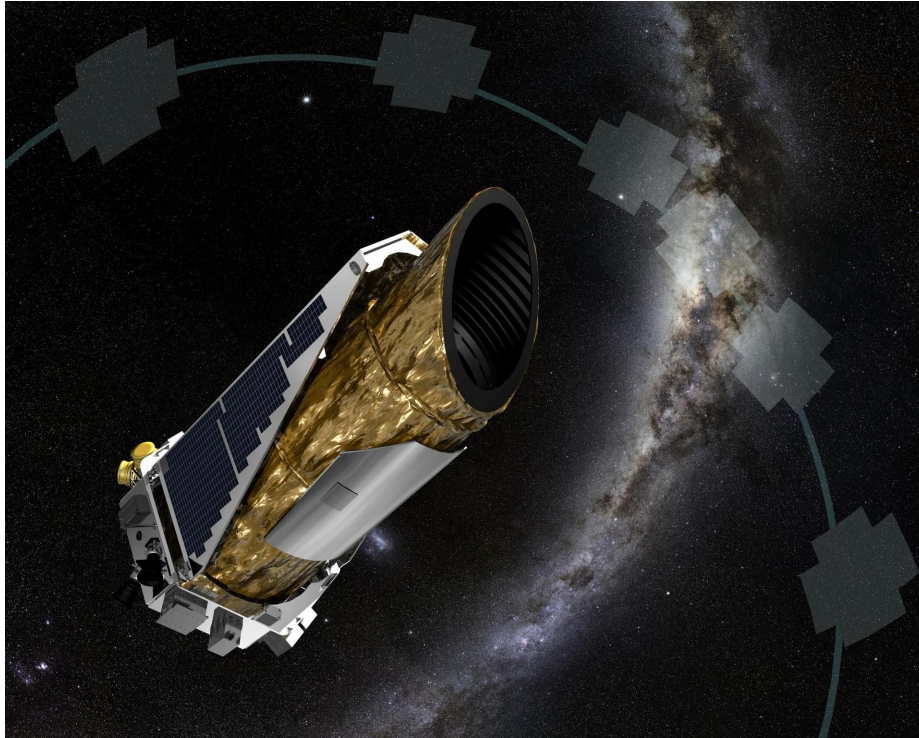# Project Report: A Machine Learning Model for Classifying Potential Exoplanets from Kepler Data
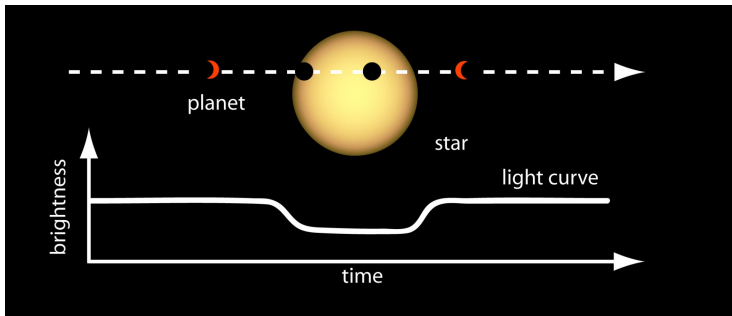


The Kepler Space Telescope. Source: NASA



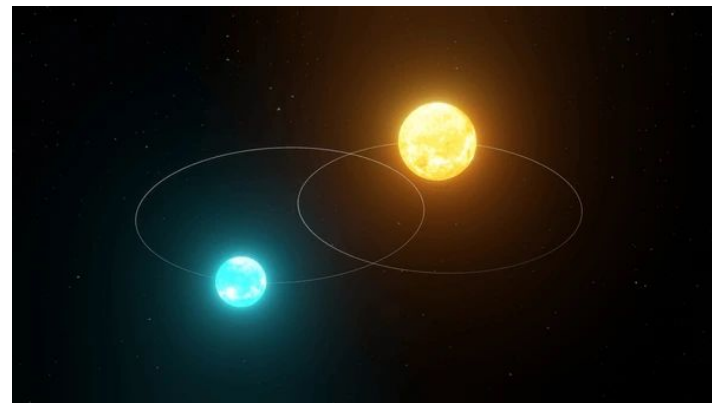Johannes Kepler (1571 - 1630)

# Background



- NASA's Kepler Space Telescope has detected thousands of exoplanets (planets orbiting stars outside of our solar system) by measuring dips in stellar brightness caused by planets passing in front of the star.

