# july 3rd, 2025

christina liu

# stuff i did

- working on poster for astroAI on july 7th (so this will be a short presentation)
- some new SHAP updates (force plots)
- literary studies on more of the astrophysics side of classification instead of just using machine learning
- some ideas on how to approach planetary systems analysis

POSTER

# new SHAP graphs: force plots

SHAP **force** plot for a positive prediction sample (exoplanet)

higher ⇄ lower

base value

f(x)

| –0.29 | –0.09 | 0.11 | 0.31 | 0.51 | 0.71 | 0.91 | **1.00** | 1.11 | 1.31 |

Surface Gravity [log10(cm/s^2)] = 0.9146   st_teff: Stellar Effective Temperature [K] = 0.08841   pl_rade: Planet Radius [Earth Radius] = 0.03308   st_rad: Stellar Radius [Solar Radius] = 0.0009595   pl_orbsmax: Planet Orbit Semi-Major Axis [au] = 0

SHAP **force** plot for a negative prediction sample (exoplanet)

higher ⇄ lower

f(x)

base value

| –0.09 | **0.00** | 0.11 | 0.31 | 0.51 | 0.71 | 0.91 | 1.11 |

st_rad: Stellar Radius [Solar Radius] = 0.008686   st_mass: Stellar Mass [Solar mass] = 0.08974   st_logg: Stellar Surface Gravity [log10(cm/s^2)] = 0.7963   st_teff: Stellar Effective Temperature [K] = 0.18

SHAP **waterfall** plot (Random Forest) for a positive prediction sample (exoplanet)

SHAP **waterfall** plot (Random Forest) for a negative prediction sample (exoplanet)

# exoplanet classification – literary studies

Framework for the architecture of exoplanetary systems (2023)
(DOI: )

**Similar** – similar regardless of distance
**Anti-ordered** – as planets further, get smaller
**Ordered** – as planets further, get larger
**Mixed** – goes back and forth

| Architecture class | Condition |
|---|---|
| Anti-ordered | $C_S(M) < -0.2$ |
| Ordered | $C_S(M) > +0.2$ |
| Similar | $|C_S(M)| \leq 0.2$ and $C_V(M) \leq \dfrac{\sqrt{n-1}}{2}$ (3) |
| Mixed | $|C_S(M)| \leq 0.2$ and $C_V(M) > \dfrac{\sqrt{n-1}}{2}$ |

coefficient of similarity – positive for ordered, negative for anti-ordered

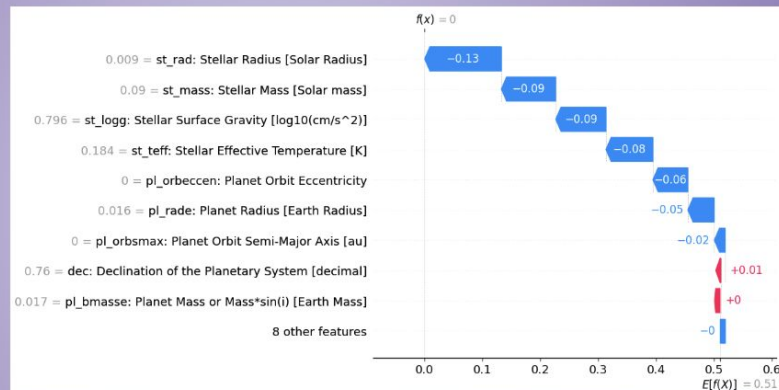$$C_s(q) = \frac{1}{n-1} \sum_{i=1}^{n-1} \left( \log \frac{q_{i+1}}{q_i} \right)$$

where $q_i$ is some planetary quantity $q$ (ex. mass, radius, orbital period, etc.) for the i[th] planet in a system.

coefficient of variation – measure magnitude of variation in a set of numbers

$$C_v(q) = \frac{\sigma(q)}{\overline{q}}$$

"while similar systems will have a low value of the coefficient of variation, mixed systems will have a high value of coefficient of variation"

# exoplanet classification – literary studies

Framework for the architecture of exoplanetary systems (2023)
(DOI: https://doi.org/10.1051/0004-6361/202243751)

**used a model called the** **GENERATION III BERN MODEL** in the process to create synthetic data (under heading *2.1 Theoretical Dataset: Bern Model*)

- system of classification they use requires ≥ 3 planets per system, thus out of their original dataset there were only 41 data points.
- gen iii bern model to generate 1000 such systems

# exoplanet classification – literary studies

Architecture Classification for Extrasolar Planetary Systems (2025)
(DOI: https://doi.org/10.1051/0004-6361/202243751)

- uses 6000 exoplanets (only real data!)
- basically just a straight-up split very similar to earlier ones we talked about

# exoplanet classification – literary studies

Planetary Population Synthesis and the Emergence of Four Classes of Planetary System Architecture (2023)
(DOI: https://doi.org/10.48550/arXiv.2303.00012)

- This paper also uses synthetic data generated using the **GENERATION III BERN MODEL** which seems to be pretty popular.

# next week

- learn more about the astrophysical side of how classifications are created (less ML, lots of literary studies)

that's all for this week. :)