# Participation in PLAsTiCC
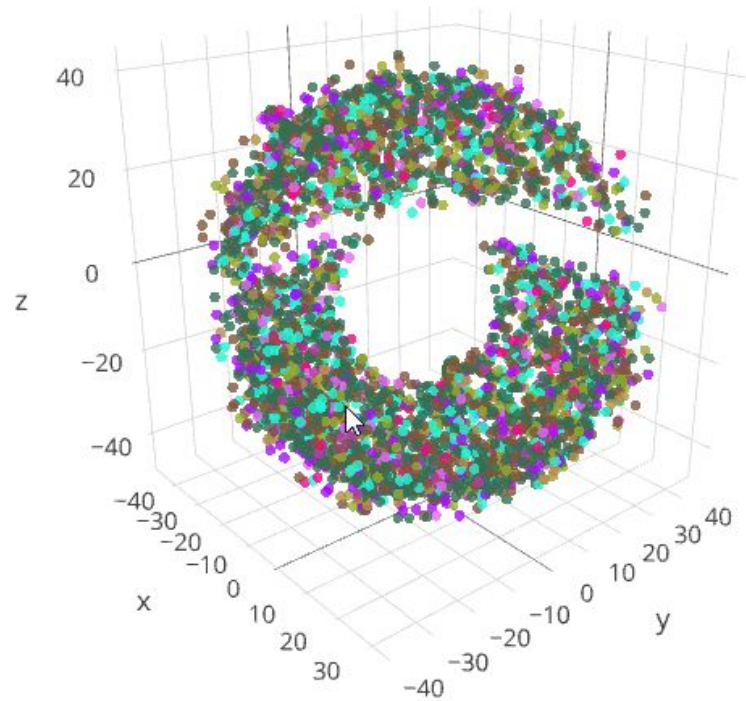
B. Chazelle
Y. Ghernouz
A. Loup
L. Rouquette

# Summary

- Context
- Dataset difficulties
- Tools
- Data transformation
- Learning models

# Context

- PLAsTiCC Contest
  - Real problem contest powered by LSST
  - Main goal : star classifications based on astrophysical data
- School specific project & Scientific synthesis
  - Mix between scientific monitoring and practical work
  - ~0.5 day a week + extra personal work
  - 4 months duration

4

# Dataset difficulties

- Very large test dataset
  - Cannot be loaded on RAM with some technologies (Java / Python)
  - Can lead to high computation time with computation consuming models
- Distinct distributions between the training data and the test data
- Unknown "rest of the world" class
- Missing & noised data on time series
- Some classes are poorly represented

# Tools

## Processing

- Java & Kotlin
  - No Framework
- Python
  - Numpy
  - Pandas
  - Scikit-learn

## Visualization

- Usage of MatPlotLib (Python)
- D3.js (Javascript)
- XChart (JVM environment)

# Data Transformation

Purposes :

- Reduce the impact of different scales (domains)
- Help the computer to handle NaN values
- Remove outliers (when required)
- Transform data into computable values
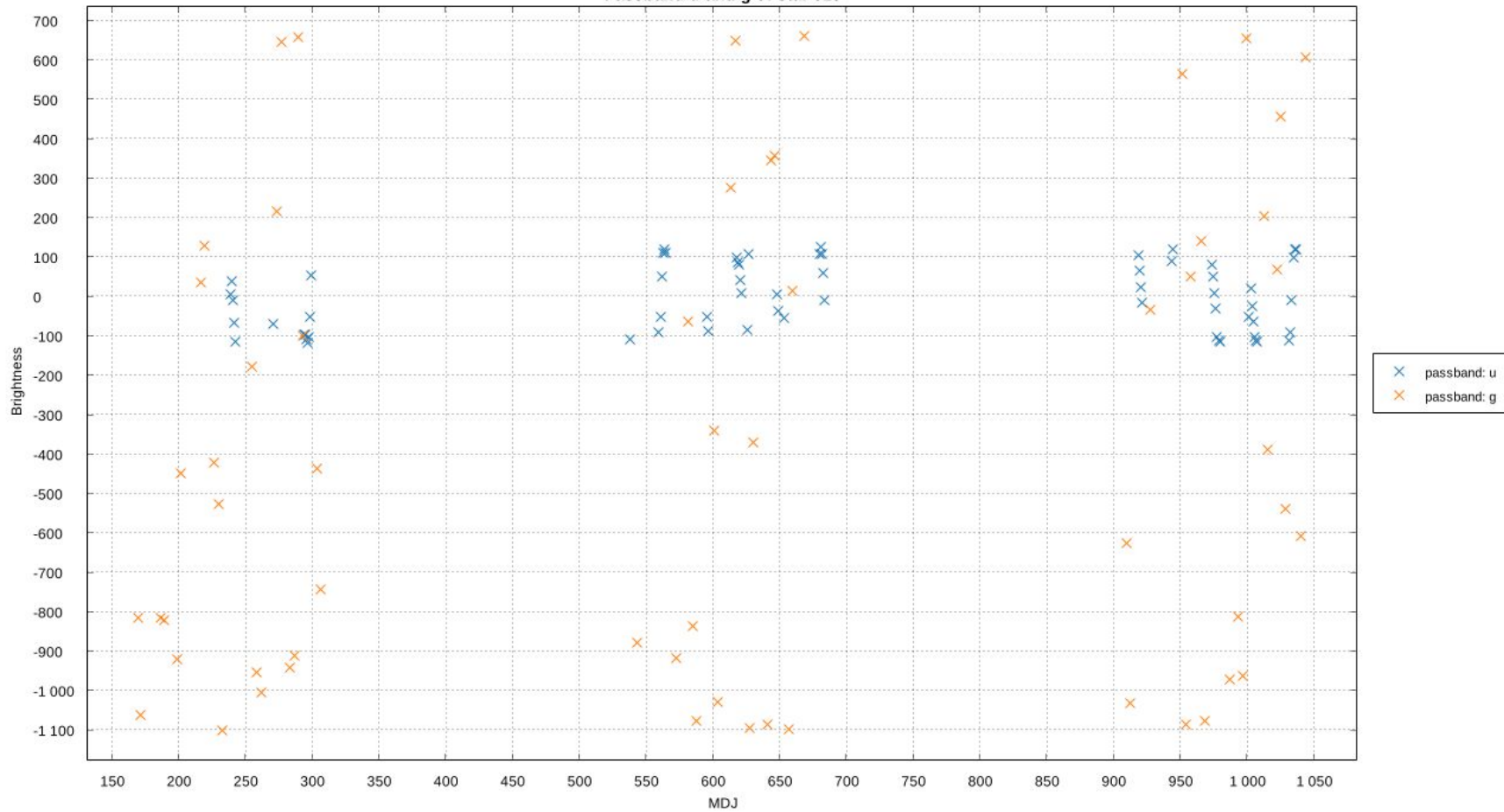
# Data Transformation

We may use different normalizations:

