

# Temperature forecast for Delhi with SARIMA-improved linear model

Feng Gu(T00751197), Yishu Liu(T00728937), Haoran He(T00749480)

December 1, 2024

## Abstract

The present study investigates regional temperature trends in Delhi. By applying time series analysis techniques, historical temperature data were analyzed, and various models, including the dynamic regression model, standard linear model, and other benchmark forecasting models, were constructed and evaluated for their performance on both training and testing datasets. The results indicate that the dynamic regression with dummy variables outperforms other models in predicting future temperatures.

Due to page limitations for formal content, most tables and figures referenced in this report are included in the appendix. Additionally, the complete workflow, a notable aspect of this study, has been made available in the accompanying GitHub repository for transparency and reproducibility.

**GitHut Repo:** <https://github.com/Gufeng-2002/Final-report-for-time-series.git>

## 1 Introduction

The subject of global warming and climate change is gradually becoming one of the significant challenges that the world must face. More frequent and intense extreme weather events, such as heat waves, dust storms, and floods, have been observed globally [1].

The issue of climate change has become particularly crucial in large, densely populated cities. One such example is Delhi, India. The effects of climate change have intensified in recent years, posing challenges to human health, agricultural production, and the environment [2]. Therefore, studying the temperature trends in Delhi holds high scientific value and practical significance. This study employs time series analysis, enabling the examination of historical temperature data, comparison of different time series models, and the prediction of future temperature trends using the best-fit model. This research aims to provide a comprehensive understanding of temperature trends in the Delhi region, and to explore and practice a more efficient way to organize and finish the paper writing work.

## 2 Data

### 2.1 Source of data

From Kaggle<sup>1</sup>, we downloaded our weather data, which is the climate data(of shape (1576,5)) about Delhi of India. Each record in the dataset contains 5 variables: date, mean temperature, humidity, wind speed and mean pressure. The mean temperature is the target variable and the other variables, except date, are predictors.

### 2.2 Preparing and processing the data

To the raw data, we process it by following the procedure below, some additional explanation and corresponding results can be found in appendix:

- Checking the missing values, if there are, replacing or removing the missing records.
- Exploring the distribution for the 4 variables, simple box-plot and hist-plot.
- Replacing the abnormal outliers<sup>2</sup> with corresponding moving average value.
- Creating dummy variables from the "date" variable: four seasons.

### 2.3 ProcessRawData module of Python

It is notable that these steps are finished in a workflow with module "*ProcessRawData.py*" [3], which has been pushed to the public Git repository. It can be easy<sup>3</sup> to repeat all these steps or make further adjustments to make it suitable for other work.

## 3 Method

### 3.1 Specifying the desired model

Before we set a specific model for forecasting *meantemp*, we decomposed the *meantemp* using TSL method[4]. Because we have daily climate data, we set the season period as 365, assuming the same day in each year should have the most similar pattern in Temperature<sup>4</sup>.

After observing the possible seasonality and trend, we create a assume the model as following:

$$y_t = \beta_{5 \times 1} X + \beta_{3 \times 1} H + \eta_t$$

in which:

$$X = \begin{bmatrix} 1 & 1 & X_{11} & X_{21} & X_{31} \\ 1 & 2 & X_{12} & X_{22} & X_{32} \\ & & \dots & \dots & \\ 1 & t & X_{1t} & X_{2t} & X_{3t} \end{bmatrix} \quad H_i = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ & & \dots \\ 0 & 0 & 0 \end{bmatrix} = \begin{cases} 1, & \text{if the season is } i \\ 0, & \text{otherwise} \end{cases}$$

---

<sup>1</sup>a machine learning community for learners

<sup>2</sup>outliers were detected using customized algorithm, which could be found in the code of Module

<sup>3</sup>only needing to point or change the directory path correctly

<sup>4</sup>But it is not rigorous, because every four-year there is one more day premium and the number of day is not an accurate "365" of interge.

