Data Science Bowl® 2019: Predicting kids performance in educational game apps

Juan Jose Garau (garau@mit.edu), Inigo de la Maza (delamaza@mit.edu)

1. Background and Objectives

- The Data Science Bowl® is the world's premier **data science for social good** competition, it brings together experts across industries to take on the world's challenges with data and technology
- The 2019 edition partners with **PBS KIDS**, a trusted name in early childhood education (3-5 years old), who aims to gain **insights** on how **media** can **help** children learn important skills for success in school and life
- Based on historical data from its app PBS KIDS Measure UP!, the challenge focuses on predicting scores on in-game assessments and understanding how to improve learning outcomes







2. Data

- The dataset contains **timestamp data** from the app activity of approximately **17,000 kids**, divided into different **sessions** (~300k) and **events** (~11M)
- Some of the sessions are **Assessments**, in which the kid must complete an evaluation task related to some measurement concept (length, weight...) • The **goal** is, given **incomplete information** of a kid's actions, **predict how well**
- he/she will do in an Assessment, by **classifying** him/her into **one** out of **4** performance groups



Time series visualization



assessment



- dataset

```
Which
```



3. Feature Engineering

For each player (installation ID), we create different instances that include all previous events of that player until the beginning of every completed

┣-+-+	<u>· · · ∤ · ·</u>	$\vdash + + + + + + + + + + + + + + + + + + +$	+ + + + + + + + + + + + + + + + + + + +	-+-+-
tivity tart	Assessment 1	Assessment 2	Assessment 3	Assessment 4
	▶ Instar	Instance 2	Instance 3	
				→ Instance 4

• We **clean** and **enrich** the data by:

Getting rid of the players who haven't completed any assessment Using completed assessments in the test data as part of the training

III. Creating a function that computes the classification label for **completed** assessments

We then extract the following features for each instance:

Assessment-related features (how hard the assessment is)



performance with logistic regression and xgboost

\mathbf{N}	letric	
Accura	ıcy	
Quadr	atic kapp	a

- surpassing *xgboost* and scoring **top 3%** in the official competition
- Models with autobalance provide more interpretable solutions at the expense of predictive power • OCT-H without autobalance proves to be the best model in terms of out-of-sample performance,

Interpretability

- OCT with autobalance makes a decision based on which assessment, the player's experience on that assessment and the **overall performance** of the player
- OCT without autobalance includes the assessment average difficulty and specific event performance Both OCT models are able to identify the hardest assessments



15.095: Machine Learning under a modern Optimization Lens

4. Performance results

• We test **OCT** and **OCT-H** models with and without class *autobalance* and compare its out-of-sample

Model									
With autobalance		Without autobalance		Python-based					
OCT	ОСТ-Н	OCT	ОСТ-Н	Logistic regression	xgboost				
0.569	0.548	0.620	0.633	0.602	0.607				
0.473	0.490	0.514	0.540	0.501	0.472				