





# Finding Home

# Finding Home

A Kevin Björk Production



What is an exoplanet? - Data - The Goal - The Process - ML - Summarizing Thoughts



What is an exoplanet? - Data - The Goal - The Process - ML - Summarizing Thoughts

LST=

Mercury  
Venus

Earth

Mars

Jupiter

Saturn

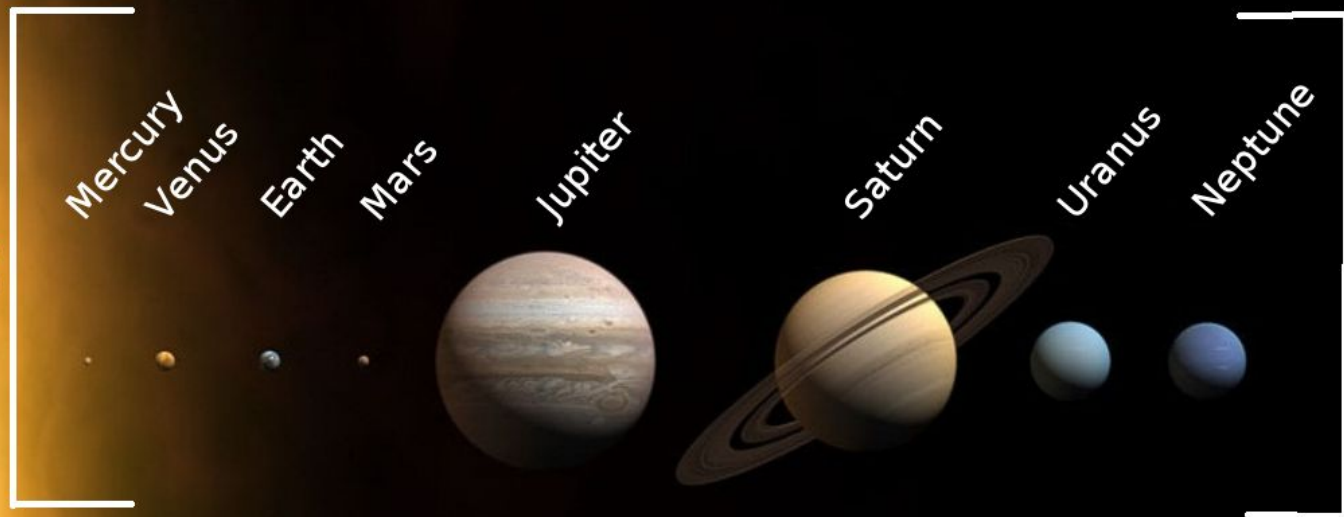
Uranus

Neptune



What is an exoplanet? - Data - The Goal - The Process - ML - Summarizing Thoughts

LST=



EXO NOT IN LST  $\Rightarrow$  TRUE

## The Data

- Before cleaning: 4048 rows x 112 columns
- After cleaning: 4046 rows x 32 columns





## The Data

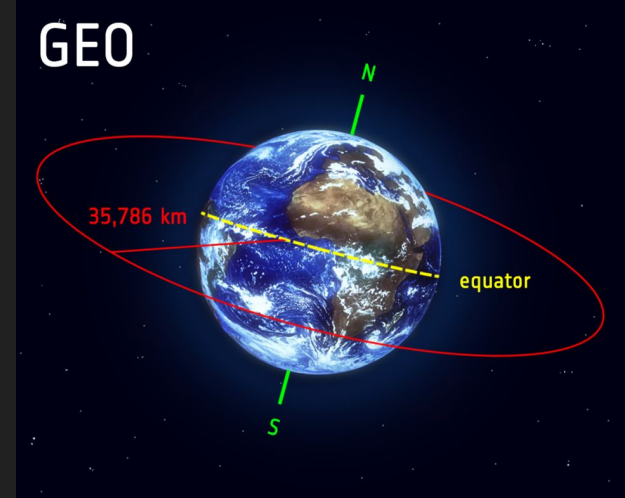
- Before cleaning: 4048 rows x 112 columns
- After cleaning: 4046 rows x 32 columns
- Originates from the Planetary Habitability Laboratory
  - Managed by the University of Puerto Rico



## The Data

- Before cleaning: 4048 rows x 112 columns
- After cleaning: 4046 rows x 32 columns
- Originates from the Planetary Habitability Laboratory
  - Managed by the University of Puerto Rico
- Contains information about exoplanets and the stars they orbits



