
WEATHER IMPACT ON DAILY CASES OF COVID-19 IN SAUDI ARABIA USING MACHINE LEARNING

A PREPRINT

Abdullah Alsuhaibani

Faculty of Computers and Information Systems
Islamic University of Madinah
Madinah, Saudi Arabia
abdullah.alsuhaibani@iu.edu.sa

Abdulrahman Alhaidari

Faculty of Computers and Information Systems
Islamic University of Madinah
Madinah, Saudi Arabia
a.alhaidari@iu.edu.sa

May 10, 2021

ABSTRACT

COVID-19 was announced by the World Health Organisation (WHO) as a global pandemic. The severity of the disease spread is determined by various factors such as the countries' health care capacity and the enforced lockdown. However, it is not clear if a country's climate acts as a contributing factor towards the number of infected cases. This paper aims to examine the relationship between COVID-19 and the weather of 89 cities in Saudi Arabia using machine learning techniques. We compiled and preprocessed data using the official daily report of the Ministry of Health of Saudi Arabia for COVID-19 cases and obtained historical weather data aligned with the reported case daily reports. We preprocessed and prepared the data to be used in models' training and evaluation. Our results show that temperature and wind have the strongest association with the spread of the pandemic. Our main contribution is data collection, preprocessing, and prediction of daily cases. For all tested models, we used cross-validation of K-fold of K=5. Our best model is the random forest that has a Mean Square Error(MSE), Root Mean Square (RMSE), Mean Absolute Error (MAE), and R2 of 97.30, 9.86, 1.85, and 82.3%, respectively.

Keywords Machine Learning · COVID-19 · Weather · Random Forest

1 Introduction

Last year, 2020, is considered to be one of the most challenging years because of the consequences of the disease that started in Wuhan, China, on December 31, 2019. The pandemic hit many countries, and many families lost their beloved ones, and companies around the world have suffered from its consequences. The World Health Organisation (WHO) reported that confirmed cases reached more than 87 million and nearly 2 million deaths on January 10, 2021.[1] It is viewed as a large-scale pandemic that invaded and disrupted all people's activities. The number of cases is still increasing, till the time of writing this paper, on a daily basis, and the mortality rate is rising substantially.

Fever, dry cough, and fatigue are the primary symptoms of COVID-19. Nevertheless, some patients who were diagnosed with COVID-19 show no symptoms, which can expedite the disease infection rate .[2]. The study in [3] shows the correlation between elderly patients and death cases for COVID-19.

On March 2, 2020, the first case was diagnosed and announced by the Saudi Ministry of Health that is located in east providence Qatif City. Saudi Arabia's government has taken some serious actions to prevent the spread of the virus, including the shutdown of borders, schools, the complete ban between regions as well as Umrah suspension [4]. Saudi Arabia is located in southwestern Asia that is closer to Africa. According to the General Authority for Statistics (GASTAT), the total area the country occupied is 2 million square kilometers with a population of 33.413.660 [5] [6]. Saudi Arabia is considered one of the largest countries in the Middle East. Riyadh is the capital of Saudi Arabia, with a total population of more than 8 million, whereas Makkah and Medina combined have a total population of approximately 10.3 million people, and both are the holy places in the Kingdom of Saudi Arabia [7].

Table 1: Daily Cases

City	Daily Cases	Daily Moralities
Makkah	3987	58
Riyadh	2980	6
Medina	2830	32
Jeddah	2738	28
Dammam	957	1
Hufuf	925	3
Taif	241	0
Jubail	211	1
Qatif	200	1
Tabuk	185	1

The weather in Saudi Arabia is mild to hot in all seasons except in some regions where the winter is cold. The goal of this study is to correlate the weather and the spread of COVID-19 in Saudi Arabia. We have applied supervised machine learning algorithms to analyze the behavioral effect of weather on the disease. The aim is to predict the daily cases of COVID-19 in Saudi Arabia, given the climate of different cities. Selecting the best state-of-the-art machine learning models and features impacts on COVID-19 spreading in Saudi Arabia is explored in this paper.

2 Related Work

In the work of [8], the authors examined the relationship between the weather in Boston, United States, and COVID-19. The result indicated that 85% out of the total number of reported cases were recorded in places where the weather is between 3C and 17C. Seven variables were used by [9] in Dhaka, Bangladesh, to correlate the climate and the spread of the pandemic. The authors concluded that only two attributes of the collected data were significant in the analysis, minimum temperature and average temperature. In another study in Canada, the authors used statistical analysis to draw a relationship between the COVID-19 cases and the climate. The result showed that there was no association between high temperature and the increase of the newly recorded cases[10].

