

1 Team

Team Members	CS id: d8z9a Student id: 53033149
Kaggle Team Name	onlyme

2 Solution Summary

The solution uses features cases, cases_14_100k, cases_100k, deaths to predict the death count of the next days. In this model, we fit a linear regression of the features, and the next prediction is done with the predictions of the next feature itself. Also the following model has hyper-parameters k (number of time-series to look behind), p (the polynomial basis), \log_base (to use log base or not) and $use_prediction$ (to use the predicted death value for the next prediction or the prediction from the feature itself).

The best solution to get is the hyper-parameters that minimizes the error function. The best hyper-parameters are chosen by training the data to the date given by phase1 and predict the next days.

3 Experiments

Feature Selection

The following are chosen as features: cases, cases_100k, cases_14_100k and deaths to predict the next death cases for next day. The deaths will also be input in our feature selection because it modifies the value for the next death number cases.

Below are some plots related to our features.

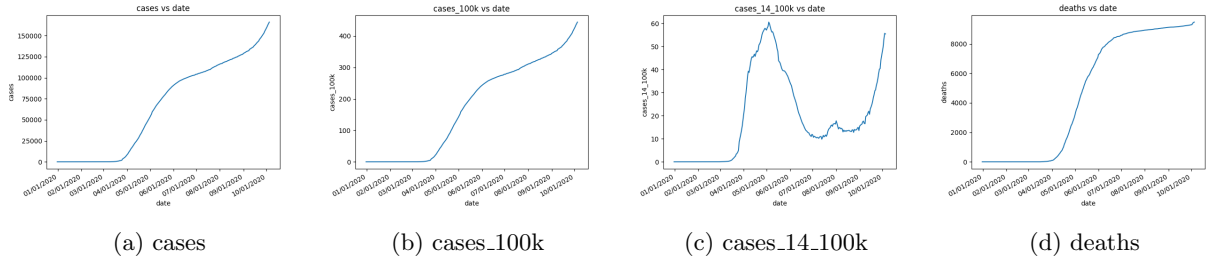


Figure 1: Plots of cases vs date

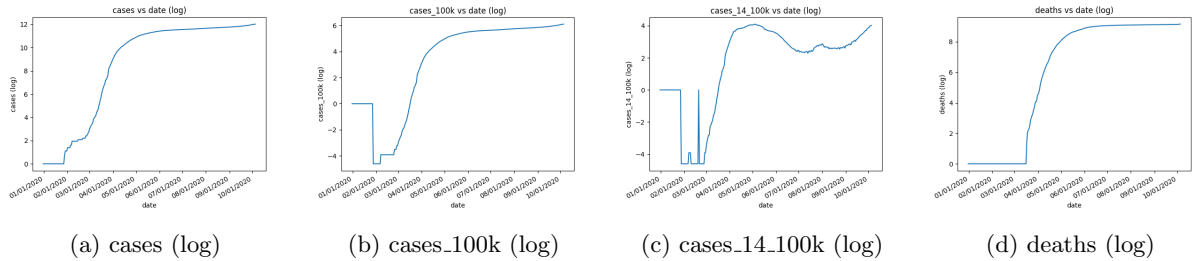


Figure 2: Plots of cases_100k vs date

The log base models are trained as well since we may be dealing with big numbers that polynomial basis cannot handle well. Also since it is a prediction with the case of a virus, we know a graph more or less has to be right skewed.

Hyperparameters

Choosing k hyperparameter

The k hyperparameter can be anything that is greater than 0 and less than the size of the feature. We chose this k value by trying bunch of different values for k . In general k can be anything that best fits our model.

Choosing p hyperparameter

In general, based on the graphs defined in Figure 1, we can see that no more than $p > 5$ is needed. Also since a greater number of p greatly influences the edges of our model (which is bad for predicting the future) having a big p will lead to extreme numbers as each time we try to predict the next values. Hence our hyperparameter will be $p < 5$ in general cases.

The value of p will also be affected greatly by the k hyperparameter. A big k will look at the k past data, so the greater k is, a greater p might generate a better fit while the lower k is, a linear fit might be better.

Other hyperparameters

The *log_base* and *use_prediction* are just binary inputs to whether we would like to use a log based data set to generate the model or use the prediction based on the fit of all data, or just keep using the fit from the death feature column.

