

Econometría II: Final Project

Time series modeling with ARIMA
Capturing time invariant correlations

Daniel Emanuel Albert Strauss

ID: 186021

Christian W. McDonald Ehrenberg

ID: 165650

November 30, 2024

Supervisor: Dr. Daniela Cortés Toto

Universidad de las Américas Puebla (UDLAP)

Abstract

Will the sun rise up tomorrow? Many people believe the sun will rise tomorrow since it rose every last day of their lives. Is that a logical conclusion? At least empirically seen, on every past day of their lives, it was more useful to predict the sun would do that thing the next morning than it did all days in their life before that past day. But is it logical to conclude that if a conclusion proved to be useful on every past day of one's life, it would be useful today? To not turn crazy, let's just make that conclusion, independent of whether it is logical or not. Carefully picking some things to assume to be time-invariant leaves us with a new ability: making predictions about the future from observations in the past. For example, one could pick, as such, the stationarity of a time series. If one wants to predict a time series W_t that takes algebraic values, such information about the future may lie in estimated autocorrelations and partial autocorrelations observed from the past. If this series is stationary, the correlations will repeat in the future. Empirically, we can not prove stationarity, but we can observe statistically significant stationarity. How can we use this knowledge to exploit autocorrelations and partial autocorrelations best for future predictions? In the course Econometria II, we learned about the theory and the practical usage of ARIMA models. ARIMA models are a family of models capable of expressing autocorrelation and partial autocorrelation functions (ACF and PACF) of a time series's stationary and invertible differentiation. This project aimed to model and validate a real-life time series with an ARIMA model.

1 Introduction

1.1 Problem to Investigate and Justification

Will the sun rise up tomorrow? Will it be colder in December than in July? At what time will it rain today? All of these phenomena have in common that they can be predicted more precisely from observations in the past. (At least so it seems.) Further examples are: Which stock should I buy? When should I sell my stock? Not only can these examples be predicted in the future from the past, but predicting them also yields practical use. There is a series of mathematical methods to do such. One of them are ARIMA models, at which we will look at this work.

1.2 Objective

The objective of this project is to fit an ARIMA model to a real-life time series we choose by ourselves and to validate the ARIMA model.

To fit the ARIMA models and predict future periods, we will follow this procedure:

- Optionally apply a transformation like BoxCox to ensure constant variance after differentiation
- Evaluate a sufficient amount of differentiations d to make the time-series stationary
- From the differentiated time series, we can extract the ACF and PACF to get an idea of the ARIMA model's necessary dimensions (p, q) . (More on p and q later)
- Use the statsmodels library to fit ARIMA models with different suggestions of (p, d, q) to the undiffer-

entiated data

- Use the fitted models to make predictions of the optionally transformed series
- If a transformation was applied, apply the inverse transformation and correct the bias if needed.

To verify the ARIMA models, they have to fulfill 8 criterions.

1. Residuals have zero mean.
2. Residuals have a constant variance.
3. Residuals are normally distributed.
4. Residuals not correlated.
5. Model parameters significantly non-zero.
6. Models are stationary and invertible
7. Model parameters have low correlation.
8. Atypical data points do not influence the result much.

The first four conditions mean that the residuals must be Gaussian noise. If they weren't, there would still be information left that was not captured by the model, even though the model could capture it, i.e., the residual criteria check if the model captured everything that the model can capture. The conditions on the parameters check the opposite; they check if the model does not contain redundancy. More on that later.

1.3 Research Hypothesis

We want to test the effectiveness of the AR and MA parts of ARIMA models. Therefore, we want to see if ARIMA models with chosen p and q perform better on the verification criteria as well as on the prediction than ARIMA models with poorly chosen p and q . Therefore, we will prosecute the above-described procedure both for well-chosen pairs of p and q as well as for bad-chosen dimensions and, in the end, derive a conclusion from our observations.

