

Discover

The covid-19 pandemic has presented us with many new challenges through the last 2 year, to adapt to these ever changing conditions, we need to be able to adapt and constantly come up with new solutions to tackle these issues. To do so, we need to determine when is the ideal time to start our ideas and seize the opportunities in this uncertain time.

Define

We start in Singapore, where we are most familiar with. It is important to know when we will be able to flatten the curve so as to open up our economy, social activities in our community given the our vacations rates, new daily cases and our government's stance.

Develop

To do so, we came up with a multiple linear regression model to calculate what would the stringency index be in Singapore given our current condions based on our past experiences with the virus. This would allow us to have a good idea of when Singapore would be able to ease our businesses and daily activities and resume our normal lives.

Deliver

We shall provide a model to allow users to key in various parameters to calculate Singapore's expected stringency index

Problem Statement

Given a rise in covid-19 cases, the Singaporean Government would always take the necessary precautions to prevent an unprecedented rise in COVID-19 transmission. This increases the Stringency Index significantly as workplaces, schools closes. In such a case, what would the stringency index of Singapore be?

Start with importing Data that outlines Singapore's Covid-19 response

Data

Link to data source: <https://github.com/owid/covid-19-data/tree/master/public/data>

Why is this data used?

1. Sources are easily documented, for example, all data under Singapore comes directly from Ministry of Health, Singapore, which ensures that data is accurate and not tweaked or false
2. Contains very little missing values and is updated daily
 - Ensures that the data given is current with little errors
 - Missing data is filled using growth trendline analysis using excel
 - Ensures a growth trend in data such as vaccinations rates using values before and after, which can only be increasing
3. Provides a range of parameters
 - Ensures that we have a diverse range of parameters to accurately train our model
 - Provides smoothed values to ensure no abnormal values
 - 0 vaccinations prior to discovery of vaccinations and delivery to Singapore on 31 Dec 2020

Model

Features considered:

1. new_cases_smoothed_per_million
 - The number of new cases directly corresponds to the community transmission of COVID-19, giving a clear indication of the prevalence of virus in the community
 - Higher number of new cases signals high transmission rates in the community, leading to higher stringency index to keep new cases low
2. total_cases_per_million
 - Total number of infected individuals is an indication of how many people in the community have natural immunity to COVID-19
 - Higher total cases will lead to herd immunity, leading to lower stringency index
3. new_deaths_smoothed_per_million
 - Indicates the number of new deaths, an important indicator that is in line with Singapore's stance to be COVID resilient with keeping daily death rates low
 - Higher death rates will signal that the virus is becoming more deadly, leading to a higher stringency index to keep people safe
4. new_tests_smoothed_per_thousand
 - Indicates the number of test done each day, a higher number of testing done will help to ensure that positive cases are picked up earlier and therefore isolated and reduce the spread of the virus
 - Higher number of new tests will lower transmissibility leading to stringency index lowering
5. positive_rate
 - Indicates the percentage of test done that are positive, will give an indication of how many cases are there in the community that is transmitted without being picked up
 - Higher positive rates with higher reservoir of undetected cases will lead to stringency index increasing
6. people_fully_vaccinated_per_hundred
 - Indicates how many people are fully vaccinated and therefore lowering the chances of death as well as having a shorter infectious period, leading to slower transmission rates
 - Higher vaccination rates gives lower death rate and lower transmissibility will likely lead to lower stringency index
- Perform various test to ensure we find the optimal alpha value (learning rate), beta value (starting weight) and number of iterations (times gradient descent is performed)

