

TRANSATLANTIC TEAM QUALITATIVE SUBMISSION

Github repo: <https://github.com/transatlantic-team/Pandemic-Prize>

I. Introduction of the research team and research directions

We are a team of 12 people from 3 continents: Europe, America, and Asia, with our leader (Martin Cepeda) currently based in Paris, France. Our team members include 5 undergraduate students, 1 post-doc, 3 senior researchers, and 3 university professors. The quantitative submission was mostly the work of the undergraduate students. Other team members contributed: infrastructure, data preparation and visualisation, and insight. In terms of infrastructure, we set up a [private competition](#) among ourselves on the Codalab platform on which team members made submissions, to monitor progress. We also set up a [Github repo](#) to share code.

While, due to time constraints, our team didn't have time to explore many avenues, even if our team members are experienced in machine learning (see e.g. the Google scholar pages of [Sergio Escalera](#), [Xavier Baro](#), [Prasanna Balaprakash](#) and [Isabelle Guyon](#)) and have studied the prediction and control of the Covid-19 epidemic (Yu et al, 2020; Cepeda, 2020). Hence we hope to make a more significant contribution in the second phase, if we are selected.

We are particularly interested in applying during the Prescriptor phase of the challenge Causal Modeling and Reinforcement Learning (RL) methods and in developing policy optimization, taking into account multiple factors. Since data availability plays a central role in obtaining good models, we started collecting additional data. Our axes for developing our models in Phase II include:

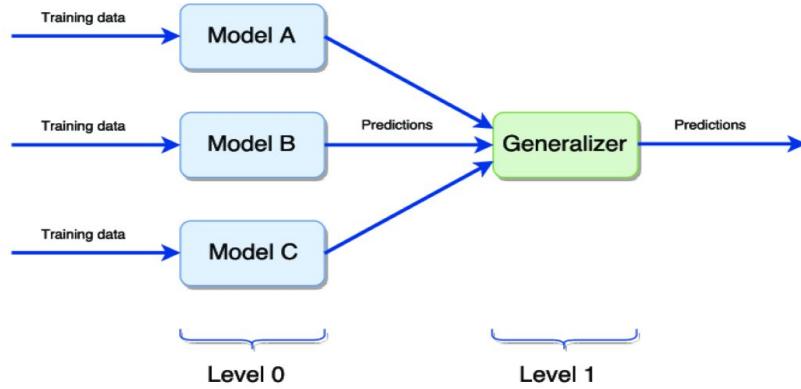
- **Causal analysis:** Policy evaluation using on-line data from previous pandemics remains a more or less open subject due to the novelty of this problem (e.g., Petersen, 2020). A notable example of this subject is a data-driven policy effect forecast (Vanderschaar, 2020) offering a counterfactual analysis framework for different countries. We are aware of the limitations of black-box predictive models based only on observable evolution of a pandemic, and the need to relate predictions to causes (e.g. lockdowns) (Goodman-Bacon, 2020). Structural equation modeling is a promising avenue, complementary to compartmental models (Pearl, 2009).
- **Economic factors:** While most authors according to our research focus on epidemiological models and neglect immediate economic impact, our interest is in blending epidemiology and econometric models, mostly short-term effects such as unemployment, business and school closures, transport reduction, which affect policy-maker decisions. Other indicators such as stock market indices, do [not necessarily correlate](#) with the economic impact of COVID-19. However, reduced morbidity and access to good healthcare facilities correlates to GDP growth (Alkire, 2018). More generally an increase in adult morbidity leads to a fall of economic growth (Javaid, 2015). Economic factors can be taken into account by appropriately defining the RL rewards, as done in (Cepeda, 2020).
- **Collateral health impact:** Another important aspect of policy optimization concerns collateral death and collateral adverse effects on public health, including untreated non-Covid related acute or chronic conditions, such as cancer, renal insufficiencies, or depression (Brodeur, 2020). For this reason, estimating secondary effects of untreated conditions due to health facilities being prioritized for COVID-19 in conjunction to pandemic response is important (Aron, 2020; Woolf, 2020). Future work could include creating a Causal/RL framework generalizing classical epidemiology models, which generally don't take into account such effects.

II. Innovation

For our pre-selection submission, our principal innovative point has been to create a meta-model based on several predictors, switching according to countries. We collected additional [country-specific time independent data](#) (see section IV). Although the individual models (Lasso, Ridge, SVR) were each trained on all the countries to gain robustness in predictive power (without country-specific data), our meta-model (figure 1) uses country-specific specialisation, making use of country-specific time independent data.

Final submission is about: Predicting 7 days smoothed daily new number of cases with robustness:

Figure 1
(source)
Schematic
of the
submitted
meta-model
with
stacking



- StackingRegressor with RidgeCV as main estimator to regularize
- Sub regressors: RidgeCV, LassoCV, BayesianRidge, LinearSVR

III. Generality

Our model performs rather well across all regions (see detailed results in appendix C.III). The following table summarizes the statistics for MAE per 100K in October test run:

Mean	Std	Median	Min	Max
1.80	4.37	0.23	0.00	32.98

However, we distinguish **3 types of regions** mutually exclusive between them according to the model's performance:

- Type 1: (MAE per 100K < 4) The model performs well across 1 month of predictions (206 out of 236 regions).
- Type 2: (MAE per 100K in [4, 10]) Bosnia and Herzegovina, Nepal, Romania, Portugal, Georgia, Italy, Costa Rica, Botswana, Moldova, Croatia, Peru, Brazil, Ireland, Iceland, Estonia, Poland, Cape Verde, Luxembourg, Oman, Slovenia, Belize, United Kingdom (country).
- Type 3: (MAE per 100K > 10, up to 33) Outliers where the model performs particularly bad: Switzerland, Netherlands, France, Spain, Bahamas, Belgium, Czech Republic, Israel and Andorra.

As the regions with the worst performance (Type 3) have in general a rather small population (Andorra, Israel, Bahamas, Czech Republic) the normalized performance per 100K inhabitants is severely penalized. For instance, even when the predictions in Belgium have roughly the same “raw” MAE as in Croatia, the MAE per 100K is 5.7 times bigger in Belgium.

