

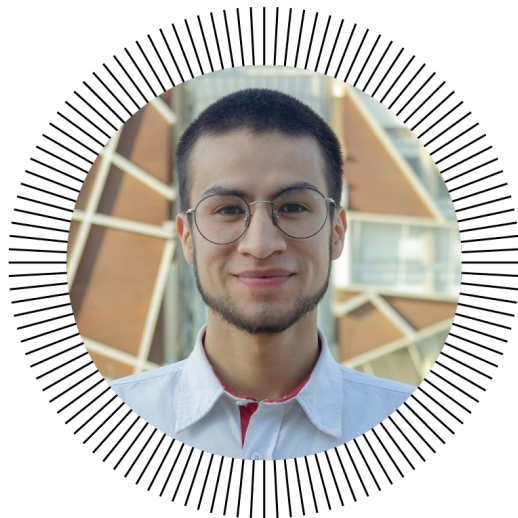


Categorization Challenge

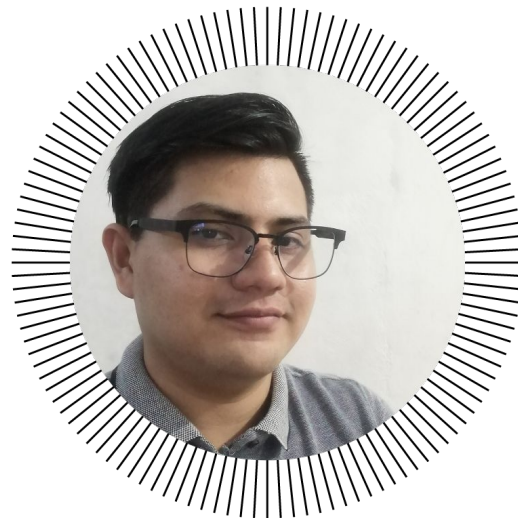
Datarangers



Datarangers Team



Edgar Steven Baquero



Jorge Daniel Gutierrez



The story of Failures Learning...

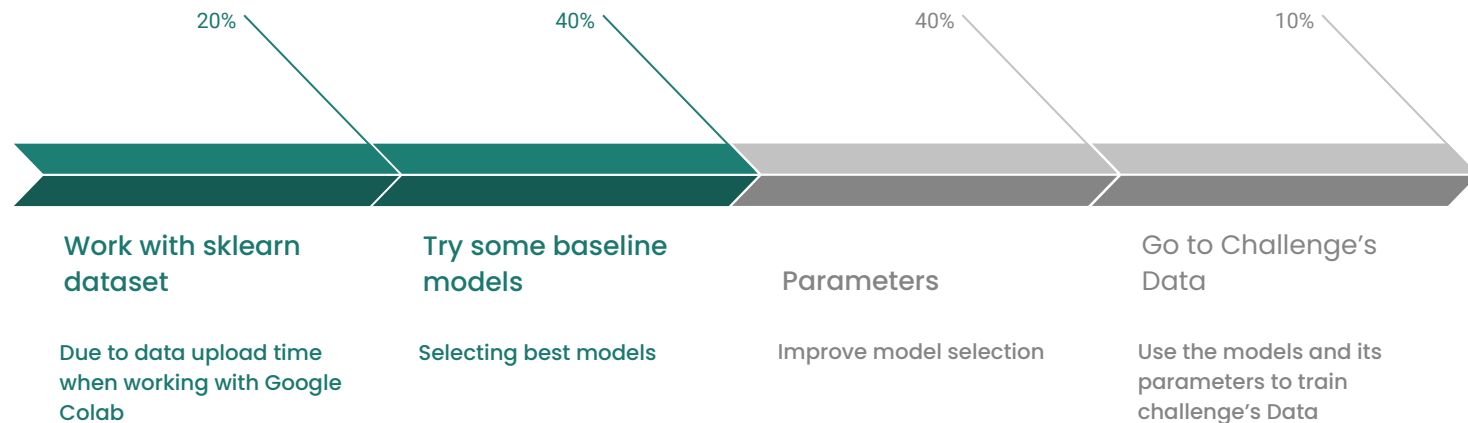
A happy story with a sad ending:

- How we handle the problem
- Results of the happy part
- Results of the sad part
- Some conclusions



