



ESCUELA SUPERIOR POLITÉCNICA DE CHIMBORAZO

FACULTAD DE CIENCIAS

ESCUELA DE FÍSICA Y MATEMÁTICA

“TÉCNICAS ESTADÍSTICAS PARA LA MODELACIÓN Y PREDICCIÓN DE LA TEMPERATURA Y VELOCIDAD DE VIENTO EN LA PROVINCIA DE CHIMBORAZO”

TRABAJO DE TITULACIÓN

TIPO: Proyecto de Investigación

Presentado para optar al grado académico de:

INGENIERO EN ESTADÍSTICA INFORMÁTICA

AUTORES: PILCO SÁNCHEZ VICTORIA KARINA

ACURIO MARTINEZ WASHINGTON DAVID

DIRECTORA: Ing. Amalia Isabel Escudero Villa.

Riobamba-Ecuador

2019

©2019, Victoria Karina Pilco Sánchez y Washington David Acurio Martínez

Autorizan la reproducción total o parcial, con fines académicos, por cualquier medio o procedimiento, incluyendo la cita bibliográfica del documento, siempre y cuando se reconozca el Derecho de Autor.

ESCUELA SUPERIOR POLITÉCNICA DE CHIMBORAZO
FACULTA DE CIENCIAS
ESCUELA DE FÍSICA Y MATEMÁTICA

El tribunal del trabajo de titulación certifica que: El trabajo de investigación: "TÉCNICAS ESTADÍSTICAS PARA LA MODELACIÓN Y PREDICCIÓN DE LA TEMPERATURA Y VELOCIDAD DE VIENTO DE LA PROVINCIA DE CHIMBORAZO", de responsabilidad de los señores Victoria Karina Pilco Sánchez y Washington David Acurio Martinez, ha sido minuciosamente revisado por los Miembros del Tribunal del Trabajo de Titulación, quedando autorizada su presentación.

Ing. Isabel Escudero


FIRMA

FECHA

09/05/2019

**DIRECTORA DEL TRABAJO
DE TITULACIÓN**

Dr. Arquímides Haro



MIEMBRO DE TRIBUNAL

09/05/2019

Nosotros Victoria Karina Pilco Sánchez y Washington David Acurio Martinez declaramos que el presente trabajo de titulación es de nuestra autoría y que los resultados del mismo son auténticos y originales. Los textos constantes en el documento que provienen de otra fuente están debidamente citados y referenciados.

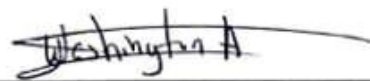
Como autores asumimos la responsabilidad legal y académica de los contenidos de este trabajo de titulación.

Riobamba, 09 de mayo 2019.



Victoria Karina Pilco Sánchez

060501798-7



Washington David Acurio Martinez

020228641-5

DEDICATORIA

El presente trabajo dedico a mis hijos Adriana y Santiago por ser mi fuente de motivación e inspiración para poder superarme cada día más.

A mis padres Héctor y Anita por estar siempre a mi lado brindándome su apoyo y confianza en todo lo necesario para cumplir mis metas profesionales, hacer de mí una mejor persona a través de sus consejos, amor y paciencia, todo lo que hoy soy es gracias a ellos.

A mi esposo Santiago por apoyarme incondicionalmente brindándome su amor, para seguir adelante.

A mi familia que de una u otra manera me brindaron su apoyo y sus palabras de aliento en esta etapa de mi vida y haber culminado con éxito un sueño tan anhelado.

Victoria

DEDICATORIA

El presente trabajo lo dedico con mucho cariño a Dios, que me ayuda cada día, cada minuto de vida porque ha estado siempre presente en cada paso que doy, guiándome, cuidándome y dándome fortaleza para continuar, ya que sin la Fe depositada en él no hubiera alcanzado mi meta; a mis padres, quienes a lo largo de mi vida me han apoyado siempre y han velado por mi bienestar, superación académica y personal siendo mi apoyo en todo momento, depositando su entera confianza en cada reto que se me presentaba sin dudar ni un solo momento en mi capacidad, a mi tío político Segundo Manuel Tamami Llugcha quien motivo en mi la tarea de superarme y de manera especial a todas y cada una de las personas que me apoyaron para lograr esta meta.

Washington

AGRADECIMIENTO

A Dios por ayudarnos a cumplir una meta más en nuestras vidas por su gran bendición amor y bondad, permitiendo que personas de gran corazón aporten con su ayuda incondicional, a nuestros padres por su constante apoyo y preocupación, a las autoridades y docentes de la ESCUELA SUPERIOR POLITÉCNICA DE CHIMBORAZO por darnos la oportunidad de superarnos profesionalmente para así poder servir de mejor manera a la sociedad, aplicando los nuevos conocimientos adquiridos.

Un agradecimiento especial a la Ing. Amalia Isabel Escudero V, Directora del trabajo de titulación y al Dr. Arquímides Haro Asesor, los mismos que nos colaboraron con sus conocimientos, tiempo y paciencia lo cual nos permitió culminar con éxito este trabajo, así como también al Centro de Energía Alternativa y Ambiente “CEAA” que nos brindaron las facilidades para la aplicación práctica de nuestro tema.

TABLA DE CONTENIDO

RESUMEN.....	xviii
SUMMARY	xix
INTRODUCCION	1
CAPÍTULO I.....	3
1 MARCO REFERENCIAL	3
1.1 Antecedentes	3
1.2 Planteamiento del Problema.....	6
1.2.1 Formulación del Problema	7
1.2.2 Sistematización del Problema	7
1.3 Justificación	7
1.3.1 Justificación Teórica	7
1.3.2 Justificación Práctica.....	8
1.4 Objetivos	9
1.4.1 Objetivos Generales	9
1.4.2 Objetivos Específicos.....	9
CAPÍTULO II	10
2 MARCO TEÓRICO REFERENCIAL	10
2.1 Definiciones básicas.....	10
2.1.1 Sensores	12
2.2 Coeficiente de Determinación y R^2 Ajustado	12
2.3 Tipos de Predicción.....	13
2.3.1 Según el horizonte.....	13
2.3.2 Según el tipo de preguntas:	14
2.4 MICE (Multiple Imputation by Chained Equations).....	14
2.4.1 Utilización del paquete MICE.....	14
2.5 Metodología Box-Jenkins (ARIMA)	17
2.5.1 Modelos Autoregresivos AR(p)	17
2.5.1.1 Proceso Autoregresivo de Orden 1: AR (1)	18
2.5.1.2 Características de un modelo AR (1)	19
2.5.2 Proceso de Medias Móviles MA (q)	21
2.5.3 Procesos Autoregresivos de Medias Móviles ARMA (p,q)	23
2.5.3.1 Características de un modelo ARMA (p,q) estacionario:.....	24
2.5.4 Proceso Autoregresivo Integrado y de media móvil ARIMA (p,d,q)	26
2.5.5 Modelos Sarima	27

2.5.6	Pruebas de Raíz Unitaria Estacional	27
2.5.7	Pasos para aplicar modelos Box-Jenkins (ARIMA).....	27
2.5.8	Supuestos	28
2.5.8.1	Prueba de Ljung-Box	28
2.5.8.2	Dickey-Fuller	29
2.5.8.3	Prueba de Kolmogorov Smirnov.....	30
2.5.8.4	Homocedasticidad (Test de Goldfeld Quandt).....	30
2.6	Teoría del Caos	31
2.6.1	Caracterización del Caos.....	32
2.6.2	Sistemas Caóticos.....	33
2.6.3	Tiempo de Retardo.....	33
2.6.4	Dimensión de Encaje.....	34
2.6.5	Reducción de Ruido	36
2.6.6	Predicción.....	38
2.7	Redes Neuronales Artificiales.....	40
2.7.1	¿Qué es una Red Neuronal?	40
2.7.1.1	Elementos básicos que componen una Red Neuronal.....	40
2.7.1.2	Ventajas que ofrecen las Redes Neuronales.....	41
2.7.1.3	Niveles o capas de una Red Neuronal.....	41
2.7.2	Redes Neuronales Recurrentes.....	42
2.7.2.1	Redes parcialmente recurrentes.....	43
2.8	Medidas de evaluación de pronósticos.....	45
2.8.1	Medidas dependientes de la escala.....	45
2.8.1.1	MSE (Mean Square Error)	45
2.8.1.2	RMSE (Root Mean Square Error).....	46
2.8.1.3	MAE (Mean Absolute Error)	47
2.8.1.4	MdAE (Median Absolute Error)	47
2.8.2	Medidas basadas en porcentajes.....	48
2.8.2.1	MAPE (Mean Absolute Percentage Error).....	48
2.8.2.2	MdAPE (Median Absolute Percentage Error)	49
2.8.2.3	RMSPE (Root Mean Square Percentage Error)	50
2.8.2.4	RMdSPE (Root Median Square Percentage Error)	50
2.8.2.5	sMAPE (Symmetric Mean Absolute Percentage Error).....	51
2.8.2.6	sMdAPE (Symmetric Median Absolute Percentage Error)	51
2.9	Criterios de información	52
2.9.1	AIC (Criterio de Información Akaike).....	52
2.9.2	BIC (Criterio de información Bayesiano)	53

2.10	Coeficiente U de Theil	54
2.11	Test de Diebold-Mariano (DM)	54
CAPITULO III.....		56
3	METODOLOGÍA	56
3.1	Tipo y diseño de la investigación.....	56
3.1.1	Descripción del área de estudio.....	56
3.1.2	Tipo de investigación	57
3.1.3	Diseño de investigación	57
3.2	Población de estudio	57
3.3	Recolección de información.....	57
3.4	Identificación de variables	58
3.5	Operacionalización de variables	58
3.6	Análisis de datos	58
3.7	Alcances de la investigación	59
CAPITULO IV.....		60
4	RESULTADOS Y DISCUSIÓN.....	60
4.1	Análisis, interpretación y discusión de resultados	60
4.1.1	Análisis estadístico de las estaciones meteorológicas.....	60
4.1.1.1	Matriz de datos.....	60
4.1.1.2	Análisis exploratorio de datos.....	60
4.1.1.3	Identificación de datos faltantes.....	60
4.1.2	Imputación de datos	63
4.1.3	Estadística Descriptiva.....	67
4.1.4	Modelación Box-Jenkins (ARIMA).....	73
4.1.5	Teoría del Caos	92
4.1.6	Redes Neuronales Recurrentes.....	103
CONCLUSIONES		119
RECOMENDACIONES		121
GLOSARIO		
BIBLIOGRAFÍA		
ANEXOS		

ÍNDICE DE TABLAS

Tabla 1-2: Intervalo de medición de los sensores.	12
Tabla 2-2: Versiones del Paquete MICE.....	16
Tabla 3-2: Técnicas de Imputación univariantes incorporadas.....	16
Tabla 4-2: Resumen de los Patrones de autocorrelación y autocorrelación parcial de los procesos de promedio móvil autoregresivos.	26
Tabla 5-2: Parámetros de la herramienta mutua	34
Tabla 6-2: Parámetros de la herramienta false_nearest.....	35
Tabla 7-2: Parámetros de la herramienta ghkss.	37
Tabla 8-2: Parámetros de la herramienta Rbf.	39
Tabla 1-3: Operacionalización de las variables.....	58
Tabla 1-4: Identificación de datos faltantes.	60
Tabla 2-4: Coeficiente de determinación de las variables imputadas (Atillo) 20104	66
Tabla 3-4: Estadísticas descriptivas de Temperatura y Velocidad de viento de cada estación meteorológica.....	71
Tabla 4-4: Test de Dickey Fuller (Estacionariedad) estaciones meteorológicas	73
Tabla 5-4: Criterios de evaluación para los posibles modelos de Temperatura ARIMA (Atillo).	75
Tabla 6-4: Criterios de información para los posibles modelos de Temperatura ARIMA (Atillo)	76
Tabla 7-4: Valores p de los supuestos del modelo $(1,0,1) (1,1,1)_{24}$ (Atillo).	77
Tabla 8-4: Criterios de evaluación para los posibles modelos de Velocidad de Viento ARIMA (Atillo).....	80
Tabla 9-4: Criterios de información para los posibles modelos de Velocidad de viento ARIMA (Atillo).....	80
Tabla 10-4: Valores p para los supuestos del modelo $(3,0,0) (1,1,1)_{24}$ (Atillo).	82
Tabla 11-4: Resumen de los modelos adecuados para cada Estación Meteorológica (Temperatura y Velocidad de viento) ARIMA.	84
Tabla 12-4: Valores p de los supuestos de cada modelo de las estaciones meteorológicas (Temperatura y Velocidad de Viento).....	87
Tabla 13-4: Criterios de evaluación para los posibles modelos de Temperatura Teoría del Caos (Atillo).....	95
Tabla 14-4: Criterios de evaluación para los posibles modelos de Velocidad de viento Teoría del Caos (Atillo).....	96

Tabla 15-4: Criterios de evaluación de cada uno de los modelos de las diferentes estaciones meteorológicas (Teoría del Caos).	98
Tabla 16-4: Criterios de evaluación de cada uno de los modelos de las diferentes Estaciones Meteorológicas con RNR.....	105
Tabla 17-4: Coeficiente U de Theil de los mejores modelos de las tres técnicas	110
Tabla 18-4: Test de Diebold-Mariano mediante la instrucción “two.sided” de los mejores modelos de las tres técnicas	111
Tabla 19-4: Test de Diebold-Mariano mediante la instrucción “greater” de los mejores modelos de las tres técnicas.....	112

ÍNDICE DE GRÁFICOS

Gráfico 1-2: Coeficientes de autocorrelación y autocorrelación parcial de los modelos AR (1).....	20
Gráfico 2-2: Coeficientes de autocorrelación y autocorrelación parcial de los modelos AR (2).	20
Gráfico 3-2: Coeficientes de autocorrelación y autocorrelación parcial de los modelos MA (1).	22
Gráfico 4-2: Coeficientes de autocorrelación y autocorrelación parcial de los modelos MA(2).	23
Gráfico 5-2: Coeficientes de autocorrelación y autocorrelación parcial de un modelo mixto ARMA(1,1).....	25
Gráfico 1-4: Gráfico de patrón de datos faltantes (Atillo 2014).....	64
Gráfico 2-4: Relación entre datos faltantes de X1 y X14 (Atillo 2014)	64
Gráfico 3-4: Gráfico de datos imputados (Atillo 2014).....	65
Gráfico 4-4: Diagramas de caja de Temperatura y Velocidad de viento para cada estación meteorológica.....	70
Gráfico 5-4: Gráfica de la serie de tiempo de la Temperatura (Atillo)	74
Gráfico 6-4: Diferencias divididas de la Temperatura (Atillo).....	74
Gráfico 7-4: Gráfica de la serie de tiempo de Temperatura eliminada la estacionalidad (Atillo)	75
Gráfico 8-4: Serie y autocorrelogramas de los residuos modelo (1,0,1) (1,1,1) ₂₄	76
Gráfico 9-4: Gráfico de normalidad de los residuos del modelo (1,0,1) (1,1,1) ₂₄	77
Gráfico 10-4: Datos reales vs predichos de Temperatura (Atillo)	78
Gráfico 11-4: Serie y autocorrelogramas de Velocidad de viento (Atillo)	78
Gráfico 12-4: Diferencias divididas de Velocidad de viento (Atillo)	79
Gráfico 13-4: Serie y autocorrelogramas sin estacionalidad de Velocidad de viento (Atillo)...	79
Gráfico 14-4: Gráfico de los errores de los modelos 3 y 5 de Velocidad de viento (Atillo).....	81
Gráfico 15-4: Serie y autocorrelogramas de los residuos del modelo (3,0,0) (1,1,1) ₂₄	81
Gráfico 16-4: Gráfico de normalidad de los residuos de Velocidad de viento (Atillo)	82
Gráfico 17-4: Datos reales vs predichos de de Velocidad de viento (Atillo)	83
Gráfico 18-4: Datos reales vs predichos de los modelos Box-Jenkins (ARIMA) para cada estación meteorológica.....	91
Gráfico 19-4: Atractor extraño de a) Temperatura y b) Velocidad de viento (Atillo).	93
Gráfico 20-4: Tiempo de retardo sin reducción de ruido de a) Temperatura y b) Velocidad de viento (Atillo).....	93

Gráfico 21-4: Dimensión de encaje con reducción de ruido de a) temperatura y b) velocidad de viento (Atillo).....	94
Gráfico 22-4: Datos reales vs predichos de a) Temperatura y b) Velocidad de viento (Atillo).	97
Gráfico 23-4: Datos reales vs predichos de los modelos de Teoría del Caos para cada estación meteorológica.....	102
Gráfico 24-4: Función de error de la red para las variables a) Temperatura y b) Velocidad de viento (Atillo).....	104
Gráfico 25-4: Datos reales vs predichos de a) Temperatura y b) Velocidad de viento (Atillo).	104
Gráfico 26-4: Datos reales vs predichos de los modelos de Redes Neuronales Recurrentes...	109
Gráfico 27-4: Datos reales vs predichos de las tres técnicas para la estación de Atillo.....	114
Gráfico 28-4: Datos reales vs predichos de las tres técnicas para cada estación meteorológica.	118

ÍNDICE DE FIGURAS

Figura 1-2: Ejemplo de una red neuronal totalmente conectada.....	40
Figura 2-2: Comparación entre una neurona biológica (izquierda) y una artificial (derecha)...	41
Figura 3-2: Tipos de Recurrencias: a) Conexión con la misma neurona, b) Conexión con neuronas de la misma capa y c) Conexión con neuronas posteriores y anteriores.....	42
Figura 4-2: Esquema de redes parcialmente recurrentes	43
Figura 5-2: Red neuronal de Elman.....	44
Figura 6-2: Red neuronal de Jordan.....	44
Figura 1-3: Ubicación de las estaciones meteorológicas (Anexo A)	56

ÍNDICE DE ANEXOS

Anexo A: Coordenadas de las estaciones meteorológicas en la provincia de Chimborazo.

Anexo B: Matriz de datos

Anexo C: Imputación de datos en R con el paquete MICE.

Anexo D: Estadísticas Descriptivas en R con el paquete ggplot.

Anexo E: Box-Jenkins (ARIMA) código en R.

Anexo F: Criterios de evaluación y supuestos de cada variable Box-Jenkins (ARIMA).

Anexo G: Gráficos de los errores de cada Estación Meteorológica Box-Jenkins (ARIMA).

Anexo H: Teoría del Caos análisis en TISEAN 3.0.1

Anexo I: Grafica de tiempo de retardo y dimensión de encaje Teoría del Caos en R-Studio.

Anexo J: Redes Neuronales Recurrentes condigo en R-Studio.

Anexo K: Redes Neuronales Recurrentes función del error de cada una de las estaciones meteorológicas.

Anexo L: Código de gráficos de Reales vs Predichos R-Studio.

Tabla 1-A: Estaciones meteorológicas en la provincia de Chimborazo.

Tabla 2-B: Matriz de datos de las estaciones Meteorológicas periodo 2014-2017.

Tabla 3-F: Criterios de evaluación e información de cada modelo Box-Jenkins (ARIMA).

Tabla 4-F: Supuestos para cada modelo de las estaciones meteorológicas Box-Jenkins (ARIMA).

Tabla 5-H: Tiempo de retardo y Dimensión de los mejores modelos.

Tabla 6-L: Criterios de evaluación de cada Red Neuronal Recurrente (Elman y Jordan).

Gráfico 1-G: Gráfico de los errores de los modelos para cada estación Box-Jenkins (ARIMA).

Gráfico 2-H: Ventana del ejecutable de reducción de ruido.

Gráfico 3-H: Ventana del ejecutable de tiempo de retardo.

Gráfico 4-H: Ventana del ejecutable de la dimensión de encaje.

Gráfico 5-H: Ventana del ejecutable para hacer predicciones.

Gráfico 6-H: Atractores extraños para cada estación meteorológica Teoría del Caos.

Gráfico 7-K: Función de error de la red neuronal para todas las estaciones meteorológicas Redes Neuronales Recurrentes.

RESUMEN

El presente trabajo de investigación tiene como objetivo determinar la técnica que proporciona mejores previsiones para modelar variables meteorológicas, período 2014-2017, para lo cual se analizó las variables Temperatura y Velocidad de Viento registradas en las estaciones meteorológicas del Centro de Energías Alternativas y Ambiente (CEAA) de la provincia de Chimborazo, se utilizó las técnicas de: Box-Jenkins (ARIMA), Teoría del Caos y Redes Neuronales Recurrentes con apoyo de softwares como R versión 3.5.1, Tisean 3.0.1 y la hoja de cálculo de Excel. Se realizó la imputación de datos faltantes de aquellas bases de datos que no superaron el 20% y considerando un coeficiente de determinación ajustado mínimo de $R^2=0.8$. Mediante las técnicas: Box-Jenkins se detectó que todas las series de tiempo presentaban estacionalidad cada 24 rezagos, identificando modelos SARIMA y un cumplimiento parcial de supuestos teóricos; con la Teoría del Caos se obtuvieron modelos con mejor ajuste para las primeras 48 horas para la temperatura, presentando una variación para velocidad de viento debido a su inestabilidad. En la modelación obtenida por redes neuronales recurrentes se obtuvo predicciones con un mayor ajuste a las técnicas anteriores y con menos variación en los datos reales versus los predichos. En conclusión, mediante el coeficiente U de Theil y el Test de Diebold-Mariano la metodología Box-Jenkins, Teoría del Caos y Redes Neuronales Recurrentes se obtuvo un U de 0.035, 0.044, 0.020 respectivamente considerando que la primera está sujeta a muchas condiciones, la segunda es adecuada para predicciones a corto plazo y la tercera a largo plazo. El 100% de los test realizados identificaron que Redes Neuronales Recurrentes tienen mayor exactitud y precisión al 95% de confiabilidad.

Palabras Claves: <ESTADÍSTICA>, <MODELACIÓN>, <BOX-JENKINS (TÉCNICA)>, <TEORÍA DEL CAOS>, <REDES NEURONALES RECURRENTE>, <METEOROLOGÍA>



ABSTRACT

The following project work is aimed at determining the technique that provides better predictions to improve weather factors, over the period 2014-2017. Therefore, it has been monitored meteorological variables with respect to temperature and wind speed recorded by weather stations within the Center for Alternative Energy and Environment (CAEE) located in Chimborazo province, there were performed: Box Jenkins (ARIMA) techniques along with Chaos theory and recurring neural networks with version R 3.5.1, Tisean 3.0.1 and excel sheet software support. It was carried the imputation of missing data of those databases not exceeding 20% and taking into consideration a minimum adjustment coefficient $R^2=0.8$. Through the method: Box-Jenkins was detected all the time series models had seasonality every 24 lags, identifying SARIMA models and a partial validity of theoretical assumptions; as of the chaos theory some models were obtained with better adjustment within the first 48 hours of the temperature process, showing wind speed variation due to its instability. In the modeling program was gathered information through recurring neural networks patterns and predictions were collected with a higher adjustment to the previous techniques and less real-data variation against all odds and predictions. In conclusion, through the U-Theil coefficient and Diebold-Mariano test of Box-Jenkins methodology, Chaos theory and neuronal networks a U of 0.035, 0.044, and 0.020 was obtained respectively taking into account the first one is subjected to many conditions, the second one is appropriate for short-term predictions and the third one is considered long-term. The 100% of the tests conducted paved the way to identify that recurring neural networks have a higher rate of 95% precision and accuracy degree of reliability.

KEY WORDS: <STADISTICS>, <MODELLING PROGRAM>, <BOX-JENKINS (METHOD)>, <CHAOS THEORY>, <RECURRING NEURAL NETWORKS>, <METEOROLOGY>



INTRODUCCION

En la actualidad las condiciones atmosféricas que se presentan en un determinado momento y lugar hablan del tiempo atmosférico, sabiendo que este es un gran condicionante para el desarrollo de cualquier actividad. Desde hace mucho tiempo el hombre ha venido estudiando los fenómenos atmosféricos, tratando de entender y explicar las causas que generan estos cambios se basaban en la magia y la religión para explicar la mayor parte de estos fenómenos meteorológicos debido a la inexistencia de tecnología e instrumentos de medición.

La ciencia ha avanzado a pasos agigantados, los estudios realizados con el fin de conocer las condiciones meteorológicas (temperatura, velocidad de viento, presión atmosférica, humedad, radiación solar, etc.) se basan en conocimientos de la física y el uso de alta tecnología, gracias a las cuales, se puede predecir el comportamiento del tiempo atmosférico hasta con una semana de antelación. Por ello es de gran importancia el estudio de la temperatura y la velocidad de viento, por su fuerte relación con el cambio atmosférico, pues, la primera es el grado de calor que la atmosfera posee y el viento es el aire en movimiento que se desplaza de manera horizontal arrastrando nubes e influyendo en la temperatura.

La variación en la temperatura y velocidad de viento hacen difícil la modelación de las mismas puesto que son variables muy complejas, se ha aplicado técnicas como: Box-Jenkins (ARIMA), esta técnica es diferente de la mayoría de métodos porque no supone ningún patrón particular en los datos históricos de las series a pronosticarse; Teoría del Caos ayuda a comprender ciertos tipos de sistemas complejos y dinámicos muy sensibles a las variaciones en las condiciones iniciales, por lo tanto, esta técnica permite deducir el orden subyacente que ocultan fenómenos aleatorios; y Redes Neuronales que crean modelos artificiales para solucionar problemas mediante técnicas algorítmicas convencionales que son de aprendizaje y simulan e imitan sistemas permitiendo crear relaciones no lineales entre capas de entrada y salida.

Se usa estas metodologías con el fin de identificar cuál de ellas presentan un mejor ajuste en el comportamiento de las variables estudiadas, utilizando el coeficiente U de Theil y el test de Diebold-Mariano (DM) para medir la precisión de pronósticos en los modelos encontrados en cada una de las técnicas empleadas. Además, el ultimo test nos ayuda a medir cuál de los modelos comparados presenta mejores pronósticos.

En el capítulo I se menciona estudios donde se han utilizado técnicas de pronósticos como son: ARIMA, Teoría del Caos y Redes Neuronales Artificiales. La importancia del desarrollo del

presente trabajo de investigación es encontrar la mejor técnica de predicción con menor error y mayor confiabilidad posible, debido a la naturaleza inestable que muestran la temperatura y velocidad de viento.

En el capítulo II se realiza una investigación teórica, comenzando desde conceptos básicos como: meteorología, climatología, etc. Se muestra información acerca de cada una de las técnicas y como va hacer evaluados los modelos de cada una respecto a los datos reales basándose en los criterios de evaluación e información, coeficiente de U de Theil y el Test de Diebold-Mariano.

En el capítulo III se muestra la ubicación de cada una de las estaciones meteorológicas en la provincia de Chimborazo y se define cada uno de los ítems de: tipo de investigación, diseño de investigación, población de estudio, recolección de información, identificación de variables, operacionalización de variables, análisis de los datos, alcance de la investigación e imputación de valores perdidos.

En el capítulo IV se analiza los resultados obtenidos luego de haber realizado la imputación de las bases de cada estación meteorológicas de temperatura y velocidad de viento. Posteriormente se analizó los modelos que proporciona cada una de las técnicas, seleccionando el mejor modelo mediante los criterios de evaluación e información. Se analiza las gráficas de los datos reales vs predichos para conocer qué tan bueno es el ajuste de cada técnica.

.

CAPÍTULO I

1 MARCO REFERENCIAL

1.1 Antecedentes

Los primeros estudios para explicar la meteorología bajo una aproximación científica se dieron en la civilización griega con: Tales Mileto (624-545 a.C.), Anaximandro (611-547 a.C.), Hipócrates de Cos (460-375 a.C.) y otros sabios helénicos. El estudio más amplio y difundido de aquellos tiempos corresponden a Aristóteles (384-322 a.C.) quien introdujo el término “Meteorología”, el cual proviene de las palabras griegas Meteoros, “alto en el cielo” y lógica, “Conocimiento, tratado” (Palomares, 2011, p. 1).

La meteorología es el estudio de la atmósfera, sus fundamentos son: la física, matemática, mecánica y química aplicada a la atmósfera; por tanto, su desarrollo ha dependido de los adelantos en estas ciencias. Autores como (Ayllon, 1996), (Pelkowski, 2000) y (Seoanez, 2002) entre otros, mencionan que desde el inicio la humanidad ha querido conocer las causas de los fenómenos meteorológicos y del estado del tiempo, para entender y adaptarse a su entorno (Aragón, 2014, p. 4).

En la actualidad existen múltiples estudios científicos de modelación y predicción con el fin de pronosticar condiciones futuras de fenómenos en diversos campos; entre los más destacados se menciona: ARIMA, Teoría del Caos y Redes Neuronales.

En la base de datos SCIELO en la revista peruana de Biología, en la publicación denominada “Previsión de la temperatura superficial del mar frente a la costa peruana mediante un modelo Autorregresivo Integrado Móvil” (ARIMA), hace referencia al evento del Niño y su conexión global con: el clima, ecosistemas y aspectos socioeconómicos. Las predicciones realizadas desde 1980 hasta 2007 mediante modelos estadísticos y dinámicos no son suficientes, por ello propusieron explorar un ARIMA mediante la identificación, estimación, verificación diagnóstica, previsión y validación del modelo. Utilizaron funciones de autocorrelación simple y parcial (FAC y FACP) para identificar y reformular los órdenes de parámetros, el criterio de información de Akaike (AIC) y de Schwartz (SC) para la verificación diagnóstica, concluyendo que un ARIMA (12,0,11) simuló las condiciones mensuales similares a los datos reales en el litoral peruano, condiciones frías a fines del 2004 y condiciones neutrales a inicios del 2005 (Purca, 2007, pp. 1-2).

En la investigación “Predicción de variables meteorológicas por medio de modelos ARIMA” elaboraron un programa de predicciones de la temperatura, radiación solar, evapotranspiración de referencia y humedad relativa con ARIMA, mediante el cual probaron la efectividad del programa para los pronósticos en condiciones de alta y baja precipitación, se evaluaron los periodos de marzo y junio de 2013 en tres Estaciones Meteorológicas Automáticas (EMAS) del Servicio Meteorológico Nacional (SMN). Donde se concluyó que fue mejor para el periodo con condiciones de baja precipitación; para este trabajo utilizaron MySQL Server y el software R Statistics (Aguado et al., 2016, pp. 1-13).

En Cali Colombia utilizaron la misma metodología para predecir la cantidad de ozono existente en la capa atmosférica, los resultados a corto plazo con el análisis univariante de series de tiempo con 2496 datos registrados en las estaciones de la Red de Monitoreo de Calidad del Aire (RMCA) correspondiente a 104 días consecutivos desde abril a julio de 2003. Los primeros 93 días se utilizaron para la estimación del modelo, y los de los 11 días restantes para su validación, se estimó el modelo con la ayuda del software EViews5 y Matlab 7, en este trabajo investigativo se tomó en cuenta los valores de correlaciones que fueron determinadas como significativas por lo cual se evaluaron cuatro modelos posibles, los mismos se tratan de modelos ARMA (24,4) de los cuales se seleccionó el modelo (4) debido a que este presentaba los mejores indicadores analizándose así lo residuales originados por este y se observó que sus coeficientes son significativamente distintos de cero, para la validación del modelo se trabajaron con 264 datos horarios que corresponden a los 11 períodos de 24 horas, al evaluar este modelo en el período de tiempo propuesto se encontró el error de pronóstico, se observaron errores pequeños lo que permitió considerar como bueno al modelo (Jaramillo et al., 2007, pp. 1-11).

German Poveda Jaramillo en su postgrado de Aprovechamiento de Recursos Hidraulicos de la Universidad Nacional de Colombia (1997)) realizó una investigación sobre “Atractores extraños (caos) en la Hidro-climatología de Colombia” utilizando datos de precipitación mensual de Bogota entre 1866 – 1992 y de Medellín entre 1908 – 1995, además una serie de caudales medios mensuales del río Magdalena en Puerto Berrío. Estimaron el Espectro de Potencias, la dimensión de HausdorffBesikovich, la dimensión de escalamiento entre distancias en los atractores y el mayor exponente de Lyapunov de las series temporales, con la finalidad de estudiar la existencia de componentes determinísticos de baja dimensionalidad y el posible comportamiento caótico. Los resultados obtenidos señalan que es posible la predicción del clima bajo escenarios de cambio global aunque falta mucho por entender acerca de la importancia de las nubes en la dinámica energética de la tierra y de la retroalimentación de los procesos oceano-atmósfera-tierra. En cualquier caso la posible existencia de caos deben jugar un papel fundamental en el terreno de la predicción climática e hidrológica (Poveda, 1997, pp. 2-14).

En la investigación “Modelación y pronóstico del potencial energético del Río Blanco usando la teoría del caos y un método convencional” cuyo objetivo fue modelar y predecir con mayor precisión el potencial energético hídrico a través del estudio de los caudales en L/s que posee el río Blanco en base a los registros almacenados en la estación dos de la Empresa Eléctrica Riobamba S.A. Aplicando los métodos ARIMA y teoría del caos, los datos obtenidos por el Simatic S5 en la central hidroeléctrica Quimiag son caudales promedios diarios de 7 años a partir de enero del 2000 hasta diciembre del 2006 en un vector con el nombre $Y=\text{caudal(L/s)}$, dejando el último mes como periodo de validación, es decir queda reservado para efectos de comprobación de la capacidad predictiva del modelo. Se utilizó el software SPSS 11 para ARIMA y TISEAN 2.1 teoría del caos, la precisión de pronósticos se halló mediante el análisis de los residuos o errores, correlación de Tukey, el período de predicción con correlación significativa es de 14 días con la teoría del caos y 17 días con ARIMA, datos reales hasta un período de 11 días los pronósticos del potencial energético son estadísticamente iguales, a partir del doceavo día la teoría del caos es más precisa que el modelo ARIMA (Escudero, 2007, pp. 65-106).

En la revista Sistemas, Cibernética e informática en la publicación denominada “Predicción de datos meteorológicos en cortos intervalos de tiempo en la ciudad de Riobamba usando la teoría del Caos”, el objetivo fue determinar y predecir los datos meteorológicos, en cortos intervalos de tiempo usando la teoría del Caos, con datos de la estación meteorológica de la ESPOCH del 2010 (Grupo de Energías Alternativas), los mismos que fueron procesados en el software TISEAN. Trabajaron 5 variables: temperatura, humedad, presión, velocidad y precipitación; concluyendo que: la precisión de la predicción depende de la variable analizada, el nivel de caoticidad calculado (valor de entropía) alcanza valores más altos en los parámetros de humedad, disminuyendo paulatinamente las demás variables, lograron determinar una predicción que se correlaciona con los datos reales, excepto para la precipitación en cortos intervalos de tiempo (24h00 a 72h00) (Arquímides et al., 2016, pp 1-7).

En el congreso VIII de Ingeniería de Organización se expone un proyecto denominado “Pronostico de la velocidad y dirección del viento mediante Redes Neuronales Artificiales”, mismo que aborda como aspecto fundamental la importancia del aprovechamiento de una instalación eólica. Realizaron pronósticos con una antelación de 24 a 36 horas la intensidad y dirección de viento mediante modelos ARIMA combinadas con las técnicas de Inteligencia Artificial como lo son las redes Neuronales Artificiales, en el proceso se analizaron a su vez la influencia de medidas anteriores del mismo lugar, así como la posible relación entre medidas en puntos relativamente alejados, con el fin de desarrollar una metodología fiable de predicción del viento para coadyuvar a la gestión de una instalación eólica, evitando futuros problemas de inestabilidad de red eléctrica (Pino et al., 2004, pp. 5-8).

En el paper publicado por Diana Lucía Poma titulado “Predicción Meteorológica mediante Redes Neuronales”, presenta los pronósticos meteorológicos considerando las variables: velocidad de viento, punto de rocío, temperatura, humedad, etc., mediante el método de alimentación hacia adelante con redes neuronales artificiales para el aprendizaje supervisado, para predecir las condiciones climáticas futuras en el Ecuador, mediante el programa de código abierto Weka y código Java da formato a los ficheros de entrada al software Weka, los datos diarios de entrada para el modelo de predicción se extrajeron de la estación meteorológica de Salinas, la misma que pertenece al proyecto Weather Underground, al subir los datos a Weka fueron clasificados mediante el algoritmo Backpropagation, para la evaluación del modelo tomó las mediciones del error, el error cuadrático medio y el error absoluto relativo. Se determinó que Redes Neuronales ofrece adecuados pronósticos para predecir los cambios climáticos mediante la utilización de un bajo número de variables (Poma, 2010, pp. 3-6).

En el trabajo de Juan Alejandro Peña Palacio titulado “Modelo para la predicción de la radiación solar a partir de redes neuronales artificiales”, analizó el comportamiento de las series de tiempo que determinan el valor de la radiación (departamento Australian Government Bureau of Meteorology), en primer lugar determinó el número de datos para el entrenamiento de la red y otras variables climáticas útiles como la temperatura máxima y mínima, y las horas de sol durante el día. Elaboró una base de datos que contenía las entradas de la red neuronal para predecir la radiación solar, incluyendo aquellas variables relacionadas con dicha variable, diseñó un modelo de predicción mediante los principios de la inteligencia computacional y seleccionó funciones de activación comúnmente utilizadas en esos casos, el modelo desarrollado consistió en una red neuronal MADALINE y un perceptrón multicapa para el aprendizaje de las variables climáticas. Se concluyó que las Redes Neuronales es un método factible para pronosticar variables climáticas, además el introducir otras variables meteorológicas asociadas con las variables a predecir ayudan a mejorar significativamente el modelo (González, 2013, pp. 17-20-46).

1.2 Planteamiento del Problema

El comportamiento de las variables meteorológicas como la temperatura y velocidad de viento son muy complejas, dado que no siguen un patrón determinístico, en este trabajo de investigación se propuso tres técnicas de modelación: Box-Jenkins, Teoría del Caos y Redes Neuronales, las mismas que por sus características son adecuadas para modelar sistemas dinámicos.

Se desea conocer el modelo que describe el comportamiento de dichas variables con un mejor ajuste; para realizar predicciones que coadyuven la toma de decisiones en los diferentes campos

que el CEAA y la comunidad en general lo requiera, a más de promover el aprovechamiento de estos recursos naturales como fuente de energía renovable y limpia.

1.2.1 Formulación del Problema

¿Qué modelación presenta un mejor ajuste en la predicción de la temperatura y velocidad del viento en las estaciones meteorológicas del Centro de Energías Alternativas y Ambiente en la Provincia de Chimborazo?

1.2.2 Sistematización del Problema

¿Qué tan fiables son los métodos de modelación para la predicción de la temperatura y velocidad de viento en la provincia de Chimborazo?

¿Qué método presenta mejor ajuste?

1.3 Justificación

1.3.1 Justificación Teórica

Existen diversos métodos de modelación para prever fenómenos naturales, sin embargo, muchos de estos se rigen a supuestos teóricos que difícilmente se ajustan a la realidad, por ello conjuntamente con el avance tecnológico surgen nuevas metodologías como: Box-Jenkins, Teoría del Caos y Redes Neuronales que permiten dar solución a aquellos inconvenientes, con el fin de evitar las múltiples transformaciones de variables perdiendo la autenticidad de la información.

La metodología Box-Jenkins abarca modelos estadísticos que predicen fenómenos mediante la identificación de características constantes (estacionarias). Su objetivo es encontrar un (o varios) modelo(s) simple(s), ajustándose al principio de parsimonia, es decir, modelos con pocos parámetros y con mayor ajuste (García, 2006, p. 5). La teoría del Caos ofrece una explicación para la mayoría de los fenómenos naturales desde el origen del Universo, estudia lo complicado, lo impredecible, lo que no es lineal, presentando propiedades como la sensibilidad a situaciones iniciales, es transitivo y las orbitas repetidas forman un conjunto denso en una región impenetrable del área física, se enfoca en fenómenos caóticos o comportamientos que depende de circunstancias inciertas (Reich, 2009, p. 1). Y las Redes Neuronales representan una metodología de modelación matemática, formadas de una estructura de neuronas unidas por enlaces que transmiten información (aplicando funciones matemáticas) a otras neuronas para entregar un resultado. Aprenden de la información histórica, adquiriendo así la capacidad de predecir

respuestas del mismo fenómeno (Sayago et al. 2011, p. 3), la ventaja de poder utilizar independientemente del cumplimiento de los supuestos teóricos (Pitarque et, al., 2000, p. 1).

En esta investigación se profundizó el estudio de las tres metodologías antes mencionadas para identificar características particulares de cada una de estas, aplicado a variables meteorológicas, con el fin de proporcionar información en cuanto a las diferencias encontradas en dichos modelos y sugerir su aplicación en otras variables meteorológicas en función de los resultados.

1.3.2 Justificación Práctica

Las generaciones pasadas argumentan que las condiciones atmosféricas eran más determinadas, pero con el pasar del tiempo son impredecibles. La meteorología busca predecir el tiempo atmosférico debido a la importancia de comprender la interacción de la atmósfera. El director del Programa Mundial de Alimentos (PMA) en el Ecuador Raphael Chuinard, asegura que mediante estudios realizados, para el año 2050 el cambio climático podrá aumentar el riesgo de hambre y la desnutrición infantil hasta un 20% (1.4 millones de niños más), siendo los agricultores dependientes de estas tierras las más afectadas, en Ecuador representan aproximadamente un millón de agricultores, debido a que los mismos proveen del 70% de los alimentos que los ecuatorianos consumen a escala nacional menciona Roberto Erreis especialista en agricultura y cambio climático. El Instituto Nacional de Meteorología e Hidrología (INAMHI) estima que la temperatura media anual en Ecuador aumenta en 0.8°C siendo las provincias más afectadas son: Cotopaxi, Tungurahua, Cañar y Azuay. Entre el 2000 y el 2010, Ecuador perdió más de 4 billones como consecuencia de la sequía, según el PMA y MAE (Comercio, 2015).

El estudio se enfoca en la modelación y predicción de variables meteorológicas (temperatura y velocidad de viento) con el objetivo de aportar con información confiable al Centro de Energías Alternativas y Ambiente para el desarrollo de sus proyectos de investigación, apoyando a instituciones gubernamentales como: Ministerio de Agricultura, Ganadería, Acuacultura y Pesca (MAGAP), Ministerio Salud Pública (MSP), Secretaria de gestión de Riesgos (SGR), entre otras; los mismos que pueden hacer uso de esta información para planificar sus actividades coadyuvando a la optimización para el uso de recursos naturales y cuidado del medio ambiente en beneficio de la ciudadanía en general.

1.4 Objetivos

1.4.1 Objetivos Generales

Modelar y predecir la temperatura y velocidad de viento en la Provincia de Chimborazo en base a los años 2014 – 2017.

