

Computer  
Science  
Center



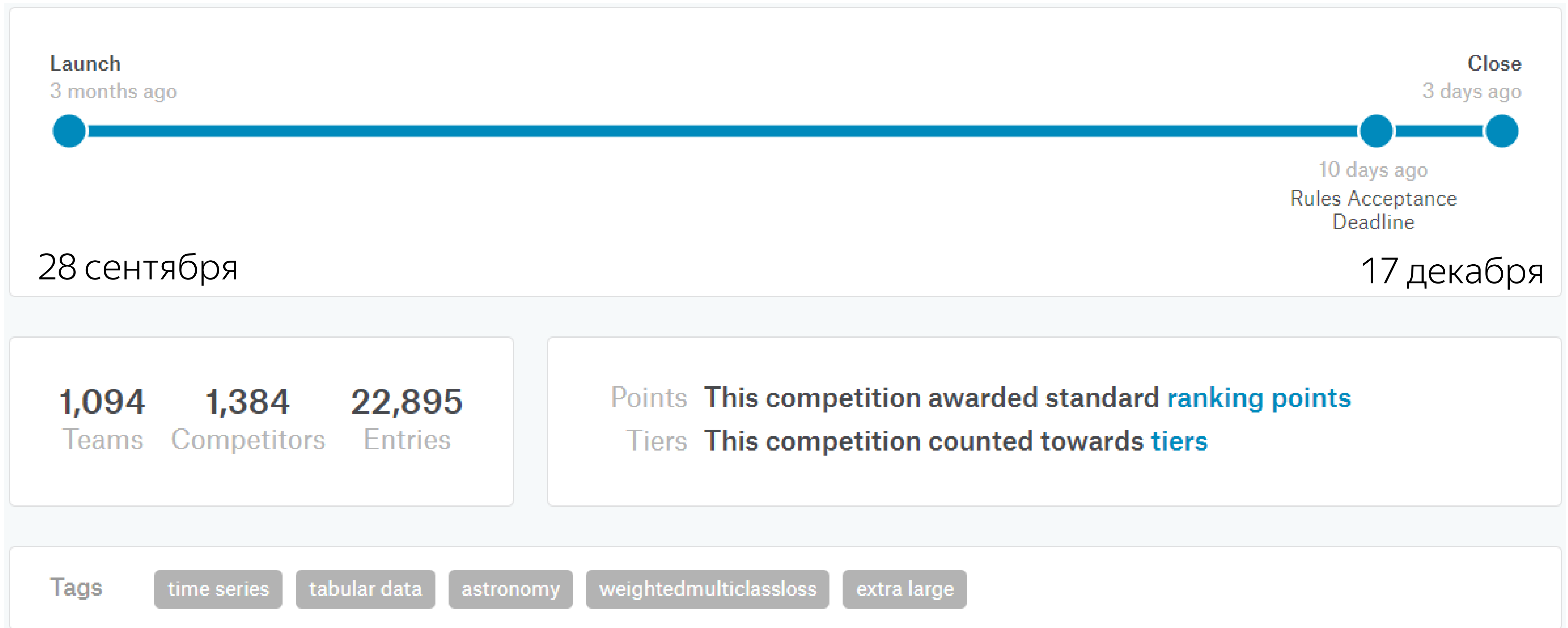
# PLAsTiCC-2018 at Kaggle

Александр Авдюшенко

24 декабря 2018

# Задача

Классификация астрономических объектов на 15 классов, одного из которых нет в train (class\_99 = неизвестные науке объекты)