Standard Score                          value = (value - mean) / std

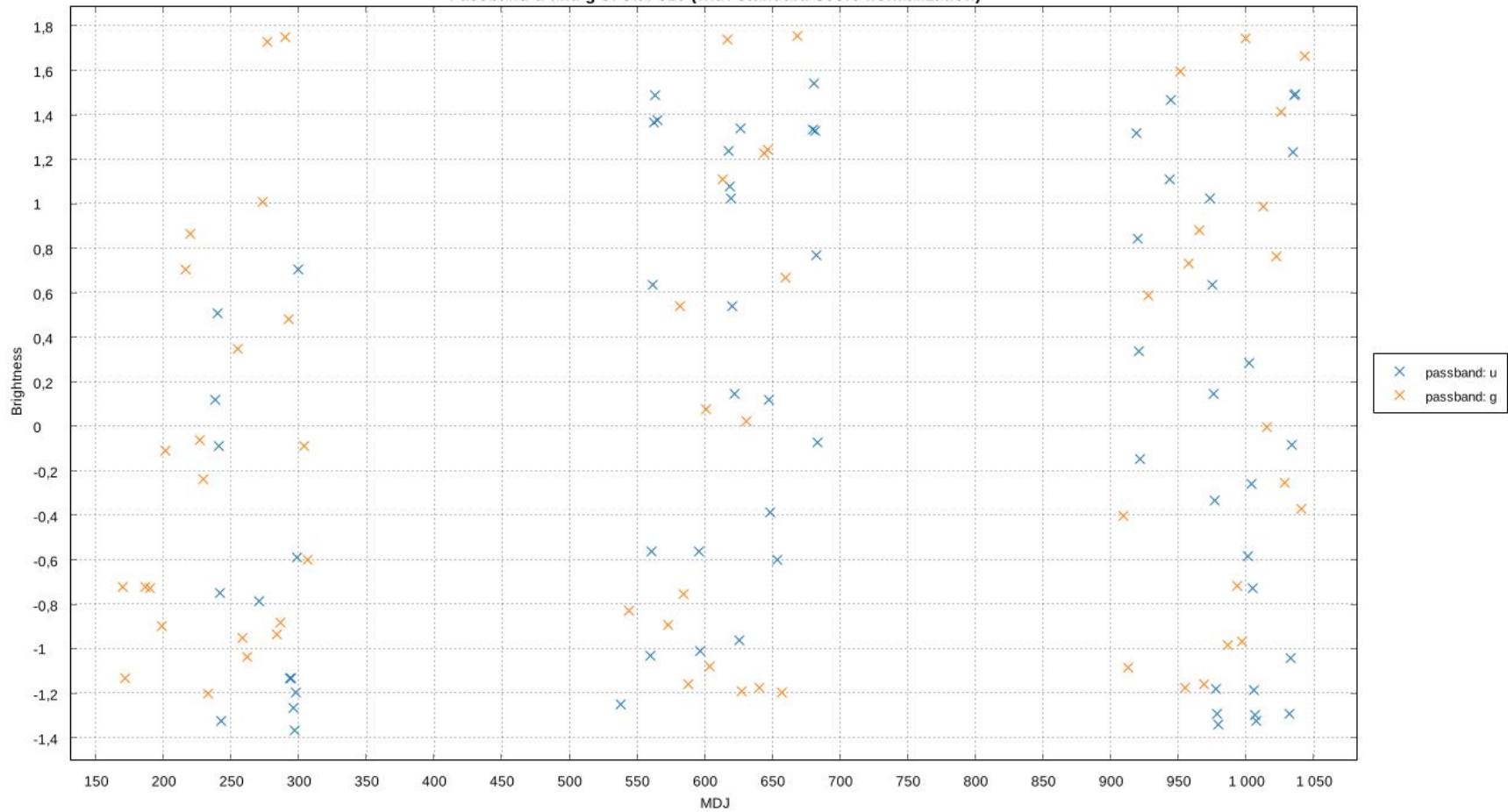Feature scaling                         value = (value - min) / (max - min)

Passband u and g of star 615

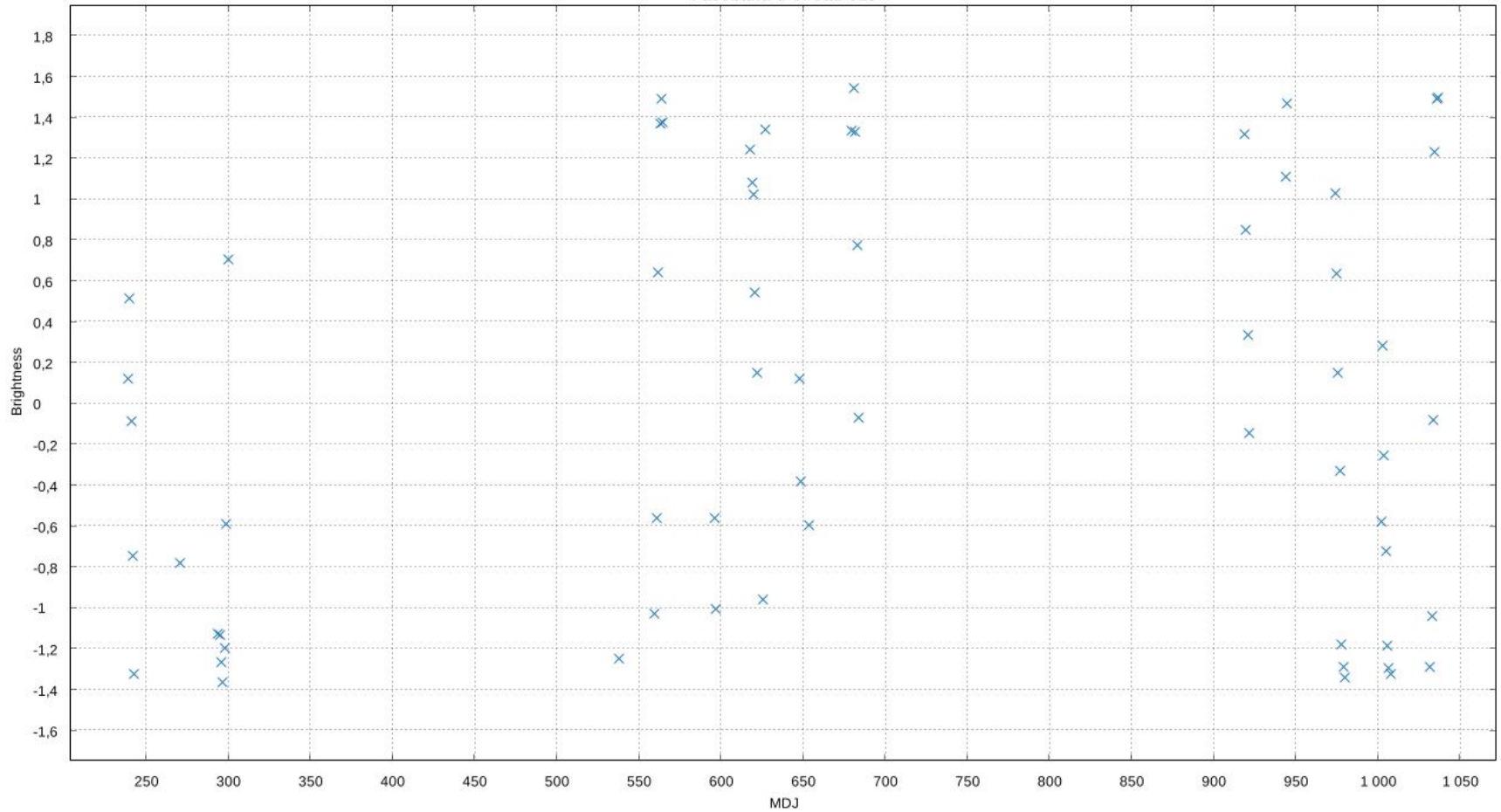Passband u and g of star 615 (with Standard Score normalization)