There are totally three  $H_i$  here to avoid multilinearity caused by including intercept. To the  $\eta_t$ , we assume it follows a SARIMA or ARIMA model, specifically:

$$\Phi^P(B^s)\phi^p(B)(1-B^s)^D(1-B)^d\eta_t = \Theta^Q(B^s)\theta^q(B)\epsilon_t$$

where, we set the  $s$  equal to 365(days). The searching for proper order of SARIMA((P,D,Q) and (p,d,q)) and the specific calculating are finished by R language.

## 3.2 Comparisions with other models

In order to assess our model properly, we build five other models in table 3 at the same time and use them fit and predict together.

## 3.3 ModelFitting module of R

To fit these models quickly and easily, we choose R to build these models and do relevant tests on them and visualize the results. Same to the Python module, there is a "*ModlFitting.R*" in the repo. To make the workload less, there are some functions that quickly convert the "table-similar" R object into Latex code of table, that's why there are many tables in the appendix, similar function is defined for saving plot objects quickly.

# 4 Result

The specific settings about models can be found in the R module, here we turn to our attention to the standard linear models and dynamic regression models, especially changes of significance level for the three dummy variables.

From table 4 and 5 in appendix, we found that models with dummy variables are always better in performance than models without that. The four variables: *time*, *season\_Autumn*, *season\_Spring*, *season\_Summer* pass the 0.05 significance level test under the  $H_0$ <sup>5</sup> assumption, however, they do not pass the corresponding tests in dynamic regression models.

To the reason why these variables become not important in dynamic regression, one explanation might be that the influences from these four variables can be captured well by the errors of SARIMA process in the model, and the long-term trend with *time* is also not important to mean temperature, based on the given sample.

Additionally, there is one counterintuitive coefficient: the coefficient for *season\_Summer* is smaller than that of *season\_Spring*, which should not be correct by checking the summary about the average feature value in table 1(appendix). It might indicate that some predictors in our model take the effect from *season\_Summer*, if we remove these interfering factors, the relationship might be shown correctly, or this is the truth of the real word.

Although some coefficients did not pass the significance level test, we can still use the model to forecast, because we are focusing on the relationship between these variables but the future values of target.

According the table 4 and 7, we compared the performance of these models on training data and testing data, finally chose the dynamic regression with dummy variables as the model to represent the forecasts, information about this model is shown in table 1.

---

<sup>5</sup> $H_0$ : the coefficient is value of 0, namely no influence from this variable

Table 1: *Summay about the dynamic regression model. Including the coefficients, tests about residuals from training data, and criteria about performance from testing data. (Note: the dynamic regression model here is called 'sarima\_dummy' in R code and tables in appendix)*

Metric	ME	RMSE	MAE	MPE
dynamic regression	0.5447	3.0829	2.6184	-0.079
	MAPE	ACF1	log_lik	AIC
	12.4737	0.8543	-2369.349	4762.699
Coefficient	Estimate	Std. Error	Statistic	P-value
ar1	0.9898	0.0041	242.1087	0.0000
ma1	-0.0953	0.0298	-3.2015	0.0014
ma2	-0.1798	0.0300	-5.9982	0.0000
humidity	-0.1363	0.0042	-32.4098	0.0000
wind_speed	-0.0291	0.0072	-4.0637	0.0001
meanpressure	-0.0322	0.0076	-4.2461	0.0000
time	0.0021	0.0045	0.4730	0.6363
season_Autumn	0.2608	0.5227	0.4990	0.6179
season_Spring	0.5930	0.5235	1.1326	0.2576
season_Summer	0.4116	0.6098	0.6751	0.4997
intercept	63.9278	8.5701	7.4594	0.0000
Other Metrics	sigma2	log_lik	AICc	BIC
dynamic regression	1.5048	-2369.349	4762.914	4826.149
	lb_stat	lb_pvalue	bp_stat	bp_pvalue
	1.5524	0.2128	1.5492	0.2133

We also visualized the forecasts for all models to make the comparisons more clear and direct in figure 1.