2 Methodology

2.1 Dataset

To choose an ARIMA-friendly data set, we used the WWWUsage dataset from [MWH08], which is often used in education examples. For testing and verification purposes, we also added code to create synthetic data, as shown in the code below.

```

1 # Variable para alternar entre datos artificiales y datos reales
2 test_with_artificial = False
3
4 # Funcion para generar datos sintéticos siguiendo un modelo ARIMA
5 def artificial_arima(p=np.array([]), d=0, q=np.array([]), f=lambda x: x, n=100, m=0):
6     """
7     Genera datos sintéticos basados en un modelo ARIMA para validar el método en datos reales.
8
9     Parámetros:
10     - p: Coeficientes del modelo AR (Autoregresivo).
11     - d: Numero de diferencias acumulativas para hacer la serie estacionaria.
12     - q: Coeficientes del modelo MA (Media Movil).
13     - f: Función de transformación aplicada a los datos generados.
14     - n: Número de puntos en la serie temporal.
15     - m: Media del ruido blanco agregado.
16
17     Retorna:
18     - Serie transformada basada en los parámetros proporcionados.
19     """
20     a = np.random.normal(0, 1, n) # Generar ruido blanco con media 0 y varianza 1
21     W = np.zeros(n) # Inicializar la serie temporal
22
23     for t in range(n):
24         if t < len(p) or t < len(q): # Manejar índices fuera de rango
25             W[t] = 0
26         else:
27             # Aplicar componentes AR y MA usando productos escalares
28             W[t] = -W[t-len(p):t] @ p[:-1] + a[t] + a[t-len(q):t] @ q[:-1]
29
30     for d_c in range(d): # Aplicar diferenciación acumulativa d veces
31         W = np.cumsum(W)
32
33     W += m # Agregar media a la serie
34     return f(W) # Aplicar la transformación final
35
36 # Si se activa `test_with_artificial`, generar datos sintéticos
37 if test_with_artificial:
38     n = 1000 # Número de puntos en la serie sintética
39     points = artificial_arima(
40         p=np.array([.5, -.4]), # Coeficientes AR
41         q=np.array([0]), # Coeficientes MA
42         d=1, # Diferenciación
43         f=lambda x: 1.01**x, # Función de transformación exponencial
44         n=n, # Tamaño de la serie
45         m=0 # Media
46     )
47     # Crear un índice de fechas mensual empezando desde 1800
48     dates = pd.date_range(start='1800-01-01', periods=n, freq='MS')
49     # Crear un DataFrame con los datos generados y el índice temporal
50     time_series_df = pd.DataFrame({'points': 1.7**points}, index=dates)

```

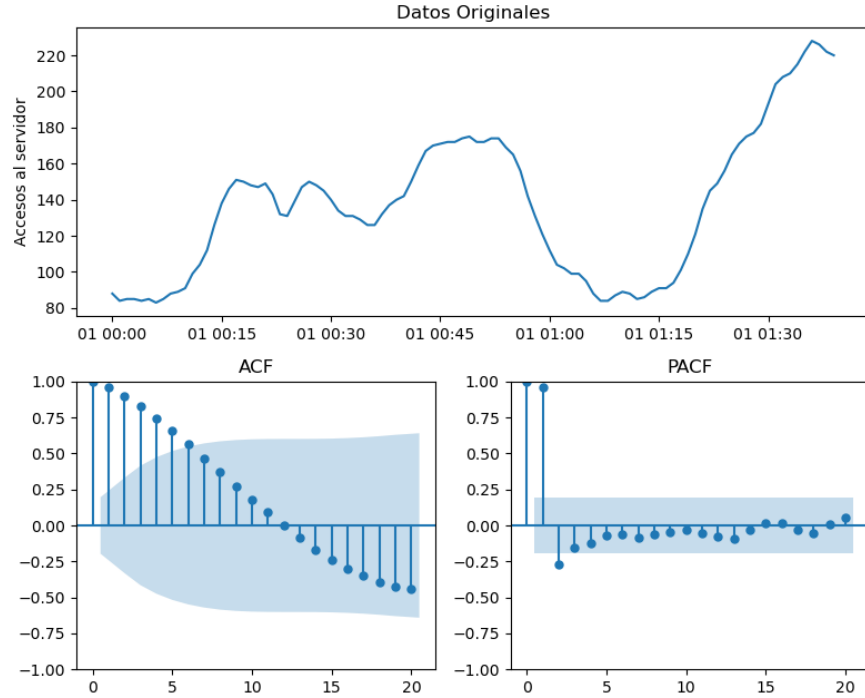


Figure 1: WWWUsage dataset

2.2 Identify a Box-Jenkins Model

2.2.1 Evaluate Stationarity

To evaluate stationarity, we used the statsmodels library. This library provides the Augmented Dickey–Fuller test, which has the null hypothesis, that a unit root is present, and as an alternative hypothesis, that the series is stationary. Since the p-value was roughly 12% we did not conclude stationarity.

```
1 import statsmodels.tsa.stattools as st # Herramientas estadísticas para series temporales
2
3 # Verificar la estacionariedad de la serie temporal
4 # La hipótesis alternativa de la prueba ADF es que la serie es estacionaria
5 adf_result = st.adfuller(time_series_df['points']) # Realizar la prueba ADF
```