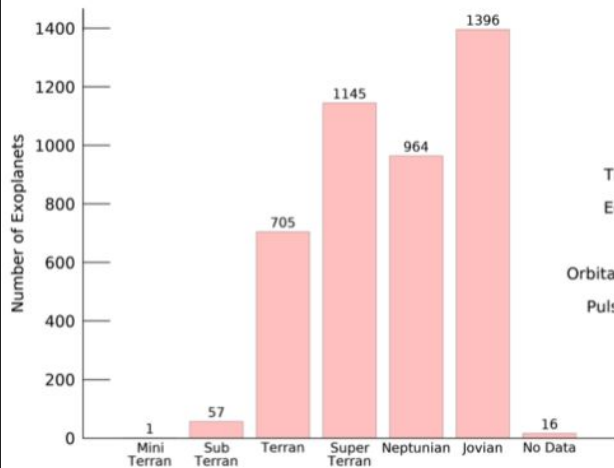


## The Data

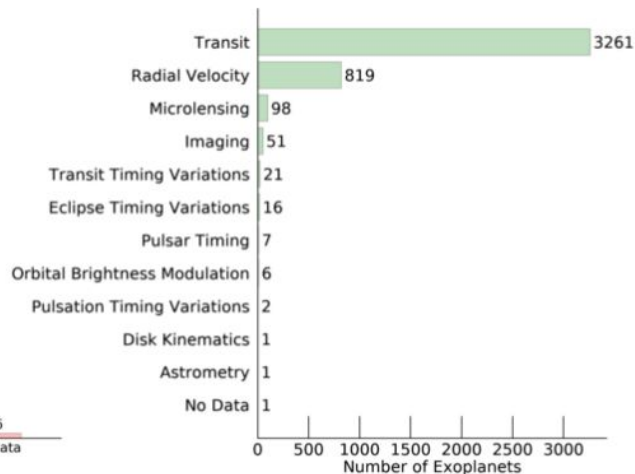
- Before cleaning: 4048 rows x 112 columns
- After cleaning: 4046 rows x 32 columns
- Originates from the Planetary Habitability Laboratory
  - Managed by the University of Puerto Rico
- Contains information about exoplanets and the stars they orbits
- Unbalanced, ca 1 % of planets are habitable



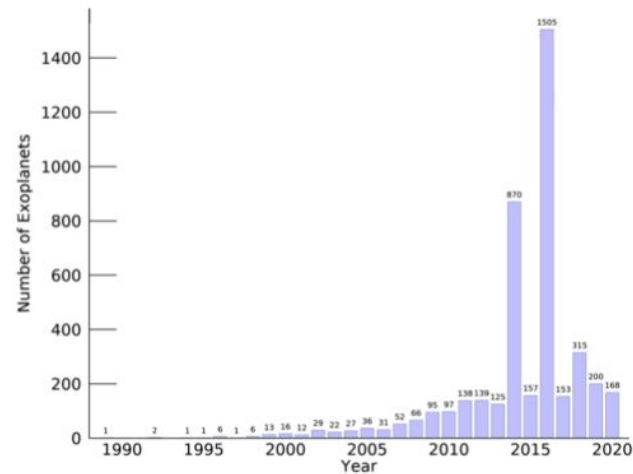
A. Planet Types (total 4284)



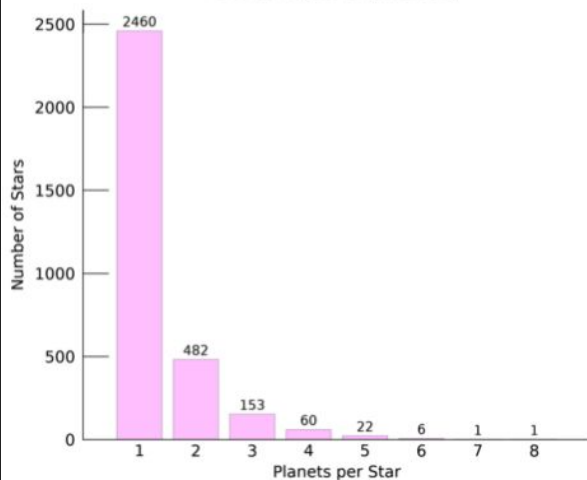
B. Planet Detection Methods (total 4284)



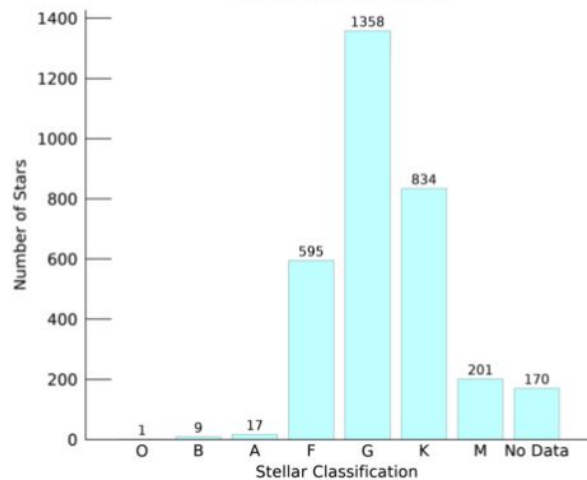
C. Planet Discovery Years (total 4284)



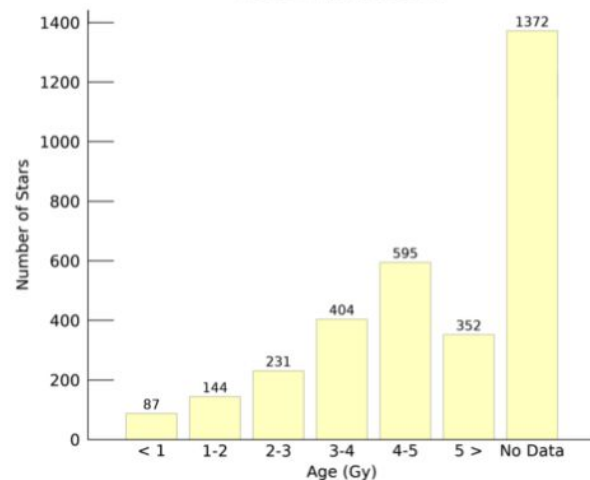
D. Stellar Systems (total 3185)



E. Stellar Types (total 3185)



F. Stellar Ages (total 3185)





What is an exoplanet? - Data - **The Goal** - The Process - ML - Summarizing Thoughts

# GOAL:

What is an exoplanet? - Data - **The Goal** - The Process - ML - Summarizing Thoughts

# GOAL:

Is the planet habitable?