Time series prediction was used in the prediction of [11] in different countries. They use a starting point of 30 cases, and then they fitted the consecutive 12 days. They foretasted that the northern hemisphere countries on the globe would have as a result of hot weather and lockdown policies.

A study by [12] investigated 21 different countries in addition to the regions administrated by French. They used publicly available data and concluded that as the temperature decreases, the COVID-19 cases do the same.

The authors in [13] conclude that weather features such as temperature and humidity predicted death rate as well as confirmed cases. KNN and Decision Tree obtained the highest performance models to predict confirmed cases and death cases. In[14] shows that a low wind speed contributes to the increase of the number of infected cases for COVID-19 in Jakarta, Indonesia, by using the Spearman correlation test. The paper in [15] proposed a support vector regression model to predict a COVID-19 prevalence across countries including (Mainland China, US, Italy, South Korea and, India) pandemic comes to an end, and analyses the growth and transmission rates. Using Pearson's correlation, humidity and temperature were important factors for increasing positive cases in both New York and Milan. In Jakarta, Indonesia [16] determined that weather particularly temperature average has correlation with spread of COVID-19 cases where 9999 ($r = 0.392$; $p < .01$)by using Spearman correlation coefficient.

The study in [17] shows the relation between COVID 19 daily cases and weather features in two major cities Riyadh and Makkah. The period for this study was 35 days starts from 8 April 2020 until 13 May 2020. The high temperature results in high daily cases when the temperature between 19 C to 42 C. However, the study was applied to two cities in Saudi Arabia. Long short-term memory (LSTM), a deep learning approach, was trained by [18] to predict total new cases, cured cases, and mortality in Saudi Arabia utilizing official data by the Saudi ministry of health. The result was measured by seven metrics, such as Mean square error (MSE). The model was used to predict the new cases of COVID-19 in six different countries [18].

3 Methodology

3.1 Data Collection

3.1.1 COVID-19 Dataset

The dataset was gathered from the Ministry of Health of the Kingdom of Saudi Arabia’s official website[19]. According to the General Authority for Statistics (GASTAT), the total area the country occupied is 2 million square kilometers with a population of 33.413.660 [5] [6]. At the time where the data was collected for this study, the dataset included 89 cities around Saudi Arabia divided among 13 regions that COVID-19 infected. The data from March 2020 until April 2020 were present in this study. The total number of records is 4860, where features as follows: date, city name, region name, cumulative cases, cumulative recoveries, cumulative mortalities, cumulative active cases, daily cases, daily recoveries, and daily mortalities. As shown in table 1, the total number of reported cases is 15254 based on the daily infections for the top 10 cities is nearly 93% of the total daily cases. Mortality in Makkah was the highest.

3.1.2 Weather Dataset

The data were obtained from the Open Weather website using their API <https://openweathermap.org/>. The retrieval information of weather is based on 13 regions along with Jeddah and Taif. Population density and climatic weather lead us to handle them as regions. The start date was on the same period of the study and came with 26 features. However, the needed features were the following: feels like, humidity, pressure, temperature, temperature by hours, maximum temperature, minimum temperature, weather condition, wind degree, wind speed, and weather main. Other features were ignored because either empty values such as sea level, snow, and rain or ineffective for our experiments like timezone and ID.

3.2 Methods

In order to establish a relationship between the pandemic and the daily new cases, we used different machine learning algorithms. These were trained to detect the pattern and the trend by machine learning models. We used python and Sklearn that allow us to train the preprocessed data and used these models on unseen, new data. Cross-validation was utilized, meaning that our data is divided into training, validation, and testing. Then, the best subset of the data is used for training the model, based on the performance of the validation. Next, the model is tested on separate testing data. The utilized algorithms are summarized in the result section. For evaluation, each model was trained separately, and the main goal is to select the best performing model on the collected data. Aiming to increase the models’ ability to detect daily case changes as the weather changing, predicting daily cases was feasible. Since the target variable we are looking at is regression, we focused on the Mean Square Error (MSE), which is the primary metric for evaluation.

For the prediction, the used features are in the table 2

Table 2: Features used in the models

Feature Name	Excluded?	Feature Name	Excluded?
Cumulative active cases	Yes	Temperature by hours	No
Cumulative cases	No	Temperature maximum	No
Cumulative mortalities	No	Temperature minmium	No
Cumulative recoveries	No	Weather main (i.e. rain, or cloudy)	Yes
Daily mortalities	No	Wind degree	Yes
Daily recoveries	No	Wind speed	Yes
Feels like	Yes	City	Yes
Humidity	Yes	Dates	No
Pressure	Yes		
Temperature	Yes		

For using features, it is determined by one of two possibilities. First, if the feature does not significantly impact the prediction, or we thought it would leakage data, we exclude it from the projection. The latter is tremendously essential because if the model received a hint about the unseen data, it can not be generalized and may fail to completely new predictions out of the testing data.

