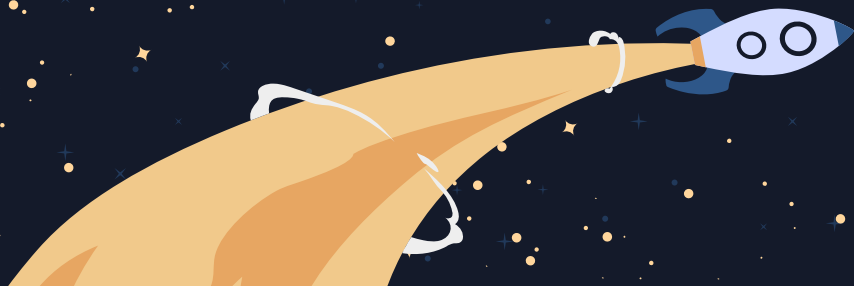


# A World Away: Hunting for Exoplanets with AI

---

Super Code Ninja Girlies

Shalisha Mathew  
Ronja Matthews  
Breanna Hernandez  
Stefany Buccat



Thousands of new exoplanets have been discovered manually from data collected through space-based exoplanet surveying missions. To overcome the manual process of exoplanet identification, artificial intelligence and machine learning (AI/ML) can be used to automatically analyze large sets of data and identify exoplanets.

**OUR GOAL:** create and train an AI/ML model using NASA's open-source exoplanet datasets to analyze data for the purpose of identifying exoplanets

NASA Exoplanet Archive

HomeAbout UsDataToolsSupportLogin

Select ColumnsDownload ToolsFilter DataDownload Data ProductsView DocumentationUser Preferences