What is an exoplanet? - Data - The Goal - **The Process** - ML - Summarizing Thoughts

# The process: DAE

## The process: DAE

- Lots and lots of columns have high correlation





$$g_p = G \frac{m_p}{r_p^2}$$

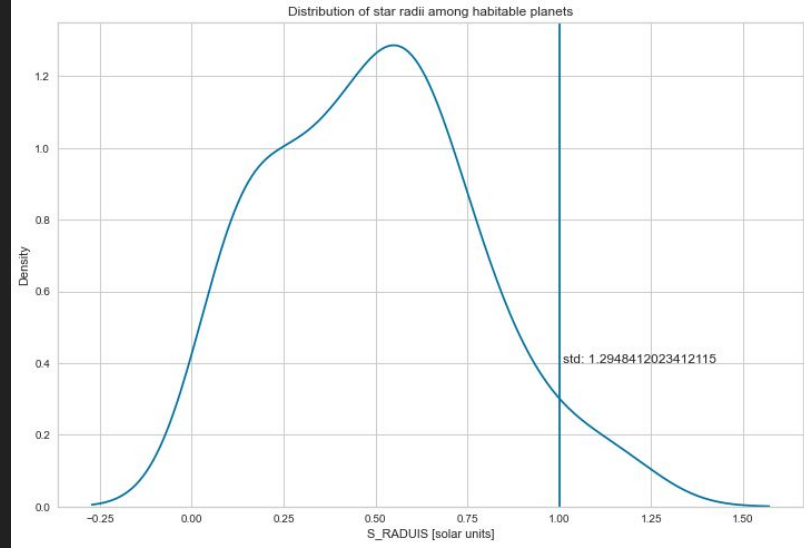
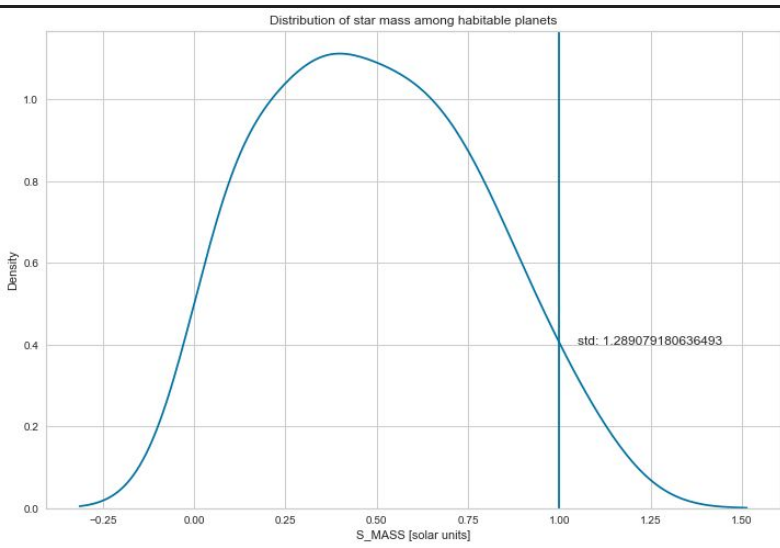
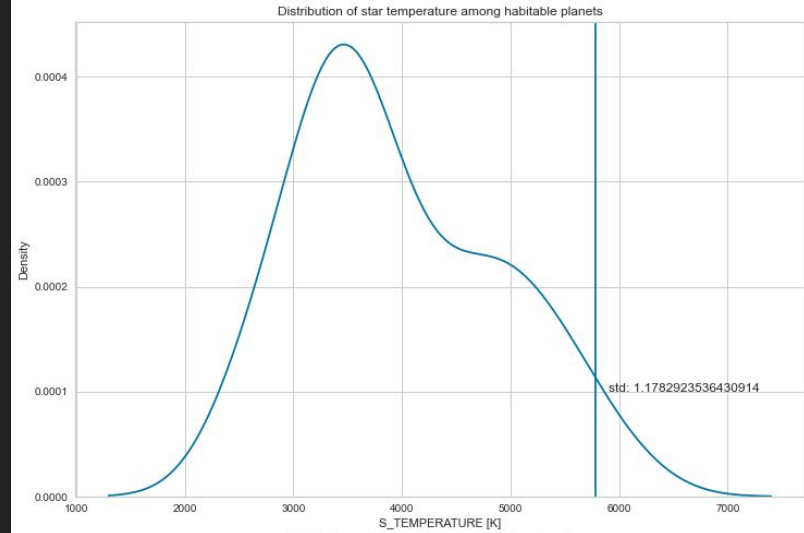
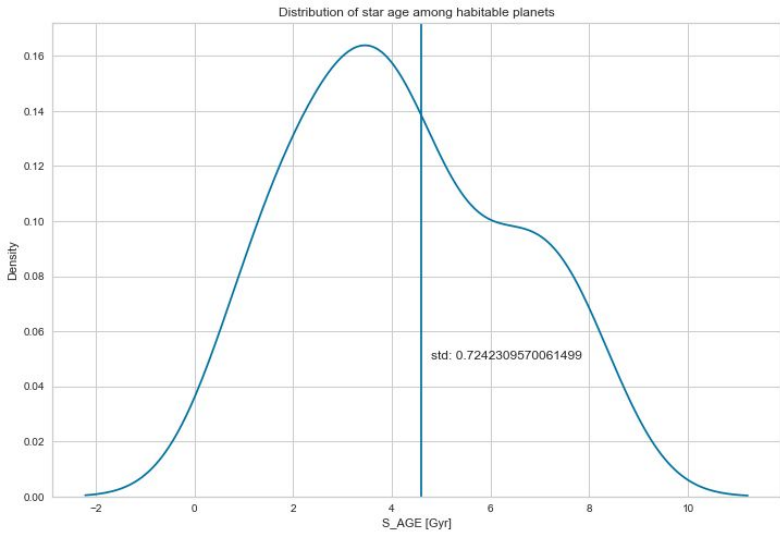
$$L_s = 4\pi r_s^2 \sigma T^4$$