# Data Transformation
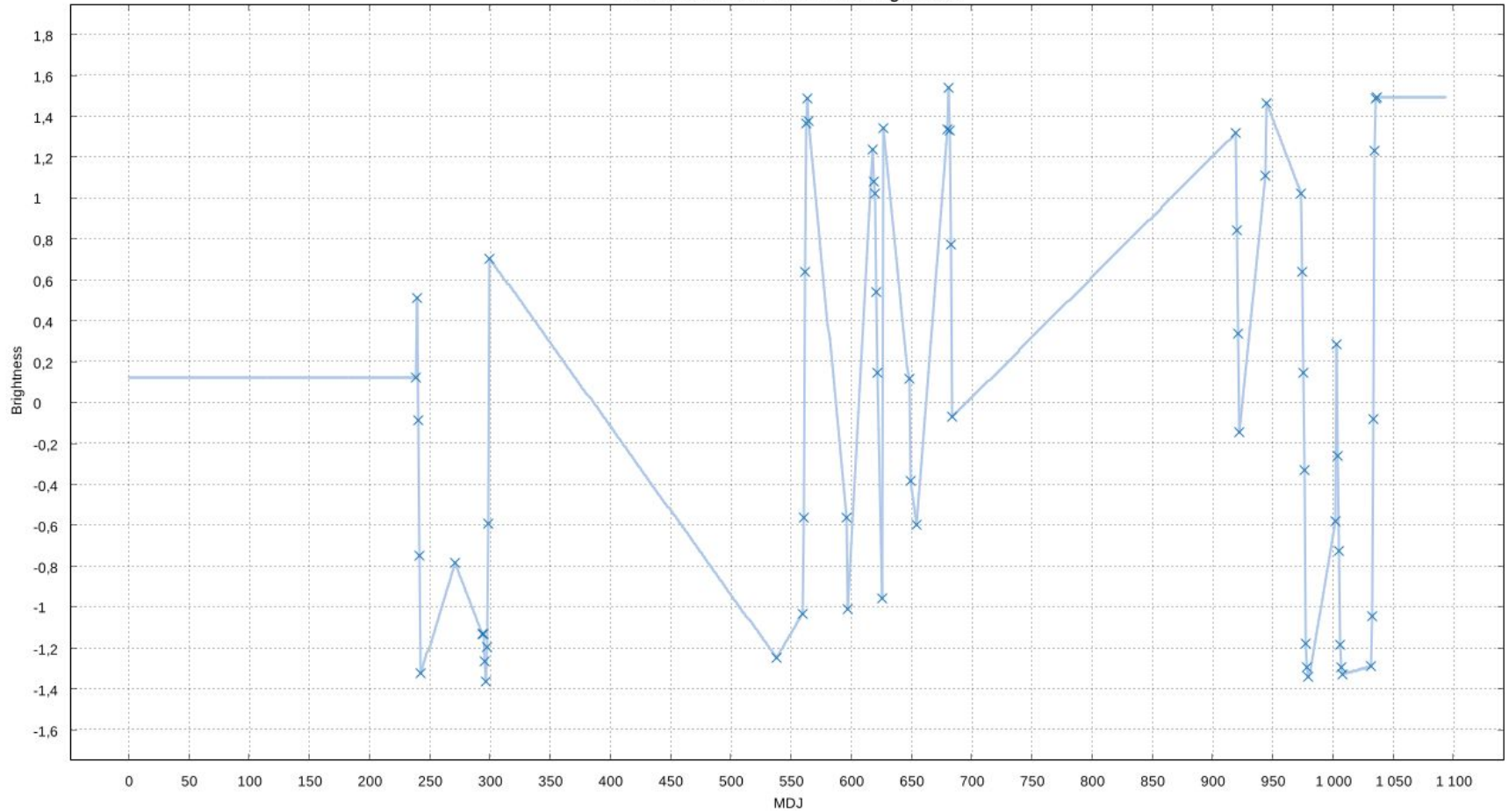
Replacing missing values

- Linear regression
- Cubic spline
- Polynomial regression
- Random point generation [1]
- Normalization

[1] T. A. Hinners, K. Tat, et R. Thorp, « Machine Learning Techniques for Stellar Light Curve Classification »,
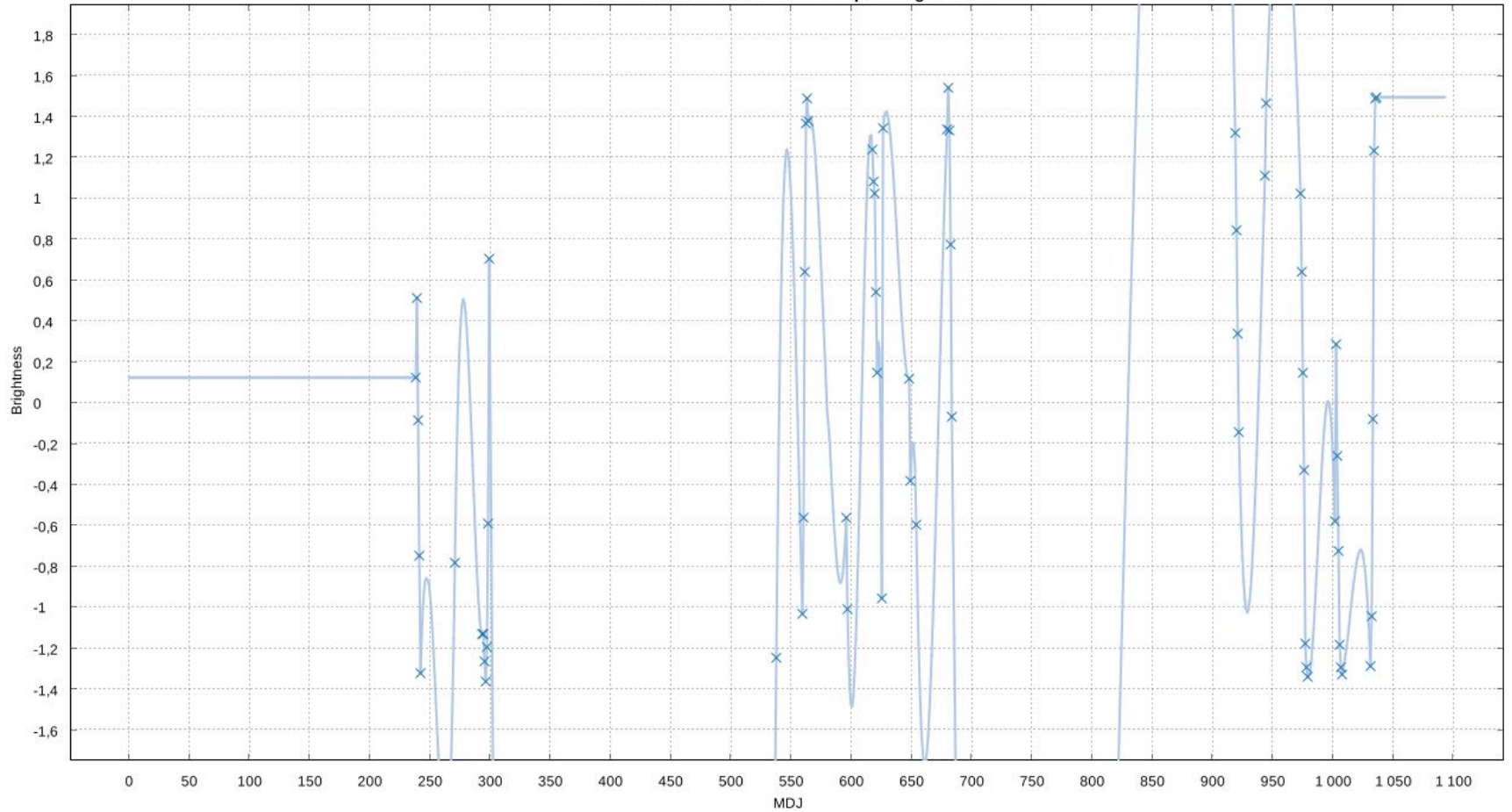*The Astronomical Journal*, vol. 156, nᵒ 1, p. 7, juin 2018. Disponible à https://arxiv.org/pdf/1710.06804.pdf
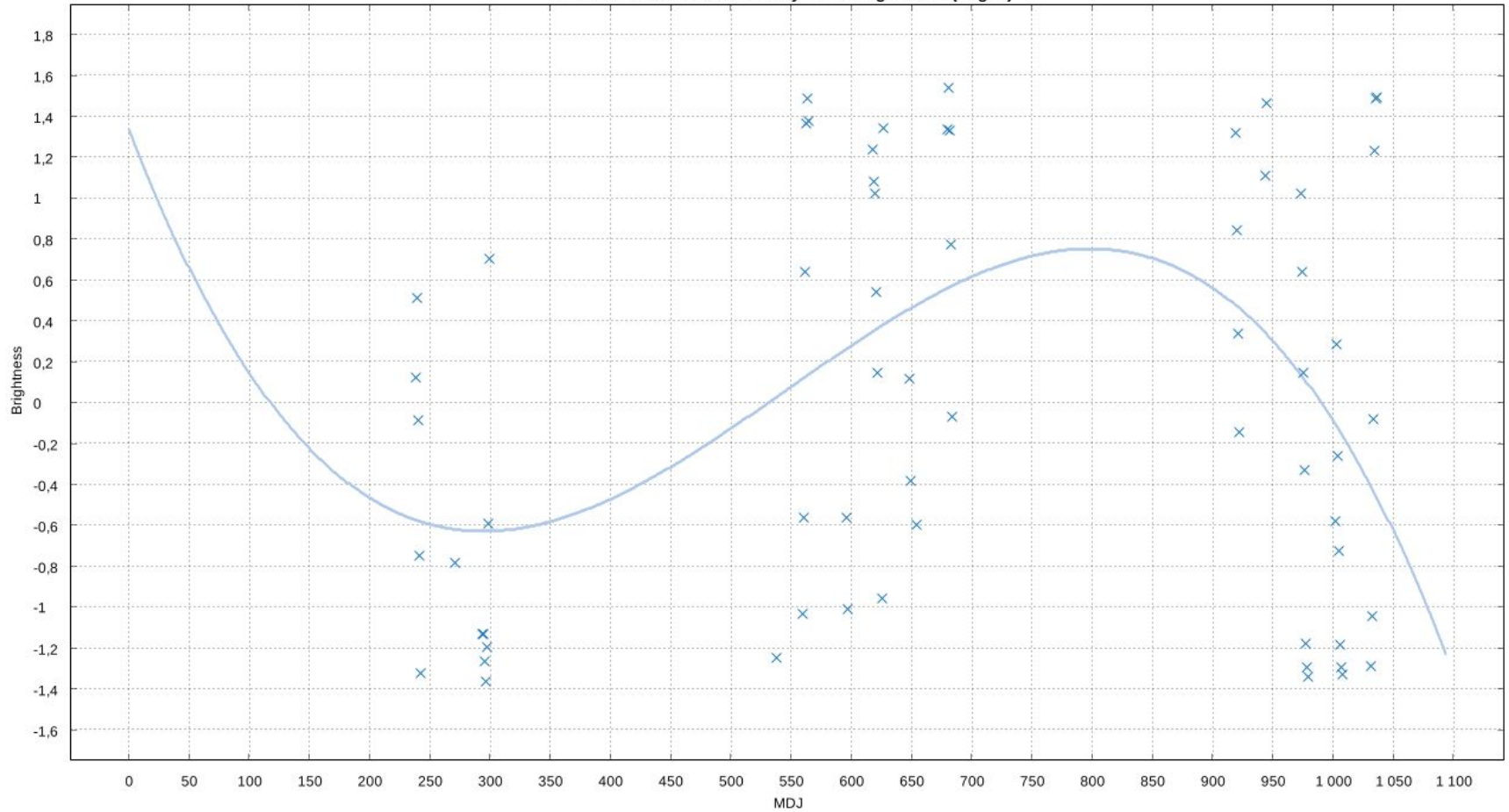
Passband u of star 615