Training

Training is done via minimizing our error function depending on our values of hyperparameters. For different values of k , p , *log_base*, *use_prediction*, we try to get the best values that minimizes the error.

Phase 1

In general for phase 1, the hyper parameters are chosen by fitting a model based on features until October 5 and then trying to minimize the error by comparing with next 11 days (October 16).

Below is a table summarizing the trials.

<i>log_base</i>	<i>use_prediction</i>	k	p	$\min(\sum(y_{pred} - y)^2)$
False	False	88	1	1177.149
False	True	127	1	1313.394
True	False	81	1	1678.265
True	True	137	2	1372.992

Table 1: Table summarizing k and p results for different *log_base* and *use_prediction* (phase 1)

For phase 1, the minimum error is given when *log_base* = *False*, *use_prediction* = *False*, $k = 88$, $p = 1$.

Phase 2

For phase 2, the values for the hyperparameters are also found using the same strategy for phase 1, now except trying to get best result from next 11 days starting October 5, the test is done with the next 20 days starting on October 5.

Below is the table summarizing the trials.

<i>log_base</i>	<i>use_prediction</i>	k	p	$\min(\sum(y_{pred} - y)^2)$
False	False	88	1	2184.482
False	True	127	1	3198.579
True	False	81	1	2210.375
True	True	79	1	1954.314

Table 2: Table summarizing k and p results for different *log_base* and *use_prediction* (phase 2)

The only difference between Table 1 and Table 2 is the last row. Phase 1's $k = 137$ and $p = 2$ could not possibly be used for further predictions because of the polynomial extreme values going faster over each prediction. Hence we got a new value $k = 79$ and $p = 1$ which we will use it to predict the next 5 days using now all data until October 25.

4 Results

Team Name	Kaggle Phase 1 Score	Kaggle Phase 2 Score
onlyme	10.13544	N/A

5 Code

```

1 class AutoRegressor:
2     def __init__(self, k, p=1,
3                 log_base=False,
4                 use_prediction=True,
5                 model=LeastSquares):
6
7         """
8         k                = the number of data sets to look behind to make next prediction
9         p                = polynomial basis
10        log_base         = change data to log space
11        use_prediction    = use prediction from X_k to test y_k+1, if false uses
12                           time series the the deaths column itself
13        model             = the linear model to use. We only use Least Squares
14        In this class, everything with variable X means it is not weighted
15        everything with variable Z means it has been transformed/weighted.
16        """
17        self.k = k
18        self.p = p
19        self.log_base = log_base
20        self.use_prediction = use_prediction
21        self.model = model
22        self.y_model = model()
23
24    def fit(self, X, y):
25        N,D = X.shape
26        k = self.k
27
28        # if k is too large don't do this.
29        self.size = N - k
30        if (self.size <= 0): return
31
32        # transform X to Z using transform method.
33        Z = self.__transform(X)
34        # get values of y for time series
35        y_k = y[self.size:N]
36
37        # set our features to predict next features.
38        self.__feature_time_series(X, Z)
39        # our general prediction for y_k+1
40        self.y_model.fit(Z, y_k)
41
42    def __feature_time_series(self, X, Z):
43        """
44        Grabs data from k to N-1 for each feature. If p > 0, also adds polynomial basis
45        columns.
46        Generates models to predict x_i_k+1 so that we can compute y_k+1.
47        The models here are by independent, meaning they will have their own
48        bias when fitting and predicting the model.
49        """
50        N, D = X.shape