## 5 Discussion

### 5.1 Explanation about the model results

According the regression results, we found that *humidity*, *wind speed* and *mean pressure* have negative effect on mean temperature, with their increases, the temperature decreases. In comparison with the winter, the other seasons have higher mean temperature, even though the coefficient of *Spring* and *Summer* might look counterintuitive, which could be the task for further exploration.

The autoregression and moving average parts show there are strong autocorrelation in the mean temperature variable, which could be explained by standard linear model that considers *time* and *seasons* variables in some extent.

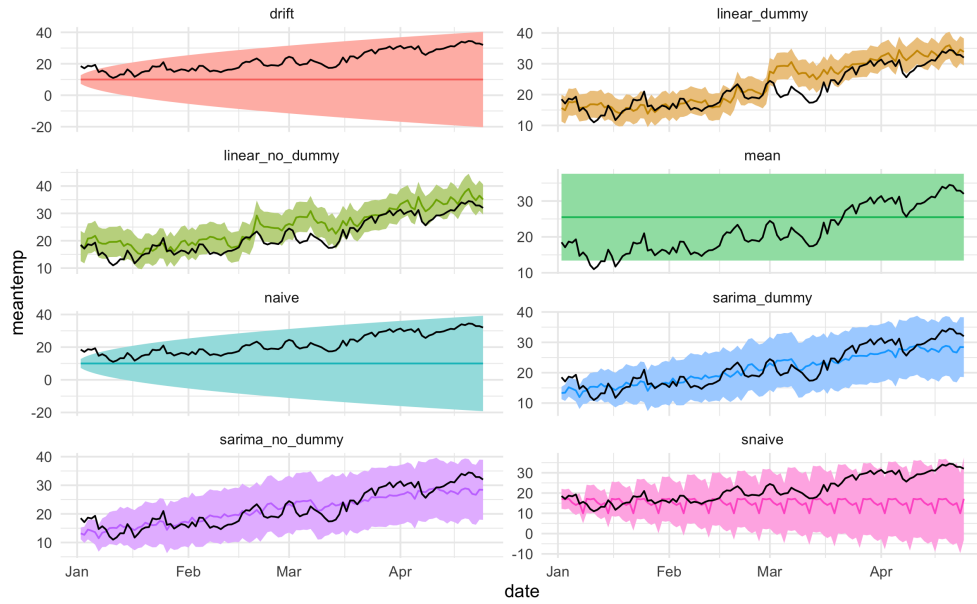


Figure 1: *Forecasts from eight models. Because of the assumptions and settings to models, we should compare the closest forecasts from dynamic regression model with forecasts from other models.*

## 5.2 Other useful work and further improvement.

We have to admit it is not a very rigorous report due to the lack of time and the limits of our skills and professional knowledge in coding and Time Series field.

However, this report is a try in using Vscode, Rstudio, Latex entention as a complete workflow, in which we manage to finish all the work in one system and make the whole process automatic as much as possible. The complete frame work could be found in the GitHub repository, including the document frame of Latex.

To make the whole workflow better, i think we can make improvement with the following aspects:

- Learn Time Series forecasting models more
- Be familiar with R and Python for Data Science
- Be familiar with using Vscode, Rstudio and GitHub for collaboration.

## References

- [1] A. Dabhade, S. Roy, M. S. Moustafa, S. A. Mohamed, R. El Gendy, and S. Barma, “Extreme Weather Event (Cyclone) Detection in India Using Advanced Deep Learning Techniques,” *2021 9th International Conference on Orange Technology (ICOT)*, Tainan, Taiwan, 2021, pp. 1–4, doi: 10.1109/ICOT54518.2021.9680663.
- [2] Hussain S., Hussain E., Saxena P., Sharma A., Thathola P., Sonwani S., “Navigating the impact of climate change in India: a perspective on climate action (SDG13) and sustainable cities and communities (SDG11),” *Frontiers in Sustainable Cities*, 2024 Jan 23. Available from: <https://research-ebsco-com.ezproxy.tru.ca>.
- [3] A. McNeil. “Financial Risk Forecasting: R Best Practice,” *Financial Risk Forecasting Notebook*. Available at: <https://www.financialriskforecasting.com/notebook/R/BestPractice.html>. Accessed: November 30, 2024.
- [4] H. Hyndman and G. Athanasopoulos. “STL Decomposition,” *Forecasting: Principles and Practice (3rd ed.)*. Available at: <https://otexts.com/fpp3/stl.html>. Accessed: November 30, 2024.