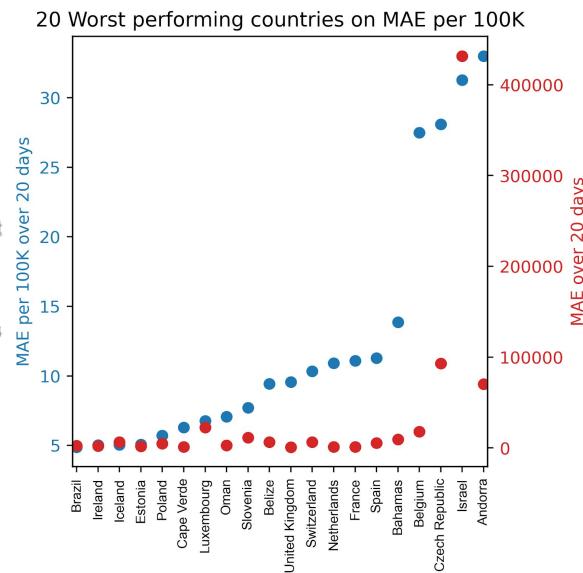


Figure 2
Countries sorted
by increasing
MAE per 100K.
Small countries
are among worst
performance
countries due to
normalization by
100K inhabitants

We can see the previous phenomena in figure 2: the gap between raw MAE and MAE per 100K is bigger when the country is smaller. For the Prescriptor phase, we'll pay special attention to the MAE in small countries.

We explain also the good performances in Type 1 and 2 countries mainly because of the stage of the pandemic: those countries have had a monotonous increase in cases and/or haven't reached a wake peak, whereas Type 3 regions are on a second, more powerful pandemic wave or never ended a first one (as in the US), which is a more complex evolution that our model fails to capture.

We show now a sample prediction for countries in each type (20 days prior to prediction period also shown):

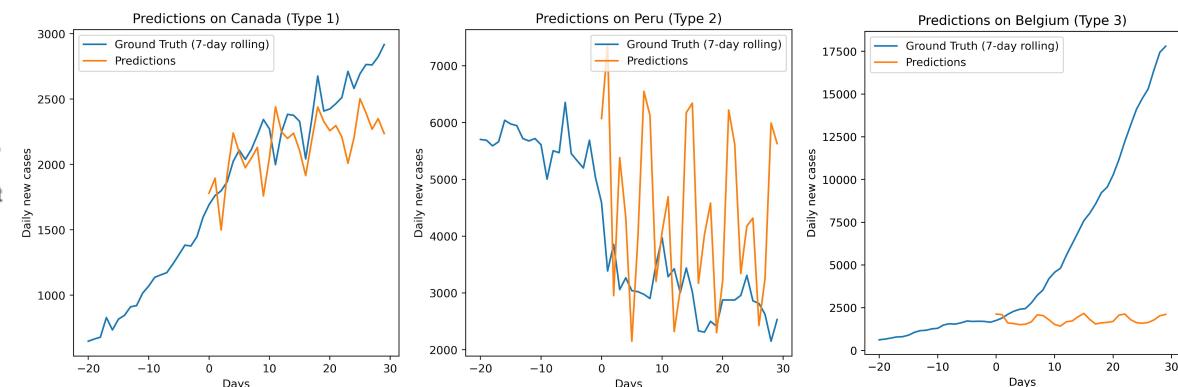


Figure 3
Sample
predictions in
countries from
Types 1, 2 and 3.
Negative values
on X axis are last
days of the
training period

Additional figures comparing different models can be found in Appendix C.IV.

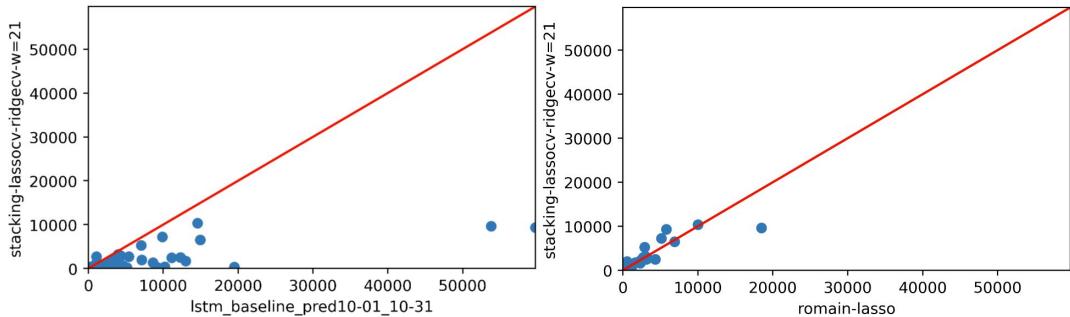
IV. Collaborative contributions

Our [code](#) and [country-specific time independent data](#) are open-sourced.

V. Consistency

We have homed in on linear predictive models because they seem most consistent on the short and long range. We are aware that this may not seem very refined, however, we compared with a variety of more complex models and found that the simplest models were the most robusts. A detailed analysis between different models can be found in annex C.II. Here we present the performance comparison between our submitter model, the LSTM baseline and a Lasso model:

Figure 4
Performance comparison between submitted model and baselines: Lasso Linear Regression and LSTM



In the plots, each point is a country and the axis represents MAE on 2 different models. The further the points are from the diagonal line (same error in both models), the greater the difference in performance. For instance, our model performs 5 times better than the given baselines over all countries in the period of 1 month of predictions (see annex C.II. for more information on model selection).

VI. Speed and resource use

Our model being based on linear predictors, it largely respects the time constraint imposed by the challenge (predictions in all countries on 180 under 1 hour) both to train and at prediction time: for prediction, running our model over all countries/regions for 180 days takes **less than 4 minutes** in the provided Sandbox environment.

For the submitted model on CPU Intel i7 6 cores: **Training up to 21st of December, 5 mins 30 seconds**. Predicting 30 days for all countries takes 19 seconds.

VII. Addressing the challenge

Prior to start modelling, we explored the available data (see Appendix C.I. and [our repo](#)). We discovered that a) new cases time series is very noisy, due to a certain periodicity (more cases are reported on Mondays) and changes in counting methodology per country since the beginning of the pandemic (which results in negative daily new cases), b) NPI data is not consistent (all countries have NaNs) and overall c) data from the first months must be taken cautiously, as COVID-19 was in an early stage and testing, data gathering and individual countries' response was not fully developed.