Multiple Linear Regression

```
In [48]: import pandas as pd
import matplotlib.pyplot as plt
import seaborn as sns
import numpy as np

# Data here is from the best growth trendline analysis using excel and it's where places are obvious
df = pd.read_csv('data/task 2/Singapore_Covid_Data.csv.csv')
display(df)
```

	location	date	new_cases_smoothed_per_million	total_cases_per_million	new_deaths_smoothed_per_million	reproduction_rate	new_tests_smoothed_per_thousand	positive_rate	people_fully_vaccinated_per_hundred
0	Singapore	1/30/2020	0.445	0.183	0.000	0.56			
1	Singapore	2/3/2020	0.498	0.550	0.000	0.54			
2	Singapore	3/3/2020	0.498	0.550	0.000	0.53			
3	Singapore	4/3/2020	0.445	0.733	0.000	0.57			
4	Singapore	5/3/2020	0.629	0.937	0.000	0.73			
...
609	Singapore	31/10/2021	673.998	15193.634	0.288	0.97			
610	Singapore	1/11/2021	655.556	15496.195	0.367	0.94			
611	Singapore	2/11/2021	661.293	15760.762	0.419	0.92			
612	Singapore	3/11/2021	617.050	16116.327	0.472	0.90			
613	Singapore	4/11/2021	655.812	16418.329	0.472	0.88			
614 rows × 10 columns									

```
In [100]: # Functions from cohort and homework week 9

def normalize_z(df):
    return (df - df.mean(axis=0))/df.std(axis=0)

def get_features_targets(df, feature_names, target_names):
    # get df of selected features
    df_feature = df[feature_names]
    # get df of selected targets
    df_target = df[target_names]
    return df_feature, df_target

def prepare_feature(df_feature):
    # number of columns in the dataframe
    cols = len(df_feature.columns)
    # convert df to numpy
    feature = df_feature.to_numpy().reshape(-1,cols)
    array = np.concatenate((np.ones((feature.shape[0],1)), feature), axis = 1)
    return array

def prepare_target(df_target):
    cols = len(df_target.columns)
    target = df_target.to_numpy().reshape(-1,cols)
    return target

def predict(df_feature, beta):
    df_feature = normalize_z(df_feature)
    preped_feature = prepare_feature(df_feature)
    return predict_norm(preped_feature, beta)

def predict_norm(X, beta):
    return np.matmul(X,beta)

def split_data(df_feature, df_target, random_state=100, test_size=0.3):
    indexes = df_feature.index
    if random_state != None:
        np.random.seed(random_state)
        k = int(test_size * len(indexes))
        test_index = np.random.choice(indexes, k, replace=False)
        indexes = set(indexes)
        test_index = set(test_index)
        train_index = indexes - test_index
        # the above indexes just help us to get random indexes within the entire data
        df_feature_train = df_feature.loc[train_index,:]
        df_feature_test = df_feature.loc[test_index,:]
        df_target_train = df_target.loc[train_index,:]
        df_target_test = df_target.loc[test_index,:]
    return df_feature_train, df_feature_test, df_target_train, df_target_test

def r2_score(y, ypred):
    ss_res = np.sum((y-ypred)**2)
    y_mean = np.mean(y)
    ss_tot = np.sum((y-y_mean)**2)
    r_2 = 1-(ss_res/ss_tot)
    return r_2

def mean_squared_error(target, pred):
    num_data = target.shape[0]
    return (1/num_data)*(np.sum((target-pred)**2))

def mean_absolute_error(target, pred):
    num_data = target.shape[0]
    return (1/num_data)*(abs(np.sum(target-pred)))

def compute_cost(X, y, beta):
    #beta is weighted values, in this case it is just choosen from random values
    J = 0
    number_of_samples = X.shape[0]
    error = np.matmul(X, beta) - y
    error_sq = np.matmul(error.T, error)
    J = (1/(2*number_of_samples)) * error_sq
    return J

def gradient_descent(X, y, beta, alpha, num_iters):
    number_of_samples = X.shape[0]
    J_storage = []
    for i in range(num_iters):
        # derivative error = (1/number of samples) * np.matmul(X.T, (np.matmul(X, beta) - y))
        beta = beta + alpha * - derivative_error
        J_storage.append(compute_cost(X, y, beta))
    return beta, J_storage

# Single Function to make the model
# 1. alpha-value (step for gradient descent)
# 2. beta (starting beta values for gradient descent)
# 3. iterations (number of iterations of gradient descent)
# 4. start (starting row)
# 5. end (last row)
# 6. feature_parameters (features used to train model)
# Return r^2 and mse values + alpha value

def make_model(alpha, beta, iterations, start=0, end=None, feature_parameters=["new_cases_smoothed_per_million","total_cases_per_million","new_deaths_smoothed_per_million","reproduction_rate","new_tests_smoothed_per_thousand","people_fully_vaccinated_per_hundred","positive_rate"], dataset = "Data/Task 2/Singapore_Covid_Data.csv.csv", target_column=["stringency_index"]):
    df = pd.read_csv(dataset)

    # Extract the features and the target
    df_features_orignal_train, df_target = get_features_targets(df.loc[start:end,:],feature_parameters,target_column)

    # Split the data set into training and test
    df_features_train, df_features_test, df_target_train, df_target_test = split_data(df_features_orignal_train,df_target_train,0.3)

    # Normalize the features using z normalization
    df_features_train_z = normalize_z(df_features_train)

    # Change the features and the target to numpy array using the prepare functions
    X = prepare_feature(df_features_train_z)
    target = prepare_target(df_target_train)

    # Call the gradient descent function
    beta, J_storage = gradient_descent(X, target, beta, alpha, iterations)

    # call the predict() method
    pred = predict(df_features_test,beta)

    target = prepare_target(df_target_test)
    r2 = r2_score(target,pred)
    mse = mean_squared_error(target, pred)
    mae = mean_absolute_error(target, pred)
    # print(f"r^2 value = {r2}, mean squared error = {mse}, mean absolute error = {mae}")
    return r2, mse, mae, beta
```