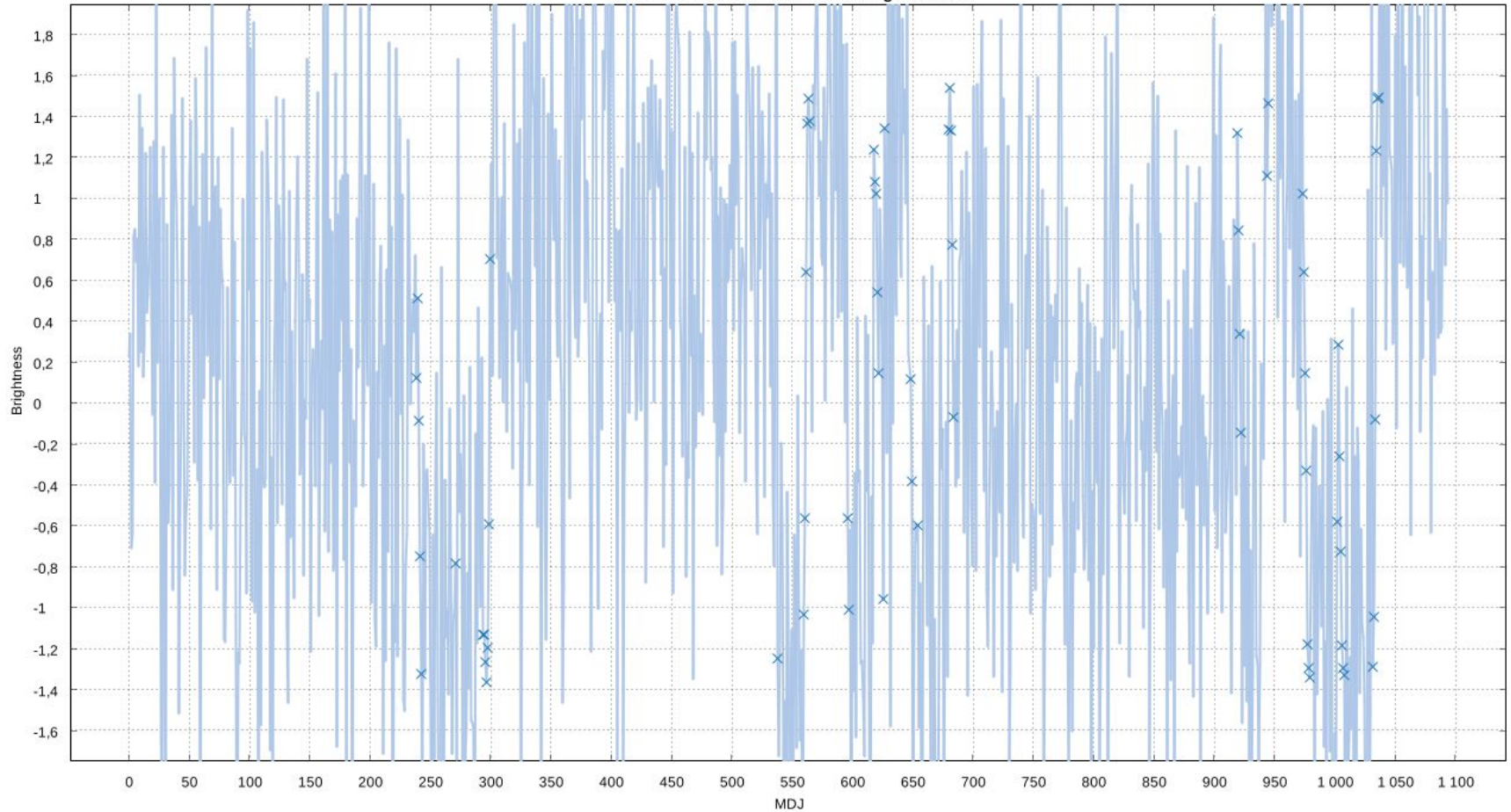Passband u of star 615 - Linear regression

