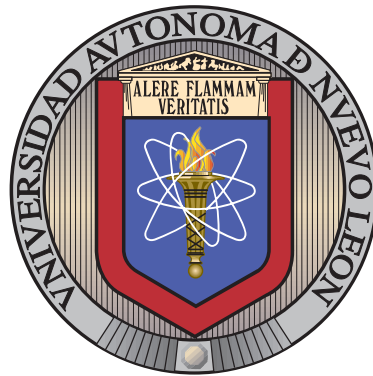


UNIVERSIDAD AUTÓNOMA DE NUEVO LEÓN

FACULTAD DE INGENIERÍA MECÁNICA Y ELÉCTRICA

SUBDIRECCIÓN DE ESTUDIOS DE POSGRADO



PRONÓSTICO DE FALLAS EN LA DISTRIBUCIÓN
DE ENERGÍA ELÉCTRICA

POR

MARIO ALBERTO GUTIÉRREZ CARRALES

EN OPCIÓN AL GRADO DE

MAESTRÍA EN CIENCIAS DE LA INGENIERÍA

CON ORIENTACIÓN EN SISTEMAS

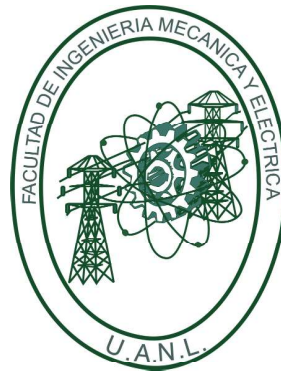
SAN NICOLÁS DE LOS GARZA, NUEVO LEÓN

JULIO 2020

UNIVERSIDAD AUTÓNOMA DE NUEVO LEÓN

FACULTAD DE INGENIERÍA MECÁNICA Y ELÉCTRICA

SUBDIRECCIÓN DE ESTUDIOS DE POSGRADO



PRONÓSTICO DE FALLAS EN LA DISTRIBUCIÓN
DE ENERGÍA ELÉCTRICA

POR

MARIO ALBERTO GUTIÉRREZ CARRALES

EN OPCIÓN AL GRADO DE

MAESTRÍA EN CIENCIAS DE LA INGENIERÍA

CON ORIENTACIÓN EN SISTEMAS

SAN NICOLÁS DE LOS GARZA, NUEVO LEÓN

JULIO 2020



UNIVERSIDAD AUTÓNOMA DE NUEVO LEÓN
FACULTAD DE INGENIERÍA MECÁNICA Y ELÉCTRICA
SUBDIRECCIÓN DE ESTUDIOS DE POSGRADO

Los miembros del Comité de Tesis recomendamos que la Tesis «Pronóstico de fallas en la distribución de energía eléctrica», realizada por el alumno Mario Alberto Gutiérrez Carrales, con número de matrícula 1549273, sea aceptada para su defensa como requisito para obtener el grado de Maestría en Ciencias de la Ingeniería con Orientación en Sistemas.

El Comité de Tesis

Dr. José Arturo Berrones Santos
Director

Dr. César Emilio Villarreal Rodríguez
Revisor

Dr. Jonás Velasco Álvarez
Revisor

Vo. Bo.

Dr. Simón Martínez Martínez

Subdirector de Estudios de Posgrado



San Nicolás de los Garza, Nuevo León, julio 2020

Dedico este trabajo con mucho amor a Dios por permitir que este proyecto se hiciera realidad, a mis padres Alfredo y Ángeles quienes me dieron la vida, educación y apoyo, a mis hermanos Guadalupe, Alicia y Alfredo por su cariño y por compartir cosas importantes conmigo y a mi novia Brenda Sotelo por alentarme y acompañarme en todo momento, incluso en los más difíciles.

ÍNDICE GENERAL

Agradecimientos	XV
Resumen	XVI
1. Introducción	1
1.1. Descripción del problema	3
1.2. Objetivo	3
1.2.1. Objetivos específicos	3
1.3. Motivación y justificación	4
1.4. Hipótesis	4
1.5. Estructura de la tesis	4
2. Revisión de literatura	6
2.1. Sector eléctrico	6
2.1.1. Procesos industriales	6
2.1.2. Industria energética	7
2.1.3. Sector eléctrico	8

2.2. Pronóstico	9
2.2.1. Métricas de evaluación	9
2.2.2. Estrategias de pronóstico	10
2.2.3. Modelos	11
3. Marco teórico	13
3.1. Series de tiempo	13
3.1.1. Descomposición	15
3.1.2. Autocorrelación	16
3.1.3. Estacionariedad	17
3.2. Pronóstico	18
3.2.1. Elementos de pronóstico	18
3.2.2. Métrica de evaluación	19
3.2.3. Modelos de regresión	22
3.2.4. Pronóstico con horizonte mayor a uno	28
4. Metodología	33
4.1. Descripción y preprocesamiento de los datos	33
4.2. Metodología de pronóstico	34
5. Experimentación y resultados	40
5.1. Características de la experimentación	40

5.2. Experimentación	42
5.2.1. Montemorelos	43
5.2.2. Zona Metropolitana Norte	51
5.2.3. Zona Metropolitana Oriente	61
5.2.4. Zona Metropolitana Poniente	70
6. Conclusiones	79
6.1. Conclusiones	79
6.2. Contribuciones	81
6.3. Trabajo a futuro	81

ÍNDICE DE FIGURAS

1.1. Divisiones de distribución de la CFE Fuente: https://www.cantilever.com.mx/	1
1.2. Proceso de generación, transmisión y distribución de electricidad . . .	2
3.1. Descomposición de la serie de tiempo del número de pasajeros en una aerolínea internacional	15
3.2. Autocorrelación de la serie de tiempo $(x_t)_{t \in N}$	17
3.3. Representación visual de la métrica MAPE	20
3.4. Visualización del ejemplo de un árbol de decisión	25
3.5. Visualización de las divisiones del espacio	25
3.6. Perceptrón multicapa	28
3.7. Series de entrenamiento y prueba de una serie de tiempo	29
3.8. Clasificación de modelos de regresión de acuerdo al número de variables	30
3.9. Ejemplo con $N = 10$ y $n = 4$	31
4.1. Preprocesamiento de datos para la creación de series de tiempo	34
4.2. Metodología de pronóstico	35

5.1. Fallas en Montemorelos	43
5.2. Fallas en Montemorelos por año	44
5.3. Descomposición aditiva de las fallas en Montemorelos	45
5.4. Autocorrelación de fallas en Montemorelos	45
5.5. Estacionariedad de fallas en Montemorelos	46
5.6. MAPE por longitud de prueba para cada familia en Montemorelos . .	46
5.7. Conjuntos de entrenamiento y prueba de fallas en Montemorelos . . .	47
5.8. MAPE mínimo por combinación de ventana de pronóstico para cada familia en Montemorelos (N=30)	48
5.9. MAPE por hiperparámetro(s) de cada familia en Montemorelos (N=30, V=(1,3))	49
5.10. Comparación de MAPE de los modelos representativos de cada familia en Montemorelos (N=30, V=(1,3))	50
5.11. Pronóstico de fallas en Montemorelos (N=30, V=(1,3), M=KNNR(10))	50
5.12. Distribución del error porcentual absoluto por intervalos en Montemorelos (N=30, V=(1,3), M=KNNR(10))	51
5.13. Fallas en la Zona Metropolitana Norte	52
5.14. Fallas en la Zona Metropolitana Norte por año	53
5.15. Descomposición aditiva de las fallas en la Zona Metropolitana Norte .	54
5.16. Autocorrelación de fallas en la Zona Metropolitana Norte	54
5.17. Estacionariedad de fallas en la Zona Metropolitana Norte	55

5.18. MAPE por longitud de prueba para cada familia en la Zona Metropolitana Norte	55
5.19. Conjuntos de entrenamiento y prueba de fallas en la Zona Metropolitana Norte	56
5.20. MAPE mínimo por combinación de ventana de pronóstico para cada familia en la Zona Metropolitana Norte (N=50)	57
5.21. MAPE por hiperparámetro(s) de cada familia en la Zona Metropolitana Norte (N=50, V=(1,1))	58
5.22. Comparación de MAPE de los modelos representativos de cada familia en la Zona Metropolitana Norte (N=50, V=(1,1))	59
5.23. Pronóstico de fallas en la Zona Metropolitana Norte (N=50, V=(1,1), M=KNNR(1))	60
5.24. Distribución del error porcentual absoluto por intervalos en la Zona Metropolitana Norte (N=50, V=(1,1), M=KNNR(1))	60
5.25. Fallas en la Zona Metropolitana Oriente	61
5.26. Fallas en la Zona Metropolitana Oriente por año	62
5.27. Descomposición aditiva de las fallas en la Zona Metropolitana Oriente	63
5.28. Autocorrelación de fallas en la Zona Metropolitana Oriente	63
5.29. Estacionariedad de fallas en la Zona Metropolitana Oriente	64
5.30. MAPE por longitud de prueba para cada familia en la Zona Metropolitana Oriente	64
5.31. Conjuntos de entrenamiento y prueba de fallas en la Zona Metropolitana Oriente	65

5.32. MAPE mínimo por combinación de ventana de pronóstico para cada familia en la Zona Metropolitana Oriente (N=40)	66
5.33. MAPE por hiperparámetro(s) de cada familia en la Zona Metropolitana Oriente (N=40, V=(14,8))	67
5.34. Comparación de MAPE de los modelos representativos de cada familia en la Zona Metropolitana Oriente (N=40, V=(14,8))	68
5.35. Pronóstico de fallas en la Zona Metropolitana Oriente (N=40, V=(14,8), M=KNNR(3))	68
5.36. Distribución del error porcentual absoluto por intervalos en la Zona Metropolitana Oriente (N=40, V=(14,8), M=KNNR(3))	69
5.37. Fallas en la Zona Metropolitana Poniente	70
5.38. Fallas en la Zona Metropolitana Poniente por año	71
5.39. Descomposición aditiva de las fallas en la Zona Metropolitana Poniente	72
5.40. Autocorrelación de fallas en la Zona Metropolitana Poniente	72
5.41. Estacionariedad de fallas en la Zona Metropolitana Poniente	73
5.42. MAPE por longitud de prueba para cada familia en la Zona Metropolitana Poniente	73
5.43. Conjuntos de entrenamiento y prueba de fallas en la Zona Metropolitana Poniente	74
5.44. MAPE mínimo por combinación de ventana de pronóstico para cada familia en la Zona Metropolitana Poniente (N=50)	75
5.45. MAPE por hiperparámetro(s) de cada familia en la Zona Metropolitana Poniente (N=50, V=(17,1))	76

5.46. Comparación de MAPE de los modelos representativos de cada familia en la Zona Metropolitana Poniente (N=50, V=(17,1))	77
5.47. Pronóstico de fallas en la Zona Metropolitana Poniente (N=50, V=(17,1), M=KNNR(9))	77
5.48. Distribución del error porcentual absoluto por intervalos en la Zona Metropolitana Poniente (N=50, V=(17,1), M=KNNR(9))	78
6.1. Serie de tiempo de las fallas en Montemorelos y la velocidad del viento en la estación Suroeste en el 2013	85

ÍNDICE DE TABLAS

2.1. Estudios relacionados con actividades económicas y el pronóstico . . .	7
2.2. Estudios relacionados con la industria energética y el pronóstico . . .	8
2.3. Clasificación de pronósticos	10
2.4. Estudios con las estrategias recursiva, directa y mezcla	11
3.1. Correlaciones de la serie de tiempo contra sus desplazos	17
3.2. Descripción de métricas	21
3.3. Ventanas de pronóstico de aprendizaje y predicción con longitud: N , longitud de prueba: n y ventana de pronóstico: (p, h)	29
3.4. Ventanas de pronóstico para el ejemplo de la figura 3.9 (b)	31
3.5. Regresión multisalida como múltiples regresiones con una salida . . .	32
5.1. Posibles valores de la longitud del conjunto de prueba y ventana de pronóstico	41
5.2. Posibles valores de los hiperparámetros en los modelos	42
5.3. Desempeño del pronóstico en Montemorelos	51
5.4. Desempeño del pronóstico en Montemorelos	61

5.5. Desempeño del pronóstico en la Zona Metropolitana Oriente	69
5.6. Desempeño del pronóstico en la Zona Metropolitana Poniente	78
6.1. Parámetros obtenidos para cada zona	80
6.2. Desempeño del pronóstico para cada zona	80

AGRADECIMIENTOS

Agradezco principalmente a la Facultad de Ingeniería Mecánica Eléctrica (FI-ME) y a la Universidad Autónoma de Nuevo León (UANL) por brindarme la oportunidad de realizar mis estudios de posgrado en sus instalaciones.

Por su parte, también agradezco al Posgrado de Ingeniería de Sistemas (PISIS) tanto a los profesores como al personal administrativo y a los compañeros. Gracias por compartir su paciencia, tiempo y dedicación conmigo.

También agradezco a los doctores Arturo Berrones, Jonás Velasco y César Villarreal por formar parte del comité de revisión de esta tesis y dedicar su tiempo en leerla y brindar aportes para su mejora.

Así mismo, agradezco al Dr. Miguel Mata por la publicación de la plantilla de tesis, así como al personal de mantenimiento de la CFE de la división de distribución Golfo Norte por proporcionar los datos de las fallas en sus zonas.

Por último, agradezco al Consejo Nacional de Ciencia y Tecnología (CONACYT) por el apoyo económico otorgado mediante una beca de estudios de tiempo completo.

RESUMEN

Mario Alberto Gutiérrez Carrales.

Candidato para obtener el grado de Maestría en Ciencias de la Ingeniería con Orientación en Sistemas.

Universidad Autónoma de Nuevo León.

Facultad de Ingeniería Mecánica y Eléctrica.

Título del estudio: PRONÓSTICO DE FALLAS EN LA DISTRIBUCIÓN DE ENERGÍA ELÉCTRICA.

Número de páginas: 91.

OBJETIVOS Y MÉTODO DE ESTUDIO: En este trabajo se consideran cuatro series de tiempo, cada una representa las fallas en la distribución de energía eléctrica en alguna de las zonas que conforman la división de distribución Golfo Norte, las cuales son Montemorelos, Zona Metropolitana Norte, Zona Metropolitana Oriente y Zona Metropolitana Poniente. Las fallas contempladas son aquellas causadas por un árbol ó una planta. El período comprendido por las series de tiempo consta de enero del 2013 a diciembre del 2018, en donde cada observación es registrada de manera semanal.

El presente estudio propone realizar un análisis estadístico para determinar la mejor configuración para pronosticar cada serie de tiempo en el que se consideran

variaciones en la longitud del conjunto de prueba, ventana de pronóstico y modelos de pronóstico. Las familias de modelos de pronóstico consideradas son: regresión lineal (RL), regresión de Ridge (RR), regresión de los k vecinos más cercanos (KNNR), árbol de decisión (DTR), bosque aleatorio (RFR) y perceptrón multicapa (MLPR) los cuales son modelos de regresión dentro del campo del aprendizaje supervisado en el aprendizaje automático. La métrica de evaluación que se utiliza para medir el desempeño de un modelo es el error absoluto porcentual medio (MAPE).

CONTRIBUCIONES Y CONCLUSIONES: Se utilizan diferentes librerías de python, entre ellas statsmodels y sklearn para estudiar características descriptivas de cada una de las series de tiempo, además, para cada una de ellas se determinan los valores adecuados para llevar a cabo el pronóstico. De esta manera, se determina el modelo que mejor desempeño tiene para pronosticar cada serie de tiempo y finalmente, se realiza un análisis para evaluar la calidad del pronóstico obtenido en el que se utilizan los errores absolutos porcentuales y diferentes métricas como error cuadrático medio (RMSE), error absoluto medio (MAE) y el coeficiente de determinación ajustado (R^2)

Firma del asesor:

Dr. José Arturo Berrones Santos

CAPÍTULO 1

INTRODUCCIÓN

La Comisión Federal de Electricidad (CFE) es una empresa industrial mexicana que brinda servicios de generación, transmisión y distribución de energía eléctrica. Sus servicios llegan a clientes, los cuales son residenciales e industriales sumando entre ellos cerca de 41 millones de beneficiarios [1].

Para eficientizar el suministro de la electricidad en el territorio mexicano existen divisiones de distribución que se encargan de brindar servicio a distintas zonas [2], estas divisiones se muestran en la figura 1.1.



Figura 1.1: Divisiones de distribución de la CFE

Fuente:<https://www.cantilever.com.mx/>

Esta investigación se centra en la división de distribución Golfo Norte, que a su vez cubre 4 zonas: Montemorelos, Zona Metropolitana Norte, Zona Metropolitana Oriente y Zona Metropolitana Poniente.

De acuerdo con información proporcionada por la división de distribución Golfo Norte el proceso que se sigue para generar, transmitir y distribuir la electricidad empieza en las centrales eléctricas, donde se produce una corriente con una tensión de 3-36 kV. Posteriormente, al salir de las centrales eléctricas, se eleva la tensión de la corriente hasta 110-380 kV. Después, en estaciones transformadoras, se varía de nuevo el voltaje de la corriente 25-132 kV. Finalmente, a través de redes de distribución la electricidad llega al cliente final. En la figura 1.2 se muestra este proceso a detalle.

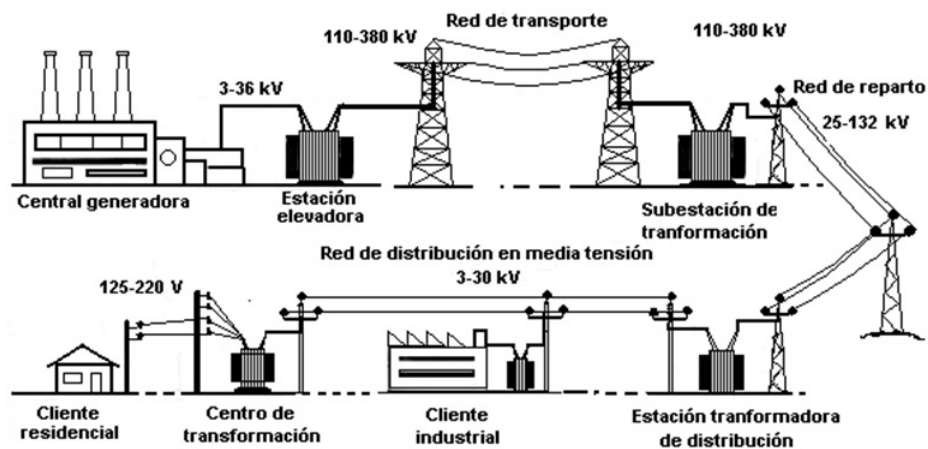


Figura 1.2: Proceso de generación, transmisión y distribución de electricidad

En [3] se define que una falla es cualquier interrupción en el proceso, ya sea en la parte de la generación, transmisión o distribución, y puede ser originada por distintos fenómenos. Para identificar la causa a la que una falla es atribuida, se tienen criterios establecidos y el personal que atiende dicha falla está capacitado para identificarla de manera adecuada.

Entre las diferentes fallas es importante identificar aquellas que son ocasionadas por árboles y plantas en la distribución de esta energía, ya que desde un punto de

vista estratégico se pueden desarrollar proyectos a corto, mediano y largo plazo. Debido a esto, esta investigación se centra en dichas fallas en la parte de la distribución de la energía eléctrica.

1.1 DESCRIPCIÓN DEL PROBLEMA

El personal de mantenimiento de la división de distribución Golfo Norte cuenta con los registros de las fallas ocasionadas por diferentes causas, específicamente las que son causadas por un árbol ó planta, pero no se han estudiado a detalle dichos registros, por lo que no se cuenta con un modelo matemático que permita pronosticar el número de fallas.

1.2 OBJETIVO

Emplear técnicas estadísticas y modelos de aprendizaje automático para pronosticar el número de fallas causadas por un árbol o planta en las redes de distribución de energía eléctrica de las zonas comprendidas por la división de distribución Golfo Norte de la CFE.

1.2.1 OBJETIVOS ESPECÍFICOS

Para lograr el objetivo de este estudio se necesita cubrir con cada uno de los siguientes puntos:

- Seleccionar modelos matemáticos para llevar a cabo el pronóstico.
- Elegir una métrica que mida el desempeño del pronóstico obtenido por un modelo específico.

- Realizar experimentación para la elección adecuada de los elementos de pronóstico y de los hiperparámetros de los modelos.
- Evaluar la calidad del pronóstico efectuado.

1.3 MOTIVACIÓN Y JUSTIFICACIÓN

La principal motivación y justificación de este trabajo es el contar con un modelo matemático para cada zona de de la división de distribución Golfo Norte que pronostique con un buen desempeño las fallas causadas por un árbol o planta.

En segundo lugar, se tiene que el personal de mantenimiento de la CFE realiza podas de manera frecuente pero sin el conocimiento de como llevar a cabo esta operación, por lo que conociendo el número de fallas en períodos futuros, se puede saber a que zona brindar mayor atención, y además, en que medida realizar las podas, ya que estas tienen un costo económico.

