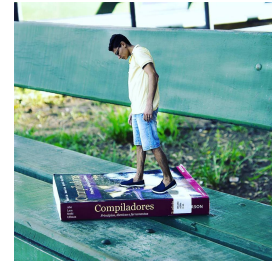# Brazil against the advance of Covid-19
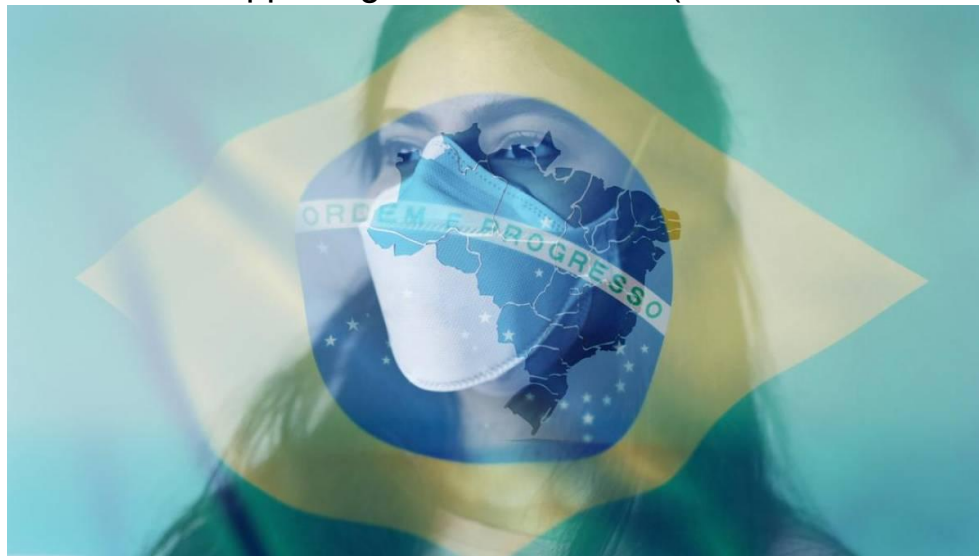
07th May

Crislânio Macêdo

# Diagnosis of COVID-19 and its clinical spectrum

AI and Data Science supporting clinical decisions(from 28th Mar to 1st Apr)



source

# Task Details

Predict confirmed **COVID-19** cases among suspected cases. Based on the results of laboratory tests commonly collected for a suspected **COVID-19** case during a visit to the emergency room, would it be possible to predict the test result for **SARS-Cov-2** (positive/negative)?
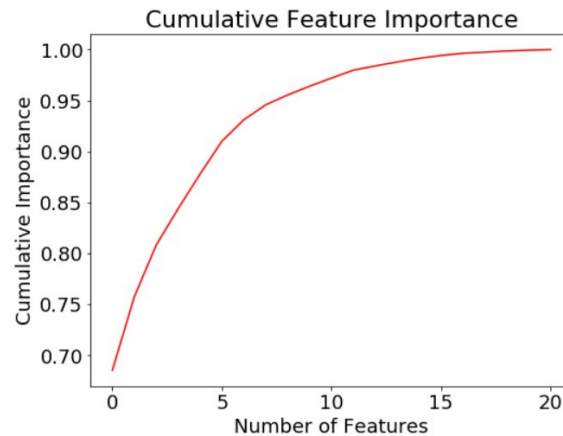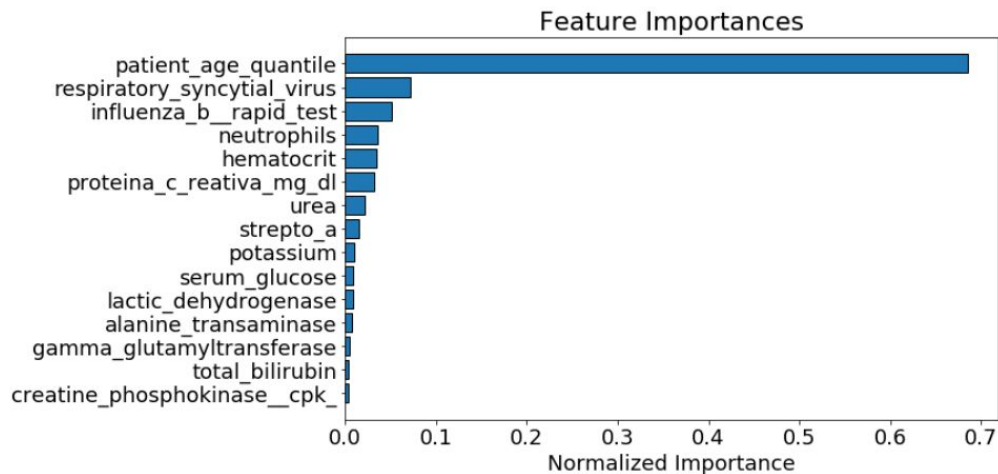
# Content

# Demo

Original features 111, after reduction 21 features.

# Research Question

What is the best way to predict positive results?

Focus on reducing false negatives. VS

Focus on reducing false positives. VS

Focus on a custom balance?

# Results -  StackNet

StackNet allows you to define all kinds of models. For example, Sklearn models, **LightGBM**, **RandomForest** and **CatBoost** can all be used with StackNet.

