

Classifying exoplanets with Kepler data

Trevor Santiago

A dark blue diagonal gradient bar that starts from the bottom left and extends towards the top right, covering the lower half of the slide.

Data set

Kepler Exoplanet Search Results

10000 exoplanet candidates examined by the Kepler Space Observatory



NASA • updated 3 years ago (Version 2)

Kepler Objects of Interest Data from Kaggle

- Target feature: koi_pdisposition
 - CONFIRMED v. FALSE POSITIVE
- Other interesting features:
 - koi_period: orbital period
 - koi_duration: transit duration
 - koi_kepmag: estimated star magnitude in Kepler band
 - koi_srad: star radius

[Data source](#)

Purpose

Space is cool!

How accurately can we predict whether one of the objects of interest will be confirmed an exoplanet or a false positive?

Which features are most important for making these predictions?

Procedure

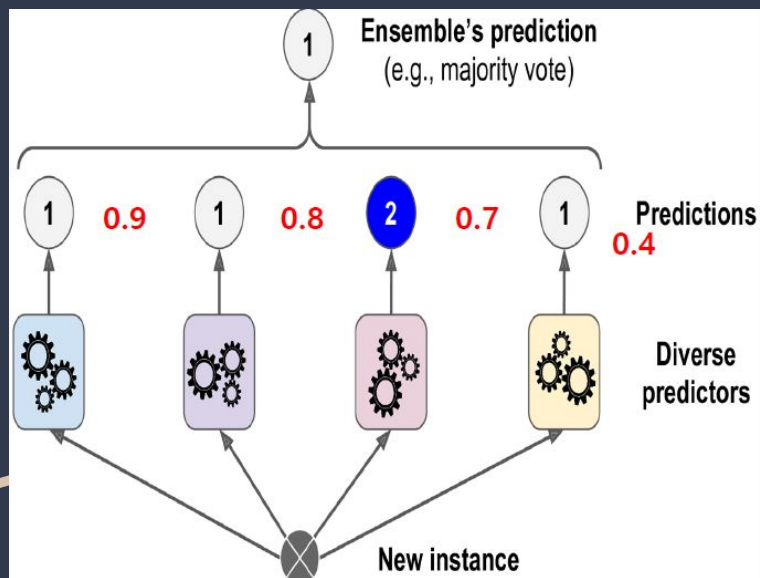
- Split data (3-way)
- Define base pipeline

```
num_prep = Pipeline([
    ('imp', SimpleImputer(strategy='median')),
    ('scale', RobustScaler())
])

# Bin koi_kepmag
binner = ColumnTransformer([('bin', KBinsDiscretizer(), [0])], remainder='passthrough')
cat_prep = Pipeline([
    ('imp', SimpleImputer(strategy='most_frequent')),
    ('bin', binner)
])

prepper = ColumnTransformer([
    ('numeric', num_prep, num_cols),
    ('categorical', cat_prep, cat_cols)
], remainder='drop')
```

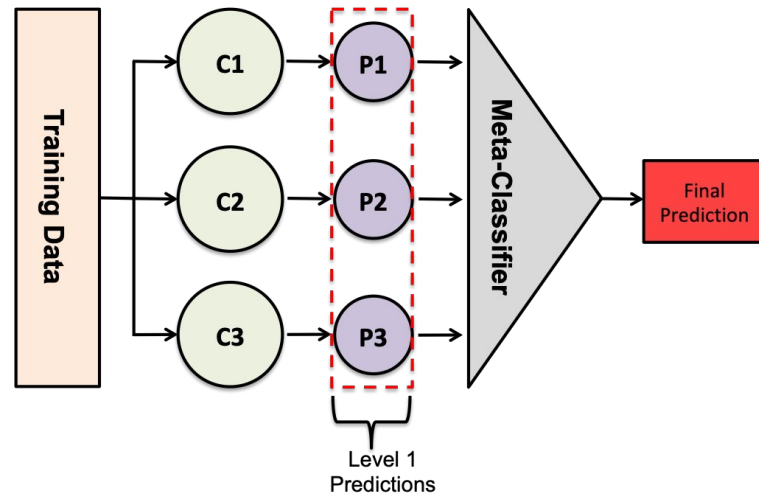
Procedure (Cont.)



- Fit Baseline model (KNN)
- Explored bagging/voting ensembling methods with additional simple models
- Best validation metrics from starter models:
 - Model: VotingClassifier with knn, decision tree, logistic regression
 - Accuracy: 80.8%
 - Log-loss: 0.43038

Procedure (Cont.)

- Explored more robust models like RandomForest and xgboost with automated hyperparameter search
- Final model: StackingClassifier
 - Base learners: same 3 simpler algorithms from voting
 - Meta-learner: Best estimator output from Cross-validated hyperparameter search



Results

- Final model metrics on hold-out test set:
 - Accuracy: 85.9%
 - Log-loss: 0.32342
- Top 3 important features using permutation importance:
 - koi_duration
 - koi_model_snr: normalized transit depth (percent flux blocked by object)
 - koi_period