Initial Model

1. Extract from data our features and targets
2. Split data into training and testing for both features and target
3. Normalize features training data using normalize_z()
4. Prepare training features and target for gradient descent to find out beta values
5. Run gradient_descent() to find optimal beta values
6. Run predict() to get target values and save into variable pred

Generate model matrices

1. Prepare target testing values to be compared with pred variable from above
2. Run r2_score() to calculate the r² value
3. Run mean_squared_error() to calculate the Mean Squared Error

- Instantiate values for alpha, beta and number of iterations
- These numbers are calculated to ensure accuracy of our model

```
In [161]: # Instantiate default values for alpha, beta and number of iterations
alpha = 0.01
iterations = 3300
beta = np.zeros(9,1))

In [162]: # Initial model
r2_value, mse_value, mae_value, beta = make_model(alpha, beta, iterations)
print(f"r2_value = {r2_value}, mse_value = {mse_value}, mae_value = {mae_value}")

r2_value = 0.8519373139537825, mse_value = 11.289873799212793, mae_value = 0.23802453443723483
```

Model Evaluation

To further improve our model, we decided to create a new function that is able to take in parameters to selectively tweak certain parameters to allow us to compare a range of values more efficiently

Here are the features that we used to evaluate our model:

1. Entire Dataset vs Selective Dataset
2. Alpha values
3. Beta values
4. Number of iterations
5. Different features used for training

Include entire data set vs using only specific values

```
In [163]: # Include entire dataset, consisting of all values
r2_value, mse_value, mae_value, beta = make_model(0.01,np.zeros(9,1),3300,0, None)
print(f"r2_value = {r2_value}, mse_value = {mse_value}, mae_value = {mae_value}")

r2_value = 0.43556226845834, mse_value = 45.638099434916836, mae_value = 1.6289599955472684

In [164]: # change in dataset to include only targeted dataset
r2_value, mse_value, mae_value, beta = make_model(0.01,np.zeros(9,1),3300,46, None)
print(f"r2_value = {r2_value}, mse_value = {mse_value}, mae_value = {mae_value}")

r2_value = 0.8519373139537825, mse_value = 11.289873799212793, mae_value = 0.23802453443723483
```