From this, we chose to rely on a) cumulated number of cases and b) 7-day smoothed new cases, as training input. Also we considered predicting (apart from daily new cases) the daily rate of change in new cases (used as a building block in autoregressive models). See Appendix C.V. for the exploratory modelling. We decided to take into account for the final modelling only the provided data and no external series such as deaths or bed occupancy, because these data would not have been updated during the testing phase. We did not attempt to exploit any loophole in the provided sandbox/predictor API or whatsoever.

VIII. Explanation

The stacking model presented in section II is trained on the whole OxCGRT dataset (01 of January to 21 of December). As it is an ensemble method (see section II), it first trains the sub-regressors (RidgeCV, LassoCV, BayesianRidge, LinearSVR) and then the meta RidgeCV regressor. Note that we also learn from NPIs as they are part of the dataset.

During the training and prediction stage, we consider a lookback window of 28 days to predict the next day. Predictions over an arbitrary time window (see section III) are computed via a rollout algorithm (Sutton 2018).

APPENDICES

A. Acknowledgements

This work builds upon internships on Covid-19 performed over the summer 2020, sponsored by ChaLearn and INRIA. We are grateful to Alain-Jacques Valleron, Sam Evans, Kristin Bennett, Paola Tubaro and John Erickson for help and advice.

B. References

(Yu et al, 2020) Yang Yu, Yu-Ren Liu, Fan-Ming Luo, Wei-Wei Tu, De-Chuan Zhan, Guo Yu, Zhi-Hua Zhou, COVID-19 Asymptomatic Infection Estimation. medRxiv, 2020
<https://www.medrxiv.org/content/10.1101/2020.04.19.20068072v1>

(Cepeda, 2020) Covid-19 risk mitigation. Master thesis.
<https://github.com/cepedus/COVID19-Risk-Mitigation>. Submitted to JDSE2021.

(Pearl, 2009) Judea Pearl. Causality. Cambridge University Press, 2009.

(Sutton 2018) Sutton, R.S. and Barto, A.G. Reinforcement Learning. An Introduction. 2nd Edition, A Bradford Book, Cambridge, 2018. incompleteideas.net/book/RLbook2020.pdf

(Petersen, 2020) Eskild Petersen, Marion Koopmans, Unyeong Go, Davidson H Hamer, Nicola Petrosillo, Francesco Castelli, Merete Storgaard, Sulien Al Khalili, Lone Simonsen, Comparing SARS-CoV-2 with SARS-CoV and influenza pandemics, The Lancet Infectious Diseases, 2020.

(Goodman-Bacon, 2020) Andrew Goodman-Bacon and Jan Marcus. Using Difference-in-Differences to Identify Causal Effects of COVID-19 Policies. Techreport, 2020.
<https://papers.ssrn.com/abstract=3603970>.

(Alkire, 2018) Blake C. Alkire, Alexander W. Peters, Mark G. Shrime, and John G. Meara. The Economic Consequences Of Mortality Amenable To High-Quality Health Care In Low- And Middle-Income Countries. Health Affairs, 2018.

(Javaid, 2015) Kiran Javaid et al. Morbidity and Economic Growth. Tech report Evans School Policy Analysis and Research Group, 2015.
<https://evans.uw.edu/policy-impact/epar/research/morbidity-and-economic-growth>.

(Brodeur, 2020) Abel Brodeur and Andrew E. Clark and Sarah Fleche and Nattavudh Powdthavee. Assessing the impact of the coronavirus lockdown on unhappiness, loneliness, and boredom using Google Trends. ArXiv preprint, 2020. <https://arxiv.org/abs/2004.12129>

(Aron, 2020) Janine Aron and John Muellbauer. A pandemic primer on excess mortality statistics and their comparability across countries. Our World in Data, 2020.
<https://ourworldindata.org/covid-excess-mortality>

(Woolf, 2020) Steven H. Woolf, Derek A. Chapman, Roy T. Sabo, Daniel M. Weinberger, and Latoya Hill. Excess Deaths From COVID-19 and Other Causes, JAMA, 2020.
<https://doi.org/10.1001/jama.2020.11787>.

(Bai, 2018) Shaojie Bai, J. Zico Kolter, and Vladlen Koltun. An Empirical Evaluation of Generic Convolutional and Recurrent Networks for Sequence Modeling. ArXiv preprint, 2018.
<https://arxiv.org/abs/1803.01271>

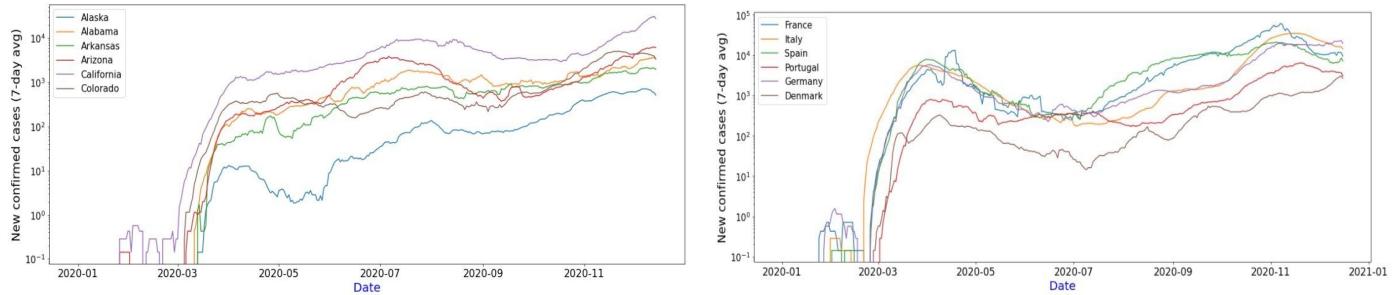
(Roser, 2020) Max Roser, Hannah Ritchie, Esteban Ortiz-Ospina and Joe Hasell. Coronavirus Pandemic (COVID-19). <https://ourworldindata.org/coronavirus>. Our World in Data, 2020.