## 6 Appendix

### 6.1 Figures

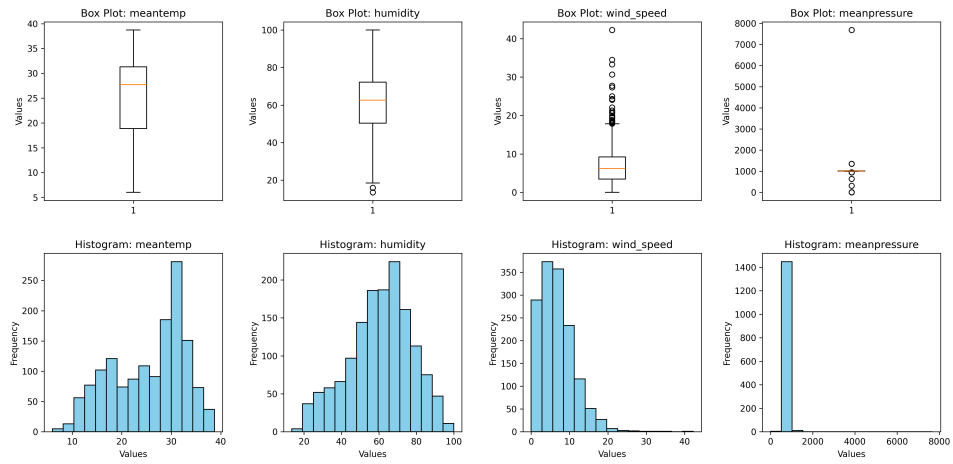


Figure 2: *Distribution of the raw data without replacing outliers*

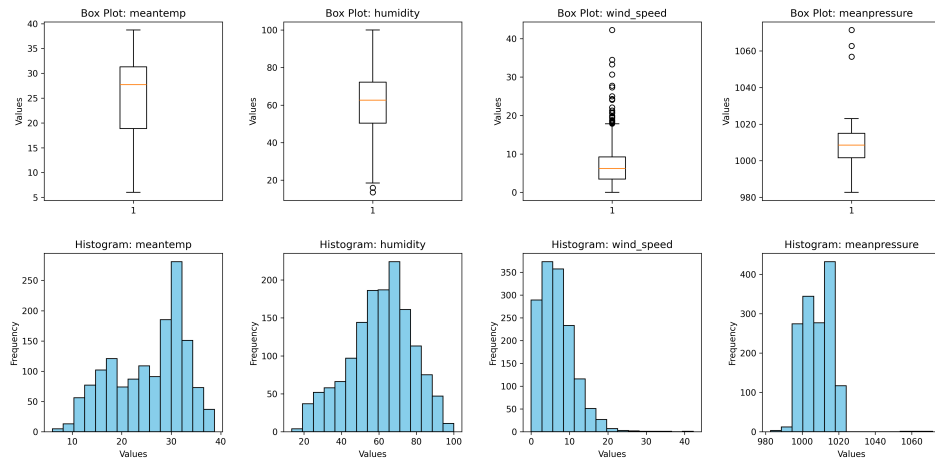


Figure 3: *Distribution of the processed data after replacing outliers*

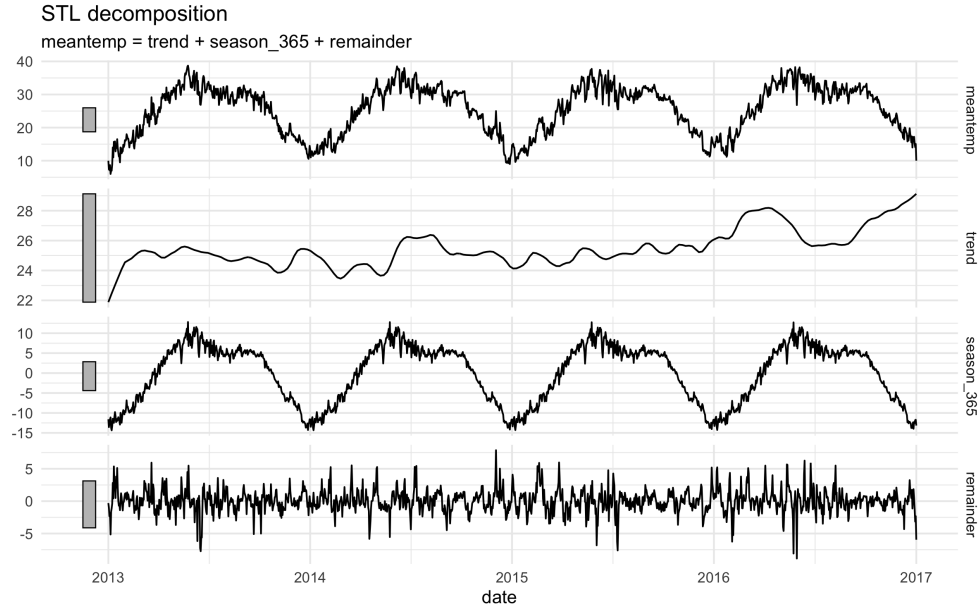


Figure 4: *STL Decomposition of the mean temperature variable (pointing period = 365 days)*