1.4 HIPÓTESIS

La implementación de modelos de aprendizaje automático permitirá pronosticar las fallas causadas por un árbol o planta en las zonas comprendidas por la división de distribución Golfo Norte con un buen desempeño en base a una métrica de evaluación.

1.5 ESTRUCTURA DE LA TESIS

Los siguientes capítulos en este trabajo de tesis se llevan a cabo de la siguiente manera:

En el capítulo 2 se realiza la revisión de literatura acerca de la energía eléctrica, así como modelos de pronóstico. En el capítulo 3 se revisan los conceptos esenciales de una serie de tiempo, además se describen los modelos que se aplican a los datos para pronosticar una serie de tiempo.

En el capítulo 4 se explica la metodología que se sigue, desde el preprocesamiento de los datos hasta aplicar los modelos de pronóstico de una manera organizada y, de este modo, identificar el modelo que es más apropiado. En el capítulo 5 se llevan a cabo los pasos descritos en la metodología para cada una de las series de tiempo de las fallas en cada una de las zonas.

CAPÍTULO 2

REVISIÓN DE LITERATURA

Este capítulo incluye una revisión literaria sobre temas relacionados con este trabajo. Por un lado se encuentra el sector eléctrico, el cual juega un papel importante en la economía de cualquier país y por otra parte el pronóstico, partiendo de las estrategias de como llevar a cabo las predicciones hasta algunos de los modelos existentes.

2.1 SECTOR ELÉCTRICO

La energía eléctrica toma un rol muy importante dentro de la economía debido a la aportación que tiene en los distintos tipos de industria. Su importancia es de tal grado que se busca mejorar la calidad y funcionamiento del sistema eléctrico y una mala gestión de las actividades relacionadas con esta energía se ve reflejado en un estancamiento en el desarrollo de la población en cuestión.

2.1.1 PROCESOS INDUSTRIALES

En [4] se define al concepto de economía como el estudio de como son utilizados los recursos para obtener distintos bienes y distribuirlos entre la sociedad para

consumirlos.

Por otra parte, en [5, 6] se hace mención sobre la composición de la estructura económica de un país, en donde se menciona que la economía se divide en tres sectores, primario, secundario y terciario.

Actividades como la agricultura y minería pertenecen al sector primario, mientras que algunas actividades del sector secundario son la automotriz y energética. Finalmente, el turismo y la banca corresponden a actividades dentro del sector terciario.

Tabla 2.1: Estudios relacionados con actividades económicas y el pronóstico

Autor	Contexto	Modelo(s)
[7]	Exportación de café	ARIMA
[8]	Almacenamiento de productos perecederos	ARIMA, HW
[9]	Suministro de oxígeno de vuelo	ARIMA
[10]	Varios	Red neuronal, regresión KNN, regresión soporte vectorial, procesos gaussianos
[11]	Ventas	Bosque aleatorio, red neuronal

2.1.2 INDUSTRIA ENERGÉTICA

De acuerdo con [12] el sector energético es determinante para fortalecer la economía, finanzas públicas, desarrollo tecnológico, balanza comercial, así como las relaciones con el exterior.

[13] menciona que la industria energética está conformada por 3 subsectores:

- Extracción de petróleo y gas.
- Generación, transmisión y distribución de energía eléctrica.

- Fabricación de productos de derivados del petróleo y carbón.

Tabla 2.2: Estudios relacionados con la industria energética y el pronóstico

Autor	Campo	Modelo(s)	Horizonte de pronóstico
[14]	Electricidad	Red neuronal	Anual
[15]	Electricidad	ARIMA	Diario
[16]	Gas	Hibrido	Anual
[17]	Gas	Hibrido	Diario
[18]	Viento	ACO, PSO	Diario
[19]	Viento	Hibrido	Diario

2.1.3 SECTOR ELÉCTRICO

[20] afirma que el sector eléctrico eficiente es de alta importancia para todos los países, ya que el mal manejo de este tipo de energía representa un obstáculo a largo plazo para el desarrollo tecnológico, tanto para el crecimiento del producto interno bruto (PIB) por habitante y, además, para el mejoramiento del bienestar de la población. Debido a esto, resulta fácil notar que una gran cantidad de países, entre ellos México, permanezcan a la vanguardia en profundas transformaciones estructurales de la industria eléctrica para lograr su modernización, fortalecer la competitividad y brindar mejores servicios a los usuarios.

Entre los trabajos que relacionan el sector eléctrico y el pronóstico se basan principalmente en el pronóstico de la demanda del consumo y el pronóstico del precio de la energía eléctrica.

2.2 PRONÓSTICO

En esta sección se realiza una búsqueda literaria de los aspectos principales para llevar a cabo el pronóstico de una serie de tiempo de una variable, los cuales son: la métrica de evaluación del pronóstico, la estrategia de pronóstico que es la manera secuencial de realizar dicho pronóstico y el modelo matemático que se utiliza, el cual proporciona una o varias ecuaciones que son utilizadas para predecir. Los elementos necesarios para llevar a cabo el pronóstico se describen a detalle en el capítulo 3.

2.2.1 MÉTRICAS DE EVALUACIÓN

Los modelos matemáticos brindan una solución a un determinado problema, por lo general, la solución que brinda un modelo es distinta a la que brindan los demás, aunque pueden contar con similitudes de acuerdo con la idea principal con la que estos son desarrollados. Debido a esto, es necesario contar con métricas de evaluación para modelos de aprendizaje automático supervisado, por ejemplo las métricas de **error cuadrático medio**, **error absoluto medio**, **error absoluto porcentual medio**, entre otras. las cuales se muestran en [21, 22], las cuales determinan que tan bueno o malo es el desempeño de un modelo al brindar solución a dicho problema. Usualmente se compara el desempeño de dos o más modelos para saber cual es el indicado resolviendo una o varias instancias en particular del problema. Los problemas que se presentan en los distintos sectores requieren de un alto nivel de precisión, ya que de él dependen cosas importantes para la humanidad, por ejemplo: determinar si un paciente tiene una enfermedad o no [23], determinar las ventas de un producto [24, 25], niveles de contaminantes presentes en una ciudad [26], entre otros. Y es evidente la necesidad de una buena precisión ya que una mala decisión puede representar pérdidas significativas de acuerdo con el contexto del problema.

2.2.2 ESTRATEGIAS DE PRONÓSTICO

Una vez que se determinan los elementos de pronóstico de manera adecuada, [27, 28] afirman que el pronóstico se clasifica en **corto**, **mediano** y **largo plazo** de acuerdo al horizonte de pronóstico, el cual representa el número de observaciones a pronosticar. En la tabla 2.3 se muestran las características de cada tipo de pronóstico.

Tabla 2.3: Clasificación de pronósticos

Tipo de pronóstico	Cantidad de variables para el aumento de precisión	Horizonte de pronóstico (observaciones)	Aplicaciones
Corto plazo	2 ó menos	Menos de 12	Razones tácticas que incluyen la planificación y el control de la producción
Mediano plazo	De 2 a 4	de 12 a 24	Toma de decisiones estratégicas menores en relación con la operación del negocio
Largo plazo	Más de 4	Más de 24	Toma de decisiones estratégicas importantes dentro de una organización, y se relacionan mucho con las implicaciones de los recursos

De acuerdo con [29, 30, 31] existen tres estrategias principales para realizar el pronóstico. La primera es la **estrategia recursiva** la cual consiste en pronosticar y utilizar dichas predicciones para pronosticar las siguientes hasta completar el horizonte, repitiendo este paso hasta pronosticar el período de prueba. La segunda es la **estrategia directa** y consiste en crear múltiples modelos de pronóstico para así pronosticar cada una de las observaciones en el horizonte, de esta manera se pronostica utilizando solamente valores conocidos. La tercera estrategia es la **directa y recursiva** o también llamada *DirRect*, esta estrategia combina aspectos de las dos anteriores.

En la tabla 2.4 se muestran distintos estudios en los que se compara el rendimiento de los modelos empleando las diferentes estrategias.

Tabla 2.4: Estudios con las estrategias recursiva, directa y mezcla

Autor	Modelo	Casos de estudio	Estrategia de pronóstico (MAE)
Antti Sorjamaa Amaury Lendasse [30]	KNN	2	Rec 3379 0.0318 Dir 1057 0.0124 DirRect 850 0.0098
Adil Ahmed Muhammad Khalid [32]	MLP	1	Rec 1.00 Dir 1.05
Amaury Lendasse [33]	SVM	1	Rec 0,00713 Dir 0,00260

2.2.3 MODELOS

Los modelos sirven para obtener los pronósticos, por lo que así como la métrica de evaluación y la estrategia de pronóstico también juegan un papel fundamental.

Están dados por una ecuación en la que para pronosticar una observación se utiliza un número dado de observaciones al pasado y cada modelo tiene su propia naturaleza de funcionamiento, ya sea por la ecuación de pronóstico, el método de obtención de los parámetros (en caso de que la ecuación tenga) o por las observaciones sobre las cuales se aplica la ecuación ya mencionada.

Los modelos clásicos de pronóstico son aquellos en los que se tiene involucrado el número el número de observaciones al pasado como hiperparámetro de la familia correspondiente. Dos de estas familias son ARIMA y Holt-Winters.

Por otra parte, los modelos de aprendizaje supervisado ó específicamente los modelos de regresión, son aquellos en los que se predice una variable numérica con ciertas variables numéricas llamadas predictoras. Si estas variables representan las observaciones pasadas de la serie de tiempo, se puede obtener una predicción para el pronóstico, permitiendo así a los modelos de regresión pronosticar una serie de tiempo. Dentro de los modelos más utilizados de regresión se encuentran: lineales,

no lineales, ensemble.

A pesar de que la expresión de un modelo lineal predice la respuesta en términos de una combinación lineal de las variables predictoras, existen diversas maneras de obtener los valores de los coeficientes, los cuales son llamados parámetros del modelo. el método más conocido es el mínimos cuadrados ordinarios, otro es el de la regresión de ridge, la cual obtiene los parámetros de manera similar a los mínimos cuadrados ordinarios, con la diferencia que agrega una penalización en los coeficientes.

Dentro de los modelos no lineales se encuentran principalmente: k vecinos más cercanos, máquinas de soporte vectorial, árboles de decisión.

Los modelos ensemble se basan en que se tienen múltiples submodelos predictores y en base a esas predicciones realizar una de manera general. Modelos que están dentro de esta categoría son: bosque aleatorio, refuerzo adaptativo (*Ada boost*, por el acrónimo de adaptive boosting, en inglés), entre otros.

Finalmente se encuentran los modelos de aprendizaje profundo, los cuales están basados en la arquitectura de una red neuronal. Para decir que un modelo pertenece a los modelos de aprendizaje profundo tiene que poseer una gran cantidad de parámetros a optimizar.

CAPÍTULO 3

MARCO TEÓRIO

En este capítulo se describen los conceptos teóricos fundamentales dentro del campo de las series de tiempo de una variable, así como las fórmulas matemáticas que se utilizan para llevar a cabo el pronóstico de manera adecuada.

3.1 SERIES DE TIEMPO

El concepto más básico del que se parte es el de una **serie de tiempo** que según [34, 35] es una sucesión finita de observaciones cuantitativas $(x_t)_{t \in T}$, donde $T = \{1, 2, \dots, N\}$, además cada una de dichas observaciones es registrada en un tiempo específico t y dicha serie de tiempo está conformada por N observaciones.

Las series de tiempo pueden modelar una gran cantidad de procesos económicos, físicos ó demográficos, entre otros. Algunos ejemplos de estos procesos son:

- El precio de un artículo.
- Temperatura máxima en una ciudad.
- Tasa de crecimiento de una determinada población.
- Número de accidentes automovilísticos en una carretera.

- Número de infectados en una epidemia.

Cada serie de tiempo tiene su propio comportamiento, es decir, si se mide una misma variable en diferentes ubicaciones o localidades se desconoce si serán iguales o no, esto se debe a la presencia de factores externos como el contexto del proceso ó las condiciones del ambiente, entre otros.

Una serie de tiempo puede ser **continua** ó **discreta** de acuerdo con el tipo de variable en que las observaciones son medidas, por ejemplo, el número de personas que visitan un supermercado representa una variable discreta, mientras que el peso de una persona representa una continua.

Otra característica de las series de tiempo tiene que ver con el periodo en que cada observación es registrada. Si el periodo es constante, se dice que la serie de tiempo es **equidistante**, mientras que en el caso en que algún periodo no sea constante, se le conoce como **serie de tiempo no equidistante**. Ejemplos de estas series de tiempo son el promedio diario de CO_2 emitido en una ciudad y la escala de Richter de un terremoto, en el caso equidistante y no equidistante, respectivamente. Esta es una característica importante ya que hay procesos que ocurren de manera periódica y otros de manera espontánea debido a su naturaleza.

En las secciones 3.1.1, 3.1.2 y 3.1.3 se estudia la descomposición, autocorrelación y estacionariedad de una serie de tiempo, respectivamente. Estos conceptos proporcionan la suficiente información sobre el comportamiento de una serie de tiempo, dicha información se utiliza para entender de qué manera se describen y relacionan temporalmente las observaciones y para auxiliar a los modelos de pronóstico para obtener un mejor rendimiento.

3.1.1 DESCOMPOSICIÓN

En [36] se mencionan los conceptos de tendencia, estacionalidad y aleatoriedad, en el que definen la **tendencia** como la representación del crecimiento o declinación en los valores a largo plazo, la **estacionalidad** como los patrones de cambio en distintas frecuencias a través del tiempo y la **aleatoriedad** como un comportamiento irregular causado por sucesos impredecibles. Además, en dicho trabajo, se afirma que se puede descomponer una serie de tiempo utilizando dichos términos, esta descomposición puede ser aditiva, multiplicativa o mixta, según como sea vista la serie de tiempo observada, sumando, multiplicando o mezclando suma y multiplicación de las tres componentes mencionadas, respectivamente. En la figura 3.1 se muestra un ejemplo de la descomposición aditiva de la serie de tiempo del número de pasajeros en una aerolínea, utilizada en diversos trabajos [37, 38]

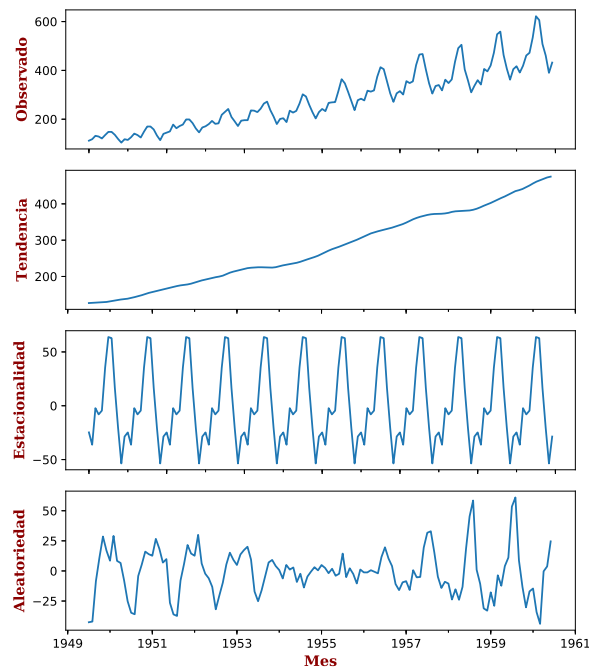


Figura 3.1: Descomposición de la serie de tiempo del número de pasajeros en una aerolínea internacional

3.1.2 AUTOCORRELACIÓN

La función de autocorrelación de una serie de tiempo mide la asociación lineal que existe entre las observaciones de una serie separadas por d unidades de tiempo (x_t, x_{t-d}) , está definida por la ecuación 3.1:

$$r(x, y) = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum_{i=1}^n (x_i - \bar{x})^2 \sum_{i=1}^n (y_i - \bar{y})^2}} \quad (3.1)$$

donde el valor de $r(x, y)$ representa el coeficiente de correlación de pearson entre las variables x e y . Es importante observar que cada valor de r está en el intervalo $[-1, 1]$ para cualquier desplazo y se cumplen las relaciones:

- si $|r| = 1$, la correlación es perfecta
- si $0.5 \leq |r| < 1$, la correlación es fuerte
- si $0 < |r| < 0.5$, la correlación es débil
- si $|r| = 0$, la correlación es nula

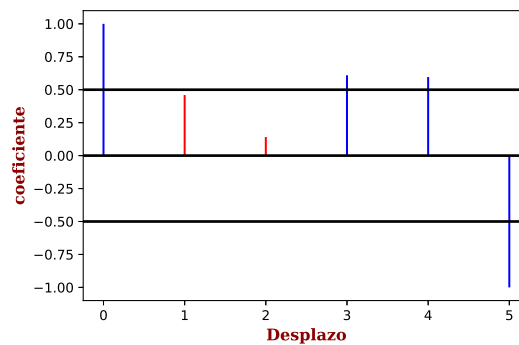
Por ejemplo, si se tiene la serie de tiempo $(x_t)_{t \in T} = (10, 14, 12, 6, 5, 8, 3)$, para calcular la función de autocorrelación se determina que los posibles desplazamientos al pasado de la serie de tiempo, son 0, 1, 2, 3, 4 y 5 y en general hasta la longitud de la serie de tiempo menos dos, ya que si d excede este número, no habrá datos para comparar o bien, solo hay un dato por variable. En la mayoría de los cuadros cuando la serie de tiempo contiene muchas observaciones, se suele visualizar sólo hasta cierto valor de desplazamiento.

En la tabla 3.1 se muestran los desplazamientos para la serie de tiempo en el ejemplo dado, así como sus coeficientes calculados.

Tabla 3.1: Correlaciones de la serie de tiempo contra sus desplazos

t	x_t	x_{t-1}	x_{t-2}	x_{t-3}	x_{t-4}	x_{t-5}
1	10	-	-	-	-	-
2	14	10	-	-	-	-
3	12	14	10	-	-	-
4	6	12	14	10	-	-
5	5	6	12	14	10	-
6	8	5	6	12	14	10
7	3	8	5	6	12	14
$\text{corr}(x_t, x_{t-d})$	1	0.46	0.14	0.61	0.60	-1

En la figura 3.2 se representa la visualización grafica de la tabla 3.1, resaltando con azul aquellos que están dentro de cierta región, en este caso los $x \in [-1, 1]$ tal que $|x| \geq 0.5$ y con rojo los que no.

Figura 3.2: Autocorrelación de la serie de tiempo $(x_t)_{t \in N}$

3.1.3 ESTACIONARIEDAD

De acuerdo con [39] la idea central de la estacionariedad de una serie de tiempo está relacionada con la variabilidad estadística de ésta, es decir, si una serie de tiempo es estacionaria entonces la media, varianza y autocovarianza son constantes a lo largo del tiempo, esto es, no están en función del tiempo. Es importante que una serie de tiempo sea estacionaria, ya que la mayoría de los modelos de pronóstico parten bajo

este supuesto, de hecho, cuando una serie no es estacionaria se aplican distintos métodos para inducir estacionariedad.

La raíz unitaria es una característica que tienen las series de tiempo, la cual proporciona información sobre la inferencia de dicha serie. Se dice que, si la serie tiene una raíz unitaria, ésta posee un patrón impredecible, en otras palabras, no es estacionaria. Mientras que si la serie no tiene raíz unitaria entonces es estacionaria.

La prueba de Dickey-Fuller Aumentada (ADF, por sus siglas en inglés) contrasta la hipótesis nula de que la serie tiene raíz unitaria contra la hipótesis alternativa de que la serie no tiene raíz unitaria. Visto de otro modo, la hipótesis nula es equivalente a que la serie de tiempo no es estacionaria y la hipótesis alternativa es que si lo es [40].

3.2 PRONÓSTICO

Existen diversas maneras de llevar a cabo el pronóstico de una serie de tiempo, en esta sección se describen los conceptos necesarios, así como los modelos de pronóstico.

3.2.1 ELEMENTOS DE PRONÓSTICO

Principalmente se definen las series de **entrenamiento** y **prueba** que sirven para que los modelos optimicen sus parámetros (ó almacenen información) y evalúen qué tan bien se desempeñan, respectivamente.

Una **familia** F con **hiperparámetros** H y **parámetros** W es una ecuación, también llamada modelo matemático $F(H, W)$ en la que se busca optimizar los parámetros W de tal manera que se minimice una métrica en función del error de pronóstico, aquellas ecuaciones que no cuentan con parámetros solamente conservan

la información de los valores al momento de entrenar tal como se muestra en la tabla 3.3.

El **número de observaciones en la serie de prueba** es el número de observaciones a ser pronosticadas en total y que van siendo pronosticadas secuencialmente de acuerdo a una **ventana de pronóstico**, la cual es una tupla (p, h) donde p es el número de las últimas observaciones al pasado que se consideran para pronosticar y h es el número de observaciones a ser pronosticadas también llamado horizonte de pronóstico.

En la sección 3.2.3 se explican los modelos de pronóstico con horizonte de pronóstico igual a uno. posteriormente en la sección 3.2.4 se explica de manera general como se hace para un horizonte mayor que uno.