1.4.2 Objetivos Específicos

- Estructurar una base de datos de la temperatura y velocidad de viento del 2014-2017.
- Diseñar el modelo para predecir la temperatura y velocidad del viento con los métodos propuestos.
- Evaluar los modelos de predicción.
- Generar las predicciones de la temperatura y velocidad del viento.

CAPÍTULO II

2 MARCO TEÓRICO REFERENCIAL

2.1 Definiciones básicas

Meteorología: estudia la atmosfera, sus propiedades y los fenómenos que suceden en la misma y que se dan a cada instante, utilizando parámetros como la temperatura, humedad, presión atmosférica, viento o las precipitaciones variando en el espacio y el tiempo. Predecir el tiempo que hará en 24 o 48 horas y en menor medida es el objetivo principal de la meteorología (Rodríguez. et al., 2004, p. 6).

Climatología: estudia las medidas que se registran de aquellos parámetros meteorológicos, y estudia el estado físico medio de la atmosfera y la variación que se da en el tiempo y el espacio. La misma clasifica los distintos tipos de clima que se dan en el planeta, según la localización geográfica y la evolución en el tiempo (Navarra., 2019, p. 3).

Clima: es un conjunto de estados de tiempo atmosférico que se dan en una determinada región (Rodríguez. et al., 2004, p. 61).

Elementos y factores del Clima: los elementos del clima son la combinación de los parámetros temperatura, precipitación, viento, humedad, presión atmosférica y nubosidad, mientras que los factores del clima son agentes que modifican o limitan los elementos del clima dando lugar a diferentes tipos de clima estos son latitud, vientos predominantes, corrientes marinas, distancia al mar, altitud y relieve (Navarra., 2019, p. 2).

Estación Meteorológica: esta debe disponer de varios instrumentos, y para que las medidas sean bien tomadas la ubicación, orientación y condiciones del entorno de los equipos deben estar bajo las normas de la Organización Meteorológica Mundial (Rodríguez. et al., 2004, p. 42).

Variable Meteorológica: es una magnitud meteorológica que alcanza valores numéricos dentro de un conjunto de números especificado.

Temperatura: es una magnitud física la cual se relaciona con la rapidez del movimiento de las partículas que compone, mientras mayor sea la agitación existan entre estas mayores será la temperatura, al ser una magnitud física tiene unidades de medida distintas en la escala que elijamos, escala Celsius (°C), escala Fahrenheit (°F), escala Kelvin (K) (Rodríguez. et al., 2004, p. 15).

Velocidad de Viento: evalúa la componente horizontal del desplazamiento del aire en un sitio y momento explícito, su medida por lo general está definida en metros sobre segundo (m/s) (Navarra., 2019, p. 1).

Predicción Meteorológica: determina los valores anticipados correspondientes a las variables meteorológicas como son: temperatura, presión, humedad, velocidad de viento que afectarán un lugar específico (Rodríguez. et al., 2004, p. 56).

Serie Temporal: es un conjunto de valores ordenados cronológicamente de un fenómeno en función del tiempo.

Tendencia: indica la trayectoria general de los datos en estudio, presenta un movimiento constante y suave a lo largo del tiempo, es decir el crecimiento o decrecimiento de la serie.

- **Ciclos:** el componente de ciclos fluctúa en onda alrededor de la tendencia.
- **Estacionalidad:** este patrón presenta cambios que se repiten año tras año, se asemeja a los ciclos, pero tiene una diferencia fundamental que es el tiempo entre las dos crestas consecutivas: en los ciclos ese tiempo es superior a un año, y en la estacionalidad, es inferior a un año es decir son de corto plazo (Hanke., 2010, pp. 169-178).
- **Series Estacionarias:** este tipo de series presentan una estabilidad a lo largo del tiempo, es decir que la media y su varianza son invariables en el tiempo. Gráficamente se dice que sus valores fluctúan alrededor de la media constante y la variabilidad con respecto a esa media también estable en el tiempo.
- **Series No estacionarias:** en este tipo de series la tendencia y la variabilidad varían en el tiempo, esto me indica que la media establece una tendencia a crecer o decrecer a largo plazo, en conclusión, la serie no fluctúa alrededor de un valor constante.

2.1.1 Sensores

La estación meteorológica automática MAWS100 está configurada con los siguientes sensores:

Tabla 1-2: Intervalo de medición de los sensores.

Nombre e identificador del sensor	Tipo de sensor	Canal de MAWS	Rango de medición	Intervalo de medición
Temperatura del aire	HMP155	CH 7	-80 ... + 60 °C	10 s
Humedad relativa (RH)	HMP155	CH 0	0 ... 100 %	10 s
Presión del aire (PA)	BARO-1	CNT	500 ... 1100 hPa	10 s
Temperatura del suelo 1-7 (TG1-7)	QMT107	CH2	-50 ... + 60 °C	10 s
Radiación solar 1 (SR 1)	CMP6	CH6	0 ... 1500 W/m ²	10 s
Radiación solar 2 (SR 2)	CMP6	CH 5	0 ... 1500 W/m ²	10 s
Voltaje de la batería	QMBATT	CH 4	0 ... 15.0 Vde	10 s

Fuente: Vaisala, 2013.

El sistema también genera una salida de voltaje al relé de control que se utiliza para controlar el estado de reinicio de encendido / apagado de los módulos de módem GPRS.

2.2 Coeficiente de Determinación y R² Ajustado

Se denominan coeficientes de determinación sea este r² en el caso de dos variables o R² (regresión múltiple), se define por:

$$R^2 = \frac{SCE}{SCT} = 1 - \frac{SCR}{SCT} \quad (1.2)$$

Dónde:

SCE es la suma de cuadrados explicados, **SCR** es la suma de cuadrados de los residuos y **SCT** es la suma de cuadrados total.

- El valor de R^2 se define como el coeficiente de determinación muestral múltiple y es la métrica más común para el ajuste de una línea de regresión. Además de ser una estadística descriptiva que halla la proporción de la varianza de la variable endógena explicada por las variables independientes, presenta dos propiedades iniciales de R^2 (Vélez et al., 2016, p.30).
- No toma valores negativos
- Toma valores entre $0 \leq R^2 \leq 1$

Las desventajas son:

- Evalúa bien la bondad dentro del modelo, más no da seguridad que lo haga fuera.
- Al comparar dos o más valores, la variable endógena debe ser la misma.
- R^2 no reduce si se añaden más variables (Vélez et al., 2016, p.30).

2.3 Tipos de Predicción

2.3.1 Según el horizonte

- **Corto Plazo:** Este tipo de pronóstico se efectúa cada mes o menos, y el tiempo de planeación es de un año, comúnmente se usa para abastecimiento, producción, asignación de mano de obra y planificación de los departamentos de fabricación.
- **Mediano Plazo:** El tiempo de validez es de seis meses a tres años, se utiliza para estimar planes de ventas, producción flujos de efectivo y elaboración de presupuestos.
- **Largo Plazo:** este tipo de pronóstico es utilizado en la planificación de nuevas inversiones, lanzamientos de nuevos productos y tendencias tecnológicas de materiales, procesos y productos, así como la elaboración de proyectos, su tiempo de duración es de tres años o más.
- **Longitud del plazo de predicción**
 1. Separar un grupo de datos según la serie antes de la modelación.
 2. Graficar los datos de la serie real y datos predichos.
 3. Utilizar un método de comparación para analizar los errores.
 4. Identificar los datos donde el modelo produce mayor error o deja de ser similar a la serie real.
 5. Definir el número de datos a predecir indicando el error de predicción.

6. Exponer las sugerencias o recomendaciones en cuanto al número de datos que el modelo permite predecir.

2.3.2 Según el tipo de preguntas:

- Resultados de un suceso: ganador de escrutinios, nota de un examen.
- Momento de un hecho: fecha de elecciones, fecha próxima recesión.
- Pronósticos de series temporales: precio acciones en meses cercanos, natalidad en los siguientes 15 años (Beyaert, 2018, pp. 1-2).

2.4 MICE (Multiple Imputation by Chained Equations).

Realiza una imputación múltiple utilizando Fully Conditionally Specification (FCS) implementado por el algoritmo MICE, cada variable posee su propio modelo de imputación, así proporciona modelos de imputación incorporados para datos continuos (pmm), datos binarios (regresión logística), datos categóricos no ordenados (regresión logística polinómica) y datos categóricos ordenados (odds proporcional), se puede utilizar como alternativa la imputación pasiva para la estabilidad entre las variables, dispone de varios gráficos de diagnóstico para examinar la calidad de las imputaciones (Castro M., 2014, p. 16).

2.4.1 Utilización del paquete MICE

Para aplicar los diferentes métodos de imputación múltiple a un conjunto de datos reales, se indican los patrones necesarios para llevar a cabo este proceso, especificar el modelo de imputación es un paso importante en la imputación múltiple, un modelo debe cumplir con lo siguiente:

- Expresar el proceso que creó los datos faltantes.
- Salvaguardar las relaciones en los datos.
- Preservar la incertidumbre sobre las relaciones.

La adhesión a estos principios genera imputaciones adecuadas, dando paso a inferencias estadísticas válidas, se pueden elegir cualquiera de estas opciones para elegir la imputación adecuada:

1. Decidir si el supuesto MAR es estimable.
2. Elegir la forma de imputación así pues la forma incluye la parte estructural y la distribución del error.

3. Dicha opción se enfoca sobre el conjunto de variables que se contienen como predictores en el modelo de imputación, esto puede ser aceptable para pequeños o medianos conjuntos de datos es decir pueden contener de 20 a 30 variables. Para imputar los datos conviene elegir un subconjunto de datos adecuados que no contenga más de 15 a 25 variables para lo cual Van Buuren ayuda con la siguiente estrategia para seleccionar las variables predictoras de una base:
 - a) Contener todas las variables existentes en el modelo de datos completos, esto indica que el modelo será aplicado los datos después de la imputación.
 - b) Incluir las variables que se relacionan con la falta de respuesta, además las razones por las cuales se dan la ocurrencia de datos faltante e pueden encontrar estudiando las correspondencias con el indicador de respuesta de la variable a ser imputadas si dicha correlación excede un cierto nivel, la variable deberá ser incluida.
 - c) Incluir a aquellas variables que presentan una proporción significativa de la varianza, debido que estas ayudan a reducir la incertidumbre de las imputaciones.
 - d) Retirar a aquellas variables seleccionadas en el paso b y c y aquellas variables que presentan demasiados valores faltantes.
- 1) La cuarta decisión es imputar variables que son función de otras variables, siendo útil para agregar las variables transformadas en el algoritmo de imputación múltiple.
- 2) Para que las variables puedan ser imputadas las mismas deben ser ordenadas, la firmeza de visita puede afectar a la tendencia del algoritmo.
- 3) Esta ocupa de la configuración de las imputaciones de partida y el número de iteraciones.
- 4) La séptima elección es m, el número de datos de imputación múltiple, si se asigna un m demasiado pequeño genera grandes errores de simulación e ineficiencia estadística, en especial si la fracción de pérdida de información es alta (Castro M., 2014, pp. 18-20).

En la actualidad existen varios paquetes del software libre R disponibles en la página web cran.r-project.org/.

El MICE en R imputa los datos multivariante incompletos mediante ecuaciones encadenadas.

Tabla 2-2: Versiones del Paquete MICE

Versión	Año	Características
MICE 1.0	2000 2001	<ul style="list-style-type: none"> • Como una librería S-PLUS. • Como R. MICE 1.0 introduciendo selección de predictores, imputación pasiva y automática.
MICE 2.9		<p>Presenta mejor funcionalidad que el MICE 1.0</p> <p>El análisis de los datos imputados es de manera general.</p> <p>Agrega una funcionalidad para la imputación de los datos:</p> <ul style="list-style-type: none"> • Varios niveles. • Selección automática de predictores. • Manejo de datos. • Valores de post-procesamiento imputados. • Rutinas de puesta común especializada. • Herramientas de selección del modelo. • Gráficos de diagnóstico

Fuente: Castro Moisés, 2014.

Elaborado por: Pilco V. Acurio W., 2019.

Dicho algoritmo trabaja con un método de imputación univariante para cada variable incompleta de manera separada. El nivel de la medida establece en su mayoría la representación del modelo de imputación univariante. La función `mice()` en R busca las diferencias de las variable numéricas, binarias, categóricas ordenadas y categóricas no ordenadas estableciendo valores por defecto (Castro M., 2014, p.17).

Tabla 3-2: Técnicas de Imputación univariantes incorporadas.

Método	Descripción	Tipo de escala
Pmm	Predictive mean matching	Numérico
Norm	Regresión lineal bayesiana	Numérico
norm.nob	Regresión lineal no bayesiana	numérico
norm.predict	Regresión lineal	Numérico
Mean	Imputación por media incondicional	Numérico
Logreg	Regresión logística	Factor, >2 niveles
Polyreg	Modelo logístico multinomial	Factor, >2 niveles
Polr	Modelo logístico ordenado	Ordenado
Lda	Análisis lineal discriminante	Factor

Cart	Arboles de clasificación y regresión	Cualquiera
Simple	Muestra aleatoria de los datos observados	Cualquiera

Elaborado por: Castro Moisés, 2014.

Al seleccionar el método de imputación, se debe tener en cuenta la frecuencia con la que las variables continuas no se distribuyen mediante una normal, un problema que se presenta al imputar dichas variables suponiendo que presentan normalidad es que la distribución de los datos imputados no se relacione con los valores observados en el caso de no normalidad de los datos. Una opción para tratar la falta de normalidad es utilizar el predictive mean matching, el pmm es un método de imputación para valores perdidos bajo la propiedad que los valores imputados obtenidos son valores observados de la variable (Castro M., 2014, p. 18).

2.5 Metodología Box-Jenkins (ARIMA)

La Metodología de los modelos ARIMA se estableció por Box y Jenkins en el año de 1976, por tal razón se denomina modelos Box-Jenkins, esta técnica fue creada con la perspectiva de que la serie temporal a predecirse la misma se genera por un proceso estocástico el cual es caracterizado a través de un modelo. Es decir, consiste en hallar un modelo matemático que figure la conducta de la serie temporal de datos y que ayude a predecir solamente ingresando el período de tiempo proporcionado (Jiménez et al., 2006, pp. 187-188).

Los modelos ARIMA no involucran variables independientes, utilizan la información de la serie misma para generar los pronósticos, depende de los patrones de autocorrelación en los datos. Refiriéndose a un conjunto de procedimientos para identificar, ajustar y verificar modelos ARIMA con los datos de la serie de tiempo. Las predicciones se derivan directamente de la forma de un modelo ajustado.

Los datos estimados con el modelo seleccionado se comparan con los datos históricos para observar la serie con exactitud, para ello se toma en cuenta si los residuos son muy pequeños, distribuidos aleatoriamente, sino es satisfactorio el modelo especificado se vuelve a repetir el proceso utilizando un nuevo modelo que mejore el original y en ese instante se considerara útil para pronosticar (Hanke et al, 2010, p. 399).

2.5.1 Modelos Autoregresivos AR(p)

Estos modelos son adecuados para trabajar con series de tiempo estacionarias y el coeficiente f_0 se relaciona con el nivel constante de la serie, no se requerirá del coeficiente ϕ_0 si los datos varían alrededor de cero o se expresan como desviaciones de la media $Y_t - \bar{Y}$.

Para trabajar con estos modelos se toma el valor actual de la serie X_t , y se explica en función de p valores pasados $X_{t-1}, X_{t-2}, \dots, X_{t-p}$, donde p establece el número de rezagos necesarios para predecir un valor actual.

Un modelo autoregresivo de orden p tiene la forma:

$$Y_t = \phi_0 + \phi_1 Y_{t-1} + \phi_2 Y_{t-2} + \dots + \phi_p Y_{t-p} + \varepsilon_t \quad (2.2)$$

Donde:

Y_t : variable respuesta en el tiempo t .

$\phi_0, \phi_1, \phi_2, \dots, \phi_p$: coeficiente a ser estimados.

$Y_{t-1}, Y_{t-2}, \dots, Y_{t-p}$: variable respuesta (independiente) en los retrasos de tiempo $t-1, t-2, \dots, t-p$.

ε_t : es el error de las variables no explicadas por el modelo.

2.5.1.1 Proceso Autoregresivo de Orden 1: AR (1)

En un modelo Autoregresivo AR (1) la variable X_t se determina por el valor pasado, es decir X_{t-1} .

$$X_t = \phi X_{t-1} + \varepsilon_t \quad (3.2)$$

Donde:

ε_t es un procedimiento de ruido blanco con media 0 y varianza constante σ^2 , ϕ es el parámetro.

Si este modelo es estacionario para cualquier valor del parámetro, debe cumplir con los siguientes requisitos:

Estacionario en media: $E(X_t) = E(X_t = \phi X_{t-1} + \varepsilon_t) = \phi E(X_{t-1}) \quad (4.2)$

Esto me indica que su media debe ser constante y finita en el tiempo, esto indica lo siguiente:

$$\begin{aligned}E(X_t) &= \phi E(X_t) \\(1 - \phi)E(X_t) &= 0 \\E(X_t) &= \frac{0}{1 - \phi} = 0\end{aligned}$$

En conclusión, será estacionario siempre que su parámetro $\phi \neq 1$.

Estacionario en Covarianza: su varianza debe ser constante y finita en el tiempo:

$$\gamma_0 = E(X_t - E(X_t))^2 = E(\phi X_{t-1} + \epsilon_t - 0)^2 = \phi^2 V(X_{t-1}) + \sigma^2 \quad (5.2)$$

Con la autocorrelación del procedimiento

$$E(X_{t-1}\epsilon_t) = E[(X_{t-1} - 0)(\epsilon_t - 0)] = \text{cov}(X_{t-1}\epsilon_t) = 0$$

Por el supuesto de un proceso estacionario,

$$E(X_{t-1})^2 = V(X_{t-1}) = V(X_t) = \gamma_0 \quad (6.2)$$

Por lo que $\gamma_0 = \sigma^2$, entonces $\gamma_0 = \frac{\sigma^2}{1 - \phi^2}$

El modelo autoregresivo AR (1) es una versión restringida de un modelo general de medias móviles.

2.5.1.1 Características de un modelo AR (1)

- Invertible
- Estacionario si $|\phi| < 1$
- Su representación visual de la función de autocorrelación, presenta una conducta amortiguada hacia cero y todos sus valores positivos, si $\phi > 0$ o a su vez alternando el signo, iniciando con negativo, si $\phi < 0$.
- Si la función de autocorrelación parcial se anula para retardos superiores a uno, es una norma.

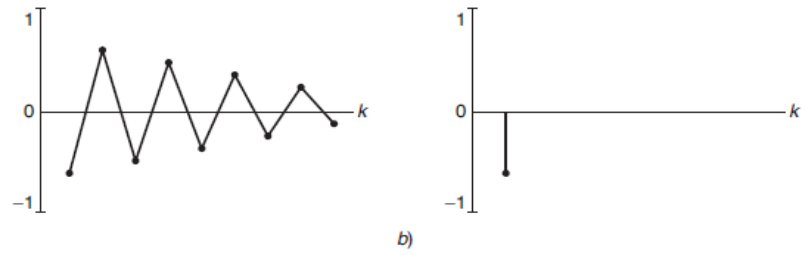
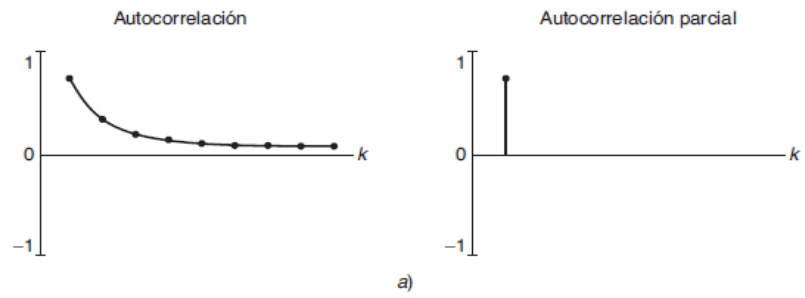


Gráfico 1-2: Coeficientes de autocorrelación y autocorrelación parcial de los modelos AR (1).

Fuente: Hanke, 2010.

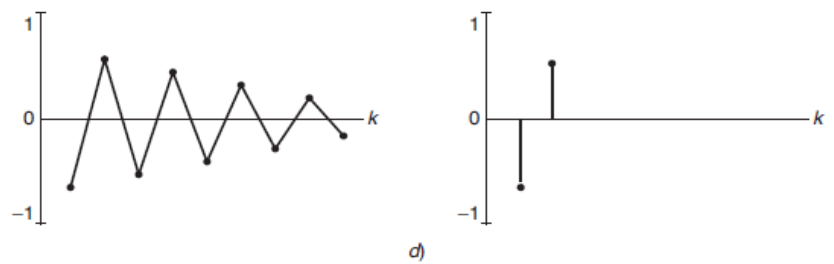
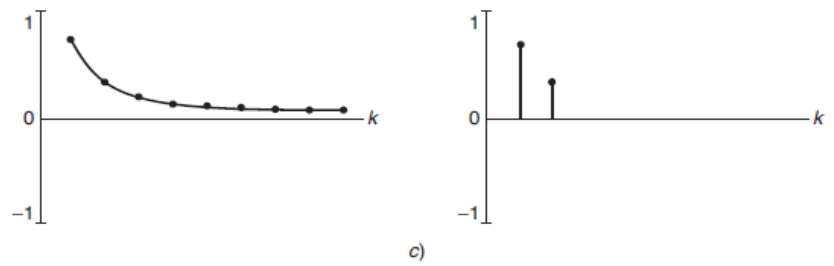


Gráfico 2-2: Coeficientes de autocorrelación y autocorrelación parcial de los modelos AR (2).

Fuente: Hanke, 2010.

2.5.2 *Proceso de Medias Móviles MA (q)*

Este tipo de modelos se establecen por una fuente externa, a su vez manejan el criterio de la linealidad, los valores de la fuente externa influyen en el valor actual X_t .

Un modelo de promedio móvil de orden q se representa por:

$$X_t = \theta_0 - \theta_1 \varepsilon_{t-1} - \theta_2 \varepsilon_{t-2} - \dots - \theta_q \varepsilon_{t-q} + \varepsilon_t \dots \quad (7.2)$$

Si se expresa en términos del operador de retardos:

$$\begin{aligned} X_t &= (1 - \theta_1 L - \theta_2 L^2 - \dots - \theta_q L^q) \varepsilon_t \\ X_t &= \theta_q(L) \varepsilon_t \end{aligned} \quad (8.2)$$

Donde:

ε_t es un ruido blanco y $\mu, \theta_1, \theta_2, \dots, \theta_q$ son los parámetros del modelo.

La diferencia que existe entre los modelos de promedio móvil (MA) y autoregresivos (AR) es que el primero realiza los pronósticos de Y_t con base en una combinación lineal de un número definido de errores ocurridos y el modelo autoregresivo pronostica Y_t con una función lineal de un número finito de valores pasados de Y_t .

Se puede decir que a los modelos de promedio móvil se los designa procesos de memoria corta, y a los modelos autoregresivos procesos de memoria larga.

Para realizar un modelo de promedio de media móvil debe cumplir los siguientes requerimientos:

- Estacionario en media, para todo valor del parámetro

$$\begin{aligned} E(X_t) &= E(\varepsilon_t - \theta \varepsilon_{t-1}) \\ E(X_t) &= E(\varepsilon_t) - \theta E(\varepsilon_{t-1}) \\ E(X_t) &= 0 \end{aligned} \quad (9.2)$$

- Estacionario en covarianza

$$\begin{aligned}\gamma_0 &= E(X_t - E(X_t))^2 = E(X_t)^2 = E(\varepsilon_t - \theta\varepsilon_{t-1})^2 \\ \gamma_0 &= E(\varepsilon_t)^2 + \theta^2 E(\varepsilon_{t-1})^2 - 2\theta E(\varepsilon_t\varepsilon_{t-1}) = \sigma^2 + \theta^2\sigma^2 - 0 \\ \gamma_0 &= (1 + \theta^2)\sigma^2 < \infty \quad (18.2)\end{aligned}$$

Condiciones de Invertibilidad

$$X_t = \varepsilon_t - \theta\varepsilon_{t-1}, \text{ entonces } X_t = (1 - \theta L)\varepsilon_t \quad (8.2)$$

El polinomio de las medias móviles está dado por: $\theta_1(L) = 1 - \theta L$, para hallar las raíces de dicho polinomio se resuelve la ecuación $1 - \theta L = 0$, obteniendo así: $L = \frac{1}{\theta}$.

Un modelo MA (1) cumple la condición de invertibilidad siguiente:

$$|L| = \left| \frac{1}{\theta} \right| > 1, \text{ esto es } |\theta| < 1.$$

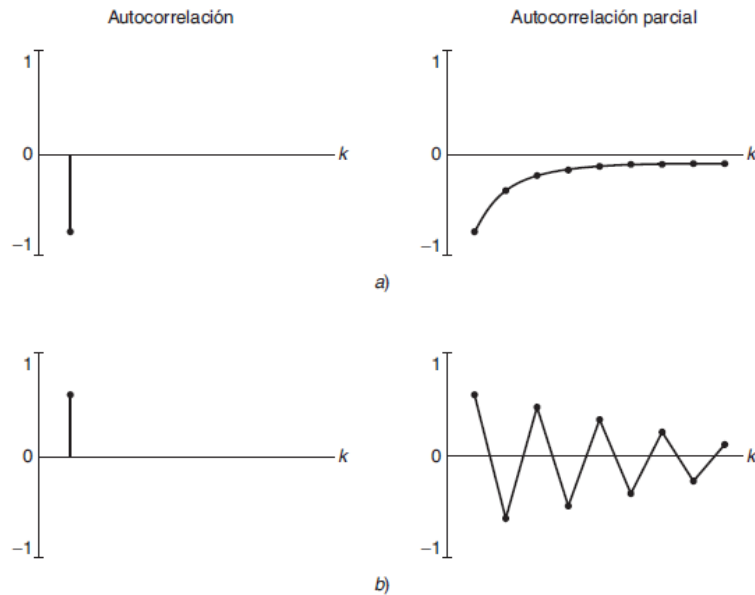


Gráfico 3-2: Coeficientes de autocorrelación y autocorrelación parcial de los modelos MA (1).

Fuente: Hanke, 2010

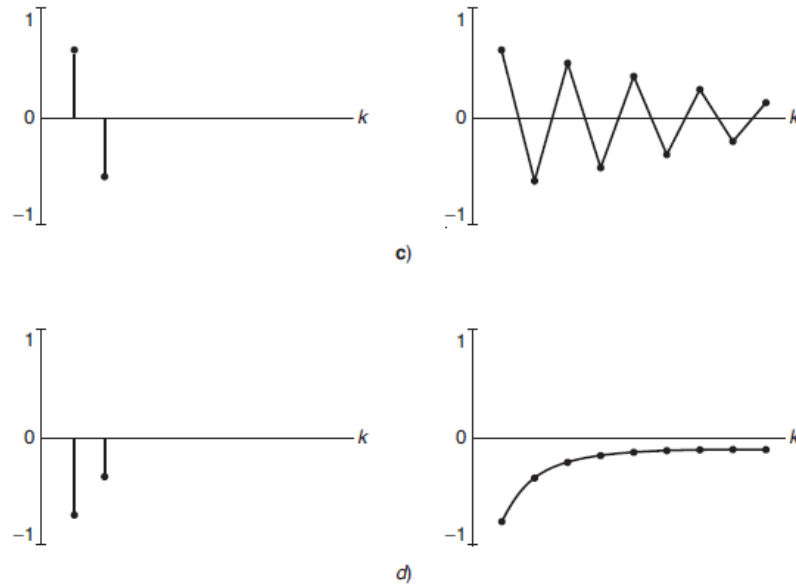


Gráfico 4-2: Coeficientes de autocorrelación y autocorrelación parcial de los modelos MA (2).

Fuente: Hanke, 2010.

2.5.3 Procesos Autoregresivos de Medias Móviles ARMA (p,q)

Es la combinación de un modelo compuesto por términos autoregresivos y un modelo con términos de promedio móvil y así se obtiene un modelo mixto de promedio móvil autoregresivo. Para dichos modelos será ventajoso utilizar la notación ARMA (p,q), en el cual se forma por p términos autoregresivos y q términos de media móvil. Un modelo ARMA (p,q) se representa de la siguiente manera:

$$X_t = c + \phi_1 x_{t-1} + \dots + \phi_p x_{t-p} + \theta_1 \varepsilon_{t-1} + \theta_2 \varepsilon_{t-2} + \dots + \theta_q \varepsilon_{t-q} + \varepsilon_t \quad (9.2)$$

Donde:

ε_t es ruido blanco y $c, \phi_1, \dots, \phi_p, \dots, \theta_1, \dots, \theta_q$ parámetros del modelo.

Un proceso ARMA (p,q) tiene la misma condición de estacionariedad que un proceso AR(p), y de invertibilidad que un promedio móvil.

El modelo ARMA (p,q) se escribe en términos del operador de retardos:

$$\begin{aligned} & (1 - \phi_1 L - \phi_2 L^2 - \dots - \phi_p L^p) X_t \\ & = (1 - \theta_1 L - \theta_2 L^2 - \dots - \theta_q L^q) \varepsilon_t \\ & \phi_p(L) X_t = \theta_q(L) \varepsilon_t \end{aligned}$$

Donde:

ϕ_p es el polinomio autoregresivo y $\theta_q(L)$ es el polinomio de medias móviles.

La representación $MA(\infty)$ cuando su proceso es estacionario es:

$$X_t = \frac{\theta_q L}{\phi_p L} \varepsilon_t, \text{ entonces } X_t = \phi_1 \varepsilon_{t-1} + \phi_2 \varepsilon_{t-2} + \phi_3 \varepsilon_{t-3} + \dots \quad (10.2)$$

Cuando el proceso es invertible se representa un modelo $AR(\infty)$ como sigue:

$$\frac{\theta_q L}{\phi_p L} X_t = \varepsilon_t, \text{ entonces } X_t = \varepsilon_t + \pi_1 Y_{t-1} + \pi_2 Y_{t-2} + \pi_3 Y_{t-3} + \dots \quad (11.2)$$

Para la representación de los pesos de $MA(\infty)$ como de $AR(\infty)$, se ven condicionados a depender del vector finito de parámetros del modelo ARMA (p,q): $\phi_1, \dots, \phi_p, \theta_1, \dots, \theta_q$.

Un procedimiento autoregresivo de medias móviles ARMA(p,q) será estacionario si y solo si el módulo de las raíces del polinomio autoregresivo $\phi_1(L)$ esta fuera del circulo unitario.

Los modelos ARMA (p,q) generan predicciones que dependen de valores actuales y pasados de la respuesta Y, así como de los valores actuales y pasados de los errores (Hanke, 2010, p. 407).

2.5.3.1 Características de un modelo ARMA (p,q) estacionario:

Para este tipo de modelos se toma en cuenta los siguientes criterios: su media es igual a cero, su varianza es constante y finita, la función de autocorrelación es infinita, es decir va decreciendo rápidamente hacia cero.

Media

$$\begin{aligned} E(X_t) &= E(\phi X_{t-1} + \varepsilon_t - \theta \varepsilon_{t-1}) = \phi E(X_{t-1}) \\ E(X_t) &= 0 \quad (12.2) \end{aligned}$$

Función de Autocovarianza

$$\gamma_k = \begin{cases} \gamma_0 = \frac{(1 + \theta^2 - 2\phi\theta)\sigma^2}{1 - \phi^2} & k = 0 \\ \gamma_1 = \phi\gamma_0 - \theta\sigma^2 & k = 1 \\ \gamma_k = \phi\gamma_{k-1} & k > 1 \end{cases}$$

Función de autocorrelación

$$\rho_k = \begin{cases} \rho_1 = \phi - \frac{\theta\sigma^2}{\gamma_0} & k = 0 \\ \rho_k = \phi\rho_{k-1} & k > 0 \end{cases}$$

Para determinar el número de términos autoregresivos y de promedio móvil (orden p y orden q) en un modelo ARMA se determina por los patrones de las autocorrelaciones y autocorrelaciones parciales de la muestra y los valores de los criterios de selección del modelo.

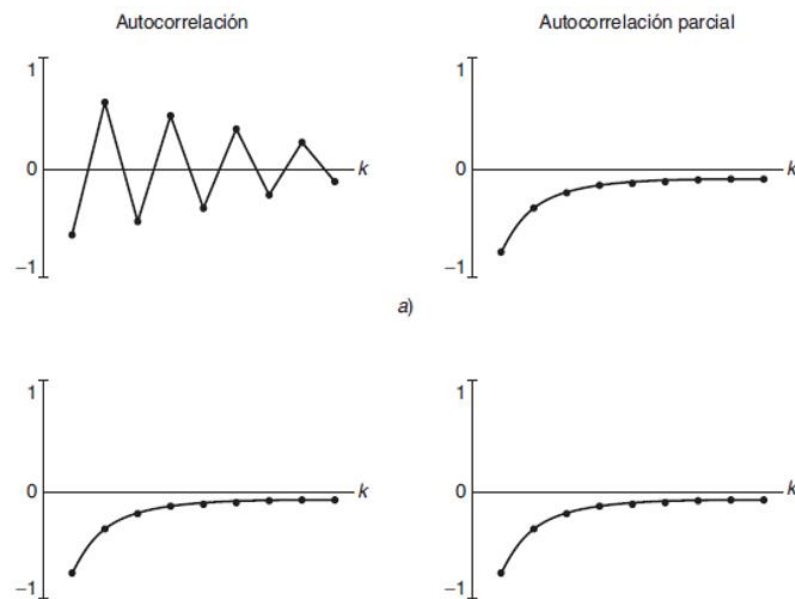


Gráfico 5-2: Coeficientes de autocorrelación y autocorrelación parcial de un modelo mixto ARMA (1,1).

Fuente: Hanke, 2010

Tabla 4-2: Resumen de los Patrones de autocorrelación y autocorrelación parcial de los procesos de promedio móvil autoregresivos.

	Autocorrelaciones	Autocorrelaciones Parciales
MA(q)	Termina después del orden q del proceso	Se desvanecen
AR(p)	Se desvanecen	Terminan después del orden p del proceso
ARMA(p,q)	Se desvanecen	Se desvanecen

Fuente: Hanke, 2010.

Realizado por: Pilco V. y Acurio W., 2019.

2.5.4 Proceso Autoregresivo Integrado y de media móvil ARIMA (p,d,q)

Este tipo de modelos se utilizan cuando ya no existe el supuesto de estacionariedad, es decir cuando su media y la varianza para una serie de tiempo ya no son constantes. Comúnmente las series económicas no son estacionarias, debido a que cambian de nivel en el tiempo o naturalmente la varianza no es constante por lo cual a estos procesos se los denomina procesos integrados.

Por lo tanto se debe realizar diferencias a la serie de tiempo d veces según sea necesario hasta que la misma presente estacionariedad y se aplica a dicha series diferenciada $ARMA(p,q)$, la serie original es $ARIMA(p,d,q)$, transformándose así en una serie de tiempo ARIMA.

Dónde:

p es el número de términos de autoregresivos, d es el número de veces que la serie es diferenciada para que sea estacionaria y q es el número de términos de la media móvil invertible.

Se expresa algebraicamente:

$$X_t^d = c + \phi_1 X_{t-1}^d + \dots + \phi_p X_{t-p}^d + \theta_1 \varepsilon_{t-1}^d + \dots + \theta_q \varepsilon_{t-q}^d + \varepsilon_t^d \quad (13.2)$$

Modelo en forma de Polinomio operador de retardos

$$\Phi(L)(1-L)^d X_t = c + \theta(L)\varepsilon_t \quad (14.2)$$

Donde:

X_t^d es la serie de las diferencias de orden d , ε_t^d

2.5.5 Modelos Sarima

Si una serie Y_t posee una componente módulo con período s es posible excluir diferenciando con un rezago de orden s , es decir, transformando Y_t a:

$$W_t = (1 - L^s)^D X_t = \Delta_s^D X_t \quad D = 0, 1, 2, \dots \quad (15.2)$$

Y se busca una estructura ARMA para W_t .

2.5.6 Pruebas de Raíz Unitaria Estacional

Existen dos posibles modelos para una serie con componente estacional:

1. Modelo con componentes determinísticas, es decir el modelo de descomposición.
2. Modelos no estacionarios, integrados estacionalmente, es decir SARIMA(p,d,q)(P,D,Q)[s], con $d \neq 0$ y $D \neq 0$.

Aunque un modelo SARIMA (p,d,q) (P,D,Q)[s], podría ser tanto estacionario como indicar características estacionales, por lo cual un tercer modelo es el anterior, que se denomina “estacionario estacional”.

El modelo con componente determinístico no es estacionario, pero eliminando la tendencia, el proceso filtrado da como resultado la suma de un proceso determinístico periódico y otro estacionario en covarianza.

2.5.7 Pasos para aplicar modelos Box-Jenkins (ARIMA)

1. Recolectar los datos.
2. Graficar y analizar la serie.
3. Graficar y analizar los autocorrelogramas simple y parcial.
4. Estimar el número de diferencias.
5. Identificar el modelo.
6. Estimar los coeficientes.
7. Validar el modelo:
 1. Graficar los autocorrelogramas simple y parcial de los residuos.
 2. Analizar la significancia de los parámetros del modelo.
 3. Analizar el ECM, AIC, BIC.

8. Verificación de los supuestos:
 1. Independencia
 2. Prueba de Dickey-Fuller (Estacionariedad)
 3. Normalidad
 4. Heterocedasticidad
9. Predecir

2.5.8 Supuestos

2.5.8.1 Prueba de Ljung-Box

El test es usado en series de tiempo para contrastar la hipótesis de independencia. La prueba de Ljung-Box se utiliza comúnmente en un modelo autoregresivo integrado de media móvil (ARIMA).

Se debe tener en cuenta que se aplica a los residuos de un modelo ARIMA equipada, no en la serie original, y en tales aplicaciones, la hipótesis de hecho objeto del ensayo es que los residuos del modelo ARIMA no tienen auto correlación (Fernández, 2017, p. 5).

Hipótesis a Probar

$$H_0: \rho_1 = \rho_2 \dots = \rho_m = 0$$

$$H_1: \rho_i \neq 0 \text{ para algún } i$$

Se debe elegir un m tal que la estimación $r(m)$ de $\rho_m = \rho(m)$ sea confiable.

Estadístico de Prueba:

$$Q = n(n+2) \sum_{k=1}^m \frac{r(k)^2}{n-k} \sim X_{m-1}^2, \quad \text{si } H_0 \text{ es cierta}$$

Región de Rechazo

Se rechaza la hipótesis nula si el valor observado es grande.

$$Q \geq X_{m-1, 1-\alpha}^2$$

$$p = P(X_{m-1}^2 \geq Q)$$

2.5.8.2 Dickey-Fuller

Esta prueba es utilizada para determinar las propiedades de estacionariedad de las series de tiempo, se puede utilizar diferentes procesos.

Se asume que y_t sigue un modelo AR (1)

$$\begin{aligned} Y_t &= \phi Y_{t-1} + \varepsilon_t \\ Y_t - Y_{t-1} &= (\phi - 1)Y_{t-1} + \varepsilon_t \end{aligned} \quad (16.2)$$

Donde:

$\rho = \phi - 1$, $\rho = 0$, me indica que existe una raíz unitaria equivalente a $\phi = 1$ convierte caminata aleatoria sin deriva, proceso estocástico no estacionario.

$$\Delta Y_t = \rho Y_{t-1} + \varepsilon_t$$

Transforma operador de primeras diferencias

$$\Delta Y_i = \beta Y_{i-1} + \varepsilon_i$$

Hipótesis a Probar

$H_0: X_t$ no es estacionaria.

$H_1: X_t$ es estacionaria.

Se rechazará H_0 si el estimador es δ es negativo y significativamente distinto de cero.

- Si $\delta = 0$ entonces: $\Delta X_t = \mu_t$.
- Las primeras diferencias son estacionarias.
- La serie no es estacionaria si $\rho = 1$.
- La regresión estimada y realizado el test sobre la significancia de δ .
- La hipótesis nula es $\delta = 0$

Este test se estima en distintas formas:

- Random Walk
- Random Walk con drift
- Random Walk con drift y tendencia (Gómez Giraldo N., 2006, pp. 116-119).

2.5.8.3 Prueba de Kolmogorov Smirnov

Es un contraste que nos permite comprobar si una distribución empírica se ajusta a una distribución teórica propuesta sea esta una distribución con parámetros o sin parámetros.

Pero cuando se utiliza esta prueba para probar normalidad la prueba de hipótesis queda planteada de la siguiente manera.

H_0 = La variable x sigue una distribución normal.

H_1 = La variable x no sigue una distribución normal.

Pasos:

1. La columna con los valores x_i debe estar ordenada de menor a mayor.
2. La función de distribución f_i se calcula acumulando las equi-probabilidades individuales i/n .
3. z_i obtenemos estandarizando la variable $z_i = \frac{x_i - \bar{x}}{\sigma}$.
4. Encontrar los $\varphi(z_i)$ distribución normal estándar de z_i .
5. Las últimas dos columnas son las distancias acumuladas entre los valores de probabilidad acumulada y los valores teóricos.

$$|F_i - \varphi(z_i)| \text{ y } |F_{i-1} - \varphi(z_i)|$$

6. Obtener el máximo de las dos columnas

$$D = \max (|F_i - \varphi(z_i)| \text{ y } |F_{i-1} - \varphi(z_i)|)$$

2.5.8.4 Homocedasticidad (Test de Goldfeld Quandt)

Es utilizada en el análisis de regresión para detectar la homocedasticidad, comparando las varianzas de dos subgrupos el primero es un conjunto de todos los valores altos y el otro un conjunto de todos los valores bajos, si las variaciones son distintas la prueba rechaza la hipótesis nula, que indica que las variaciones de los errores no son constantes. Goldfeld y Quandt describieron dos tipos de pruebas en su artículo paramétricas y no paramétricas, el supuesto para la prueba es que los datos se distribuyen normalmente.

El estadístico de prueba es la proporción de errores residuales cuadrados medios para las regresiones de los dos subconjuntos de datos, pertenece a la F-test para la igualdad de varianzas. Se puede usar la prueba para una o dos colas.

Pasos:

1. Se ordenan los datos en forma ascendente.
2. Dividir en tres partes los datos.
3. Ubicar el punto medio de los datos.
4. Realizar el análisis de regresión por separado en la parte superior e inferior (valores altos y valores bajos), luego hallar la suma de cuadrados residual.
5. Calcular la relación de la suma de los cuadrados residual, donde RSS_2 es el grupo de valores altos y RSS_1 el grupo de valores bajos $(RSS_2/df)/(RSS_1/df)$.
6. Utilizar la regla de decisión convencional para una prueba F, comúnmente los valores F grandes indican que las variaciones son diferentes.

$$F(n_2 - k, n_1 - k) = \frac{\frac{RSS_2}{n_2} - k}{\frac{RSS_1}{n_1} - k}$$

Dónde:

n_1 y n_2 es el número de observaciones en las regresiones superior e inferior, k es el número de parámetros en el modelo (Goldfeld et al., 1965; citado en Stephanie, 2016, p.1).

