/L<sub>o</sub>), Radius (R/R<sub>o</sub>), Absolute magnitude (M<sub>v</sub>), Star color, and Spectral Class.

**Data Diversity:** This dataset encompasses a wide range of star types across the Hertzsprung-Russell diagram, reflecting the diversity in temperature, size, and luminosity inherent to stars in the universe.

**Preprocessing Steps:** I standardized the numerical features and variable for consistent analysis and applied encoding techniques to categorical data, preparing our dataset for the anomaly detection process.

**Objective:** The prepared dataset served as the foundation for applying unsupervised learning models to detect stellar anomalies, aiming to identify unusual or rare celestial phenomena.

# Exploratory Data Analysis





# Methodology

**Model Selection:** Given the nature of cosmology data, i thought it would useful to approach from different methodology. I selected the Isolation Forest and Local Outlier Factor (LOF) models for their efficiency in detecting anomalies within high-dimensional data, such as our cosmology dataset. Anomalies in astronomical data often have the more impactful insights.

**Anomaly Detection Process:** The models were applied to the preprocessed dataset to identify stars with properties significantly deviating from the norm, indicating potential anomalies.

**Model Evaluation:** I next evaluated the models based on visualization techniques and silhouette scores to assess the effectiveness of the anomaly detection and the distinctiveness of identified anomalies.

**Insights & Interpretation:** The analysis aimed to interpret the detected anomalies in the context of astrophysics, seeking to uncover potentially rare or previously unidentified celestial phenomena.



## Model Development and Tuning

**Model Training:** We trained both the Isolation Forest and LOF models on the cosmology dataset, focusing on accurately identifying anomalies based on stellar properties."

**Parameter Optimization:** Tuning the models involved adjusting parameters such as the number of trees in Isolation Forest and the number of neighbors in LOF to enhance model sensitivity and specificity.

**Evaluation Metrics:** For this, I utilized silhouette scores to evaluate model performance, alongside visual assessments of anomaly detection results, ensuring a balanced approach to identifying true anomalies.

**Challenges Overcome:** I addressed overfitting and ensuring the models' generalizability across unseen data were key focuses during the tuning process.

## Model Development and Tuning

```
lof_anomalies_avg, overall_avg, common_anomalies
```

```
(Temperature (K)          12891.394737
 Luminosity(L/Lo)         173731.381828
 Radius(R/Ro)             401.987558
 Absolute magnitude(Mv)    5.753526
 dtype: float64,
 Temperature (K)          10497.462500
 Luminosity(L/Lo)         107188.361635
 Radius(R/Ro)             237.157781
 Absolute magnitude(Mv)    4.382396
 dtype: float64,
 14)
```

```
silhouette_iso, silhouette_lof
```

```
(0.4564553845578292, 0.14671093515624603)
```

```
[16]: LOF_Anomaly
      1    202
     -1     38
      Name: count, dtype: int64
```

Both models effectively identified anomalies, with the Isolation Forest showing a higher silhouette score, indicating a more distinct separation of anomalies. The LOF seeing 38 anomalies and the isolated forest seeing more at 53 with 14 being identified by both.



## Key Insights and Recommendations

Anomalies identified include stars with temperatures exceeding the dataset's 95th percentile, pointing towards extremely hot stars possibly indicative of young, massive stars or rare stellar remnants. For instance, a subset of anomalies demonstrated temperatures over 20,000 K, significantly higher than the average, suggesting unusual stellar processes or evolutionary stages.

The comparison between Isolation Forest and LOF models revealed a distinct set of anomalies, with the Isolation Forest identifying a higher proportion of high-luminosity outliers. This discrepancy highlights the diverse nature of stellar anomalies and underscores the value of multiple detection methods in capturing the full spectrum of celestial diversity.