C. Supplementary material

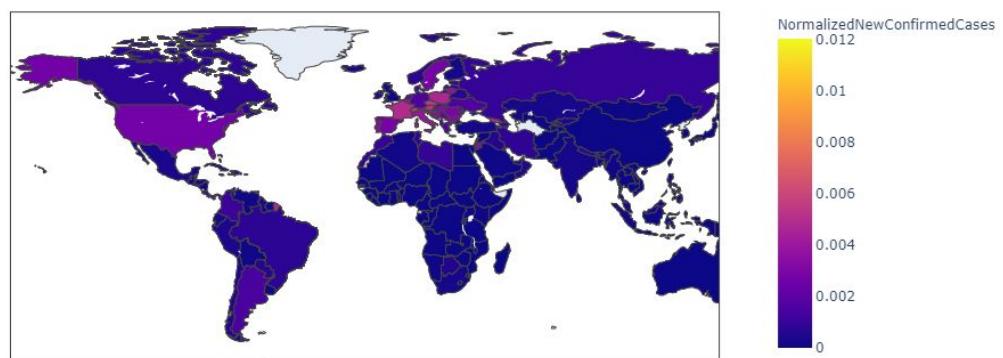
I. Data Understanding

Similar evolution in geographically close countries/regions (log scale)

We observed two principal waves of outbreak with shifted starting periods depending on the countries. Countries geographically close or with a lot of commercial/tourism mobility will most likely show similar behaviors.



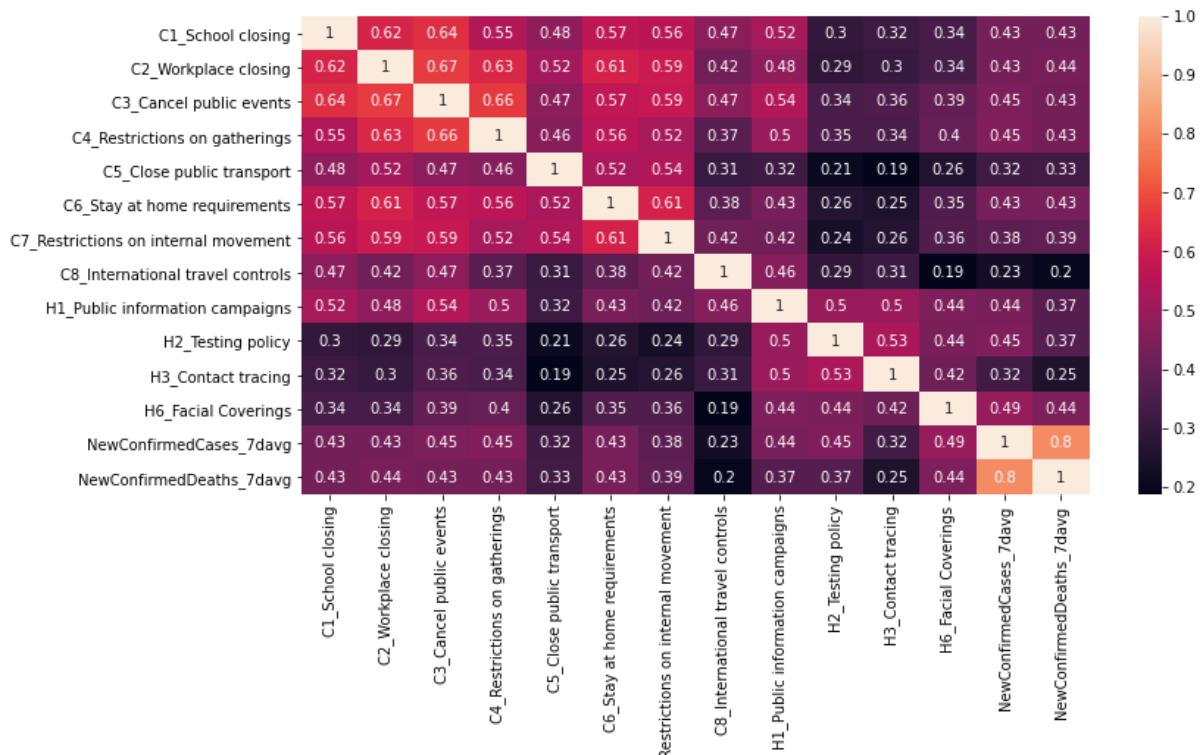
2nd wave visualization (ipynb available [here](#))



New cases normalized by country population.

Correlations when considering NPIs as ordinal data (using Kendall's Tau coefficient)

This correlation plot was to give some insight about the possible relations between the NPIs and the evolution of cases and deaths. Correlations between the NPIs themselves are understandable as their evolution is correlated but the correlation with the number of cases and/or deaths is not that obvious. Indeed, we would expect negative correlations at least for the PIs imposed even with a delayed impact.



Features Importance

In order to understand better the relation between variables, we performed a Random Forest Regression model and then selected the feature importances of the variables. There was no important relevance on any of the variables in the OxCGRT dataset.

II. Model selection:

To select the final submission, we did the following:

- 1) TrainRidgeCV, LassoCV, BayesianRidge, LinearSVR and stack them with RidgeCV over data up to september.
- 2) Compute prediction over 1 month (October 2020)
- 3) Compute MAE over 20 days (length of testing phase) taking as reference the 7-day smoothed number of cases in reality (same as the testing phase). These MAE scores are computed per-country to see on which of them a particular model performs the best.