2.6 Teoría del Caos

Esta teoría da a conocer ciertos tipos de sistemas complejos y dinámicos que son sensibles a las variaciones en las condiciones iniciales, en meteorología el clima no alcanza un patrón fijo y previsible, presenta las siguiente propiedades: sensibilidad a las condiciones iniciales, es transitivo y sus orbitas periódicas forman un conjunto denso en una región compacta del espacio físico, por lo tanto, se comporta de manera caótica, los procesos y conductas depende de circunstancias inciertas para poder predecir el comportamiento del clima (Reich, 2009, p. 2).

2.6.1 Caracterización del Caos

Permite deducir el orden subyacente que ocultan fenómenos aparentemente aleatorios, se conoce que ecuaciones totalmente deterministas (como el set de Lorenz), muestran las siguientes características que definen el Caos:

- a. Deterministas, existe una “ley” que rige la conducta del sistema (¿Qué es lo contrario de “determinista”? ¿“Aleatorio” o “con libre albedrío”? ¿Existe el Libre Albedrío para las Ciencias Duras o es sólo una ilusión?).
- Para expresar un fenómeno se lo hace por “compresión” en lugar de hacerlo por “extensión”.
 - Para generar los mismos datos observados que el sistema original existe una simulación de menor tamaño (Kb), según Chaitin (1994) un sistema es aleatorio cuando el algoritmo que crea su propia serie ocupa más Kb que el sistema original, en conclusión, expresar el sistema por “extensión” será lo más eficiente y no por medio de un algoritmo.
- b. Son muy sensibles a las condiciones iniciales:
- La situación que provoca una divergencia exponencial en la trayectoria del Espacio de Fase es una desviación infinitesimal en el punto de inicio, lo que se puede cuantificar con el “Exponente de Lyapunov”.
 - La extrema sensibilidad a las condiciones iniciales involucra que el comportamiento del sistema se sugiere partir de cierto “Horizonte de Predictibilidad”, debido a que la incerteza tecnológica asociada a los datos de entrada siempre va a ser mayor que el concepto de “infinitesimal matemático”.
 - A pesar de la impredecibilidad de una trayectoria particular del espacio de Fase, se pueden hallar “Atractores” o zonas del espacio de fase que tienden a ser “visitadas” con mayor frecuencia que otras.

NOTA: la trayectoria en el Espacio de fase de un sistema caótico genera una curva fractal (de dimensión fraccionaria)

- c. Aparenta ser desordenados o aleatorios, pero no lo son:

- Persiguen ecuaciones deterministas
- Presentan atractores

Ecuación determinista pero caótica es:

$$Y_t = 4Y_{t-1}(1 - Y_{t-1}) \quad (17.2)$$

Para ilustrar el efecto mariposa comparando los gráficos que resultan cuando se utilizan las siguientes condiciones iniciales.

- Sistema A: $X_0 = 0.399999$
- Sistema A + una mariposa: $X_0 = 0.400000$ (apenas una millonésima de diferencia).

2.6.2 *Sistemas Caóticos*

En los sistemas reales es común encontrar señales que aparentemente tienen un comportamiento casual, representado por una elevada sensibilidad a las condiciones iniciales e imprevisibilidad a través del tiempo.

Normalmente se define como caóticos, y son caracterizados mediante ciertas variables en el espacio de fases donde sus dimensiones representan variables dinámicas.

Según los autores Hegger Rainer, Kantz Holger y Schreiber Thomas de la página web oficial de Tisean 3.0.1 https://www.pks.mpg.de/~tisean/Tisean_3.0.1/index.html proporciona definiciones para los siguientes comandos.

2.6.3 *Tiempo de Retardo*

- **Incrustaciones y secciones de Poincaré**

El concepto de espacio de las fases es el centro de todos los métodos no lineales en este paquete, el tiempo de retardo y las incrustaciones se utilizan dentro de la mayoría de los otros programas, será importante contar con estas técnicas para la visualización de datos, la selección de parámetros, etc.

- **Información Mutua de Datos**

Esta ayuda a calcular el tiempo de retraso de la información mutua de los datos, es el proceso más sencillo, puesto que utiliza una malla fija de cajas, no implementa correcciones de muestras finitas hasta el momento.

- **Uso: [Opciones] mutuas**

Al no ser válida una opción se interpretará como un nombre de archivo de datos potencial, dado que no existe un archivo de datos, esto indica que es una lectura estándar, significa también stdin. Su posible opción es:

Tabla 5-2: Parámetros de la herramienta mutua

Opción	Descripción	Defecto
-l#	Número de datos a utilizar	Archivo completo
-x#	Número de líneas a ignorar	0
-c#	Columna para leer	1
-b#	Número de cajas para la partición	16
-D#	Retraso de tiempo máximo	20
-o[#]	Nombre del archivo de salida	Sin nombre dado: 'datafile'.mut (o stdin.mut si los datos se leyeron desde stdin) sin -o los resultados se escriben en stdout.
-V#	Nivel de verbosidad 0: solo mensajes de pánico 1: añadir mensajes de entrada / salida	1
-h	Mostrar estas opciones	Ninguna

Fuente: Hegger et al., 2007

- **Descripción de salida:**

En la primera línea está contenido el número de casillas ocupadas, en la segunda la entropía de Shannon (normalizada al número de casillas ocupadas), la última alinea la información mutua (primera columna: retraso, segunda columna: información mutua).

2.6.4 Dimensión de Encaje

- **Descripción del comando: false_nearest**

El programa false ayuda a buscar a los vecinos más cercanos de todos los puntos de datos en m dimensiones e itera a un paso a dichos puntos, es decir retrasa los pasos en el futuro, cuando la distancia de la iteración y del punto más cercano excede un umbral fracción de vecinos falsos para las dimensiones de incrustación especificadas.

Se implementó un segundo criterio que explica si la distancia al vecino más cercano se vuelve más pequeña que la desviación estándar de los datos divididos por el umbral, se omitirá un punto, resultando, así como un criterio más estricto, pero mostrando el efecto de que al aumentar las

dimensiones de incrustación el número de puntos que ingresan a las estadísticas es tan pequeño que las estadísticas completas carecen de significado.

- **Uso: false_nearest [opciones]**

Si presenta una opción que no sea válida se interpretara como un nombre de archivo de datos potencial, puesto que no existe ningún archivo de datos es una lectura estándar y significa stdin, sus posibles opciones son:

Tabla 6-2: Parámetros de la herramienta false_nearest

Opción	Descripción	Defecto
-l#	Número de datos a utilizar	Archivo completo
-x#	Ignorar las primeras # filas	0
-c#	Columnas para leer	1
-m#	Dimensiones mínimas de incrustación de los vectores.	1
-M#,#	Número de componentes, max. dimensión de incrustación de los vectores	1,5
-d#	Retraso de los vectores	1
-f#	Factor de relación	2.0
-t#	Ventana de Theiler	0
-o[#]	Nombre del archivo de salida	Sin nombre de archivo: 'datafile'.fnn (o stdin.del si se leyó stdin) Si no se le da -o se usa la salida estándar
-V#	Nivel de verbosidad 0: solo mensajes de panico 1: añadir mensajes de entrada / salida 2: añadir información sobre el estado actual del programa	3
-h	Mostrar estas opciones	Ninguna

Fuente: Hegger et al., 2007.

- **Construcción de los vectores**

En el caso de una entrada multivariable, los vectores se construyen de la siguiente manera (n=número de componentes):

$$(x_1(i), \dots, x_n(i), x_1(i + delay), \dots, x_n(i + delay), \dots, x_1(i + (maxemb - 1) * delay), \dots, x_n(i + (maxemb - 1) * delay))$$

La dimensión de incrustación mínima dada por la marca -m solo se refiere a la incrustación de los componentes, es decir que si empieza con un vector de tres componentes y das -ml -M3,3 -dl, el programa inicia con

$$(x_1(t), x_2(t), x_3(t))$$

Se prueba la siguiente ecuación:

$$(x_1(t), x_2(t), x_3(t), x_1(t + 1), x_2(t + 2), x_3(t + 3),) \quad (18.2)$$

Y para el siguiente caso realiza lo mismo y así sucesivamente

- **Descripción de la salida:**

La salida se presenta en stdout (o en el archivo):

- Primera columna: la dimensión (contada como se muestra arriba)
- Segunda columna: la fracción de falsos vecinos más cercanos.
- Tercera columna: el tamaño medio del barrio.
- Cuarta columna: el promedio del tamaño al cuadrado del barrio.
- Salida en stderr: Una estadística sobre cuántos puntos se encontraron hasta el tamaño del vecindario dado.

2.6.5 Reducción de Ruido

Debido a que los filtros espectrales son problemáticos con señales caóticas de banda ancha, se han necesitado nuevas técnicas, y que han sido implementadas para utilizar proyecciones de espacio de fase para reducir el ruido, este programa maneja localmente aproximaciones constantes

de la dinámica. La librería ghkss realiza proyecciones localmente lineales. Para realizar pruebas es probable que desee agregar ruido a los datos y comparar el resultado de limpieza con la señal real.

- **Descripción del comando: ghkss**

El ghkss realiza una reducción de ruido, inicialmente plasma una proyección ortogonal en q-dimensional usando una métrica especial, en el caso que se establezca el parámetro -2, se usara una métrica euclidiana esto se aplica en Cawley et al. Así como Sauer y en ocasiones es útil para sistemas de flujo.

- **Uso: ghkss [Opciones]**

Para toda opción que no sea válida se interpretara como un nombre de archivo de datos potencial, debido a que no existe ningún archivo de datos, significa también stdin las posibles opciones son:

Tabla 7-2: Parámetros de la herramienta ghkss.

Opción	Descripción	Defecto
-l#	Número de puntos a utilizar	Archivo completo
-x#	Número de líneas a ignorar	0
-c#	Columna para leer	1, ..., dimensión de los vectores
-m#,#	Número de componentes, dimensión de inserción.	1,5
-d#	Retraso de la incrustación	1
-q#	Dimensión de la variedad para proyectar a	2
-k#	Número mínimo de vecinos	30
-r#	Tamaño mínimo del barrio	(intervalo de datos) / 1000
-i#	Número de iteraciones	1
-2	Usa la métrica euclidiana en lugar de la complicada.	Métrica complicada
-o#	Nombre del archivo de salida	Sin nombre de archivo: 'archivo de datos'.opt.n, donde n es la iteración (o stdin.opt.n si los datos se leyeron desde stdin) sin -o la última iteración también se escribe en la salida estándar

-V#	Nivel de verbosidad 0: solo mensajes de panico 1: añadir mensajes de entrada / salida 2: añadir corrección media y tendencia 4: sumar cuantos puntos se corrigieron para que épsilon	7
-h	Mostrar estas opciones	Ninguna

Fuente: Hegger et al., 2007.

- **Descripción de la salida:**

Los archivos que se obtienen van a contener la serie de tiempo filtrada como una columna, el cuadro de error estándar muestra algunas estadísticas para cada iteración:

- i. el número de vectores corregidos hasta el valor real del tamaño del vector.
- ii. El cambio promedio
- iii. La corrección promedio.
- iv. Indica cuantos puntos la corrección fue irrazonablemente grande y su última línea indica el archivo donde se escribieron los datos corregidos.

2.6.6 Predicción

En Tisean se han implementado varias técnicas de predicción fundadas en el espacio de fases, en las cuales se van diferenciando en la forma en que se aproxima la dinámica proporcionando modelos de orden cero, uno, funciones de base radial y ajustes polinomiales.

- **Descripción del comando: Rbf**

Este programa permite modelar los datos utilizando una función de base radial (rbf) ansatz, las funciones básicas que se utilizan son gaussianas, con puntos centrales para los datos de la serie temporal, si no se facilita la opción $-X$, se usa una clase de fuerza de Coulomb para desplazarse un poco y su distribución más uniforme, donde su varianza se establece en la distancia promedio entre los centros, además prueba el ansatz calculando el error de pronóstico promedio del modelo o efectúa una predicción de i-step utilizando el indicador $-L$, el ansatz es:

$$x_{n+1} = a_0 + SUM_{a_i f_i}(x_n) \quad (19.2)$$

Donde:

x_n es el n -ésimo vector de retardo y f_i es un gaussiano centrado en el punto i central.

- **Uso: rbf [Options]**

Se interpretará como un nombre de archivo de datos potencial a todo lo que no sea una opción válida. Debido a que no existe ningún archivo de datos, es lectura estándar significa también stdin, las posibles opciones son:

Tabla 8-2: Parámetros de la herramienta Rbf.

Opción	Descripción	Defecto
-l#	Número de datos a utilizar	Archivo completo
-x#	Número de líneas a ignorar	0
-c#	Columna para leer	1
-m#	Dimensión de incrustación	2
-d#	Retrasar	1
-p#	Número de centros	10
-X	Desactivar la deriva (fuerza de Coulomb)	Activado
-s#	Pasos para pronosticar (para el error de pronóstico)	1
-n#	Número de puntos para el ajuste; Los otros puntos se utilizan para estimar el error de la muestra	Número de datos
-L#	Determina la longitud de la serie predicha.	Ninguna
-o#	Nombre del archivo de salida; sin -o se usa stdout	'Archivo de datos'.rbf
-V#	Nivel de verbosidad 0: solo mensajes de pánico 1: añadir mensajes de entrada / salida	1
-h	Mostrar estas opciones	Ninguna

Fuente: Hegger et al., 2007.

- **Descripción de salida**

Al observar el archivo de salida se nota que este contiene: las coordenadas de los puntos centrales, la varianza utilizada para los gaussianos, los coeficientes (pesos) de las funciones básicas utilizadas para el modelo, los errores de pronóstico, y se estableció la bandera -L, los puntos predichos.

2.7 Redes Neuronales Artificiales

2.7.1 ¿Qué es una Red Neuronal?

La red neuronal tiene un proceso concreto y paralelo distribuido con la propensión natural de almacenar procedimientos prácticos y hacerlos valederos para su empleo. Es comparable al cerebro del ser humano en dos perspectivas, la primera, la red obtiene conocimientos mediante un procedimiento de preparación; la segunda, los enlaces inter-neuronales llamadas cargas sinápticas muestran gran estabilidad en almacenar los conocimientos (Alves, 2010, pp. 27-28).

2.7.1.1 Elementos básicos que componen una Red Neuronal.

En la figura 1-2 se muestra un bosquejo de red neuronal:

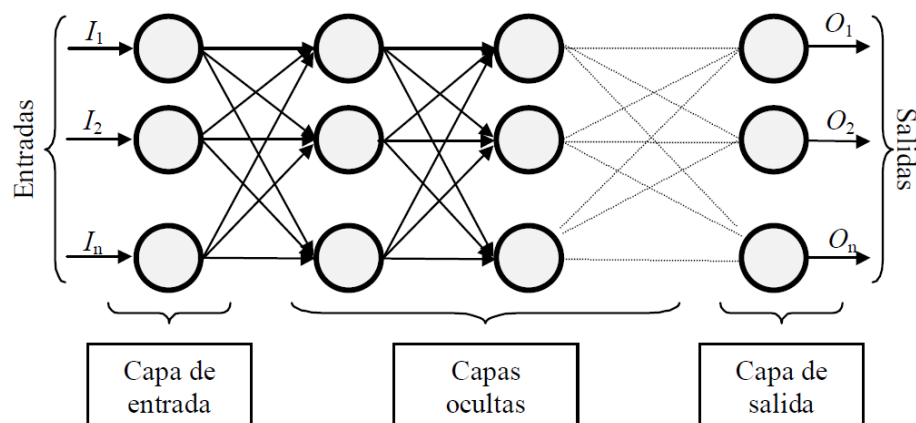


Figura 1-2: Ejemplo de una red neuronal totalmente conectada.

Fuente: Ruiz et al., 2001.

La cual está formada por neuronas interrelacionadas y organizadas en 3 capas (pudiendo cambiar). Los datos son ingresados mediante la capa de entrada, luego prosiguen a la capa oculta la misma que puede ser formada por algunas capas y su respuesta es obtenida a través de la capa de salida.

En la figura 2-2 se contrasta una neurona biológica y una artificial, las cuales tienen semejanzas (con entradas, usan pesos y producen salidas). Una neurona puede ser muy pequeña en sí misma, pero si se mezclan cientos o millones logran resolver problemas muy complicados.

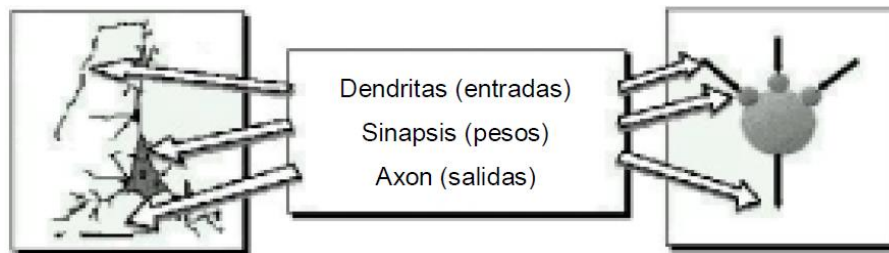


Figura 2-2: Comparación entre una neurona biológica (izquierda) y una artificial (derecha).

Fuente: Ruiz et al., 2001.

2.7.1.2 Ventajas que ofrecen las Redes Neuronales.

Las redes neuronales dado a su naturaleza y fundamentos poseen un gran número de rasgos similares a las neuronas del cerebro. Tienen la facultad de instruirse de la experiencia, como de sistematizar sucesos pasados y para poder llevarlos a nuevos sucesos, etc.; dicha metodología se está utilizando en muchas áreas de la ciencia por sus múltiples ventajas como:

- **Aprendizaje Adaptativo:** pueden aprender a resolver problemas apoyadas en un entrenamiento o en una experiencia inicial.
- **Auto-organización:** una red neuronal es capaz de establecer su propia ordenación de la información que se desarrolla en la fase de aprendizaje.
- **Tolerancia a fallos:** la red posee algunas capacidades que se pueden detener e incluso sufren un gran daño, esto se debe a la pérdida parcial de la red llevando a la degradación de su estructura.
- **Operación en tiempo real:** para los cálculos neuronales se diseñan y elaboran aparatos con hardware específico para realizar dichas operaciones.
- **Fácil inserción dentro de la tecnología existente:** para la inserción modular de sistemas existentes se utiliza chips especializados para optimizar el desenvolvimiento de las redes neuronales en ciertos problemas (Ruiz et al., 2001, p. 8).

2.7.1.3 Niveles o capas de una Red Neuronal

Las neuronas son distribuidas en la red creando niveles o capas, dentro de cada una con un número específico de neuronas, se distinguen 3 tipos de capas:

- **Capa de entrada:** esta capa recibe directamente la información procedente de fuentes externas a la red.
- **Capa oculta:** como su nombre lo indica son capas que se encuentran ocultas en la red. Las capas ocultas pueden tomar valores entre 0 y un número elevado, estas capas pueden

conectarse de diferentes maneras junto con el valor establecido, dando origen a las diversas topologías de redes neuronales.

- Capa de salida: es la última capa, que nos entrega el resultado del procesamiento de la información (Ruiz et al., 2001, p. 16).

2.7.2 *Redes Neuronales Recurrentes*

Existen varias definiciones, de las cuales podemos mencionar que:

- Este tipo de red neuronal tiene dos tipos de conexiones *feedforward* y *feedback* entre neuronas, estas conexiones logran que la información transite hacia delante como hacia atrás durante el trabajo de la red (Sepúlveda, 2011, p. 21).
- Estas redes tienen un diseño que implementa una cierta memoria, consecuentemente un sentido temporal. Para lograr esto se debe implementar algunas neuronas que reciben como entrada el resultado de una de las capas e introducen su salida en una de las capas de un nivel anterior a ella (Diaz, 2016, p. 1).

Constan con características de modelar tanto la no linealidad como las componentes dinámicas de un sistema. Además, este tipo de redes son muy adecuadas para modelar series de tiempo.

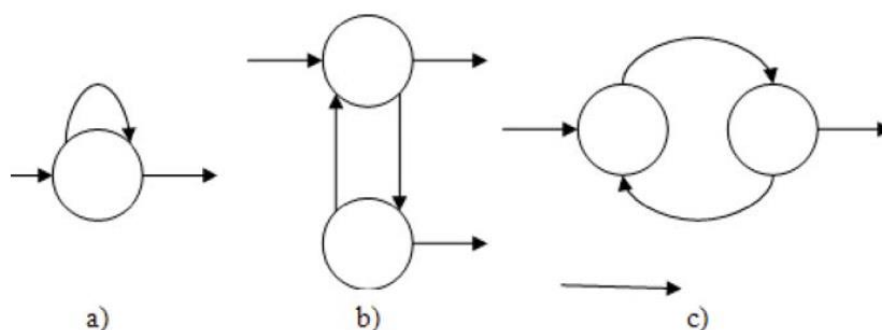


Figura 3-2: Tipos de Recurrencias: **a)** Conexión con la misma neurona, **b)** Conexión con neuronas de la misma capa y **c)** Conexión con neuronas posteriores y anteriores.
Fuente: Sepúlveda, 2011.

Hay diferentes tipos de redes neuronales recurrentes, como son las redes completamente recurrentes y redes parcialmente recurrentes. En estas se encuentran las propuestas por Elman y Jordan, que tienen enlaces *feedforward*, a las mismas que se les ha complementado algunos enlaces hacia atrás (Sepúlveda, 2011, pp. 22-23).

2.7.2.1 Redes parcialmente recurrentes

Estas redes se caracterizan porque existe un conjunto de neuronas situadas en la capa de entrada llamadas neuronas de contexto, neuronas de estado o neuronas auxiliares, trabajando como neuronas receptoras de los enlaces recurrentes (Sepúlveda, 2011, p. 24), siendo este tipo de redes adecuadas para modelar series temporales (Díaz, 2016, p. 1).

Hay dos clases de neuronas de entrada a la red:

- Las que sirven de entrada tomando señales del exterior.
- Las neuronas de contexto se desempeñan como la memoria de la red en donde se recopilan las tareas de las neuronas de cierta capa en la red en un momento o instante anterior.

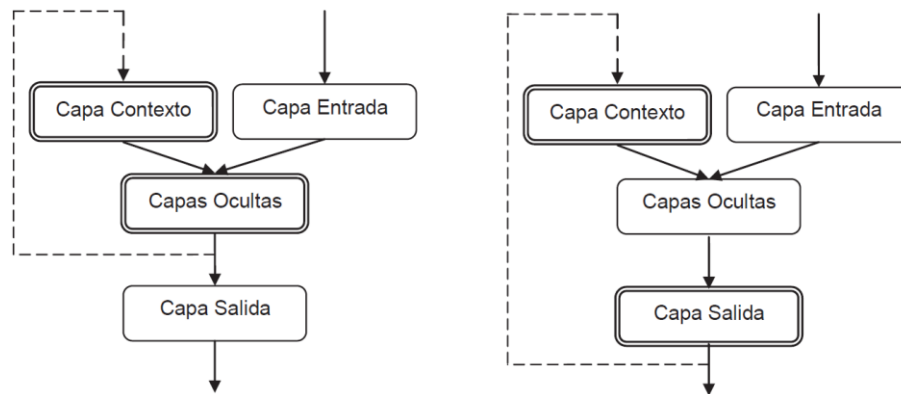


Figura 4-2: Esquema de redes parcialmente recurrentes.

Fuente: Sepúlveda; 2011.

Estas redes se crean mediante una topología tipo *feedforward*, donde se tiene en cuenta la retroalimentación desde las capas de salida o de las capas ocultas, a las capas de entrada. En este tipo de redes se consideran las redes de Elman, Jordan y variantes de estas clases (Sepúlveda, 2011, pp. 24-25).

- **RED DE ELMAN**

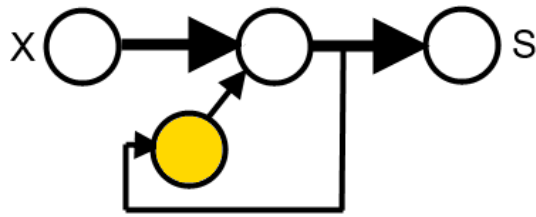


Figura 5-2: Red neuronal de Elman.
Fuente: Diaz; 2016.

La red de Elman también es conocida como red recurrente simple, considero retroalimentación desde las capas ocultas que se dirigen a la capa de contexto y no considero realimentaciones locales (Sepúlveda, 2011, p. 26).

En otras palabras, *las entradas de estas neuronas, se toman desde las salidas de las neuronas de una de las capas ocultas, y sus salidas se conectan de nuevo en las entradas de esta misma capa, lo que proporciona una especie de memoria sobre el estado anterior de dicha capa* (Diaz, 2016, p. 1).

- **RED DE JORDAN**

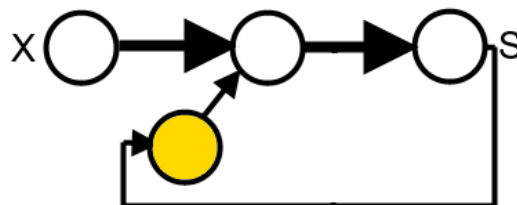


Figura 6-2: Red neuronal de Jordan.
Fuente: Diaz; 2016.

En la red de Jordan, *la entrada de las neuronas de la capa de contexto se toma desde la salida de la red* (Diaz, 2016, p. 1). El diseño de la capa de entrada de esta red se encuentra separada en dos partes: *el conjunto de entradas externas y la retroalimentación de la activación de la capa de salida a través de conexiones de valor fijo*. Este tipo de red se especifica debido a que las neuronas de contexto toman una copia de las neuronas de salida de la red y de ellas mismas (Sepúlveda, 2011, pp. 25-26).

Diferencia entre Elman y Jordan

La diferencia primordial reside en que las redes de Elman las neuronas de contexto reciben copias de las neuronas ocultas, mientras que las redes de Jordan las neuronas de contexto reciben copias de las neuronas de la capa de salida y de sí mismas (Sepúlveda, 2011, p. 25).

2.8 Medidas de evaluación de pronósticos

2.8.1 Medidas dependientes de la escala

2.8.1.1 MSE (Mean Square Error)

Es la medida de e_t^2 , es el promedio de los errores entre el estimador y lo que se estima al cuadrado:

$$MSE = \frac{1}{n} \sum_{t=1}^n (y_t - \hat{y}_t)^2 \quad (20.2)$$

Dónde:

n es el número de muestras y \hat{y}_t es el estimador de y_t .

Características:

- Es el segundo momento del error relacionando la varianza del estimador y el sesgo, el MSE entre más pequeño mejor será el ajuste del estimador de mejor a los datos reales.
- Los datos positivos y negativos al ser elevados al cuadrado hacen que no se cancelen entre sí, otorgando así mayor peso a los errores de mayor tamaño pues son de naturaleza cuadrática.
- Las unidades del MSE no son las mismas de los datos ni del estimador.

Esta medida toma valores entre 0 y ∞ (Vélez et al., 2016, p. 38).

Desventajas

- Es una medida sensible a valores estimados atípicos es decir con poca frecuencia de ocurrencia, esto proyectará un valor superior en la diferencia con el dato real que aumentará con la potencia cuadrada y por lo cual dará un peso en el cálculo del promedio haciéndolo poco fiable.
- Ashley (1983) impulsó un teorema en el cual demuestra que si el MSE es mayor que la varianza de la variable explicativa incluir en un modelo de pronóstico devolverá resultados peores que cuando se la omite.

2.8.1.2 RMSE (Root Mean Square Error)

Raíz cuadrada de la media de los errores al cuadrado

$$RMSE = \sqrt{MSE} = \sqrt{\frac{1}{n} \sum_{t=1}^n (y_t - \hat{y}_t)^2} \quad (21.2)$$

Dónde:

n es el número de muestras, \hat{y}_t pertenece a los valores observados de y_t y se modela en valores en tiempo o lugar t . La función de la medida es la cuadrática o error medio al cuadrado (Vélez et al., 2016, p. 38).

Particularidades

- Se lo llama también desviación cuadrática media se lo usa frecuentemente en la diferencia entre valores pronosticados por un modelo y los datos reales observados, el RMSE ayuda a agregar en una sola medida la capacidad de predicción.
- Esta medida es más adecuada que el MAE para representar el rendimiento del modelo cuando se espera que la distribución de error sea gaussiano acorde con Chai y Draxler (2014), la divergencia con el MSE es que sus resultados está en las unidades originales de la información histórica.
- El uso del RMSE se da cuando existe una función de pérdida cuadrática.

Desventajas

- Se debe evitar el uso del RMSE cuando al evaluar el error hay
- valores atípicos pues estos valores tendrán un efecto muy fuerte en la medida en que se elevan al cuadrado, se evita su uso debido a que es muy sensible a valores atípicos.
- Willmott y Matsuura (2005) sugieren que el RMSE no es un buen indicador de promedio del rendimiento de un modelo y puede ser engañoso este indicador.

La previsión de la varianza de error varía a través del tiempo, a causa de no linealidad en el modelo y a su variación en las variables exógenas si se incluyen en el modelo (Vélez et al., 2016, p. 38).

2.8.1.3 MAE (Mean Absolute Error)

Es el promedio de los valores absolutos de los errores calculados:

$$MAE = \frac{1}{n} \sum_{t=1}^n |y_t - \hat{y}_t| \quad (22.2)$$

Dónde:

n es el número de muestras y \hat{y}_t es el estimador de y_t , la función de pérdida de la medida es la del error absoluto.

Características

- Esta medida entrega un número que puede ser directamente interpretado ya que la pérdida se halla en las unidades de la variable de salida.
- Al utilizar el MAE, los resultados podrían ser afectados por una gran cantidad de valores de error promedio sin necesariamente reflejar errores de gran tamaño.
- Al incluir los valores absolutos en el cálculo del MAE, significa dificultades en los cálculos de la sensibilidad del valor de MAE con respecto a los parámetros del modelo en cuestión (Vélez et al., 2016, pp. 41).

2.8.1.4 MdAE (Median Absolute Error)

Es la mediana de los errores absolutos calculados:

$$MdAE = \text{Mediana}(|y_t - \hat{y}_t|) \quad \text{para } t = 1, \dots, n \quad (23.2)$$

Esta medida de precisión depende de la escala de los datos, es muy útil para comparar métodos del mismo conjunto de datos, pero con la misma escala, no se ve afectada por valores extremos, esta ventaja es a su vez una debilidad ya que no maximiza el uso de la información disponible sobre los errores, un rasgo que debe compartir toda medida “robusta” según Swanson, Tayman y Bryan (2011) (Vélez et al., 2016, p. 42).

2.8.2 Medidas basadas en porcentajes

2.8.2.1 MAPE (Mean Absolute Percentage Error)

Se obtiene como el promedio de los errores porcentuales y no toma en cuenta el signo, este mide el tamaño del error en términos porcentuales por lo tanto independiente de la escala.

$$MAPE = \frac{1}{n} \sum_{t=1}^n \left| \frac{y_t - \hat{y}_t}{y_t} \right| \quad (24.2)$$

Dónde:

n es el número de muestras, y_t es su valor actual y \hat{y}_t es su estimador, el error absoluto es su función de pérdida.

Características

- Para usar esta medida se asocia a series de tiempo homogéneas e equivalentemente espaciadas.
- Donde existe un bajo nivel de datos no se debe usar esta métrica puesto que es sensible a la escala.
- El MAPE es una medida fácil de entender debido a que sus términos están dados en porcentajes, por lo cual es común utilizarla para comparar diferentes modelos de pronósticos en conjunto de datos distintos.
- Posee propiedades estadísticas muy importantes por lo que utiliza las observaciones y posee la variabilidad más pequeña de muestra a muestra (Levy & Lemeshow, 1991).

Desventajas

- Si existen valores en cero, me indica una división entre cero, pero se lo puede solucionar desestimando a los datos de valor cero, surgen varios problemas si existiesen valores de y_t muy pequeños y \hat{y}_t grande puesto esto redundaría en valores de MAPE grandes que generan comparaciones infructuosas. En los casos donde el valor cero sea de gran importancia, es decir casos como la temperatura en grados Fahrenheit y Centígrados, esta métrica no será la más adecuada para la comparación de modelos.
- Si y_t fuese mayor que \hat{y}_t nos genera un porcentaje más pequeño que cuando y_t es menor que \hat{y}_t .
- Pocos valores atípicos pueden alejar su cálculo y así no ajustarse a varias distribuciones por lo cual hace que esta medida no sea tan robusta.

- Según el National Research Council (1980) la validez del MAPE es discutible ya que la distribución de porcentajes los valores absolutos son frecuentemente asimétrica y sesgada a la derecha.
- Sus valores pueden entre 0 y ∞ (Vélez et al., 2016, p. 43).

2.8.2.2 MdAPE (Median Absolute Percentage Error)

Esta medida se calcula como la mediana de los errores porcentuales sin tomar en cuenta el signo, mide el tamaño del error en términos porcentuales:

$$MdAPE = Mediana \left(\left| \frac{y_t - \hat{y}_t}{y_t} \right| \right) \text{ para } t = 1, \dots, n \quad (25.2)$$

Dónde:

n es el número de muestras, y_t es su valor actual y \hat{y}_t es su estimador, el error absoluto es su función de pérdida.

Características:

- Es similar al MAPE (siendo esta regular o simétrica) pero ahora calcula su promedio en vez de resumir los errores en un porcentaje absoluto (APE). Donde toda la APE son ordenadas de menor a mayor y el PAE central (existe un número par de APE entonces el promedio de los dos medios se calcula) se usa para denotar la mediana.
- No es afectada por valores atípicos y no existen dependencia de la escala, por lo tanto, se puede utilizar para diferentes conjuntos de datos.

Desventajas

- Estas medidas al ser basadas en porcentajes de error tienden a ser infinitas o indefinidas si $y_t = 0$ para t en cualquier instante en un período de interés, y presenta una distribución sesgada muy amplia cuando cualquier y_t esta cerca de cero. Según Hyndman, 2006 significa que la MAPE es esencialmente mayor que la MdAPE.
- Su significado es menos intuitivo, la utilización de la simétrica APE minimiza las posibilidades de los valores atípicos y reduce la necesidad de utilizar MdAPE, es difícil combinar MdAPE a través de horizontes y/o series según y cuando las nuevas bases de datos de que se disponga (Makridakis S., 1993; citados en Vélez et al., 2016, p. 44).

2.8.2.3 RMSPE (Root Mean Square Percentage Error)

Es la raíz cuadrada del promedio de los errores en términos porcentuales al cuadrado:

$$RMSPE = \sqrt{\frac{1}{n} \sum_{t=1}^n \left(\frac{y_t - \hat{y}_t}{y_t} \right)^2} \quad (26.2)$$

Dónde:

n es el número de muestras, y_t es su valor actual y \hat{y}_t es su estimador, el error al cuadrado es su función de pérdida de medida.

Características:

Posee las mismas propiedades que RMSE solo que sus resultados se expresan en porcentajes (Swanson, Tayman & Bryan, 2011; citado en Vélez et al., 2016, p. 44).

Desventajas

- Toma un valor infinito o indefinido si y_t es cero o presentar una distribución sesgada cuando y_t es muy próximo a cero, por tal razón el uso de esta medida en modelos donde sus datos son pequeños conteos no es utilizada.
- Puede tomar valores entre 0 y ∞ (Vélez et al., 2016, p. 44).

2.8.2.4 RMdSPE (Root Median Square Percentage Error)

Se define como la raíz cuadrada de la mediana de los errores en términos porcentuales al cuadrado:

$$RMdSPE = \sqrt{\text{Mediana} \left(\left(\frac{y_t - \hat{y}_t}{y_t} \right)^2 \right)} \text{ para } t = 1, \dots, n \quad (27.2)$$

Dónde:

n es el número de muestras, y_t es su valor actual y \hat{y}_t es su estimador, el error al cuadrado es su función de pérdida de medida

Características:

Es independiente de la escala de datos, por lo cual es utilizada con frecuencia para comparar el rendimiento de previsión a través de diferentes conjuntos de datos.

2.8.2.5 sMAPE (Symmetric Mean Absolute Percentage Error)

Esta medida es independiente de la escala, mide el tamaño del error en términos porcentuales, es calculado como el promedio de los valores absolutos de los errores y se divide entre los promedios de los datos reales y sus pronósticos.

$$sMAPE = \frac{1}{n} \sum_{t=1}^n \frac{|y_t - \hat{y}_t|}{\frac{(|y_t| + |\hat{y}_t|)}{2}} \quad (28.2)$$

Dónde:

n es el número de muestras, y_t es su valor actual y \hat{y}_t es su estimador, el error absoluto es su función de pérdida de medida, existen variaciones a la ecuación y se mencionan dos:

$$sMAPE = \frac{1}{n} \sum_{t=1}^n \frac{|y_t - \hat{y}_t|}{\frac{(y_t + \hat{y}_t)}{2}} \quad (29.2)$$

$$sMAPE = \frac{1}{n} \sum_{t=1}^n \frac{|y_t - \hat{y}_t|}{(y_t + \hat{y}_t)} \quad (30.2)$$

En la ecuación 29.2 nos permite el ingreso de valores negativos y positivos al igual que la ecuación 30.2 pero este último limita a valores de 0 a 100%.

El MAPE simétrico se diseñó para solucionar ciertas limitaciones del MAPE, las revisiones iniciales involucran a Armstrong (1985), Flores (1986) y finalmente a Makridakis (1993).

Características:

- Explora que los valores extremos tengan menor influencia, así como neutralizar la asimetría que se genera con el término $|y_t - \hat{y}_t|$ cuando el pronóstico \hat{y}_t sea mayor o menor al real y_t .
- Sus valores a tomar están entre 0 y 200% (Vélez Julián et al., 2016, p. 47).

2.8.2.6 sMdAPE (Symmetric Median Absolute Percentage Error)

Está definido como la mediana de los valores absolutos de los errores divididos entre los promedios de los datos reales y los pronosticados.

$$sMdAPE = \text{Mediana} \frac{|y_t - \hat{y}_t|}{\frac{(|y_t| + |\hat{y}_t|)}{2}} \text{ para } t = 1, \dots, n \quad (31.2)$$

Dónde:

n es el número de muestras, y_t es su valor actual y \hat{y}_t es su estimador, el error absoluto es su función de pérdida de medida.

Características

- Adopta valores entre 0% y 200% al igual que *sMAPE*.
- Los problemas que generan los valores pequeños de y_t pueden ser menos rigurosos para *sMdAPE*.
- Hyndman y Koehler (2005) explican brevemente esta medida y citan en su literatura a Makridakis (1993) el mismo que señala que, tanto MAPE como MdAPE presentan la misma desventaja de que se sancionan más los errores positivos que los negativos, así esta observación condujo a la utilización de las denominadas medidas simétricas (Vélez et al., 2016, p. 48).

2.9 Criterios de información**2.9.1 AIC (Criterio de Información Akaike)**

Este criterio impone una penalización para añadir variables independientes, se debe tener en cuenta que este criterio no busca el modelo correcto ya que parte del indicio que el modelo verdadero puede no estar dentro del conjunto de modelos a evaluar, por tanto, su objetivo es elegir el modelo que genere mejores predicciones (Peña, 2002, pp. 346-348). Por último, el criterio AIC puede ser utilizado para el desarrollo de predicciones dentro de la muestra, para predicciones fuera de la muestra, modelos anidados y no anidados (Gujarati & Porter, 2010; citado en Vélez et al., 2016, p.33).

Y el AIC se lo calcula como sigue:

$$AIC = \ln(\hat{\sigma}^2) + \frac{2}{n}r \quad (32.2)$$

2.9.2 BIC (Criterio de información Bayesiano)

Conocido también como el criterio de Schwarz (SBC también SBIC) es un criterio para seleccionar modelos entre un conjunto finito de modelos, se basa en la función de probabilidad y se encuentra relacionado con el criterio de información de Akaike (AIC).

Es posible aumentar la probabilidad mediante la adición de parámetros, pero si lo hace puede resultar en sobreajuste, el AIC y BIC resuelven este problema mediante la introducción de un término de la penalización para el número de parámetros en el modelo, el término de penalización es mayor en el BIC que en AIC.

Se toma en cuenta la suposición de que los errores del modelo o perturbaciones son independientes e idénticamente distribuidos según una distribución normal, y que la condición límite que la derivada de la probabilidad de registro con respecto a la varianza real es cero:

$$BIC = \ln(\hat{\sigma}^2) + \frac{\ln(n)}{n}r \quad (33.2)$$

Dónde:

$\hat{\sigma}_e^2$ es la varianza del error, la varianza del error se define como:

$$\hat{\sigma}_e^2 = \frac{1}{n} \sum_{i=1}^n (x_i - \hat{x}_i)^2 \quad (34.2)$$

Dónde:

ln logaritmo natural, $\hat{\sigma}^2$ suma residual de cuadrados dividida entre el número de observaciones, **n** número de observaciones (residuos) y **r** es el número total de parámetros (incluyendo el término constante) en el modelo ARIMA.

En series de tiempo según otros autores se lo puede calcular como sigue:

$$BIC = n * \ln\left(\frac{SSE}{n}\right) + k * \ln(n) \quad (35.2)$$

Dónde:

k es el número de parámetros en el modelo, para un modelo sin el término constante se utiliza $k=p+q+1$ (incluyendo $\phi_1, \dots, \phi_p, \theta_1, \dots, \theta_p, \sigma^2$); en si existiese un término constante, $k=p+q+2$ (incluido ϕ_0). Para elegir el mejor modelo se debe elegir el valor de AIC y BIC más pequeño.

Según el libro de “Pronósticos en los negocios de Hanke” en series de tiempo después de la estimación y verificación los dos modelos pueden representar adecuadamente los datos.

- Si los modelos tienen el mismo número de parámetros, el modelo con error cuadrático medio más pequeño s^2 , es el preferido.
- Si los modelos tienen diferente número de parámetros, se selecciona el mejor bajo el principio de parsimonia.
- El modelo con más parámetros puede tener un error cuadrático medio considerablemente menor (Vélez et al., 2016, p.34).