- Not all changes in brightness are caused by exoplanets; some are caused by other stars or measurement noise.

- Exoplanet candidates are called Kepler Objects of Interest (KOIs), and are characterized by many measured variables and categorical features.

- The NASA Exoplanet Archive contains a table of past KOIs which have been marked as **CONFIRMED** (true exoplanet) or **FALSE POSITIVE** based on follow-up observations.

Project Goal: Accelerate the process of exoplanet discovery by training a machine learning model to classify KOIs from the NASA Exoplanet Archive as <span style="color:green">CONFIRMED</span> or <span style="color:red">FALSE POSITIVE</span>

- Utilize tabulated data from the KOI Cumulative Table.

- Try both supervised and unsupervised methods.

- Produce an accurate, *interpretable* model with visualizations of results.

- Discover insights about which features of the dataset are most useful.

# Related Work

- Morton et. al (2016) developed a statistical modeling tool called VESPA for calculating false positive probabilities from several combined data sources [1].

- Shallue & Vanderburg (2018) achieved high accuracy (>98%) using convolutional neural networks (CNNs) trained on time-series Kepler light curve data [2].

- Our project complements these approaches with a more interpretable classifier based on tabular data.

- The Kepler False Positive Working Group gained insight into some of the most common causes of false positives:

  - Orbiting binary star
  - Fluctuation in brightness
  - Electronic crosstalk

1. Morton, T.D. et al. *False positive probabilities for all Kepler objects of interest.* (2016). The Astrophysical Journal.
2. Shallue, Christopher J and Vanderburg, Andrew. *Identifying Exoplanets with Deep Learning.* (2018). The Astronomical Journal.

# Data Understanding

- Studied the features of the KOI Cumulative Table and selected the most useful ones, omitting irrelevant or redundant features.

- Removed features that gave away information about the target variable, such as "false positive flags" indicating the reason for a false positive.

- Removed features that were derived from combinations of other features, such as "planetary radius", which was estimated based on star and telescope data already in the table.

# Data Inspection and Cleaning

- After feature selection, data was imported to Pandas dataframe containing 9564 rows and 18 columns.

- Surprisingly, features are all numerical.

- Rows without a positive or negative target label (marked as CANDIDATE) were removed, leaving 7585 entries with a 36%-64% positive-negative split.

- Most columns had some missing values.