Handling the problem

Handling the problem – Initial insights



Handling the problem - Sklearn dataset

- We flatten image data
- It lead us to a classical regression problem

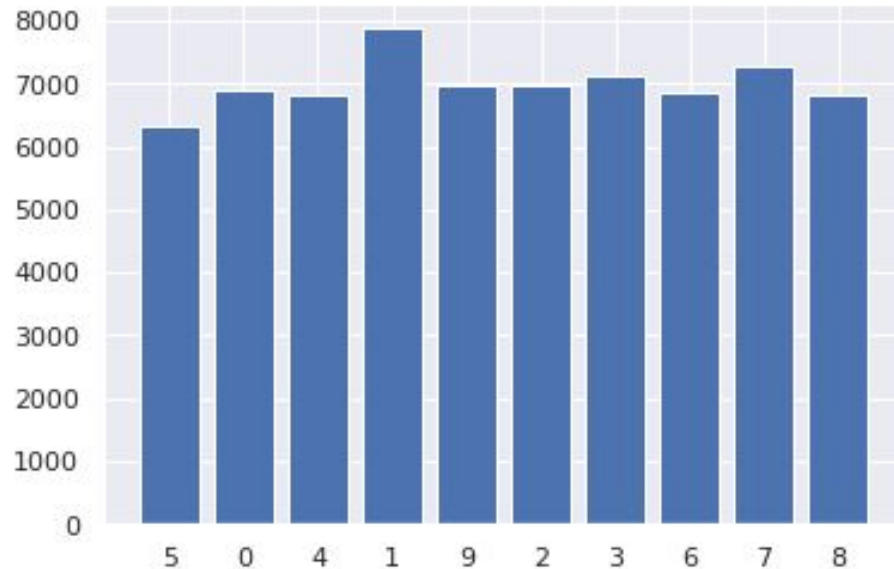


We splitted the Train dataset
into

Train	Test
80%	20%

Handling the problem – Distribution of categories

It seems like data is balanced:



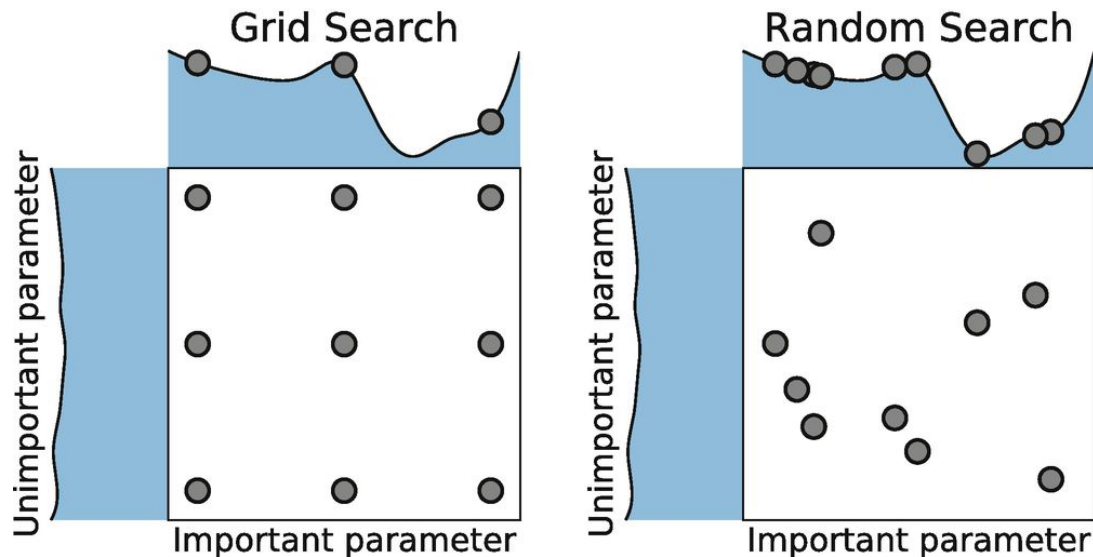
Handling the problem – Models



1	Mult. Regression	5	Tree
2	Linear Discriminant Analysis	6	Quadratic Discrimination Analysis
3	Logistic Regression	7	SVM
4	Neural Network	8	Boosting