## The process: DAE

- Lots and lots of columns have high correlation
- Feature Creation: SSI

$$SSI = \sum_i \sqrt{(S_{exo,i} - S_{sun,i})^2}$$

$$S_i = [Mass, Radius, \dots]$$



What is an exoplanet? - Data - The Goal - **The Process** - ML - Summarizing Thoughts

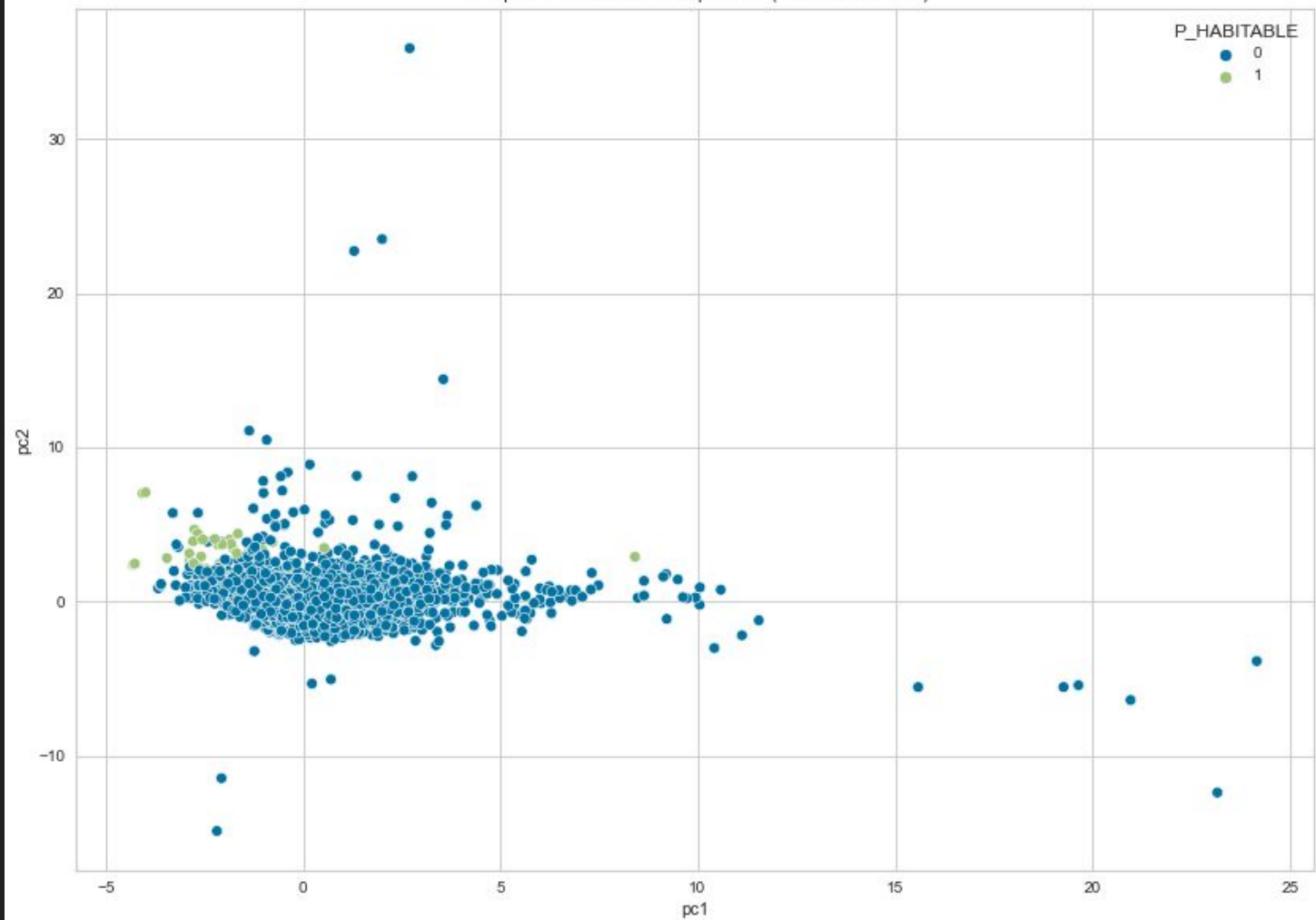
## The process: EDA - Habitable vs Inhabitable



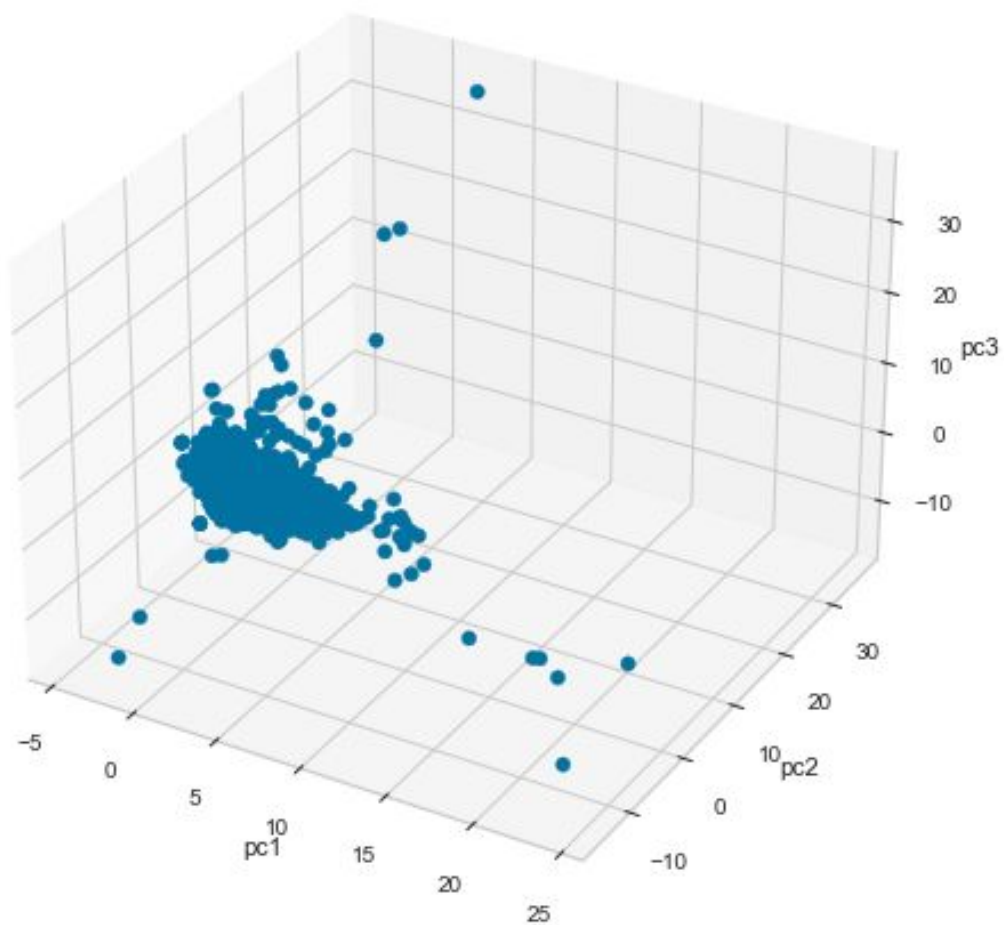
## The process: EDA - Habitable vs Inhabitable