- Missing values were imputed using the median of the non-null values.

```
df['koi_disposition'].value_counts()
```
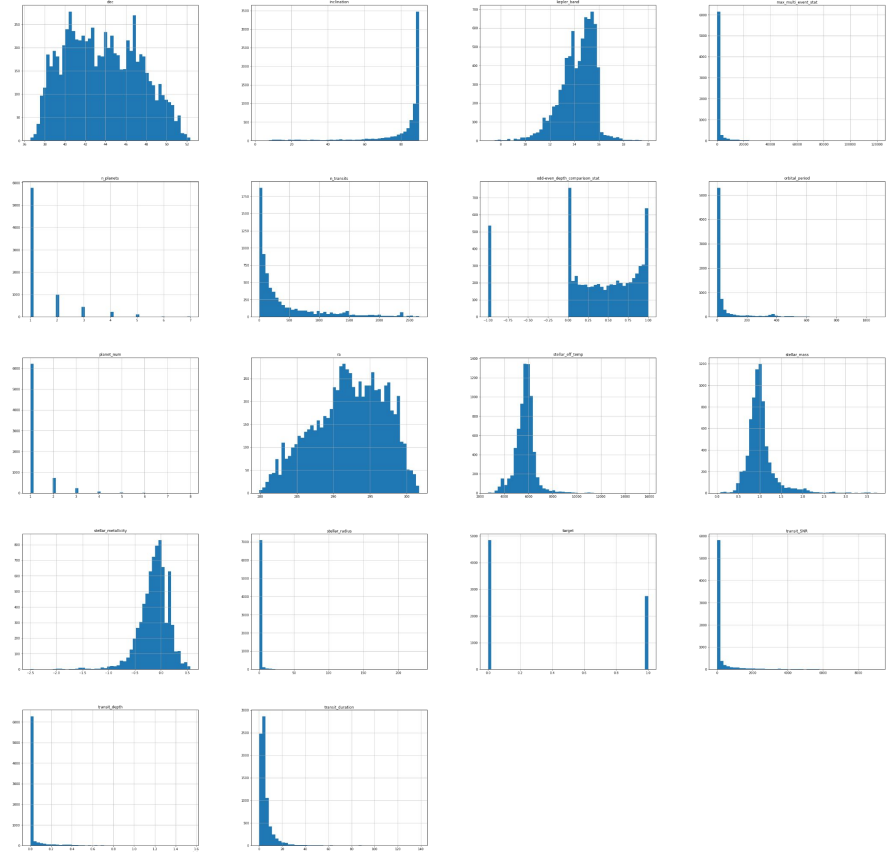
```
FALSE POSITIVE    4839
CONFIRMED         2746
CANDIDATE         1979
Name: koi_disposition, dtype: int64
```

```
data.info()
```

```
<class 'pandas.core.frame.DataFrame'>
Int64Index: 7585 entries, 0 to 9563
Data columns (total 18 columns):
 #   Column                        Non-Null Count  Dtype
---  ------                        --------------  -----
 0   orbital_period                7585 non-null   float64
 1   transit_duration              7585 non-null   float64
 2   transit_depth                 7326 non-null   float64
 3   inclination                   7325 non-null   float64
 4   max_multi_event_stat          6904 non-null   float64
 5   transit_SNR                   7326 non-null   float64
 6   n_planets                     7585 non-null   int64
 7   n_transits                    6904 non-null   float64
 8   planet_num                    7299 non-null   float64
 9   odd-even_depth_comparison_stat 6675 non-null   float64
 10  stellar_eff_temp              7326 non-null   float64
 11  stellar_metallicity           7306 non-null   float64
 12  stellar_radius                7326 non-null   float64
 13  stellar_mass                  7326 non-null   float64
 14  ra                            7585 non-null   float64
 15  dec                           7585 non-null   float64
 16  kepler_band                   7584 non-null   float64
 17  target                        7585 non-null   int64
dtypes: float64(16), int64(2)
memory usage: 1.1 MB
```

# Exploratory Data Analysis

- Plotted distributions of each variable using histograms.

- Some variables are normally distributed, others exponentially distributed.

- One variable codes missing data as -1, treated same as null

# Exploratory Data Analysis

- Plotted correlation matrix using Seaborn heatmap function.

- Not many strong correlations between variables.

- Statistical significance variable strongly correlated with magnitude of signal - perhaps redundant

- Star mass and temperature strongly correlated

# Data Preprocessing

- Normally-distributed variables transformed by Z-score.

- Log transformation for exponentially-distributed variables, followed by normalization.

- Principal Component Analysis (PCA) further transformed data

# Exploratory Data Analysis

- Scatterplot of first two PCA components.

- True exoplanets in yellow, false positives in blue.

- 49% of explained variance

Plot of explained variance vs. number of components:

# Exploratory Data Analysis

- 3-D Plotly Scatterplot of first three PCA components.

- 58% of explained variance

- Less overlap than 2-D plot

- Looks promising for clustering model.

# Data Modeling - Decision Tree

- Split data into training set and testing set with 80%-20% split.

- Used non-transformed data, since decision tree uses only thresholds.