# Dataset

Количество строк  
7.8K – training\_set\_metadata  
1.4M – training\_set

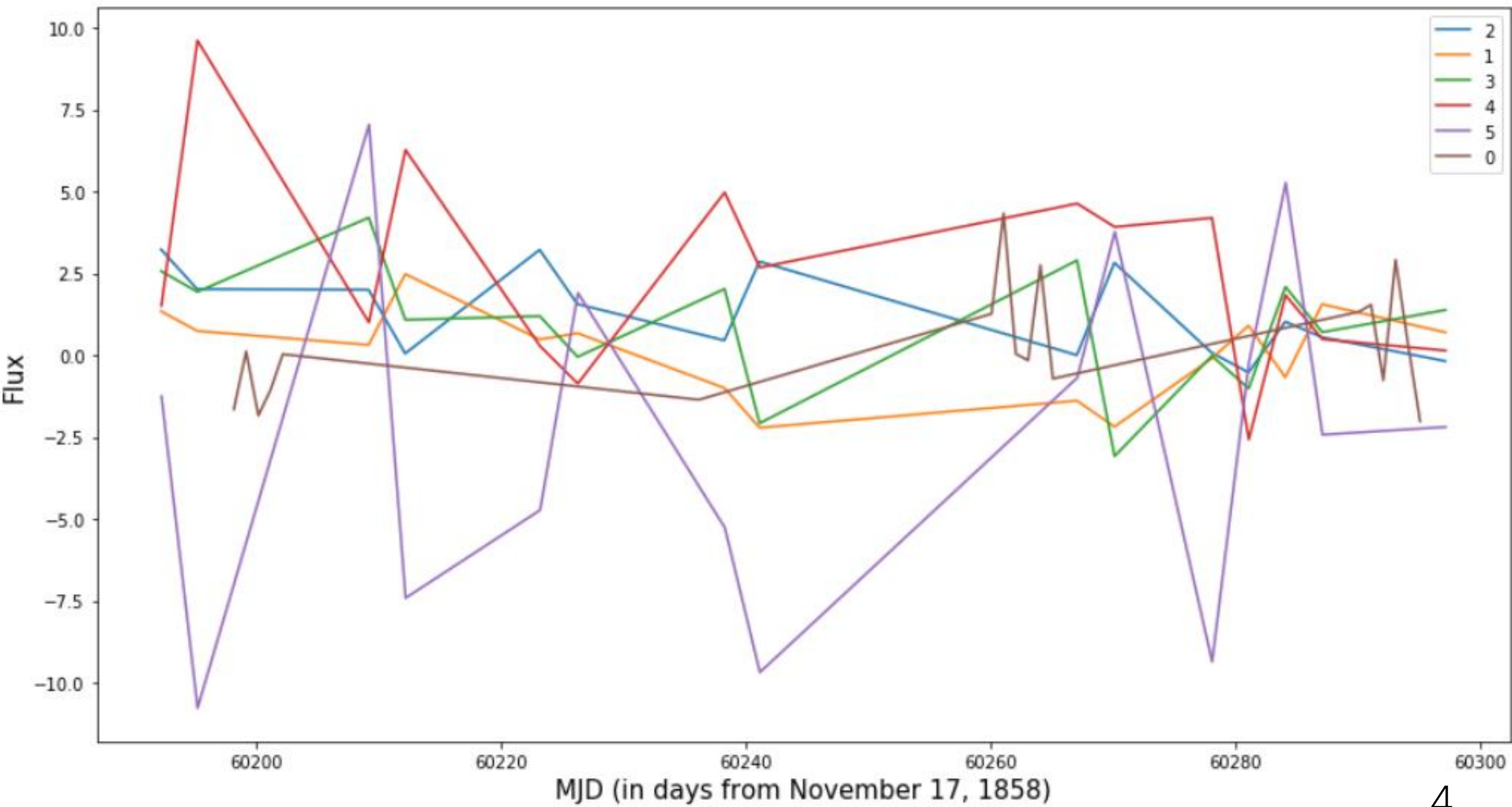