Cumulative KOI Data

KOI#D	KOI Name	Kepler Name	Engaged	Active	Disposition	Using Data	Disposition Score	Not Trans-Listed	Star	Eclipse	Central Offset	Ephemeris	Orbit Period (days)	Transit Epoch	Impact Parameter	Transit Duration	Transit Depth (ppm)	Planetary Radius (Earth radii)	Equilibrium Temperature (K)				
10787469	KOI2572.01	Kepler-227.0	CONFIRMED	CANDIDATE	1.0000	0	0	0	0	0	0	0	6.48805572-0.779e-05	170.538766-00218	0.146	2.395701-0819	615.8195	2.26	783	93.56			
10787469	KOI2572.02	Kepler-227.0	CONFIRMED	CANDIDATE	0.8950	0	0	0	0	0	0	0	6.49439372-0.00274e-05	12.35845-0007	0.068	4.9507-105	874.6505	2.83	443	91.9			
10787469	KOI2572.03	Kepler-227.0	CONFIRMED	CANDIDATE	0.8950	0	0	0	0	0	0	0	1.99891399e-01-0.00017e-05	170.080232-0000	0.099	1.7522-10-01	1920.71-01	1.61	710	105			
10844649	KOI678.01	Kepler-44.0	FALSE POSITIVE	FALSE POSITIVE	0.0000	0	0	0	0	0	0	0	1.75804524-0.24e-07	170.390556-0001	1.276	2.4064-0710	809.2122	3.16	535	89.19			
10844649	KOI678.02	Kepler-44.0	CONFIRMED	CANDIDATE	1.0000	0	0	0	0	0	0	0	1.52087773-0.73e-04	171.595561-0011	0.701	1.6544-04-02	60.36-19e	2.76	1400	90.21			
10844649	KOI678.03	Kepler-44.0	CONFIRMED	CANDIDATE	0.9900	0	0	0	0	0	0	0	1.62026326-0.00023e-05	171.212561-0001	0.018	1.5717-02-04	157.542-04	3.48	185	104.19			
10872992	KOI2367.01	Kepler-228.0	CONFIRMED	CANDIDATE	0.9900	0	0	0	0	0	0	0	1.54331431-0.004e-05	170.32797-0015	0.342	1.3426-07-03	686.66-07	2.17	621	92.76			
10872992	KOI2367.02	Kepler-228.0	CONFIRMED	CANDIDATE	0.9900	0	0	0	0	0	0	0	2.56680871-0.78e-05	170.35947-0001	0.061	1.2420-10-15	256.15e-10	1.50	1300	80.77			
872123	KOI149.11	Kepler-20.0	FALSE POSITIVE	FALSE POSITIVE	0.0000	0	0	0	0	0	0	0	1.26011958-0.28e-05	12.220595-0003	1.169	5.0224-10-18	23.70-18e	39.51	142	78.72			
10872992	KOI2367.03	Kepler-228.0	CONFIRMED	CANDIDATE	0.9900	0	0	0	0	0	0	0	1.95466481-0.10e-05	171.573197-0001	0.018	1.5545-03-01	494.543-03	1.61	530	102.64			
10872992	KOI2367.04	Kepler-228.0	CONFIRMED	CANDIDATE	0.9900	0	0	0	0	0	0	0	2.47061377-0.14e-08	172.303556-0e-06	0.180-01	1.7474-01-07	162.924-07	13.36e-01	1330	70.14			
10872992	KOI2367.05	Kepler-228.0	CONFIRMED	CANDIDATE	0.9900	0	0	0	0	0	0	0	2.25475841-0.14e-08	172.385411-0e-06	0.224-19	3.8886-03-01	642.917-16	1.61	506	104.06			
892224	KOI149.01	Kepler-9.0	FALSE POSITIVE	FALSE POSITIVE	0.0000	0	0	0	0	0	0	0	1.52440064-0.96e-07	171.174922-04-01	0.610-01	3.5948-03-03	649.74e-6	14.05e-11	1521	130.4			
10844649	KOI678.04	Kepler-44.0	CONFIRMED	CANDIDATE	0.9900	0	0	0	0	0	0	0	1.78871475-0.00e-05	170.335616-0001	0.041	2.6304-07-07	191.14-10	1.9e-10	528	90.74			
1041211	KOI224.01	Kepler-61.0	CONFIRMED	FALSE POSITIVE	0.0000	0	0	0	0	0	0	0	1.52524868-0.00e-05	170.335616-0001	2.482	3.3016-07-04	17984.301-01	10.5e-10	733	75.88			
1048478	KOI249.01	Kepler-11.0	FALSE POSITIVE	FALSE POSITIVE	0.0000	0	0	0	0	0	0	0	19.40377762-0.00e-06	172.482324-0000	0.004	12.2159-05-06	189.7533-10	7.23e-10	523	117.80			
10872992	KOI2367.04	Kepler-228.0	CONFIRMED	FALSE POSITIVE	0.0000	0	0	0	0	0	0	0	1.23713985-0.12e-06	174.162153-04-01	0.101	4.7964-03-03	724.24-10	49.2e-10	888	55.97			
1048478	KOI249.02	Kepler-11.0	FALSE POSITIVE	FALSE POSITIVE	0.0000	0	0	0	0	0	0	0	1.847378-0.00e-06	180.781738-0000	0.202	3.4278-03-06	16978.338-06	7.8e-06	566	28.81			
1025464	KOI249.01	Kepler-60.0	CONFIRMED	CANDIDATE	1.0000	0	0	0	0	0	0	0	0.72738717-0.07e-06	172.251150-0007	0.337	3.2875-03-09	1308.26-18	2.40	678	61.85			
1053368	KOI274.01	Kepler-61.0	CONFIRMED	CANDIDATE	1.0000	0	0	0	0	0	0	0	0.62093263-0.00e-06	171.850299-0007	0.258	1.5521-03-11	1937.73-04	2.87	446	54.04			
1053368	KOI274.02	Kepler-61.0	CONFIRMED	CANDIDATE	1.0000	0	0	0	0	0	0	0	1.26670973-0.15e-06	170.72994-0023	0.044	3.51e-06-06	397.86-05	1.8e-10	598	39.5			
1005124	KOI124.01	Kepler-20.0	CONFIRMED	CANDIDATE	1.0000	0	0	0	0	0	0	0	0.54391818-0.00e-06	171.869684-0007	0.002	3.0276-07-01	851.61-28	1.40e-10	919	106.13			
1005124	KOI124.02	Kepler-20.0	CONFIRMED	CANDIDATE	1.0000	0	0	0	0	0	0	0	1.84952211-0.00e-06	170.860624-0003	0.228	2.9466-07-07	363.31-19	2.7e-10	108	255.1			
1005124	KOI124.03	Kepler-20.0	CONFIRMED	CANDIDATE	0.9710	0	0	0	0	0	0	0	0.10094408-0.78e-05	12.82110-0028	0.005	3.5966-10-18	217.44e-18	1.00	97.80				
1007576	KOI249.01	Kepler-61.0	CONFIRMED	CANDIDATE	1.0000	0	0	0	0	0	0	0	3.8640325-0.00e-06	170.860624-0003	0.016	6.8907-10-14	2402.64-04-05	1.41e-10	267	114			
1056232	KOI249.01	Kepler-60.0	CONFIRMED	CANDIDATE	1.0000	0	0	0	0	0	0	0	6.21699914-0.38e-05	171.535544-0008	0.098	2.2155-05-06	887.227-1	2.86	520	103.64			
1056241	KOI275.01	Kepler-60.0	CONFIRMED	CANDIDATE	1.0000	0	0	0	0	0	0	0	6.49877881-0.38e-05	171.535544-0008	0.171	2.1717-07-02	834.7287-2	2.76	487	106.6			