- Trained a baseline decision tree using SciKit-learn library with default parameters

  - No limit for number of variables or leaves in the tree
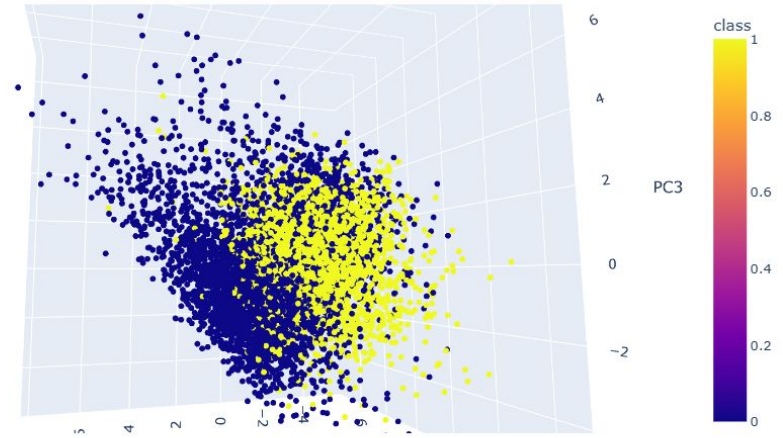
- Achieved baseline testing accuracy of 88%.

- Precision (accuracy on predicted positives) and recall (percentage of true positives identified) were lower, at 81% and 83%, respectively.



Confusion Matrix

|  | False Positive | Confirmed Exoplanet |
|---|---|---|
| **False Positive** | 902.0000 | 99.0000 |
| **Confirmed Exoplanet** | 88.0000 | 428.0000 |

# Data Modeling - Decision Tree

- Optimized decision tree model with SciKit-learn GridSearchCV() function

  - Hyperparameter sweep with 5-fold cross-validation

  - Optimized for ROC-AUC score to try to improve recall

- Accuracy: 92.4%

- Precision: 88.9%

- Recall: 90.7%



Confusion Matrix

```
print("Best Parameters:", grid_search.best_params_)
print("Best Cross-Validation Score:", grid_search.best_score_)

Best Parameters: {'criterion': 'entropy', 'max_depth': 10, 'min_samples_leaf': 4, 'min_samples_split': 10}
Best Cross-Validation Score: 0.9374086565335495
```

# Data Modeling - Random Forest

- Optimized with GridSearchCV()

  - Best AUC: ensemble of 200 decision trees randomly assigned 4 features each

- Required 15 minutes to run grid search, vs. 30 seconds for decision tree

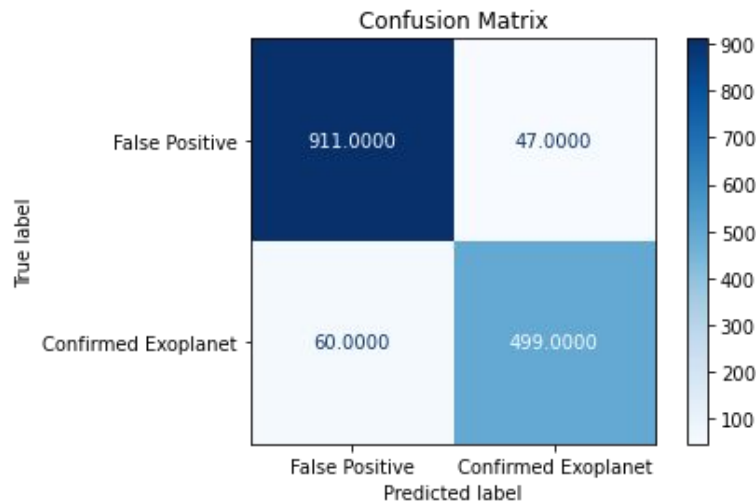- Accuracy: 92.9%

- Precision: 91.4%

- Recall: 89.3%



Confusion Matrix

```
print("Best Parameters:", grid_search.best_params_)
print("Best Cross-Validation Score:", grid_search.best_score_)

Best Parameters: {'bootstrap': True, 'max_features': 'sqrt', 'min_samples_leaf': 2, 'min_samples_spli
t': 5, 'n_estimators': 200}
Best Cross-Validation Score: 0.9730332920978869
```

# Data Modeling - Agglomerative Clustering

- Unsupervised learning model - no train-test split

- Used PCA-transformed data.