3.2.2 MÉTRICA DE EVALUACIÓN

Una métrica de evaluación sirve para determinar que tan bien o mal un algoritmo o modelo desempeña una tarea.

Dado que el tipo de problema que se aborda es de pronóstico se debe contar con una métrica de regresión. Se puede utilizar una métrica ya definida o definir una de acuerdo con las necesidades del problema. Para utilizarlas se necesita de la variable de los datos reales en un período de prueba y otra variable que son las predicciones o pronósticos, ambas series, es decir, la serie de prueba y la serie pronosticada deben tener la misma cantidad de observaciones para efectuar una comparación entre lo que sucedió realmente y lo pronosticado. La variable que representa los pronósticos se obtiene de manera secuencial de acuerdo a la ventana de horizonte explicado en la sección 3.2.1.

La principal métrica que se utiliza en este estudio es el error porcentual absoluto medio (MAPE, por sus siglas en inglés) definida por la ecuación 3.2

$$MAPE = \frac{\sum_{i=1}^N \frac{|y_i - y_i^*|}{y_i}}{N}. \quad (3.2)$$

Partiendo de las suposiciones en que se pronostica una serie de prueba con N observaciones, y las variables y_i e y_i^* representan el valor real y el pronóstico de la i -ésima observación en la serie de prueba, respectivamente.

La métrica MAPE está definida como el promedio de los errores porcentuales absolutos. En la figura 3.3 se muestra una representación visual de esta métrica.

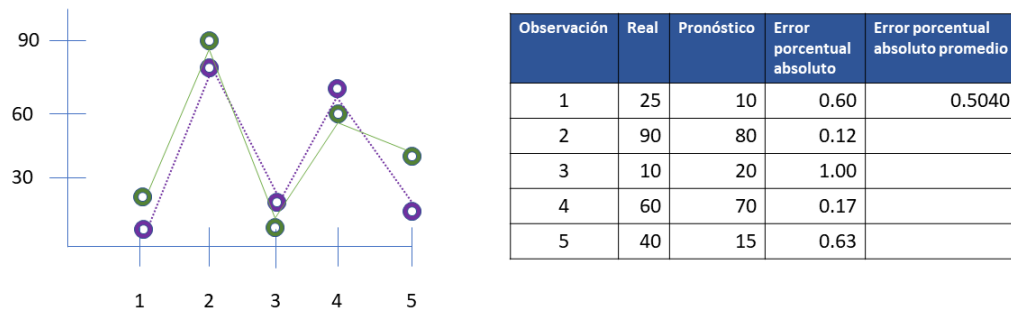


Figura 3.3: Representación visual de la métrica MAPE

Además de la métrica MAPE se consideran otras métricas como: error cuadrático medio (RMSE) y error absoluto medio (MAE), y el coeficiente de determinación ajustado (R^2) las cuales se calculan utilizando las ecuaciones 3.3, 3.4 y 3.5 respectivamente.

$$RMSE = \sqrt{\frac{1}{n} \sum_{i=1}^n (y_i - y_i^*)^2}, \quad (3.3)$$

$$MAE = \frac{1}{n} \sum_{i=1}^n |y_i - y_i^*|, \quad (3.4)$$

$$R^2 = 1 - \frac{\sum_{i=1}^n (y_i - y_i^*)^2}{\sum_{i=1}^n (y_i - \bar{y})^2}. \quad (3.5)$$

Donde las variables y_i e y_i^* representan el valor real y el pronóstico de la i -ésima observación en la serie de prueba, mientras que \bar{y} representa el la media muestral de la variable y , la cual es conocida al saber cual es la serie de prueba.

En la tabla 3.2 se muestran los posibles valores para dichas métricas en base a los pronósticos realizados, además las unidades de éstas y el objetivo deseable en el valor obtenido utilizando dicha métrica.

Tabla 3.2: Descripción de métricas

Métrica	Posibles valores	Unidades	Objetivo
MAPE	$[0, \infty)$	Porcentaje	Minimizar
RMSE	$[0, \infty)$	Variable observada	Minimizar
MAE	$[0, \infty)$	Variable observada	Minimizar
R-cuad	$(-\infty, 1]$	Porcentaje	Maximizar

De acuerdo con [41] la métrica R-cuad puede tomar valores negativos cuando el modelo generalmente tiene un bajo desempeño, es decir, el pronóstico está alejado de el valor real, mientras que un valor cercano a cero representa que los pronósticos son todos cercanos a una constante. Además, existe la posibilidad de que el R-cuad se indetermina cuando el resultado es de la forma $1 - \frac{0}{0}$, en tal caso se toma R-cuad igual a 1.

La métrica MAE da un promedio del error absoluto por lo que con ella se pueden establecer intervalos de estimación del pronóstico. La métrica RMSE está basada en una idea similar, pero obteniendo la raíz cuadrada de del error al cuadrado.

3.2.3 MODELOS DE REGRESIÓN

En este trabajo las ecuaciones de predicción de aprendizaje supervisado que se utilizan son:

1. **Combinación lineal.** En esta ecuación se considera que la salida depende linealmente de las entradas, en este caso las p observaciones al pasado de la ventana de tiempo y está dada por la ecuación 3.6

$$y_t = \sum_{n=1}^p a_n y_n + \epsilon. \quad (3.6)$$

2. **Promedio parcial.** Esta ecuación está enfocada en obtener una predicción de acuerdo al promedio de otras predicciones tal como se muestra en la ecuación 3.7.

$$y_t = \frac{\sum_{n=1}^k y_n}{k} + \epsilon. \quad (3.7)$$

En donde las k predicciones que se utilizan son aquellas en la serie de entrenamiento que cumplen con ciertas condiciones de acuerdo con familia y los hiperparámetros que se están utilizando.

3.2.3.1 FAMILIA DE REGRESIÓN LINEAL

La familia de Regresión lineal (RL) [42] es el modelo más simple de los modelos de regresión, depende linealmente de las variables predictoras y los parámetros son obtenidos por el método de mínimos cuadrados ordinarios (MCO), el cual resuelve el problema dado por la expresión 3.8. Este modelo está dado por la ecuación 3.6 y no tiene ningún hiperparámetro.

$$\min_{a_k} \|Xa_k - y\|_2^2. \quad (3.8)$$

Donde a_k es el k -ésimo parámetro de la ecuación 3.6, X es el vector de entradas, es decir, las observaciones utilizadas para pronosticar e y es la variable de respuesta.

El problema dado por la expresión 3.8 tiene como objetivo determinar los valores de los parámetros a_k que minimizan el error de pronóstico cuadrático.

3.2.3.2 FAMILIA RIDGE

La familia Ridge (RR) [42] al igual que la familia MCO tiene como ecuación de regresión la ecuación 3.6 con la diferencia que, obtiene el valor de los parámetros por el método de mínimos cuadrados ordinarios agregando una penalidad en los parámetros, así como se muestra en el problema dado por la expresión 3.9.

$$\min_{a_k} \|Xa_k - y\|_2^2 + \rho \|Xa_k\|_2^2. \quad (3.9)$$

El problema dado por la expresión 3.9 tiene como objetivo determinar los valores de los parámetros a_k que minimizan el error de pronóstico cuadrático, considerando una penalidad en dichos parámetros.

El hiperparámetro de complejidad ρ controla la cantidad de contracción: cuanto mayor es el valor de ρ , mayor es la cantidad de contracción y, por lo tanto, los coeficientes se vuelven más robustos a la colinealidad.

3.2.3.3 FAMILIA DE K VECINOS MÁS CERCANOS

La familia de K vecinos más cercanos (KNNR) [43] utiliza la ecuación 3.7 para predecir por lo que es necesario definir que observaciones tomar para promediar su respuesta.

Para implementar este modelo primero se define una métrica de distancia en

el hiperespacio \mathbb{R}^n , usualmente se utiliza la distancia euclídeana, la cual está dada por la ecuación 3.10.

$$d(X, Y) = \sqrt{(x_1 - y_1)^2 + \dots + (x_n - y_n)^2}. \quad (3.10)$$

Se identifica cada una de las observaciones en el hiperespacio y para pronosticar una observación dada se localizan las K observaciones más cercanas de acuerdo a la métrica establecida, a estas observaciones se les llama K vecinos más cercanos. El valor pronosticado es el promedio de las variables independientes de los K vecinos más cercanos.

El valor adecuado del hiperparámetro K para realizar el pronóstico se obtiene mediante experimentación computacional.

3.2.3.4 FAMILIA DE ÁRBOLES DE DECISIÓN

La familia de árboles de decisión (DTR) [44] al igual que la familia de K vecinos más cercanos predice de acuerdo a la ecuación 3.7 por lo que también es necesario definir sobre que observaciones se promedian las respuestas.

Funciona dividiendo el hiperespacio de características en regiones rectangulares simples llamadas hiperbloques, divididas por líneas paralelas a los ejes que son representados por condiciones en las variables.

En la figura 3.4 se muestra la visualización de un árbol para una instancia en particular, mientras que en la figura 3.5 se muestra la división del espacio (en este caso \mathbb{R}^2).

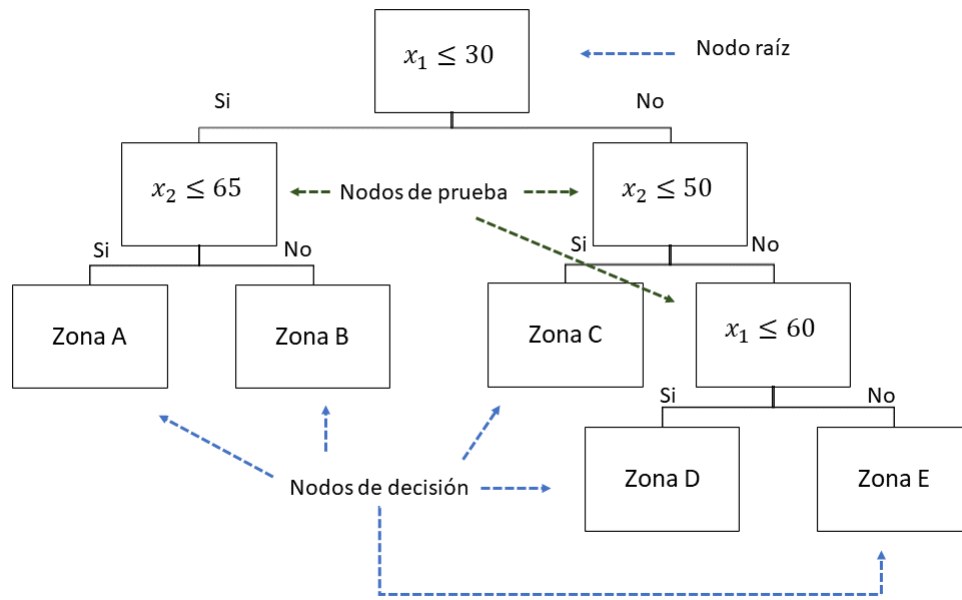


Figura 3.4: Visualización del ejemplo de un árbol de decisión

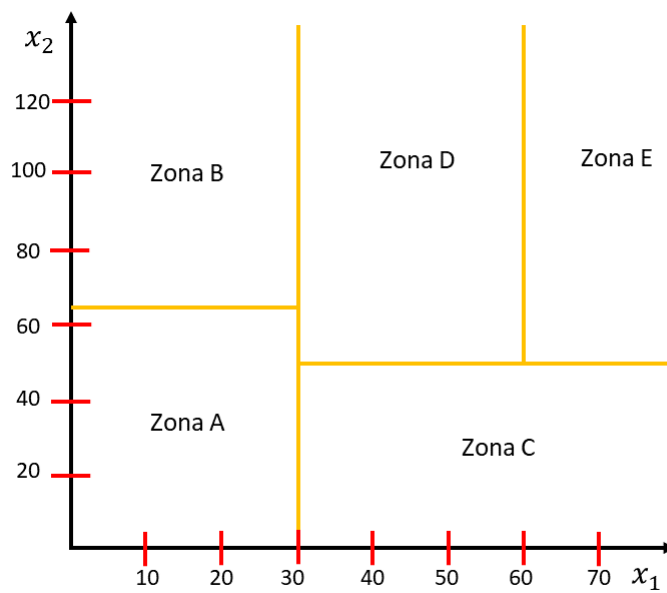


Figura 3.5: Visualización de las divisiones del espacio

Para obtener una predicción a una observación nueva, se utiliza el promedio

de las variables de respuesta de aquellas observaciones en la región en la que ésta nueva observación pertenece.

El hiperparámetro considerado en este modelo es el número mínimo de observaciones que conforman en cada grupo, lo cual beneficia a que no haya exceso de grupos con pocas observaciones.

Una de las desventajas de este algoritmo de regresión es que sobreentrena los datos de entrenamiento, por lo que, para predecir nuevas observaciones se necesita que posean características similares a los datos de entrenamiento para obtener una respuesta precisa.

3.2.3.5 FAMILIA DE BOSQUE ALEATORIO

La familia de bosque aleatorio (RFR) [45] se basa en la familia de árbol de decisión, ya que su funcionamiento generaliza la idea de este último y además también tiene por ecuación de predicción la ecuación 3.7.

El funcionamiento de este modelo es el siguiente: se generan múltiples árboles de decisión diferentes entre sí, y para predecir la variable de respuesta de una observación nueva se evalúa en cada uno de los árboles y finalmente, se promedian los valores obtenidos.

De esta manera, se puede decir que se promedia sobre la respuesta que brinda cada árbol de decisión, y como cada uno es un promedio, entonces se observa que esta ecuación es un promedio sobre promedios.

Los hiperparámetros contemplados son: el número de árboles y el que se contempla para cada árbol de decisión que es, el número mínimo de observaciones en cada bloque creado.

3.2.3.6 PERCEPTRÓN MULTICAPA

El perceptrón multicapa (MLPR) [46] está compuesto por una capa de entrada, una capa de salida y una o varias capas ocultas intermedias. Se caracteriza por tener salidas disjuntas pero relacionadas entre sí, de tal manera que la salida de una neurona es la entrada de la siguiente. Cada capa está conformada por neuronas y un nodo extra que representa el error. Además, las aristas representan un peso los cuales, para los cuales en este trabajo se utiliza el método de propagación hacia atrás para obtener su estimación.

El funcionamiento de este modelo recae en la arquitectura de la red. Partiendo de un conjunto de pesos sinápticos aleatorio, el proceso de aprendizaje busca un conjunto de pesos que permitan a la red predecir de manera precisa. Durante el proceso de aprendizaje se va refinando iterativamente la solución hasta alcanzar un nivel de operación suficientemente bueno.

Para fines prácticos, en este trabajo se utiliza sólo la capa de entrada y de salida, variando el número de neuronas en la capa de salida. De esta manera, el modelo de esta familia con hiperparámetro n está dado por la ecuaciones 3.11, 3.12 y 3.13, que representan la capa oculta, capa de salida y predicción, respectivamente.

$$z_{1,i} = f\left(\sum_{j=1}^p w_{1,ji}x_j + b_1\right) \quad \forall i \in \{1, \dots, n\}, \quad (3.11)$$

$$z_{2,i} = f\left(\sum_{j=1}^n w_{2,ji}z_{1,j} + b_2\right) \quad \forall i \in \{1, \dots, n\}, \quad (3.12)$$

$$y_t = f\left(\sum_{j=1}^n w_{3,j}z_{2,j} + b_3\right). \quad (3.13)$$

Donde p es el número de variables predictoras, f es la función de activación, para este trabajo se considera la ReLu, definida por $f(x) = \max(x, 0)$. Por otra parte, $w_{1,ji} \quad \forall i \in \{1, \dots, n\} \quad \forall j \in \{1, \dots, p\}$, $w_{2,ji} \quad \forall i \in \{1, \dots, n\} \quad \forall j \in \{1, \dots, n\}$ y $w_{3,j} \quad \forall j \in \{1, \dots, n\}$ son los parámetros a estimar en el entrenamiento y $b_i \quad \forall i \in \{1, 2, 3\}$ son

los errores (o sesgo) asignados a la capa de entrada, capa oculta y capa de salida, respectivamente.

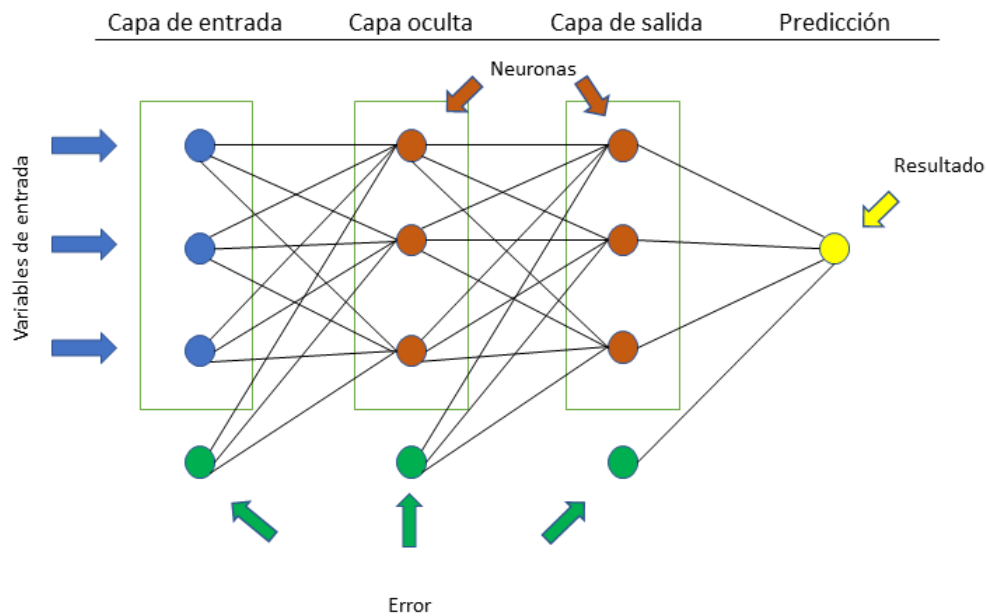


Figura 3.6: Perceptrón multicapa

3.2.4 PRONÓSTICO CON HORIZONTE MAYOR A UNO

Es importante resaltar que en la sección 3.2.3 se describen las ecuaciones de los modelos para las distintas familias con determinados hiperparámetros y a su vez considerando otros aspectos, entre ellos el número de observaciones que se utilizan para pronosticar, pero el horizonte de pronóstico es uno para todas. En esta sección se explica como estos modelos brindan un pronóstico con horizonte mayor a uno siguiendo la estrategia directa de pronóstico.

Suponiendo que se tiene una serie de tiempo con N observaciones, dicha serie de tiempo se particiona en dos, de tal manera que en la primera serie quedan $N - n$ observaciones y en la segunda n , éstas series son de entrenamiento y de prueba, respectivamente. En la figura 3.7 se muestra un ejemplo de una serie de tiempo, con

su respectivas series de entrenamiento y prueba.

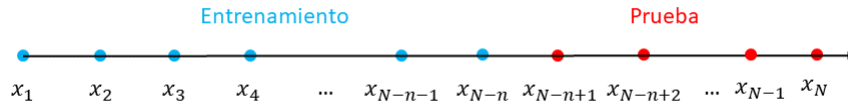


Figura 3.7: Series de entrenamiento y prueba de una serie de tiempo

Una vez obtenidas las series de entrenamiento y prueba se procede a identificar dicha serie de tiempo a manera de una tabla utilizando las observaciones al pasado y el horizonte, tal como se muestra en la tabla 3.3.

Tabla 3.3: Ventanas de pronóstico de aprendizaje y predicción con longitud: N , longitud de prueba: n y ventana de pronóstico: (p, h)

	Pasado					Futuro				
	1	2	...	$p - 1$	p	1	2	...	$h - 1$	h
Entrenamiento	x_1	x_2	...	x_{p-1}	x_p	x_{p+1}	x_{p+2}	...	x_{p+h-1}	x_{p+h}
	x_2	x_3	...	x_p	x_{p+1}	x_{p+2}	x_{p+3}	...	x_{p+h}	x_{p+h+1}
	x_3	x_4	...	x_{p+1}	x_{p+2}	x_{p+3}	x_{p+4}	...	x_{p+h+1}	x_{p+h+2}
	\vdots	\vdots	\vdots	\vdots	\vdots	\vdots	\vdots	\vdots	\vdots	\vdots
	$x_{N-n-h-p+1}$	$x_{N-n-h-p+2}$...	$x_{N-n-h-1}$	x_{N-n-h}	$x_{N-n-h+1}$	$x_{N-n-h+2}$...	x_{N-n-1}	x_{N-n}
Prueba	$x_{N-n-p+1}$	$x_{N-n-p+2}$...	x_{N-n-1}	x_{N-n}	x_{N-n+1}	x_{N-n+2}	...	$x_{N-n+h-1}$	x_{N-n+h}
	$x_{N-n-p+1+h}$	$x_{N-n-p+2+h}$...	$x_{N-n-1+h}$	x_{N-n+h}	$x_{N-n+1+h}$	$x_{N-n+2+h}$...	$x_{N-n+2h-1}$	x_{N-n+2h}
	$x_{N-n-p+1+2h}$	$x_{N-n-p+2+2h}$...	$x_{N-n-1+2h}$	x_{N-n+2h}	$x_{N-n+1+2h}$	$x_{N-n+2+2h}$...	$x_{N-n+3h-1}$	x_{N-n+3h}
	\vdots	\vdots	\vdots	\vdots	\vdots	\vdots	\vdots	\vdots	\vdots	\vdots

Con dicha tabla se procede a entrenar el modelo $F(H, W, P)$ con la serie de entrenamiento, es decir, se obtienen los parámetros del modelo (si tiene, sino simplemente se guardan las observaciones con sus respuestas ya que algunas familias no tienen parámetros, pero utiliza información de dichas las observaciones con sus respectivas salidas).