Showing records 1 of 27 (64 of 954 total)

DOI:10.26333/4E

Clear CacheCheck AllReset Rows

# OUR APPROACH

## 1 Initial Approach: Random Forest + CNN on Light Curves

- At first, we tried combining **Random Forest models on 4 basic features** with **1D and 2D CNNs** using the raw light curve data (lightcurve)
- The idea was to use the CNN to capture patterns in the light curve shapes, and the Random Forest to combine that with numeric features
- However, this required **folding the light curves for each star and period** using a for loop
- It quickly became **inefficient** because there were thousands of stars to process
- Random Forest was ideal because:
  - It **handles both numeric and categorical features** easily
  - It's **robust to outliers and noisy features**
  - It provides **feature importance** automatically, letting us identify which features contribute most

## 2 Focus on Outliers and Data Cleaning

- Next, we decided to **handle outliers** in the numeric features (like transit depth, duration, planetary radius)
- We used the **IQR method**: removed values outside  $1.5 \times$  the interquartile range
- After cleaning, the Random Forest model performed **better on the confusion matrix**:
  - More true positives were correctly classified
  - False positives decreased
  - Overall accuracy improved

```
# Handle Outliers

# List of numeric columns to check for outliers
numeric_cols = ["koi_period", "koi_duration", "koi_prad", "koi_depth"]

# Create copies to avoid modifying original X and y
X_clean = X.copy()
y_clean = y.copy()

for col in numeric_cols:
    # Calculate IQR
    Q1 = X_clean[col].quantile(0.25)
    Q3 = X_clean[col].quantile(0.75)
    IQR = Q3 - Q1
    lower_bound = Q1 - 1.5 * IQR
    upper_bound = Q3 + 1.5 * IQR

    # Remove outliers
    mask = (X_clean[col] >= lower_bound) & (X_clean[col] <= upper_bound)
    X_clean = X_clean[mask]
    y_clean = y_clean[X_clean.index]
```