## Key Insights and Recommendations

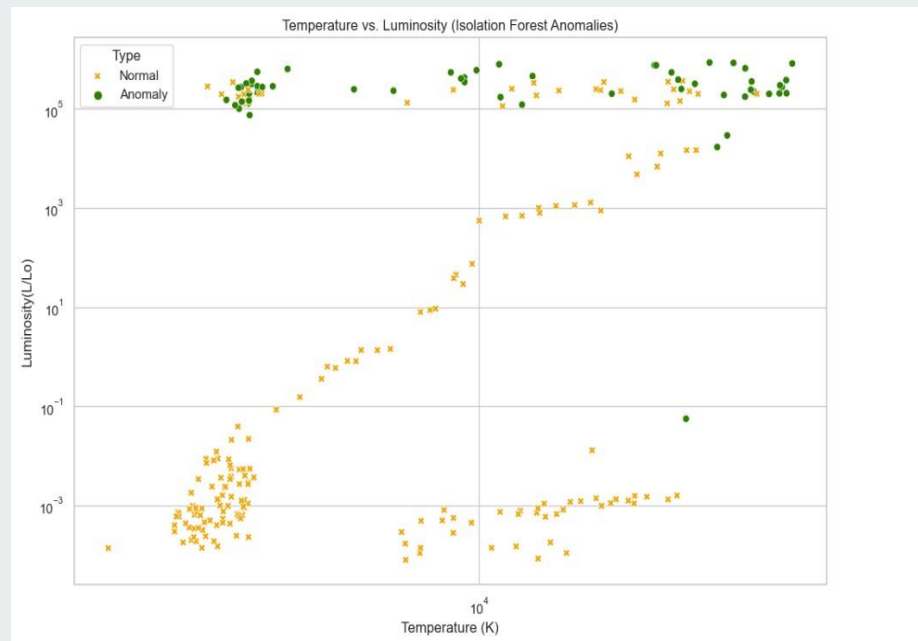
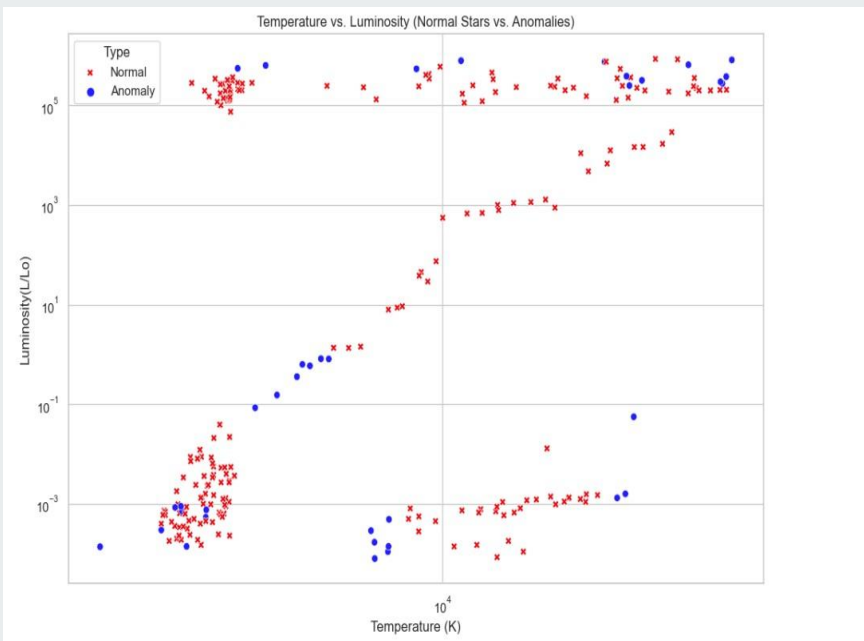
**Astronomical Validation:** Initiate collaborations with observatories to conduct targeted observations of the identified anomalies, especially those with extreme temperatures and luminosities. This step is vital for confirming their astronomical significance and understanding their physical characteristics.

**Feature Expansion:** Incorporate advanced spectral analysis into the dataset to differentiate between anomalies that represent new types of stars and those that are statistical outliers. This could lead to the discovery of previously unrecognized patterns in stellar evolution.

**Model Exploration:** Develop a hybrid model that combines the strengths of Isolation Forest and LOF with deep learning techniques to improve the detection of subtle anomalies. Testing these models on a wider dataset could reveal more about the universe's complexity.