Passband u of star 615 - Cubic spline regression

Passband u of star 615 - Polynomial regression (deg=3)
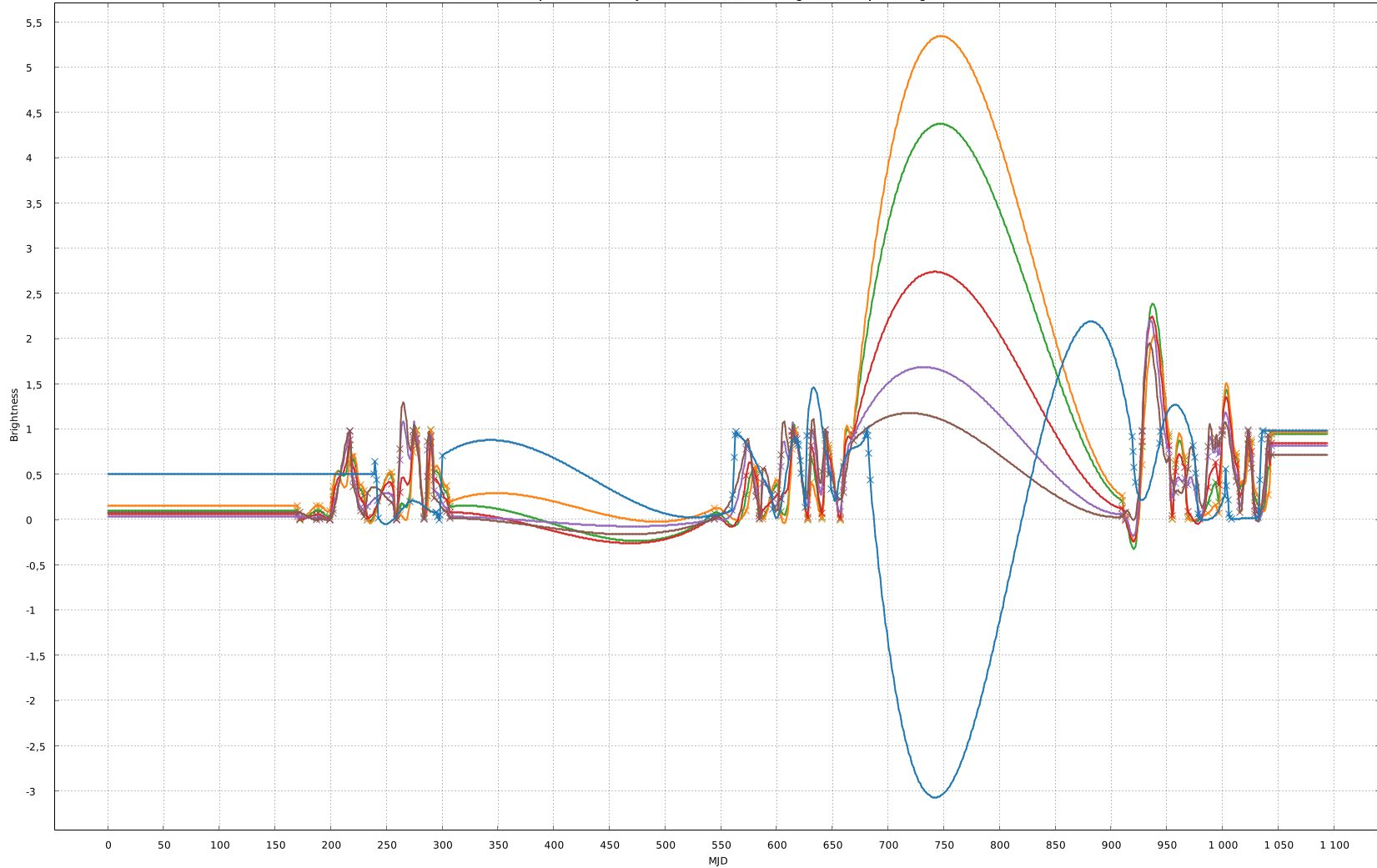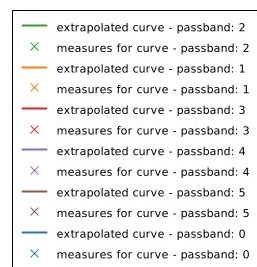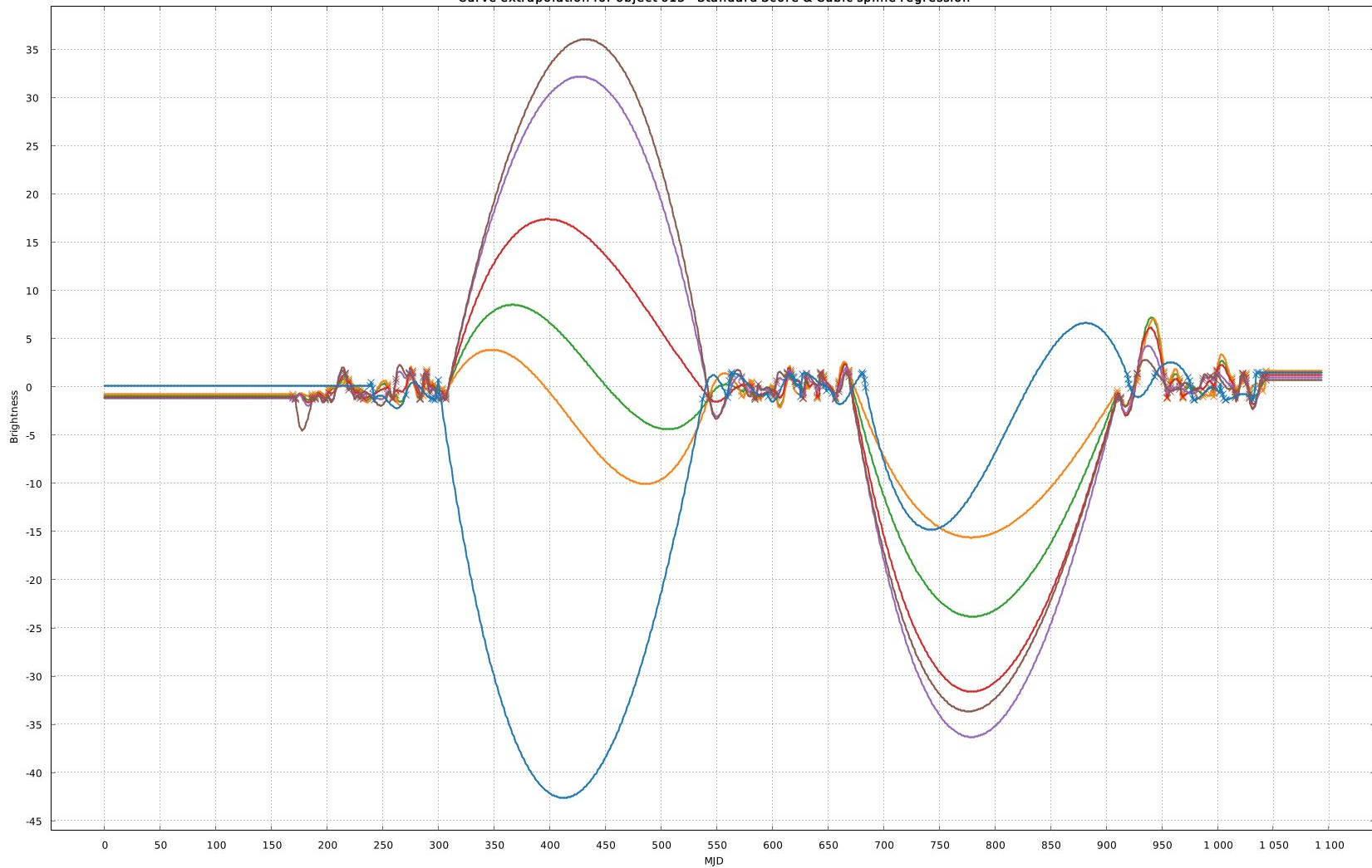
Passband u of star 615 - Random generation