En la figura 3.8 se muestran los diferentes modelos matemáticos, en donde se

resaltan el número de observaciones al pasado (variables predictoras) y el horizonte de pronóstico (variables dependientes) y como está dada la ecuación con notación escalar y vectorial.

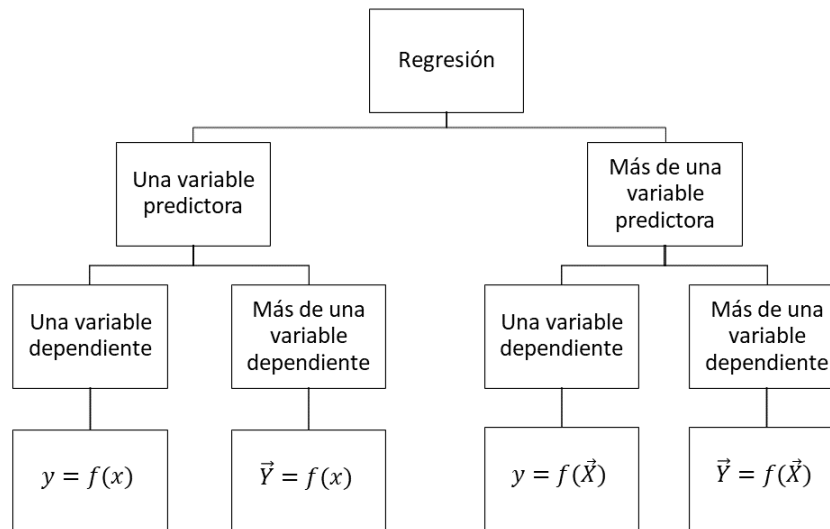


Figura 3.8: Clasificación de modelos de regresión de acuerdo al número de variables

En el caso en que el horizonte de pronóstico es mayor a uno es en los bloques correspondientes a más de una variable independiente y esto es posible en los casos en que se tiene una sola variable predictora o más de una, que representan ventanas de pronóstico $V_1 = (1, h)$ y $V_2 = (p, h)$ respectivamente.

En la figura 3.9 se observa un ejemplo para el cual se muestra una ventana de pronóstico cuyo valor es de V_1 y también para V_2 . Para ambos casos se supone que $N = 10$, $n = 4$ y $h = 2$, donde n es la longitud de la serie de prueba.

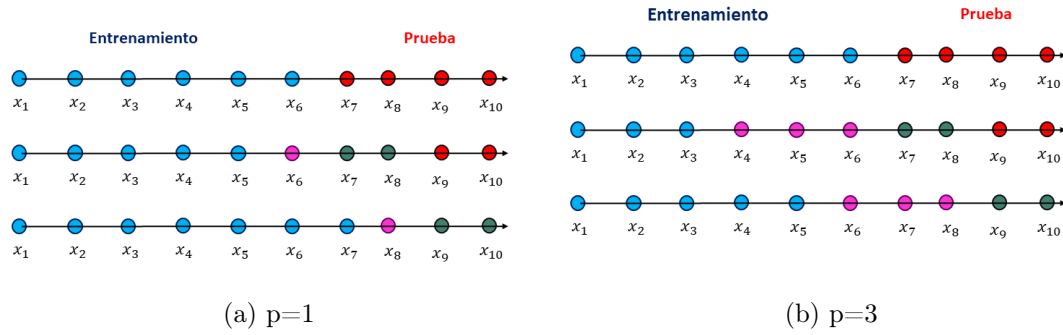


Figura 3.9: Ejemplo con $N = 10$ y $n = 4$

La idea principal es tratar al problema de regresión multisalida como múltiples problemas de regresión con una sola salida, en la que se conservan las mismas variables predictoras pero con una variable independiente para cada problema. En la figura 3.4 se observa como realizar la representación de las ventanas de pronóstico del ejemplo con múltiples salidas a múltiples problemas con una salida, utilizando el ejemplo de la figura 3.9 para $p = 3$.

Tabla 3.4: Ventanas de pronóstico para el ejemplo de la figura 3.9 (b)

Pasado			Futuro	
1	2	3	1	2
x_1	x_2	x_3	x_4	x_5
x_2	x_3	x_4	x_5	x_6
x_4	x_5	x_6	x_7	x_8
x_6	x_7	x_8	x_9	x_{10}

Se observa que cada ventana de pronóstico es una observación al problema de regresión. Además las primeras dos observaciones son utilizadas para entrenar y las últimas para llevar a cabo el pronóstico. En la tabla 3.5 se muestra como el problema de regresión multisalida se puede tratar desde un enfoque en que se tienen múltiples problemas de regresión de una salida.

Tabla 3.5: Regresión multisalida como múltiples regresiones con una salida

(a) primer componente				(b) segunda componente			
Pasado			Futuro	Pasado			Futuro
1	2	3	1	1	2	3	2
x_1	x_2	x_3	x_4	x_1	x_2	x_3	x_5
x_2	x_3	x_4	x_5	x_2	x_3	x_4	x_6
x_4	x_5	x_6	x_7	x_4	x_5	x_6	x_8
x_6	x_7	x_8	x_9	x_6	x_7	x_8	x_{10}

Ya que se tienen problemas de regresión con una sola salida, se puede aplicar cualquier modelo que se ha descrito anteriormente. Es importante observar que la predicción para cada salida se puede llevar a cabo de diferentes maneras (ordenada), De esta manera se dice que para cada componente del vector de salidas se crea una ecuación que sigue un modelo específico. En este trabajo todas las componentes (del horizonte de pronóstico) utilizan el mismo modelo (Familia e hiperparámetros).

En la tabla 3.5 se observa que cada componente cuenta con un modelo de pronóstico y, así, ya se puede realizar el pronóstico para una ventana de pronóstico general $V = (p, h)$.

CAPÍTULO 4

METODOLOGÍA

En este capítulo se describen los datos, su preprocesamiento para llevar a cabo la evaluación computacional. Además, se presenta la metodología a seguir para efectuar el pronóstico, en la que resalta la selección de valores para diferentes parámetros, así como selección de modelo entre distintos posibles y la evaluación del desempeño del pronóstico utilizando dicho modelo.

4.1 DESCRIPCIÓN Y PREPROCESAMIENTO DE LOS DATOS

Los datos con los que se trabaja en este proyecto fueron proporcionados por personal de la CFE división Golfo Norte, dichos datos están en 6 archivos en formato `.xlsx` y en cada uno de ellos están los reportes de fallas para las 4 zonas para un año en específico, desde el 2013 al 2018

En el reporte de fallas una observación corresponde a una falla reportada y que fue atendida por el personal de mantenimiento y se cuenta con diferentes características para las fallas, tales como fecha, lugar en el que ocurrió (subestación, circuitos y ramal), causa que la origino, etc.

Dado que los datos estaban separados inicialmente por año, se procede a juntarlos o concatenarlos para así tener en un mismo archivo las fallas en todas las zonas y con todos los años. Posteriormente se particiona por zona con el objetivo de tener las fallas en el período dado para cada zona. Para crear la serie de tiempo de cada zona se selecciona una frecuencia para realizar el conteo de fallas, esta puede ser por día, mes, año, etc. En la figura 4.1 se muestra una representación visual del preprocesamiento de los datos para obtener las series de tiempo de fallas en cada zona.



Figura 4.1: Preprocesamiento de datos para la creación de series de tiempo

4.2 METODOLOGÍA DE PRONÓSTICO

En la figura 4.2 se muestran todos los pasos a seguir desde que se tiene una serie de tiempo hasta que se selecciona un modelo perteneciente a una determinada familia de pronóstico para finalmente, evaluar su desempeño utilizando los errores porcentuales absolutos de las observaciones pronosticadas. La selección de valores de los parámetros es basada en minimizar la métrica MAPE.

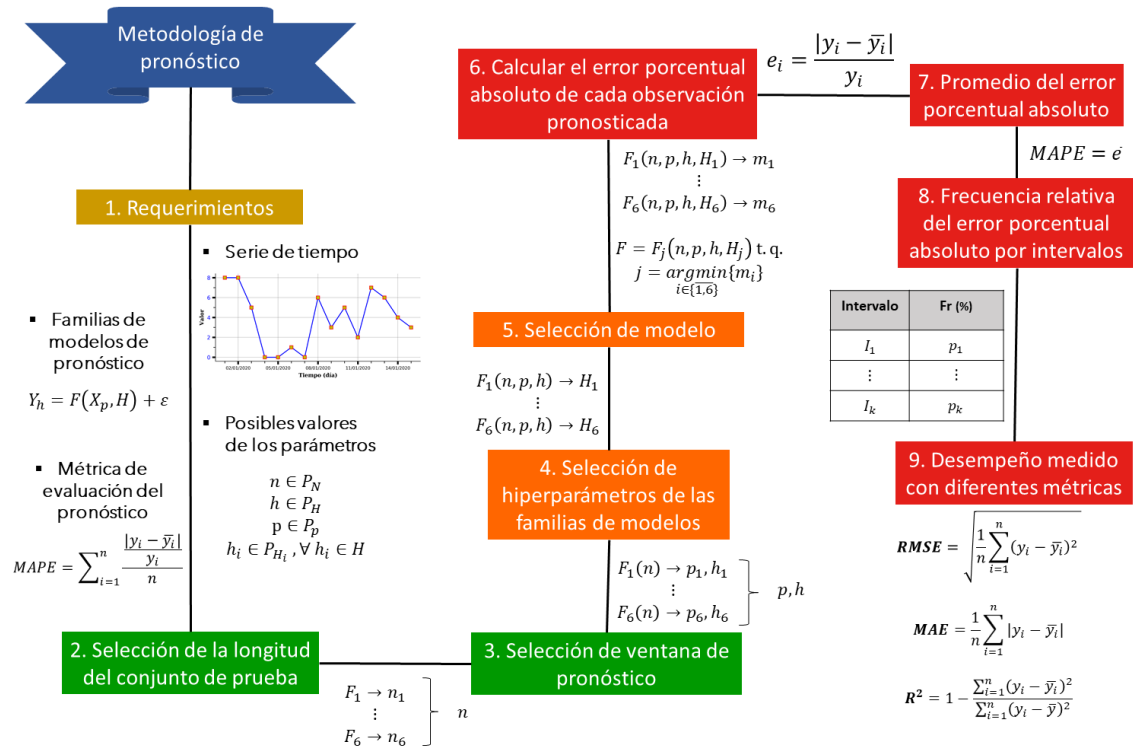


Figura 4.2: Metodología de pronóstico

1. **Requerimientos.** Los elementos necesarios para llevar a cabo satisfactoriamente el pronóstico de una serie de tiempo son:

- Serie de tiempo univariada con N observaciones, $(x_t)_{t \in T}$
- Familias de modelos, en este caso las familias que se utilizan son modelos de regresión dentro del aprendizaje automático. La ecuación de predicción y los hiperparámetros de dichos modelos son descritos en la sección ??.
- Posibles valores para cada uno de los parámetros con los que el pronóstico puede ser efectuado, dichos parámetros son: longitud del conjunto de prueba (n), ventana de pronóstico ($V = (p, h)$), Familia de modelos (F) con su respectiva configuración de hiperparámetros (H). Al conjunto de posibles valores se les asigna el nombre de $P_n, P_p, P_h, P_F, P_{H_i}$, en donde H_i hace referencia al i -ésimo hiperparámetro de la familia F .
- La métrica de evaluación permite conocer el desempeño del pronóstico de

un modelo en particular. De esta manera, se puede conocer el de pronósticos brindados por diferentes modelos para determinar cual obtiene mayor precisión basandose en que para esta métrica, un menor valor proporciona un mejor desempeño.

La serie de tiempo se pronostica con todos los posibles valores en cada parámetro para analizar la variación del desempeño el cual es evaluado con la métrica MAPE y seleccionar los valores adecuados de estos, en cada uno de los siguientes pasos se especifica dicha selección para los modelos deterministas y no deterministas, para estos últimos se considera una cantidad r de réplicas.

2. Selección de longitud del conjunto de prueba. La selección del parámetro n se divide en dos fases.

La primer fase consiste en determinar aquel valor $n_i \in P_n$ para cada familia de modelo, de acuerdo a la naturaleza de estos.

- Modelos deterministas. Para una familia de modelos deterministas F_i se consideran los pronósticos efectuados con todas las combinaciones en los valores de los parámetros, analizando la variación del MAPE por los posibles valores de n en la longitud del conjunto de prueba. Se selecciona n_i como el elemento de P_n con el que se obtiene la menor mediana en el MAPE.
- Modelos no deterministas. A diferencia de las familias con modelos deterministas, para la familia que proporciona modelos no deterministas F_i , dado que se cuenta con réplicas, también se consideran los pronósticos efectuados con todas las combinaciones en los valores de los parámetros, pero en este caso se analiza la variación de la mediana del MAPE encontrado en las réplicas realizadas de acuerdo a los posibles valores de n . Se selecciona n_i como el elemento de P_n con el que se obtiene la menor mediana en las medianas del MAPE.

La fase dos consta de seleccionar de manera general un valor para el parámetro n . Con la finalidad de que la mayor cantidad de familias tengan la configuración que minimiza el MAPE, se selecciona n como el valor con mayor frecuencia en el conjunto $\{n_1, \dots, n_{C_F}\}$, donde C_F es el número de familias consideradas para realizar el pronóstico. En caso de ser necesario, se define el criterio de desempate como n igual al mayor valor, ya que brinda mayor repeticiones utilizando la ventana de pronóstico.

3. Selección de ventana de pronóstico. Para seleccionar la ventana de pronóstico $V = (p, h)$ se consideran los pronósticos con longitud en el conjunto de prueba n obtenido en el paso 2, similar a dicho paso, se consideran dos fases para seleccionar el valor de los parámetros p y h .

La fase uno consiste en la selección específica de la ventana de pronóstico para cada familia.

- Modelos deterministas. Para una familia de modelos deterministas F_i se considera el MAPE de los pronósticos efectuados para cada combinación $(p_{i_j}, h_{i_k}) \in P_p \times P_h$ con longitud en el conjunto de prueba n y a cada (p_{i_j}, h_{i_k}) se le asigna el menor valor posible para el MAPE, ya que, como es un modelo determinista, se puede elegir la configuración de hiperparámetro(s) que brinda dicho error.
- Modelos no deterministas. De la misma manera en que se consideran los pronósticos con longitud en el conjunto de prueba n para los modelos deterministas, una familia con modelos no deterministas considera dichos pronósticos con la diferencia de que efectuando el pronóstico en múltiples ocasiones no hay garantía de obtener el mismo MAPE, por lo que a cada configuración de hiperparámetro(s) se considera la mediana del MAPE obtenido en las réplicas. A cada (p_{i_j}, h_{i_k}) se le asigna el menor valor posible para la medianas del MAPE, ya que es posible determinar la configuración de hiperparámetro(s) que tiene la menor mediana en el MAPE.

Para la familia F_i se considera la ventana de pronóstico (p_i, h_i) con menor valor de la variable en cuestión (MAPE o la mediana de este, según sea el caso). Si hay empate entre dos o más ventanas de pronóstico, el primer criterio de desempate es elegir aquella con mayor valor en el parámetro h , mientras que si sigue ocurriendo empate, el segundo criterio es seleccionar la de menor valor en el parámetro p .

La segunda fase se basa en determinar una ventana de pronóstico (p, h) de manera general. Con el motivo de aprovechar el mejor desempeño de las familias, se considera aquella ventana de pronóstico que obtiene el menor valor de la variable medida de manera global, es decir, aquella ventana de pronóstico que proporciona el menor valor en comparación de todas las ventanas de pronóstico posibles para todas las familias posibles.

4. **Selección de hiperparámetros de las familias de modelos.** La configuración de hiperparámetro(s) $H_i \in P_{H_i}$ que se selecciona para la familia F_i depende de la naturaleza de este.

- Modelos deterministas. Para una familia de modelos deterministas F_i se consideran los pronósticos efectuados con la longitud del conjunto de prueba n y ventana de pronóstico (p, h) , se selecciona la configuración de hiperparámetros H_i , de tal manera que se minimiza el MAPE.
- Modelos no deterministas. De manera semejante a los modelos deterministas, para una familia de modelos no deterministas F_i se consideran los mismos pronósticos (con réplicas), seleccionando la configuración de hiperparámetros H_i como aquella que minimiza la mediana del MAPE.

De esta manera, se tiene un modelo representativo para cada familia de modelos F_i para el pronóstico de la serie de tiempo con longitud en el conjunto de prueba n y ventana de pronóstico (p, h) .

5. **Selección de modelo.** Para elegir el modelo con el que se pronosticará la serie de tiempo con longitud en el conjunto de prueba n y ventana de pronóstico

(p, h) se compara el valor del MAPE que proporciona cada modelo representativo de las diferentes familias. Se selecciona aquel modelo que minimiza el MAPE.

6. **Calcular el error porcentual absoluto de cada observación pronosticada.** Una vez que se tiene el modelo $Y_h = F(X_p, H) + \epsilon$ se pronostica el conjunto de prueba y se calcula el error porcentual absoluto de cada observación en el conjunto de prueba, las cuales permitirán evaluar el desempeño del pronóstico obtenido por el modelo en cuestión.
7. **Promedio del error porcentual absoluto.** Como se menciona en la sección 2 la métrica MAPE es una de las más utilizadas en la literatura por lo que se procede a evaluar el pronóstico.
8. **Error porcentual absoluto por intervalos.** Después que se evalúa el desempeño el pronóstico con la métrica MAPE, se calcula la distribución de las observaciones pronosticadas en intervalos deseados para obtener información más detallada de los pronósticos.

Si $I = \{I_1, \dots, I_k\}$ es una partición finita de $[0, \infty)$, entonces se puede calcular la frecuencia relativa en la que el error porcentual absoluto pertenezca a I_j la cual se calcula mediante la ecuación 4.1,

$$fr_j = \frac{|I_j|}{|I|}, \forall j \in \overline{1, k}. \quad (4.1)$$

9. **Desempeño medido con diferentes métricas.** Finalmente, con el objetivo de tener más información del desempeño del pronóstico por el modelo seleccionado, se aplican distintas métricas, entre ellas:

- RMSE
- MAE
- R^2

En la sección 3 se muestran las ecuaciones necesarias para la evaluación del desempeño de los pronósticos.

EXPERIMENTACIÓN Y RESULTADOS

Los experimentos realizados en esta investigación se realizaron en una PC con un sistema operativo Windows 10 de 64 bits, Intel Core i7-8550U @ 1.80 GHz 2.00 GHz y 8 GB de memoria RAM.

El lenguaje de programación que se utiliza es `python` [47] y algunas de las principales librerías que son necesarias para completar la experimentación son `matplotlib` [48], `pandas` [49], `statsmodels` [50] y `sklearn` [41] que sirven principalmente para la manipulación de datos, creación y prueba de los modelos estadísticos, así como el ambiente gráfico necesario.

5.1 CARACTERÍSTICAS DE LA EXPERIMENTACIÓN

Se cuenta con cuatro series de tiempo, cada una representa el número de fallas de cada zona que comprende la cobertura de la División Golfo Norte a lo largo del período de enero del 2013 a diciembre del 2018 (314 semanas). El conteo de las fallas es llevado a cabo de manera fija, para el cual se consideran lapsos semanales.

De acuerdo a la metodología presentada en el capítulo 4, el pronóstico se realiza para cada una de las series de tiempo y los parámetros a considerar son, principalmente, la longitud del conjunto de prueba, la ventana de pronóstico, así como

también los hiperpárametros de cada modelo. Además, dado que los valores de las observaciones de las series de tiempo son enteros, al obtener el pronóstico de cada modelo, cada observación predicha se redondea al entero más cercano.

En la tabla 5.1 se muestran los posibles valores que se consideran para la longitud de prueba y ventana de pronóstico, mientras que en la tabla 5.2 se muestran aquellos posibles valores para los hiperparámetros de los diferentes modelos.