2.2.2 Why we did not need a transformation to stabilize the variance

Since the data looked homoscedastic after two differentiations, it seemed like it was not necessary to apply a transformation to stabilize the variance. Furthermore, we tested fitting ARIMA models after applying the box Cox transformation with automated λ parameter and received not as good results when verifying the 8 criteria as without the transformation. Therefore we did not use any transformation. When we tested how the box-cox transformation would turn out, we used the Scipy stats library to automatically find λ .

```
1 from scipy import stats # Importar herramientas estadísticas de SciPy
2
3 # Variable para controlar si se aplica la transformación Box-Cox
4 use_trafo = False
5
6 # Verificar si se aplica la transformación Box-Cox
```

```

7 if use_trafo:
8     # Aplicar la transformación Box-Cox y estimar lambda automáticamente si no se especifica
9     lambda = None # Valor inicial de lambda; si es None, se estima automáticamente
10    if lambda is None:
11        # Aplicar Box-Cox y estimar lambda óptimo
12        transformed_prices, lambda = stats.boxcox(time_series_df["points"])
13    else:
14        # Aplicar Box-Cox con un valor específico de lambda
15        transformed_prices = stats.boxcox(time_series_df["points"], lambda=lambda)
16
17    # Almacenar los valores transformados en el DataFrame
18    time_series_df['T_points'] = transformed_prices
19    print("Lambda usado: ", lambda) # Mostrar el valor de lambda utilizado
20 else:
21     # Si no se aplica la transformación, conservar la serie original
22     lambda = None
23     time_series_df['T_points'] = time_series_df['points']
24
25 # Graficar la serie temporal (transformada o sin transformar)
26 plot_series(time_series_df['T_points'], f"Serie Temporal Transformada con Box-Cox (lambda={lambda})")

```

2.2.3 Differentiate the series

In order to make the series stationary we had to differentiate the series. Nicely, the Pandas' data frame provides a function 'diff()', which applies differentiation. To ensure the minimum necessary amount of differentiations to achieve stationarity, we differentiated the series until the Augmented Dickey-Fuller test significantly ($\alpha = 5\%$) rejected the null hypothesis.

```

1
2 # Nivel de significancia para la prueba ADF
3 alpha = 0.05
4
5 # Seleccionar la serie transformada sin valores nulos
6 current_series = time_series_df['T_points'].dropna()
7 d = 0 # Contador de diferenciaciones aplicadas
8
9 # Función para obtener el p-valor de la prueba ADF
10 def get_adf_p_value(series):
11     """
12     Realiza la prueba Dickey-Fuller Aumentada (ADF) y retorna el p-valor.
13
14     Parámetros:
15     - series: Serie temporal a analizar.
16
17     Retorna:
18     - p-valor de la prueba ADF.
19     """
20     adf_result = st.adfuller(series.dropna()) # Realizar la prueba ADF
21     return adf_result[1] # Retornar el p-valor
22
23 # Iterar hasta que la serie sea estacionaria o se alcance el máximo de diferenciaciones
24 while get_adf_p_value(current_series) >= alpha:
25     d += 1 # Incrementar el contador de diferenciaciones

```

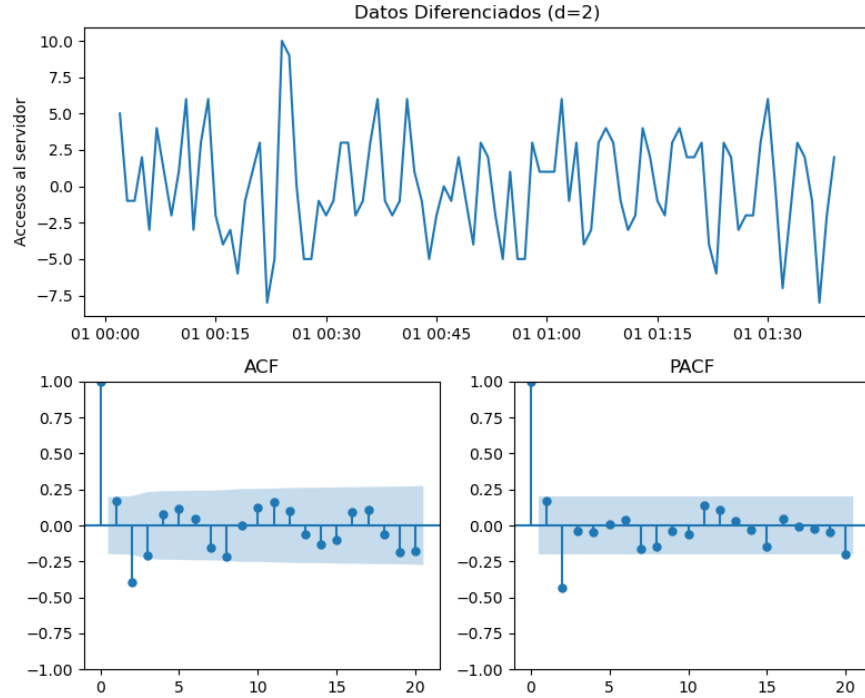


Figure 2: WWWUsage dataset after two differentiations

```

26 # Aplicar diferenciación de primer orden
27 current_series = current_series.diff() # Calcula la diferencia entre valores consecutivos
28
29 # Imprimir el progreso y el p-valor después de cada diferenciación
30 print(f"Después de {d} diferenciación(es), el p-valor de ADF es: {get_adf_p_value(
31     current_series)}")
32
33 # Agregar la serie diferenciada al DataFrame original
34 time_series_df['diff_points'] = current_series.copy()
35
36 # Verificar si la serie es estacionaria después de la diferenciación
37 if get_adf_p_value(current_series) < alpha:
38     print(f"La serie es estacionaria después de {d} diferenciación(es).")
39 else:
40     print("Se alcanzó el máximo de diferenciaciones sin lograr estacionariedad.")