Passband u of star 615 - Random generation

Curve extrapolation for object 615 - Feature scaling & Cubic spline regression

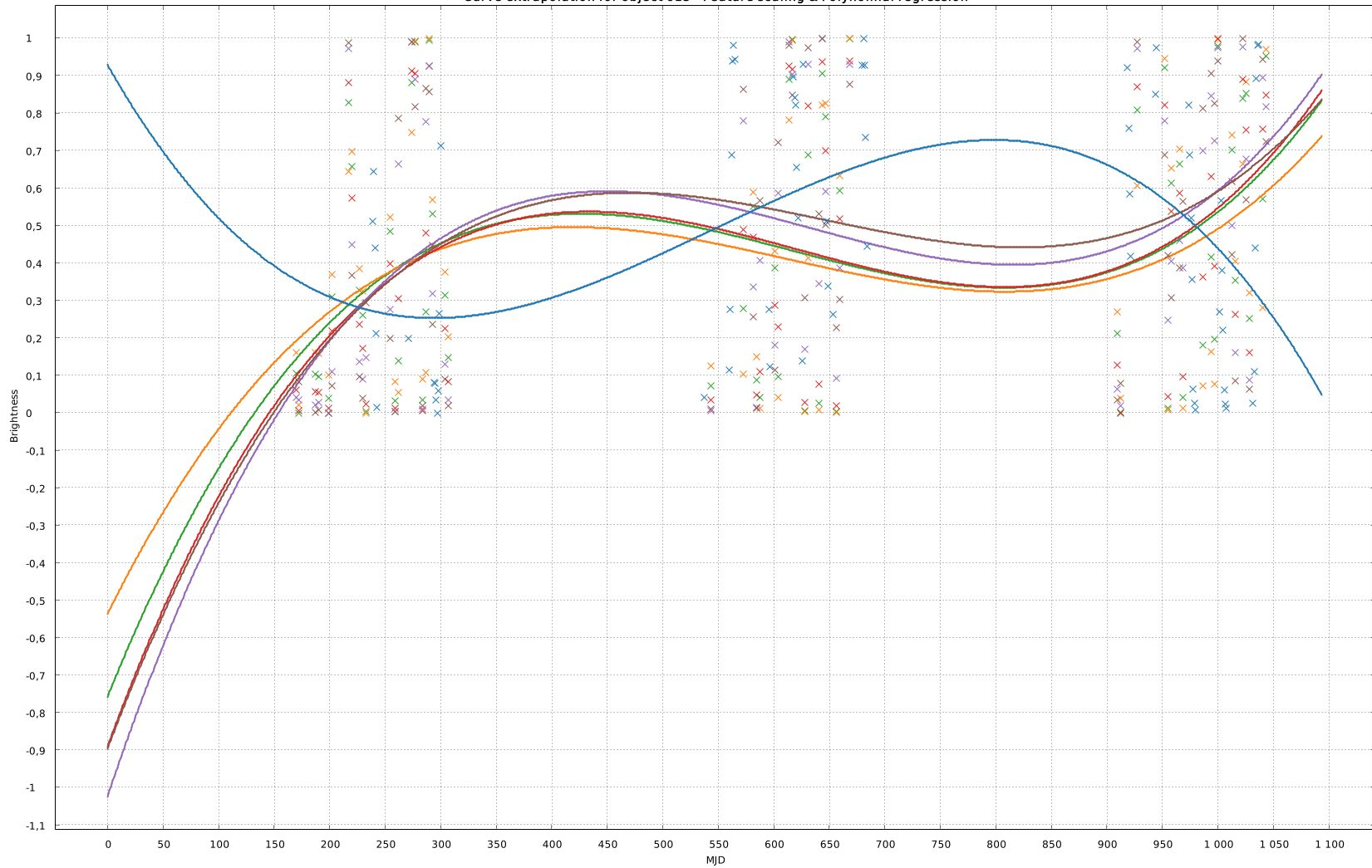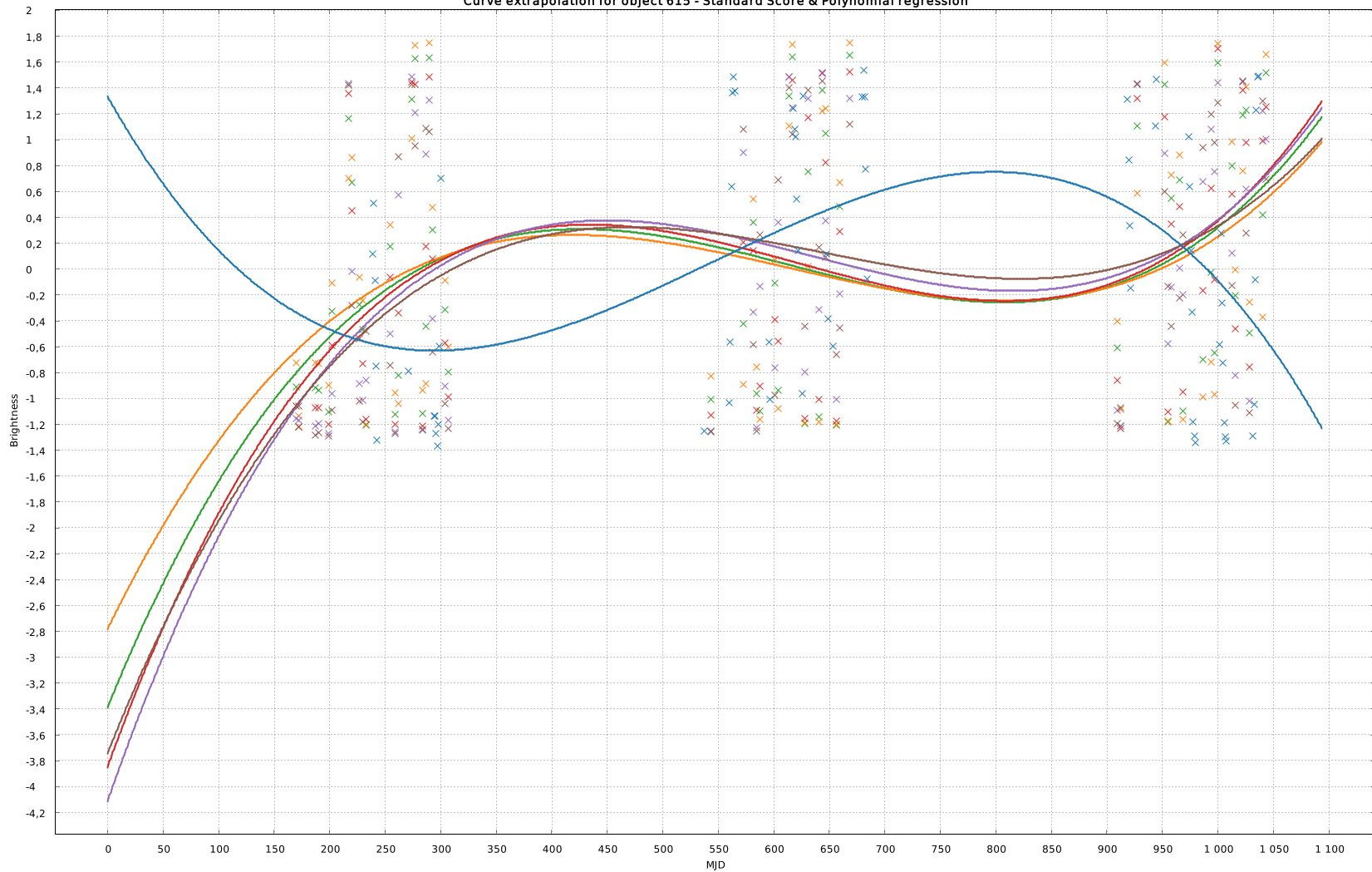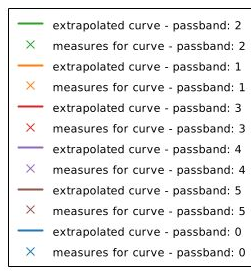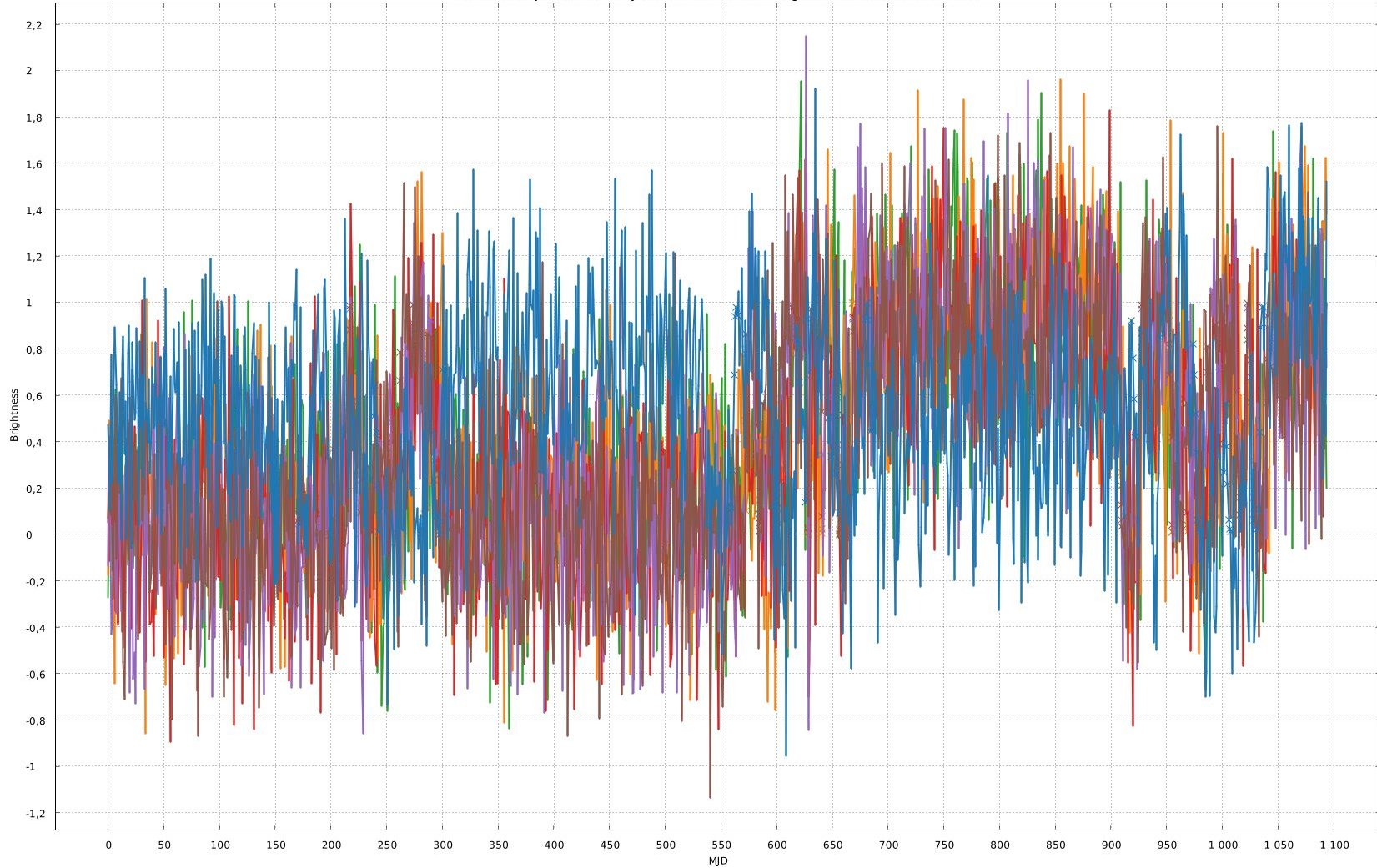Curve extrapolation for object 615 - Standard Score & Cubic spline regression