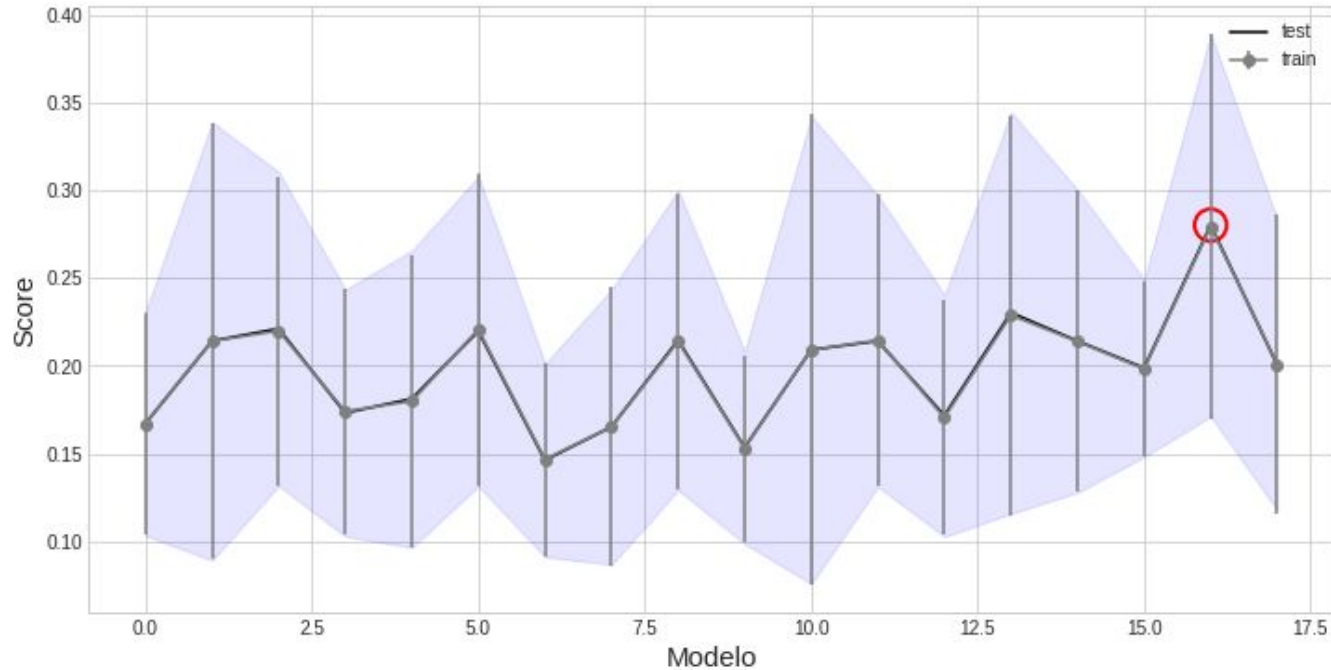
Handling the problem - Model Selection

We used Grid Search to select models



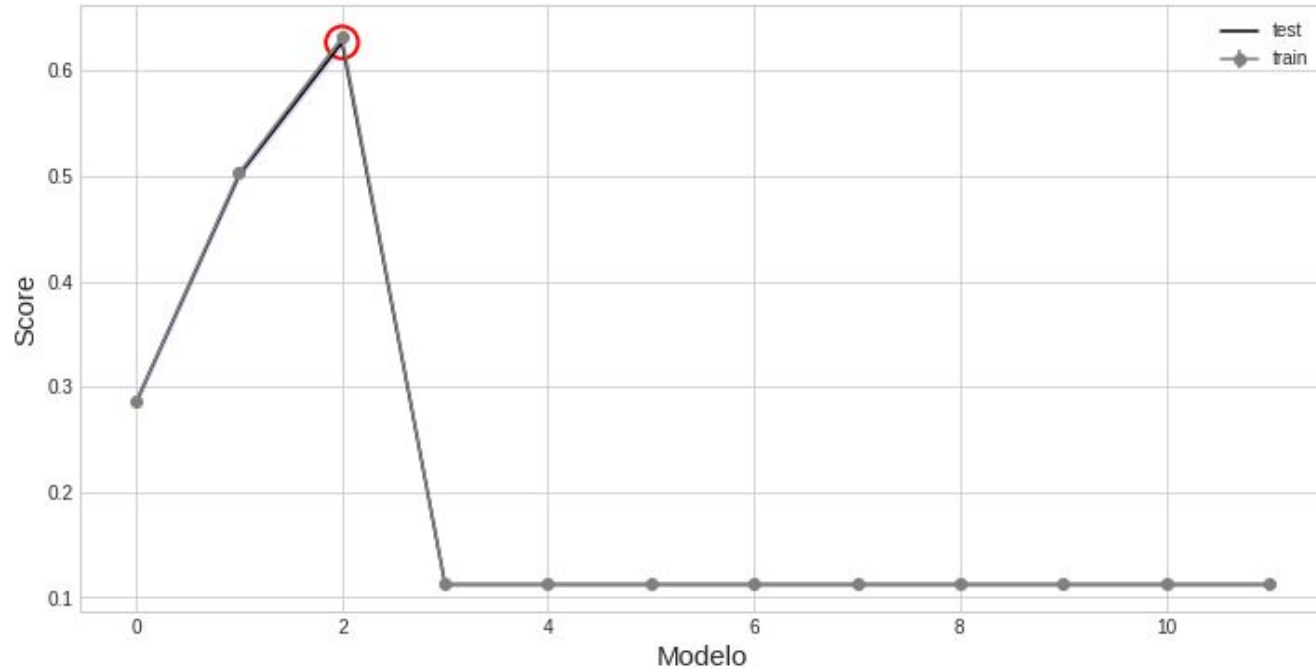
Model Selection – Results

Grid Search CV - Neural Network



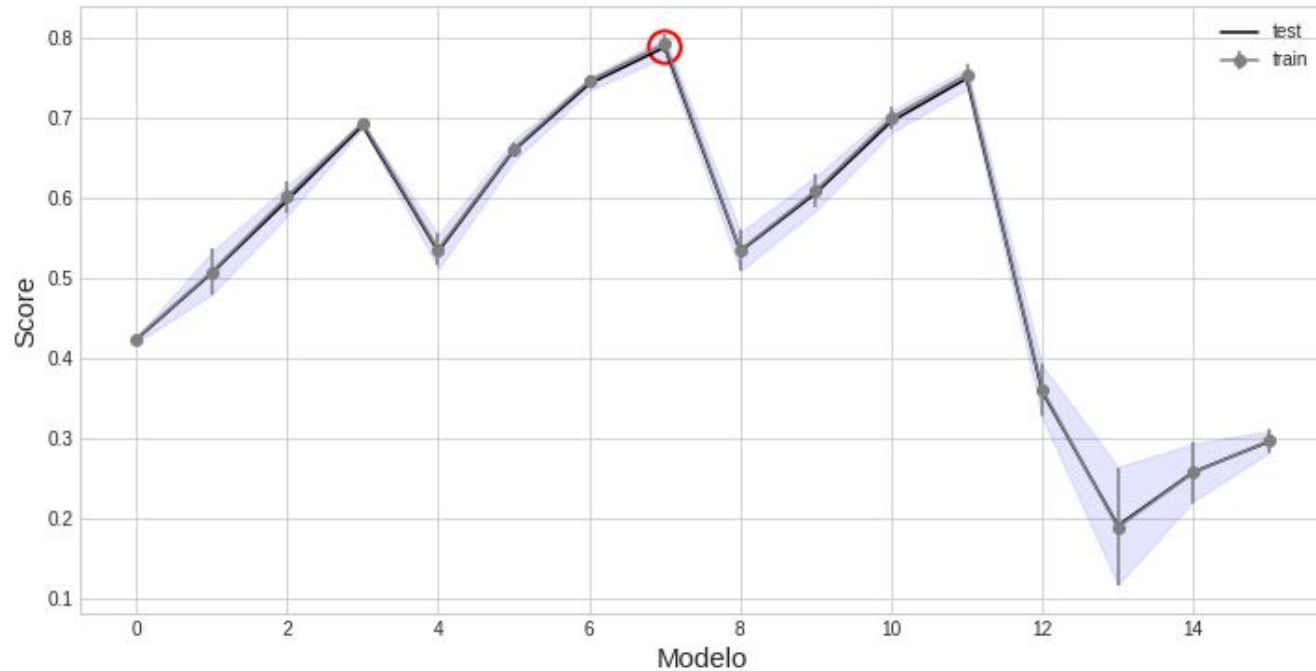
Model Selection – Results

Grid Search CV - Decision Tree

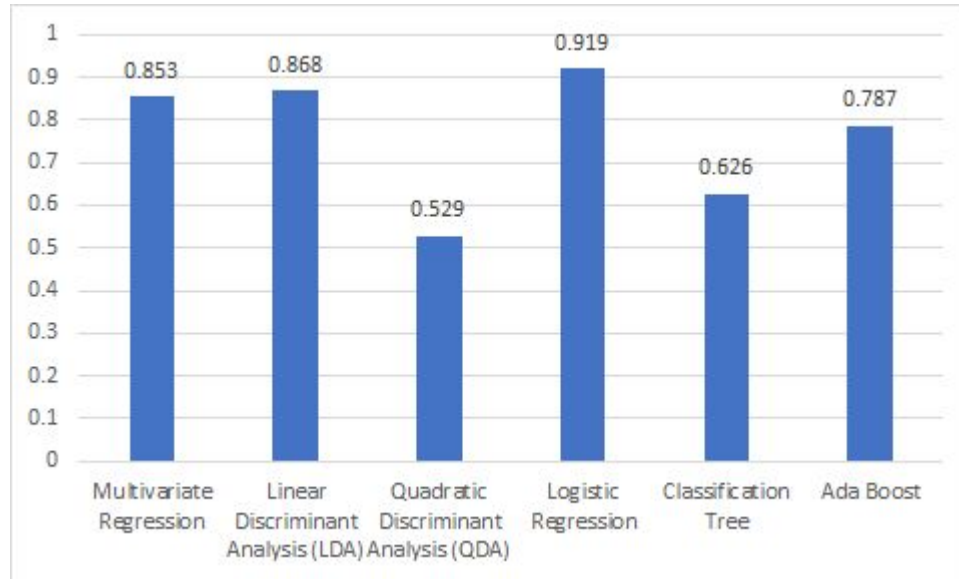