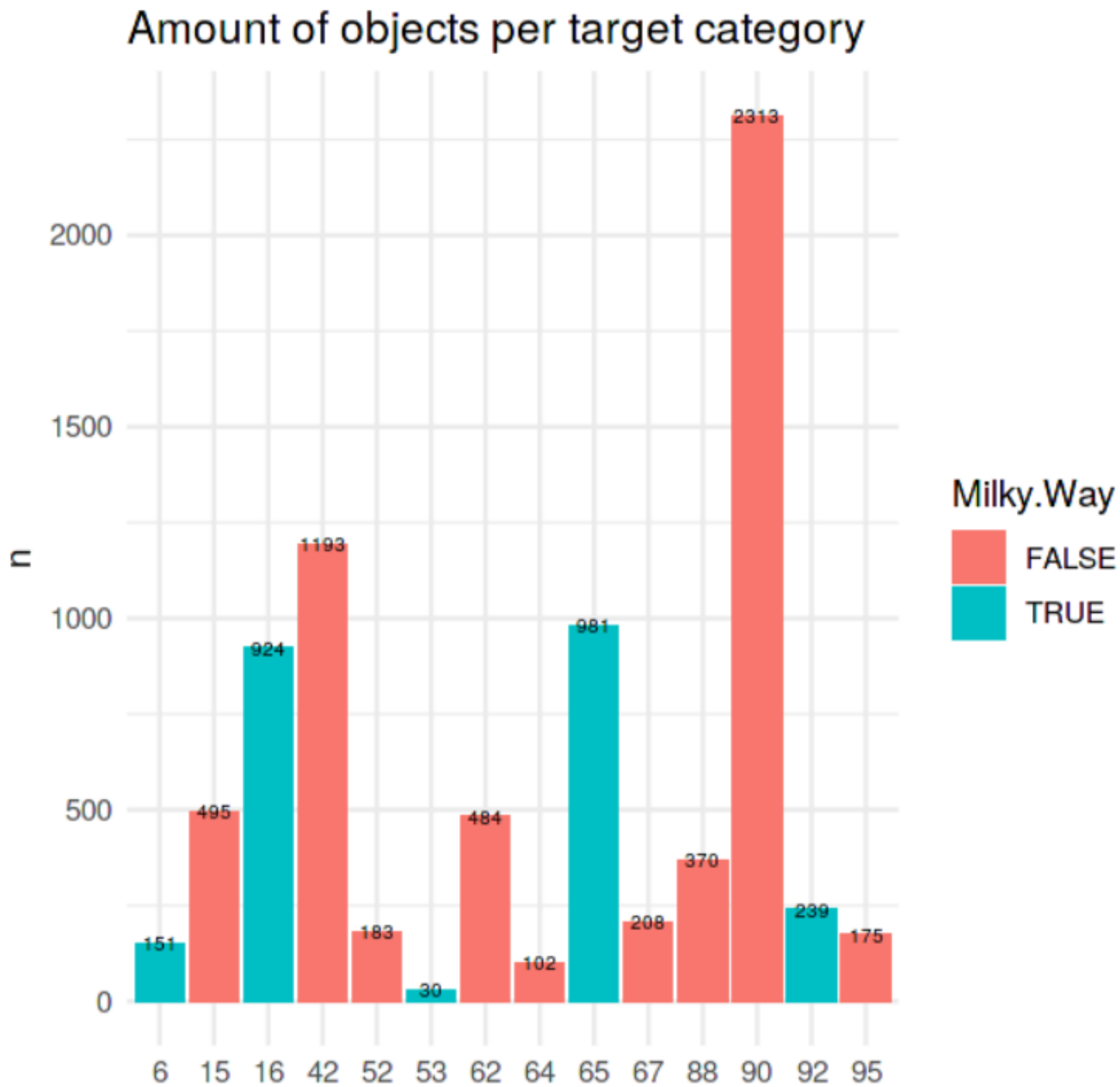
3.5M – test\_set\_metadata  
454M – test\_set