```

We achieved for this dataset, significant stationarity after $d = 2$ differentiations.

2.2.4 Identify significant autocorrelations and partial autocorrelations

The statsmodels library provides methods to obtain the significant partial autocorrelations:

```

1 import statsmodels.tsa.stattools as st
2 # Calculate ACF with confidence intervals
3 acf_values, acf_confint = st.acf(data, alpha=alpha, fft=True, nlags=nlags, adjusted=True)
4 pacf_values, pacf_confint = st.pacf(data, alpha=alpha, nlags=nlags)

```

If the confidence interval of a given lag did not include 0 this would mean, that the corresponding correlation would be significantly non zero.

In the end, we obtained these significant ACF and PACF:

```

1 Significant lags at 5.0% significance level:
2 -----
3 Lag | ACF Value | Confidence Interval
4 -----
5 0 | 1.000 | [ 1.000, 1.000]
6 2 | -0.399 | [-0.603, -0.195]
7 24 | -0.339 | [-0.640, -0.038]
8
9 Significant lags at 5.0% significance level:
10 -----
11 Lag | PACF Value | Confidence Interval
12 -----
13 0 | 1.000 | [ 1.000, 1.000]
14 2 | -0.443 | [-0.641, -0.245]
15 20 | -0.266 | [-0.464, -0.068]
16 23 | -0.258 | [-0.456, -0.061]
17 24 | -0.263 | [-0.461, -0.065]

```

Listing 1: Significant Correlations

2.2.5 Proposal of Box-Jenkins Models

As seen in Table 1, the only significant autocorrelations and partial autocorrelations lie at lag 2. The correlations after lag 20 are at a very big lag and introduce a lot of parameters, which is why we choose not to model these correlations. Since the correlations at lag 1 are not significantly not zero, we are relatively sure that an ARIMA model with $p, q \leq 1$ would not be able to capture these autocorrelations. Thus, we choose, as reasonable models, the models ARIMA(0,2,2), ARIMA(2,2,0), and ARIMA(2,2,2). As unreasonable models (to control the reasonable models), we choose ARIMA(0,2,1), ARIMA(1,2,0), ARIMA(1,2,1), and the white noise model ARIMA(0,2,0).

To fit these models, we used the ARIMA class of statsmodels. Specifically, we used the method 'statsmodels.tsa.arima.model.ARIMA.fit()':

```

1 from statsmodels.tsa.arima.model import ARIMA # Importar el modelo ARIMA
2
3 # Seleccionar la serie transformada y eliminar valores nulos
4 data = time_series_df['T_points'].dropna()
5
6 # Definir los modelos ARIMA sugeridos
7 suggested_models = np.array([
8     [0, d, 0], # Modelo ARIMA(0,d,0)
9     [0, d, 1], # Modelo ARIMA(0,d,1)
10    [1, d, 0], # Modelo ARIMA(1,d,0)
11    [1, d, 1], # Modelo ARIMA(1,d,1)
12    [0, d, 2], # Modelo ARIMA(0,d,2)
13    [2, d, 0], # Modelo ARIMA(2,d,0)
14    [2, d, 2], # Modelo ARIMA(2,d,2)

```

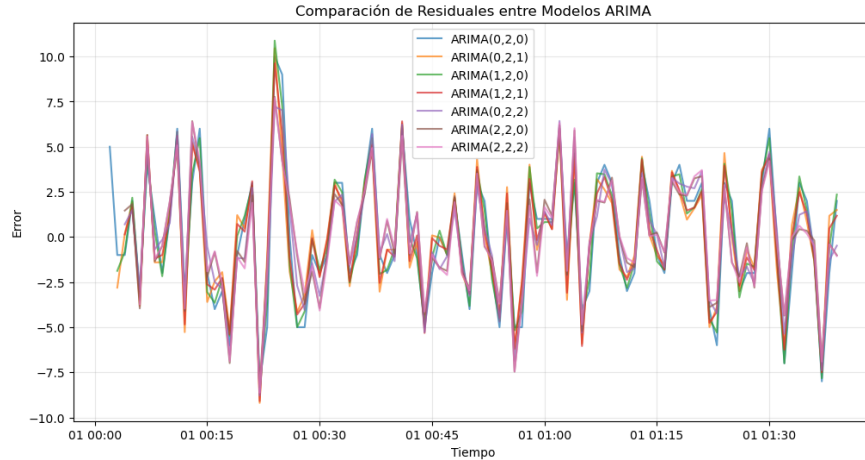



Figure 3: Residuals of the different fitted ARIMA models

```

15 1)
16
17
18 # Iterar sobre los modelos sugeridos y ajustar cada uno
19 for p, d, q in suggested_models:
20     # Ajustar el modelo ARIMA con los parámetros actuales (p, d, q)
21     model = ARIMA(data, order=(p, d, q))
22     fitted = model.fit()

```

2.3 Verification of the Models

2.3.1 Model verification from the residuals

Four of the eight criteria can be extracted from the residuals. They are

1. Residuals have zero mean.
2. Residuals have a constant variance.
3. Residuals are normally distributed.
4. Residuals not correlated.