2.10 Coeficiente U de Theil

Henry Theil expuso dos métricas de error en distintos momentos (Vélez, 2016, pp. 54-56), el coeficiente U mide la precisión del modelo, si este coeficiente se acerca a 0 buena precisión en cambio sí se acerca a 1 el modelo no es confiable (Sepúlveda, 2011, p. 35).

$$U = \frac{\sqrt{\frac{1}{n} \sum_{t=1}^n (y_t - \hat{y}_t)^2}}{\sqrt{\frac{1}{n} \sum_{t=1}^n y_t^2 + \frac{1}{n} \sum_{t=1}^n \hat{y}_t^2}} \quad (36.2)$$

Dónde:

n número de muestras, \hat{y}_t valores predichos y y_t valores reales.

2.11 Test de Diebold-Mariano (DM)

Este test lo presentaron Diebold y Mariano en 1995 es asintótico y es válido para condiciones generales, contrasta la igualdad en precisión entre dos conjuntos de pronósticos, este test puede utilizarse aun si se incumple los supuestos con respecto a los errores de predicción, es decir: media cero, distribución normal, falta de autocorrelación y de correlación.

Este contraste compara si la diferencia existente entre las funciones de perdida de errores de predicción de ambos modelos es significativamente distinta de cero (entonces los modelos son diferentes en su calidad predictiva es decir que uno será mejor que el otro) o no.

Contraste

$$DM = L(e_{t+h|t}^b) - L(e_{t+h|t}^m) \quad (37.2)$$

Siendo esta la pérdida diferencial, d_t , donde $L(\cdot)$ es la función de pérdida se puede usar la función de pérdida absoluta o la función de pérdida cuadrática entre otras en el modelo correspondiente.

Hipótesis a Probar

$$H_0: L(e_{t+h|t}^b) = L(e_{t+h|t}^m)$$

Utilizando resultados estándar, Diebold y Mariano proponen contrastar $H_0^{DM} = E(d_t) = 0$ a través de:

$$S_1 = \frac{\bar{d}}{\sqrt{\frac{2\pi\hat{f}_d(0)}{T}}} \sim N(0,1) \quad (38.2)$$

Dónde:

\bar{d} representa la media muestral del diferencial de pérdidas y $\hat{f}_d(\mathbf{0}) = \frac{1}{2\pi} \sum_{T=-\infty}^{\infty} \gamma_d(r)$ densidad espectral del diferencial de pérdidas en la frecuencia cero siendo $\gamma_d(r)$ la autocovarianza de orden del diferencial.

Será unilateral o bilateral dicho contraste según se tenga a priori cual es el conjunto de predicciones superior (Eransus, 2010, p. 84).

CAPITULO III

3 METODOLOGÍA

3.1 Tipo y diseño de la investigación

3.1.1 Descripción del área de estudio

El presente estudio se ejecutó en el Centro de Energía Alternativa y Ambiente “CEAA” de la Escuela Superior Politécnica de Chimborazo. El mismo que cuenta con 11 estaciones meteorológicas ubicada en: Alao, Atillo, EsPOCH, Matus, Multitud, Quimiag, Cumandá, San Juan, Tixán, Tunshi y Urbina.



Figura 6-3: Ubicación de las estaciones meteorológicas (Anexo A).
Elaborado por: GEAA, 2018.

3.1.2 Tipo de investigación

La investigación según:

- El grado de abstracción es de tipo aplicada porque se centra en encontrar mecanismos o estrategias que permitan abordar el problema en específico.
- La naturaleza de estudio es exploratorio puesto que se ha logrado un primer acercamiento a la problemática mediante la búsqueda de artículos relacionados, y predicativa porque mediante modelación estadística se pretende predecir fenómenos futuros.
- La naturaleza de los datos es cuantitativa puesto que las variables a analizar son temperatura y velocidad de viento.
- El periodo temporal es transversal dado que se manejará datos en un periodo de 3 años.
- El método de estudio es inductivo porque se modelará la temperatura y velocidad de viento con el fin de sugerir predicciones con mejor ajuste (Mimenza, 2018, p. 1).

3.1.3 Diseño de investigación

Es no experimental debido a que permite realizar estudios sin manipular las variables y se puede observar los fenómenos en su ambiente natural, para después poder analizarlos, es decir, que no se genera ninguna situación, más bien se observan situaciones ya existentes no provocadas intencionalmente en la investigación. Debido a que se analizó cambios a través de un periodo de tiempo es transversal y de tendencia (Hernández et al., 2010, pp. 151-160).

3.2 Población de estudio

El estudio se realizará de los 33,600 datos registrados durante el período 2014-2017 de las 11 estaciones meteorológicas instaladas en la provincia de Chimborazo.

3.3 Recolección de información

Los datos fueron obtenidos de los diferentes ordenadores que posee cada estación meteorológica y proporcionados por el CEAA para la realización del presente estudio.

3.4 Identificación de variables

Para el presente estudio de investigación se propone trabajar con las variables meteorológicas temperatura y velocidad de viento, debido a la importancia y necesidad de conocer y prever el comportamiento de dichas variables que aportan a los diversos proyectos de investigación que se desarrollan dentro del CEEA.

3.5 Operacionalización de variables

Tabla 1-3: Operacionalización de las variables.

VARIABLE	TIPO	ESCALA	UNIDAD DE MEDIDA	DESCRIPCIÓN
X ₁ : Temperatura	Cuantitativa Continua	Intervalo	°C	Es una magnitud física que refleja la cantidad de calor, ya sea de un cuerpo, de un objeto o del ambiente.
X ₁₄ : Velocidad de Viento	Cuantitativa Continua	Razón	m/s	Es la velocidad con la que el aire de la atmósfera se mueve sobre la superficie de la tierra.

Realizado por: Pilco V. y Acurio W., 2019.

3.6 Análisis de datos

En la presente investigación se trabajó con 33,600 datos cuantitativos, de los cuales se realizó un análisis previo obteniendo estadísticas descriptivas como: media, mediana, moda, porcentaje de datos faltantes, etc. Una vez analizada las bases de datos por años de cada estación meteorológica se hizo el relleno mediante la librería MICE, con la única restricción de aplicar el algoritmo a bases de datos que tenga menos del 20% de los datos faltantes (Argote & López, 2014, p. 19). Luego, se aplicó técnicas estadísticas de pronósticos como: Box-Jenkins, Teoría del Caos y Redes Neuronales obteniendo de cada uno un modelo para realizar los pronósticos, los cuales fueron desarrolladas en R-Studio y TISEAN 3.0.1. Para redes neuronales se hizo a priori una transformación de escala a las bases de datos (temperatura y velocidad de viento) de 0 a 1. Se utilizó los criterios de evaluación e información para escoger el mejor modelo de cada una de las técnicas estadísticas, para la selección de la mejor técnica se usó el coeficiente de U de Theil y el test de Diebold-Mariano.

3.7 Alcances de la investigación

El presente trabajo de investigación tiene alcances descriptivo, correlacional y predictivo el primero indica la descripción de las variables temperatura y velocidad de viento, para identificar el comportamiento de las mismas en la Provincia de Chimborazo y el segundo ayudó a estudiar la relación existente entre los datos.

Tiene como eficacia realizar pronósticos del tiempo en corto y largo plazo mediante la utilización de Box-Jenkins, Teoría del Caos y Redes Neuronales dado que no se cuenta con un modelo definido para variables meteorológicas, buscando así conocer cuan acertado es el modelo de predicción encontrado.

CAPITULO IV

4 RESULTADOS Y DISCUSIÓN

4.1 Análisis, interpretación y discusión de resultados

ANÁLISIS ESTADÍSTICO DESCRIPTIVO

4.1.1 Análisis estadístico de las estaciones meteorológicas

En la presente investigación, se analizó los registros de las 11 estaciones meteorológicas de la provincia de Chimborazo información que proporcionada por el “CEAA” de los años 2014-2017 con lo que se indica los modelos obtenidos con cada una de las técnicas estadísticas.

4.1.1.1 Matriz de datos

En la presente investigación se abordó las variables temperatura y velocidad de viento.

4.1.1.2 Análisis exploratorio de datos

Se necesita conocer a priori el estado de las bases de datos para lo cual se realizó un análisis exploratorio de los registros de las 11 estaciones meteorológicas, para conocer si en todos los años se puede realizar el relleno de datos (MICE) y se efectuó el cálculo de: medidas de tendencia central, dispersión, posición, gráficos.

4.1.1.3 Identificación de datos faltantes

Tabla 1-4: Identificación de datos faltantes.

Estación	Variable	Años	Total	Faltantes	% Faltantes	R ajustado
Cumandá	X ₁	2014	8760	166	1.89%	87.81%
		2015	8760	74	0.84%	83.74%
		2016	8784	3500	39.85%	---

		2017	7296	2189	30.00%	---
	X ₁₄	2014	8760	372	4.25%	15.3%
		2015	8760	74	0.84%	16.86%
		2016	8784	3500	39.85%	---
		2017	7296	2189	30.00%	---
San Juan	X ₁	2014	8760	24	0.27%	94.21
		2015	8760	0	0.00%	---
		2016	8784	0	0.00%	---
		2017	7296	1708	23.41%	---
	X ₁₄	2014	8760	1158	13.22%	40.14%
		2015	8760	4843	55.29%	---
		2016	8784	0	0.00%	---
		2017	7296	1708	23.41%	---
Tixán	X ₁	2014	8760	292	3.33%	89.95%
		2015	8760	0	0.00%	---
		2016	8784	0	0.00%	---
		2017	7296	1637	22.44%	---
	X ₁₄	2014	8760	2365	27.00%	---
		2015	8760	1	0.01%	---
		2016	8784	0	0.00%	---
		2017	7296	1638	22.45%	---
Tunshi	X ₁	2014	8760	228	2.60%	94.97%
		2015	8760	0	0.00%	---
		2016	8784	629	7.16%	80.12%
		2017	7296	1812	24.84%	---
	X ₁₄	2014	8760	230	2.63%	61.29%
		2015	8760	0	0.00%	---
		2016	8784	629	7.16%	52.34%
		2017	7296	1874	25.69%	---
Urbina	X ₁	2014	8760	31	0.35%	93.4%
		2015	8760	0	0.00%	---
		2016	8784	0	0.00%	---
		2017	7296	1687	23.12%	---
	X ₁₄	2014	8760	253	2.89%	38.41%
		2015	8760	0	0.00%	---
		2016	8784	1	0.01%	---
		2017	7296	1687	23.12%	---
Alao	X ₁	2014	8760	67	0.76%	94.43%
		2015	8760	0	0.00%	---

		2016	8784	0	0.00%	---
		2017	7296	221	3.03%	85.78%
	X ₁₄	2014	8760	75	0.86%	52.82%
	X ₁₄	2015	8760	0	0.00%	---
	X ₁₄	2016	8784	0	0.00%	---
	X ₁₄	2017	7296	221	3.03%	49.86%
Atillo	X ₁	2014	8760	549	6.30%	45.81%
		2015	8760	0	0.00%	---
		2016	8784	0	0.00%	---
		2017	7296	221	3.00%	82.76%
	X ₁₄	2014	8760	1637	18.70%	56.20%
		2015	8760	32	0.40%	53.59%
		2016	8784	1	0.00%	---
		2017	7296	221	3.00%	53.55%
Espoch	X ₁	2015	8760	168	1.92%	97.23%
		2016	8784	3	0.03%	95.98%
		2017	7296	173	2.37%	86.51%
	X ₁₄	2015	8760	201	2.29%	75.63%
		2016	8784	4	0.05%	72.22%
		2017	7296	175	2.40%	62.56%
Matus	X ₁	2014	8760	7938	90.62%	---
		2015	8760	107	1.22%	55.99%
		2016	8784	8569	97.55%	---
		2017	7296	221	3.03%	96.32%
	X ₁₄	2014	8760	7939	90.63%	---
		2015	8760	30	0.34%	72.47%
		2016	8784	8569	97.55%	---
		2017	7296	221	3.03%	70.62%
Multitud	X ₁	2014	8760	473	5.40%	81.48%
		2015	8760	551	6.29%	87.00%
		2016	8784	0	0.00%	---
		2017	7296	3622	49.64%	---
	X ₁₄	2014	8760	473	5.40%	2.14%
		2015	8760	921	10.51%	6.46%
		2016	8784	4164	47.40%	---
		2017	7296	3623	49.66%	---
Quimiag	X ₁	2014	8760	13	0.15%	95.27%
		2015	8760	0	0.00%	---
		2016	8784	0	0.00%	---

		2017	7296	221	3.03%	87.57%
	X ₁₄	2014	8760	14	0.16%	70.13%
		2015	8760	0	0.00%	---
		2016	8784	0	0.00%	---
		2017	7296	221	3.03%	64.59%

Realizado por: Pilco V. y Acurio W., 2019.

Fuente: CEAA.

Algunos registros (Tabla 1-4) de las bases de datos se encuentran completas y otras tienen datos faltantes superando el 20%, por lo cual en esos años no se realizó la imputación de datos.

4.1.2 Imputación de datos

Para la imputación de datos se utilizó la librería MICE (Multiple Imputation by Chained Equations), para imputar datos faltantes dado que trabaja en un conjunto de datos reales e imputa datos multivariantes incompletos bajo la utilización de ecuaciones encadenadas, este algoritmo requiere de una especificación de un método de imputación univariante separado para cada variable (Castro M., 2014, p.17). Se requiere la instalación de otros paquetes como el “VIM”, “lattice”, “Rcpp” los mismos que se encuentran disponibles en la página web cran.r-project.org/, el método que se empleó fue el “pmm” el cual realiza un análisis de datos de escala numérica y su ventaja es que algunos de sus valores imputados coinciden con valores observados en la misma variable preservando relaciones no lineales así sea que su parte estructural del modelo de imputación sea incorrecta.

Se mostrará el procedimiento de cómo se realizó la imputación de datos en R-Studio con la librería MICE para la estación de Atillo año 2014, obteniendo los siguientes resultados:

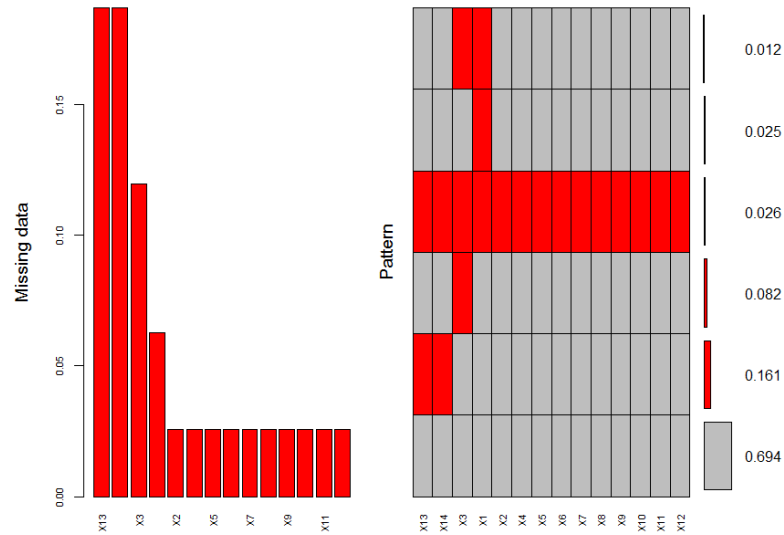


Gráfico 1-4: Gráfico de patrón de datos faltantes (Atillo 2014).
Realizado por: Pilco V. y Acurio W., 2019.

Se observa (Gráfico 1-4) que existen dos variables (X_{13} y X_{14}) que tienen prácticamente el doble de datos faltantes a diferencia de las demás, el siguiente gráfico indica la relación de los valores faltantes entre dos variables:

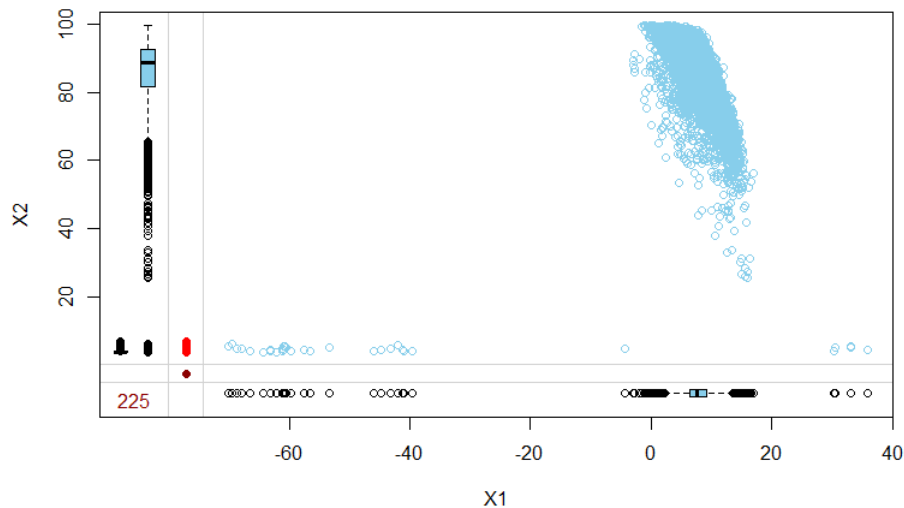


Gráfico 2-4: Relación entre datos faltantes de X1 y X14 (Atillo 2014).
Realizado por: Pilco V. y Acurio W., 2019.

Se evidencia (Gráfico 2-4) que existe una co-ocurrencia de datos faltantes entre las variables, los puntos rojos indican la existencia de los mismos y a su vez se analiza el box-plot el cual indica la existencia de una distribución asimétrica de cola derecha para los mencionados y a su vez para los datos observados, los azules señalan los datos observados, existe una co-ocurrencia de 225 datos faltantes entre las dos variables, también se observan datos outliers.

Se procede a realizar la imputación múltiple, utilizando el paquete MICE, para verificar el modelo de imputación empleado en cada una de las variables, así pues, las variables fueron imputadas mediante el método “pmm” (predictive mean matching esta es la imputación que por defecto realiza el MICE).

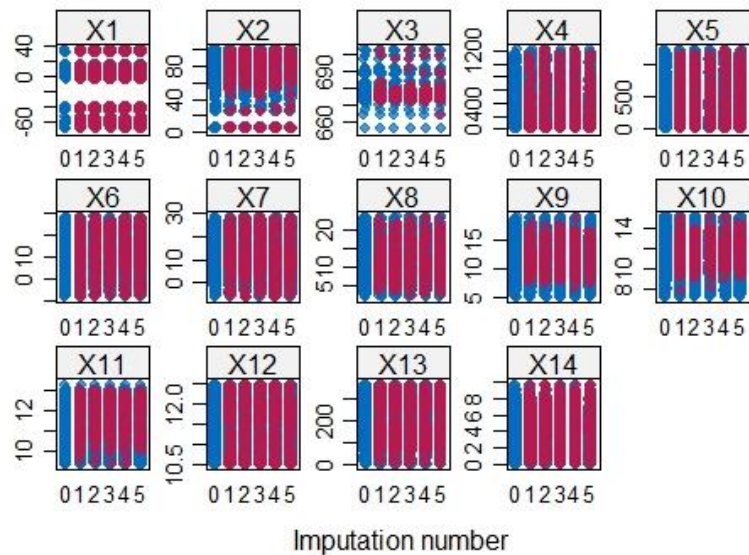


Gráfico 3-4: Gráfico de datos imputados (Atillo 2014).
Realizado por: Pilco V. y Acurio W., 2019.

Los puntos azules (Gráfico 3-4) representan los datos observados mientras que los de color rojo corresponden a los datos imputados, siendo estos valores posibles de acuerdo con el modelo de regresión generado, además, los puntos rojos tienen una correlación de 1 ya que se encuentran sobre una línea y a su vez también se encuentran combinados con los azules lo cual indica que su correlación crece.

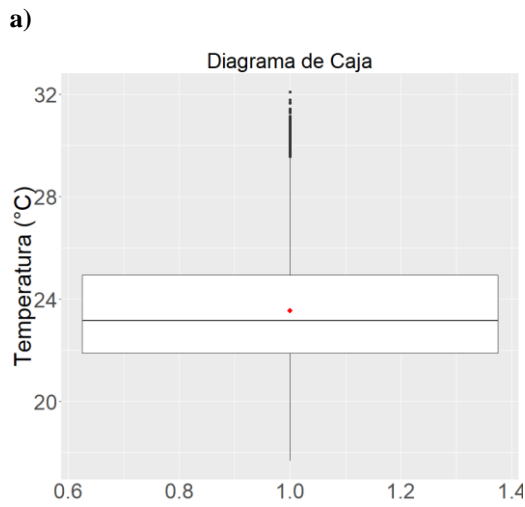
Tabla 29-4: Coeficiente de determinación de las variables imputadas (Atillo 2014).

Resultados del relleno de las Variable meteorológicas con el paquete MICE		
Variabes	R²	RSE
X ₁	45.81%	8.317
X ₂	50.18%	13.93
X ₃	15.29%	1.372
X ₄	82.37%	84.02
X ₅	91.51%	75.17
X ₆	99.78%	0.264
X ₇	99.87%	0.2
X ₈	99.34%	0.2809
X ₉	98.46%	0.2566
X ₁₀	95.36%	0.288
X ₁₁	93.84%	0.2216
X ₁₂	81.70%	0.2415
X ₁₃	35.29%	60.54
X ₁₄	56.20%	1.328

Realizado por: Pilco V. y Acurio W., 2019.

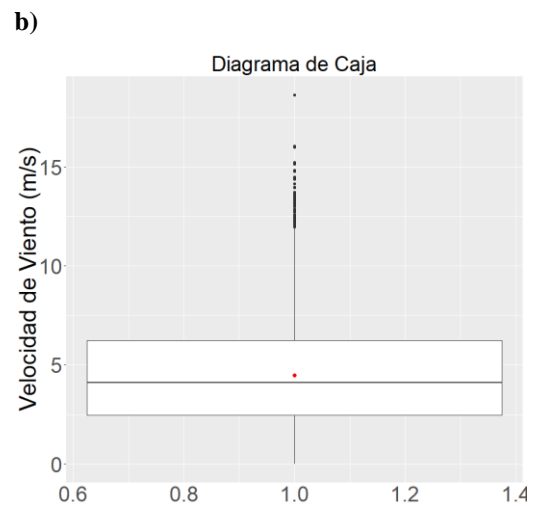
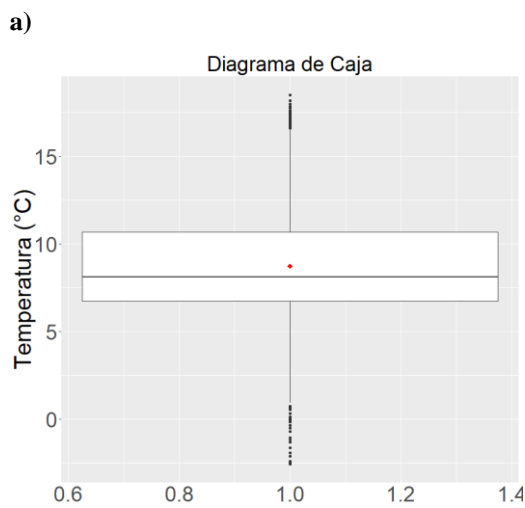
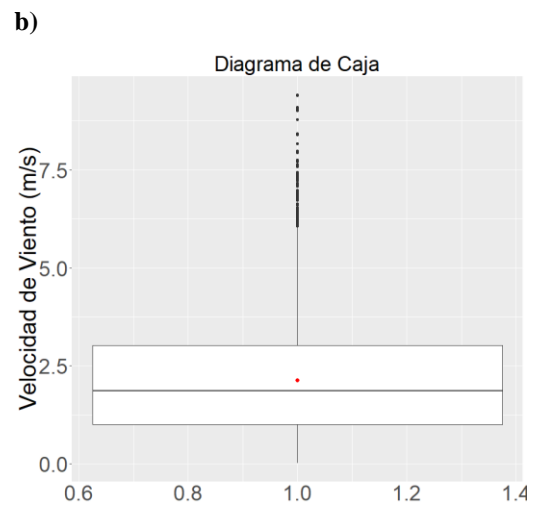
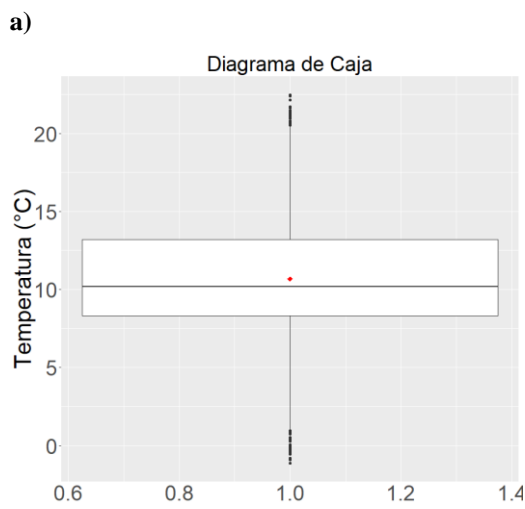
Una vez realizada la imputación de datos se analizó el coeficiente de determinación ajustado (Tabla 2-4) teniendo como condición que si presenta un valor mayor al 79% (Acuña E., 2016, p. 19) la imputación es adecuada para realizar las predicciones con los datos.

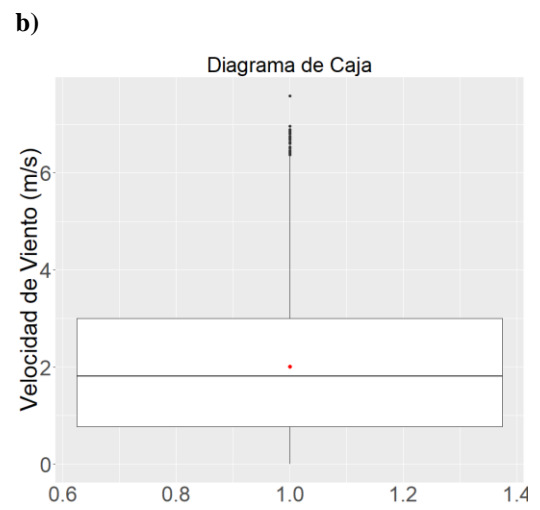
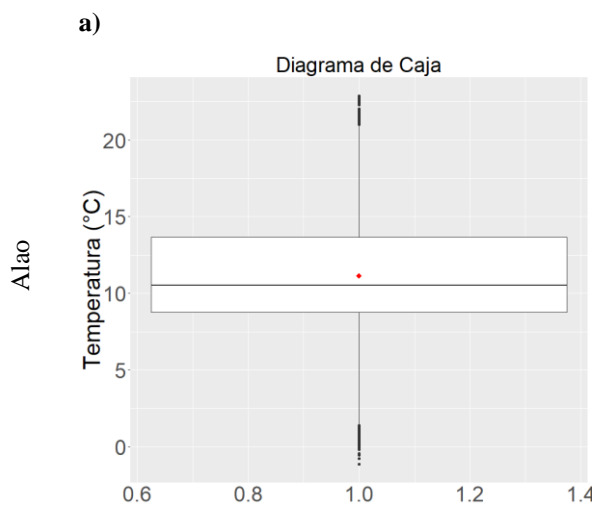
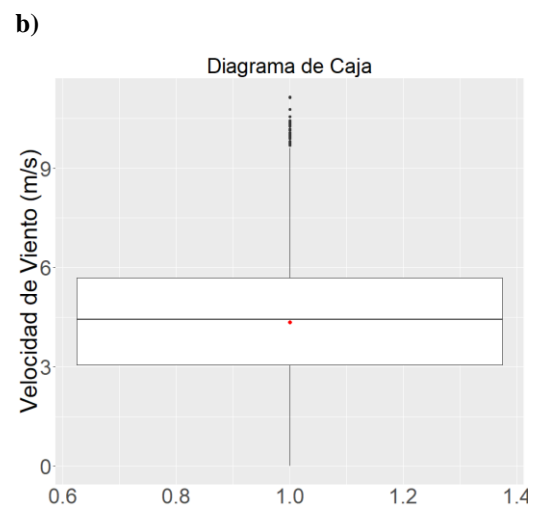
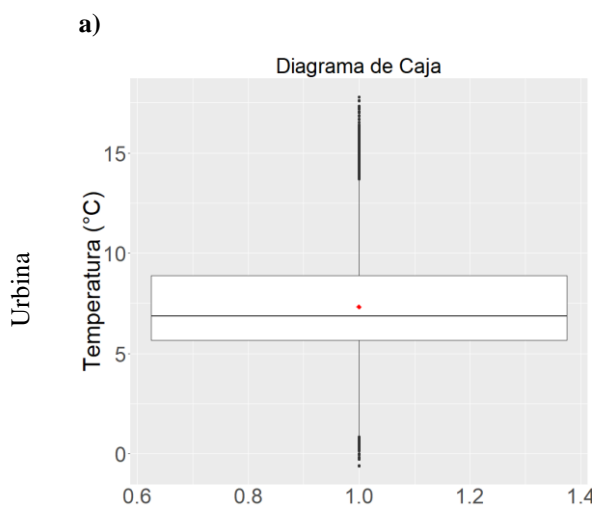
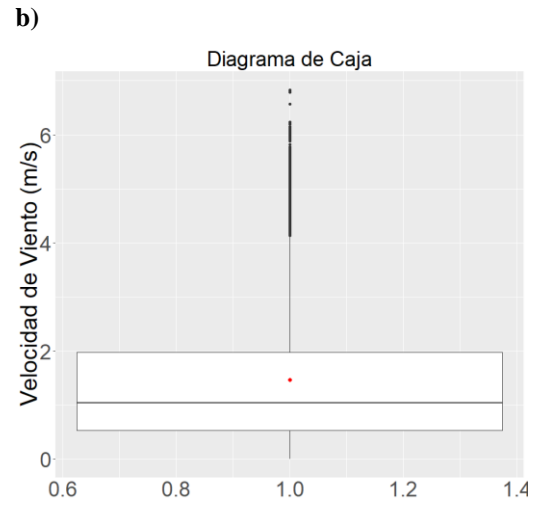
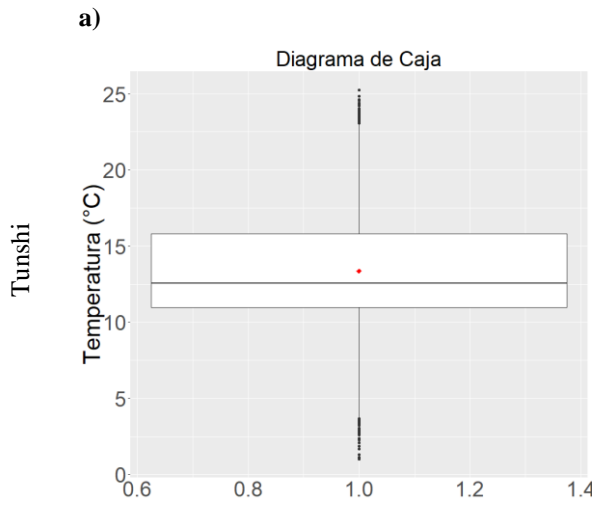
4.1.3 Estadística Descriptiva

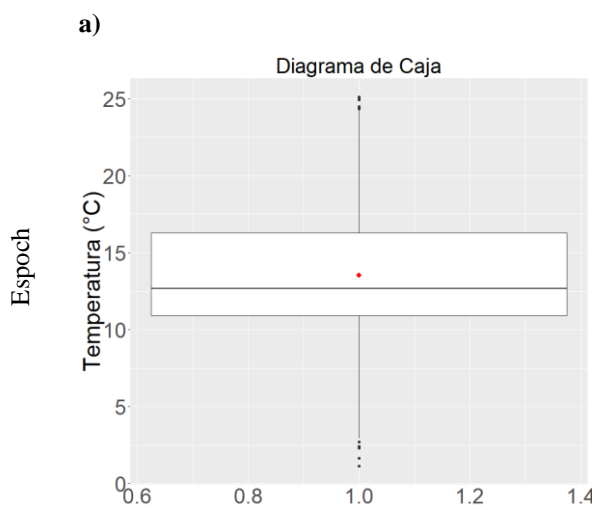
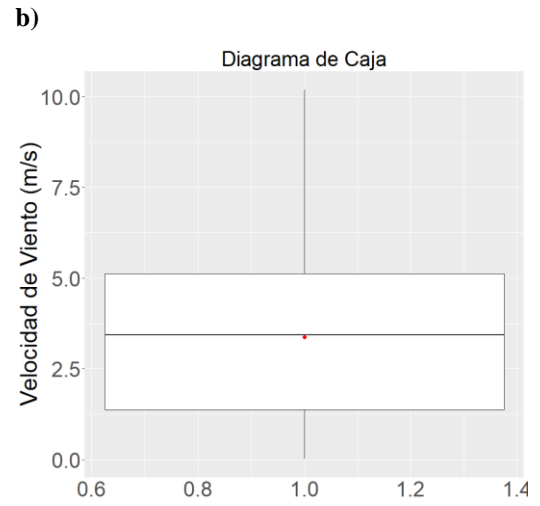
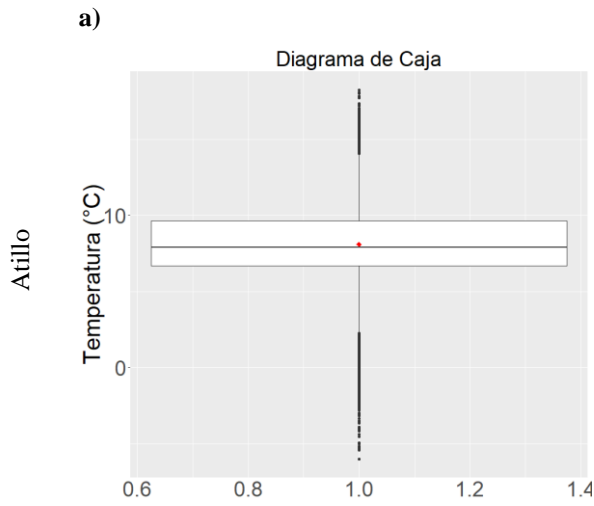


b)

No se realizó el análisis en esta variable, dado que el coeficiente de determinación ajustado es menor al 79% en la imputación y por lo tanto no se considera adecuada.

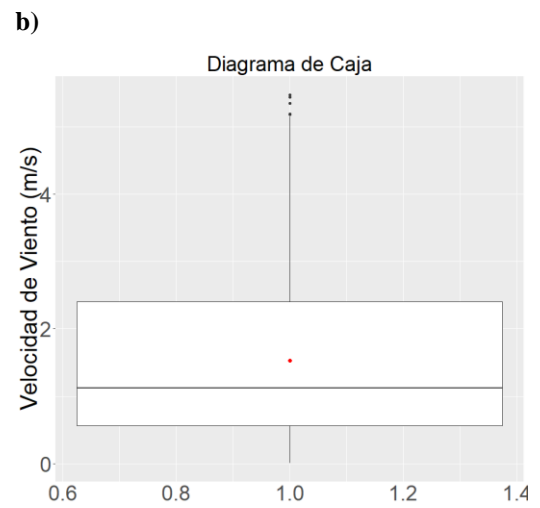
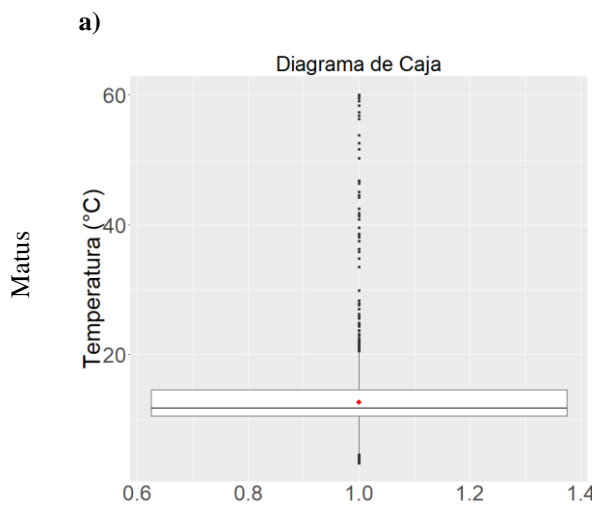


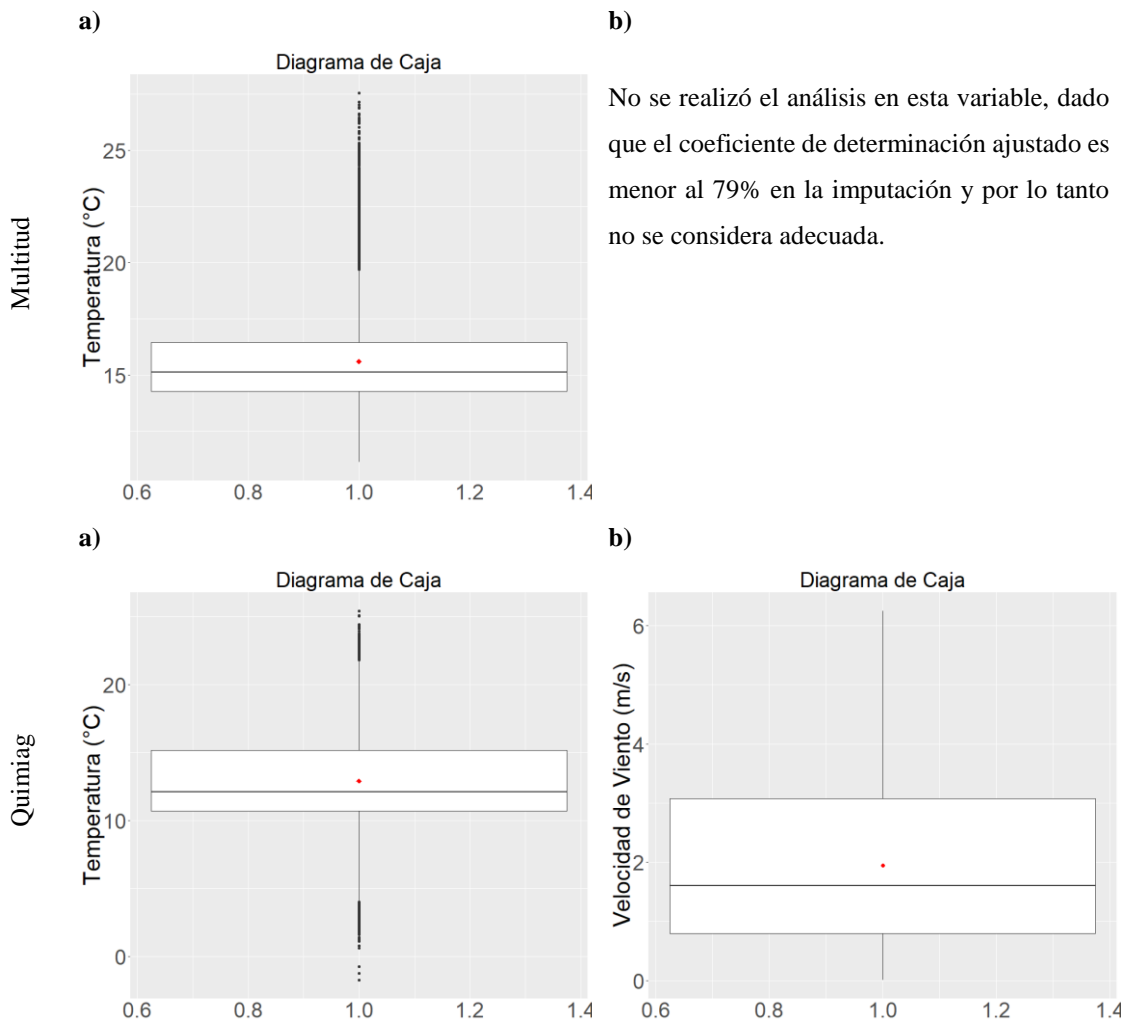




b)

No se realizó el análisis en esta variable, dado que el coeficiente de determinación ajustado es menor al 79% en la imputación y por lo tanto no se considera adecuada.





No se realizó el análisis en esta variable, dado que el coeficiente de determinación ajustado es menor al 79% en la imputación y por lo tanto no se considera adecuada.

Gráfico 4-4: Diagramas de caja de Temperatura y Velocidad de viento para cada estación meteorológica.

Realizado por: Pilco V. y Acurio W., 2019.

En los diagramas de caja (Gráfico 4-4) la mayor parte de las estaciones muestra la existencia de datos sospechosos, los cuales al ser contrastados mediante la Tabla 1-2 propuesta por la OMM (Organización Meteorológica Mundial) se encuentran dentro de los intervalos indicados, por lo tanto, no serán eliminados de la base de datos y se considerarán para el análisis. Además, podemos acotar que estos datos se deben a la variación que existe en las variables temperatura y velocidad de viento.

El punto rojo refleja la media, la misma que se encuentra dispersa de la mediana en todos los diagramas de caja. Se puede concluir para temperatura que en las estaciones de: Cumandá, San Juan, Tixán, Tunshi, Urbina, Alao, Atillo, Epoch, Matus, Multitud y Quimiag siguen una distribución asimétrica sesgada a la derecha de igual manera para Velocidad de Viento excepto en Urbina y Atillo en las que se evidencio una distribución sesgada a la izquierda.

Tabla 3-4: Estadísticas descriptivas de Temperatura y Velocidad de viento de cada estación meteorológica.

Estación	Variable	Estadísticas Descriptivas											
		Período	Media	Mediana	Moda	Desviación estándar	Curtosis	Coefficiente de asimetría	Rango	Mínimo	Máximo	1 ^{er} Cuartil	3 ^{er} Cuartil
Cumandá	X ₁	2014-2015	23.56	23.15	22.02	2.3	-0.01	0.66	14.41	17.68	32.09	21.88	24.95
San Juan	X ₁	2014-2015-2016	10.67	10.17	8.79	3.43	-0.25	0.22	23.59	-1.12	22.47	8.31	13.19
	X ₁₄	2016	2.13	1.86	1.07	1.44	0.74	0.91	9.39	0.02	9.41	1	3.02
Tixán	X ₁	2014-2015-2016	8.74	8.13	7.89	2.85	-0.19	0.47	21.05	-2.54	18.5	6.74	10.69
	X ₁₄	2015-2016	4.47	4.12	2.85	2.55	-0.1	0.58	18.63	0	18.63	2.45	6.24
Tunshi	X ₁	2014-2015-2016	13.35	12.55	12.07	3.4	-0.16	0.42	24.23	1.02	25.25	10.98	15.81
	X ₁₄	2015	1.47	1.04	0.56	1.27	1.29	1.37	6.82	0.01	6.83	0.54	1.98
Urbina	X ₁	2014-2015-2016	7.32	6.85	6.38	2.56	0.32	0.54	18.4	-0.61	17.8	5.66	8.87
	X ₁₄	2015-2016	4.34	4.44	4.54	1.93	-0.47	-0.06	11.14	0.01	11.15	3.06	5.68
Alao	X ₁	2014-2015-2016-2017	11.14	10.51	9.54	3.49	-0.2	0.28	24.02	-1.14	22.88	8.76	13.66
	X ₁₄	2015-2016	2.01	1.81	0.21	1.43	-0.42	0.61	7.59	0	7.59	0.77	3
Atillo	X ₁	2015-2016-2017	8.1	7.9	7.79	2.7	1.45	-0.17	24.24	-6	18.24	6.68	9.64
	X ₁₄	2016	3.37	3.43	0.48	2.14	-1.04	0.18	10.17	0.02	10.19	1.37	5.13
Epoch	X ₁	2015-2016-2017	13.54	12.65	11.78	3.58	-0.39	0.42	23.95	1.15	25.1	10.93	16.28
Matus	X ₁	2015-2017	12.7	11.76	11.05	3.6	30.42	3.15	56.67	3.33	60	10.58	14.58
	X ₁₄	2015	1.53	1.13	0.38	1.19	-0.45	0.81	5.47	0.01	5.48	0.56	2.41
Multitud	X ₁	2015-2016	15.61	15.13	18.21	2.04	3.41	1.52	16.43	11.14	27.56	14.27	16.45
Quimiag	X ₁	2014-2015-2016-2017	12.91	12.13	11.28	3.29	0.13	0.45	27.15	-1.73	25.43	10.7	15.15
	X ₁₄	2015-2016	1.95	1.6	0.4	1.34	-0.93	0.51	6.24	0.01	6.25	0.8	3.08

Realizado por: Pilco V. y Acurio W., 2019.