- PCA

PCA plot with the first 2 components (ca 19 % variance)

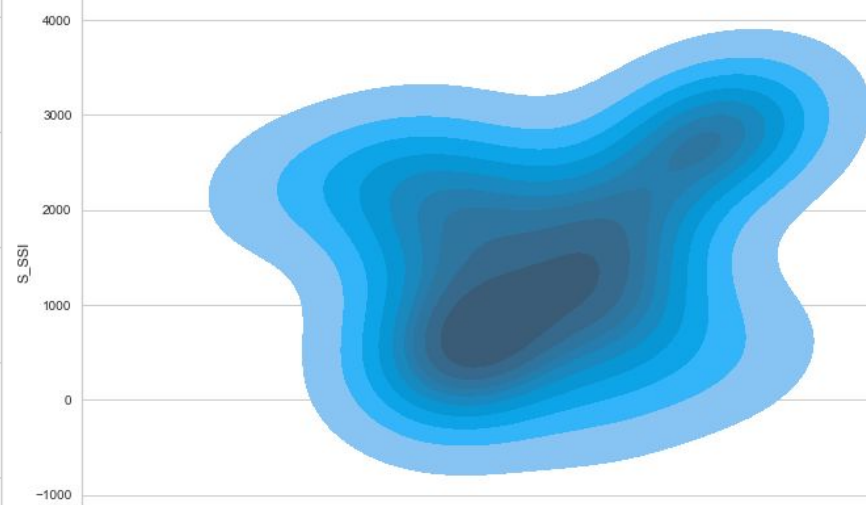
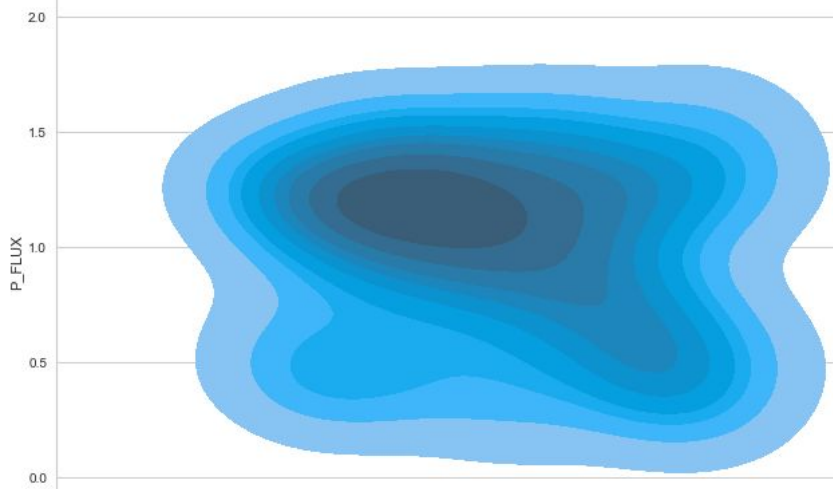
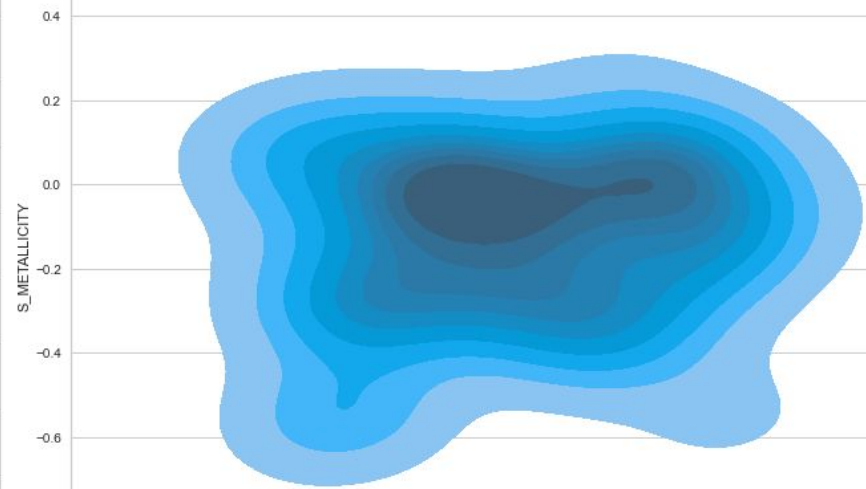
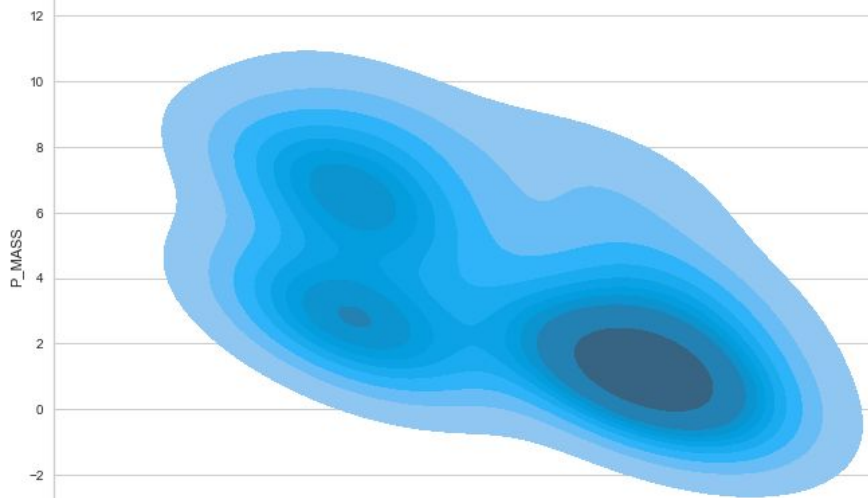