If one of the criteria of 1,2,4 is violated, it may indicate that there is information that can be captured by box Jenkins models but has not been captured. ARIMA models with high enough dimensions, for example, should be able to capture all autocorrelations, which would have, as a consequence, that the residuals do not have any autocorrelation.

Note that for all of these criteria, there exists a hypothetical time series, that is, ϵ -close to fulfill this criterion, without actually fulfilling it. For example, the mean of the residuals could be $\epsilon > 0$ instead of 0, or the residuals follow a distribution that is almost a normal distribution but not the normal distribution itself.

That's why it is statistically impossible without further assumptions to observe that any of these criteria are fulfilled with a significance level smaller than 1. What we do instead is to try to reject the criteria with the standard significance level of 5%. Given that we have enough data and that the actual properties of the residual are close to the properties in the criteria, we hope that this does not lead to the rejection of the null hypothesis and that the opposite scenario does.

Mean of Zero: To evaluate whether the residuals have a zero mean, we apply a one-sample t-test, which is included in the Scipy stats library. Note that a one-sample t-test assumes normality and independence of residuals. Normality, we will test later, and for independence, we will test if they are not correlated.

```
1 # 1. Prueba de media cercana a 0 (t-test)
2 p_value = ttest_1samp(residuals, 0).pvalue
3 results['mean_close_to_0'] = p_value > alpha # La media está cerca de 0 si p-valor > alpha
4 print(f"{model_name}: Media de residuales = {np.mean(residuals):.4f}, p-valor = {p_value:.4f}")
```

Constant Variance: To evaluate whether the residuals are homoscedastic, we use the Breusch-Pagan test, which is part of the statsmodels.stats.diagnostic package. Its null hypothesis is homoscedasticity, so we opt again for no-reject.

```
1 # 2. Prueba de varianza constante (homocedasticidad) usando Breusch-Pagan
2 _, pvalue, _, _ = het_breuschpagan(residuals, sm.add_constant(np.arange(len(residuals))))
3 results['constant_variance'] = pvalue > alpha # Pasa si p-valor > alpha
4 print(f"{model_name}: Homocedasticidad (p-valor de Breusch-Pagan) = {pvalue:.4f}")
```

Normally Distributed: To evaluate whether the residuals are normally distributed, we consider two tests both implemented in the scipy.stats library. Namely, we use the Shapiro-Wilk and Jarque-Bera Tests, which both have a zero hypothesis for the normality of the residuals.

```
1 # 3. Pruebas de normalidad (Shapiro-Wilk y Jarque-Bera)
2 _, shapiro_pvalue = shapiro(residuals) # Prueba Shapiro-Wilk
3 jb_stat, jb_pvalue = jarque_bera(residuals) # Prueba Jarque-Bera
4 results['normal_distribution'] = shapiro_pvalue > alpha and jb_pvalue > alpha # Ambas deben pasar
5 print(f"{model_name}: Normalidad (p-valor Shapiro-Wilk) = {shapiro_pvalue:.4f}")
6 print(f"{model_name}: Normalidad (p-valor Jarque-Bera) = {jb_pvalue:.4f}")
```

Uncorrelated Residuals: To evaluate whether the residuals are not correlated, we use the Ljung-Box Test, implemented in the statsmodels.stats.diagnostic library. For n lags into the past, it can tell for each of these lags if the according to estimated autocorrelation is significantly not zero. If none of the autocorrelations are significantly not zero, we do not reject the null hypothesis that the residuals are uncorrelated.

```
1 # 4. Prueba de independencia usando Ljung-Box
2 lb_test = acorr_ljungbox(residuals, lags=[10], return_df=True)
3 pvalue_ljungbox = lb_test['lb_pvalue'].values[0]
4 results['independent_errors'] = pvalue_ljungbox > alpha # Pasa si p-valor > alpha
5 print(f"{model_name}: Independencia (p-valor de Ljung-Box) = {pvalue_ljungbox:.4f}")
```

In the end, we obtained these results:

Model	Mean Close to 0	Constant Variance	Normal Distribution	Not Correlated Errors
ARIMA(0,2,0)	Pass	Pass	Pass	Fail
ARIMA(0,2,1)	Pass	Pass	Pass	Fail
ARIMA(1,2,0)	Pass	Pass	Pass	Fail
ARIMA(1,2,1)	Pass	Pass	Pass	Fail
ARIMA(0,2,2)	Pass	Pass	Pass	Pass
ARIMA(2,2,0)	Pass	Pass	Pass	Pass
ARIMA(2,2,2)	Pass	Pass	Pass	Pass

We see that for ARIMA models with $p, q \leq 1$, we see that the residuals have significant autocorrelations. This indicates they did not capture all autocorrelations from the original stationary data. Thus, this was to be expected as both the ACF and the PACF had a peak at 2. We see an example of why it is useful to check the 8 criteria.