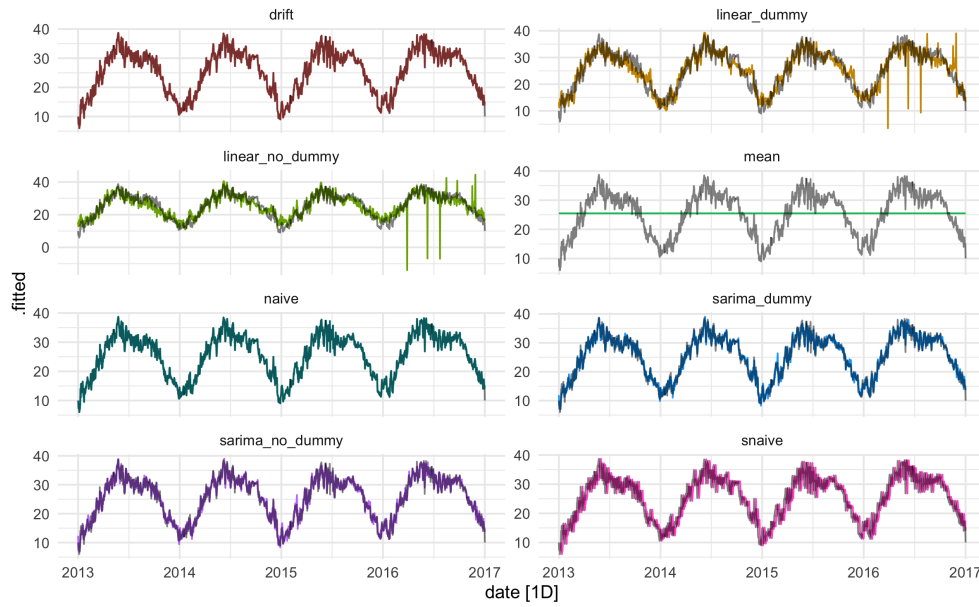


Figure 5: *Fitted values and the true values on training data*



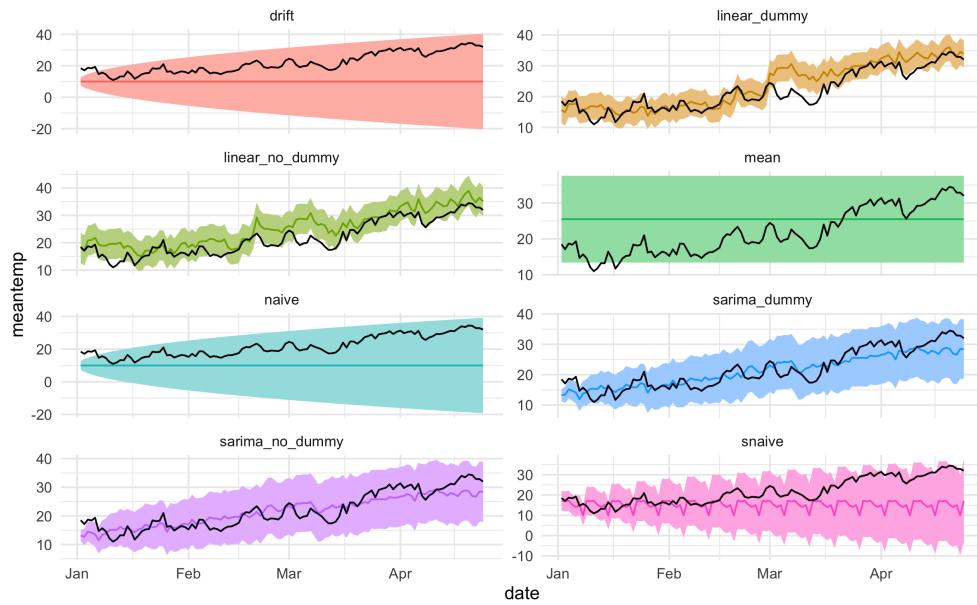


Figure 6: *Forecasts with 90% confidence interval of all models*

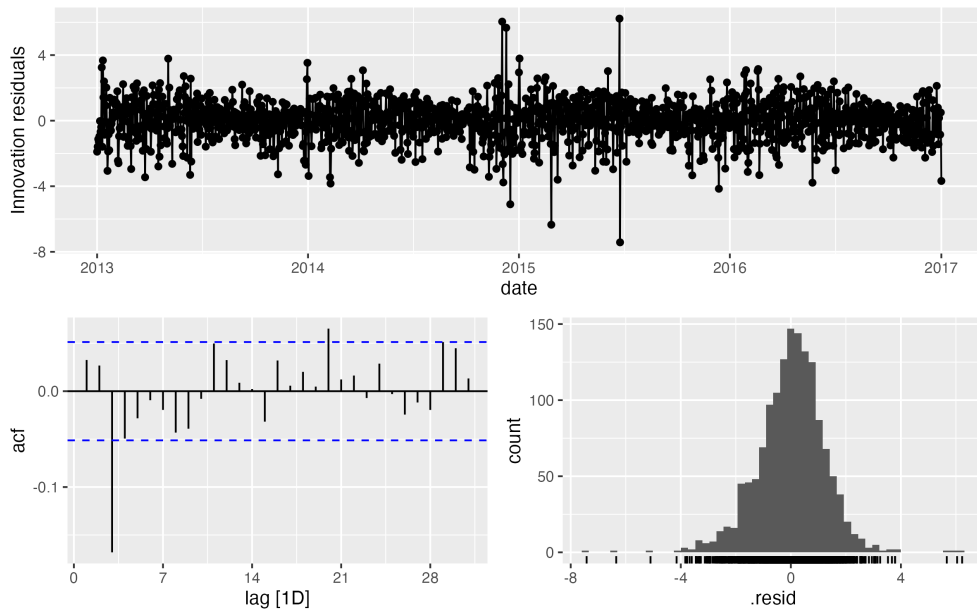


Figure 7: *Residual diagnostic plot for SARIMA with dummy variables*

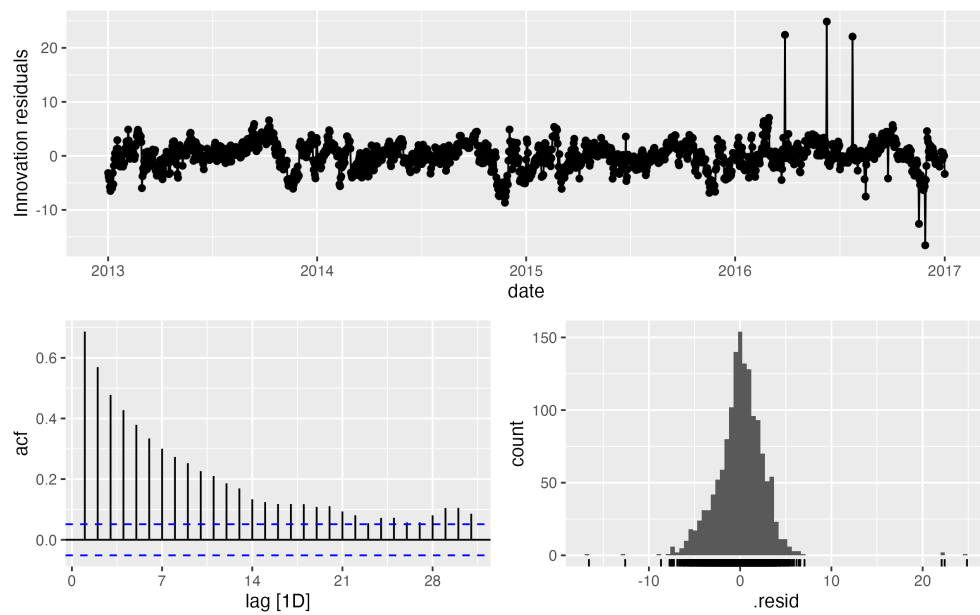


Figure 8: *Residual diagnostic plot for standard linear model with dummy variables*

## 6.2 Tables

Table 2: *Mean features for different seasons*

season	mean_temp	mean_pressure	mean_w_speed	mean_humidity
Autumn	26.0792	1009.7134	5.4029	60.8571