```

```

50     p = self.p
51
52     feature_models = []
53     for i in range(D):
54         pos = i*p+1
55         # grab the data from Z and add a bias 1 or e.
56         x_i = self.__add_bias(Z[:,pos:pos+p])
57         x_k = X[self.size:N,i]
58         # create model and fit.
59         model = self.model()
60         model.fit(x_i, x_k)
61         # save model to predict x_i_k+1
62         feature_models.append(model)
63
64     self.feature_models = feature_models
65
66     def __transform(self, X):
67         """
68         Transforms X into Z. Adds the bias 1 if non logbase, e if logbase.
69         Adds polynomial basis given p.
70         """
71         N,D = X.shape
72         k = self.k
73         p = self.p
74
75         # new data set
76         Z = np.ones(shape=(k, D*p+1))
77         k_data = X[self.size-1:N-1, :]
78
79         for i in range(D):
80             pos = i*p+1
81             # add polynomial basis
82             Z[:,pos:pos+p] = self.__polyBasis(k_data[:,i])
83
84         # handle log case
85         if self.log_base == True:
86             Z[:,0] = np.exp(1)
87
88         return Z
89
90     def __polyBasis(self, X):
91         """
92         General function to change to poly basis.
93         NOTE: this does not add 1 as a bias.
94         """
95         N = X.size
96         p = self.p
97
98         Z = np.ones((N, p))
99         for j in range(1, p + 1):
100             Z[:,j-1] = X**j
101
102         return Z
103
104     def __add_bias(self, X):
105         """
106         Adds a bias of e if log_base, 1 otherwise.
107         """
108         N,D = X.shape
109         Z = np.ones(shape=(N,D+1))
110         Z[:,1:] = X
111
112         if self.log_base == True:
113             Z[:,0] = np.exp(1)
114

```

```

115         return Z
116
117     def predict(self, X, num_times):
118         """
119         Predict the next sequence of num_times.
120         """
121         def append_time_series(Z, y_k):
122             """
123             Predict x_i_k+1, removes x_i_1 from Z, and appends the new value.
124             If use_prediction is on, it will use the predicted value y_k+1 and its
125             polynomial biases for next predict. Else it will use the values from
126             its features series x_i_k+1.
127             """
128             N,D = Z.shape
129             p = self.p
130             feature_models = self.feature_models
131
132             i = 0
133             time_preds = np.ones(shape=(1,D))
134             for model in feature_models:
135                 # predict the new features
136                 pos = i*p+1
137                 x_i = self.__add_bias(Z[:,pos:pos+p])
138                 pred = model.predict(x_i)[-1]
139                 i += 1
140
141                 time_preds[:,pos:pos+p] = self.__polyBasis(pred)
142
143             # determine log base
144             if self.log_base == True:
145                 time_preds[:,0] = np.exp(1)
146
147             # remove first row, append to last row.
148             Z = Z[1:]
149             Z = np.append(Z, time_preds, axis=0)
150
151             # use prediction will change the deaths column to predicted value and its
152             # polynomial basis.
153             if self.use_prediction == True:
154                 Z[:,pos:pos+p] = self.__polyBasis(y_k)
155
156             return Z
157
158         Z = self.__transform(X)
159         y_k = self.y_model.predict(Z)[-1]
160
161         y_pred = np.zeros(num_times)
162
163         for i in range(num_times):
164             # predict, and save results.
165             Z = append_time_series(Z, y_k)
166             y_k = self.y_model.predict(Z)[-1]
167             y_pred[i] = y_k
168
169         return y_pred

```

Listing 1: AutoRegressor class implementation

```

1 if __name__ == "__main__":
2     # modify variables here
3     phase = "phase2"           # the phase to test
4     num_data = 5               # the number of data to predict
5     log_base = True            # variable if our model is in log base or not
6     use_prediction = True      # variable to use predictions in general over time series
7