2.3.2 Model verification from the parameter estimators

Significance of Parameters: One of the criteria was to check that the estimators for the parameters of the ARIMA model are significantly different from zero. If, for example, a model's highest MA degree parameter would not be significantly different from zero, this would mean that this parameter is redundant and that redundancy could be decreased by lowering the MA dimension. The statsmodels ARIMAResults class provides confidence intervals for the estimated parameters of the ARIMA models. If the confidence interval does not include 0, then the according parameter is significantly different, not zero.

We wrote some code to check for all proposed models if all parameters would be significantly non-zero and then printed which models passed this test and which ones did not:

```

1
2 Resumen de Significancia de Parámetros de Modelos ARIMA Basado en Intervalos de Confianza:
3 =====
4 Modelos que Aprobaron (Intervalos de Confianza NO contienen cero):
5 ['ARIMA(0,2,0)', 'ARIMA(0,2,1)', 'ARIMA(2,2,0)']
6
7 Modelos que Fallaron (Intervalos de Confianza contienen cero):
8 ['ARIMA(1,2,0)', 'ARIMA(1,2,1)', 'ARIMA(0,2,2)', 'ARIMA(2,2,2)']

```

For the model ARIMA(0,2,2) it is important to note that it did not pass this test because the AR(1) parameter was not significantly not zero. The AR(2) parameter was on the other hand. Thus changing the ARIMA(0,2,2) model into an ARIMA(0,2,1) model would remove more than just redundancy as far as I can see it. Thus, one could argue that the ARIMA(0,2,2) model actually passed the test.

Admissible Models: Another criterion that the estimated models have to fulfill is stationarity and invertability. This was easy to check since the Numpy library provides a method "roots" to evaluate the roots of a polynomial. From the moving average and from the autoregressive parameters the according polynomials could be deducted and inserted into numpy.roots. If, for both polynomials, the roots would lie outside the unit circle, the model was admissible, i.e., stationarity and invertible.

In our case, we obtained all modules to be admissible.

```

1
2
3 Resumen de la Prueba de Admisibilidad Basada en Raíces de los Polinomios MA y AR:
4 =====
5 Modelos que Pasaron (Modelos Admisibles con Raíces Fuera del Círculo Unitario):
6 ['ARIMA(0,2,0)', 'ARIMA(0,2,1)', 'ARIMA(1,2,0)', 'ARIMA(1,2,1)', 'ARIMA(0,2,2)', 'ARIMA(2,2,0)', '
   ARIMA(2,2,2)']
7
8 Modelos que Fallaron (Modelos No Admisibles con Raíces Dentro o en el Círculo Unitario):
9 []

```

Correlations of the Parameters: It is better if the estimators for the model parameters have low correlations since then, I believe the model carries less redundancy.

The statsmodels ARIMAResults class offers the covariance matrix of the parameter estimators. Also the statsmodels library offers the function to turn a covariance matrix into a correlation matrix. Thus we calculated the correlation matrixes like this:

```

1 # Extraer la matriz de varianza-covarianza de los parámetros estimados
2 cov_matrix = fitted_model.cov_params() # Matriz de covarianza de los parámetros
3
4 # Convertir la matriz de covarianza a una matriz de correlación
5 correlation_matrix = cov2corr(cov_matrix)

```

We obtained these correlation matrixes (I only print the matrixes for models with more than one parameter):

```

1 Matriz de Correlación de Parámetros para ARIMA(1,2,1):
2           ar.L1      ma.L1      sigma2
3 ar.L1      1.000000 -0.882550  0.084942
4 ma.L1     -0.882550  1.000000 -0.043799
5 sigma2     0.084942 -0.043799  1.000000
6
7 Matriz de Correlación de Parámetros para ARIMA(0,2,2):
8           ma.L1      ma.L2      sigma2
9 ma.L1      1.000000  0.318789  0.135840
10 ma.L2      0.318789  1.000000 -0.095333
11 sigma2     0.135840 -0.095333  1.000000
12
13 Matriz de Correlación de Parámetros para ARIMA(2,2,0):
14           ar.L1      ar.L2      sigma2
15 ar.L1      1.000000 -0.154842  0.235739
16 ar.L2     -0.154842  1.000000 -0.038835
17 sigma2     0.235739 -0.038835  1.000000
18
19 Matriz de Correlación de Parámetros para ARIMA(2,2,2):
20           ar.L1      ar.L2      ma.L1      ma.L2      sigma2
21 ar.L1      1.000000 -0.141192 -0.900061 -0.092983  0.230904
22 ar.L2     -0.141192  1.000000  0.125011 -0.870712  0.121231
23 ma.L1     -0.900061  0.125011  1.000000  0.073469 -0.144525
24 ma.L2     -0.092983 -0.870712  0.073469  1.000000 -0.216743
25 sigma2     0.230904  0.121231 -0.144525 -0.216743  1.000000

```

For the ARIMA(2,2,0) models, we can see a lower correlation of the parameters as for the ARIMA(0,2,2) model. In the ARIMA(2,2,2) model we see low correlations in general but high correlations between the AR(1) and MA(1) and the AR(2) and MA(2) part, indicating that both an autoregressive as well as a moving average part may not be necessary.