- Performed parameter sweep of distance metrics, cluster-linkage metrics, and number of PCA components included.

- Best model used 8 PCA components

  - About 85% of explained variance

- Only 68% accurate



Confusion Matrix

# Data Modeling - Dense Neural Network

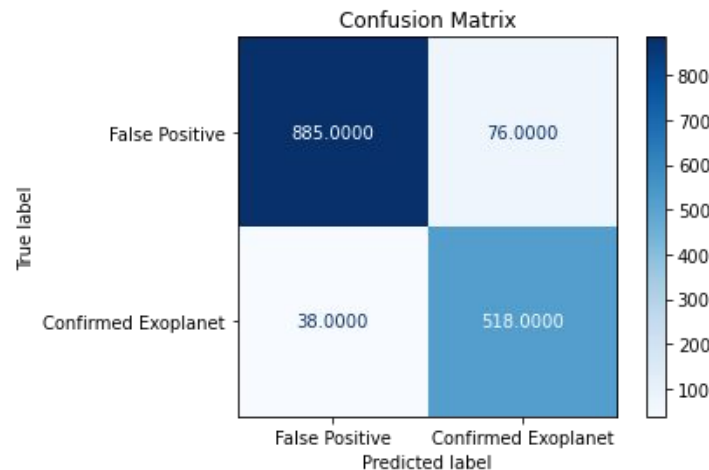- Used all components of PCA-transformed data.

- 3 dense layers plus dropout layer

- Accuracy: 92.5%

- Precision: 87.2%

- Recall: 93.2%



```
model = Sequential([
    Dense(256, activation='relu', input_shape=(17,)),
    BatchNormalization(),
    Dense(128, activation='relu'),
    BatchNormalization(),
    Dense(64, activation='relu'),
    BatchNormalization(),
    Dropout(0.3),
    Dense(1, activation='sigmoid')
])
```

```
model.compile(optimizer='adam', loss='binary_crossentropy', metrics=['accuracy'])
```