In order to discover the most important features that affect the response, Relief, which is an algorithm that measures the sensitive features to response, was used. Also, We used different machine learning algorithms to associated the

COVID-19 pandemic with the weather. After data preprocessing, multiple machine learning models were created. We based our work on the classical machine learning algorithms, as the main goal is finding a relation between COVID-19 spread and weather fluctuations. We include not only temperature but also other weather's properties such as wind speed and degree. The collected data has 16 features for each day in 89 regions of Saudi Arabia as time series. The target for the created models is the daily cases. However, to prevent data leakage, we removed some features to prevent the model from inferring the number of newly infected cases from an existing feature. In our approach, we use time-series as we have each date's infection data. Then, the data items were sorted based on the dates range sequence. Subsequently, the dates were removed after sorting. The reason is that we encoded the dates using one-hot encoding. For instance, if a case was reported in Madinah on March 20, 2020, then only this date will have a value of one, and the proceeding and the succeeding dates will have zeros. Therefore, numerous columns were created. This led to not only minimization of the training time but also a slight improvement of the model performance was sought.

It turned out that not all features are essential when it comes to the final response. For instance, when both daily cases and mortality were included, the model's performance degraded substantially, as we use ranked and sorted features, they have outlined in the table 2. Additionally, cumulative deaths were excluded because we do not want the model to have information related to daily cases, our prediction.

4 Results and Discussion

This section summarises our findings and contributions made. Our findings on COVID-19 and its relation to the disease spread at least hint that as the weather changes, there is a change in the number of daily cases. In the table 3 we report the models that were trained on the weather and COVID-19 data.

Table 3: Models' evaluation

Model	MSE	RMSE	MAE	R^2
Random forest	97.30	9.86	1.85	82.3%
AdaBoost	121.65	11.03	1.95	76.7%
Decision tree	167.14	12.92	2.38	69.2 %
KNN	249.21	15.78	2.80	54.3 %
Neural Network	356.88	18.89	4.41	34.5 %

The results show that the best performing model among all is random forest. It is vivid that the neural network was not able to capture the variability of the training data rather than testing samples. Therefore, it overfit the training; consequently, it has a mediocre performance on the testing set. All features were used in the final prediction except the dates of when the cases were reported. This was determined to reduce the training time, and more importantly, the results of the experiments were not significantly affected by removing the dates. We ranked the features according to their influence on the result. The most significant was wind degree followed by temperature.

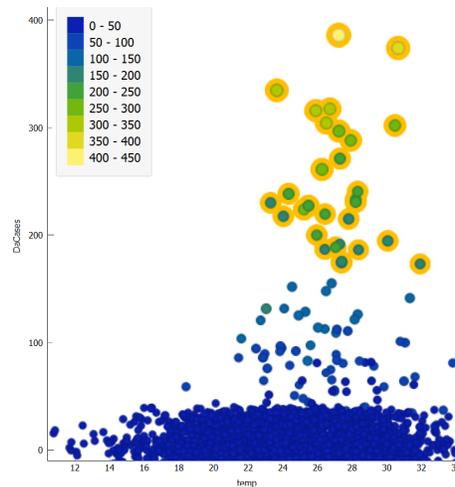


Figure 1: Temperature and daily cases correlation

In figure 4, when the temperature is between 25 and 32 C, we see from the data that the daily cases increase. When it reached 32.5 C, it starts to decrease. Moreover, we see that the more the heat increase does not necessarily affect the reported case. Weather main such as raining is effecting the cases after the temperature.

One concern about the findings was that the data used is from the beginning of the COVID-19 spread in Saudi Arabia, where there are fewer cases than in other countries. It remains unclear that in the extreme cold climate, which is not the case of Saudi Arabia, to which degree weather is attributed to increasing or decreasing in COVID-19 daily cases.

5 Conclusion

The findings of this study can be understood as there is a correlation between weather and the infection of COVID-19. Also, our experiments concluded that the most critical weather feature is temperature, wind speed, feels like, wind degree and humidity. For future research, we suggest using different AI techniques such as deep learning with more diverse data from other parts of the world. Besides, we recommend that using COVID-19 and weather data of a country where the pandemic has been around for several months, and the number of cases is significantly high. This likely makes the correlation clearer as the number of daily cases may significantly change with the weather fluctuation.