2.3.3 Criterion 8

As visible in Figure 1, our dataset does not contain any noticeable outliers. Therefore no further analysis is needed.

2.3.4 Which model to choose?

From our proposed models ARIMA(0,2,2), ARIMA(2,2,0), and ARIMA(2,2,2), the model ARIMA(2,2,0) looked the best with respect to the 8 criteria. Whilst all models passed all criteria for the residuals and all models are admissible, ARIMA(2,2,0) has the lowest correlations of the parameters and it was the only model, that just had parameters that were significantly not zero. Therefore we choose the ARIMA(2,2,0) model to presecute predictions.

2.4 Prediction

Here again, the statsmodels class ARIMAResults provides the function predictions to allow us to obtain the predictions of a fitted ARIMA model. If the requested times of the predictions lie within that time training data for the model stems from, then the model will predict one step into the future. For both prediction tasks (predict the last ten datapoints points and predict the next ten data points points), we also plotted the confidence intervals.

```
1 # Realizar pronósticos
2 forecast = fitted_model.get_prediction(start=start, end=end)
3 forecast_mean = forecast.predicted_mean # Valores pronosticados
4 forecast_ci = forecast.conf_int(alpha=1 - confidence_level) # Intervalos de confianza ajustados
```

In figure 4 we see the predictions of our chosen model ARIMA(2,2,0) and the control random model ARIMA(0,2,0) one step into the future. In table 2, we see the results of the forecast compared to the original data of the ARIMA(2,2,0) model, and in table 3, we see the results of the ARIMA(0,2,0) model.

In figure 5 we see the predictions of our chosen model ARIMA(2,2,0) and the control random model ARIMA(0,2,0) 10 steps into the future after the original data has ended. In table 4 we see the results of this forecastthe ARIMA(2,2,0) model and in table 5 we see the results of the ARIMA(0,2,0) model.

Time	Forecast	Original Data	Difference	Lower CI	Upper CI
01:30:00	188.66	193.00	4.34	182.42	194.89
01:31:00	204.23	204.00	-0.23	197.99	210.46
01:32:00	212.36	208.00	-4.36	206.12	218.59
01:33:00	210.19	210.00	-0.19	203.96	216.43
01:34:00	214.57	215.00	0.43	208.33	220.81
01:35:00	221.66	222.00	0.34	215.42	227.89

01:36:00	228.19	228.00	-0.19	221.96	234.43
01:37:00	232.86	226.00	-6.86	226.62	239.10
01:38:00	222.38	222.00	-0.38	216.14	228.61
01:39:00	221.01	220.00	-1.01	214.77	227.25

Table 2: Prediction of ARIMA(2,2,0) with 95% Confidence Interval (CI) in the last 10 periods for the next period.

Time	Forecast	Original Data	Difference	Lower CI	Upper CI
01:30:00	187.00	193.00	6.00	179.93	194.07
01:31:00	204.00	204.00	-0.00	196.93	211.07
01:32:00	215.00	208.00	-7.00	207.93	222.07
01:33:00	212.00	210.00	-2.00	204.93	219.07
01:34:00	212.00	215.00	3.00	204.93	219.07
01:35:00	220.00	222.00	2.00	212.93	227.07
01:36:00	229.00	228.00	-1.00	221.93	236.07
01:37:00	234.00	226.00	-8.00	226.93	241.07
01:38:00	224.00	222.00	-2.00	216.93	231.07
01:39:00	218.00	220.00	2.00	210.93	225.07

Table 3: Prediction of ARIMA(0,2,0) with 95% Confidence Interval (CI) in the last 10 periods for the next period.

Time	Forecast	Lower CI	Upper CI
01:40:00	219.40	213.16	225.63
01:41:00	218.27	202.87	233.68
01:42:00	216.40	191.48	241.32
01:43:00	214.56	180.10	249.02
01:44:00	213.06	168.23	257.90
01:45:00	211.64	155.16	268.11
01:46:00	210.08	140.95	279.21
01:47:00	208.46	125.94	290.97
01:48:00	206.87	110.30	303.45
01:49:00	205.33	93.97	316.70

Table 4: Forecast for ARIMA(2,2,0) with 95% Confidence Interval 10 periods into the future.