Mediante las estadísticas descriptivas (Tabla 3-4) se puede decir que:

Para temperatura en las estaciones de: Cumandá, San Juan, Tixán, Tunshi, Urbina, Alao, Atillo, Espoch, Matus, Multitud y Quimiag la media es mayor que la mediana por lo cual se deduce que sigue una distribución asimétrica sesgada a la derecha de igual manera para Velocidad de Viento excepto en Urbina y Atillo donde la media es menor que la mediana lo que indica que sigue una distribución sesgada a la izquierda, corroborando lo mencionado en los diagramas de caja.

Para temperatura se encontró un coeficiente de asimetría mayor a cero en las estaciones de: Cumandá, San Juan, Tixán, Tunshi, Urbina, Alao, Espoch, Matus, Multitud y Quimiag lo que indica que la curva es asimétricamente positiva, además, que los datos tienden agruparse más en la parte izquierda de la media; excepto para Atillo lo que muestra que la curva es asimétricamente negativa, asimismo, los datos tienden agruparse más en la parte derecha de la media.

En cuanto a curtosis para temperatura se encontró un valor mayor a cero en las estaciones de: Urbina, Atillo, Matus, Multitud y Quimiag lo que indica una distribución leptocúrtica; excepto para Cumandá, San Juan, Tixán, Tunshi, Alao y Espoch donde se nota un valor menor de cero, lo que muestra una distribución platicúrtica.

Para velocidad de viento se encontró un coeficiente de asimetría mayor a cero en las estaciones de: Cumandá, San Juan, Tixán, Tunshi, Alao, Atillo, Espoch, Matus, Multitud y Quimiag lo que indica que la curva es asimétricamente positiva, además, que los datos tienden agruparse más en la parte izquierda de la media; excepto para Urbina lo que muestra que la curva es asimétricamente negativa, asimismo, los datos tienden agruparse más en la parte derecha de la media.

En cuanto a curtosis para velocidad de viento se encontró un valor mayor a cero en las estaciones de: San Juan y Tunshi lo que indica una distribución leptocúrtica; excepto para Tixán, Urbina, Alao, Atillo, Matus y Quimiag donde se nota un valor menor de cero, lo que muestra una distribución platicúrtica.

Los registros de todas las estaciones meteorológicas no se encuentran normalmente distribuidos.

4.1.4 Modelación Box-Jenkins (ARIMA)

Para implementar esta técnica se requiere conocer a priori que las variables que se van a estudiar son series estacionarias.

Tabla 4-4: Test de Dickey Fuller
(Estacionariedad) estaciones meteorológicas.

Estación	Variable	Valor p
Cumandá	X_1	0.01
San Juan	X_1	0.01
	X_{14}	0.01
Tixán	X_1	0.01
	X_{14}	0.01
Tunshi	X_1	0.01
	X_{14}	0.01
Urbina	X_1	0.01
	X_{14}	0.01
Alao	X_1	0.01
	X_{14}	0.01
Atillo	X_1	0.01
	X_{14}	0.01
EsPOCH	X_1	0.01
Matus	X_1	0.01
	X_{14}	0.01
Multitud	X_1	0.01
Quimiag	X_1	0.01
	X_{14}	0.01

Realizado por: Pilco V. y Acurio W., 2019.

A un nivel de significancia de 0.05 se rechaza la hipótesis nula (Tabla 4-4) y se concluye que las series de tiempo de las distintas estaciones meteorológicas son estacionarias. Posteriormente, se prosigue a modelar cada serie de tiempo en este caso se analiza la estación de Atillo:

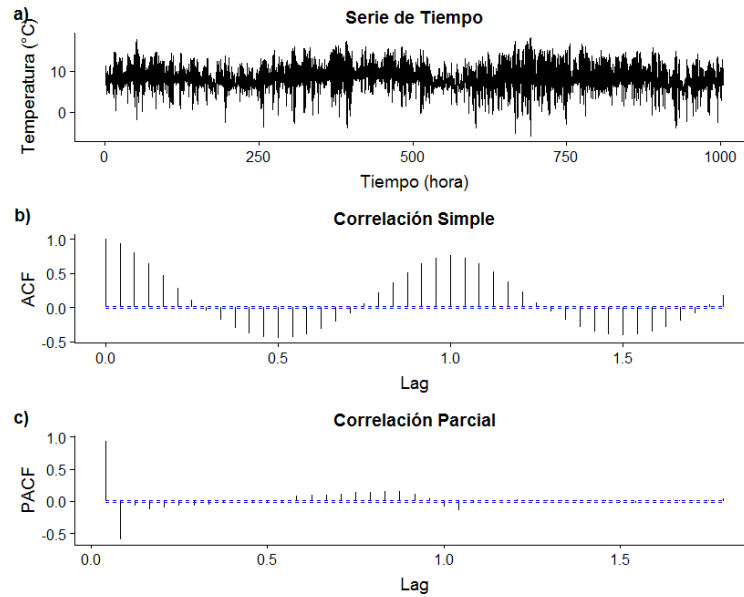


Gráfico 5-4: Gráfica de la serie de tiempo de la Temperatura (Atillo).

Realizado por: Pilco V. y Acurio W., 2019.

Se puede constatar que la serie de tiempo (Gráfico 5-4) a) cumple el requerimiento de estacionariedad y no presenta tendencia, la correlación simple b) muestra estacionalidad cada 24 rezagos lo que muestra indicios de un modelo Sarima y en la autocorrelación parcial c) se evidencia crecimiento exponencial. Se prosigue a eliminar la estacionalidad, para lo cual es necesario aplicar un test de diferencias divididas.

```
> # Número de Diferencias
> ndiffs(datos, test = "pp")
[1] 0
> ndiffs(datos, test = "kps")
[1] 1
> ndiffs(datos, test = "adf")
[1] 0
```

Gráfico 6-4: Diferencias divididas de la Temperatura (Atillo).

Realizado por: Pilco V. y Acurio W., 2019.

Se encuentra el número de diferencias divididas (Gráfico 6-4) con el test “Kps”, el mismo que recomienda realizar una diferencia dividida a la serie de tiempo para eliminar la estacionalidad.

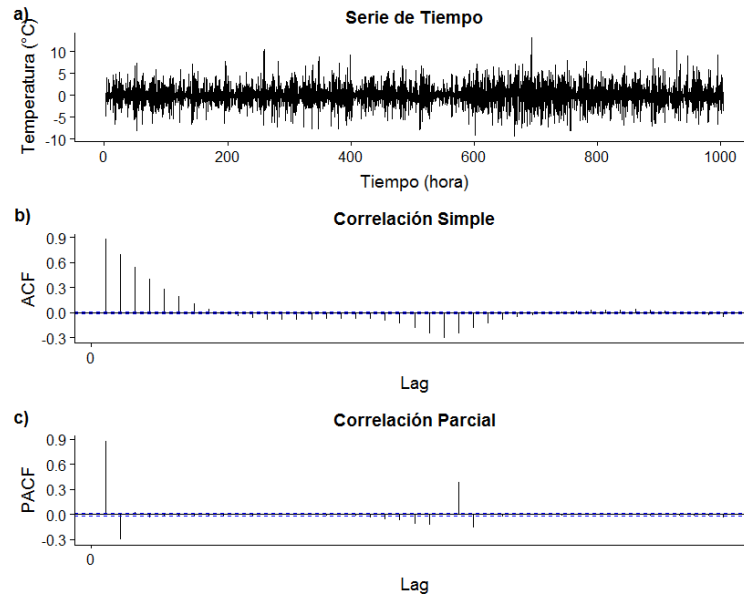


Gráfico 7-4: Gráfica de la serie de tiempo de Temperatura eliminada la estacionalidad (Atillo).
Realizado por: Pilco V. y Acurio W., 2019.

Se puede constatar que la serie de tiempo (Gráfico 7-4) a) cumple el requerimiento de estacionariedad y no presenta tendencia, en la correlación simple b) no se observa el patrón de estacionalidad y en la autocorrelación parcial c) se evidencia crecimiento exponencial, por lo tanto, se proceder a buscar los modelos.

Teniendo en cuenta estas características mostradas se han propuesto 5 modelos de los cuales mediante los criterios de información y evaluación se seleccionará el más idóneo.

Tabla 5-4: Criterios de evaluación para los posibles modelos de Temperatura ARIMA (Atillo).

Nº	Modelo	Medidas de Escala				Medidas Basada en Porcentajes					
		MSE	RMSE	MAE	MdAE	MAPE	MdAPE	RMSPE	RMdSPE	Smape	sMdAPE
1	(4,0,1) (2,1,0) ₂₄	38.1846	6.1794	6.0542	6.3482	8.259	3.4179	22.2041	3.4179	1.2189	1.2617
2	(2,0,0) (2,1,0) ₂₄	36.9783	6.081	5.9566	6.2371	8.1164	3.3917	21.8366	3.3917	1.2126	1.2581
3	(2,0,0) (1,1,0) ₂₄	34.4638	5.8706	5.7315	6.0463	7.9768	3.2604	21.6502	3.2604	1.1951	1.2396
4	(1,0,1) (1,1,1) ₂₄	26.2898	5.1274	4.9992	5.0543	6.5193	2.8379	17.4911	2.8379	1.1491	1.1732

Realizado por: Pilco V. y Acurio W., 2019.

Tabla 6-4: Criterios de información para los posibles modelos de Temperatura ARIMA (Atillo).

Nº	Modelo	Criterios de Información	
		AIC	BIC
1	(4,0,1) (2,1,0) ₂₄	53251.81	53316.53
2	(2,0,0) (2,1,0) ₂₄	53469.41	53509.85
3	(2,0,0) (1,1,0) ₂₄	55729.76	55762.11
4	(1,0,1) (1,1,1) ₂₄	48534.18	48574.62

Realizado por: Pilco V. y Acurio W., 2019.

De los modelos propuestos tanto por los criterios de evaluación (Tabla 5-4) e información (Tabla 6-4) se evidencia que el mejor modelo es (1,0,1) (1,1,1)₂₄, por lo cual se analizara los supuestos en dicho modelo.

Análisis de los supuestos para Temperatura del modelo (1,0,1) (1,1,1)₂₄ de Atillo.

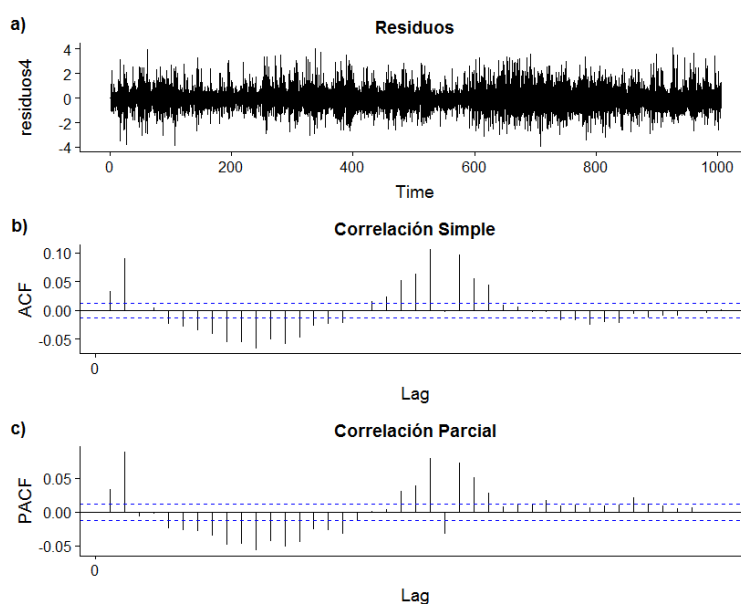


Gráfico 8-4: Serie y autocorrelogramas de los residuos de Temperatura del modelo (1,0,1) (1,1,1)₂₄ (Atillo)

Realizado por: Pilco V. y Acurio W., 2019.

Se puede constatar que los residuos (Gráfico 8-4) del modelo a) cumplen con el requerimiento de estacionariedad y no presenta tendencia, en la correlación simple b) y parcial c) se evidencia que los rezagos no se encuentran dentro de la banda de confianza, por lo tanto, los residuos probablemente no cumplan con algunos de los supuestos como: normalidad, estacionariedad, independencia y varianza constante.

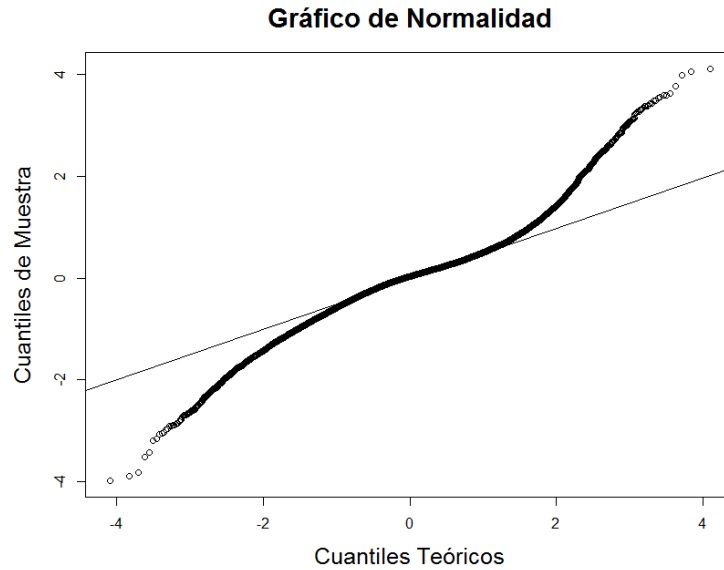


Gráfico 9-4: Gráfico de normalidad de los residuos de Temperatura del modelo (1,0,1) (1,1,1)₂₄ (Atillo).
Realizado por: Pilco V. y Acurio W., 2019.

Tabla 7-4: Valores p de los supuestos del modelo (1,0,1) (1,1,1)₂₄ (Atillo).

Valores p			
Independencia	Estacionariedad	Normalidad	Homocedasticidad
2.2e-16	0.01	2.2e-16	2.2e-16

Realizado por: Pilco V. y Acurio W., 2019.

Se realiza un análisis de los residuos (Tabla 7-4) se verifica independencia mediante el test de Box-Ljung encontrando un valor p de 2.2e-16, para estacionariedad se utiliza el test de Dickey-Fuller obteniendo un valor p de 0.01, se prueba normalidad con el test de Jarque Bera mostrando un valor p de 2.2e-16 y homocedasticidad con el test Goldfeld-Quandt identificando un valor p de 2.2e-16, por lo que se rechaza en todos los supuestos la hipótesis nula, concluyendo que los residuos no son independientes (están correlacionados), son estacionarios, no siguen un patrón de normalidad y no son homogéneos (heterocedásticos).

Predicción

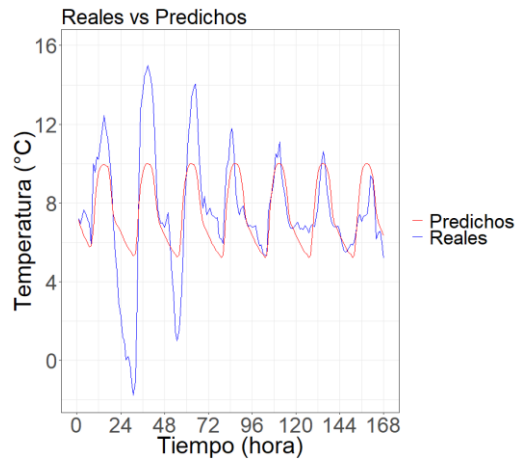


Gráfico 10-4: Datos reales vs predichos de Temperatura (Atillo).

Realizado por: Pilco V. y Acurio W., 2019.

Al comparar los datos reales vs predichos (Gráfico 10-4) presentan un comportamiento casi parecido, evidenciándose un mejor ajuste en los 150 primeros pronósticos, a partir de aquí se observa muchos desfases. Se obtiene un coeficiente de determinación 51.53% lo que indica la existencia de una correlación moderada, considerando que no es un buen modelo de predicción.

Variable: *Velocidad de Viento*

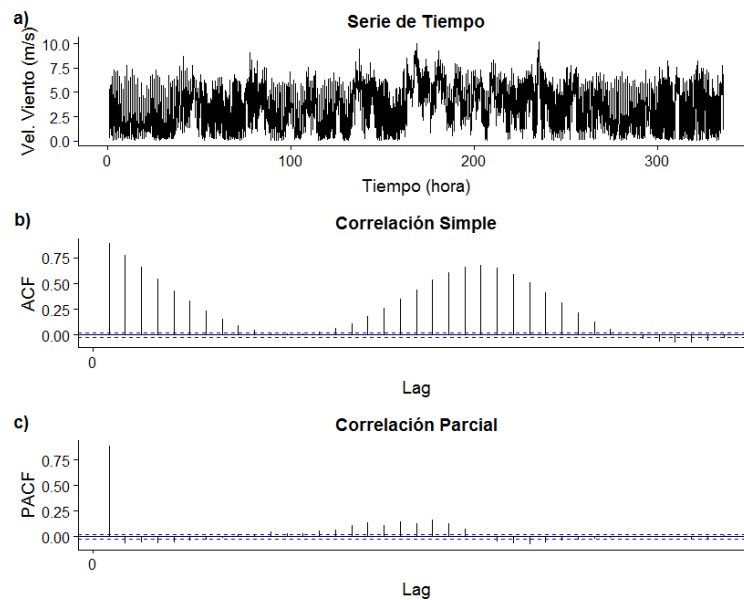


Gráfico 11-4: Serie y autocorrelogramas de Velocidad de Viento (Atillo).

Realizado por: Pilco V. y Acurio W., 2019.

Se puede constatar (Gráfico 11-4) que la serie de tiempo a) cumple el requerimiento de estacionariedad y no presenta tendencia, en la correlación simple b) se observa estacionalidad cada 24 rezagos lo que muestra indicios de un modelo Sarima y en la autocorrelación parcial c) se evidencia crecimiento exponencial. Se prosigue a eliminar la estacionalidad, para lo cual es necesario aplicar un test de diferencias divididas.

```
> ndiffs(datos, test = "pp")
[1] 0
> ndiffs(datos, test = "kps")
[1] 1
> ndiffs(datos, test = "adf")
[1] 0
```

Gráfico 12-4: Diferencias divididas de Velocidad de viento (Atillo).

Realizado por: Pilco V. y Acurio W., 2019.

Se encuentra el número de diferencias divididas (Gráfico 12-4) con el test “Kps” el mismo que recomienda realizar una diferencia dividida a la serie de tiempo para eliminar la estacionalidad.

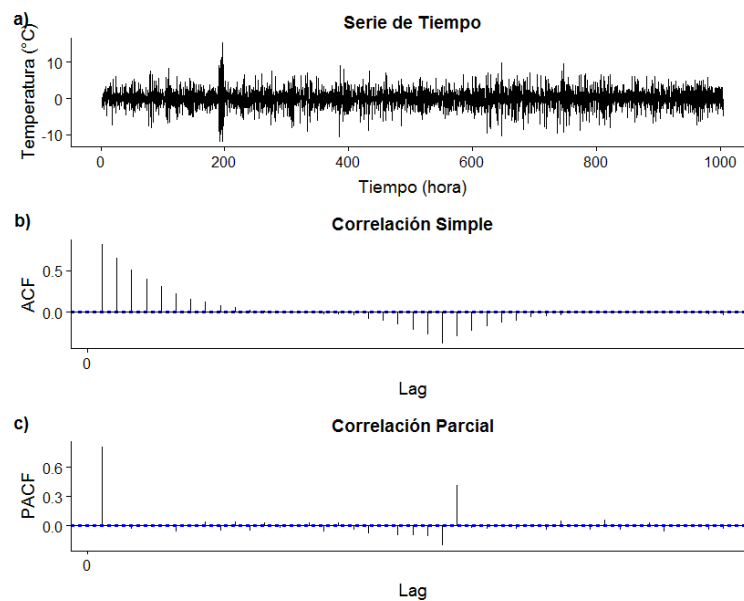


Gráfico 13-4: Serie y autocorrelogramas sin estacionalidad de Velocidad de viento (Atillo).

Realizado por: Pilco V. y Acurio W., 2019.

Se puede constatar que la serie de tiempo (Gráfico 13-4) a) cumple el requerimiento de estacionariedad y no presenta tendencia, en la correlación simple b) no se observa el patrón de estacionalidad y en la autocorrelación parcial c) se evidencia crecimiento exponencial, por lo tanto, se proceder a buscar los modelos.

Teniendo en cuenta estas características mostradas se han propuesto 5 modelos de los cuales mediante los criterios de información y evaluación se seleccionará el más idóneo.

Tabla 8-4: Criterios de evaluación para los posibles modelos de Velocidad de Viento ARIMA (Atillo).

Nº	Modelo	Medidas de Escala				Medidas Basada en Porcentajes					
		MSE	RMSE	MAE	MdAE	MAPE	MdAPE	RMSPE	RMdSPE	Smape	sMdAPE
1	(2,0,1) (2,1,0) ₂₄	10.2557	3.2025	2.811	3.0473	4.9185	1.587	14.0848	1.587	0.9003	0.8849
2	(2,0,2) (2,1,0) ₂₄	9.9633	3.1565	2.7674	3.003	4.8498	1.5511	13.8736	1.5511	0.8948	0.8736
3	(3,0,0) (1,1,1) ₂₄	1.6477	1.2836	1.0883	1.0809	1.7074	0.5028	5.2894	0.5028	0.5816	0.4467
4	(1,0,0) (1,1,0) ₂₄	12.0721	3.4745	2.9847	3.1443	5.4033	1.6695	15.5318	1.6695	0.9088	0.9099
5	(1,0,1) (1,1,1) ₂₄	1.6433	1.2819	1.0866	1.081	1.7048	0.4996	5.2889	0.4996	0.5812	0.4418

Realizado por: Pilco V. y Acurio W., 2019.

Tabla 9-4: Criterios de información para los posibles modelos de Velocidad de viento ARIMA (Atillo).

Nº	Modelo	Criterios de Información	
		AIC	BIC
1	(2,0,1) (2,1,0) ₂₄	22709.95	22751.89
2	(2,0,2) (2,1,0) ₂₄	22710.04	22758.96
3	(3,0,0) (1,1,1) ₂₄	20940.46	20982.39
4	(1,0,0) (1,1,0) ₂₄	23651.27	23672.23
5	(1,0,1) (1,1,1) ₂₄	20995.85	21030.8

Realizado por: Pilco V. y Acurio W., 2019.

De los modelos propuestos por los criterios de evaluación (Tabla 8-4) se evidencia que el mejor modelo es (1,0,0) (1,1,1)₂₄, por otra parte, por los criterios de información (Tabla 9-4) el modelo más adecuado es (3,0,0) (1,1,1)₂₄. Para poder seleccionar el mejor modelo entre los dos, procedemos a valorar los errores que presenta cada uno.

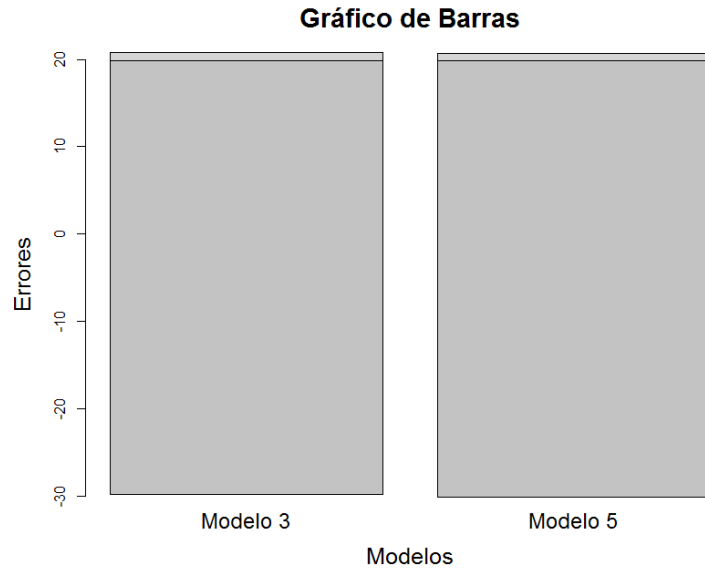


Gráfico 14-4: Gráfico de los errores de los modelos 3 y 5 de Velocidad de viento (Atillo).

Realizado por: Pilco V. y Acurio W., 2019.

En la evaluación de los errores (Gráfico 14-4) de los modelos 3 y 5 se observa que presentan el mismo error de predicción, por ello se aplica el test de exactitud de Diebold-Mariano obteniendo un valor p de 0.216 al ser menor al nivel de significancia (0.05) se rechaza la hipótesis nula y se concluye que los modelos no tienen la misma exactitud. Por lo que se puede decir mediante los criterios de información y el test que el mejor modelo es $(3,0,0)$ $(1,1,1)_{24}$.

Análisis de los supuestos para Velocidad de viento del modelo $(3,0,0)$ $(1,1,1)_{24}$ de Atillo.

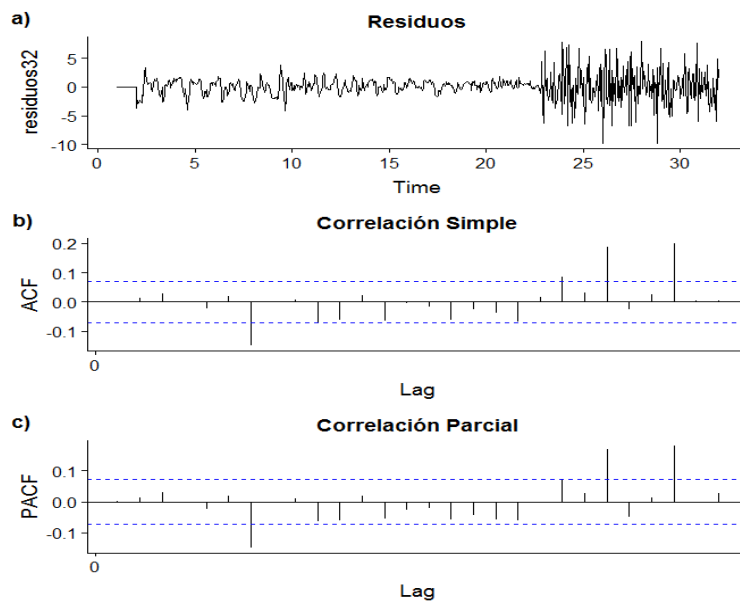


Gráfico 15-4: Serie y autocorrelogramas de los residuos de Velocidad de viento del modelo $(3,0,0)$ $(1,1,1)_{24}$ (Atillo).

Realizado por: Pilco V. y Acurio W., 2019.

Se puede constatar que los residuos (Gráfico 15-4) del modelo a) cumplen el requerimiento de estacionariedad y no presenta tendencia, en la correlación simple b) y parcial c) se evidencia que no todos los rezagos se encuentran dentro de la banda de confianza, por lo tanto, los residuos probablemente no cumplan con algunos de los supuestos como: normalidad, estacionariedad, independencia y varianza constante.

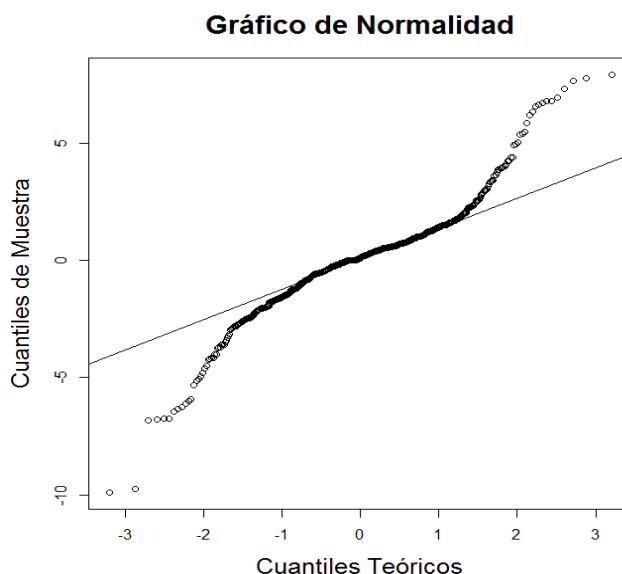


Gráfico 16-4: Gráfico de normalidad de los residuos de Velocidad de Viento del modelo (3,0,0) (1,1,1)₂₄ (Atillo)
Realizado por: Pilco V. y Acurio W., 2019.

Tabla 10-4: Valores p para los supuestos del modelo (3,0,0) (1,1,1)₂₄ (Atillo).

Valor p			
Independencia	Estacionariedad	Normalidad	Homocedasticidad
7.149e-08	0.01	2.2e-16	2.2e-16

Realizado por: Pilco V. y Acurio W., 2019.

Se realiza un análisis de los residuos (Tabla 10-4): se verifica independencia mediante el test de Box-Ljung encontrando un valor p de 7.149e-08, para estacionariedad se utiliza el test de Dickey-Fuller obteniendo un valor p de 0.01, se prueba normalidad con el test de Jarque Bera mostrando un valor p de 2.2e-16 y homocedasticidad con el test Goldfeld-Quandt identificando un valor p de 2.2e-16, por lo que se rechaza en todos los supuestos la hipótesis nula, concluyendo que los residuos no son independientes (están correlacionados), son estacionarios, no siguen un patrón de normalidad y no son homogéneos (heterocedásticos).

Predicción

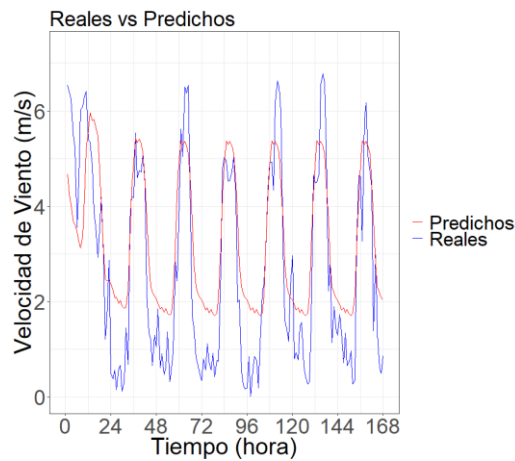


Gráfico 17-4: Datos reales vs predichos de Velocidad de viento (Atillo).

Realizado por: Pilco V. y Acurio W., 2019.

Al comparar los datos reales vs predichos (Gráfico 17-4) presentan un comportamiento casi parecido, evidenciándose un mejor ajuste en los 150 primeros pronósticos, a partir de aquí se observa muchos desfases. Se obtiene un coeficiente de determinación 60.01% lo que indica la existencia de una correlación moderada, considerando que no es un buen modelo de predicción.

Tabla 11-4: Resumen de los modelos adecuados para cada Estación Meteorológica (Temperatura y Velocidad de Viento) ARIMA.

Estación	Variable	Modelo	Medidas de Escala				Medidas Basada en Porcentajes						Criterios de Información	
			MSE	RMSE	MAE	MdAE	MAPE	MdAPE	RMSPE	RMdSPE	Smape	sMdAPE	AIC	BIC
Cumandá	X ₁	(2,0,0) (0,1,5) ₂₄	2.8339	1.6834	1.3257	1.1119	0.0527	0.0468	0.0652	0.0468	0.0541	0.0476	30114.62	30176.43
San Juan	X ₁	(1,0,1) (3,1,2) ₂₄	3.2125	1.7923	1.4349	1.2878	0.1535	0.1208	0.2118	0.1208	0.1472	0.1231	56644.29	56709.47
	X ₁₄	(1,0,2) (1,1,1) ₂₄	1.7098	1.3076	1.0715	0.9728	1.3273	0.4682	3.0997	0.4682	0.5564	0.4684	21263.24	21305.18
Tixán	X ₁	(2,0,0) (1,1,2) ₂₄	2.2568	1.5023	1.1800	0.9864	0.1633	0.1230	0.2381	0.1230	0.1488	0.1205	59868.1	59916.99
	X ₁₄	(2,1,2) (0,0,2) ₂₄	7.5134	2.7411	2.4246	2.5032	1.7137	0.8935	3.3698	0.8935	0.6670	0.6176	56214.59	56268.69
Tunshi	X ₁	(2,0,1) (2,1,2) ₂₄	23.1682	4.8133	3.7199	2.9016	0.2966	0.2196	0.4244	0.2196	0.2671	0.2224	56301.63	56366.81
	X ₁₄	(1,0,1) (3,1,2) ₂₄	0.5578	0.7469	0.5326	0.3639	0.5378	0.2961	1.0987	0.2961	0.3983	0.3194	15829.6	15885.49
Urbina	X ₁	(1,0,1) (2,1,1) ₂₄	2.1338	1.4608	1.1727	1.0452	0.1978	0.1393	0.3047	0.1393	0.1767	0.1411	53415.48	53464.36
	X ₁₄	(1,1,2) (1,0,1) ₂₄	5.0792	2.2537	1.8331	1.6735	1.5961	0.6201	3.1588	0.6201	0.5939	0.4734	46521.02	46567.39
Alao	X ₁	(5,0,1) (0,1,9) ₂₄	11.5799	3.4029	2.3804	1.6234	0.2279	0.1526	0.3694	0.1526	0.2235	0.1591	74541.99	74676.38
	X ₁₄	(2,0,1) (1,1,3) ₂₄	1.2847	1.1334	0.8961	0.7550	2.2322	0.4467	9.3118	0.4467	0.5892	0.4539	42285.8	42347.62
Atillo	X ₁	(1,0,1) (1,1,1) ₂₄	26.2898	5.1274	4.9992	5.0543	6.5193	2.8379	17.4911	2.8379	1.1491	1.1732	48534.18	48574.62
	X ₁₄	(3,0,0) (1,1,1) ₂₄	1.6477	1.2836	1.0883	1.0809	1.7074	0.5028	5.2894	0.5028	0.5816	0.4467	20940.46	20982.39
EsPOCH	X ₁	(2,0,4) (2,1,1) ₂₄	7.2879	2.6996	1.9032	1.4519	0.1534	0.1088	0.2449	0.1088	0.1471	0.1129	55962.54	56043.42
Matus	X ₁	(2,0,1) (1,1,1) ₂₄	6.5318	2.5557	1.8125	1.2658	0.1583	0.0988	0.2512	0.0988	0.1474	0.0997	54221.86	54267.95
	X ₁₄	(1,0,1) (1,1,1) ₂₄	0.5719	0.7562	0.5910	0.4824	0.9837	0.3542	4.4115	0.3542	0.4730	0.3792	11695.38	11730.31
Multitud	X ₁	(4,0,0) (1,1,1) ₂₄	2.9426	1.7154	1.2594	0.8372	0.0850	0.0557	0.1153	0.0558	0.0791	0.0543	35685.53	35739.62
Quimiag	X ₁	(4,0,0) (1,1,1) ₂₄	7.3202	2.7056	1.9110	1.2229	0.1535	0.0983	0.2335	0.0983	0.1463	0.0975	78449.75	78508.7
	X ₁₄	(3,0,1) (1,1,1) ₂₄	0.6468	0.8043	0.5876	0.4596	0.7589	0.2685	1.8880	0.2685	0.4032	0.2830	27068.24	27122.33

Realizado por: Pilco V. y Acurio W., 2019.

Para las estaciones restantes se aplicó el mismo procedimiento anteriormente indicado, para el análisis de las series de cada estación meteorológica, todas presentaron estacionalidad (Tabla 5-4). Se realizó la diferencia dividida correspondiente debido a que en la correlación parcial simple todas las series evidencian estacionalidad cada 24 rezagos, indicios de modelos SARIMA.

Se indica el mejor modelo (Tabla 11-4) para cada estación meteorológica, los mismos que fueron seleccionados bajo el criterio de evaluación, de información, mediante los errores más bajos de cada modelo, en algunos se debió aplicar el test de Diebold-Mariano para saber si presentaban o no la misma exactitud entre modelos y en el caso de tener igual precisión se eligió por el criterio de parsimonia.

Estación: Cumandá

Para la variable temperatura (Tabla 11-4) el mejor modelo es $(2,0,0) (0,1,5)_{24}$, seleccionado mediante los criterios de información y por la gráfica de los errores (Anexo G) debido a que muestra menor error de predicción.

Estación: San Juan

Para la variable temperatura (Tabla 11-4) el mejor modelo es $(1,0,1) (3,1,2)_{24}$, seleccionado mediante los criterios de evaluación (RMSPE), de información y por la gráfica de los errores (Anexo G) debido a que muestra menor error de predicción.

En cuanto, para la variable velocidad de viento (Tabla 11-4) el mejor modelo es $(1,0,2) (1,1,1)_{24}$, seleccionado mediante los criterios de información, además, en la gráfica de los errores (Anexo G) se nota que la altura de las barras son similares, por lo que se procede a realizar el test de Diebold-Mariano obteniendo un valor p de 0.2208 mayor al nivel de significancia, por lo tanto, se concluye que ambos modelos tienen la misma exactitud de predicción y mediante el principio de parsimonia se seleccionó dicho modelo.

Estación: Tixán

Para la variable temperatura (Tabla 11-4) el mejor modelo es $(2,0,0) (1,1,2)_{24}$, seleccionado mediante los criterios de evaluación (MdAPE, RMSPE, RMdSPE), de información (BIC), además, en la gráfica de los errores (Anexo G) se nota que la altura de las barras son similares, por lo que se procede a realizar el test de Diebold-Mariano obteniendo un valor p de 0.3346 mayor

al nivel de significancia, por lo tanto, se concluye que ambos modelos tienen la misma exactitud de predicción y mediante el principio de parsimonia se seleccionó dicho modelo.

En cuanto, para la variable velocidad de viento (Tabla 11-4) el mejor modelo es (2,1,2) (0,0,2)₂₄, seleccionado mediante los criterios de evaluación, de información y por la gráfica de los errores (Anexo G) se nota que la altura de las barras son similares, por lo que se procede a realizar el test de Diebold-Mariano obteniendo un valor p de 0.1361 mayor al nivel de significancia, por lo tanto, se concluye que ambos modelos tienen la misma exactitud de predicción y mediante las características propuestas se seleccionó dicho modelo.

Estación: Tunshi

Para la variable temperatura (Tabla 11-4) el mejor modelo es (2,0,1) (2,1,2)₂₄, seleccionado mediante los criterios de información y por la gráfica de los errores (Anexo G) debido a que muestra menor error de predicción.

En cuanto, para la variable velocidad de viento (Tabla 11-4) el mejor modelo es (1,0,1) (3,1,2)₂₄, seleccionado mediante los criterios de evaluación, de información y por la gráfica de los errores (Anexo G) debido a que muestra menor error de predicción.

Estación: Urbina

Para la variable temperatura (Tabla 11-4) el mejor modelo es (1,0,1) (2,1,1)₂₄, seleccionado mediante los criterios de evaluación (MSE, RMSPE, MAE, MAPE, RMSPE), de información y por la gráfica de los errores (Anexo G) debido a que muestra menor error de predicción.

En cuanto, para la variable velocidad de viento (Tabla 11-4) el mejor modelo es (1,1,2) (1,0,1)₂₄, seleccionado mediante los criterios de evaluación (MAE, MdAE, MAPE, RMSPE, Smape, sMdAPE), de información y por la gráfica de los errores (Anexo G) debido a que muestra menor error de predicción.

Estación: Alao

Para la variable temperatura (Tabla 11-4) el mejor modelo es (5,0,1) (0,1,9)₂₄, seleccionado mediante los criterios de evaluación (RMSPE), de información y por la gráfica de los errores (Anexo G) debido a que muestra menor error de predicción.

En cuanto, para la variable temperatura (Tabla 11-4) el mejor modelo es (2,0,1) (1,1,3)₂₄, seleccionado mediante los criterios de evaluación, de información y por la gráfica de los errores (Anexo G) debido a que muestra menor error de predicción.

Tabla 12-4: Valores p de los supuestos de cada modelo de las estaciones meteorológicas (Temperatura y Velocidad de Viento).