Curve extrapolation for object 615 - Feature scaling & Polynomial regression

Curve extrapolation for object 615 - Standard Score & Polynomial regression

Curve extrapolation for object 615 - Feature scaling & Normal distribution filler

Curve extrapolation for object 615 - Standard Score & Normal distribution filler

# Learning models

Deep learning

| Pros | Cons |
|---|---|
| ● May integrate some noise in the model<br>● Do not require important data transformation and analysis | ● Cannot handle properly the "rest of the world" class<br>● Sensible to class distribution<br>● No explainable<br>● Need very large sets |

# Learning models

Decision Tree

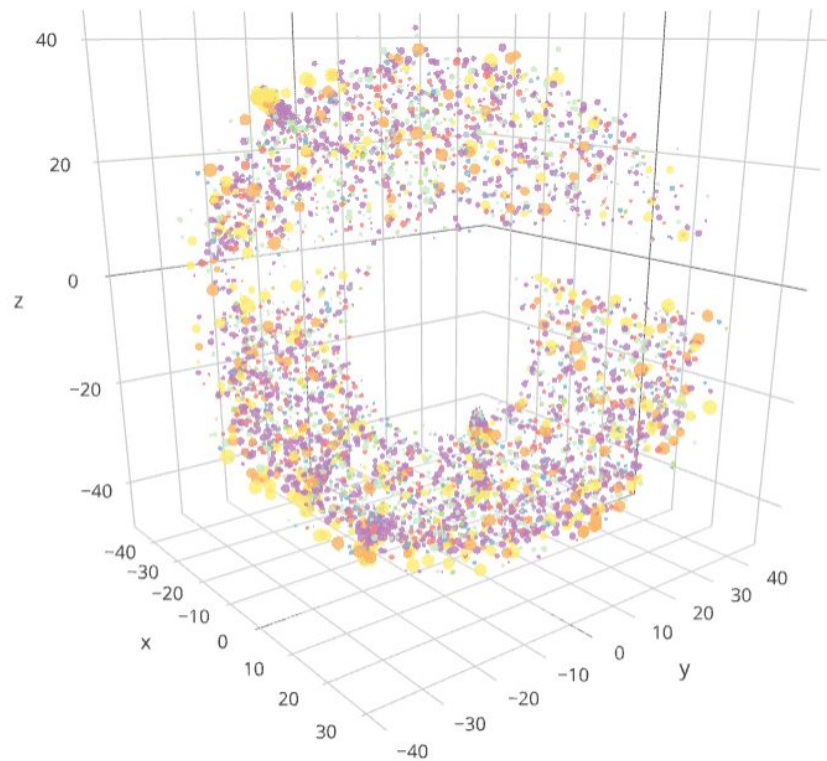| Pros | Cons |
|------|------|
| <ul><li>Can handle missing data (axis are not required to be continued)</li><li>Can manage nominal properties</li><li>Explainable if the size is not big</li></ul> | <ul><li>Overfit quickly</li><li>Cannot manage an unknown class</li></ul> |

# Learning models

KNN + Features engineering + cos similarity



DATA 1
id = 615, ra = 349.04, decl = -61.94, gal_l = 320.79, gal_b = -51.75, ddf = 1, ... , target = 92

Properties → Features

DDF
ra
decl

mwebv

is ddf == true ?
is distmod NaN ?
is distmod < 34 ?
is 34 < distmod <= 38 ?

is 34 < distmod <= 38 ?