Model Selection – Results

Grid Search CV - Ada Boost



Model Selection – General Performance



Plot twist: Challenge's Data



mnist sklearn dataset

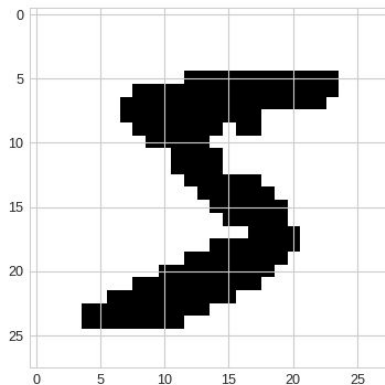
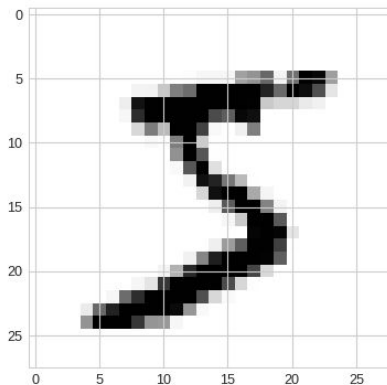
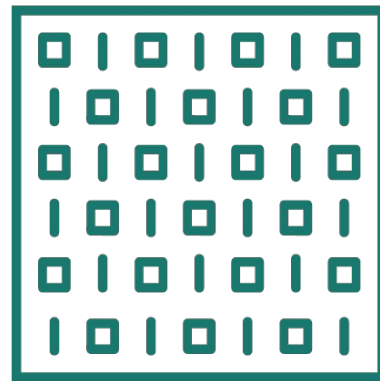


mnist bhchallenge dataset

"Everything went wrong when we changed dataset"

One way to handle challenge data

- We noticed that $\min(\text{train_image}) = 22$, while $\min(\text{test_image}) = 0$
- Also noticed that $\max(\text{train_image}) = 220$, while $\max(\text{test_image}) = 255$
- Standardize with Boolean Masks!



To improve for the next Challenge



- Display the data before the challenge dataset
- Convolutional neural networks
- Take naps