The **StackNetClassifier** will perform cross-validation (CV) and will output the CV scores for each model. To make sure we can output a probability of pacient result for sarscov2 exam result we specify "use_proba=True".

# Results -  StackNet

ROC Plot for StackNet Baseline

A model with high recall  achieves good results in finding positive patients among those true positive patients. But, the hospital may not have enough resources to apply the necessary procedures for all patients assigned with a positive label if that number is too high.

Hence, an ideal model is one that is well-balanced, i.e., one that has high recall but it does not over-assign patients with positive labels.

Classification report train

```
Confusion matrix:
 [[3558    1]
 [ 363   28]]
Classification report:
              precision    recall  f1-score   support

           0       0.91      1.00      0.95      3559
           1       0.97      0.07      0.13       391

    accuracy                           0.91      3950
   macro avg       0.94      0.54      0.54      3950
weighted avg       0.91      0.91      0.87      3950
```
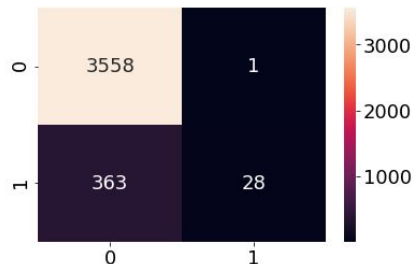


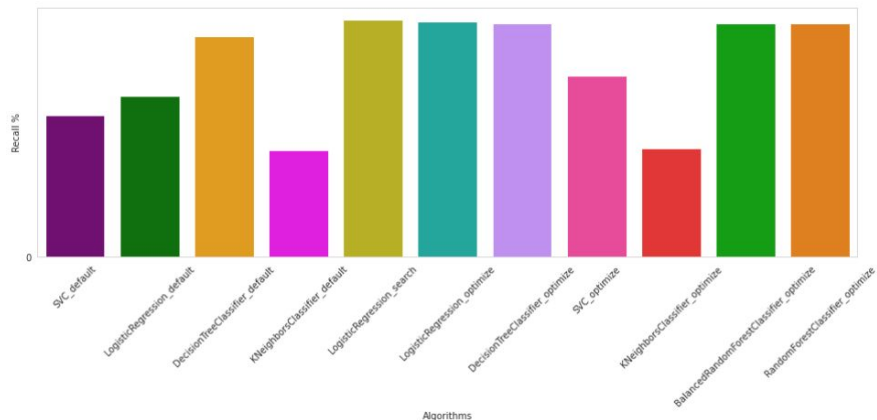What proportion of actual positives was identified correctly?

Our model has a recall of 7% —in other words, it correctly identifies 7% of all positive tests results for SARS-Cov-2.

What proportion of positive identifications was actually correct? 97%

Our model has a precision of 97%—in other words, when it predicts a positive test result for SARS-Cov-2, it is correct 97% of the time.

# Results

| | Train_Recall | Test_Recall | Test_Specificity | Optimize | Mean_RecSpe |
|---|---|---|---|---|---|
| SVC_default | 0.689608 | 0.546218 | 0.484158 | 0.333806 | 0.515188 |
| LogisticRegression_default | 0.77478 | 0.621849 | 0.452475 | 0.387974 | 0.537162 |
| DecisionTreeClassifier_default | 0.820591 | 0.848739 | 0.335644 | 0.546120 | 0.592192 |
| KNeighborsClassifier_default | 0.582707 | 0.411765 | 0.645545 | 0.258521 | 0.528655 |
| LogisticRegression_search | 1 | 0.915966 | 0.140594 | 0.560892 | 0.52828 |
| LogisticRegression_optimize | 1 | 0.907563 | 0.167327 | 0.559516 | 0.537445 |
| DecisionTreeClassifier_optimize | 0.988869 | 0.89916 | 0.262376 | 0.571803 | 0.580768 |
| SVC_optimize | 1 | 0.697479 | 0.370297 | 0.432043 | 0.533888 |
| KNeighborsClassifier_optimize | 0.582707 | 0.420168 | 0.646535 | 0.265441 | 0.533351 |
| BalancedRandomForestClassifier_optimize | 0.988869 | 0.89916 | 0.262376 | 0.571803 | 0.580768 |
| RandomForestClassifier_optimize | 0.988869 | 0.89916 | 0.262376 | 0.571803 | 0.580768 |

Recall: the proportion of positive cases that were identified correctly.

Specificity: the proportion of negative cases that have been correctly identified.

Better custom balance model is Balanced Random Forest with: **recall: 89.9%, specificity: 26%**

Focus on reducing false positives: log. regression **91.5%**

This result helps the hospital to release patients who obtained a pre-analysis diagnosed as NEGATIVE by the artificial intelligence system.

A situation in Brazil will improve #StayHome

# References

- COVID19: Recall
- Brazil against the advance of Covid-19
- EDA - first try / python lgb / shap
- XGB | LGB | CB
- LigthGBM simple FE
- Optimizing Imbalanced Classification
- Feature Selection
- Ensembling With StackNet
- EDA and Prediction

# References

- **<u>Kaggle</u> - @CaesarLupum**
- **Github - @crislanio**
- **Linkedin - @crislanio**
- **Medium - @crislanio.ufc**
- **Twitter - @crs_macedo**

Thanks !