Testing Alpha Values

```
In [165]: # test for changes to alpha value
alpha_range = [0.001, 0.01, 0.05, 0.1]
r2_value_range = []
mse_value_range = []
mae_value_range = []
for i in alpha_range:
    r2_value, mse_value, mae_value, beta = make_model(i,np.zeros(9,1),3300,46, None)
    mse_value_range.append(r2_value)
    mse_value_range.append(mse_value)
    mae_value_range.append(mae_value)

data = pd.DataFrame(list(zip(alpha_range, r2_value_range, mse_value_range, mae_value_range)))
data.columns = ["alpha_range", "r2_value_range", "mse_value_range", "mae_value_range"]
display(data)

sns.lmplot(data = data,x="alpha_range", y="r2_value_range", label = "r^2")
sns.lmplot(data = data,x="alpha_range", y="mse_value_range", label = "Mean Squared Error")
sns.lmplot(data = data,x="alpha_range", y="mae_value_range", label = "Mean Absolute Error")
cs = sns.factorplot(alpha_range, r2_value_range, mse_value_range, mae_value_range)

# r^2 and mse values become significantly more accurate after 0.01 alpha value
```

	alpha_range	r2_value_range	mse_value_range	mae_value_range
0	0.001	0.374238	18.265238	2.17206
1	0.01	0.851937	11.289074	0.230025
2	0.05	0.851984	11.285510	0.230025
3	0.100	0.851984	11.285510	0.230025

```
Out[165]: [Text(0.5, 1.8, 'Model Metrics against alpha value'), Text(0, 0.5, '')]
```

Testing Beta Values

```
In [166]: beta_range = [np.zeros(9,1),np.ones(9,1),np.full((9,1),5),np.full((9,1),30)]
r2_value_range = []
mse_value_range = []
mae_value_range = []
for i in beta_range:
    r2_value, mse_value, mae_value, beta = make_model(0.01,1,3300,46, None)
    mse_value_range.append(r2_value)
    mse_value_range.append(mse_value)
    mae_value_range.append(mae_value)

data = pd.DataFrame(list(zip(beta_range, r2_value_range, mse_value_range, mae_value_range)))
data.columns = ["beta_range", "r2_value_range", "mse_value_range", "mae_value_range"]
display(data)

# starting beta value makes little difference
# mean squared error is the lowest when beta_range is np.zeros(8,1))
```

	beta_range	r2_value_range	mse_value_range	mae_value_range
0	[0. 0. 0. 0. 0. 0. 0. 0. 0.]	0.851937	11.289074	0.230025
1	[1. 1. 1. 1. 1. 1. 1. 1. 1.]	0.851927	11.289851	0.230025
2	[5. 5. 5. 5. 5. 5. 5. 5. 5.]	0.851886	11.28977	0.230025
3	[30. 30. 30. 30. 30. 30. 30. 30. 30.]	0.851834	11.286927	0.230025

Testing Optimal Number of Iterations

```
In [167]: iterations_range = [180,300,590,700,900,1100,1300,1500,1700,1900,2100,2300,2500,2700,2900,3100,3300,3500,3700,3900,4100,4300,4500,4600,4800,5000,5200,5400,5600]
# change in number of iterations
r2_value_range = []
mse_value_range = []
mae_value_range = []
for i in iterations_range:
    r2_value, mse_value, mae_value, beta = make_model(0.01,np.zeros(9,1),1,46, None)
    mse_value_range.append(r2_value)
    mse_value_range.append(mse_value)
    mae_value_range.append(mae_value)

data = pd.DataFrame(list(zip(iterations_range, r2_value_range, mse_value_range, mae_value_range)))
data.columns = ["iterations_range", "r2_value_range", "mse_value_range", "mae_value_range"]
display(data)

sns.lmplot(data = data,x="iterations_range", y="r2_value_range", label = "r^2")
sns.lmplot(data = data,x="iterations_range", y="mse_value_range", label = "Mean Squared Error")
sns.lmplot(data = data,x="iterations_range", y="mae_value_range", label = "Mean Absolute Error")
cs = sns.factorplot(iterations_range, r2_value_range, mse_value_range, mae_value_range)

# r^2 and mse values become significantly more accurate after 300 iterations
# to reach mse value with accuracy of up to 0.0, you will need to reach 3300 iterations
# iterations_range[16] = 3300 is where the value of mse becomes accurate to 2dp, r2 and mae is similar before that
```