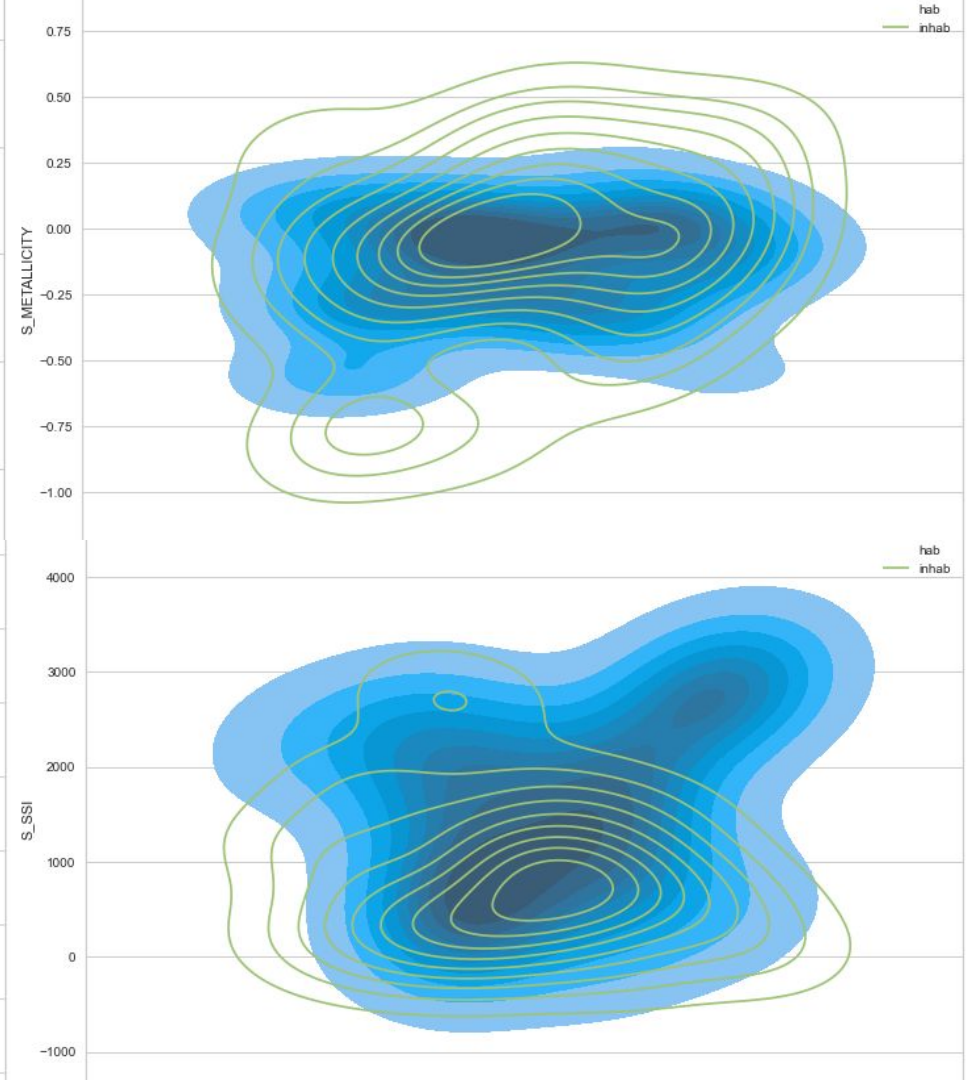
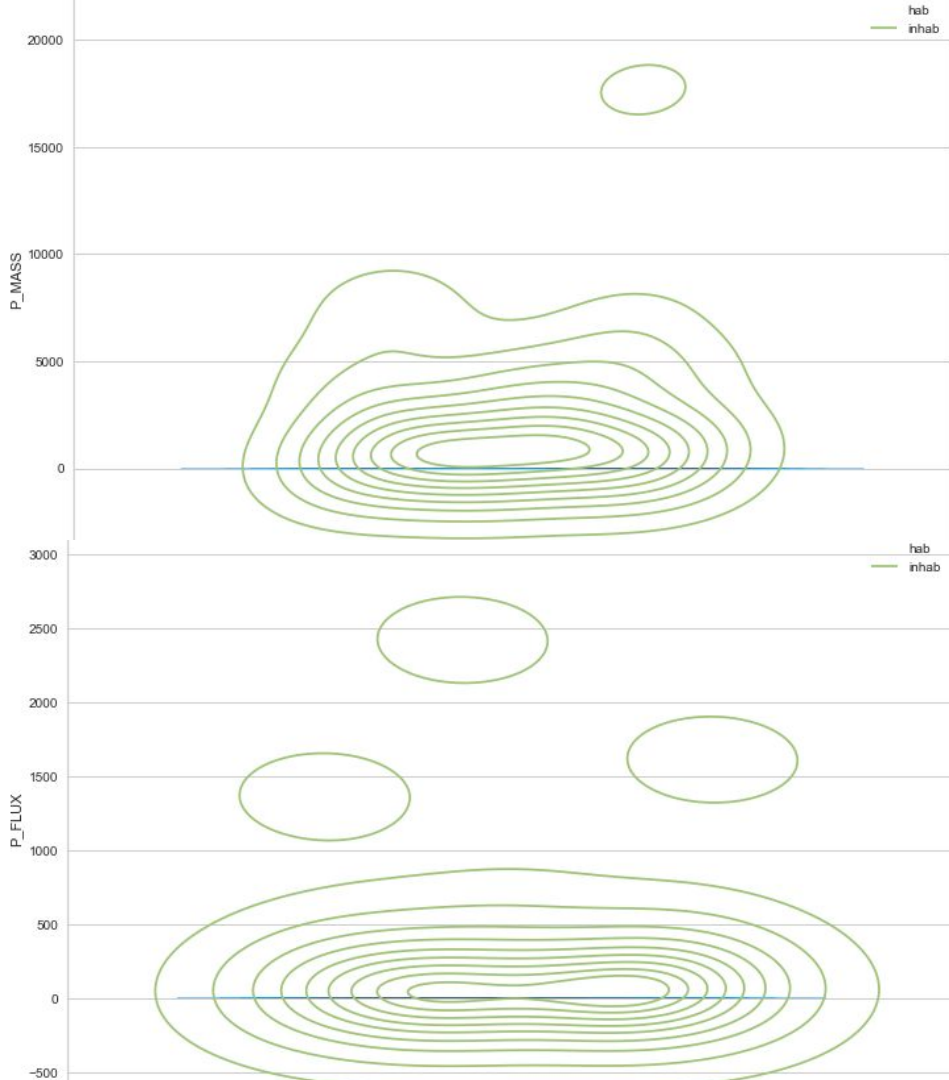
PCA plot with the first 3 components (ca 27 % variance)