Обучение 20-30 м

Применение 3-4 ч

Metadata  
object\_id, координаты,  
redshift, distmod, **target**  
Data

unix\_time, flux (brightness) x6  
flux\_err, detected




# #169 solution – @ovalur

54

▲ 49

Aleksandr Avdyushenko




1.755

10

36m

Your Best Entry ↑

Your submission scored 1.755, which is an improvement of your previous score of 2.158. Great job!

 Tweet this!

	private	public	
<b>xgb_lgb_2_blend.csv</b> a month ago by Aleksandr Avdyushenko xgb_lgb_2_blend	1.12522	1.09950	<input type="checkbox"/>
<b>single_subm_0.656146_2018-11-21-09-33.csv</b> a month ago by Aleksandr Avdyushenko lgb_2_sum_ideas	1.10791	1.08183	<input type="checkbox"/>
<b>single_subm_0.904055_2018-11-20-05-55.csv</b> a month ago by Aleksandr Avdyushenko lgb_sum_ideas	1.41797	1.38728	<input type="checkbox"/>
<b>single_subm_0.684427_2018-11-20-04-15.csv</b> a month ago by Aleksandr Avdyushenko xgb_sum_ideas	1.15567	1.13203	<input type="checkbox"/>
<b>single_predictions.csv</b> a month ago by Aleksandr Avdyushenko lgb_with_time_series	2.06563	1.98198	<input type="checkbox"/>
<b>single_predictions.csv</b> 2 months ago by Aleksandr Avdyushenko lgb_full	1.79270	1.75516	<input type="checkbox"/>

	private	public	
<b>single_subm_pred_extragalactic_gen_unk.csv</b> 14 days ago by Aleksandr Avdyushenko pred_1_037_and_extraGal_genUnk	1.03363	1.02195	<input type="checkbox"/>
<b>single_subm_pred_extragalactic.csv</b> 14 days ago by Aleksandr Avdyushenko pred_1_037_and_extragal	1.04733	1.03429	<input type="checkbox"/>
<b>single_subm_lgb_genUnknown_extragalactic.csv</b> 14 days ago by Aleksandr Avdyushenko lgb_genUnknown_extragalactic	1.09448	1.07006	<input type="checkbox"/>
<b>single_subm_lgb_and_gen_unknown.csv</b> 14 days ago by Aleksandr Avdyushenko lgb_and_gen_unknown	1.10202	1.07725	<input type="checkbox"/>
<b>single_subm_0.652208_2018-12-09-05-14.csv</b> 14 days ago by Aleksandr Avdyushenko lgb_3	1.10568	1.08044	<input type="checkbox"/>

Submission and Description	Private Score	Public Score	Use for Final Score
<b>025_lgb_1036_075_all_1021.csv</b> 7 days ago by Aleksandr Avdyushenko 025_lgb_1036_075_all_1021	1.02831	1.01663	✓
<b>single_subm_0.118646_0.872081_2018-12-17-05-26.csv</b> 7 days ago by Aleksandr Avdyushenko 2_lgbm_smote	1.04803	1.03659	✓

## Непроверенные идеи

- catboost
- lstm
- умный blending с oof, stacking

В целом очень позитивные впечатления!  
Тратил примерно 5-10 ч времени в неделю.



# Part of #26 solution

Simple RNN baseline (in pytorch), kernel <https://www.kaggle.com/johnfarrell/plasticc-2018-emb-gru>

Only uses 5 raw series: 'mjd', 'flux', 'flux\_err', 'detected', 'passband' (embedding dim is 16) & meta features: 'ddf', 'hostgal\_photoz', 'hostgal\_photoz\_err', 'distmod', 'mwebv'.

The score result is **oof0.650/pub0.938/pri0.970 (~ 90 place, bronze)**.

And it has very low correlation with stats features so stacking works very well with other models using stats features.