	iterations	r2_value_range	mse_value_range	mae_value_range
0	100	0.439631	407.121423	19.534997
1	300	0.697973	23.028140	2.816431
2	500	0.820591	13.068085	0.576550
3	700	0.842572	12.003109	0.276452
4	900	0.847267	11.645144	0.236245
5	1100	0.846994	11.403469	0.238848
6	1300	0.850475	11.400575	0.230136
7	1500	0.851061	11.355872	0.230039
8	1700	0.851395	11.330430	0.230027
9	1900	0.851593	11.315337	0.230025
10	2100	0.851715	11.306026	0.230025
11	2300	0.851793	11.300005	0.230025
12	2500	0.851845	11.296116	0.230025
13	2700	0.851881	11.293299	0.230025
14	2900	0.851906	11.291485	0.230025
15	3100	0.851924	11.290100	0.230025
16	3300	0.851937	11.289074	0.230025
17	3500	0.851947	11.288300	0.230025
18	3700	0.851955	11.287709	0.230025
19	3900	0.851961	11.287251	0.230025
20	4100	0.851966	11.286893	0.230025
21	4300	0.851970	11.286611	0.230025
22	4500	0.851974	11.286295	0.230025
23	4600	0.851976	11.286138	0.230025
24	5000	0.851977	11.286012	0.230025
25	5200	0.851979	11.285912	0.230025
26	5400	0.851980	11.285831	0.230025
27	5600	0.851981	11.285787	0.230025

```
Out[167]: [Text(0.5, 1.8, 'Model Metrics against number of iterations'), Text(0, 0.5, '')]
```

Testing Optimal Features to Consider

```
In [168]: # use different features to decide if it would improve MSE and r^2 values
features_evaluated = ["total_cases","new_cases","new_cases_smoothed","total_deaths","new_deaths","new_deaths_smoothed","total_cases_per_million","new_cases_per_million","new_cases_smoothed_per_million","total_deaths_per_million","new_deaths_per_million","new_deaths_smoothed_per_thousand","positive_rate","tests_per_case","total_cases"]
r2_value_range = []
mse_value_range = []
mae_value_range = []
for i in range(1,len(features_evaluated)):
    # print(features_evaluated[i])
    r2_value, mse_value, mae_value, beta = make_model(0.01,np.zeros((1+1,1)),3300,46, None, features_evaluated[0:i])
    mse_value_range.append(r2_value)
    mse_value_range.append(mse_value)
    mae_value_range.append(mae_value)

# create new array for new col for use change
r2_value_change = ["NIL"]
for value in range(1,len(r2_value_range)):
    prev_value = r2_value_range[value-1]
    current_value = mse_value_range[value]
    difference_r2 = current_value - prev_value
    if difference_r2 < 0.1:
        difference_r2 = "NSR"
        r2_value_change.append(difference_r2)
    else:
        r2_value_change.append("Significant")

# create new array for new col for use change
mse_value_change = ["NIL"]
for value in range(1,len(mse_value_range)):
    prev_value = mse_value_range[value-1]
    current_value = mse_value_range[value]
    difference_mse = current_value - prev_value
    if difference_mse > 0.1:
        difference_mse = "NSR"
        mse_value_change.append(difference_mse)
    else:
        mse_value_change.append("Significant")

# create new array for new col for use change
mae_value_change = ["NIL"]
for value in range(1,len(mae_value_range)):
    prev_value = mae_value_range[value-1]
    current_value = mae_value_range[value]
    difference_mae = current_value - prev_value
    if difference_mae > 0.1:
        difference_mae = "NSR"
        mae_value_change.append(difference_mae)
    else:
        mae_value_change.append("Significant")

data = pd.DataFrame(list(zip(features_evaluated, r2_value_range, mse_value_range, "mae_value_range", r2_value_change, "mse_value_change", mae_value_change)))
data.columns = ["features_evaluated", "r2_value_range", "mse_value_range", "mae_value_range", "r2_value_change", "mse_value_change", "mae_value_change"]
display(data)

# NSR = No Significant Reduction
```