Time	Forecast	Lower CI	Upper CI
01:40:00	218.00	210.93	225.07

01:41:00	216.00	200.20	231.80
01:42:00	214.00	187.56	240.44
01:43:00	212.00	173.29	250.71
01:44:00	210.00	157.59	262.41
01:45:00	208.00	140.59	275.41
01:46:00	206.00	122.39	289.61
01:47:00	204.00	103.07	304.93
01:48:00	202.00	82.70	321.30
01:49:00	200.00	61.34	338.66

Table 5: Forecast for ARIMA(0,2,0) with 95% Confidence Interval 10 periods into the future.

3 Conclusion

In this work, we fitted ARIMA models to the WWWUsage dataset. First, we found out how many differentiations are necessary for stationarity. Then, from the obtained stationary series, we could deduct that either q or p must be at least 2. Then we fitted some ARIMA models with different dimensions. some ARIMA models, fulfilled the criterion of having q or p at the value 2. Others had smaller values, and one was even the model that assumed noise after differentiation, e.g., it was an ARIMA(0,2,0) model. The models with badly chosen parameters indeed performed worse when it came to verifying the 8 criteria of the model. This was prominent when it was about verifying the criteria for the residuals, which indicated that the models with lower dimensions were indeed not able to capture the to-be-captured autocorrelations. This suggests that, on the other hand, the AR or MR components could indeed express relevant parts of the probability distribution of the dataset. However, when comparing the predictions of the ARIMA(0,2,2) model to the predictions of the ARIMA(0,2,0) model, the predicted data points were very close. This suggests that, in this case, the differentiation contributed the most to the prediction and that the MA component was only of secondary matter. It is important to note that the obtained results hold only for this dataset and that one should not be surprised when, from observing other datasets, different results are observed.

References

- [MWH08] Spyros Makridakis, Steven C Wheelwright, and Rob J Hyndman. *Forecasting methods and applications*. John wiley & sons, 2008.

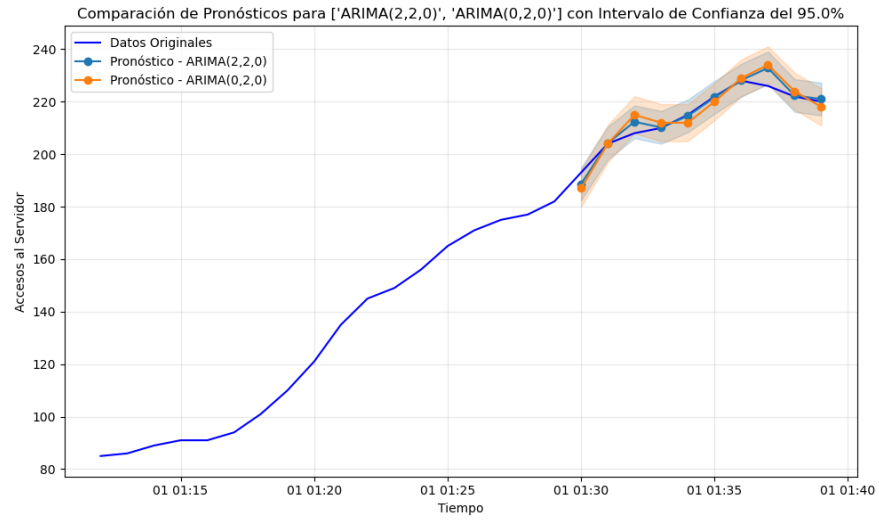


Figure 4: Prediction for ARIMA(2,2,0) and ARIMA(0,2,0) in the last 10 periods for the next period.

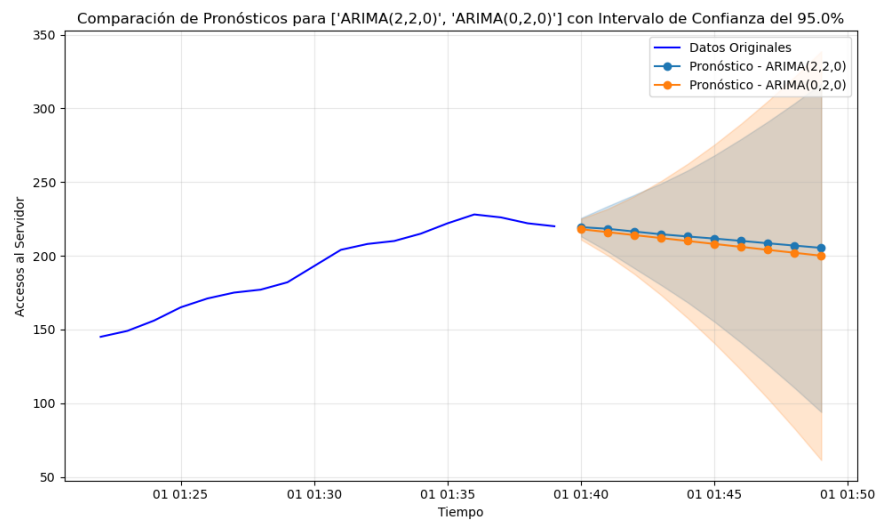


Figure 5: Prognosis ten periods into the future