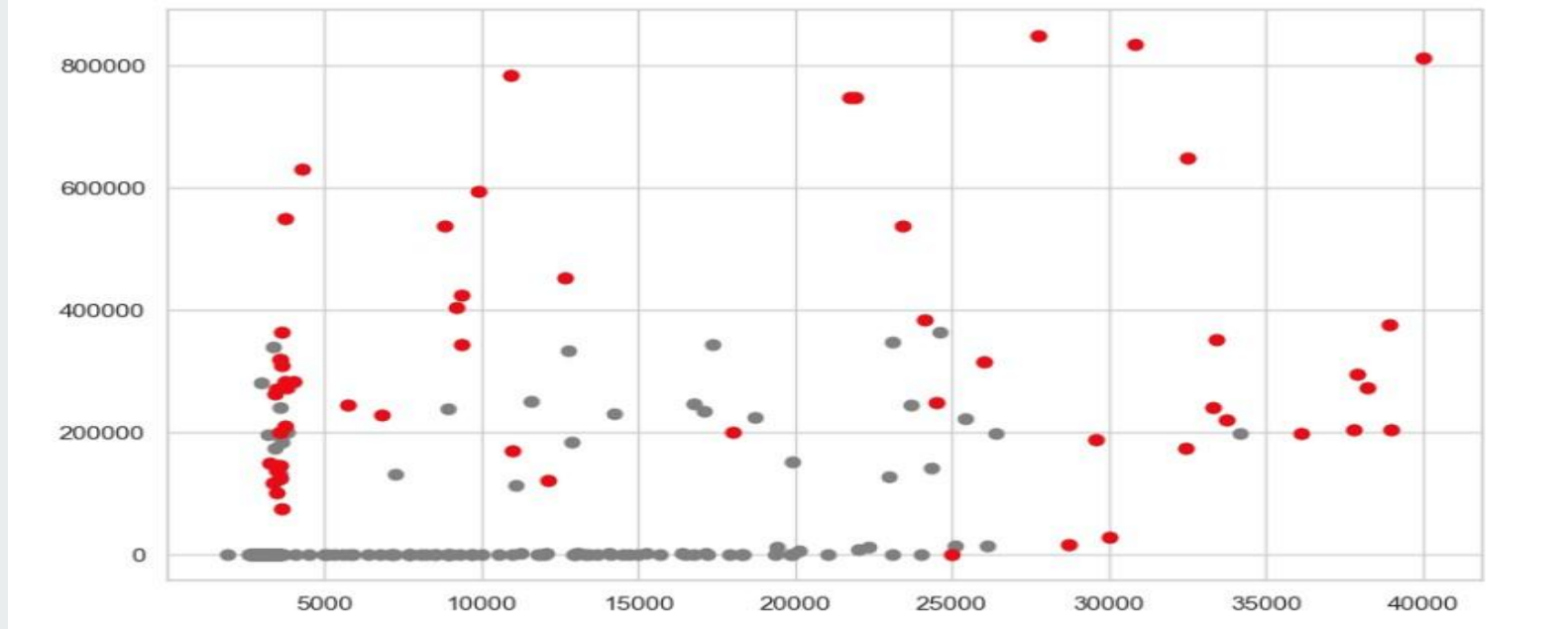
## Key Insights and Recommendations

Let's take a look at the anomalies in conjunction with the typical stars from LOF vs Isolation Forest



## Key Insights and Recommendations

Let's take a look at the anomalies in conjunction with the anomalies LOF vs Isolation Forest for temperature vs luminosity





## Conclusions

**Summary of Achievements:** I highlighted the application of advanced anomaly detection techniques, identifying significant outliers within the stellar dataset. This analysis pinpointed stars with extreme attributes, offering potential new insights into celestial classifications.

**Impact of Findings:** The discovery of anomalies with extreme temperatures and luminosities challenges existing understandings of stellar behavior, suggesting the presence of rare or unusual star types. I can see how this would lead to reevaluation of current models of stellar evolution.

**Future Potential:** This work demonstrates the critical role of data science in astrophysics. This could be the start of paving the way for future studies that could leverage more sophisticated models and datasets to explore further.