Tabla 5.1: Posibles valores de la longitud del conjunto de prueba y ventana de pronóstico

Característica	Posibles valores
Longitud del conjunto de prueba	$\{30,40,50\}$
Ventana de pronóstico	$([1, 20] \times [1, 12]) \cap \mathbb{Z}^2$

Tabla 5.2: Posibles valores de los hiperparámetros en los modelos

Modelo	Hiperparámetro(s)	Valores a probar
Regresión lineal	-	-
Regresión Ridge	β : penalización en los coeficientes	{5,10,15,20,25}
K vecinos más cercanos	K : número de vecinos para promediar	{3,4,...,20}
Árbol de decisión	m : Mínima cantidad de elementos por bloque	{5,10,15,20}
Bosque aleatorio	A : Cantidad de árboles	{30,40,50,60,70}
	m : Mínima cantidad de elementos por bloque	{5,10,15,20}
Perceptrón multicapa (1 capa oculta, FA: ReLu)	N : Neuronas	{5,6,7,8,9,10}

5.2 EXPERIMENTACIÓN

En las secciones 5.2.1, 5.2.2, 5.2.3 y 5.2.4 se realiza un análisis de cada zona en específico que abarca el estudio, correspondientes a Montemorelos, Zona Metropolitana Norte, Zona Metropolitana Oriente y Zona Metropolitana Poniente, respectivamente.

Para el conteo de las fallas es importante saber que cada zona cuenta con la presencia de vegetación distinta, así como diferentes aspectos ambientales que pueden influenciar a que un árbol o planta provoque una falla.

5.2.1 MONTEMORELOS

En la figura 5.1 se observa la serie de tiempo para el número de fallas en Montemorelos, resaltando que en el 2015 y 2016 se presentan observaciones con el mayor número de fallas, mientras que en el 2013 y 2014 se presentan observaciones con menor cantidad, siendo inferiores a 10 fallas en cada semana contemplada.

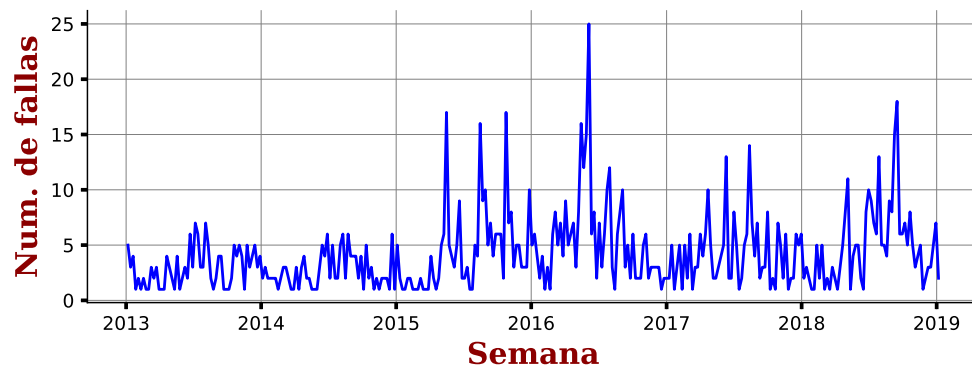


Figura 5.1: Fallas en Montemorelos

En la figura 5.2 se muestra la serie de tiempo en Montemorelos particionada por año, en la que se observa que en cada año, el máximo es alcanzado entre la semana 20 y la 45, por otra parte, para el valor mínimo de fallas no hay un período fijo.

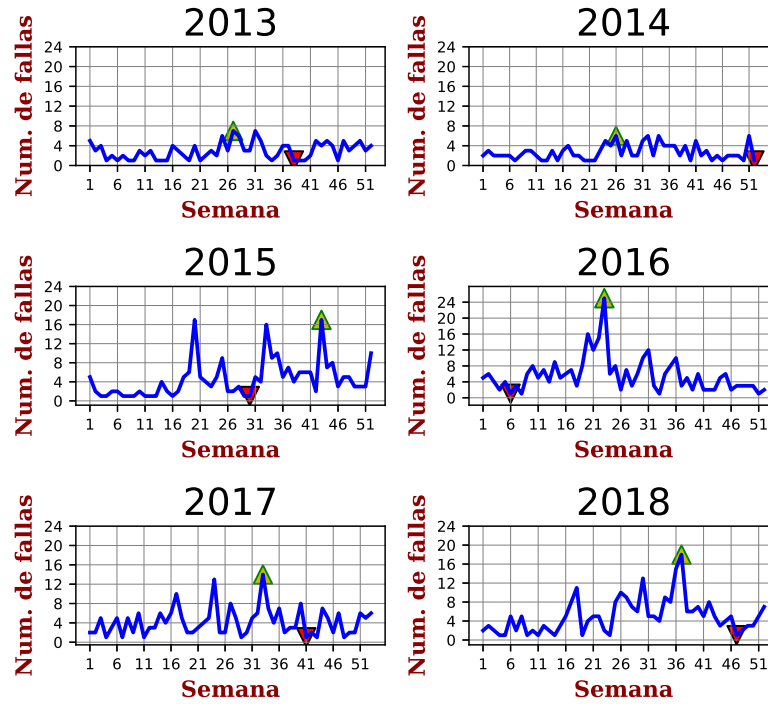


Figura 5.2: Fallas en Montemorelos por año

En la figura 5.3 se observa la descomposición aditiva de la serie de tiempo para las fallas en Montemorelos, en la que se observa que hay una tendencia creciente entre los años 2013 al 2015, mientras que del 2016 al 2018 es decreciente. Por otra parte, en la misma figura se observa que basándose en la estacionalidad hay un patrón parabólico, donde el vértice es aproximadamente la mediación de año y se presentan bajas y altas alternantes.

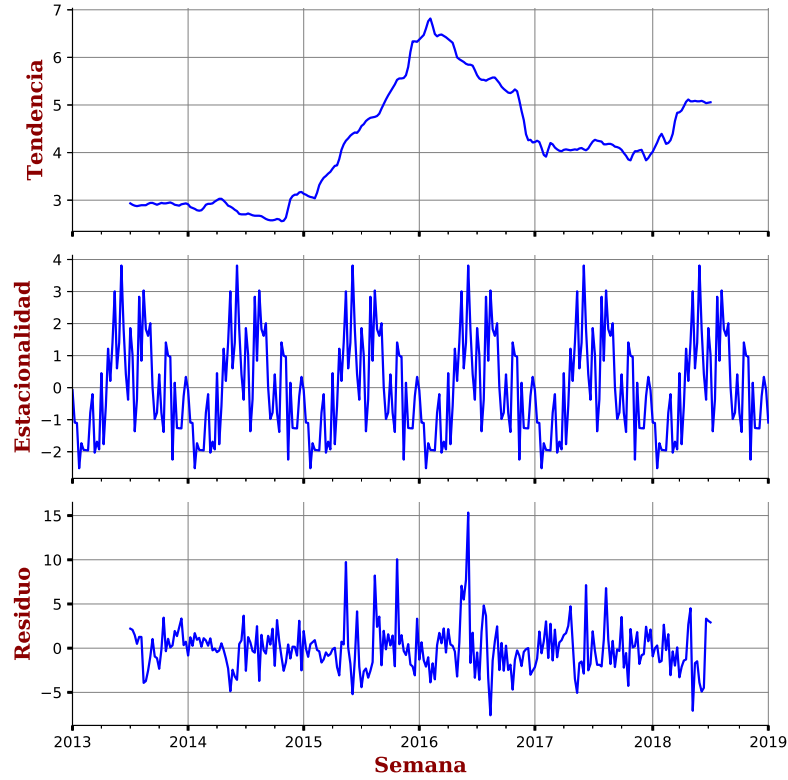


Figura 5.3: Descomposición aditiva de las fallas en Montemorelos



Figura 5.4: Autocorrelación de fallas en Montemorelos

En la figura 5.4 se observa que dadas las fallas de una semana específica en Montemorelos, en general, las semanas con la cantidad de fallas similar a dicha semana son la primera, segunda y novena anteriores a esta.

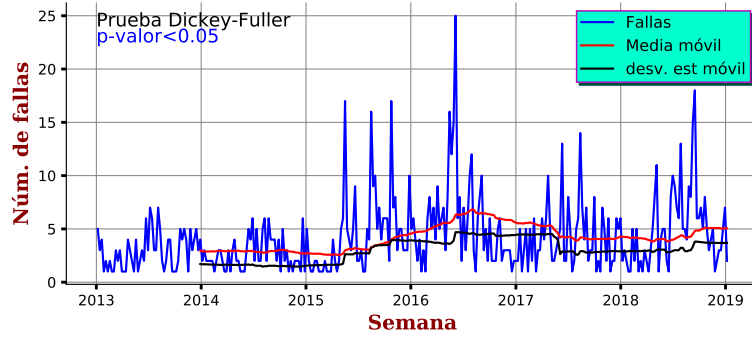


Figura 5.5: Estacionariedad de fallas en Montemorelos

En la figura 5.5 se observa que el p-valor obtenido al realizar la prueba es menor a $\alpha = 0,05$ por lo que se rechaza la hipótesis nula de que la serie de tiempo tiene raíz unitaria, en favor de que esta serie de tiempo sea estacionaria.

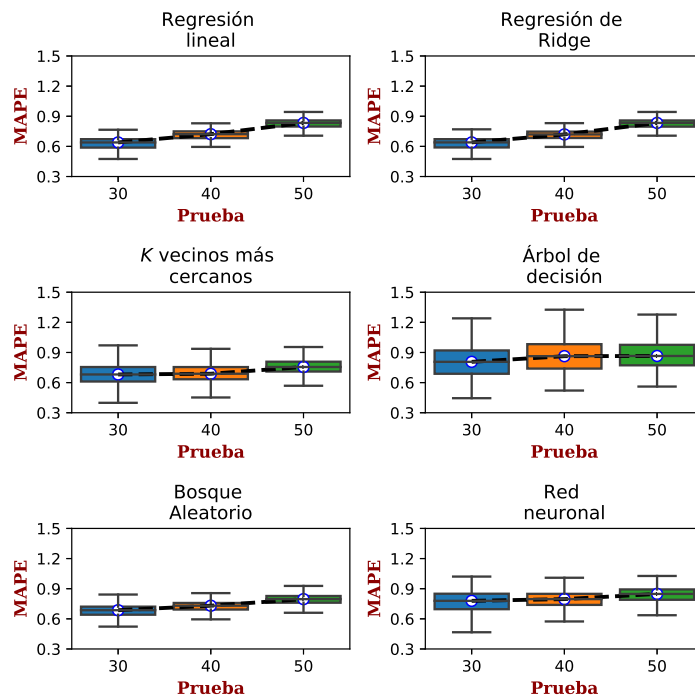


Figura 5.6: MAPE por longitud de prueba para cada familia en Montemorelos

En la figura 5.6 se observa el MAPE obtenido de acuerdo a la variación en la longitud en el conjunto de prueba para cada familia. En esta figura se observa que

el valor que presenta menor mediana en el MAPE es $n_i = 30 \forall i \in 1, 2, 3, 4, 5, 6$, por lo que de manera general se contempla $n = 30$.

En la figura 5.7 se muestra la serie de tiempo de fallas en Montemorelos dividida en el conjunto de entrenamiento y de prueba. El conjunto de prueba está conformado por $n = 30$ observaciones.

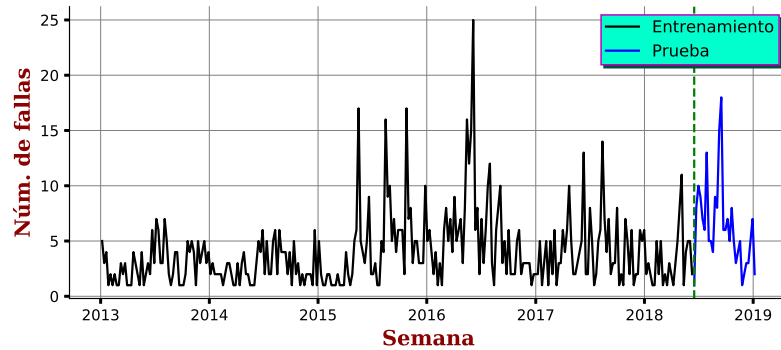


Figura 5.7: Conjuntos de entrenamiento y prueba de fallas en Montemorelos

En la figura 5.8 se muestra el MAPE obtenido variando las ventanas de pronóstico para cada familia dada la longitud del conjunto de prueba de $n = 30$. Se observa que la familia que tiene menor MAPE posible para sus ventanas de pronóstico es K vecinos más cercanos y con $V = (1, 3)$, por lo que se elige dicha ventana de pronóstico para obtener mayor ventaja de este menor MAPE.

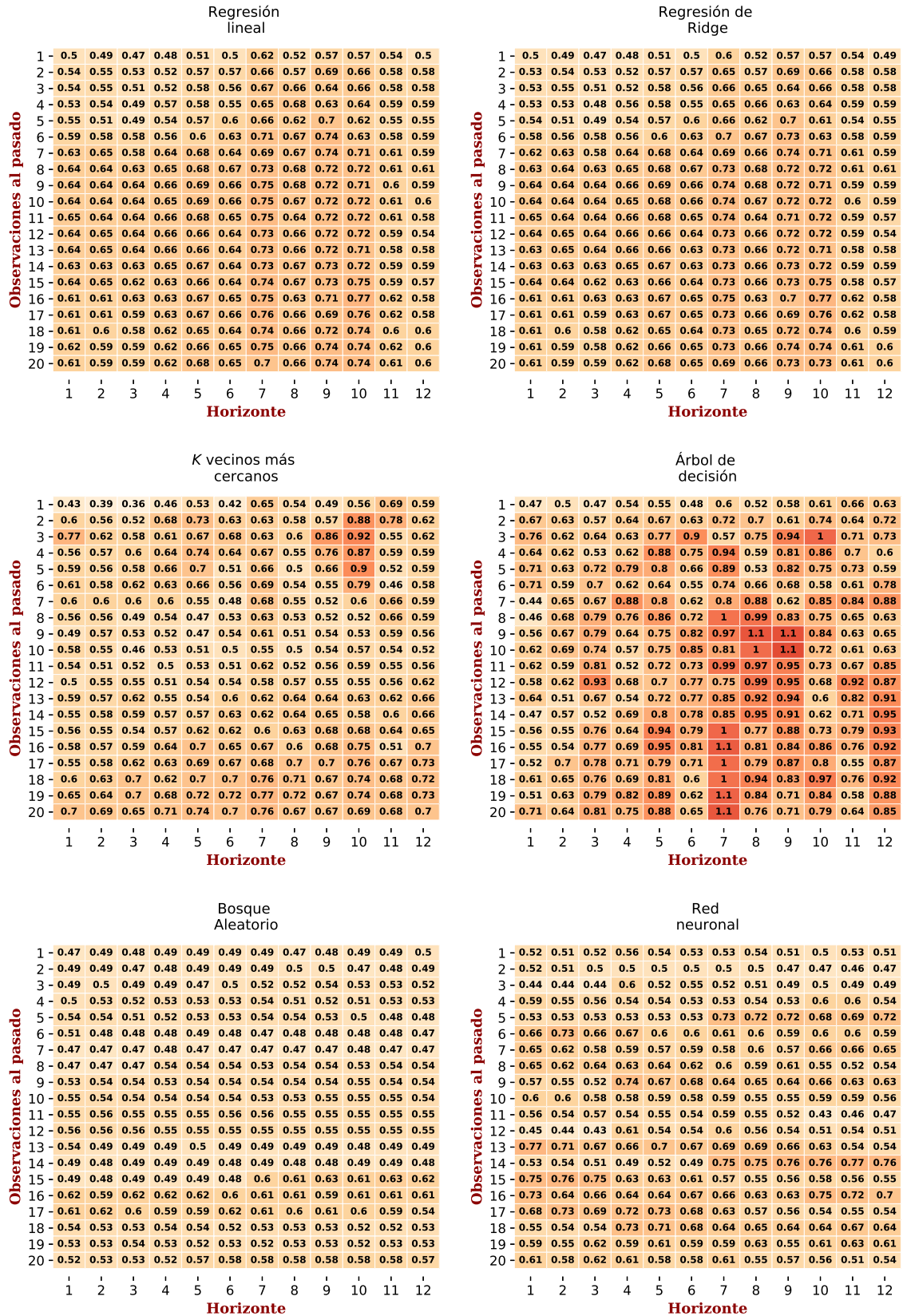


Figura 5.8: MAPE mínimo por combinación de ventana de pronóstico para cada familia en Montemorelos (N=30)

En la figura 5.9 se muestra el valor del MAPE obtenido variando la configuración de hiperparámetro(s) para cada familia de modelos en el pronóstico de fallas en Montemorelos. De esta manera, se obtiene aquel modelo que representa cada familia dada la longitud del conjunto de prueba y la ventana de pronóstico.

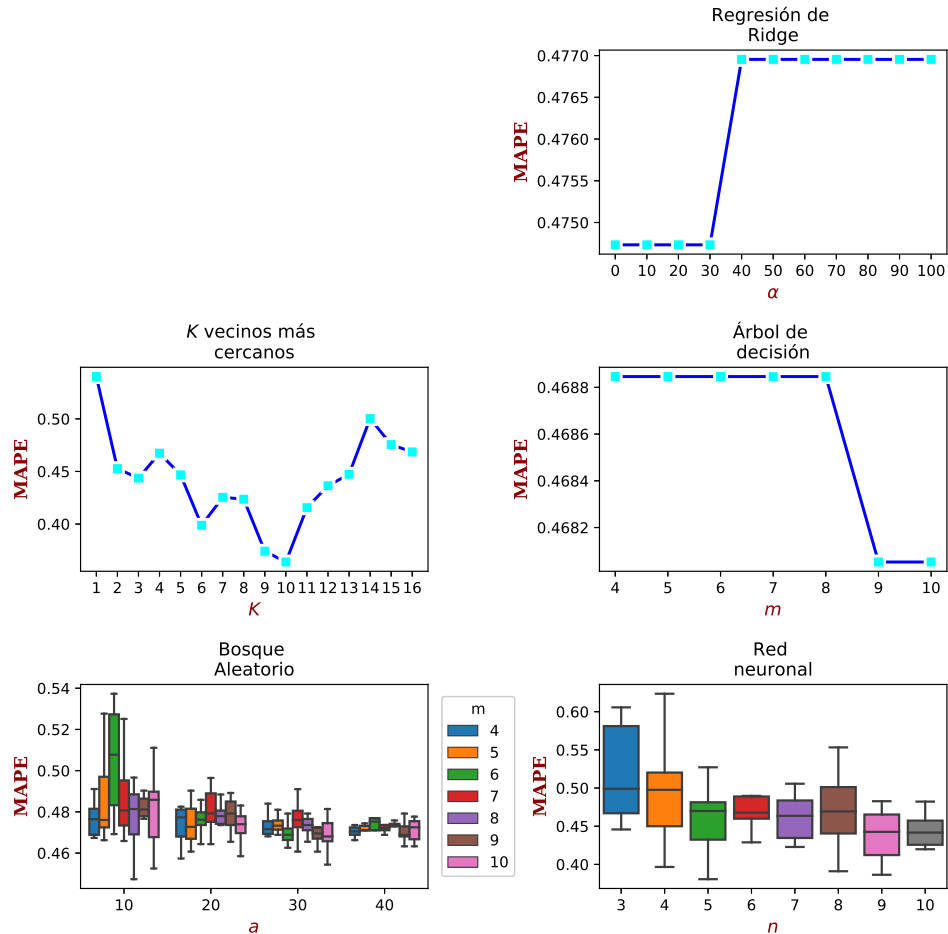


Figura 5.9: MAPE por hiperparámetro(s) de cada familia en Montemorelos ($N=30$, $V=(1,3)$)

Ya que se han elegido los hiperparámetros correspondientes de las familias de modelos para la configuración de la longitud del conjunto de prueba $n = 30$ y la ventana de pronóstico $V = (1, 3)$, en la figura 5.10 se compara el MAPE obtenido por cada uno para seleccionar aquel modelo con menor valor.

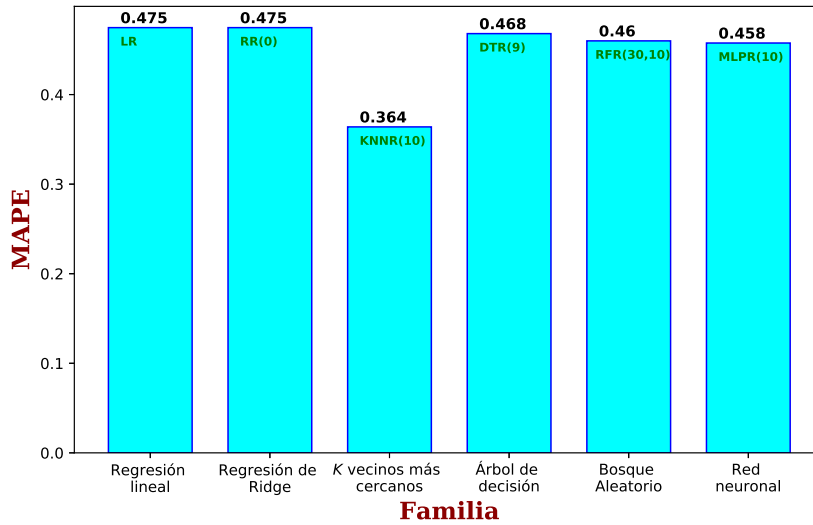


Figura 5.10: Comparación de MAPE de los modelos representativos de cada familia en Montemorelos (N=30, V=(1,3))

En la figura 5.11 se muestran las observaciones del conjunto de prueba de la serie de tiempo de fallas en Montemorelos y el pronóstico brindado por el modelo con menor MAPE para los parámetros seleccionados.