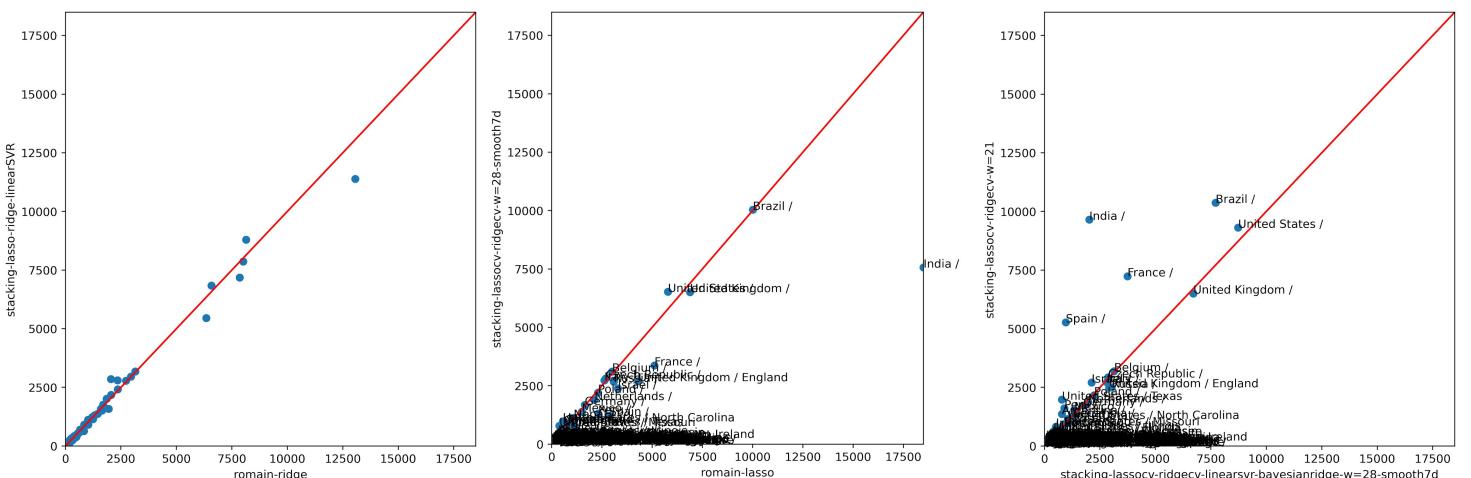
Table: cumulative 20MAE on all countries, alongside with worst performing country/region

romain-lasso	113188.24285 7	'India / '
romain-ridge	114758.81401 6	'Brazil / '
stacking-lasso-ridge-linearSVR	111055.40453 0	'Brazil / '
stacking-lassocv-ridgecv-linearsvr-bayesianridge-w=28-smooth7d	84126.300332	'United States / '
stacking-lassocv-ridgecv-w=21	107380.49242 1	'Brazil / '
stacking-lassocv-ridgecv-w=28-smooth7d	86885.818699	'Brazil / '

- 4) Comparing models side by side: we scatter the countries over their 20MAE on 2 models (axes) to visualize which countries make the difference:

Figures: model name on each axis, diagonal indicates equi-error. How to read: models in the upper triangle perform better on the x axis model and vice versa. The more aligned the points, the more similar both models perform. The axes are scaled to the max error across all models. Two very similar models halves error in countries in the upper triangle model on the y axes halves the error of India (worst error country) w.r.t model on x axis.

As we measure in this case the "raw" MAE over 20 days, the bigger countries such as Brazil, India and the United States have the biggest error and their points are in the upper-right region of the plots.

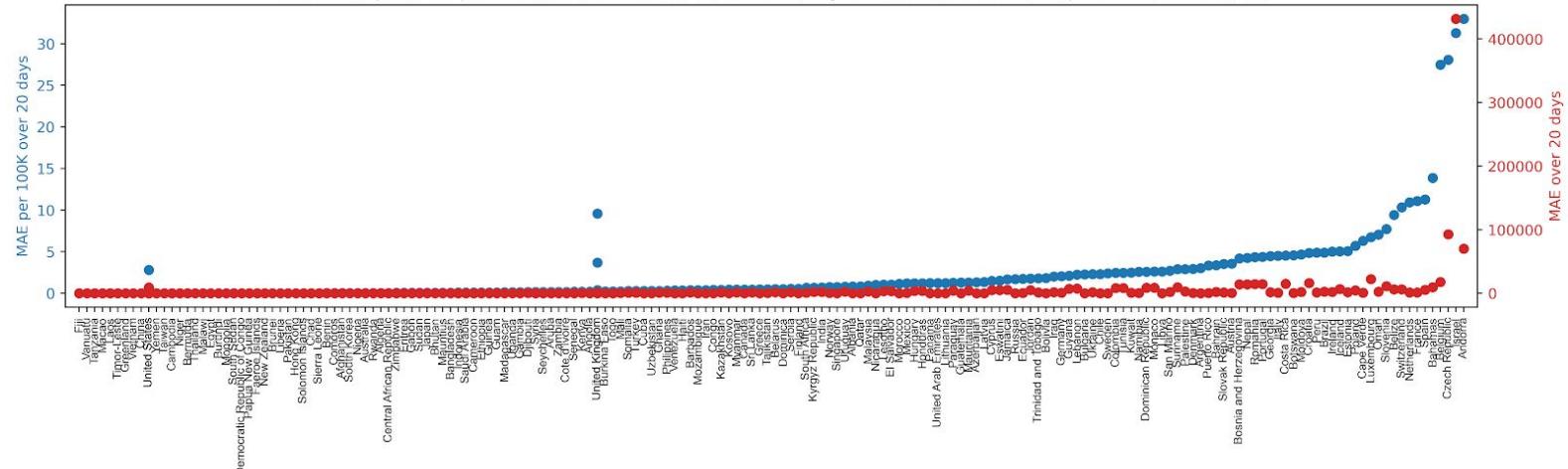


At the left plot: 2 very similar models (points more or less aligned over the equi-error diagonal). At the center plot: Model in the y axis halves the error in India (worst performing country in the x axis model). At the right plot: Overall improvement of model on the x axis, specifically in India, Brazil, France, Spain and the US (farthest points from diagonal).

- 5) Final submission: as we work with general statistical models, it makes little sense to have specific countries for each one of them. We go for a stacking approach, where the final model has no specific regions. For more complex models such as SEIR/SEIRD we could give a super-well fitted to only one country, but it loses the interest of the challenge itself to help predictions and generate prescriptions on a global scale.