# OUR APPROACH

## ③ Hyperparameter Tuning

- We **increased the number of trees** from 100 → 300 in the Random Forest
  - More trees improve model **stability** and **reduce variance**.
- We also adjusted parameters like `max_depth`, `min_samples_split`, and `min_samples_leaf` to prevent overfitting

```
# Train Random Forest

# Split into training and test sets
X_train, X_test, y_train, y_test = train_test_split(X_clean, y_clean, test_size=0.2, random_state=42)

#rf_model = RandomForestClassifier(n_estimators=100, random_state=42)
#rf_model.fit(X_train, y_train)
rf_model = RandomForestClassifier(
    n_estimators=300,      # more trees = better stability
    max_depth=15,         # prevent overfitting
    min_samples_split=5,  # don't split tiny nodes
    min_samples_leaf=3,   # ensures each leaf has enough data
    random_state=42
)

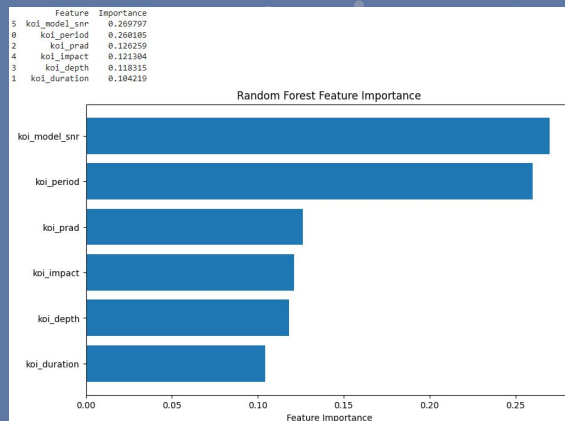
rf_model.fit(X_train, y_train)
```

## ④ Adding More Features

- Initially, we used only 4 features
- We realized that **adding more features could help**, but adding **too many irrelevant features** could actually decrease model performance
- After research, we chose **10 features** that made the most sense physically (transit parameters + stellar properties)

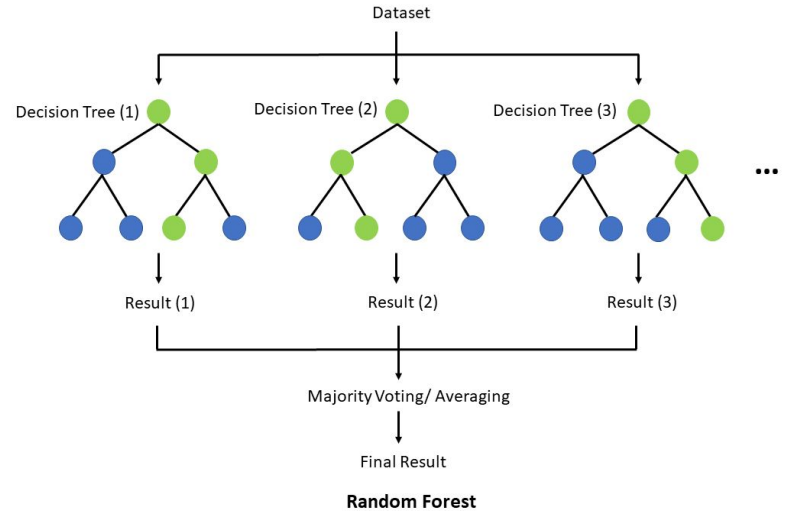
## ⑤ Feature Importance Analysis

- We used Random Forest's **built-in feature importance**:
  - Each tree in the forest splits data based on features to reduce uncertainty (impurity)
  - Features that **consistently help separate classes** get higher importance scores
- We plotted a **bar chart** to visualize which features mattered most
  - The top features (like `koi_model_snr`, `koi_period`) had the highest scores
  - The lowest-ranked features were removed, but the model's performance **didn't change much**, confirming they weren't very useful



# WHAT IS THE RANDOM FOREST CLASSIFIER?

The Random Forest classifier is the AI/ML model's core part. It is an ensemble of 300 decision trees, where each tree will focus and analyze its own data. The trees will then create a prediction (false positive, candidate, or confirmed exoplanet), and the forest will average all the individual predictions for a final prediction. This process was used due to its ability to ensure high accuracy and reliability in determining the most influential features for exoplanet identification, as well as analyzing complex data obtained from space agency open-data, while being robust to noise or outliers.