	features_evaluated	r2_value_range	mse_value_range	mae_value_range	r2_value_change	mse_value_difference	mae_value_change
0	total_cases	0.254793	53.669781	0.031019	NSR	NSR	NSR
1	new_cases	0.634051	26.337903	0.031019	0.379757	-27.331278	NSR
2	new_cases_smoothed	0.700002	21.523307	0.031019	NSR	-4.804137	NSR
3	total_deaths	0.763001	16.902118	0.031019	NSR	-4.541249	NSR
4	new_deaths	0.769923	16.990543	0.031019	NSR	NSR	NSR
5	new_deaths_smoothed	0.769989	16.990951	0.031019	NSR	NSR	NSR
6	new_cases_per_million	0.770817	16.484357	0.031019	NSR	NSR	NSR
7	new_deaths_per_million	0.770738	16.500106	0.031019	NSR	NSR	NSR
8	new_cases_smoothed_per_million	0.770948	16.502036	0.031019	NSR	NSR	NSR
9	total_deaths_per_million	0.772648	16.362630	0.031019	NSR	NSR	NSR
10	new_deaths_per_million	0.772957	16.367006	0.031019	NSR	NSR	NSR
11	new_deaths_smoothed_per_million	0.771745	16.427603	0.031019	NSR	NSR	NSR
12	positive_rate	0.776161	16.347761	0.031019	NSR	NSR	NSR
13	new_tests_smoothed_per_thousand	0.842361	11.345302	0.031019	NSR	-4.034674	NSR
14	new_tests_smoothed_per_thousand	0.842794	11.341450	0.031019	NSR	NSR	NSR
15	positive_rate	0.825968	12.544563	0.031019	NSR	NSR	NSR
16	tests_per_case	0.626209	12.520785	0.031019	NSR	NSR	NSR

As shown in the table above, certain features does not result in a significant difference. Having additional features also adds to the computational complexity of the model, thus we remove such features.

Features that result in a significant difference are:

1. total_cases
2. new_cases
3. new_cases_smoothed
4. total_deaths
5. new_deaths_smoothed

In such cases where many features result in novery little reduction in the mean squared error of the model, we do not include it into our model as every feature considered will result in additional computational complexity, which would require a significantly higher amount of resources as the model gets larger.

Moreover, using only these 5 features also resulted in a Mean Absolute Value (0.031) that is lower than our initial model (0.230).

Finally, we decided that it is important that we optimise our model using the least amount of features as only those that will significantly affect our model and to reduce our computational complexity

```
In [170]: r2_value, mse_value, mae_value, beta = make_model(0.01,np.zeros(6,1),3300,46, None, ["total_cases","new_cases","new_cases_smoothed","total_deaths","new_deaths_smoothed","total_cases_per_million","new_cases_per_million","new_cases_smoothed_per_million","total_deaths_per_million","new_deaths_per_million","new_deaths_smoothed_per_thousand","positive_rate",12.03],col_normalized=["reproduction_rate",1.93],col_normalized=["new_tests_smoothed_per_thousand",12.03],col_normalized=["people_fully_vaccinated_per_hundred",51.33],col_normalized=["positive_rate",0.082])

alpha = 0.7842477846629241, mse_value = 15.52777208198358084, mae_value = 0.83101821470374917
```

```
In [172]: r2_value = 0.01
iterations = 3300
beta = np.zeros(6,1))

r2_value, mse_value, mae_value, beta = make_model(alpha, beta, iterations)
print(f"r2_value = {r2_value}, mse_value = {mse_value}, mae_value = {mae_value}")

r2_value = 0.8519373139537825, mse_value = 11.2898
```