III. Best predictor performance over different countries:

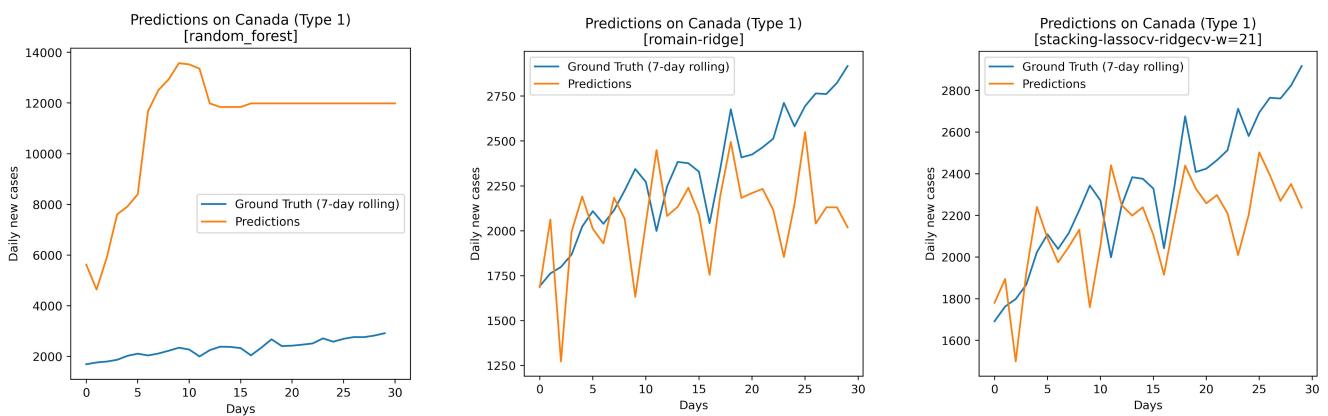
Submitted predictor performance (sorted) over all countries/regions (train until 31/09/20, predict until 20/10/20)



Across all regions we reach a rather good performance (MAE per 100K less than 10). Particularly bad predictions are Bahamas, Belgium, Czech Republic, Israel and Andorra. Some of these countries have either a) very strong 2nd waves of cases and/or c) haven't reached the peak of a wave.

IV. Comparison test run in October

From left to right: Random Forest, Ridge Regression, Submitted model (Stacking ensemble). All models were trained until 30/09.

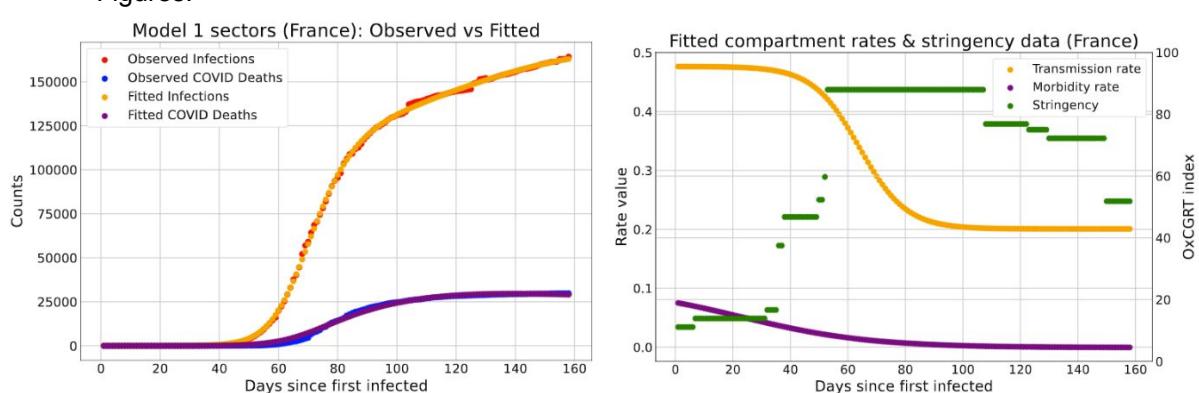


V. Other modelling choices:

On epidemiological models:

- Advantages:

- 1) Disease dynamics is accurate and has already been studied, in particular SEIR/SEIRD models are adapted to COVID-19 as the incubation period of the virus is important (4 days).
- 2) Comparatively few parameters to fit w.r.t. NN approaches and each parameter has a direct interpretation (contagion rate, death rate, etc.)
- 3) Once fitted, simulation is fairly fast depending on the ODE solver
- Disadvantages:
 - 1) It contains several latent variables (S, E, R) that are difficult to estimate from data, so to plug in a trained model on an arbitrary date implies estimating all latent variables also. Best option then is to fit the model from 'time 0' and run it also from time 0, and start considering as a prediction when we surpass the end of the training data
 - 2) Highly country dependent: depends on context and time-shift (when the first infected person is detected). Country dependence is necessary as that's what we're trying to predict and in COVID-19 international transport is severely reduced.
 - 3) For COVID-19, the whole game is how to map NPIs to changes in transmission/death/incubation rates (as they're clearly not constant, see case of the studies on R_0). To do that, the rate modelling can be arbitrarily complex (piecewise linear, piecewise constant, arbitrary continuous function, etc.)
 - 4) To learn rates over time, the loss is not straightforward to compute (as the rate is the input of the actual model) and highly non linear, due to the coupled differential equations for the EPI model. Thus any training is bottlenecked by the ODE solver.
- Figures:



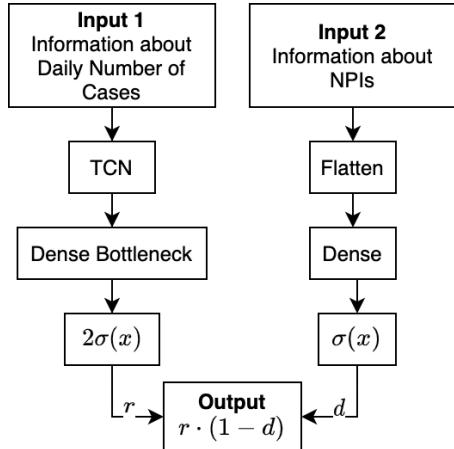
Data and fitted model in France from 25 January 2020 to 1 July 2020. At left, Infected and Dead compartments and at right, fitted transmission and morbidity rates compared with the value of the stringency index (sum of NPIs and non-NPIs). The accuracy on the model depends greatly on the assumption on transmission and morbidity rate. In this case, it was an affine sigmoid function.

On Neural Network-based models:

We tried NN models to predict the daily new number of cases (based on the linear baseline) or the normalized ratio of new daily cases (based on the LSTM baseline). In general these models performed worse than linear models. We noticed that TCN (Bai, 2018) models had a great potential to filter the input signal of the ratio of new daily cases. But, due to overfitting we did not use this model as our final submission.