## The process: EDA - Habitable vs Inhabitable

- PCA
- What specifically makes the hab/inhab planets differ?





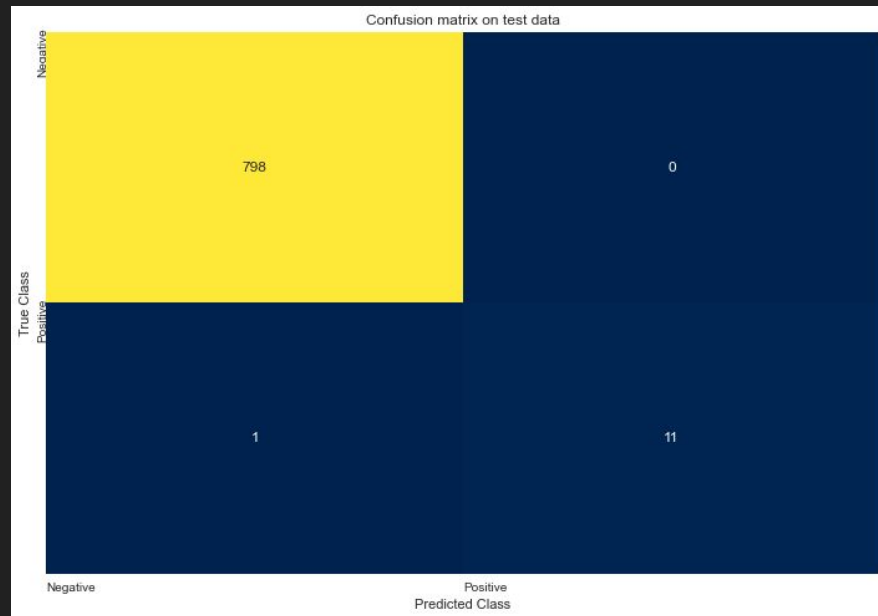
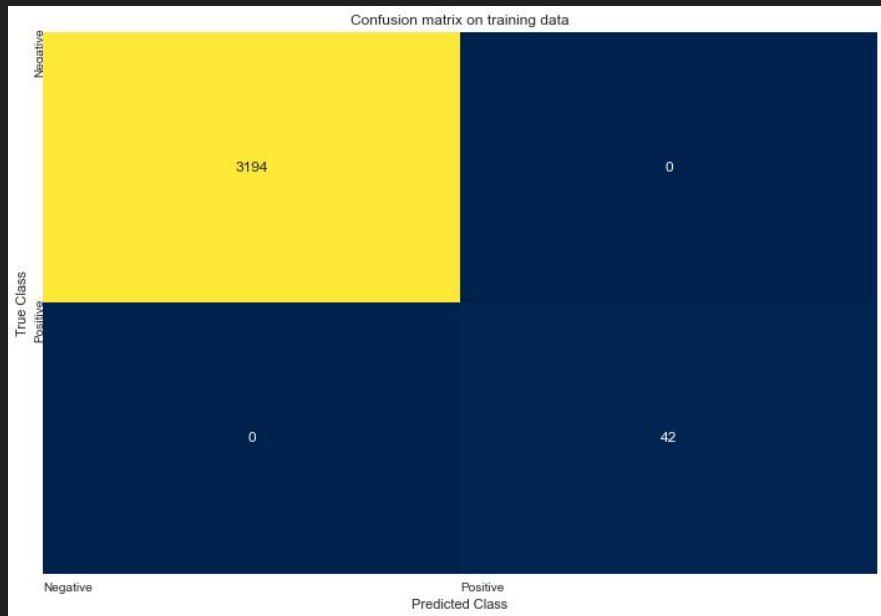
## The process: EDA - Habitable vs Inhabitable