Features vector

| 1 |
| 0 |
| 0 |
| 1 |

Target 92

Class Vector

| 5 |
| 2 |
| 0 |
| 3 |

Sum of features vector of each object of the class

# Learning models

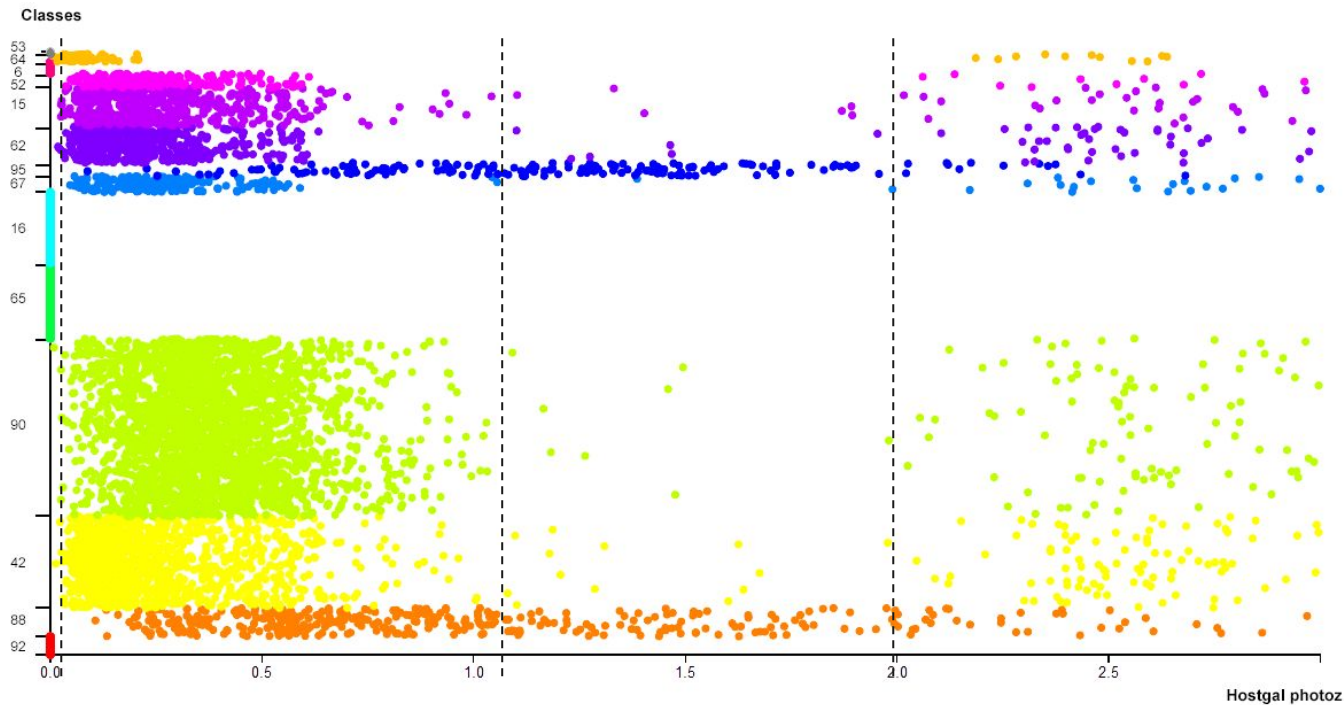KNN + Features engineering + cos similarity

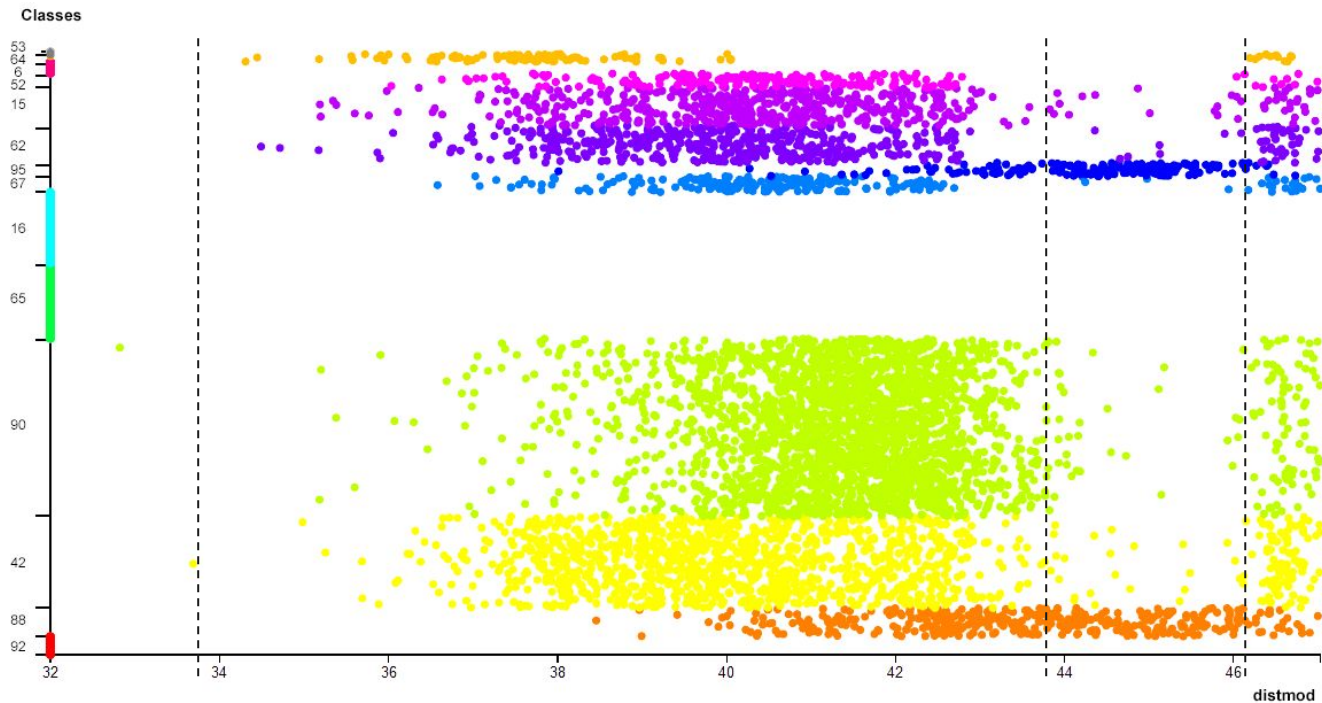| Pros | Cons |
|---|---|
| ● Easy to create<br>● Can manage the unknown class | ● Required to extract meaningful features (very complex task) |

# Features engineering

# Features engineering

# Features engineering

# Our results on PLAsTiCC

| Submission and Description | Public Score | Use for Final Score |
|---|---|---|
| **knn.zip**<br>5 days ago by Loïc Rouquette<br><br>Knn based solver | 24.793 | ☐ |
| **submission_vector.zip**<br>7 days ago by Loïc Rouquette<br><br>Test cos similarity between features vectors | 22.077 | ☐ |
| **submission.zip**<br>7 days ago by Loïc Rouquette<br><br>DecisionTree + custom features over flux (p90, p10) and some informations about metadata (specz, photoz, etc). | 30.947 | ☐ |
| **submission.zip**<br>11 days ago by Loïc Rouquette<br><br>Base line classifier (simple DecisionTreeClassifier based on extracted Features) | 31.754 | ☐ |

# Thanks for your attention