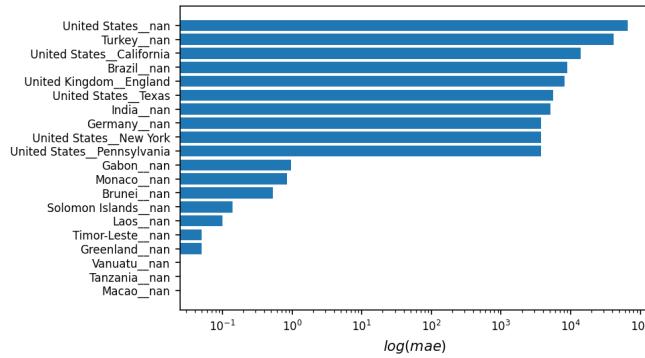
- Advantages:
 - 1) One of the main advantages is to be able to build a unique model which can learn from the data of all countries.
 - 2) A NN model can learn non-linear relationships between input and output variables.
- Disadvantages:
 - 1) It is hard to tune the hyperparameters and architecture of a NN and optimize its weights.
 - 2) Many overparameterized NN can have overfitting effects that can be hard to deal with.
 - 3) We noticed that adding a TCN on NPIS data was blocking all kinds of learning! However, adding it on the timeserie of new cases enables learning but massive overfitting.

- Figures:

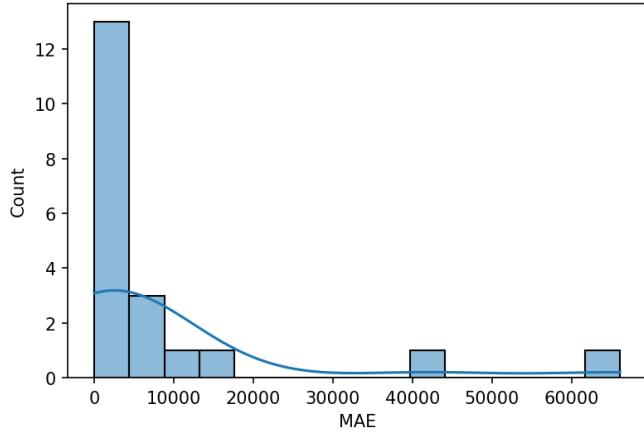


The TCN based model is inspired from the LSTM baseline. For the first input, instead of having an LSTM we filter the signal with a TCN, then a serie of dense layers (64, 32, 16, 8, 4, 2, and 1 units) to feed a rescale sigmoid function (to bound the output and avoid outliers).

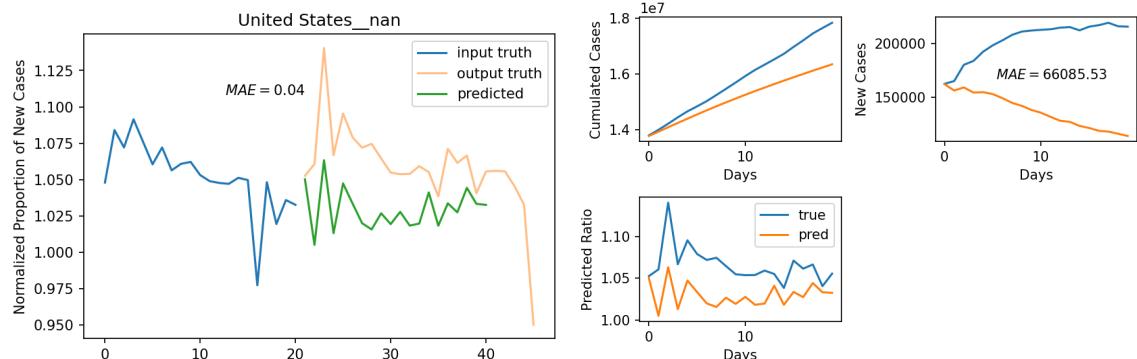
On the second input, we noticed simply use a single dense layer. In general, we were careful about doing a temporal split between our train/validation data (which is not done in the baselines), to avoid introducing a bias in our validation error.



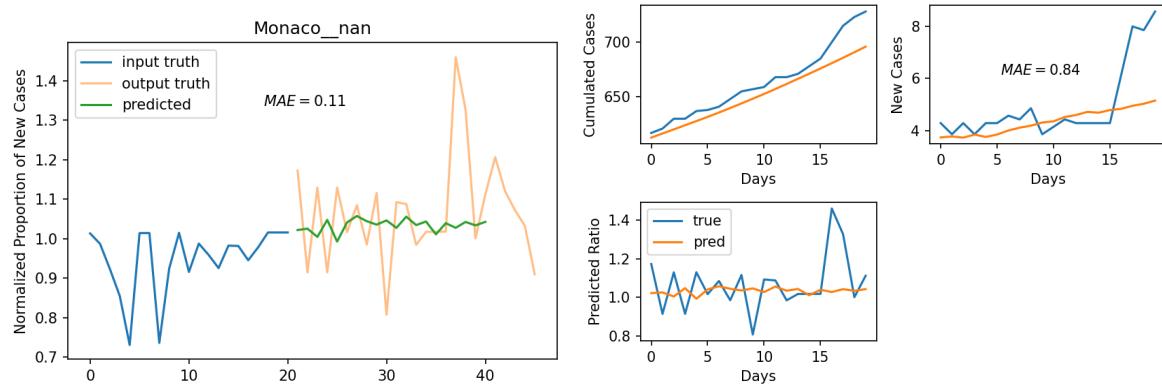
Prediction errors from 2020-12-01 to 2020-12-20, top-10 biggest errors (top) and top-10 smallest errors (bottom) are shown.



Histogram and estimated distribution of errors from 2020-12-01 to 2020-12-20, we can see that some countries are clearly outliers.