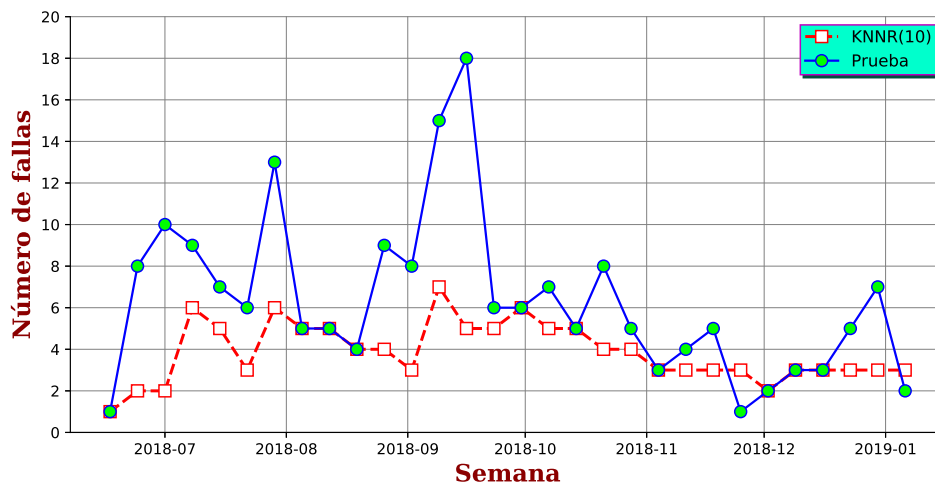


Figura 5.11: Pronóstico de fallas en Montemorelos (N=30, V=(1,3), M=KNNR(10))

En la figura 5.12 se observa la frecuencia relativa porcentual del error porcentual absoluto obtenido por el pronóstico con el respectivo modelo y las características previamente seleccionadas, en la que se resalta que el grupo con mayor frecuencia relativa porcentual es el de las observaciones pronosticadas con 0% de error absoluto porcentual.

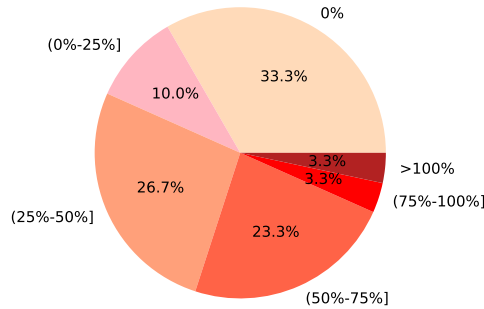


Figura 5.12: Distribución del error porcentual absoluto por intervalos en Montemorelos ($N=30$, $V=(1,3)$, $M=KNNR(10)$)

En la tabla 5.3 se muestran los valores de distintas métricas, las cuales representan el desempeño del pronóstico obtenido para la serie de tiempo de las fallas en Montemorelos con $n = 30$, $V = (1, 3)$ y empleando el modelo de regresión de los k vecinos más cercanos con hiperparámetro $k = 10$.

Tabla 5.3: Desempeño del pronóstico en Montemorelos

MAPE	RMSE	MAE	R-cuad
0.3640	3.1300	2.6666	0.7876

5.2.2 ZONA METROPOLITANA NORTE

En la figura 5.13 se observa la serie de tiempo para el número de fallas en la Zona Metropolitana Norte, resaltando que el año con la observación con mayor

número de fallas es el 2016 con 12, posteriormente el año 2017 con 8, mientras que el 2013 y el 2014 son los años con observaciones con menor número de fallas, siendo inferiores a 6 fallas. Además esta serie de tiempo se caracteriza por tener muchas observaciones con valor de uno.

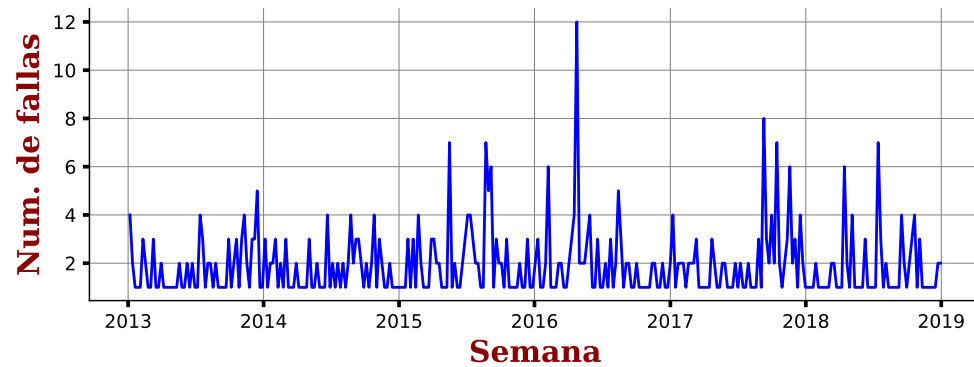


Figura 5.13: Fallas en la Zona Metropolitana Norte

En la figura 5.14 se muestra la serie de tiempo en la Zona Metropolitana Norte particionada por año, en la que se observa que en cada año, el máximo es alcanzado después de la 15ava semana, por otra parte, para el valor mínimo de fallas no hay un período fijo.

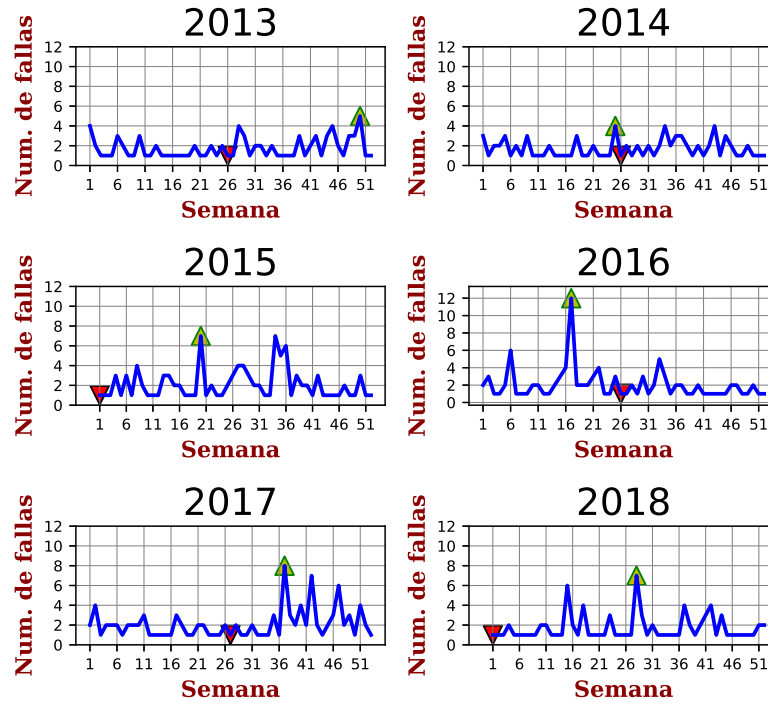


Figura 5.14: Fallas en la Zona Metropolitana Norte por año

En la figura 5.15 se observa la descomposición aditiva de la serie de tiempo para las fallas en la Zona Metropolitana Norte, en la que se observa que hay una tendencia que oscila alrededor de 2, en donde se tiene que para el 2015 y casi todo el 2016, así como algunos pequeños períodos del 2017 y 2018 están por encima, mientras que el resto están por debajo. Por otra parte, en la misma figura se observa que basándose en la estacionalidad el mayor número de fallas se obtiene entre marzo y mediados de abril, mientras que el menor se alcanza entre mayo y junio .

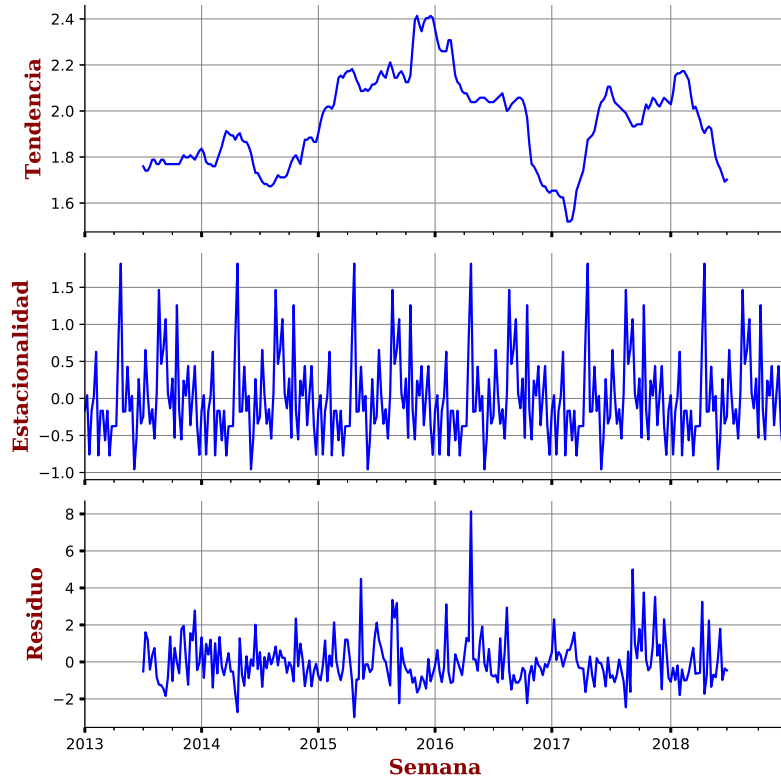


Figura 5.15: Descomposición aditiva de las fallas en la Zona Metropolitana Norte

En la figura 5.16 se observa que dadas las fallas de una semana específica en la zona Metropolitana Norte, en general, ninguna de las 52 semanas anteriores presenta una similitud estadística.



Figura 5.16: Autocorrelación de fallas en la Zona Metropolitana Norte

En la figura 5.17 se observa que el p-valor obtenido al realizar la prueba es menor a $\alpha = 0,05$ por lo que se rechaza la hipótesis nula de que la serie de tiempo tiene raíz unitaria, en favor de que esta serie de tiempo sea estacionaria.

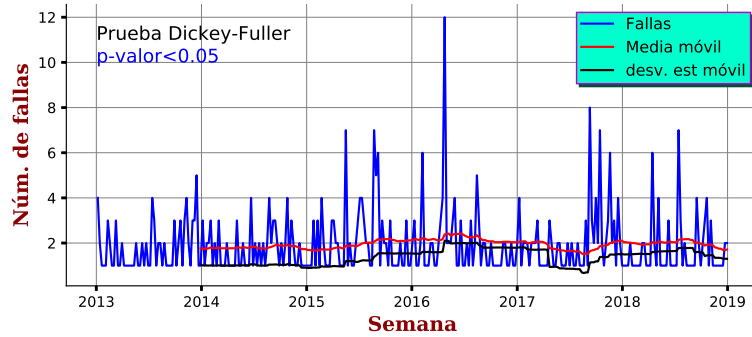


Figura 5.17: Estacionariedad de fallas en la Zona Metropolitana Norte

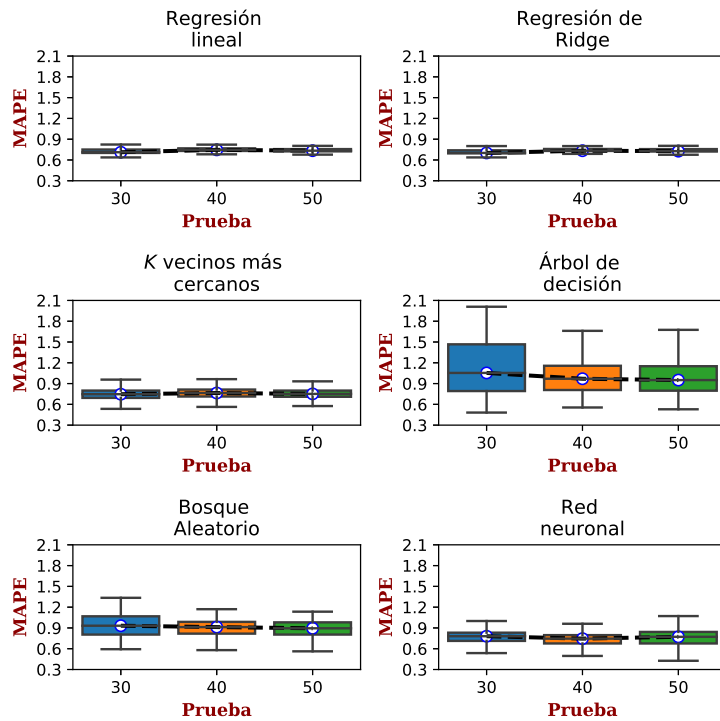


Figura 5.18: MAPE por longitud de prueba para cada familia en la Zona Metropolitana Norte

En la figura 5.18 se observa el MAPE obtenido de acuerdo a la variación en

la longitud en el conjunto de prueba para cada familia. En esta figura se observa que el valor que presenta menor mediana en el MAPE es $n_i = 30 \forall i \in \{1, 2\}$, $n_i = 50 \forall i \in \{3, 4, 5\}$ y $n_6 = 40$, de manera general se contempla $n = 50$.

En la figura 5.19 se muestra la serie de tiempo de fallas en la Zona Metropolitana Norte dividida en el conjunto de entrenamiento y de prueba. El conjunto de prueba está conformado por $n = 50$ observaciones.

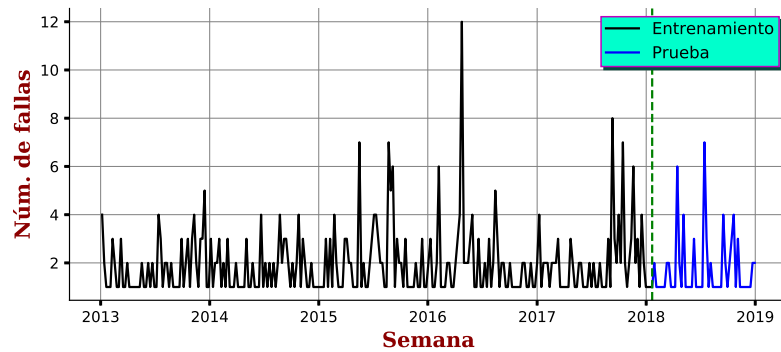


Figura 5.19: Conjuntos de entrenamiento y prueba de fallas en la Zona Metropolitana Norte

En la figura 5.20 se muestra el MAPE obtenido variando las ventanas de pronóstico para cada familia dada la longitud del conjunto de prueba de $n = 50$. Se observa que la familia que tiene menor MAPE posible para sus ventanas de pronóstico es K vecinos más cercanos y con $V = (1, 1)$, por lo que se elige dicha ventana de pronóstico para obtener mayor ventaja de este menor MAPE.

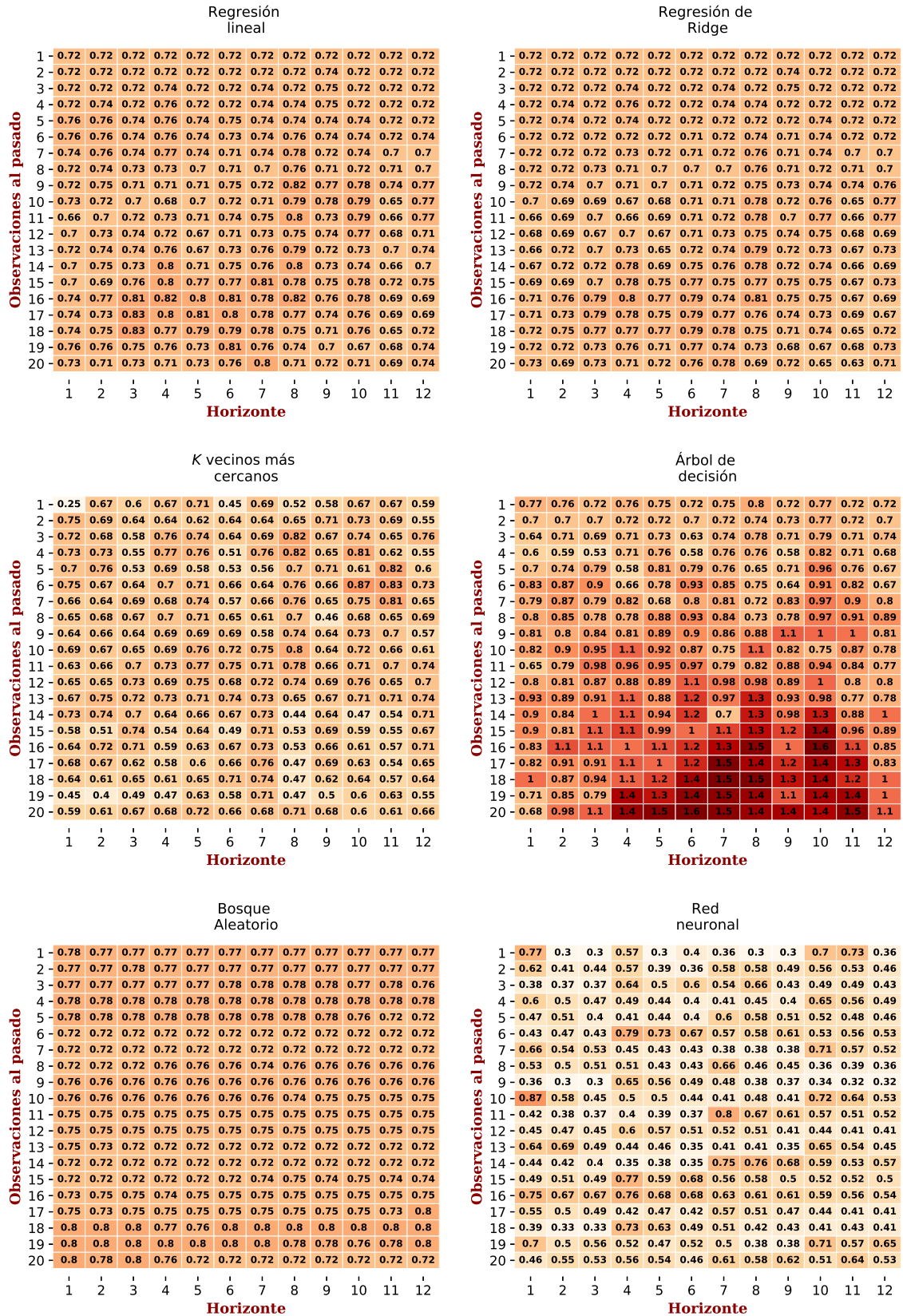


Figura 5.20: MAPE mínimo por combinación de ventana de pronóstico para cada familia en la Zona Metropolitana Norte (N=50)

En la figura 5.21 se muestra el valor del MAPE obtenido variando la configuración de hiperparámetro(s) para cada familia de modelos en el pronóstico de fallas en la Zona Metropolitana Norte. De esta manera, se obtiene aquel modelo que representa cada familia dada la longitud del conjunto de prueba y la ventana de pronóstico.

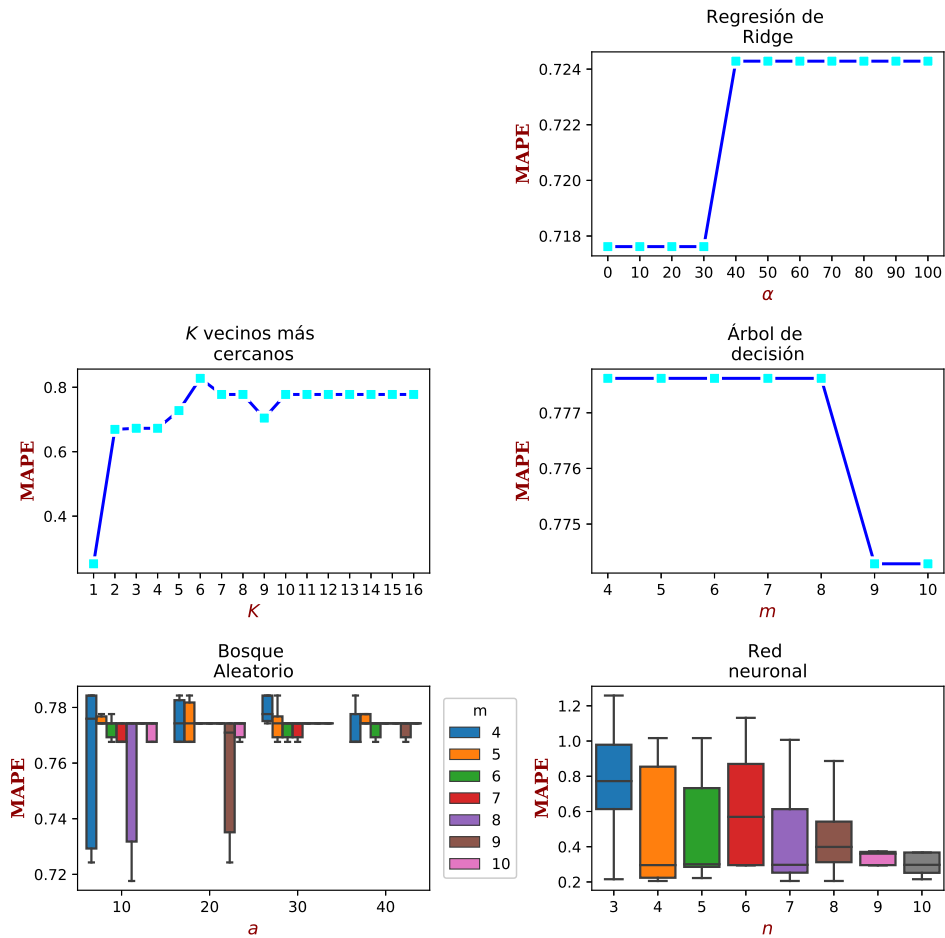


Figura 5.21: MAPE por hiperparámetro(s) de cada familia en la Zona Metropolitana Norte ($N=50$, $V=(1,1)$)

En la figura 5.21 se muestra el valor del MAPE obtenido variando la configuración de hiperparámetro(s) para cada familia de modelos en el pronóstico de fallas en la Zona Metropolitana Norte. De esta manera, se obtiene aquel modelo que representa cada familia dada la longitud del conjunto de prueba y la ventana de

pronóstico.