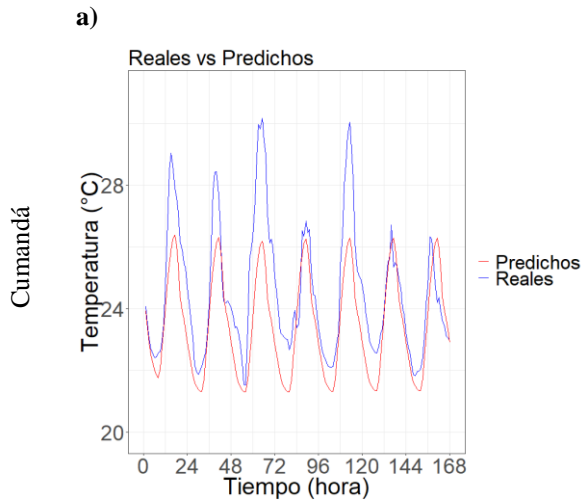
Estación	Variable	Modelo	Supuestos			
			Independencia	Estacionariedad	Normalidad	Homocedasticidad
Cumandá	X ₁	(2,0,0) (0,1,5) ₂₄	0.0002855	0.01	2.20E-16	1
San Juan	X ₁	(1,0,1) (3,1,2) ₂₄	2.20E-16	0.01	2.20E-16	0.9319
	X ₁₄	(1,0,2) (1,1,1) ₂₄	1.47E-14	0.01	2.20E-16	9.24E-05
Tixán	X ₁	(2,0,0) (1,1,2) ₂₄	2.20E-16	0.01	2.20E-16	1
	X ₁₄	(2,1,2) (0,0,2) ₂₄	2.20E-16	0.01	2.20E-16	0.1704
Tunshi	X ₁	(2,0,1) (2,1,2) ₂₄	2.20E-16	0.01	2.20E-16	1
	X ₁₄	(1,0,1) (3,1,2) ₂₄	2.20E-16	0.01	2.20E-16	0.9903
Urbina	X ₁	(1,0,1) (2,1,1) ₂₄	2.20E-16	0.01	2.20E-16	0.2986
	X ₁₄	(1,1,2) (1,0,1) ₂₄	3.18E-04	0.01	2.20E-16	0.1737
Alao	X ₁	(5,0,1) (0,1,9) ₂₄	2.20E-16	0.01	2.20E-16	1
	X ₁₄	(2,0,1) (1,1,3) ₂₄	2.20E-16	0.01	2.20E-16	0.1105
Atillo	X ₁	(1,0,1) (1,1,1) ₂₄	2.20E-16	0.01	2.20E-16	2.20E-16
	X ₁₄	(3,0,0) (1,1,1) ₂₄	7.15E-08	0.01	2.20E-16	2.20E-16
Epoch	X ₁	(2,0,4) (2,1,1) ₂₄	2.20E-16	0.01	2.20E-16	1
Matus	X ₁	(2,0,1) (1,1,1) ₂₄	2.20E-16	0.01	2.20E-16	1
	X ₁₄	(1,0,1) (1,1,1) ₂₄	2.20E-16	0.01	2.20E-16	1
Multitud	X ₁	(4,0,0) (1,1,1) ₂₄	2.20E-16	0.01	2.20E-16	2.20E-16
Quimiag	X ₁	(4,0,0) (1,1,1) ₂₄	2.20E-16	0.01	2.20E-16	2.20E-16
	X ₁₄	(3,0,1) (1,1,1) ₂₄	2.24E-09	0.01	2.20E-16	0.9727

Realizado por: Pilco V. y Acurio W., 2019.

Se analizó los supuestos (Tabla 12-4) a un nivel de significancia de 0,05 de cada una de las variables, por lo tanto, en ninguna estación meteorológica se encontró independencia y normalidad; para la variable temperatura en las estaciones de: Cumandá, San Juan, Tixán, Tunshi, Urbina, Alao, Epoch y Matus se cumple los supuestos de estacionariedad y homocedasticidad, mientras que en: Atillo, Multitud y Quimiag solo se verifica la estacionariedad.

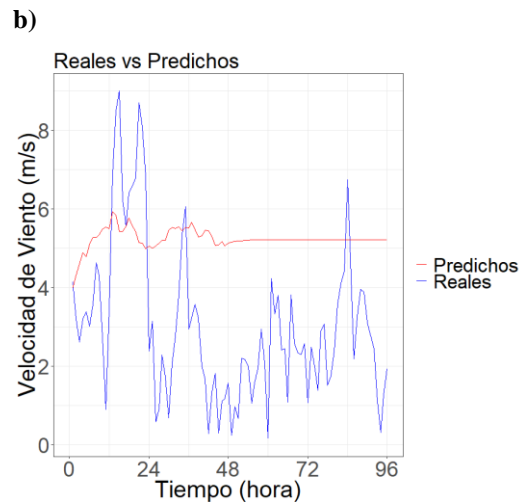
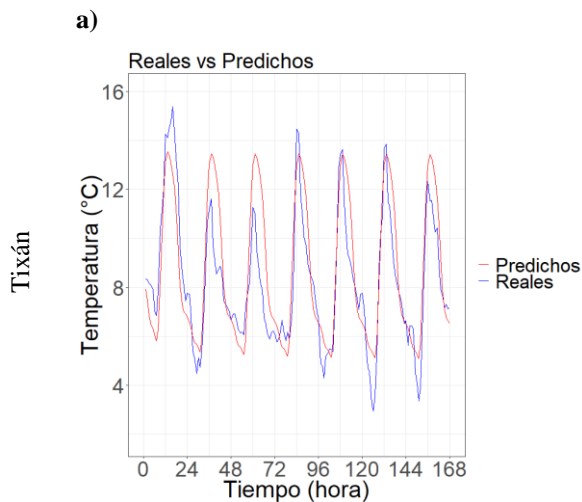
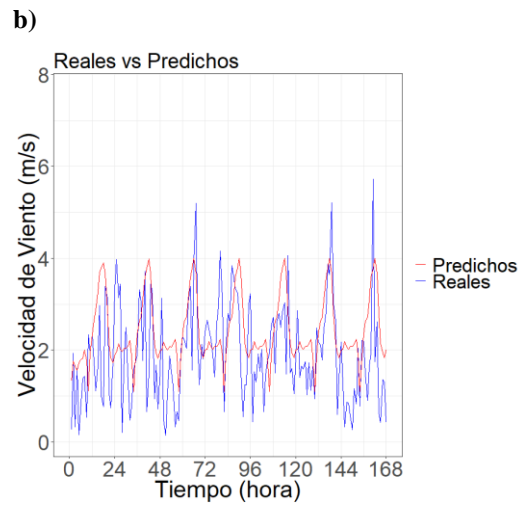
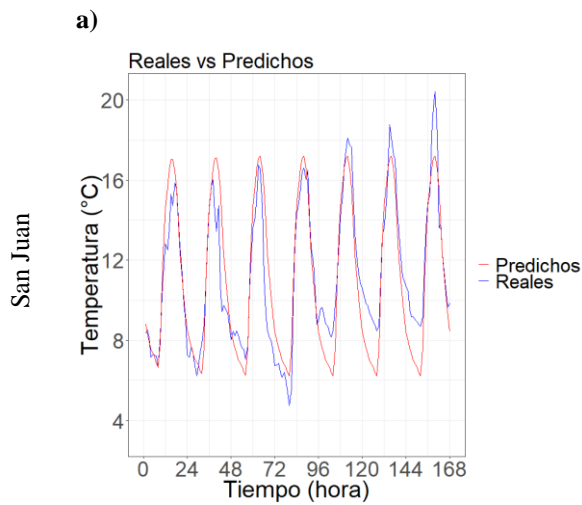
Para la variable velocidad de viento en las estaciones de: Tixán, Tunshi, Urbina, Alao, Matus y Quimiag se cumple los supuestos de estacionariedad y homocedasticidad, mientras que en: San Juan y Atillo solo se verifica la estacionariedad.

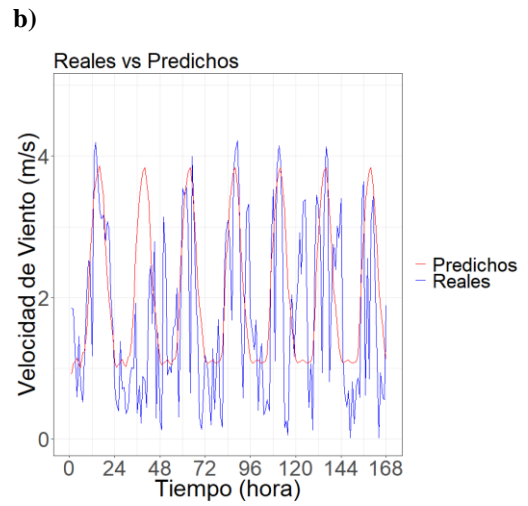
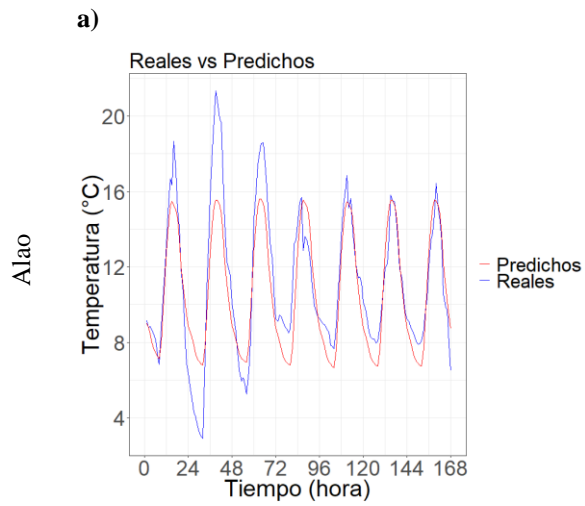
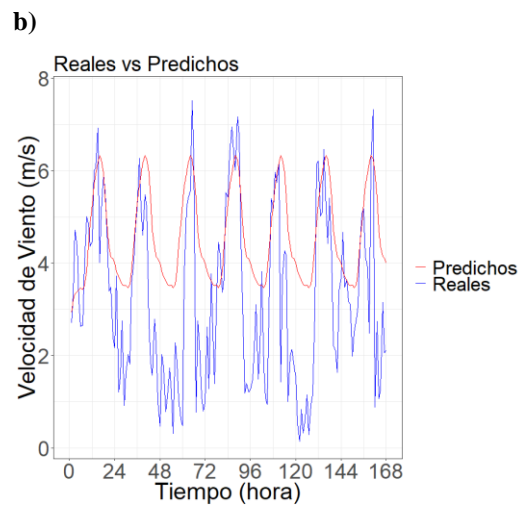
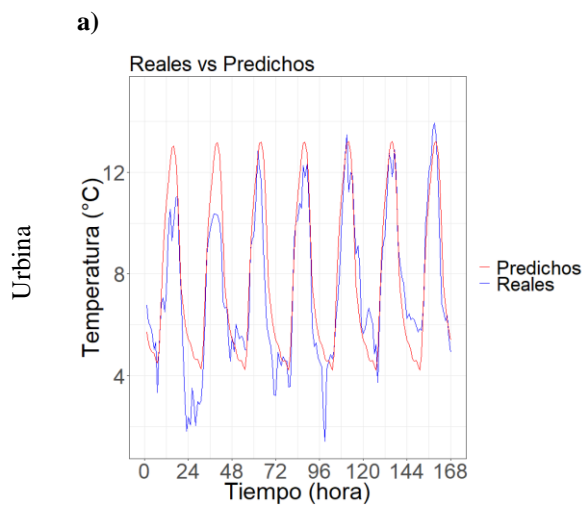
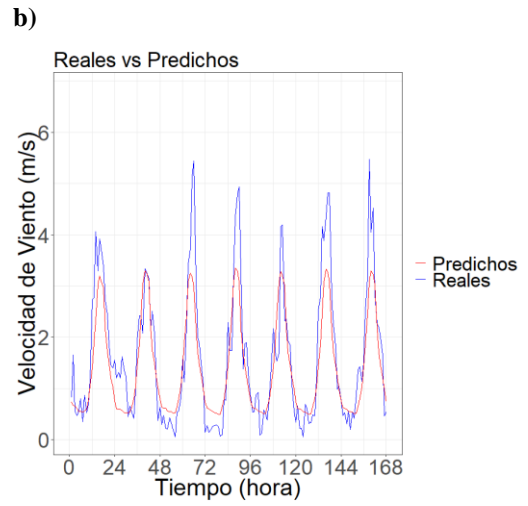
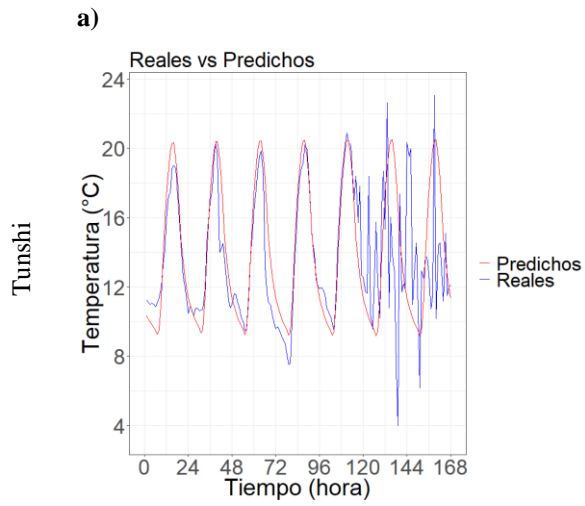
Pronósticos con los modelos seleccionados Box-Jenkins (ARIMA).

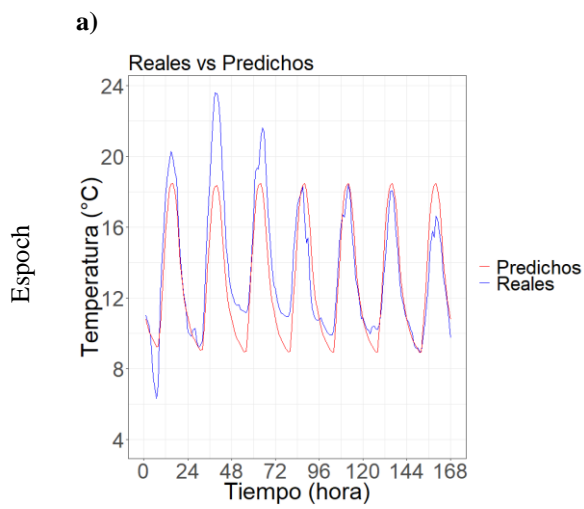
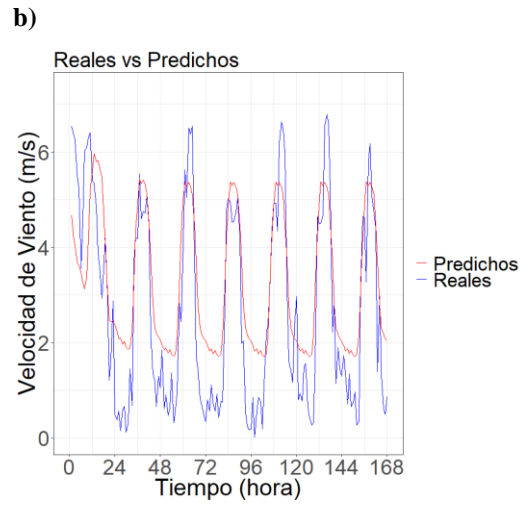
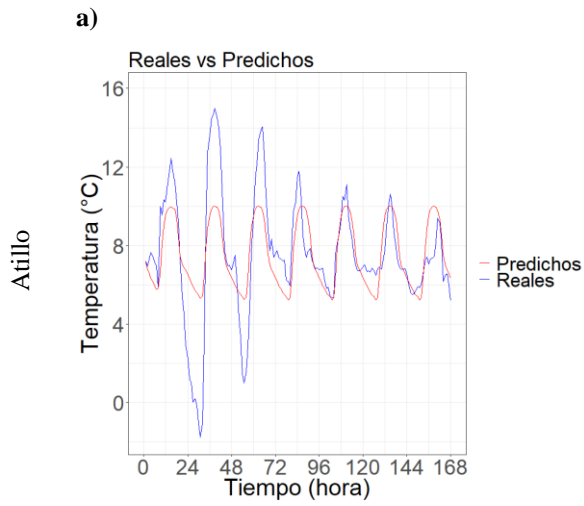


b)

No se realizó el análisis en esta variable, dado que el coeficiente de determinación ajustado es menor al 79% en la imputación y por lo tanto no se considera adecuada.

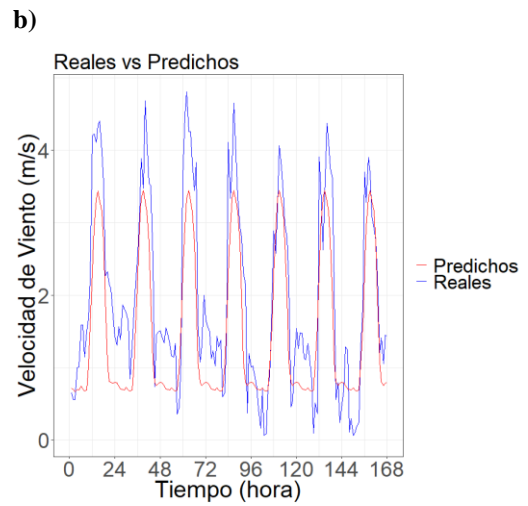
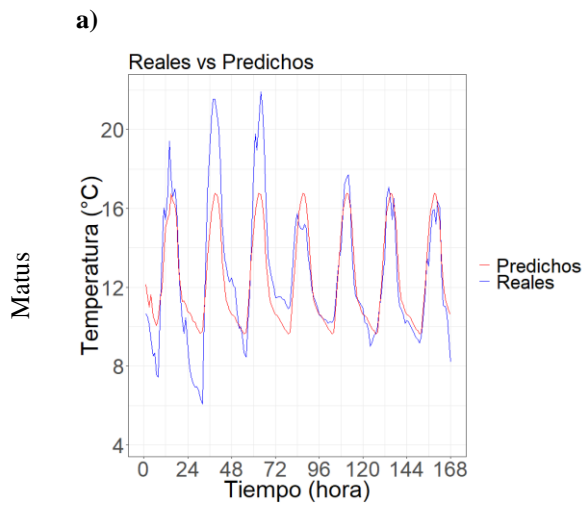


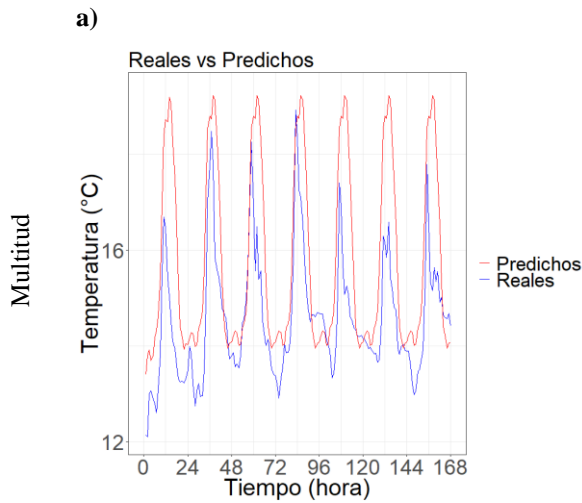




b)

No se realizó el análisis en esta variable, dado que el coeficiente de determinación ajustado es menor al 79% en la imputación y por lo tanto no se considera adecuada.





b)

No se realizó el análisis en esta variable, dado que el coeficiente de determinación ajustado es menor al 79% en la imputación y por lo tanto no se considera adecuada.

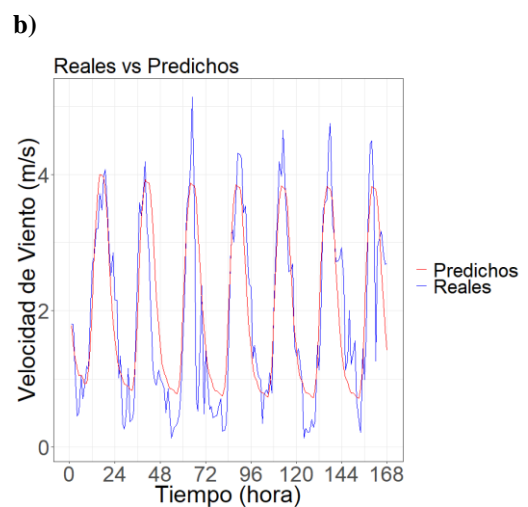
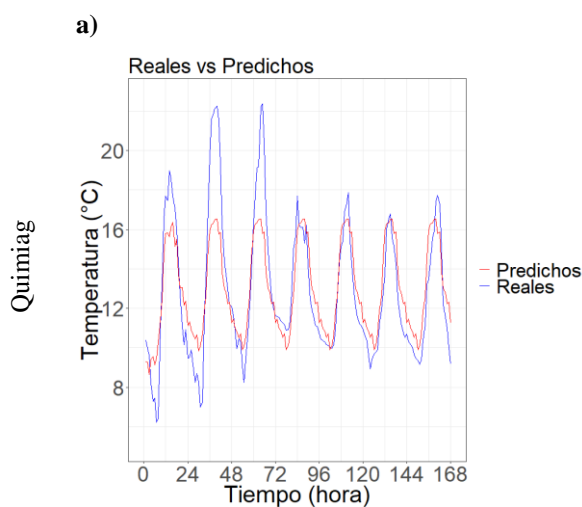


Gráfico 18-4: Datos reales vs predichos de los modelos Box-Jenkins (ARIMA) para cada estación meteorológica.

Realizado por: Pilco V. y Acurio W., 2019.

Al seleccionar el modelo más adecuado y luego de evaluar los supuestos del mismo, se procedió a realizar las debidas predicciones con los modelos ARIMA de cada estación (Gráfico 18-4), se puede observar que existe variabilidad en las series, por lo tanto, los valores predichos son casi similares a las observaciones originales. Se tomaron los datos correspondientes a una semana (168 horas), se notó que en la mayoría de las estaciones los datos siguen el mismo patrón, pero no presenta un buen ajuste ya que las predicciones no logran llegar hasta los repuntes más altos, al inicio se observa que tienden a ser similares las dos series, por lo tanto, la modelación ARIMA es válida para predicciones a corto plazo.

4.1.5 *Teoría del Caos*

Para procesar los datos se usó TISEAN 3.0.1 (Rainer H. et al., 2007), el mismo que permite analizar series de tiempo no lineales, es un software libre, en el cual se obtuvieron parámetros como: tiempo de retardo en el programa *mutuas*, dimensión de encaje en *false_nearest*, reducción de ruido en *ghkss* y las respectivas predicciones en *Rbf*.

Para la reducción de ruido se aplicó 10 iteraciones para la variable temperatura y velocidad de viento, de cada una se obtuvo un modelo analizados bajo los criterios de evaluación para seleccionar el más adecuado. Las leyes del caos explican en gran parte aquellos fenómenos que tienen su origen en la naturaleza y supone que el 90% son de carácter no lineal, dependiendo de las condiciones iniciales (Sánchez & Garduño, 2007, pp. 184-185).

Los fenómenos naturales presentan comportamiento inestable y difícil de predecir como: el clima, tormentas marinas, variación de velocidad de viento en la atmósfera, temperatura, etc., los cuales presentan no linealidad.

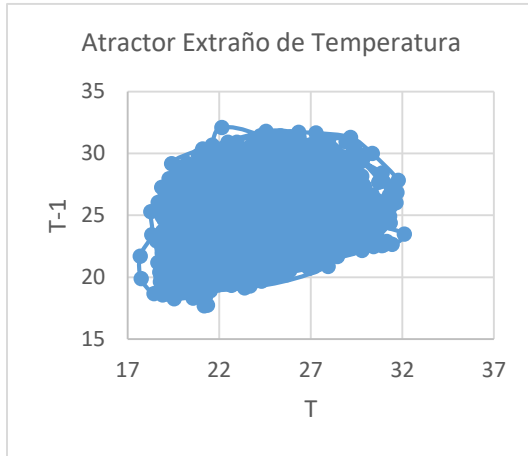
La reducción de ruido se debe utilizar dependiendo del tipo de datos a estudiar, teniendo en cuenta, por ejemplo: en el caso de una serie como temperatura solo limpia la señal sin que el conjunto de datos se altere, mientras que en una serie de datos electrocardiográficos hace que la misma se altere. A medida que el ruido aumenta la serie tiende a ser más definida.

Leith (1974) y Lorenz (1993) indica que la condición inicial se vuelve crítica en pronósticos meteorológicos que van más allá de dos semanas (Sánchez & Garduño, 2007, pp. 184-185). Se puede predecir series temporales caóticas sin una precisión considerable y si el intervalo de predicción es grande aumenta la desviación en cuanto a los datos originales.

El procedimiento es el siguiente:

Mediante el ejecutable (Anexo H) se ingresan los datos para obtener los parámetros necesarios como: el número de retardos (d) y la dimensión de encaje (m), y l es el número de datos de cada variable existente en cada una de las estaciones (Tabla 5-H), se muestra el análisis en la estación de Atillo:

a)



b)

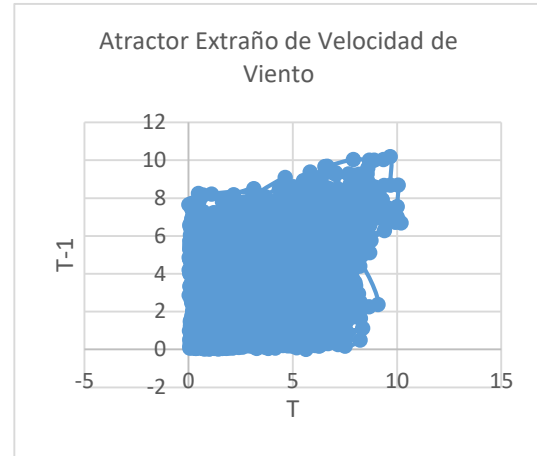
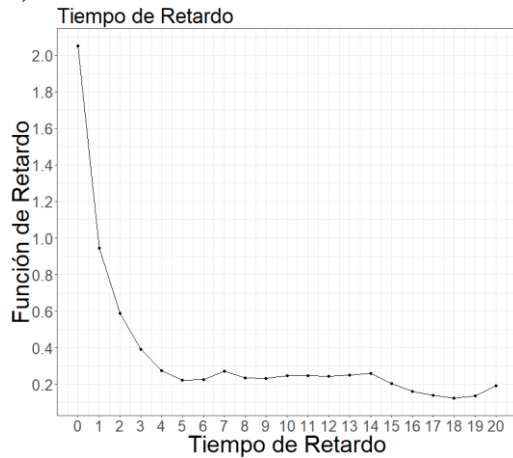


Gráfico 19-4: Atractor extraño de a) Temperatura y b) Velocidad de viento (Atillo).
Realizado por: Pilco V. y Acurio W., 2019.

El atractor extraño (Gráfico 19-4) regularmente muestra alguna simetría, se notó que el atractor de temperatura y velocidad de viento no evidencia tendencia por lo cual se procede a realizar reducción de ruido (Anexo H).

a)



b)

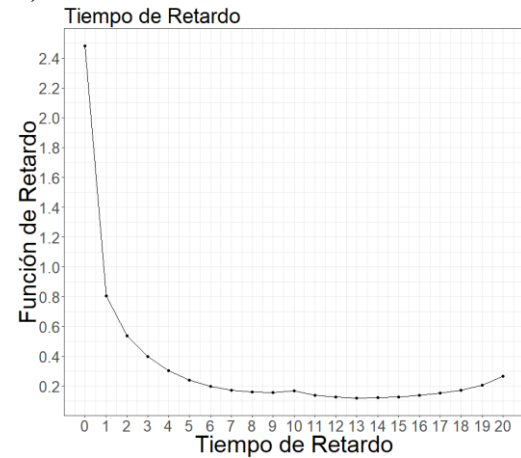
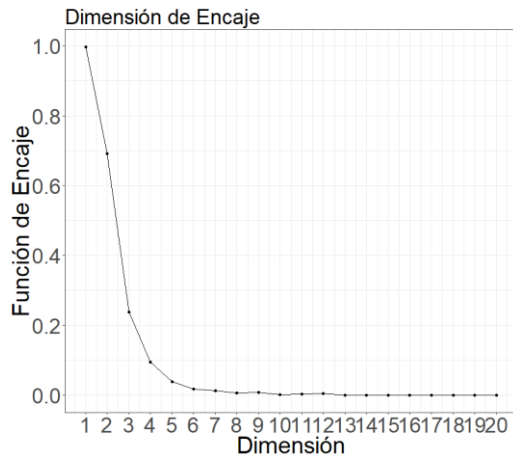


Gráfico 20-4: Tiempo de retardo sin reducción de ruido de a) Temperatura y b) Velocidad de viento (Atillo).

Realizado por: Pilco V. y Acurio W., 2019.

Para determinar el tiempo de retardo (Gráfico 20-4) se debe observar el primer valor menor de la función de correlación de retardo, en este caso para temperatura es 5 y para velocidad de viento es 9. Se tomará como referencia este valor para encontrar la dimensión de encaje (Anexo I).

a)



b)

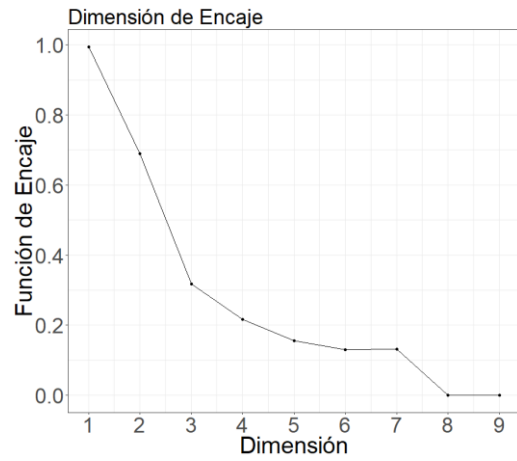


Gráfico 21-4: Dimensión de encaje con reducción de ruido de a) temperatura y b) Velocidad de Viento (Atillo).

Realizado por: Pilco V. y Acurio W., 2019.

La dimensión de encaje (Gráfico 21-4) es determinado por el valor igual a cero, para la temperatura el valor es 13 y para velocidad de viento es 8 (Anexo I).

Tabla 13-4: Criterios de evaluación para los posibles modelos de Temperatura Teoría del Caos (Atillo).

Reducción de ruido	Parámetros del Caos		Medidas de Escala				Medidas Basada en Porcentajes					
	Tiempo de Retardo	Dimensión de Encaje	MSE	RMSE	MAE	MdAE	MAPE	MdAPE	RMSPE	RMdSPE	Smape	sMdAPE
0	7	7	10.5563	3.2490	2.5481	2.0063	0.9305	0.2431	6.6260	0.2431	0.3306	0.2424
1	7	8	8.3030	2.8815	2.1087	1.3989	0.9938	0.1712	9.3882	0.1712	0.2843	0.1683
2	7	10	11.2511	3.3543	2.7300	2.5600	1.0417	0.3047	9.0833	0.3047	0.3603	0.3022
3	6	15	9.4140	3.0682	2.4314	2.0442	1.0766	0.2397	10.0728	0.2397	0.3276	0.2479
4	6	11	21.6095	4.6486	3.9648	3.6774	0.8063	0.5023	4.4845	0.5023	0.6562	0.6624
5	6	11	21.5298	4.6400	3.9619	3.7216	0.8199	0.4974	4.6493	0.4974	0.6544	0.6524
6	6	11	8.6641	2.9435	2.1867	1.5307	0.9805	0.1856	8.5261	0.1856	0.2926	0.1825
7	6	11	8.8925	2.9820	2.3709	2.0126	0.9347	0.2545	8.7252	0.2545	0.3392	0.2670
8	6	11	8.9741	2.9957	2.3597	1.9806	0.9376	0.2476	8.4588	0.2476	0.3342	0.2444
9	5	13	7.0857	2.6619	1.8903	1.3159	0.8783	0.1576	8.0087	0.1576	0.2695	0.1623
10	5	10	10.4847	3.2380	2.5282	2.1265	1.0748	0.2588	7.3468	0.2588	0.3322	0.2401

Realizado por: Pilco V. y Acurio W., 2019.

Tabla 14-4: Criterios de evaluación para los posibles modelos de Velocidad de Viento Teoría del Caos (Atillo).

Reducción de ruido	Parámetros del Caos		Medidas de Escala				Medidas Basada en Porcentajes					
	Tiempo de Retardo	Dimensión de Encaje	MSE	RMSE	MAE	MdAE	MAPE	MdAPE	RMSPE	RMdSPE	Smape	sMdAPE
0	10	---	---	---	---	---	---	---	---	---	---	---
1	13	7	8.7560	2.9591	2.4514	2.2891	3.6543	1.0127	12.7700	1.0127	0.9174	0.8916
2	13	8	9.3042	3.0503	2.5223	2.2299	4.3231	0.9861	14.1459	0.9861	0.8561	0.7962
3	9	8	7.8407	2.8001	2.3519	2.3449	3.0945	0.9421	10.1479	0.9421	1.0067	1.0612
4	8	9	9.7597	3.1241	2.6669	2.5265	4.6140	1.3628	13.7692	1.3628	0.8885	0.8702
5	8	11	8.4897	2.9137	2.4323	2.2251	4.0299	0.8092	12.9612	0.8092	0.9342	0.8990
6	8	12	9.3385	3.0559	2.5785	2.4617	4.4155	1.3261	12.8874	1.3261	0.8697	0.8323
7	8	12	9.0878	3.0146	2.5514	2.4554	4.2102	1.3111	11.5952	1.3112	0.8727	0.8347
8	8	12	9.3712	3.0612	2.5901	2.4590	4.2349	1.3897	12.1136	1.3897	0.8783	0.8609
9	8	12	8.9828	2.9971	2.5313	2.3266	4.0990	1.2771	12.0713	1.2771	0.8755	0.8673
10	8	12	8.3131	2.8832	2.4154	2.1940	3.9474	1.1038	11.7449	1.1038	0.8568	0.8497

Realizado por: Pilco V. y Acurio W., 2019.

En los espacios en blanco no se trabajó dado que no se encontraron valores del parámetro dimensión de encaje (m) debido a que no existe autocorrelación entre cada uno de los datos de la serie, por lo cual se elegirá el mejor modelo entre los presentados para temperatura y velocidad de viento.

Analizando los criterios de evaluación puedo concluir que el mejor modelo para temperatura (Tabla 14-4) se encuentra en la iteración 1 y en cuanto a velocidad de viento (Tabla 15-4) no se encuentra necesario realizar reducción de ruido.

Predicción

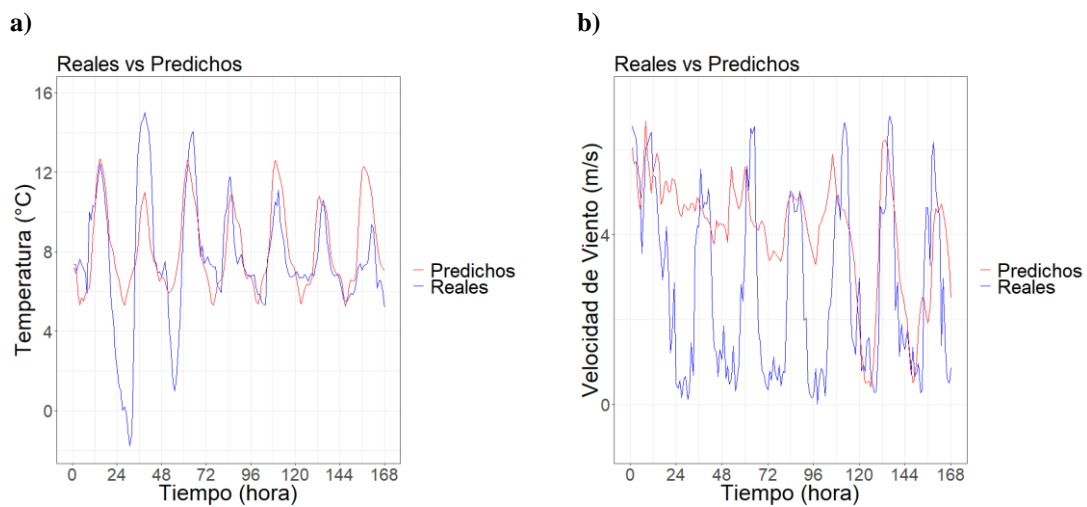


Gráfico 22-4: Datos reales vs predichos de a) Temperatura y b) Velocidad de viento (Atillo).
Realizado por: Pilco V. y Acurio W., 2019.

Se observa (Gráfico 22-4) que los datos reales de a) temperatura se ajustan con los predichos obtenidos por el modelo de Teoría del caos al comienzo de la serie y para b) velocidad de viento datos reales, recordando que estos modelos obtenidos son con reducción de ruido.

Tabla 15-4: Criterios de evaluación de cada uno de los modelos de las diferentes estaciones meteorológicas (Teoría del Caos).

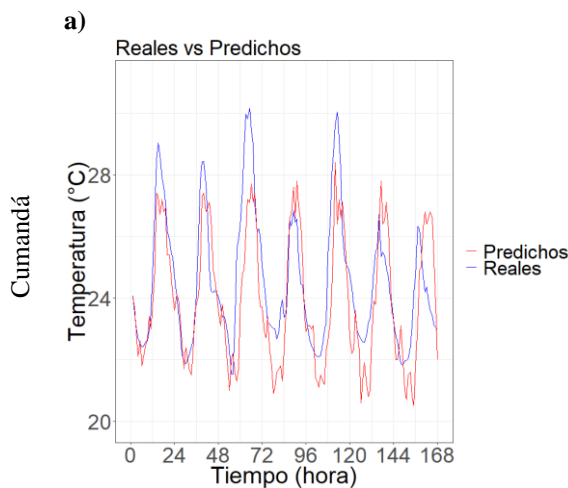
Estación	Variable	Reducción de ruido	Parámetros del Caos		Medidas de Escala				Medidas Basada en Porcentajes					
			Tiempo de Retardo	Dimensión de Encaje	MSE	RMSE	MAE	MdAE	MAPE	MdAPE	RMSPE	RMdSPE	Smape	sMdAPE
Cumandá	X ₁	9	5	12	6.285	2.507	1.785	1.114	0.070	0.048	0.094	0.048	0.074	0.048
San Juan	X ₁	9	5	12	4.013	2.003	1.542	1.158	0.163	0.109	0.261	0.109	0.150	0.111
	X ₁₄	9	5	8	2.044	1.430	1.032	0.780	0.986	0.433	2.318	0.433	0.560	0.477
Tixán	X ₁	0	6	8	3.895	1.974	1.530	1.214	0.206	0.158	0.291	0.158	0.190	0.161
	X ₁₄	6	8	10	3.778	1.944	1.494	1.188	0.964	0.387	2.428	0.387	0.501	0.384
Tunshi	X ₁	1	7	11	11.548	3.398	2.835	2.550	0.232	0.196	0.328	0.196	0.210	0.191
	X ₁₄	5	5	10	0.762	0.873	0.638	0.445	0.581	0.392	1.108	0.392	0.519	0.429
Urbina	X ₁	1	6	8	3.848	1.962	1.598	1.395	0.236	0.189	0.320	0.189	0.250	0.204
	X ₁₄	2	9	7	3.833	1.958	1.572	1.343	0.864	0.503	1.745	0.503	0.594	0.558
Alao	X ₁	6	7	11	11.236	3.352	2.346	1.583	0.232	0.140	0.384	0.140	0.216	0.146
	X ₁₄	7	5	9	2.735	1.654	1.343	1.130	3.811	0.649	20.361	0.649	0.751	0.639
Atillo	X ₁	9	5	13	7.0857	2.6619	1.8903	1.3159	0.8783	0.1576	8.0087	0.1576	0.2695	0.1623
	X ₁₄	3	9	8	7.8407	2.8001	2.3519	2.3449	3.0945	0.9421	10.1479	0.9421	1.0067	1.0612
Espoch	X ₁	4	6	11	7.697	2.774	1.964	1.408	0.161	0.107	0.252	0.107	0.152	0.109
Matus	X ₁	4	4	12	8.449	2.907	2.189	1.646	0.189	0.130	0.277	0.130	0.179	0.133
	X ₁₄	5	6	8	0.899	0.948	0.733	0.574	0.988	0.416	3.155	0.416	0.560	0.477
Multitud	X ₁	5	6	14	1.749	1.322	1.002	0.760	0.068	0.052	0.087	0.052	0.067	0.051
Quimiag	X ₁	5	6	12	9.930	3.151	2.357	1.702	0.186	0.136	0.257	0.136	0.182	0.137
	X ₁₄	1	5	9	0.782	0.885	0.641	0.478	0.822	0.299	2.034	0.299	0.422	0.294

Realizado por: Pilco V. y Acurio W., 2019.

Para las estaciones restantes se aplicó el mismo procedimiento anteriormente indicado, se realizó las gráficas de los atractores (Gráfico 6-H), se determinó el tiempo de retardo y la dimensión de encaje. Después se utilizó estos parámetros (Tabla 5-H) se procedió a realizar los pronósticos, los cuales fueron analizados mediante los criterios de información y se seleccionó el mejor modelo (Tabla 15-4) para cada estación meteorológica.

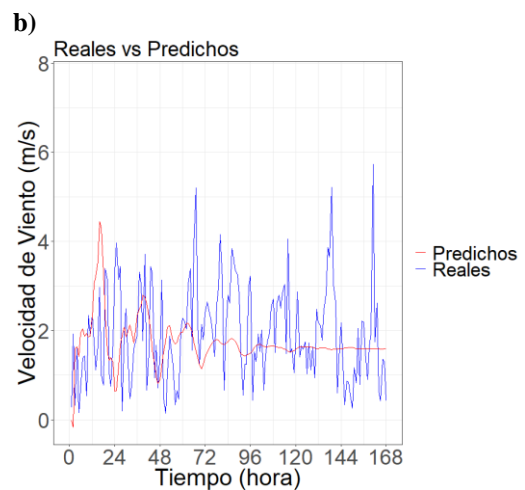
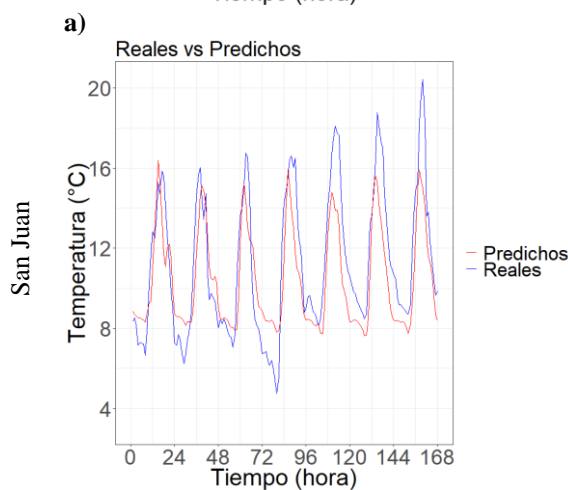
En casi todas las estaciones meteorológicas se obtuvieron mejores predicciones aplicando reducción de ruido, pero para la variable temperatura tanto en Tixán y Quimiag se obtuvieron mejores pronósticos sin reducción de ruido.