Решение	Score	Место
<ul style="list-style-type: none"> <li>• LightGBM</li> <li>• Bayesian approach to removing noise:</li> <li>• <math>\text{flux} = (\text{flux}/\text{flux\_err}^{**2} + \text{flux\_mean}/\text{flux\_std}^{**2}) / (1/\text{flux\_err}^{**2} + 1/\text{flux\_std}^{**2})</math></li> <li>• adding <b>features based on scaled flux</b></li> <li>• features to <b>capture the behaviour around the peak</b></li> </ul>	0.84070	<u>14</u>
<ul style="list-style-type: none"> <li>• LGBM model inspired by Oliver kernel, feature design and selection, and augmentation</li> <li>• magnitudes, peak widths</li> <li>• <b>parametric curve fittings</b></li> <li>• k-correction</li> <li>• <b>augmented train 10x times</b> using flux variation within normal distribution of the corresponding flux_err, similar for photoz value</li> <li>• Feature selection: manual approach and eli5 (permutation importance)</li> <li>• 50/50 blend of two LGBM models with slightly modified features and parameters</li> </ul>	0.82691	<u>13</u>
Stacking of <ul style="list-style-type: none"> <li>• lgb as in <u>Chia-Ta's kernels</u> with whole series</li> <li>• lgb with features from detected==1 and train for galaxy and ex-galaxy separately</li> <li>• MLP (multilayer perceptron) as in <u>Siddhartha's kernel</u> with same features</li> <li>• RNN feature extraction with attention and <b>pseudo labels 0.95 -&gt; 0.80</b></li> </ul>	0.80905	<u>8</u>

Решение	Score	Место
<p>NN consisting of 3 components:</p> <p>Meta Encoder - taking as input meta features hostgal photoz, hostgal photoz err, is galactic and summary stats for flux and flux_err as min/max/mean etc.</p> <p>Light Curve Encoder - bidirectional GRU taking as input grouped by day flux, flux_err, detected and time difference (autoencoder on test data)</p> <p>Spectroscopic Redshift Predictor – two fully-connected layers for predicting hostgal specz (on test data too)</p> <p>Outputs of these 3 components are fed into two last fully-connected layers for predicting class probabilities.</p> <p>Augmentations:</p> <p>flux as a Gaussian with standard deviation flux_err, hostgal photoz as a Gaussian with standard deviation hostgal photoz err, randomly dropping observations</p> <p>Top submission is an average of 5 cross-validation runs of 3 neural network variations.</p>	0.80173	<u>6</u>
<p>Blend of LGB, NN and several stacking models:</p> <ul style="list-style-type: none"> <li>• probing class_99</li> <li>• predict hostgal specz using training set+ test set with hostgalspecz</li> <li>• using normal values and log transformed values together on NN</li> <li>• Bazin: this is a light curve fit method</li> <li>• Log Ensemble: Instead of averaging predictions, I have averaged the logarithm of the predictions because in the end we try to regress the log values</li> </ul>	0.70423	<u>4</u>

Решение	Score	Место
<p>NNs worked much better than LGB models. Final ensemble (stacking) included 9 models, 7 of which were NNs that scored as low as 0.75x individually. The two LGB models were much weaker, each scoring about 0.90x, but they did provide a bit of diversity to the ensemble</p> <p>Hostgal <i>specz pseudo-labeling</i></p> <p>Flux adjustments</p> <p>Augmentation (a little bit)</p>	0.69933	<u>2</u>
<p>(an astronomer studying supernova cosmology)</p> <ul style="list-style-type: none"> <li>augmented the training set by degrading the well-observed lightcurves in the training set to match the properties of the test set</li> <li>use Gaussian processes to predict the lightcurves</li> <li>measured 200 features on the raw data and Gaussian process predictions</li> <li>trained a single LGBM model with 5-fold cross-validation</li> <li>lots of different features</li> <li>class 99 objects was a weighted average of the predictions for classes 42, 52, 62 and 95</li> </ul>	0.68503	<u>1</u>

«Чуваки усреднили решения топ-5 команд и пробили скор 0.6 на паблице!!»

<https://www.kaggle.com/c/PLAsTiCC-2018/discussion/75179>

We finally reached 0.5x using 1st, 2nd, 3rd, 4th, 5th subs!  
Really Great :)

Submission and Description	Private Score	Public Score	Use for Final Score
<a href="#">last_sub.csv.gz</a> 2 minutes ago by <a href="#">mamas</a> 1st * 0.33 + 2nd * 0.33 + (3rd + 4th + 5th)/3 * 0.33	0.61146	0.59828	<input type="checkbox"/>