Spring	28.5264	1007.4021	8.4979	45.2249
Summer	31.7559	999.5255	7.8920	64.0544
Winter	15.4633	1016.7316	5.3775	73.1533

Table 3: *Models built in this report*

Basic Model	With Dummy	Without Dummy
Drift	Linear Dummy	Linear No Dummy
Mean	Sarima Dummy	Sarima No Dummy
Naive		
SNaive		

Table 4: *Performance of models*

.model	adj_r_squared	sigma2	log_lik	AICc	BIC	df.residual
linear_no_dummy	0.7928	11.1880	-3837.231	3537.543	3569.210	1457
linear_dummy	0.8709	6.9699	-3489.786	2848.720	2896.184	1454
naive	NA	2.7938	NA	NA	NA	NA
snaive	NA	8.9231	NA	NA	NA	NA
drift	NA	2.7938	NA	NA	NA	NA
mean	NA	53.9946	NA	NA	NA	NA
sarima_dummy	NA	1.5048	-2369.349	4762.914	4826.149	NA
sarima_no_dummy	NA	1.5381	-2387.348	4790.795	4832.996	NA

Table 5: *Comparisons of criteria for forecasting*

.model	RMSE	MAE	MPE	MAPE	ACF1
sarima_dummy	3.0829	2.6184	-0.0790	12.4737	0.8543
sarima_no_dummy	3.1777	2.7123	-1.6050	13.1857	0.8641
linear_dummy	3.6619	2.7407	-9.3986	13.7899	0.8673
linear_no_dummy	4.2402	3.5275	-17.8293	18.4615	0.7798
mean	7.3533	6.5861	-27.4088	36.5698	0.9525
snaive	9.5219	7.3749	24.3618	29.8995	0.8207
drift	13.3623	11.7644	50.0270	50.0270	0.9525
naive	13.3623	11.7644	50.0270	50.0270	0.9525