Ya que se han elegido los hiperparámetros correspondientes de las familias de modelos para la configuración de la longitud del conjunto de prueba $n = 50$ y la ventana de pronóstico $V = (1, 1)$, en la figura 5.22 se compara el MAPE obtenido por cada uno para seleccionar aquel modelo con menor valor.

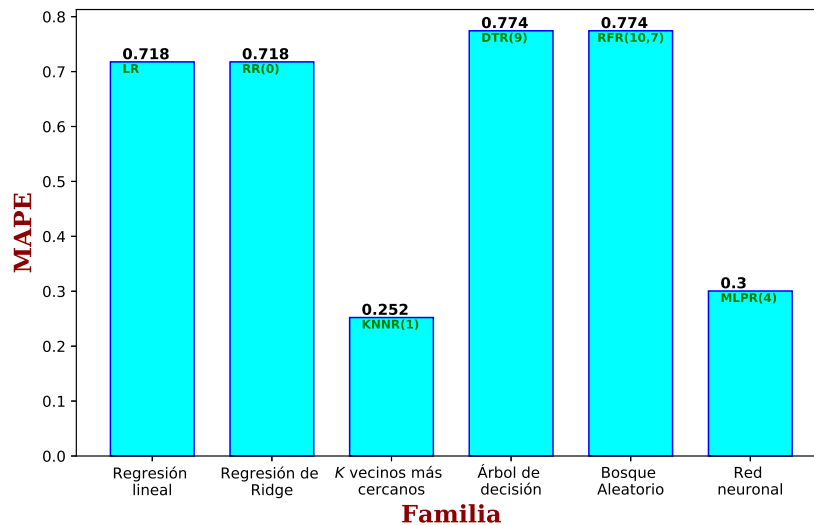


Figura 5.22: Comparación de MAPE de los modelos representativos de cada familia en la Zona Metropolitana Norte ($N=50$, $V=(1,1)$)

En la figura 5.23 se muestran las observaciones del conjunto de prueba de la serie de tiempo de fallas en la Zona Metropolitana Norte y el pronóstico brindado por el modelo con menor MAPE para los parámetros seleccionados.

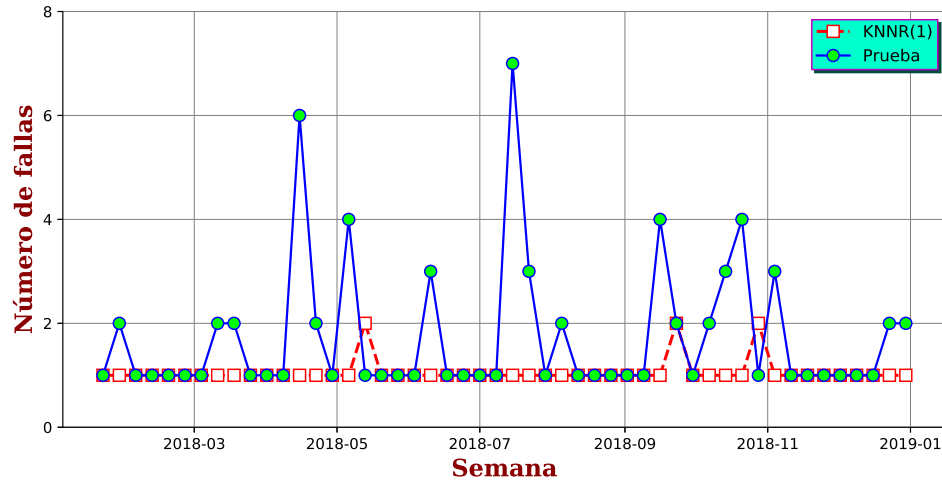


Figura 5.23: Pronóstico de fallas en la Zona Metropolitana Norte (N=50, V=(1,1), M=KNNR(1))

En la figura 5.24 se observa la frecuencia relativa porcentual del error porcentual absoluto obtenido por el pronóstico con el respectivo modelo y las características previamente seleccionadas, en la que se resalta que el grupo con mayor frecuencia relativa porcentual es el de las observaciones pronosticadas con 0% de error absoluto porcentual, además se resalta que no hay observaciones con un error absoluto porcentual mayores al 100%.

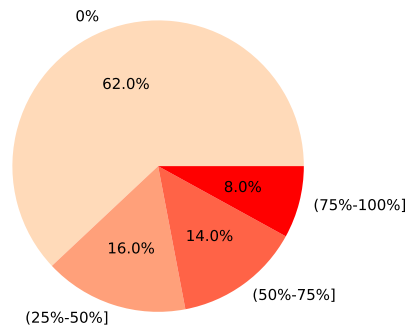


Figura 5.24: Distribución del error porcentual absoluto por intervalos en la Zona Metropolitana Norte (N=50, V=(1,1), M=KNNR(1))

En la tabla 5.4 se muestran los valores de distintas métricas, las cuales representan el desempeño del pronóstico obtenido para la serie de tiempo de las fallas en la Zona Metropolitana Poniente con $n = 50$, $V = (1, 1)$ y empleando el modelo de regresión de los k vecinos más cercanos con hiperparámetro $k = 1$.

Tabla 5.4: Desempeño del pronóstico en Montemorelos

MAPE	RMSE	MAE	R-cuad
0.2552	1.8800	0.7600	-0.0979

5.2.3 ZONA METROPOLITANA ORIENTE

En la figura 5.25 se observa la serie de tiempo para el número de fallas en la Zona Metropolitana Oriente, resaltando que el año con la observación con mayor número de fallas es el 2016 con 16, posteriormente el año 2014 y 2017 con 10 aunque en el 2014 se observa que la mayoría de las observaciones son menores que 6, mientras que el 2017 supera en varias observaciones las 6 fallas. Por otra parte, todas las observaciones del 2013 están por debajo de las 6 fallas.

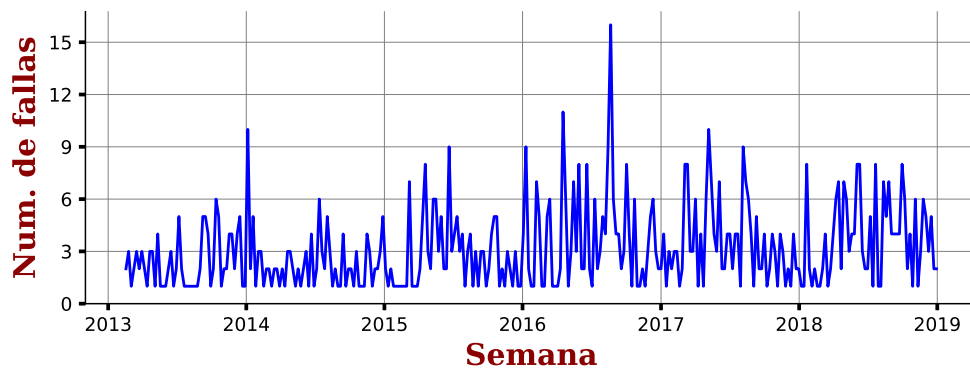


Figura 5.25: Fallas en la Zona Metropolitana Oriente

En la figura 5.26 se muestra la serie de tiempo en la Zona Metropolitana Oriente particionada por año, en la que se observa que en cada año, el máximo es alcanzado

antes de la 41ava semana, por otra parte, para el valor mínimo de fallas no hay un período fijo.

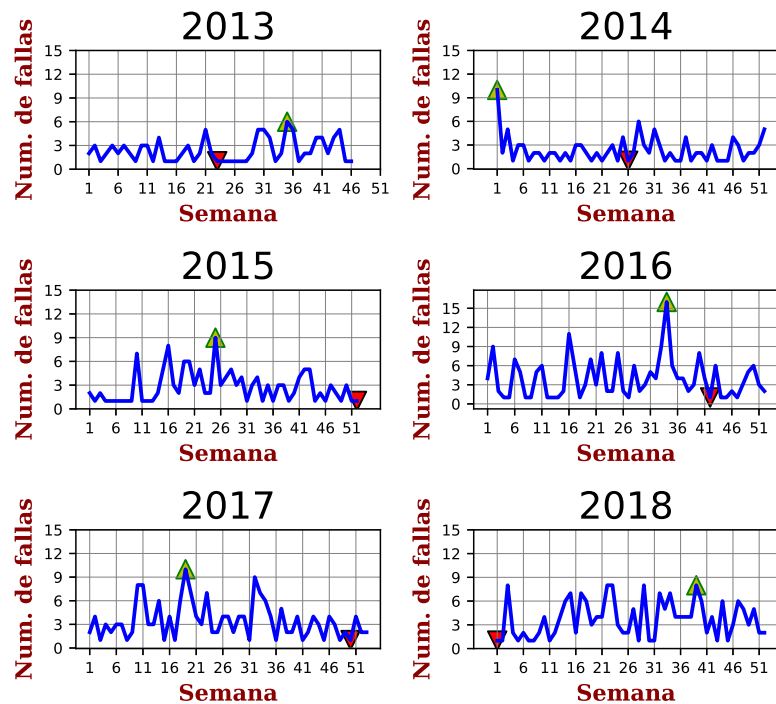


Figura 5.26: Fallas en la Zona Metropolitana Oriente por año

En la figura 5.27 se observa la descomposición aditiva de la serie de tiempo para las fallas en la Zona Metropolitana Oriente, en la que se observa que hay una tendencia creciente entre la mitad (aproximadamente) del 2014 y al 2016, mientras que del 2017 en adelante se observa un decremento. Es importante destacar que estos cambios en la tendencia son en valores relativamente pequeños, en la misma figura se observa que basándose en la estacionalidad, entre febrero y abril se obtiene el máximo número de fallas.

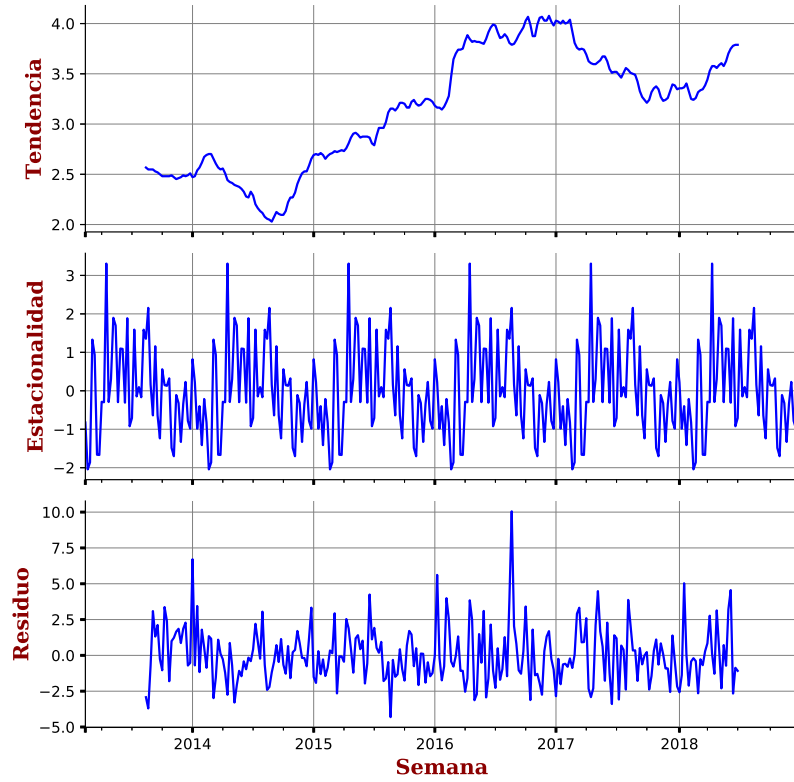


Figura 5.27: Descomposición aditiva de las fallas en la Zona Metropolitana Oriente

En la figura 5.28 se observa que dadas las fallas de una semana específica en la Zona Metropolitana Oriente, en general, las semanas con la cantidad de fallas similar a dicha semana son la primera, quinta, novena y la 52ava anteriores a esta.



Figura 5.28: Autocorrelación de fallas en la Zona Metropolitana Oriente

En la figura 5.29 se observa que el p-valor obtenido al realizar la prueba es menor a $\alpha = 0,05$ por lo que se rechaza la hipótesis nula de que la serie de tiempo tiene raíz unitaria, en favor de que esta serie de tiempo sea estacionaria.

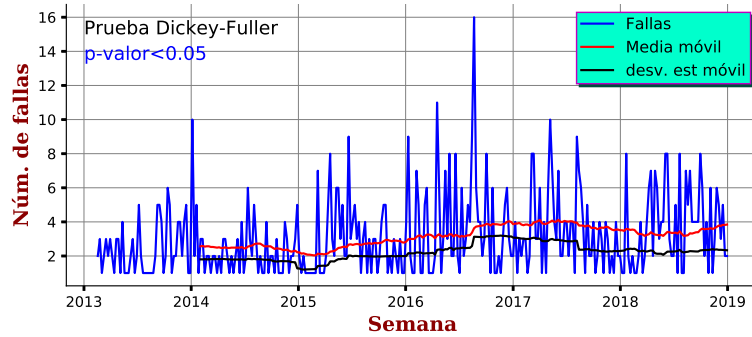


Figura 5.29: Estacionariedad de fallas en la Zona Metropolitana Oriente

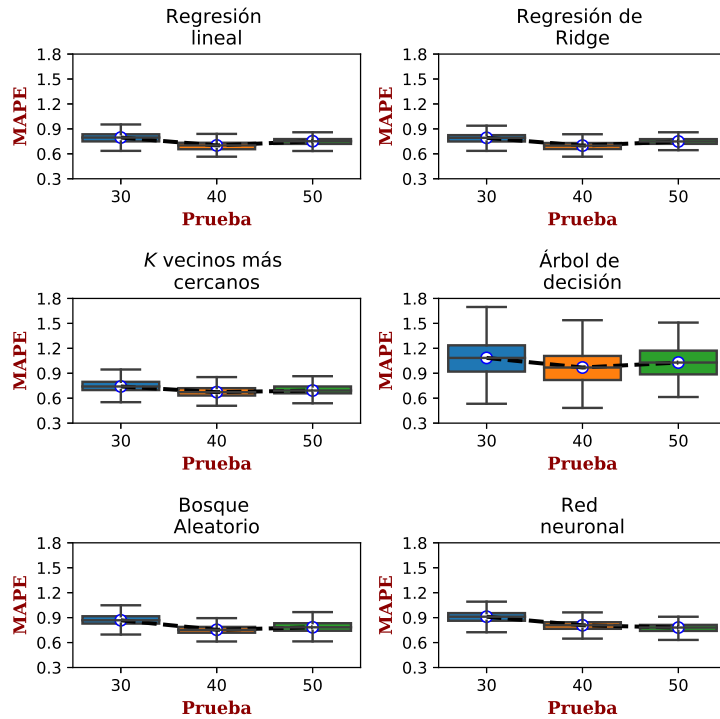


Figura 5.30: MAPE por longitud de prueba para cada familia en la Zona Metropolitana Oriente

En la figura 5.30 se observa el MAPE obtenido de acuerdo a la variación en la

longitud en el conjunto de prueba para cada familia. En esta figura se observa que el valor que presenta menor mediana en el MAPE es $n_i = 40 \forall i \in \{1, 2, 3, 4, 5\}$ y $n_6 = 50$, por lo que de manera general se contempla $n = 40$.

En la figura 5.31 se muestra la serie de tiempo de fallas en la Zona Metropolitana Oriente dividida en el conjunto de entrenamiento y de prueba. El conjunto de prueba está conformado por $n = 40$ observaciones.

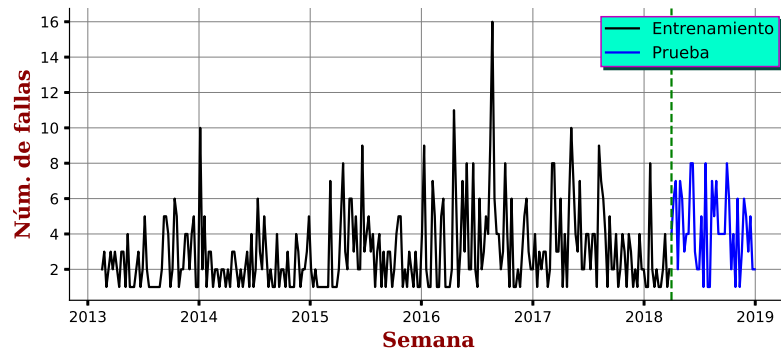


Figura 5.31: Conjuntos de entrenamiento y prueba de fallas en la Zona Metropolitana Oriente

En la figura 5.32 se muestra el MAPE obtenido variando las ventanas de pronóstico para cada familia dada la longitud del conjunto de prueba de $n = 40$. Se observa que la familia que tiene menor MAPE posible para sus ventanas de pronóstico es K vecinos más cercanos y con $V = (14, 8)$, por lo que se elige dicha ventana de pronóstico para obtener mayor ventaja de este menor MAPE.

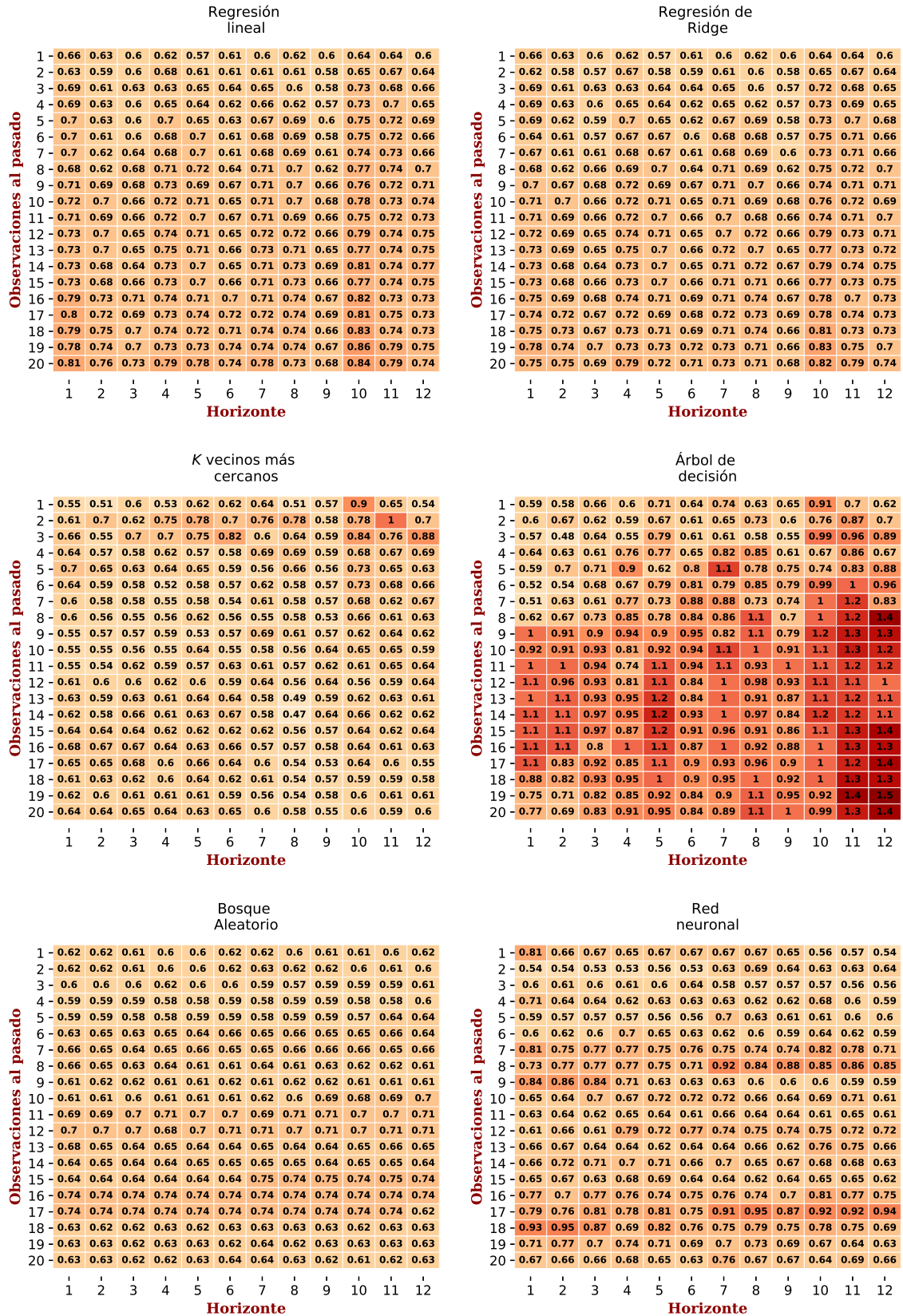


Figura 5.32: MAPE mínimo por combinación de ventana de pronóstico para cada familia en la Zona Metropolitana Oriente (N=40)

En la figura 5.33 se muestra el valor del MAPE obtenido variando la configuración de hiperparámetro(s) para cada familia de modelos en el pronóstico de fallas en la Zona Metropolitana Oriente. De esta manera, se obtiene aquel modelo que representa cada familia dada la longitud del conjunto de prueba y la ventana de pronóstico.