References

- [1] Who coronavirus disease (covid-19) dashboard | who coronavirus disease (covid-19) dashboard. <https://covid19.who.int/>. (Accessed on 01/26/2021).
- [2] Hai-Yang Wang, Xue-Lin Li, Zhong-Rui Yan, Xiao-Pei Sun, Jie Han, and Bing-Wei Zhang. Potential neurological symptoms of covid-19. *Therapeutic Advances in Neurological Disorders*, 13:1756286420917830, 2020. PMID: 32284735.
- [3] Fei Zhou, Ting Yu, Ronghui Du, Guohui Fan, Ying Liu, Zhibo Liu, Jie Xiang, Yeming Wang, Bin Song, Xiaoying Gu, et al. Clinical course and risk factors for mortality of adult inpatients with covid-19 in wuhan, china: a retrospective cohort study. *The lancet*, 395(10229):1054–1062, 2020.
- [4] Ismail Anil and Omar Alagha. The impact of covid-19 lockdown on the air quality of eastern province, saudi arabia. *Air Quality, Atmosphere & Health*, 14(1):117–128, 2021.
- [5] General information about the kingdom of saudi arabia | general authority for statistics. <https://www.stats.gov.sa/en/4025>. (Accessed on 02/14/2021).
- [6] Statistical yearbook of 2019 | issue number: 55 | general authority for statistics. <https://www.stats.gov.sa/en/1006>. (Accessed on 02/14/2021).
- [7] Emirate of Makkah region. <https://www.makkah.gov.sa/en>. (Accessed on 05/2/2021).
- [8] Qasim Bukhari, Joseph M Massaro, Ralph B D'Ágostino, and Sheraz Khan. Effects of weather on coronavirus pandemic. *International journal of environmental research and public health*, 17(15):5399, 2020.
- [9] M Mofijur, IM Rizwanul Fattah, ABM Saiful Islam, MN Uddin, SM Ashrafur Rahman, MA Chowdhury, Md Asraful Alam, Md Uddin, et al. Relationship between weather variables and new daily covid-19 cases in dhaka, bangladesh. *Sustainability*, 12(20):8319, 2020.
- [10] Teresa To, Kimball Zhang, Bryan Maguire, Emilie Terebessy, Ivy Fong, Supriya Parikh, and Jingqin Zhu. Correlation of ambient temperature and covid-19 incidence in canada. *Science of the Total Environment*, 750:141484, 2020.
- [11] Alessio Notari. Temperature dependence of covid-19 transmission. *arXiv preprint arXiv:2003.12417*, 2020.
- [12] Jacques Demongeot, Yannis Flet-Berliac, and Hervé Seligmann. Temperature decreases spread parameters of the new covid-19 case dynamics. *Biology*, 9(5):94, 2020.
- [13] Zohair Malki, El-Sayed Atlam, Aboul Ella Hassanien, Guesh Dagnew, Mostafa A Elhosseini, and Ibrahim Gad. Association between weather data and covid-19 pandemic predicting mortality rate: Machine learning approaches. *Chaos, Solitons & Fractals*, 138:110137, 2020.
- [14] Muhammad Rendana. Impact of the wind conditions on covid-19 pandemic: A new insight for direction of the spread of the virus. *Urban climate*, 34:100680, 2020.
- [15] Milind Yadav, Murukessan Perumal, and M Srinivas. Analysis on novel coronavirus (covid-19) using machine learning methods. *Chaos, Solitons & Fractals*, 139:110050, 2020.

- [16] Ramadhan Tosepu, Joko Gunawan, Devi Savitri Effendy, Hariati Lestari, Hartati Bahar, Pitrah Asfian, et al. Correlation between weather and covid-19 pandemic in jakarta, indonesia. *Science of The Total Environment*, page 138436, 2020.
- [17] EA Babeker. Correlation between some climatic factors and covid-19 epidemic in two cities in kingdom of saudi arabia.
- [18] Ammar H Elsheikh, Amal I Saba, Mohamed Abd Elaziz, Songfeng Lu, S Shanmugan, T Muthuramalingam, Ravinder Kumar, Ahmed O Mosleh, FA Essa, and Taher A Shehabeldeen. Deep learning-based forecasting model for covid-19 outbreak in saudi arabia. *Process Safety and Environmental Protection*, 149:223–233, 2020.
- [19] Covid 19 dashboard: Saudi arabia. <https://covid19.moh.gov.sa/>. (Accessed on 02/14/2021).