```

```

8 # load file + set data
9 df = pd.read_csv(os.path.join("..", "data", phase + "_training_data.csv"))
10
11 X_codes = df.country_id.unique()
12 X_codes = X_codes[X_codes != "CA"]
13
14 X_ca = df.loc[df["country_id"] == "CA"]
15
16 ca_cases = X_ca["cases"].values
17 ca_100k = X_ca["cases_100k"].values
18 ca_14_100k = X_ca["cases_14_100k"].values
19 ca_deaths = X_ca["deaths"].values
20 dates = X_ca["date"].values
21
22 # combine column
23 X = np.column_stack((ca_cases, ca_100k, ca_14_100k, ca_deaths))
24
25 # if we are using log space
26 if log_base == True:
27     X = np.log(X)
28     X[X == -np.inf] = 0
29
30     ca_deaths = np.log(ca_deaths)
31     ca_deaths[ca_deaths == -np.inf] = 0
32
33 # run test only in phase 1 mode
34 # this only works phase 1 data.
35 # uncomment next line to get best k and p. NOTE: num_data has to be either 11 or 20
36 # _find_min_poly_basis(X, ca_deaths, num_data, log_base, use_prediction)
37
38 # code for plots
39 # dates_obj = [dt.datetime.strptime(d, "%m/%d/%Y").date() for d in dates]
40 # labels = ["cases vs date", "cases_100k vs date", "cases_14_100k vs date", "deaths vs
    date"]
41 # ylabels= ["cases", "cases_100k", "cases_14_100k", "deaths"]
42 # fnames = ["cases", "cases_100k", "cases_14_100k", "deaths"]
43
44 # for i in range (4):
45 #     if log_base == True:
46 #         __plot(dates_obj, X[:,i], fnames[i]+"_log", labels[i] + " (log)", "date",
    ylabels[i] + " (log)")
47 #     else:
48 #         __plot(dates_obj, X[:,i], fnames[i], labels[i], "date", ylabels[i])
49
50
51 # fit to autogressor model
52 model = AutoRegressor(k=79, p=1, log_base=log_base, use_prediction=use_prediction)
53 model.fit(X, ca_deaths)
54 y_preds = model.predict(X, num_data)
55
56 if log_base == True:
57     y_preds = np.exp(y_preds)
58
59 y_preds = y_preds.astype('int64')
60 print(y_preds)
61
62 # to csv
63 output = pd.DataFrame({'deaths': y_preds,
64                        'id': range(num_data)})
65 out_path = os.path.join("..", "data", phase + "_out.csv")
66 output.to_csv(path_or_buf=out_path, index=False)

```

Listing 2: Main function

```

1 def __find_min_poly_basis(X, y, num_data, log_base=False, use_prediction=True):
2     # first is for predicting 11 days in the future
3     # second is for predicting 20 days in the future
4     # comment/uncomment as needed
5     # y_test = [9504,9530,9541,9557,9585,9585,9585,9627,9654,9664,9699]
6     y_test = [9504,9530,9541,9557,9585,9585,9585,9627,9654,9664,9699,
7               9722,9746,9760,9778,9794,9829,9862,9888,9922]
8
9     min_error = np.inf
10
11     for p in range(1,5):
12         for k in range(1, 280):
13             model = AutoRegressor(k=k, p=p, log_base=log_base, use_prediction=use_prediction
14             )
15             try:
16                 model.fit(X, y)
17             except:
18                 continue
19             y_preds = model.predict(X, num_data)
20             if log_base == True:
21                 y_preds = np.exp(y_preds)
22             error = np.sum(np.square(y_preds-y_test))
23             # print("p: %d, k: %d, error: %.3f" % (p, k, error))
24             if (error < min_error):
25                 p_low = p
26                 k_low = k
27                 min_error = error
28
29     print("lowest p: %d, k: %d, error: %.3f" % (p_low, k_low, min_error))

```

Listing 3: Brute force to search for k and p

```

1 def __plot(x, y, fname, label, xlabel, ylabel):
2     plt.gca().axis.set_major_formatter(mdates.DateFormatter('%m/%d/%Y'))
3     plt.gca().axis.set_major_locator(mdates.MonthLocator())
4     plt.plot(x, y)
5     plt.gcf().autofmt_xdate()
6     plt.title(label)
7     plt.xlabel(xlabel)
8     plt.ylabel(ylabel)
9     plt.savefig(os.path.join(".", "figs", fname))
10    plt.clf()

```

Listing 4: Plotting function

6 Conclusion

In this project, I have learned how to do auto regressive models in general. I learned how to implement and interpret them. I have learned the importance of how to look at a data and transform it into a model. This was the hardest part of the project since we had to start barebone with just brief hints. Finally, I learned that when fitting a model, one should go from a simple model to complex, since trying to make something complex may not work well, which I ended up wasting a lot of time.

If I were given more time, I would have tried to get nearest neighbors of canada using KNN and also add them to the feature space. Also I would have tried doing multiple transformations into on feature space. For example, right now the model I implemented it either uses all polynomial basis of 3, or all basis of 2. I could have mixed this up for each feature, having different polynomial basis depending on the model which I would do the same case for the log space. Finally I would try using more linear transformations instead of just polynomial.

7 Additional Notes

Code style from: https://www.overleaf.com/learn/latex/code_listing