- PCA
- What specifically makes the hab/inhab planets differ?
  - Answer: hab are more earth-like



# ML model: target var = P\_HABITABLE

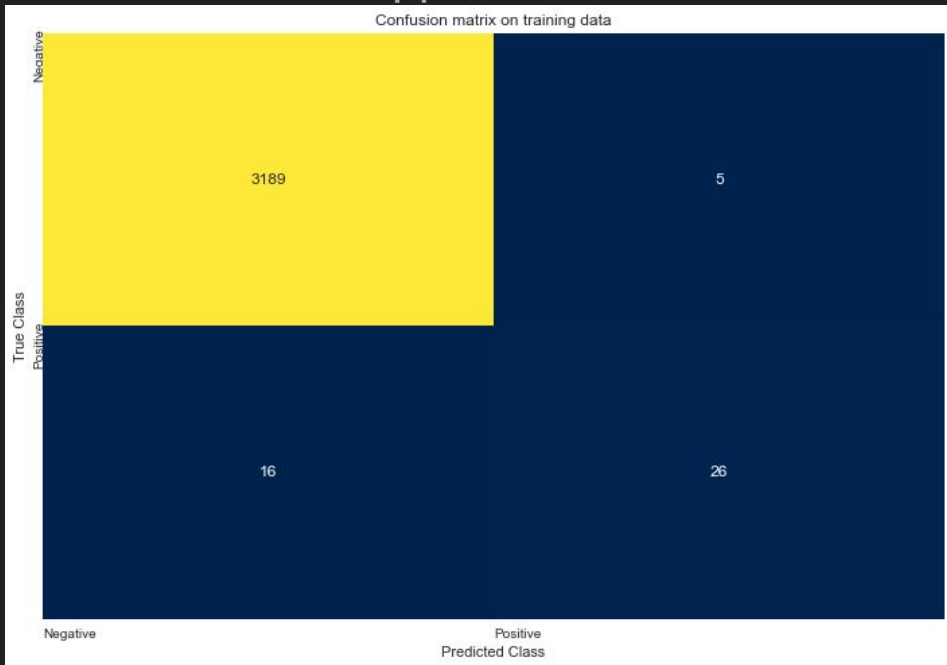
- Estimator: LogisticalRegression
- Scorer: kappa



What is an exoplanet? - Data - The Goal - The Process - **ML** - Summarizing Thoughts

## ML model: target var = P\_HABITABLE (only star data)

- Estimator: RandomForestClassifier
- Scorer: kappa

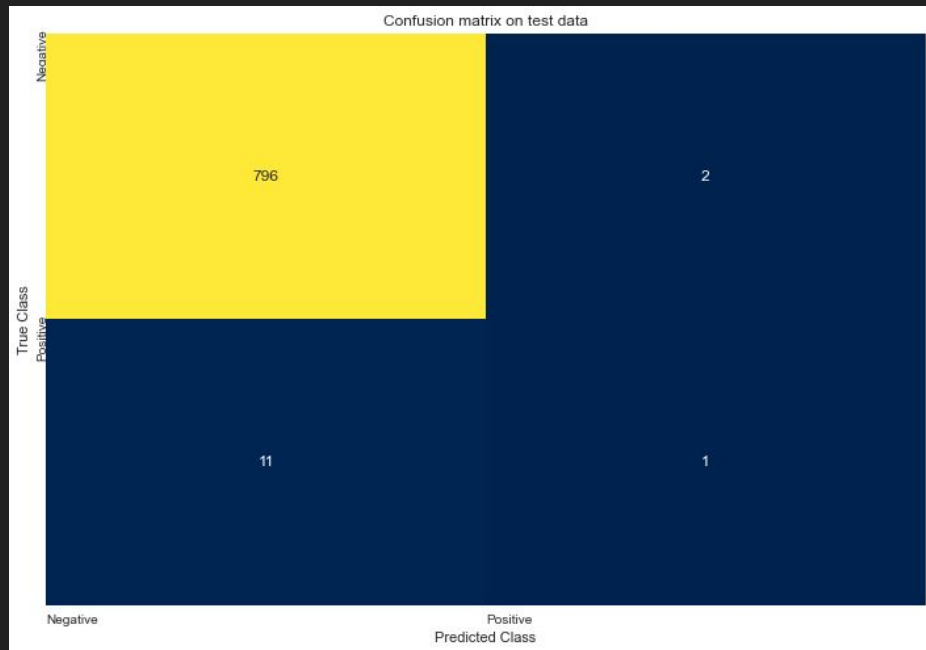


Precision: 84 %

## ML model: target var = P\_HABITABLE (only star data)

- Estimator: RandomForestClassifier
- Scorer: kappa

Precision: 33 %



What is an exoplanet? - Data - The Goal - The Process - ML - **Summarizing Thoughts**

# Summarizing thoughts

## Summarizing thoughts

- A model such as this might not be as inapplicable as one might think

## Summarizing thoughts

- A model such as this might not be as inapplicable as one might think
- Reducing complexity is important

## Summarizing thoughts

- A model such as this might not be as inapplicable as one might think
- Reducing complexity is important
- Problem with feature creation but it was useful

## Summarizing thoughts

- A model such as this might not be as inapplicable as one might think
- Reducing complexity is important
- Problem with feature creation but it was useful
- Try resampling next time



The End