Table 6: *Tests on residuals from models' fitted values*

.model	kpss_stat	kpss_pvalue	bp_stat	bp_pvalue	lb_stat	lb_pvalue
drift	0.1668	0.1000	37.3479	0.0000	37.4246	0.0000
linear_dummy	0.1640	0.1000	688.8547	0.0000	690.2692	0.0000
linear_no_dummy	0.1899	0.1000	473.1878	0.0000	474.1595	0.0000
mean	0.5774	0.0247	1378.7251	0.0000	1381.5562	0.0000
naive	0.1668	0.1000	37.3479	0.0000	37.4246	0.0000
sarima_dummy	0.1538	0.1000	1.5492	0.2133	1.5524	0.2128
sarima_no_dummy	0.0711	0.1000	0.0115	0.9145	0.0116	0.9144
snaive	0.2894	0.1000	672.6339	0.0000	674.0218	0.0000

Table 7: *Specific coefficients and statistics*

.model	term	estimate	std.error	statistic	p.value
linear_no_dummy	(Intercept)	700.3787	11.8561	59.0733	0.0000
linear_no_dummy	humidity	-0.1567	0.0058	-27.0830	0.0000
linear_no_dummy	wind_speed	-0.0409	0.0211	-1.9380	0.0528
linear_no_dummy	meanpressure	-0.6612	0.0118	-56.0071	0.0000
linear_no_dummy	time	0.0021	0.0002	10.3031	0.0000
linear_dummy	(Intercept)	419.8945	14.5850	28.7894	0.0000
linear_dummy	humidity	-0.1399	0.0056	-24.8596	0.0000
linear_dummy	wind_speed	-0.0106	0.0169	-0.6284	0.5298
linear_dummy	meanpressure	-0.3888	0.0144	-26.9365	0.0000
linear_dummy	time	0.0017	0.0002	10.1253	0.0000
linear_dummy	season_Autumn	5.9208	0.2251	26.3036	0.0000
linear_dummy	season_Spring	5.6261	0.2601	21.6269	0.0000
linear_dummy	season_Summer	8.2651	0.3086	26.7856	0.0000
drift	b	0.0000	0.0437	0.0000	1.0000
sarima_dummy	ar1	0.9898	0.0041	242.1087	0.0000
sarima_dummy	ma1	-0.0953	0.0298	-3.2015	0.0014
sarima_dummy	ma2	-0.1798	0.0300	-5.9982	0.0000
sarima_dummy	humidity	-0.1363	0.0042	-32.4098	0.0000
sarima_dummy	wind_speed	-0.0291	0.0072	-4.0637	0.0001
sarima_dummy	meanpressure	-0.0322	0.0076	-4.2461	0.0000
sarima_dummy	time	0.0021	0.0045	0.4730	0.6363
sarima_dummy	season_Autumn	0.2608	0.5227	0.4990	0.6179
sarima_dummy	season_Spring	0.5930	0.5235	1.1326	0.2576
sarima_dummy	season_Summer	0.4116	0.6098	0.6751	0.4997
sarima_dummy	intercept	63.9278	8.5701	7.4594	0.0000
sarima_no_dummy	ar1	0.9821	0.0054	180.2766	0.0000
sarima_no_dummy	ma1	-0.0234	0.0329	-0.7128	0.4761
sarima_no_dummy	humidity	-0.1355	0.0042	-32.3052	0.0000
sarima_no_dummy	wind_speed	-0.0305	0.0070	-4.3544	0.0000
sarima_no_dummy	meanpressure	-0.0321	0.0075	-4.2714	0.0000
sarima_no_dummy	time	-0.0001	0.0039	-0.0339	0.9730
sarima_no_dummy	intercept	66.1415	8.2258	8.0408	0.0000