Confusion Matrix

# Evaluation Summary

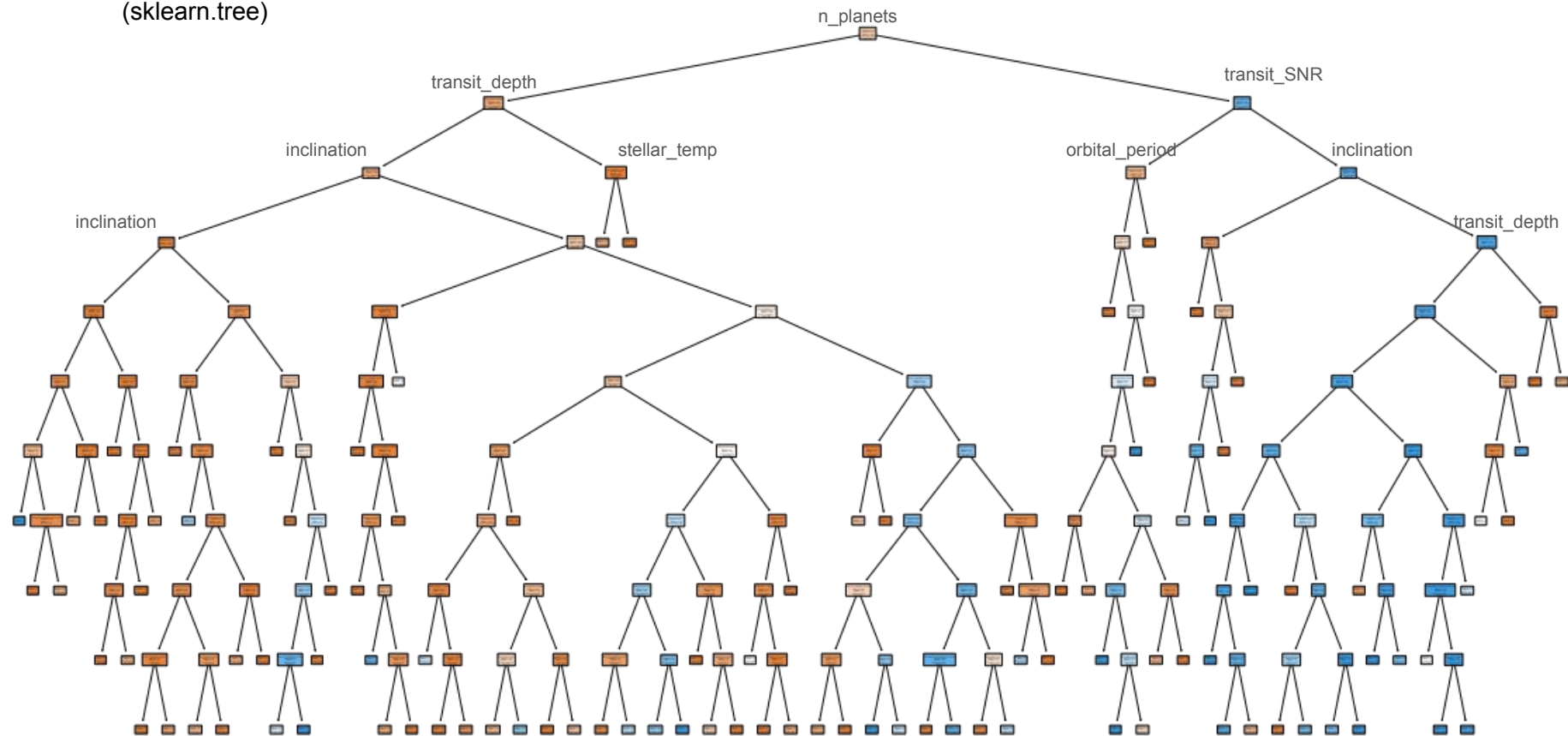| Model | Accuracy | Precision | Recall | F1 Score |
|---|---|---|---|---|
| Decision Tree | 92.4% | 88.9% | 90.7% | 0.898 |
| Random Forest | 92.9% | 91.4% | 89.3% | 0.903 |
| Agglomerative Clustering | 67.8% | 54.5% | 67.0% | 0.601 |
| Dense Neural Network | 92.5% | 87.2% | 93.2% | 0.901 |

- Supervised models (Decision Tree, Random Forest, Dense Neural Network ) greatly outperformed the unsupervised model (Agglomerative Clustering).

- A single optimized decision tree performed just about as well as more powerful methods.

- The interpretability of the decision tree model makes it the ideal choice for this project.

# Decision Tree Visualization
(sklearn.tree)

# Project Timeline

- Phase 1 - Data Cleaning and Understanding, including EDA Visualizations (2 days)

- Phase 2 - Data Preprocessing and Feature Selection (1 day)

- Phase 3 - Baseline Model Development (1 day)

- Phase 4 - Model Evaluation and Hyperparameter Tuning (2 days)

- Phase 5 - Final Presentation of Results (1 day)

# Future Work

- Try to understand why the clustering approach didn't work.

  - Try different feature selection, data transformation, feature engineering

- Optimize a model for subset of data with no other discovered planets in the star system

  - Number of planets in the system is highly predictive - could be a limitation of the model

# Key Takeaways

- The project succeeded in building an accurate, efficient, interpretable classifier for detecting exoplanets from the NASA Exoplanet Archive's KOI data.

- Data Understanding phase proved to be crucial.

- The feature space ended up being smaller and more numerical than expected.

    - Single decision tree model was sufficient for this feature space.

- Evaluation plan shift - less emphasis on accuracy due to class imbalance.

    - High recall ensures potential exoplanets don't often get missed.