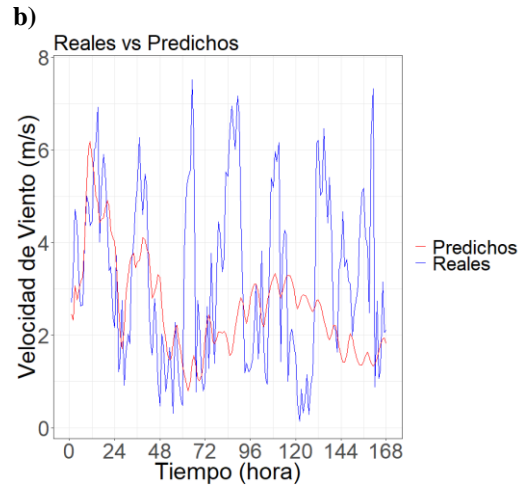
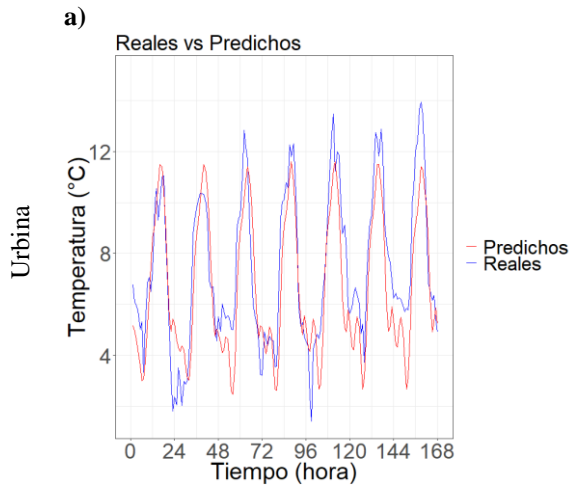
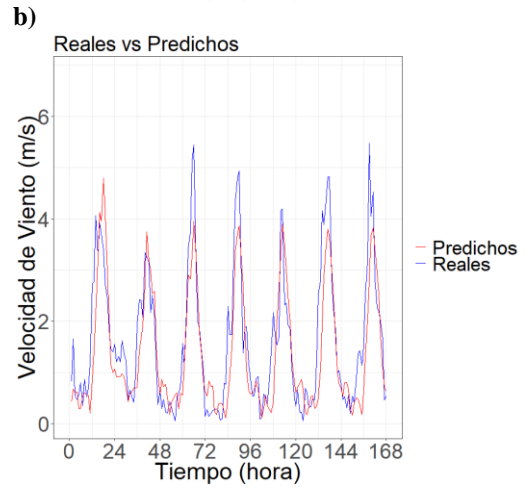
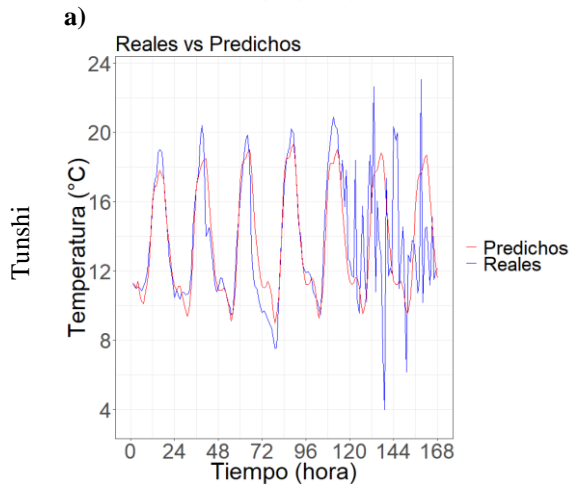
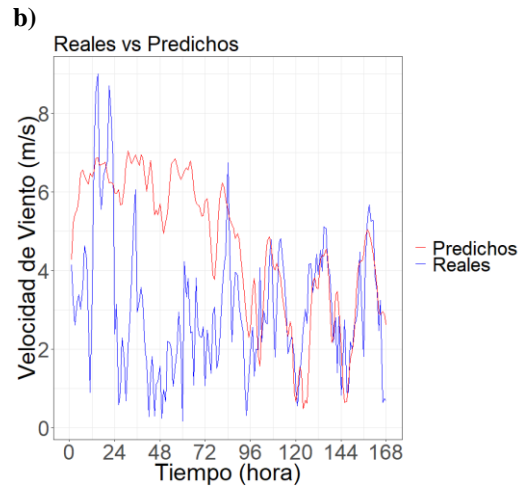
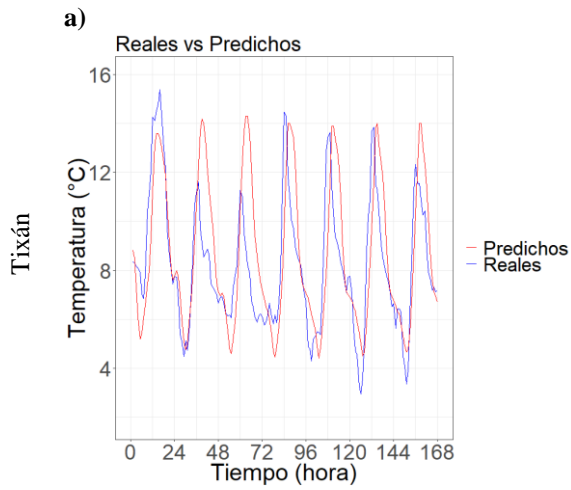
Pronósticos con los modelos seleccionados de Teoría del Caos

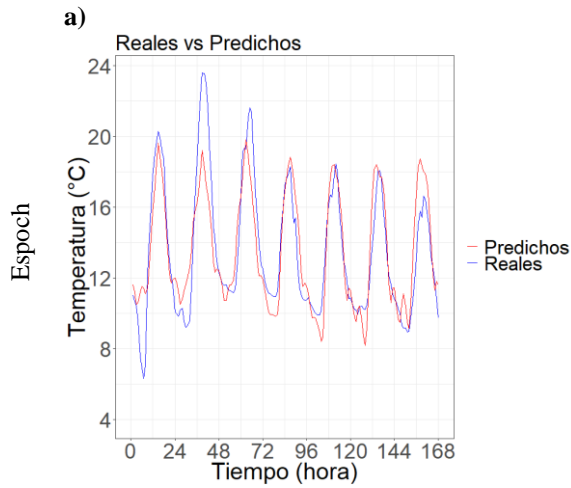
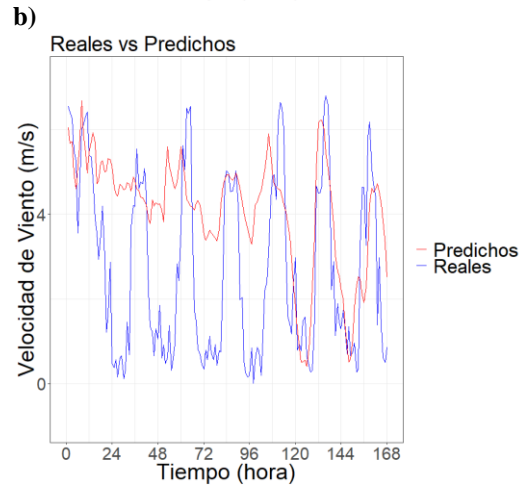
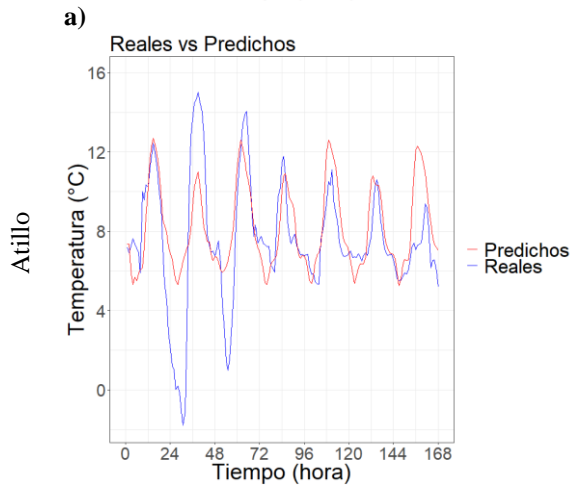
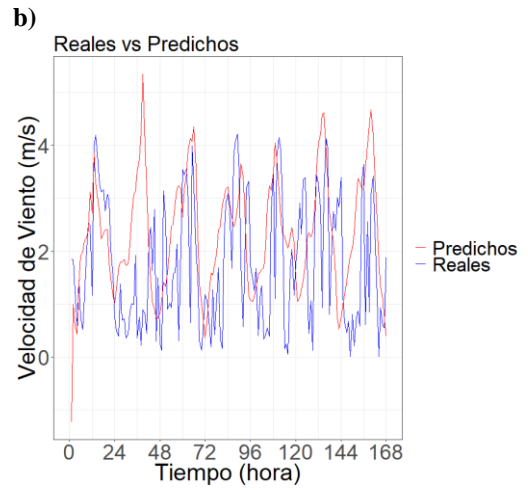
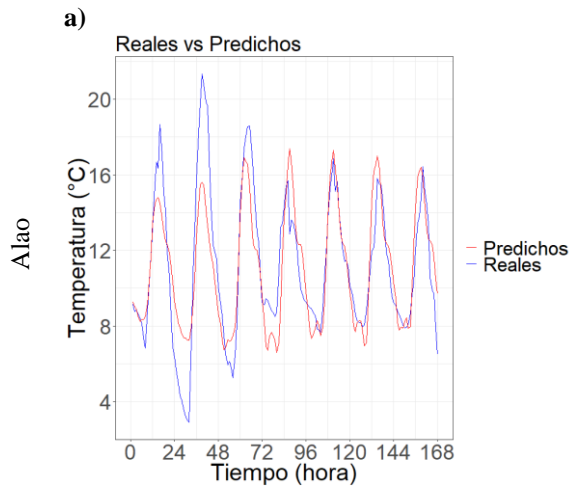


b)

No se realizó el análisis en esta variable, dado que el coeficiente de determinación ajustado es menor al 79% en la imputación y por lo tanto no se considera adecuada.







b)

No se realizó el análisis en esta variable, dado que el coeficiente de determinación ajustado es menor al 79% en la imputación y por lo tanto no se considera adecuada.

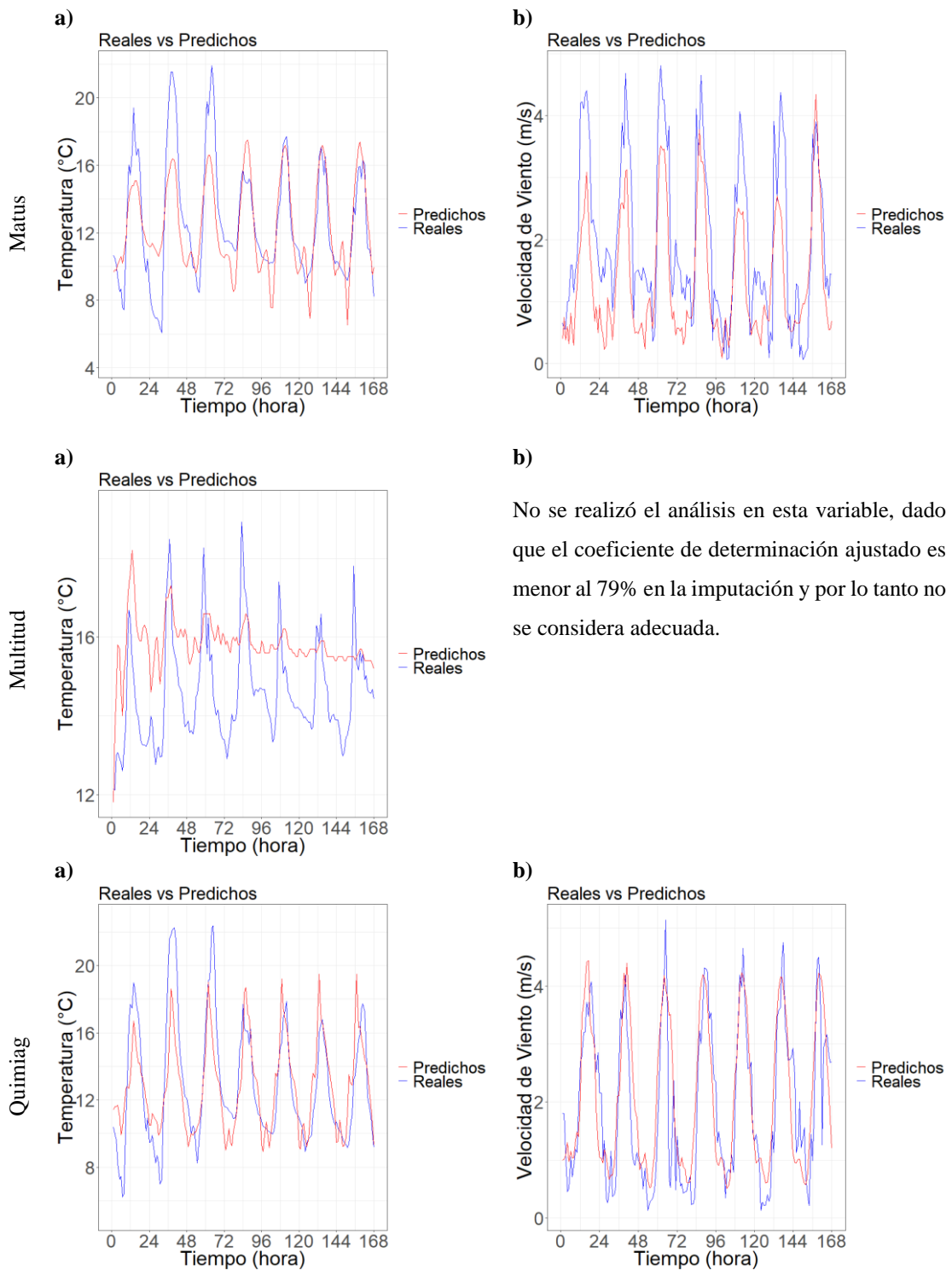


Gráfico 23-4: Datos reales vs predichos de los modelos de Teoría del Caos para cada estación meteorológica.

Realizado por: Pilco V. y Acurio W., 2019.

Al seleccionar el modelo más adecuado por Teoría del Caos se realiza sus gráficas (Gráfico 23-4), se observa que existe variabilidad entre los datos reales vs predichos para las dos variables en estudio. El análisis se realizará en un lapso de tiempo de una semana (168 horas), al comienzo el comportamiento de las dos series es aproximadamente similar.

Para temperatura se observa un mejor ajuste que para velocidad de viento, aunque en ambos casos no logran llegar a los repuntes más altos. Generando mejores pronósticos visualmente para temperatura debido a que velocidad de viento es una variable muy inestable.

4.1.6 *Redes Neuronales Recurrentes*

Las Redes Neuronales Recurrentes son adecuadas para modelar series temporales, se trabajó con redes Elman y Jordan, el análisis se realizó en el software libre R-Studio. Una de las ventajas que ofrece esta técnica es que no se necesita conocer a priori los supuestos como otros modelos (Qi y Zhang, 2001; citado en Sánchez, 2012, p. 22), su eficiencia en aprender el comportamiento y extrapolarlo ha logrado obtener mejores predicciones (De Gooijer y Kumar, 1992; citado en Sánchez, 2012, p. 22). Para dicha técnica no es obligatorio describir supuestos como la distribución de probabilidad, patrón de comportamiento de los datos de la serie para pronosticar de manera eficiente, además al ser utilizada para realizar predicciones presenta una característica significativa y es que los datos no serán desarrollados para probar supuestos como tendencia o estacionalidad en la serie previo a la realización de los pronósticos (Ruelas & Laguna, 2013, pp. 13-14).

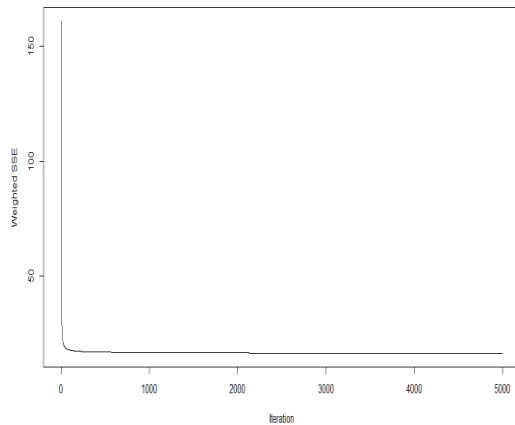
Cuando los supuestos de los métodos estadísticos paramétricos no se cumplen, no pueden ser demostrados, los supuestos de linealidad no se conservan, hace falta métodos no paramétricos que permitan suavizar los supuestos que necesitan los datos muestrales o predictores apoyándose en una deducción de linealidad, para esta problemática nacen las técnicas de Redes Neuronales como una opción para resolver problemas de regresión en las condiciones indicadas (Roldán, 2002, p. 2).

Es necesario realizar un cambio de escala a los datos originales debido a que los valores de la ecuación logística se encuentran en el rango de (0,1), se aplica la siguiente formula:

$$\hat{y} = \frac{x_i - x_{min}}{x_{max} - x_{min}}$$

Luego se procede a definir la periodicidad de los datos históricos, en nuestro caso cada 24 horas. La red de Elman está constituida por: una capa de entrada, capas ocultas y una capa de salida; tenemos dos capas ocultas, una de tres neuronas y la otra de dos neuronas, mientras tanto la red de Jordan está formada por 4 capas ocultas. Se ha establecido un ritmo de aprendizaje de 0.07 (González et, al., 2017, pp. 4-5), con un número máximo de 5000 iteraciones, tanto para temperatura y velocidad de viento. Se realizará el análisis en la estación de Atillo:

a)



b)

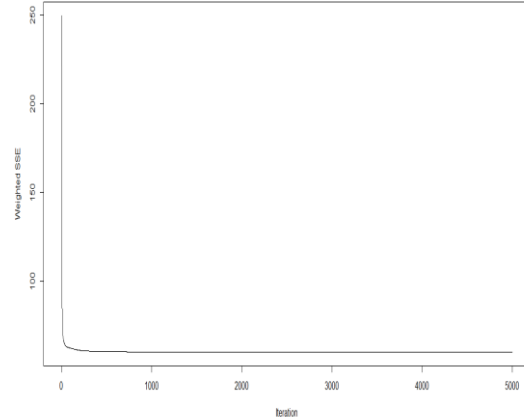


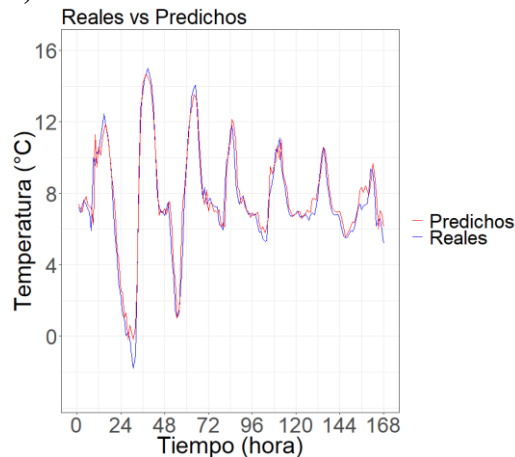
Gráfico 24-4: Función de error de la red para las variables a) Temperatura y b) Velocidad de viento (Atillo).

Realizado por: Pilco V. y Acurio W., 2019.

La función del error (Gráfico 24-4) para temperatura y velocidad de viento converge de manera muy rápida a cero, indicando el error con que la red neuronal realizara predicciones. Se presentan los modelos (Tabla 17-4) más adecuados para las respectivas predicciones según los criterios de evaluación.

Predicción

a)



b)

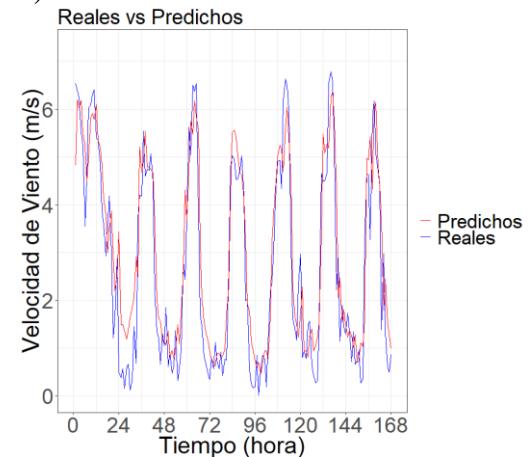


Gráfico 25-4: Datos reales vs predichos de a) Temperatura y b) Velocidad de viento (Atillo).

Realizado por: Pilco V. y Acurio W., 2019.

Al evaluar los pronósticos mediante los criterios de evaluación (Tabla 17-4) el mejor modelo generado por la red de Elman para las dos variables en estudio, logrando obtener un buen ajuste con los datos originales esto debido a que la red neuronal aprendió correctamente la función logística y puede proporcionar un siguiente valor aproximadamente con una buena precisión.

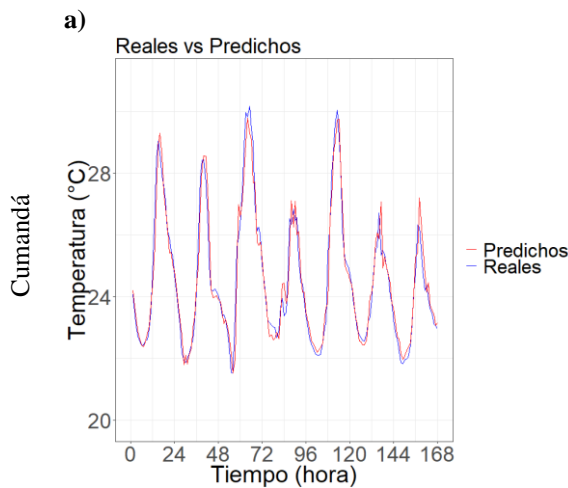
Tabla 16-4: Criterios de evaluación de cada uno de los modelos de las diferentes Estaciones Meteorológicas con RNR.

Estación	Variable	Red Neuronal	Medidas de Escala				Medidas Basada en Porcentajes					
			MSE	RMSE	MAE	MdAE	MAPE	MdAPE	RMSPE	RMdSPE	Smape	sMdAPE
Cumandá	X ₁	ELMAN 0.07	1.0560	1.0276	0.5754	0.3011	0.0235	0.0124	0.0428	0.0124	0.0234	0.0125
San Juan	X ₁	ELMAN 0.07	0.6781	0.8235	0.5967	0.4247	0.0578	0.0410	0.0827	0.0410	0.0565	0.0409
	X ₁₄	ELMAN 0.07	1.0897	1.0439	0.7856	0.6253	0.8745	0.3312	2.1916	0.3312	0.4511	0.3451
Tixán	X ₁	ELMAN 0.07	0.9015	0.9495	0.7543	0.6411	0.1077	0.0782	0.1567	0.0782	0.0990	0.0764
	X ₁₄	ELMAN 0.07	1.2051	1.0978	0.8421	0.6749	0.5320	0.2205	1.3743	0.2205	0.3271	0.2212
Tunshi	X ₁	ELMAN 0.07	18.0147	4.2444	3.1410	2.3735	0.2493	0.1895	0.3685	0.1895	0.2307	0.1865
	X ₁₄	ELMAN 0.07	0.3790	0.6156	0.4559	0.3460	0.4780	0.2668	0.9727	0.2668	0.3577	0.2859
Urbina	X ₁	ELMAN 0.07	0.5971	0.7727	0.5924	0.4928	0.0908	0.0654	0.1379	0.0654	0.0893	0.0658
	X ₁₄	ELMAN 0.07	1.2267	1.1076	0.8418	0.6926	0.5813	0.2447	1.2033	0.2447	0.3596	0.2385
Alao	X ₁	ELMAN 0.07	7.5944	2.7558	1.5846	0.6264	0.1443	0.0586	0.2897	0.0586	0.1368	0.0588
	X ₁₄	ELMAN 0.07	0.9379	0.9684	0.7401	0.5513	2.0055	0.3276	8.6194	0.3276	0.5211	0.3485
Atillo	X ₁	ELMAN 0.07	7.0448	2.6542	1.4252	0.5158	0.3498	0.0643	2.0106	0.0643	0.2168	0.064
	X ₁₄	ELMAN 0.07	0.9374	0.9682	0.7295	0.5723	0.9238	0.2754	2.6142	0.2754	0.4244	0.2749
Espoch	X ₁	ELMAN 0.07	5.0976	2.2578	1.5165	0.9741	0.1252	0.0762	0.2177	0.0762	0.1208	0.0776
Matus	X ₁	ELMAN 0.07	6.2541	2.5008	1.4388	0.6252	0.124	0.0504	0.2347	0.0504	0.1155	0.0511
	X ₁₄	ELMAN 0.07	2.738	1.6547	1.1988	0.8	0.794	0.4668	2.487	0.4668	0.6188	0.5511
Multitud	X ₁	ELMAN 0.07	0.3584	0.5987	0.4121	0.2653	0.0275	0.0182	0.0387	0.0182	0.0271	0.0182
Quimiag	X ₁	ELMAN 0.07	7.2437	2.6914	1.508	0.6248	0.1221	0.0476	0.2531	0.0476	0.1161	0.0476
	X ₁₄	ELMAN 0.07	0.4059	0.6371	0.4715	0.3768	0.6648	0.2139	1.6926	0.2139	0.3473	0.2208

Realizado por: Pilco V. y Acurio W., 2019.

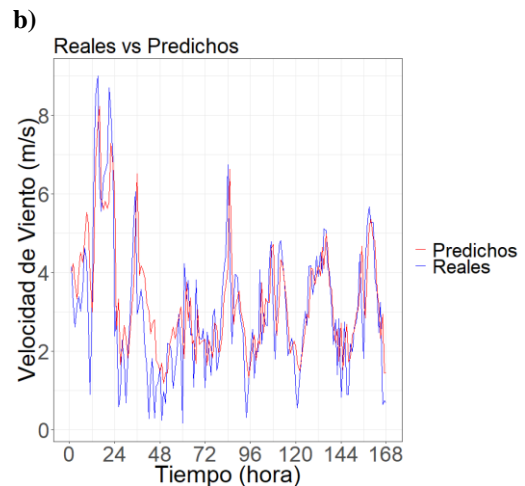
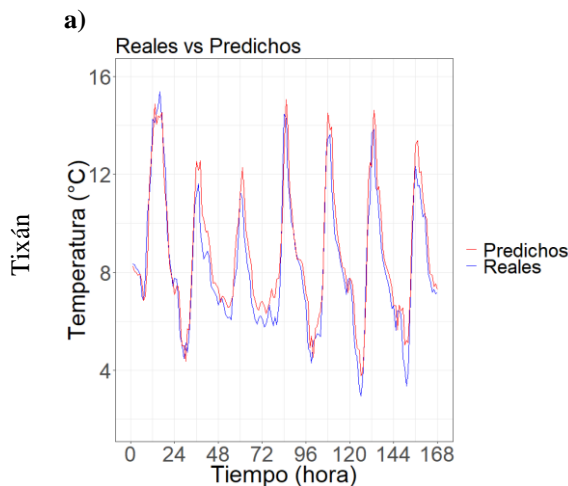
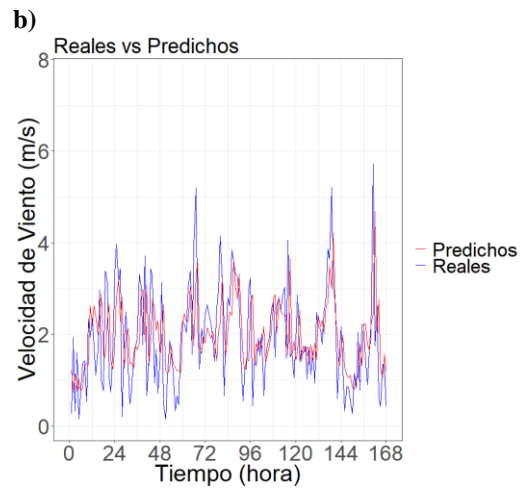
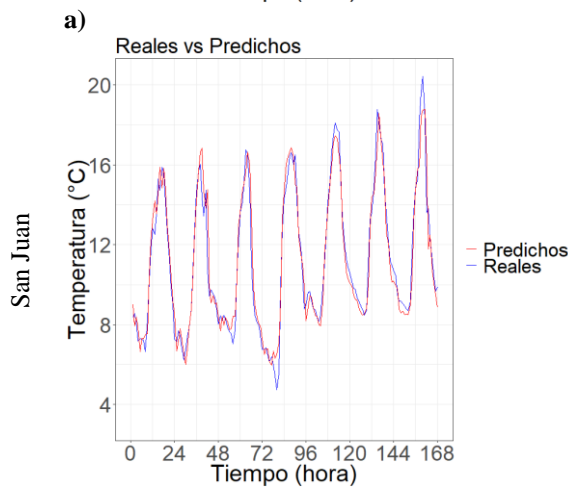
Se realizó pronósticos mediante la red de Elman y de Jordan (Anexo L) para las dos variables en estudio, los cuales fueron analizados mediante los criterios de evaluación. La mayoría de criterios de evaluación o en algunos casos todos indicaban que los pronósticos generados por la red de Elman eran más adecuados que los de la red de Jordan.

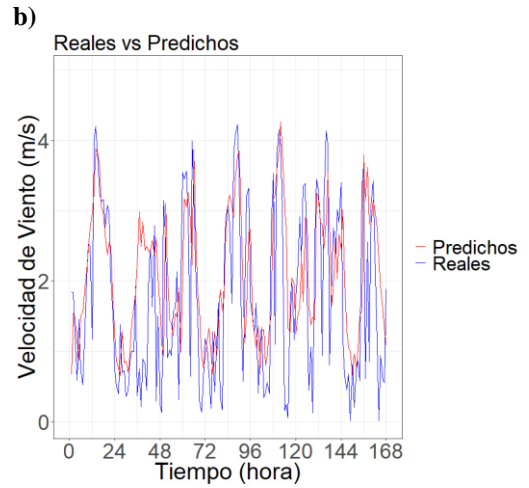
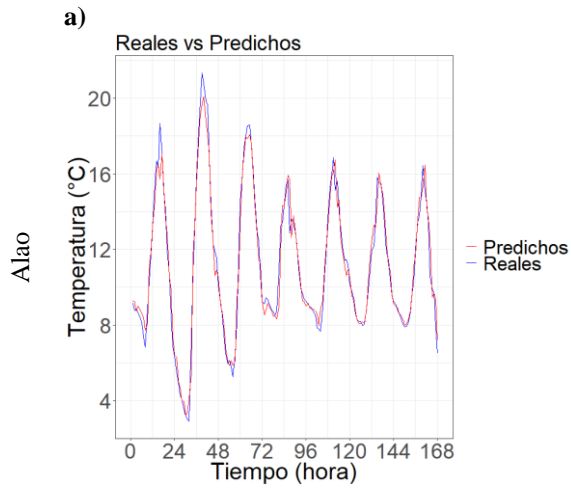
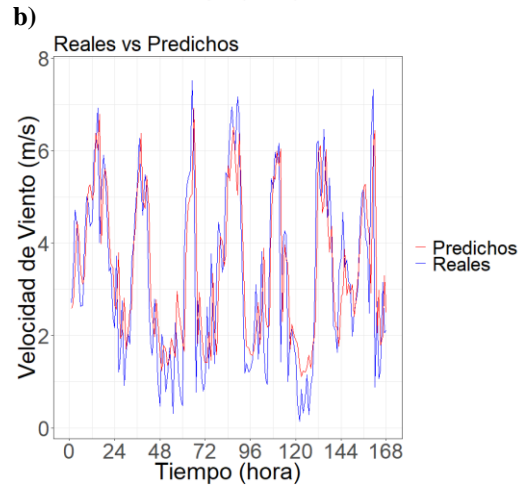
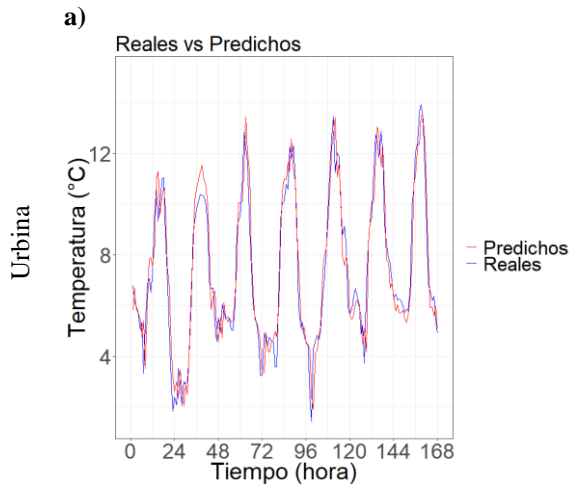
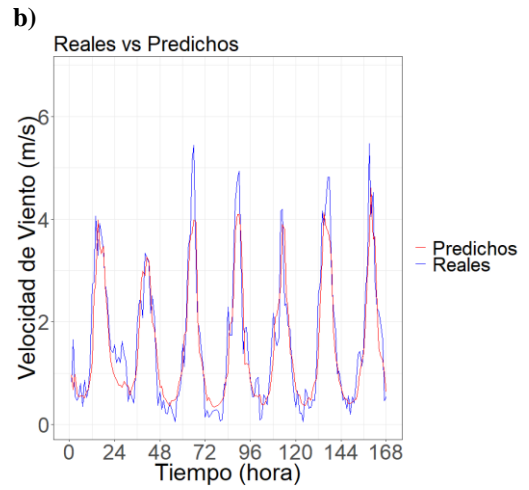
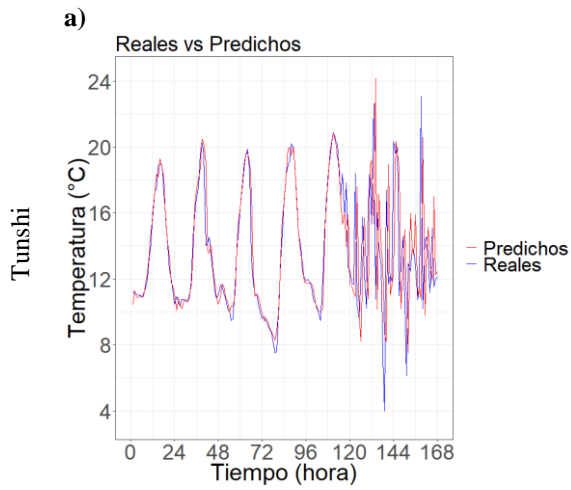
Pronósticos con los modelos seleccionados de Red de Elman

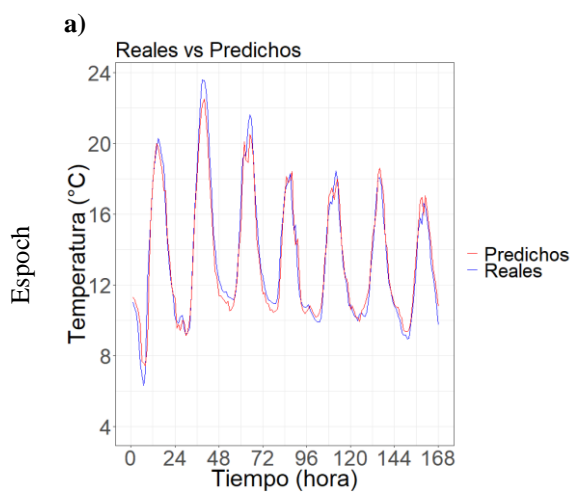
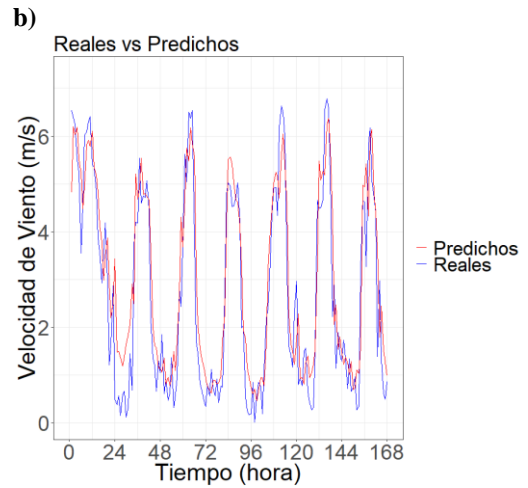
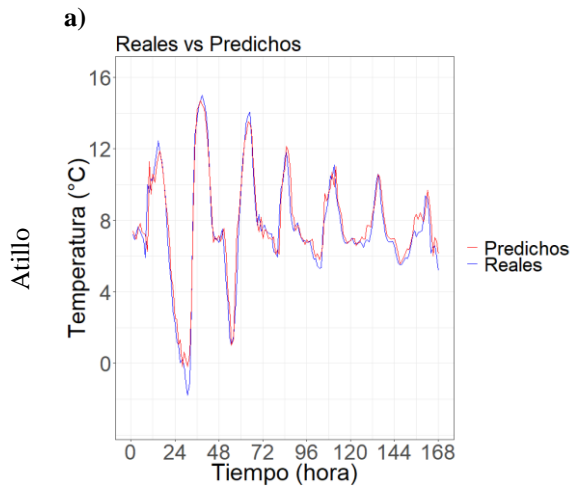


b)

No se realizó el análisis en esta variable, dado que el coeficiente de determinación ajustado es menor al 79% en la imputación y por lo tanto no se considera adecuada.

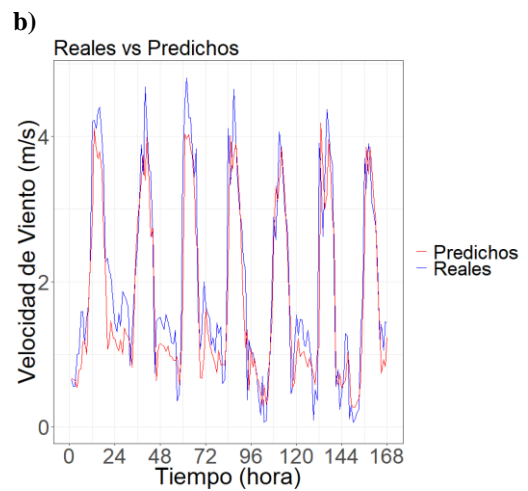
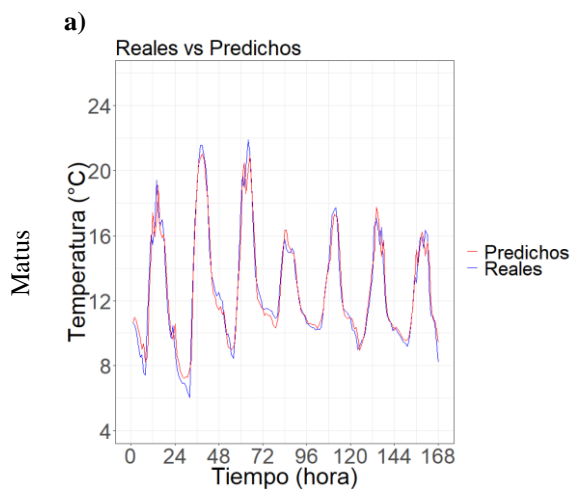


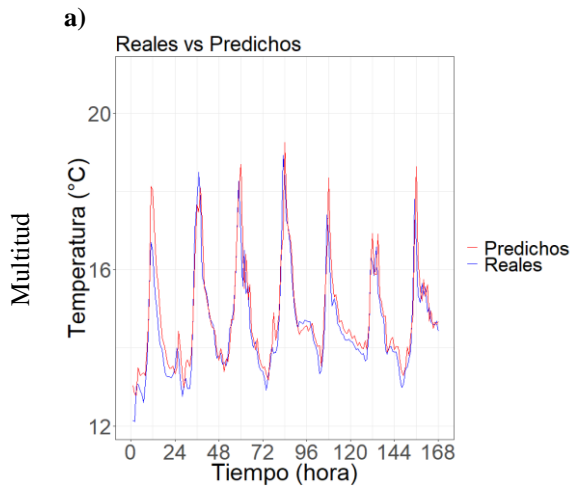




b)

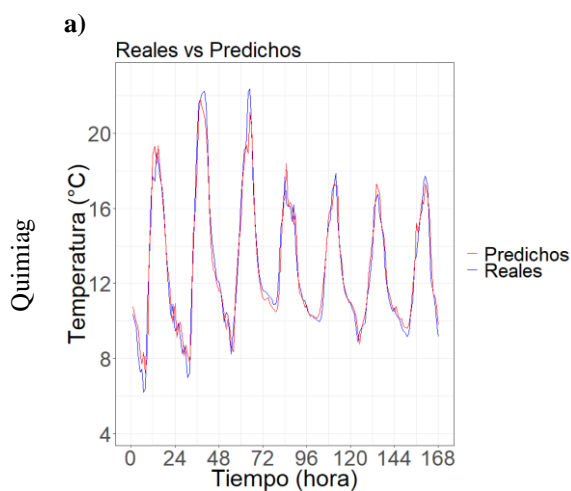
No se realizó el análisis en esta variable, dado que el coeficiente de determinación ajustado es menor al 79% en la imputación y por lo tanto no se considera adecuada.





b)

No se realizó el análisis en esta variable, dado que el coeficiente de determinación ajustado es menor al 79% en la imputación y por lo tanto no se considera adecuada.



b)

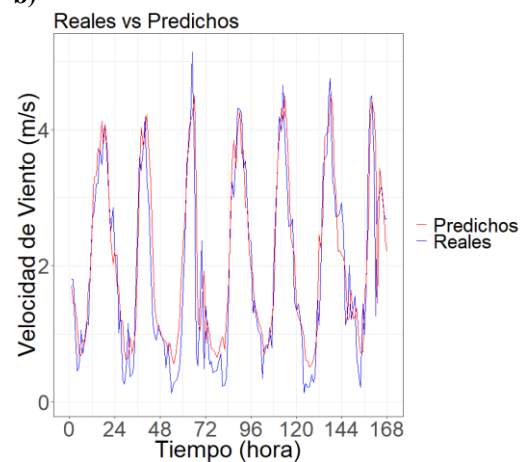


Gráfico 26-4: Datos reales vs predichos de los modelos de Redes Neuronales Recurrentes para cada estación.

Realizado por: Pilco V. y Acurio W., 2019.

Una vez seleccionado el mejor modelo de Redes Neuronales Recurrentes se realiza sus gráficas (Gráfico 26-4), se observa que no existe mucha variabilidad entre los datos originales versus los predichos para las dos variables en estudio mostrando comportamiento aproximadamente similar. Las gráficas fueron realizadas en un lapso de tiempo de una semana (168 horas).

Para temperatura y velocidad de viento se aprecia un buen ajuste a lo largo de la serie, logrando llegar casi hasta los repuntes más altos, generando muy buenos pronósticos para ambas variables.

4.1.7 Análisis de Precisión de los modelos

Se realiza un análisis de comparación de los pronósticos en un período de 31 días para conocer si en este período de tiempo se puede apreciar la precisión de los modelos encontrados, con cada una de las técnicas de modelación; para ello se calculó el coeficiente U de Theil el cual permite

hallar la precisión de pronósticos frente a los datos reales, mientras más cercano a cero se encuentre este valor la precisión será mejor.

Tabla 17-4: Coeficiente U de Theil de los mejores modelos de las tres técnicas.

Estación	Variable	Coeficiente U de Theil		
		ARIMA	T. CAOS	RNR
Cumandá	X ₁	0.0351	0.0526	0.0209
San Juan	X ₁	0.0768	0.0886	0.0354
	X ₁₄	0.2650	0.3521	0.2259
Tixán	X ₁	0.0851	0.1129	0.0539
	X ₁₄	0.3178	0.2673	0.1616
Tunshi	X ₁	0.1693	0.1565	0.1518
	X ₁₄	0.1954	0.2223	0.1537
Urbina	X ₁	0.0866	0.1265	0.0468
	X ₁₄	0.2760	0.3312	0.1606
Alao	X ₁	0.1454	0.1399	0.1142
	X ₁₄	0.2552	0.3510	0.2207
Atillo	X ₁	0.1577	0.1823	0.1568
	X ₁₄	0.1937	0.3825	0.1489
Espoch	X ₁	0.0988	0.099	0.0768
Matus	X ₁	0.0998	0.1131	0.0962
	X ₁₄	0.192	0.2464	0.1479
Multitud	X ₁	0.0558	0.0444	0.0201
Quimiag	X ₁	0.1007	0.1172	0.0986
	X ₁₄	0.1767	0.1979	0.138

Realizado por: Pilco V. y Acurio W., 2019.