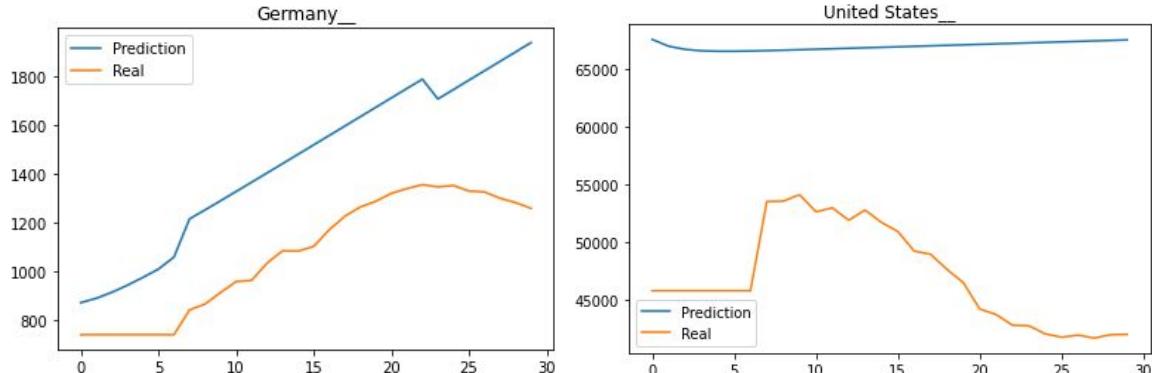
The United States is one of the worst performing countries with our TCN based NN. We can see that the ratio of new daily cases is clearly underestimated.



Monaco is one of the best performing countries, we can see that the general trend is followed without overfitting to the noise.

On ARIMA/X models:

- Advantages:
 - 1) Very fast to train and predict
 - 2) Adapted to auto-regressive time series such as the number of cumulated cases
 - 3) Idea: local context and look-back days will condition future evolution of the pandemic
- Disadvantages:
 - 1) On current data, it fails to capture trends and seasonality. (both in general and on specific countries)
 - 2) NPI context/NPI shifted context does not improve performance
- Figures: Prediction of daily new cases over 1 month (November 2020) in Germany and the US:



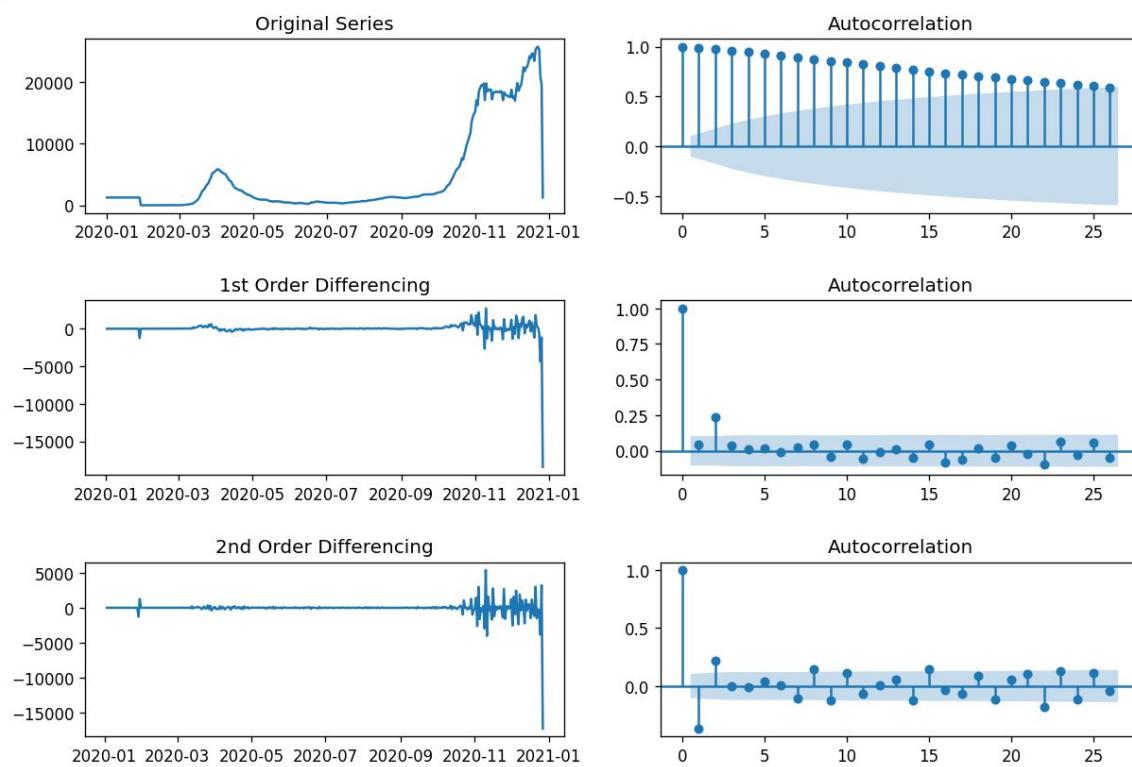
In the case of Germany, the ARIMAX model captures only the growth of cases, not the scale nor reaching the peak of a wave. In all, there's an MAE of almost 600. In the second figure, the model fails completely to predict new cases in the US, returning an almost constant prediction over the Month.

We also did an [Augmented Dickey Fuller test](#) (ADF) and Auto Correlation Function (ACF) plots analysis used in an attempt to fit Seasonal ARIMAX:

For instance for Germany:

ADF Statistic: -2.193192
p-value: 0.208726

Since p-value is greater than the significance level ($p\text{-value}>0.05$), we looked at the differencing and autocorrelation plots, to find the proper order of differencing.



From the plots above we see that a first order differencing makes the series of New Cases (smoothed over 7 days) in Germany stationary enough, no need for a second differencing.

On Random forest models:

We have tried one Random forest regressor model to try to predict daily new cases for the different countries. We have used the basic preprocessing. Hyperparameters selection includes the number of estimators, the maximum depth of the tree, and the maximum number of leaf nodes. Changing the criterion from “mse” to “mae” increased the training and testing time a lot (from minutes to several hours).

- Advantages:
 - 1) Can be very fast to train and predict.
 - 2) They are more understandable. We can use them to obtain feature importance information.
- Disadvantages:
 - 3) Hyperparameters have a big impact on the model; both in terms of results and time performance.
 - 4) The model quickly overfits the training data.
 - 5) They are not very good at capturing fluctuating trends in the medium/long term (as it can be observed in the figures).
- Figures: Results with our final model in 3 different countries. The model uses 100 estimators and “mse” as criterion. We train it with the data until 2020-09-30, and try to predict daily new cases from 2020-10-01 to 2020-10-31. The model takes 2'30" to train and 2' to make the predictions. The training/test set split is 0.8/0.2, the MAE error in the training set is 48 and in the test set is 113.