# RESULTS AND PERFORMANCE OF THE AI/ML MODEL

- Our AI/ML model analyzed NASA's Kepler exoplanet data using 6 features:
  - koi\_period: orbital period
  - koi\_depth: transit depth
  - koi\_duration: transit duration
  - koi\_impact: impact parameter
  - koi\_prad: planetary radius
  - koi\_model\_snr: transit signal-to-noise
- The model removes outliers in koi\_period, koi\_duration, koi\_prad, and koi\_depth. With the cleaned data, it will train a Random Forest classifier and evaluate its accuracy. A visual representation of each feature's importance will be displayed through a bar graph, to show which features helped the most to separate classes
- Finally, the model will display:
  - A classification report to summarize the identification of false positives, candidates, and confirmed exoplanets
  - A confusion matrix to compare the predicted labels and the true labels
- After adding more features, the model's accuracy improved from 69% to 78%

Classification Report:

	precision	recall	f1-score	support
0	0.76	0.67	0.71	366
1	0.50	0.42	0.46	245
2	0.71	0.84	0.77	459
accuracy			0.69	1070
macro avg	0.66	0.64	0.65	1070
weighted avg	0.68	0.69	0.68	1070

Confusion Matrix:

```
[[244 55 67]
 [ 50 103 92]
 [ 25 48 386]]
0.6850467289719626
```

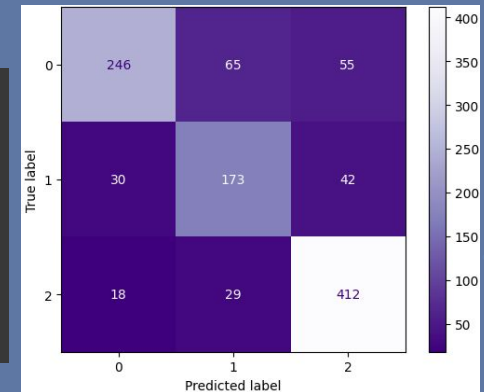


Classification Report:

	precision	recall	f1-score	support
0	0.84	0.67	0.75	366
1	0.65	0.71	0.68	245
2	0.81	0.90	0.85	459
accuracy			0.78	1070
macro avg	0.76	0.76	0.76	1070
weighted avg	0.78	0.78	0.77	1070

Confusion Matrix:

```
[[246 65 55]
 [ 30 173 42]
 [ 18 29 412]]
0.7766355140186916
```



# USER EXPERIENCE

A web interface platform was created with Streamlit to provide users with accessibility to the AI/ML model. The platform allows users to upload Kepler dataset files from space agency open data. The model will then analyze the uploaded data and display exoplanet identification results. With this interactive platform, scientists, researchers, and students can learn and explore exoplanet data.

Settings

Outliers

☒

Sigma filter (z-score)

1

Sigma threshold

2.00

1

IQR filter

IQR factor

1.50

Train/Test split

Test size

0.20

42

random\_state

--

Random Forest

n\_estimators

300

12

max\_depth (0 = None)

☒

class\_weight="balanced"

Prediction

"False Positive" : 0.6814

"Candidate" : 0.6376

"Confirmed" : 0.681

Predicted: Confirmed

Kepler Exoplanet Identifier

Upload Kepler KDI tabular data, clean outliers (sigma/IQR), train a Random Forest, and predict disposition.

Data

Upload your kepler\_data.csv (KDI table). If your file has header notes, set skipsrows accordingly.

Upload CSV

Drag and drop file here

Limit: 200MB per file • CSV

Browse files

skipsrows (for header notes)

53

No CSV uploaded yet.

Train model & evaluate

Test Accuracy

0.748

+4 -0.106 vs train

Train accuracy: 0.944

Classification report

	precision	recall	f1-score	support
False Positive	0.7727	0.6631	0.7137	282
Candidate	0.6964	0.6959	0.6481	217
Confirmed	0.8391	0.8354	0.8272	403
accuracy	0.7478	0.7478	0.7478	0.7478
macro avg	0.7827	0.7315	0.7297	950
weighted avg	0.7533	0.7478	0.7484	950