El coeficiente U de Theil (Tabla 18-4) en las tres técnicas se aproxima a cero, sin embargo, los coeficientes que destacan son los modelos proporcionados por redes neuronales recurrentes (Red Elman), concluyéndose así que mejor precisión presenta la técnica de redes neuronales recurrentes.

En el software estadístico R-Studio se encuentra el test de Diebold-Mariano con el comando `dm.test` el cual nos presenta tres alternativas de hipótesis dependiendo de la instrucción y de los objetivos que tenga el investigador, estas son: `less`, `greater` y `two.sided`. Mencionaremos sus contrastes de hipótesis:

- **less**
 H_0 : Los dos métodos tienen la misma precisión de pronóstico.
 H_1 : el método 2 es menos preciso que el método 1.
- **greater**
 H_0 : Los dos métodos tienen la misma precisión de pronóstico.
 H_1 : el método 2 es más preciso que el método 1.
- **two.sided**
 H_0 : Los dos métodos tienen la misma precisión de pronóstico.
 H_1 : No presentan la misma precisión de pronóstico.

Tabla 18-4: Test de Diebold-Mariano mediante la instrucción “two.sided” de los mejores modelos de las tres técnicas.

Estación	Variable	Test de Diebold-Mariano		
		ARIMA-T. CAOS	ARIMA-RNR	T. CAOS-RNR
Cumandá	X ₁	1.75E-15	2.20E-16	2.20E-16
San Juan	X ₁	0.04704	2.20E-16	2.20E-16
	X ₁₄	0.2524	2.20E-16	1.39E-14
Tixán	X ₁	6.47E-16	2.20E-16	2.20E-16
	X ₁₄	2.20E-16	2.20E-16	2.20E-16
Tunshi	X ₁	0.0003	1.31E-05	0.0972
	X ₁₄	2.63E-11	3.88E-09	2.20E-16
Urbina	X ₁	2.20E-16	2.20E-16	2.20E-16
	X ₁₄	0.0005	2.20E-16	2.20E-16
Alao	X ₁	0.3541	2.20E-16	2.20E-16
	X ₁₄	2.20E-16	3.73E-12	2.20E-16
Atillo	X ₁	3.73E-13	1.99E-08	2.20E-16
	X ₁₄	2.20E-16	2.20E-16	2.20E-16
EsPOCH	X ₁	0.0492	2.20E-16	2.20E-16
Matus	X ₁	5.58E-11	4.36E-07	2.20E-16
	X ₁₄	2.20E-16	2.20E-16	2.20E-16
Multitud	X ₁	2.30E-06	2.20E-16	2.20E-16
Quimiag	X ₁	4.93E-12	3.00E-07	2.20E-16
	X ₁₄	2.20E-16	2.87E-10	2.20E-16

Realizado por: Pilco V. y Acurio W., 2019.

A un nivel de significancia de 0.05 para las técnicas:

Box-Jenkins (ARIMA) con Teoría del Caos (Tabla 19-4) se rechaza la hipótesis nula para las variables temperatura y velocidad de viento en las estaciones meteorológicas de: Cumandá, San

Juan, Tixán, Tunshi, Urbina, Atillo, Espoch, Matus, Multitud y Quimiag concluyéndose que las dos técnicas no presentan la misma precisión de pronósticos.

En la estación de Alao para temperatura se acepta la hipótesis nula y se concluye que las dos técnicas tienen la misma precisión de pronósticos, en cuanto a velocidad de viento se puede decir que las dos técnicas no presentan la misma precisión de pronósticos.

Box-Jenkins (ARIMA) con Redes Neuronales Recurrentes (Elman) (Tabla 19-4) se rechaza la hipótesis nula para las variables temperatura y velocidad de viento en todas las estaciones meteorológicas concluyéndose que las dos técnicas no presentan la misma precisión de pronósticos.

Teoría del Caos con Redes Neuronales Recurrentes (Elman) (Tabla 19-4) se rechaza la hipótesis nula para las variables temperatura y velocidad de viento en las estaciones meteorológicas de: Cumandá, San Juan, Tixán, Urbina, Alao, Atillo, Espoch, Matus, Multitud y Quimiag concluyéndose que las dos técnicas no presentan la misma precisión de pronósticos.

En la estación de Tunshi para temperatura se acepta la hipótesis nula y se concluye que las dos técnicas tienen la misma precisión de pronósticos, en cuanto a velocidad de viento se puede decir que las dos técnicas no presentan la misma precisión de pronósticos.

Tabla 19-4: Test de Diebold-Mariano mediante la instrucción “greater” de los mejores modelos de las tres técnicas.

Estación	Variable	Test de Diebold-Mariano		
		ARIMA-T. CAOS	ARIMA-RNR	T. CAOS-RNR
Cumandá	X ₁	1	2.20E-16	2.20E-16
San Juan	X ₁	0.9765	2.20E-16	2.20E-16
	X ₁₄	0.1262	2.20E-16	6.94E-15
Tixán	X ₁	1	2.20E-16	2.20E-16
	X ₁₄	2.20E-16	2.20E-16	2.20E-16
Tunshi	X ₁	0.0001	6.54E-06	0.0486
	X ₁₄	1	1.94E-09	2.20E-16
Urbina	X ₁	1	2.20E-16	2.20E-16
	X ₁₄	0.0002	2.20E-16	2.20E-16
Alao	X ₁	0.1771	2.20E-16	2.20E-16
	X ₁₄	1	1.87E-12	2.20E-16
Atillo	X ₁	1	9.94E-09	2.20E-16
	X ₁₄	1	2.20E-16	2.20E-16

Epoch	X ₁	0.0246	2.20E-16	2.20E-16
Matus	X ₁	1	2.18E-07	2.20E-16
	X ₁₄	1	2.20E-16	2.20E-16
Multitud	X ₁	1.15E-06	2.20E-16	2.20E-16
Quimiag	X ₁	1	1.50E-07	2.20E-16
	X ₁₄	1	1.43E-10	2.20E-16

Realizado por: Pilco V. y Acurio W., 2019.

A un nivel de significancia de 0.05 para las técnicas:

Box-Jenkins (ARIMA) con Teoría del Caos (Tabla 20-4) se acepta la hipótesis nula para las variables temperatura y velocidad de viento en las estaciones meteorológicas de: Cumandá, San Juan, Alao, Atillo, Epoch, Matus y Quimiag concluyéndose que las dos técnicas presentan la misma precisión de pronósticos.

En la estación de Tixán para temperatura se acepta la hipótesis nula y se concluye que las dos técnicas tienen la misma precisión de pronósticos, en cuanto a velocidad de viento se puede decir que Teoría del Caos es más preciso que ARIMA.

En la estación de Tunshi para temperatura se rechaza la hipótesis nula y se concluye que Teoría del Caos es más preciso que ARIMA, en cuanto a velocidad de viento se puede decir que las dos técnicas tienen la misma precisión de pronósticos.

En la estación de Urbina para temperatura se acepta la hipótesis nula y se concluye que las dos técnicas tienen la misma precisión de pronósticos, en cuanto a velocidad de viento se puede decir que Teoría del Caos es más preciso que ARIMA.

En la estación de Epoch y Multitud para temperatura se rechaza la hipótesis nula y se concluye que Teoría del Caos es más preciso que ARIMA.

En cuanto a la comparación de ARIMA y Teoría del Caos con respecto a la Redes Neuronales Recurrentes se rechaza la hipótesis nula y se concluye que en todas las estaciones meteorológicas presenta mayor precisión tiene la red de Elman generando mejores pronósticos.

Predicción

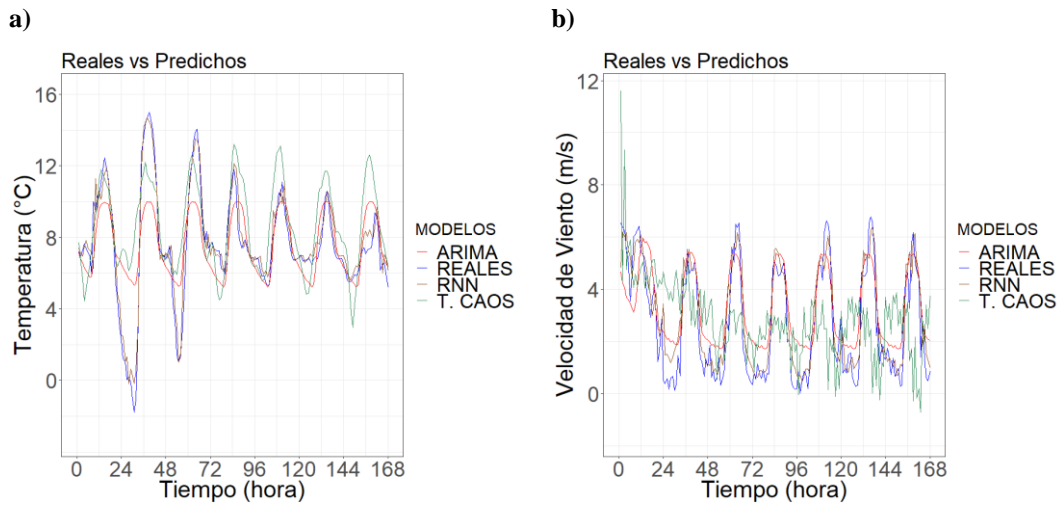
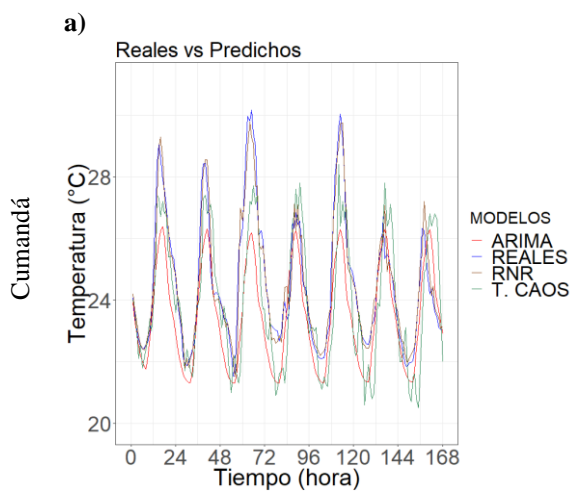


Gráfico 27-4: Datos reales vs predichos de las tres técnicas para la estación de Atillo.

Realizado por: Pilco V. y Acurio W., 2019.

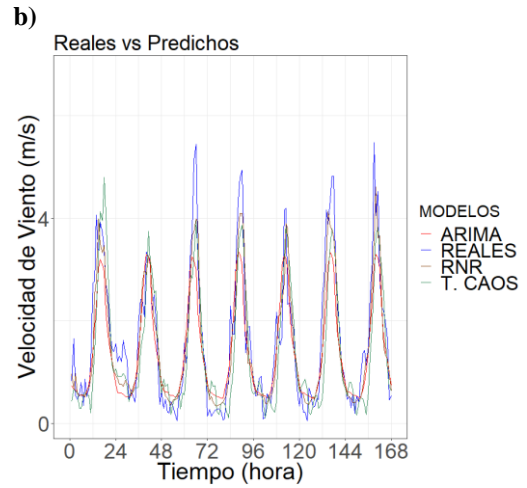
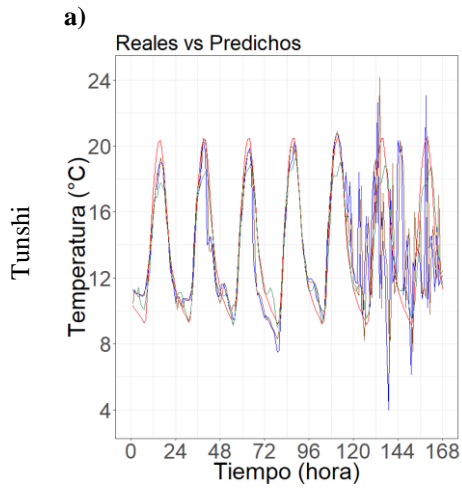
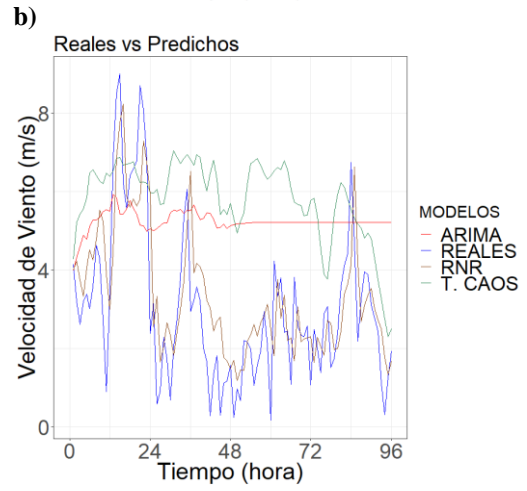
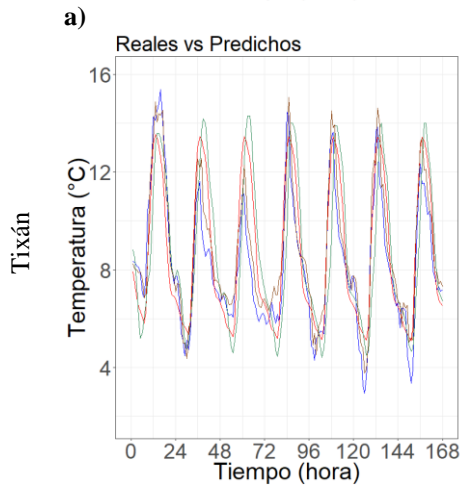
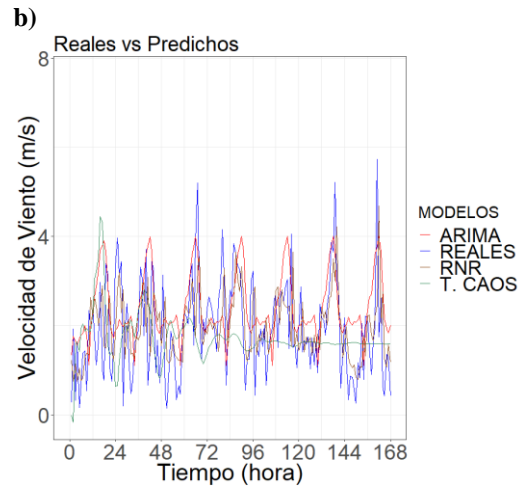
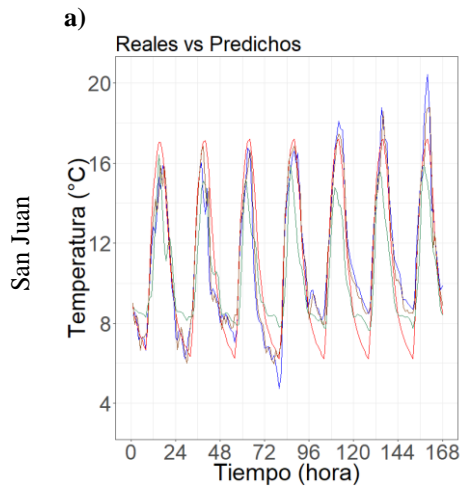
Podemos observar (Gráfica 27-4) de forma gráfica que las tres técnicas muestran semejanza a la serie original, pero la red de Elman presenta un mejor ajuste y menor variabilidad.

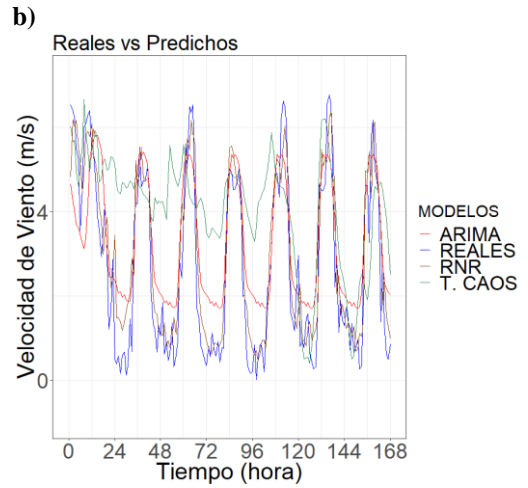
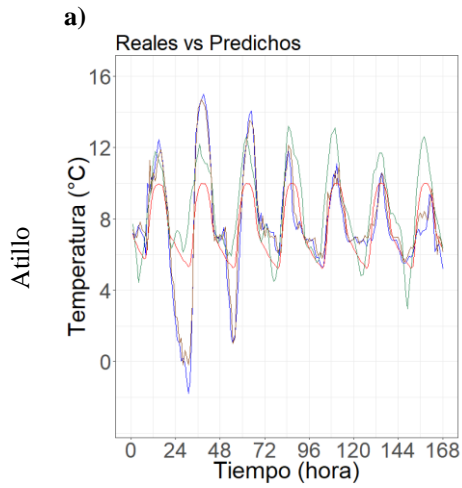
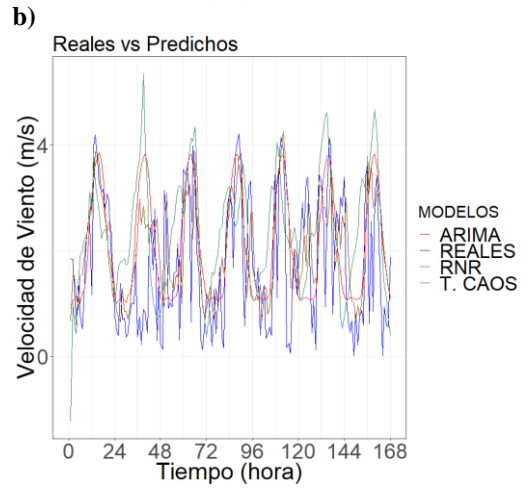
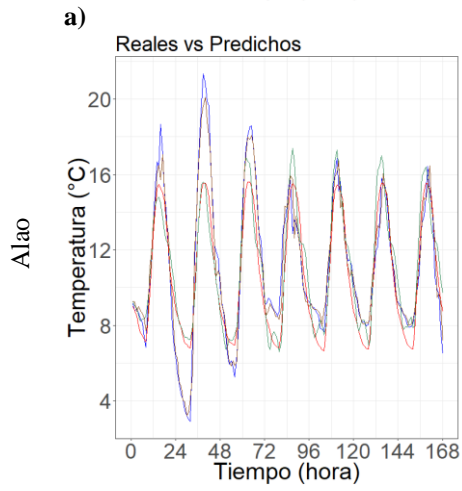
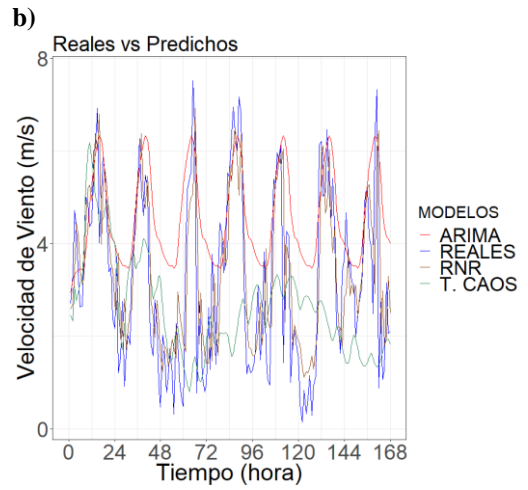
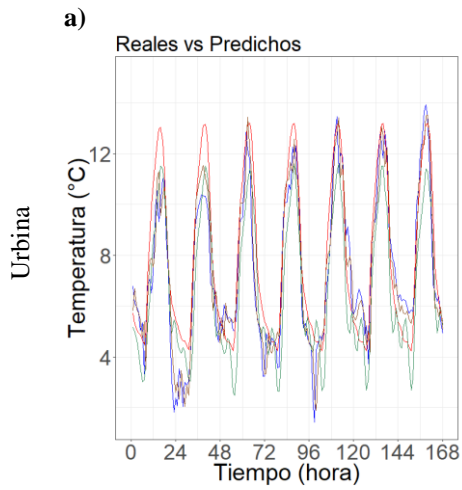
Comparación de los datos predichos obtenidos con el modelo ARIMA, Teoría del Caos y Redes Neuronales Recurrentes durante 30 días.

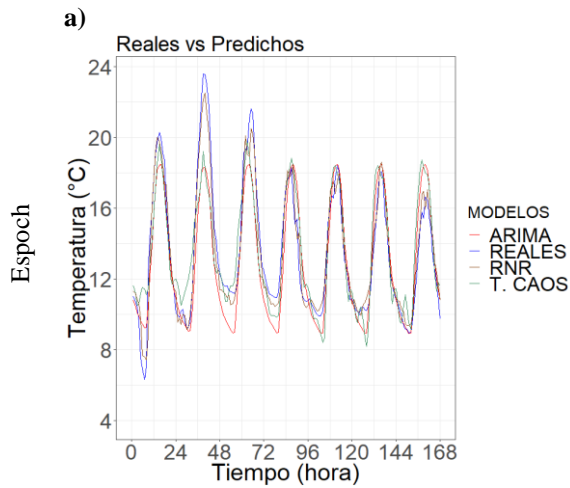


b)

No se realizó el análisis en esta variable, dado que el coeficiente de determinación ajustado es menor al 79% en la imputación y por lo tanto no se considera adecuada.

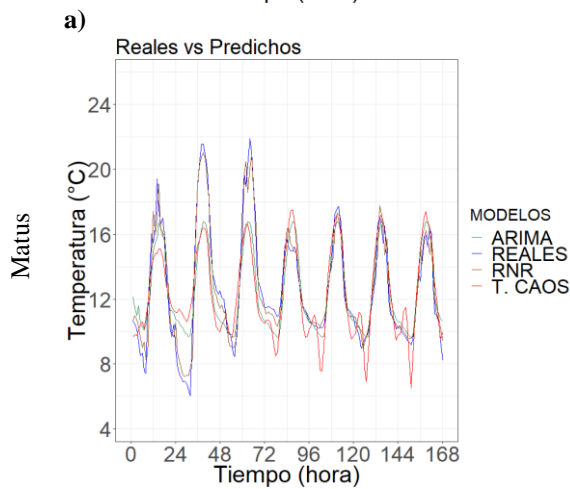




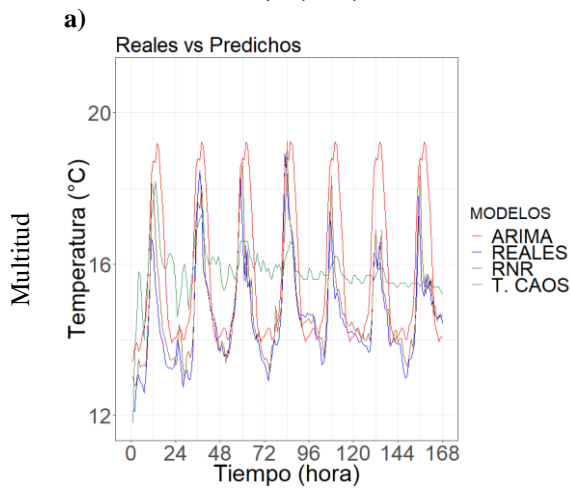
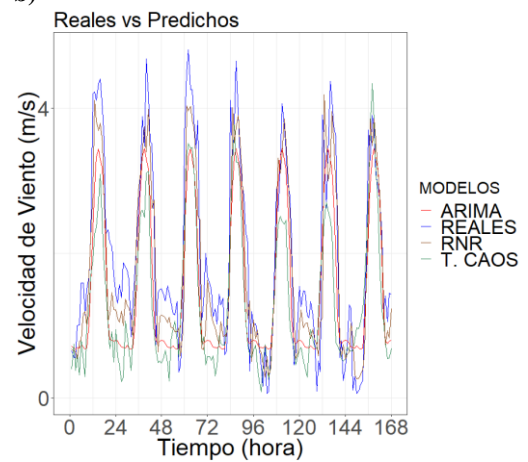


b)

No se realizó el análisis en esta variable, dado que el coeficiente de determinación ajustado es menor al 79% en la imputación y por lo tanto no se considera adecuada.



b)



b)

No se realizó el análisis en esta variable, dado que el coeficiente de determinación ajustado es menor al 79% en la imputación y por lo tanto no se considera adecuada.

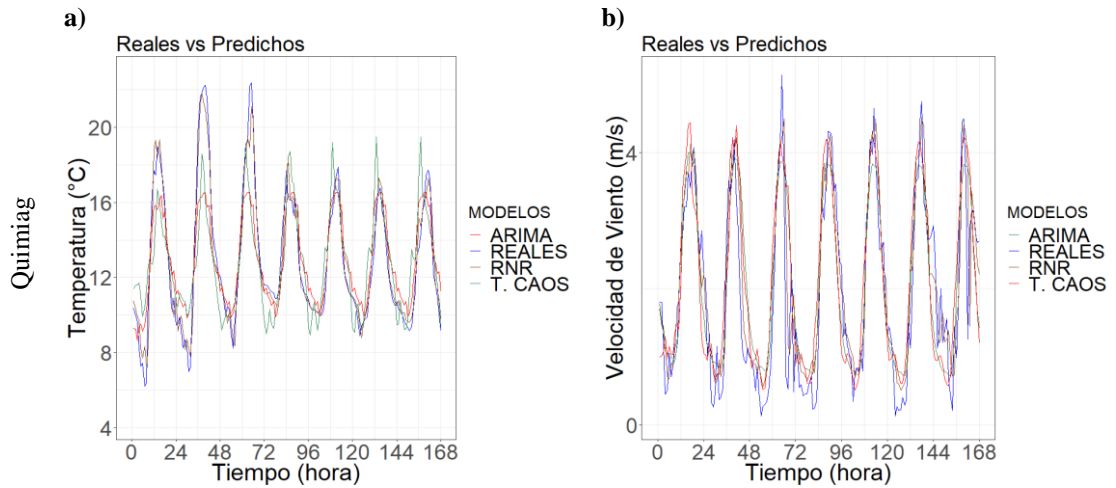


Gráfico 28-4: Datos reales vs predichos de las tres técnicas para cada estación meteorológica.
Realizado por: Pilco V. y Acurio W., 2019.

Los mejores pronósticos (Gráfico 28-4) son los generados por la red Elman (RNR), presentando mejor ajuste a la serie de datos reales y predicciones más precisas y exactas.

CONCLUSIONES

Se estructuró una base de datos por cada estación, sin considerar aquellas series con más del 20% de faltantes y que presentaron un coeficiente de determinación ajustado inferior al 79% en su imputación.

En cada técnica aplicada en esta investigación se identificó varias propuestas de modelos, de los cuales se hizo una selección de los mejores mediante los criterios de evaluación (MSE, MdAE, MdAPE, RMSPE, RMdSPE, SMAPE, SMdAPE y MAE). Al aplicar la metodología Box-Jenkins (ARIMA) se analizaron autocorrelogramas simples y parciales de cada una de las estaciones meteorológicas con sus respectivas variables (velocidad de viento y temperatura), notando que todas las series presentaron estacionalidad cada 24 rezagos; también se aplicó los criterios de información (BIC y AIC), test de Diebold Mariano para medir exactitud y el principio de parsimonia en la elección de los modelos; sin embargo los errores no cumplieron con independencia ni normalidad en general; en cuanto al análisis gráfico no mostraron un buen ajuste, debido a que las predicciones no logran alcanzar los puntos extremos de las observaciones reales. Con la Teoría del Caos los atractores en todas las estaciones no ostentaron simetría ni tendencia, así también en la estación de Atillo no existió autocorrelación, por lo que se realizó reducción de ruido hasta 10 iteraciones con el fin de mantener las características reales de los fenómenos; en los resultados gráficos se observó que las predicciones no llegan a los repuntes más altos y con el pasar el tiempo la variabilidad aumenta, a pesar de ello en un lapso de tiempo de 48 horas las previsiones obtenidas se ajustan a los datos reales incrementando su confiabilidad. Y con Redes Neuronales Recurrentes de Elman y de Jordan previamente se empleó un cambio de escala debido a que trabajan con una función logística, donde Elman obtuvo errores más bajos con aproximadamente 93% de confiabilidad y mediante la función de error cuadrático medio ponderado se visualizó que converge rápidamente a "0" presentando así un ajuste muy similar a los datos reales.

Se evaluó las predicciones obtenidas mediante coeficiente U de Theil y el test unilateral de Diebold-Mariano y se halló que las previsiones obtenidas con redes neuronales recurrentes de Elman presentaron error más cercano a 0, con un 95% de confiabilidad.

Se obtuvieron predicciones de velocidad de viento y temperatura con las 3 técnicas para un periodo hasta 31 días para las 11 estaciones meteorológicas instaladas en la provincia de Chimborazo logrando un mejor ajuste a las observaciones reales con la red neuronal recurrente de Elman.

RECOMENDACIONES

Para modelar con la técnica de Box-Jenkins (ARIMA) se debe tomar en cuenta el tamaño de la base de datos, como el procesador de la computadora a trabajar dado que influye en el tiempo de análisis de los modelos tardándose de 5 a 8 horas en un procesador Core i7 y 7 a 10 horas en un procesador Core i5.

Evitar el uso del RMSE dado que no es un buen indicador de promedio de rendimiento de un modelo, ya que al elevar dichos valores al cuadrado tiene un efecto muy fuerte en su medida y su cálculo suele ser engañoso pues el mismo es sensible a valores atípicos.

En la variable temperatura el cálculo de la métrica MAPE subestima los valores registrados en cero no es recomendable debido a que el valor de cero representa un estado, por lo cual no es adecuado el uso de la misma.

La metodología Box-Jenkins y Teoría del Caos se puede usar para predicciones a corto plazo solo para temperatura, dado que en velocidad de viento no se ajustan a los datos reales, mientras que Redes Neuronales Recurrentes proporciona mejores previsiones a largo y corto plazo para dichas variables.

Se debería considerar dar un mantenimiento adecuado a los sensores que se maneja en las estaciones meteorológicas, y corroborar que los mismos están funcionando correctamente para evitar pérdida de información.

Se debería incluir en la malla curricular de la carrera de Estadística, métodos y técnicas de imputación de datos, dado que al momento de trabajar con datos reales se presentan situaciones en las que se existen observaciones faltantes y se requiere de información confiable que al ser manipulada no se aleje de la realidad. Además, dentro de las Electivas una asignatura sobre el estudio de Redes Neuronales Artificiales debido a su alto potencial de modelación y predicción.

GLOSARIO

CEAA: Centro de Energías Alternativas y Ambiente.

INAMHI: Instituto Nacional de Meteorología e Hidrología.

PMA: Programa Mundial de Alimentos.

MAE: Ministerio del Ambiente.

MAGAP: Ministerio de Agricultura, Ganadería, Acuacultura y Pesca.

MSP: Ministerio Salud Pública.

SGR: Secretaria de Gestión de Riesgos.

RNR: Redes Neuronales Recurrentes.

OMM: Organización Meteorológica Mundial

MSE: Mean Square Error

RMSE: Root Mean Square Error

MAE: Mean Absolute Error

MdAE: Median Absolute Error

MAPE: Mean Absolute Percentage Error

MdAPE: Median Absolute Percentage Error

RMSPE: Root Mean Square Percentage Error

RMdSPE: Root Median Square Percentage Error

sMAPE: Symmetric Mean Absolute Percentage Error

sMdAPE: Symmetric Median Absolute Percentage Error

AIC: Criterio de Información Akaike

BIC: Criterio de información Bayesiano

BIBLIOGRAFÍA

Acuña E. *Regresión Lineal* [En Línea]. Puerto Rico: Universidad de Puerto Rico, 2016. [Consulta: 15 septiembre 2018]. Disponible en: <https://docplayer.es/22233458-Capitulo-9-regresion-lineal.html>

Agencia Estatal de Meteorología. *Meteorología y climatología de Navarra* [En Línea]. Navarra: Gobierno de Navarra, 2018. [Consulta: 14 septiembre 2018]. Disponible en: <http://meteo.navarra.es/estaciones/estacion.cfm?IDestacion=405>

Alves Monteiro, Carlos Cristiano. *Controlador predictivo para el proceso de lodos activados* [En línea] (Tesis). (Diplomado). Universidad Central “Marta Abreu” de las Villas, Santa Clara. 2009-2010. pp. 1-61 [Consulta: 25 enero 2018]. Disponible en: <http://dspace.uclv.edu.cu/bitstream/handle/123456789/4518/Carlos%20Cristiano%20Alves%20Monteiro.pdf?sequence=1&isAllowed=y>

Aragón Moreno, J. A. *La meteorología: Orígenes, evolución y desarrollo* [blog]. Consulta: 15 octubre 2018]. Disponible en: <http://cambioclimaticoyciudades.weebly.com/blog/la-meteorologia-origenes-evolucion-y-desarrollo>

Jaramillo Ayerbe, M. et al. *Análisis de series de tiempo univariante aplicando metodología de Box-Jenkins para la predicción de ozono en la ciudad de Cali, Colombia.* Revista Facultad de Ingeniería Universidad de Antioquia, vol. 4, N°39, (2007), (Perú) pp. 1-11

Bejar, J. *Caracterización de datos electrocardiográficos mediante la teoría del caos.* 2001, pp. 1-82.

Beyaert A. *Ejemplos de Predicción* [en línea]. Murcia, 2018. [Consulta: 10 febrero 2019]. Disponible en: <https://docplayer.es/20848629-Ejemplos-de-prediccion.html>

Castro Cacabelos, Moisés. *Imputación de datos faltantes en un modelo de tiempo de fallo acelerado* [En Línea] (Tesis). (Maestría) Universidad de Vigo, España. 2014. pp.17-18. [Consulta: 27 enero 2019]. Disponible en: http://eio.usc.es/pub/mte/descargas/ProyectosFinMaster/Proyecto_940.pdf

Díaz Kusztrich, Miguel. *Redes neuronales recurrentes y series temporales* [blog]. 2016. [Consulta: 18 diciembre 2018]. Disponible en: <http://software-tecnico-libre.es/es/articulo-por-tema/todas-las-secciones/todos-los-temas/todos-los-articulos/redes-neuronales-recurrentes-y-series-de-tiempo>

Ruelas Santoyo, Edgar Augusto. & Laguna González, José Antonio. *Comparación de predicción basada en redes neuronales contra métodos estadísticos en el pronóstico de ventas.* Revista Ingeniería Industrial Actualidad y Nuevas Tendencias, Vol. IV, No. 12, ISSN: 1856-8327 (2013), (México) pp. 91-105.

El Comercio Ecuador. *Abundante producción de granos en la sierra.* [Consulta: 12 enero 2018]. Disponible en: <https://www.elcomercio.com/actualidad/granos-sierra-ceniza-lluvias-tungurahua.html>

El Comercio Ecuador. *El cambio climático amenaza a la agricultura.* [Consulta: 12 enero 2018]. Disponible en: <http://especiales.elcomercio.com/planeta-ideas/planeta/11-octubre-del-2015/el-cambio-climatico-amenaza-a-la-agricultura>

Eransus Armendáris, Francisco J. *Evaluación y comparación de capacidad predictiva bajo funciones de pérdidas discretas* [En línea] (Doctorado) Universidad Complutense de Madrid, Madrid. 2010. p. 84. [Consulta: 03 noviembre 2018]. Disponible en: <https://eprints.ucm.es/11217/1/T32091.pdf>

Escudero Villa, Amalia Isabel. *Modelación y pronóstico del potencial energético del río Blanco usando la teoría del caos y un método convencional* [En línea] (Tesis). (Posgrado). Escuela Superior Politécnica de Chimborazo, Ecuador. 2007. pp. 1-114. [Consulta: 24 enero 2018]. Disponible en: <http://dspace.esPOCH.edu.ec/bitstream/123456789/1302/1/226T0005.pdf>

Fernández Casal, Rubén. *Diagnosis de la Independencia* [blog]. [Consulta: 08 febrero 2019]. Disponible en: <https://rubenfcasal.github.io/post/diagnosis-de-la-independencia/>

Giraldo, N. *Series de Tiempo en R.* Colombia: ISBN, 2006, pp. 1-157

Gómez Suárez, Mónica. *Espacio ocupado en el lineal por las marcas del distribuidor: estimación mediante redes neuronales vs regresión múltiple.* ScienceDirect [En línea], 2009, (Madrid) Vol. 12, pp. 37-68. [Consulta: 11 noviembre 2018]. Issue 41. Disponible en: <https://www.sciencedirect.com/science/article/pii/S1138575809700477>

González Avella, J. C.; Tudurí, J. M. & Rul-lan, G. *Análisis de series temporales usando redes neuronales recurrentes* [blog]. [Consulta: 11 noviembre 2018]. Disponible en: <https://www.apsl.net/blog/2017/06/14/analisis-de-series-temporales-usando-redes-neuronales-recurrentes/>

González Rodríguez, Andrés Felipe. *Modelo para la predicción de la radiación solar a partir de redes neuronales artificiales* [En línea] (Tesis). (Posgrado). Escuela de Ingeniería de Antioquia.

2013. pp. 1-80 [Consulta: 18 enero 2018]. Disponible en: https://repository.eia.edu.co/bitstream/11190/326/7/GonzalezAndres_2013_ModeloPrediccionRadiacion.pdf

Hanke, J. E. & Wichern, D. W. *Pronósticos en los Negocios*. Novena edición. México: Pearson Educación, 2010, pp.1-399.

Haro, Arquimides X.; Limaico, Cecilia T.; & Llosas, Yolanda E. *Predicción de datos meteorológicos en cortos intervalos de tiempo en la ciudad de Riobamba usando la teoría del Caos*. *Sistemas, Cibernética e Informática*, Vol. 13, No. 1, ISSN: 1690-8627 (Ecuador) pp. 35-41.

Hegger, Rainer.; Kantz, Holger & Schreiber, Thomas. *Nonlinear Time Series Analysis* [En línea]. 1998-2007. [Consulta: 14 febrero 2018]. Disponible en: https://www.pks.mpg.de/~tisean/Tisean_3.0.1/index.html

Bonet Cruz, Isis. et al. *Redes neuronales recurrentes para el análisis de secuencias*. *Revista Cubana de Ciencias Informáticas*, Vol. 1, No. 4 (2007), (Cuba) pp. 1 – 11.

Jiménez, José F.; Gázquez, Juan C. & Sánchez, R. *La capacidad predictiva en los métodos Box-Jenkins y Holt-Winters: una aplicación al sector turístico*. *Revista Europea de Dirección y Economía de la Empresa*, vol. 15, N° 3, (2006), (España) pp. 187-188.

López, Danilo A.; García, Nancy. & Herrera, Jhon F. *Developing a Predictive Model to Estimate the Behavior of Variables in a Network Infrastructure*. *Revista Información Tecnológica*, vol. 26, N°5, (2015), (Chile) pp. 143-154

López, E. A. *Estadística con aplicaciones en Agronomía y Ciencias Forestales*. 2008, p. 1-226

Matich, Damián Jorge. *Redes Neuronales: Conceptos Básicos y Aplicaciones*. 2001, pp. 1-55.

Morales Edwin, Palomino. *Trazando series de tiempo con ggplot2 y ggfortify* [En línea]. 2017. [Consulta: 14 enero 2019]. Disponible en: https://rstudio-pubs-static.s3.amazonaws.com/333961_44287c618bfe4f91ba01b71391cb06da.html

Ortiz, Raúl. *GGPLOT2. Modificar nombre del gráfico y de los ejes* [En línea]. 2015. [Consulta: 15 enero 2019]. Disponible en: <https://rpubs.com/Rortizdu/140201>

Ortiz, Raúl. *Manipulación de los ejes* [En línea]. 2015. [Consulta: 15 enero 2019]. Disponible en: http://rstudio-pubs-static.s3.amazonaws.com/140203_b85429a724b341e8886315c191106c16.html

Ovando, Gustavo.; Bocco, Mónica. & Sayago, Silvina. *Redes neuronales para modelar predicción de heladas.* Agricultura Técnica, Vol. 65, No. 1 (2005), (Argentina) pp. 65-73.

Peña D. *Análisis de Datos Multivariantes.* 2002, pp. 1-495

Pino Diez, Raúl. et al. *Pronóstico de la velocidad y dirección del viento mediante Redes Neuronales Artificiales.* [Consulta: 10 enero 2018]. Disponible en: https://www.researchgate.net/publication/228859782_Pronostico_de_la_Velocidad_y_Direccion_del_Viento_mediante_Red_Neuronales_Artificiales

Pitarque, Alfonso.; Ruiz, Juan Carlos. & Roy, Juan Francisco. *Las redes neuronales como herramientas estadísticas no paramétricas de clasificación.* Psicothema, Vol. 12, No. 2 (2000), (Spain) pp. 459-463.

Poma Lima, Diana Lucía. *Predicción Meteorológica mediante Redes Neuronales.* [Consulta: 11 enero 2018]. Disponible en: <https://dlpoma.wordpress.com/2010/06/01/prediccion-meteorologica-mediante-redes-neuronales/>

Poveda Jaramillo, Germán. *Atractores extraños (caos) en la hidro-climatología de Colombia.* Rev. Acad. Colomb, Vol. XXI, No. 81 (1997), (Colombia) pp. 431-444.

Purca S., Quispe C. *Forecast of sea surface temperature of the Peruvian coast using an autoregressive integrated moving average model.* Revista Peruana de Biología, vol. 4, N°1, (2007), (Perú) pp. 109-114

Reich, Carlos S. *La teoría del caos: definición y ejemplo* [blog]. 2009. [Consulta: 24 enero 2018]. Disponible en: <https://shreich.wordpress.com/2009/10/05/la-teoria-del-caos-definicion-y-ejemplo/>

Reich, Carlos. *Teoría de Sistemas* [blog]. [Consulta: 04 diciembre 2018]. Disponible en: <https://shreich.wordpress.com/tag/teoria-de-sistemas/>

Rodríguez, Aguado J. et al. *Meteorological Variables Prediction Through Arima Models.* Revista Agrociencia, vol. 50, N°1, (2016), (México) pp. 1-13

Rodríguez, Jiménez R.; Benito, Capa Águeda. & Portela, Lozano Adelaida. *Meteorología y Climatología,* (2004), (España) p. 56.

Roldán Quintero, Raúl Eduardo. *Soluciones por modelos estadísticos y redes neuronales artificiales* (Tesis). (PhD). Tenaca American University, Caracas. 2002. pp. 1-147.