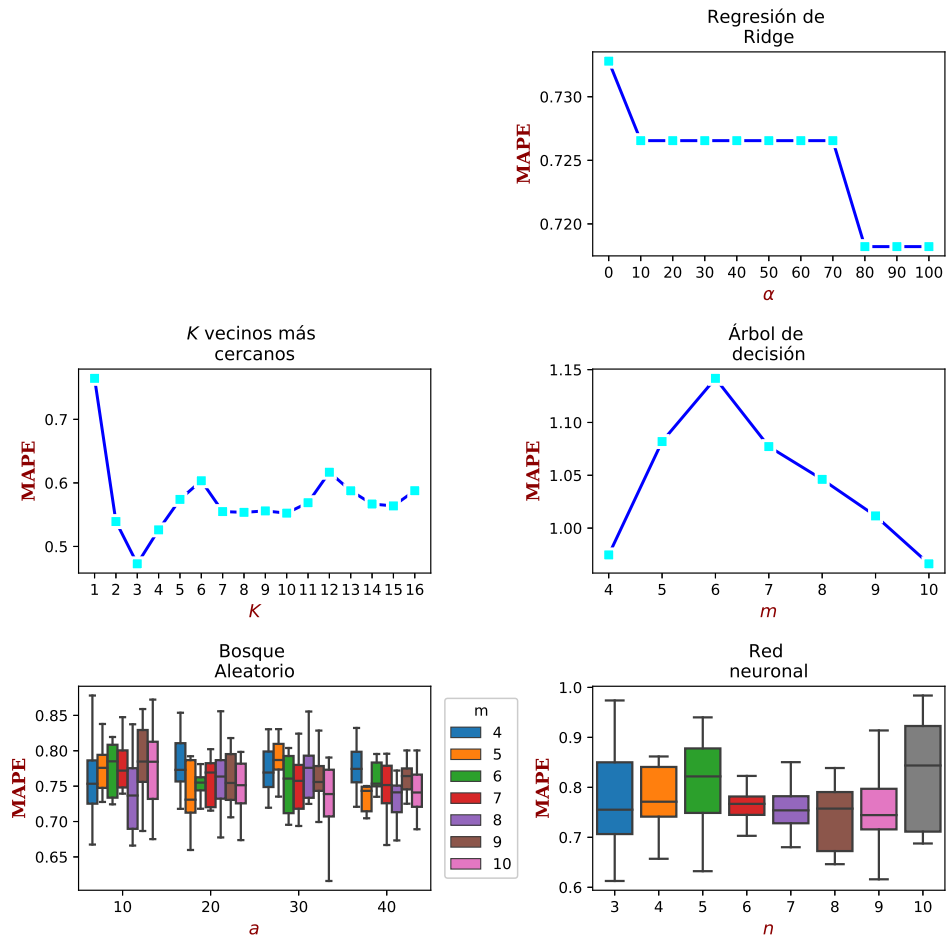


Figura 5.33: MAPE por hiperparámetro(s) de cada familia en la Zona Metropolitana Oriente ($N=40$, $V=(14,8)$)

Ya que se han elegido los hiperparámetros correspondientes de las familias de modelos para la configuración de la longitud del conjunto de prueba $n = 40$ y la ventana de pronóstico $V = (14, 8)$, en la figura 5.34 se compara el MAPE obtenido por cada uno para seleccionar aquel modelo con menor valor.

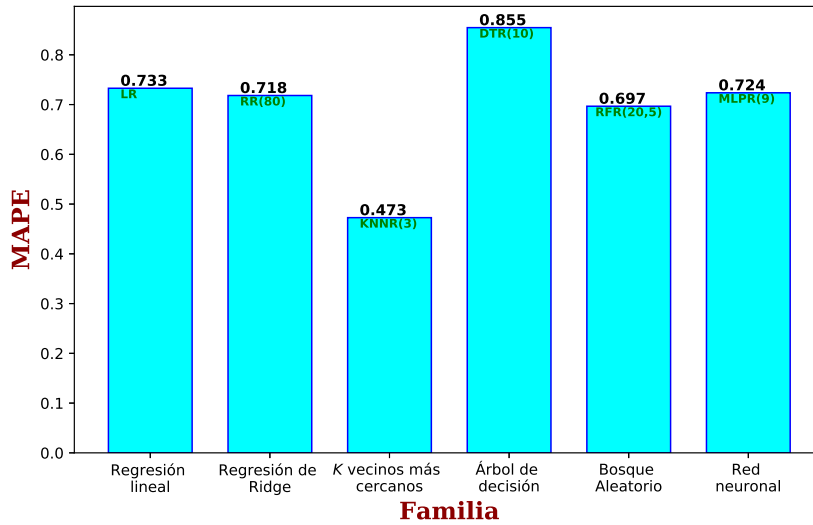


Figura 5.34: Comparación de MAPE de los modelos representativos de cada familia en la Zona Metropolitana Oriente (N=40, V=(14,8))

En la figura 5.35 se muestran las observaciones del conjunto de prueba de la serie de tiempo de fallas en la Zona Metropolitana Oriente y el pronóstico brindado por el modelo con menor MAPE para los parámetros seleccionados.

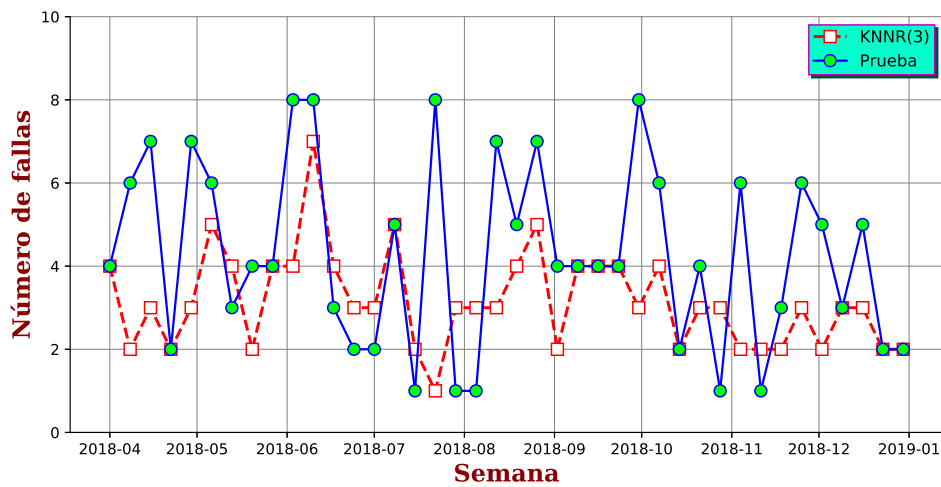


Figura 5.35: Pronóstico de fallas en la Zona Metropolitana Oriente (N=40, V=(14,8), M=KNNR(3))

En la figura 5.36 se observa la frecuencia relativa porcentual del error porcentual absoluto obtenido por el pronóstico con el respectivo modelo y las características previamente seleccionadas, en la que se resalta que el grupo con mayor frecuencia relativa porcentual es el de las observaciones pronosticadas con un error absoluto porcentual en el intervalo $(25\% - 50\%]$, seguido por las observaciones con un error absoluto porcentual de 0% .

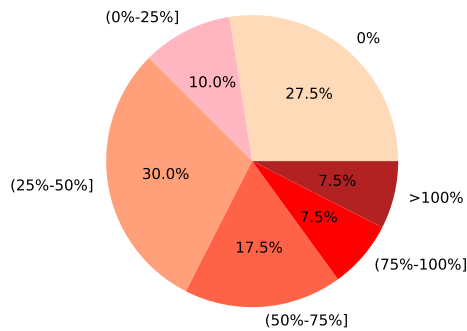


Figura 5.36: Distribución del error porcentual absoluto por intervalos en la Zona Metropolitana Oriente ($N=40$, $V=(14,8)$, $M=KNNR(3)$)

En la tabla 5.5 se muestran los valores de distintas métricas, las cuales representan el desempeño del pronóstico obtenido para la serie de tiempo de las fallas en la Zona Metropolitana Poniente con $n = 40$, $V = (14, 8)$ y empleando el modelo de regresión de los k vecinos más cercanos con hiperparámetro $k = 3$.

Tabla 5.5: Desempeño del pronóstico en la Zona Metropolitana Oriente

MAPE	RMSE	MAE	R-cuad
0.4730	1.4250	1.7250	0.7060

5.2.4 ZONA METROPOLITANA PONIENTE

En la figura 5.37 se observa la serie de tiempo para el número de fallas en la Zona Metropolitana Poniente, resaltando que el máximo valor en las observaciones es de 9 y se presenta en los años 2013 y 2014, de hecho, dichos años junto con el 2015 son los que presentan mayor cantidad de fallas, mientras que del 2016 al 2018 se tiene que las observaciones son inferiores a 6 con excepción de un caso en el 2016.

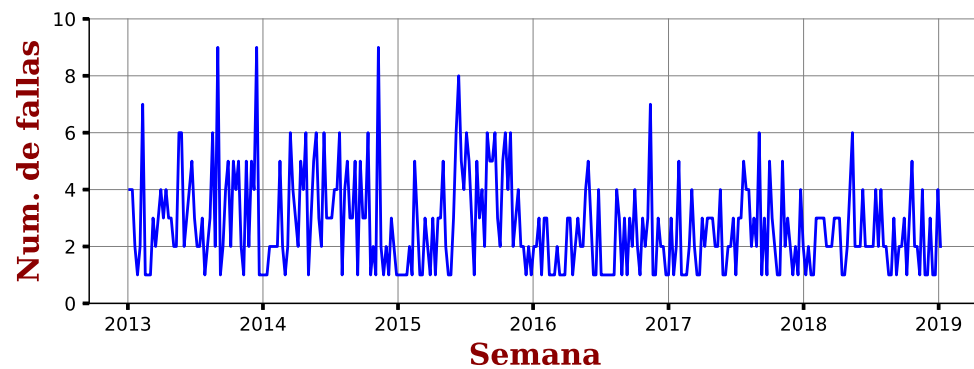


Figura 5.37: Fallas en la Zona Metropolitana Poniente

En la figura 5.38 se muestra la serie de tiempo en la Zona Metropolitana Poniente particionada por año, en la que se observa que en el 2013, 2014 y 2016 la mayor cantidad de fallas ocurre después de la 40ava semana, mientras que para el resto de los años ocurre entre la 20ava y la 40ava semana. Por otra parte, para el valor mínimo de fallas no hay un período fijo.

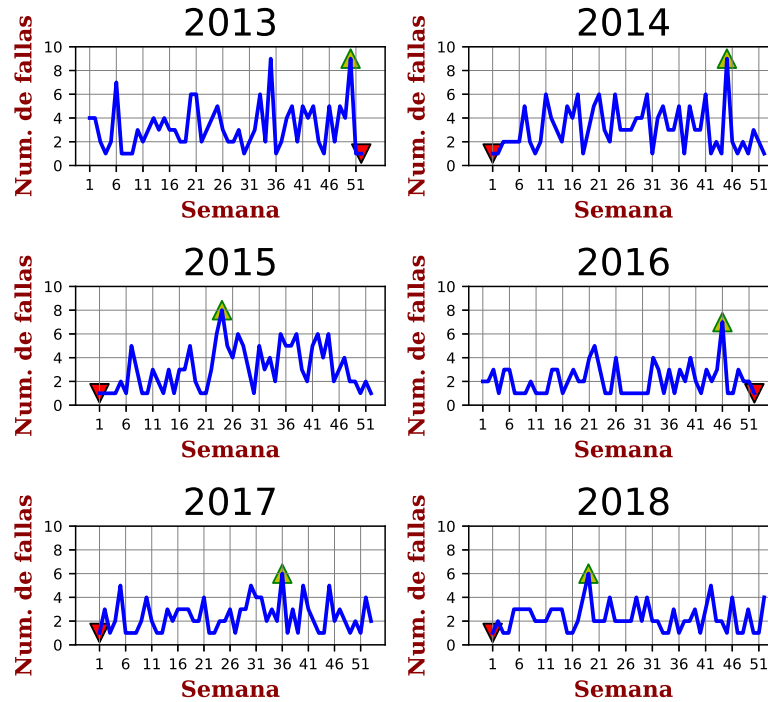


Figura 5.38: Fallas en la Zona Metropolitana Poniente por año

En la figura 5.39 se observa la descomposición aditiva de la serie de tiempo para las fallas en la Zona Metropolitana Poniente, en la que se observa que a grandes razgos, hay una tendencia decreciente con mayor notoriedad en el año 2016. Por otra parte, en la misma figura se observa que basándose en la estacionalidad, entre los meses de agosto y noviembre es cuando generalmente hay mayor presencia de fallas.

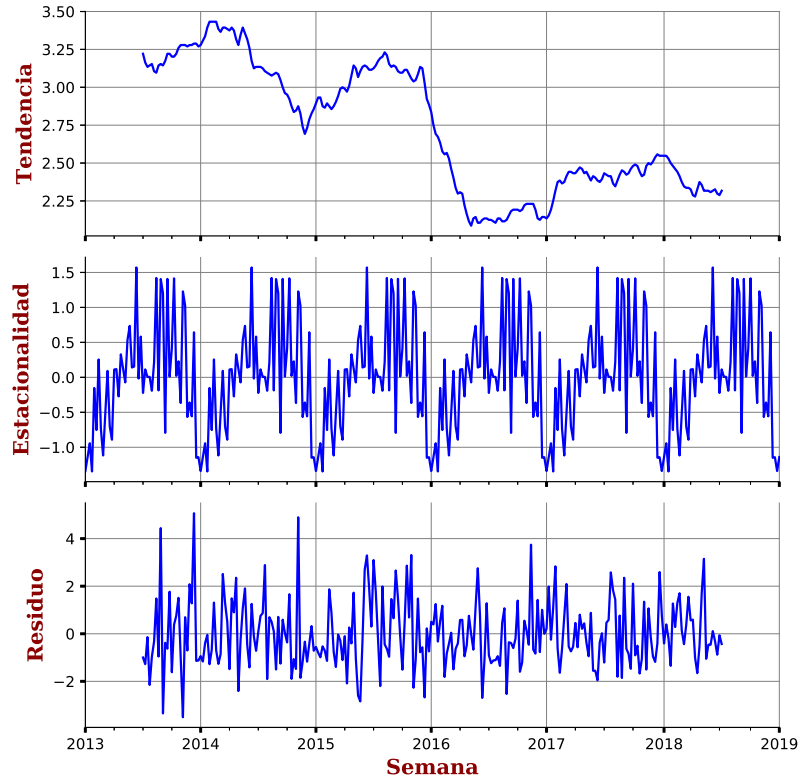


Figura 5.39: Descomposición aditiva de las fallas en la Zona Metropolitana Poniente

En la figura 5.40 se observa que dadas las fallas de una semana específica en la Zona Metropolitana Oriente, en general, las semanas con la cantidad de fallas similar a dicha semana son la cuarta, y la 47ava anteriores a esta.

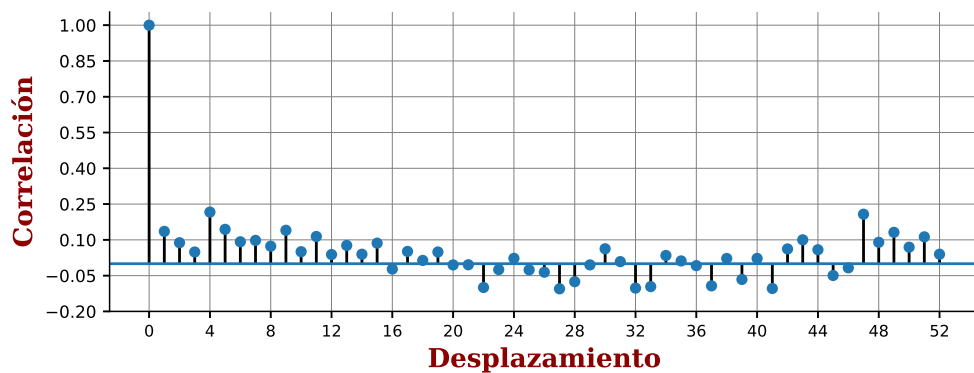


Figura 5.40: Autocorrelación de fallas en la Zona Metropolitana Poniente

En la figura 5.41 se observa que el p-valor obtenido al realizar la prueba es menor a $\alpha = 0,05$ por lo que se rechaza la hipótesis nula de que la serie de tiempo tiene raíz unitaria, en favor de que esta serie de tiempo sea estacionaria.

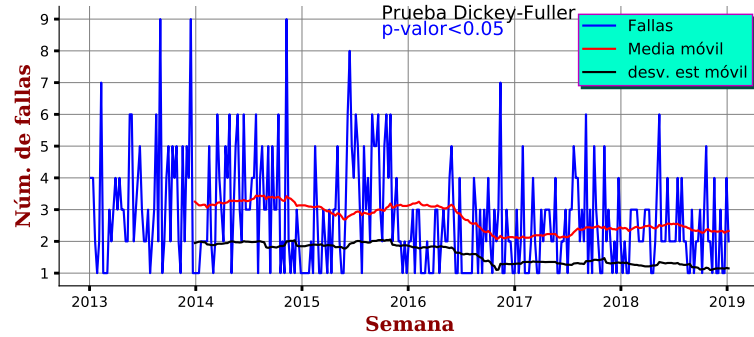


Figura 5.41: Estacionariedad de fallas en la Zona Metropolitana Poniente

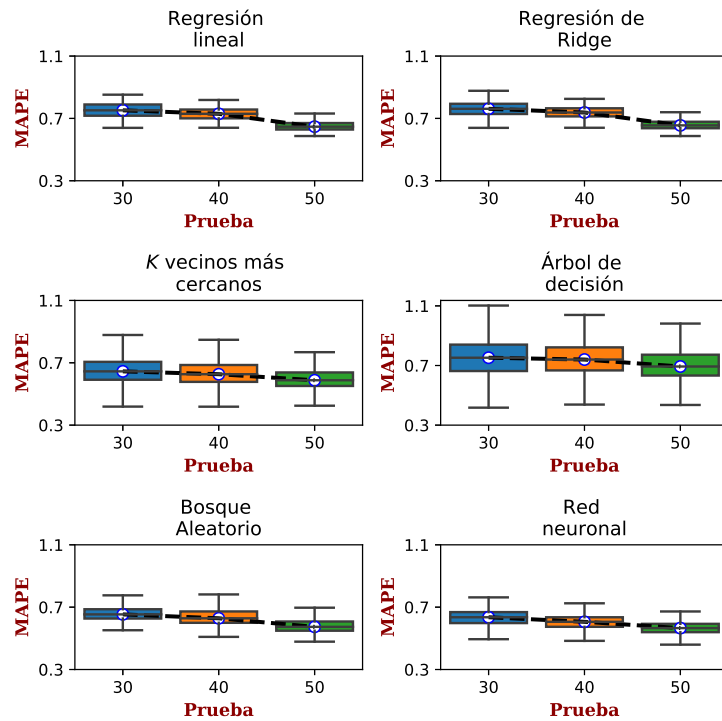


Figura 5.42: MAPE por longitud de prueba para cada familia en la Zona Metropolitana Poniente

En la figura 5.42 se observa el MAPE obtenido de acuerdo a la variación en la

longitud en el conjunto de prueba para cada familia. En esta figura se observa que el valor que presenta menor mediana en el MAPE es $n_i = 50 \forall i \in \{1, 2, 3, 4, 5, 6\}$, por lo que de manera general se contempla $n = 50$.

En la figura 5.43 se muestra la serie de tiempo de fallas en la Zona Metropolitana Poniente dividida en el conjunto de entrenamiento y de prueba. El conjunto de prueba está conformado por $n = 50$ observaciones.

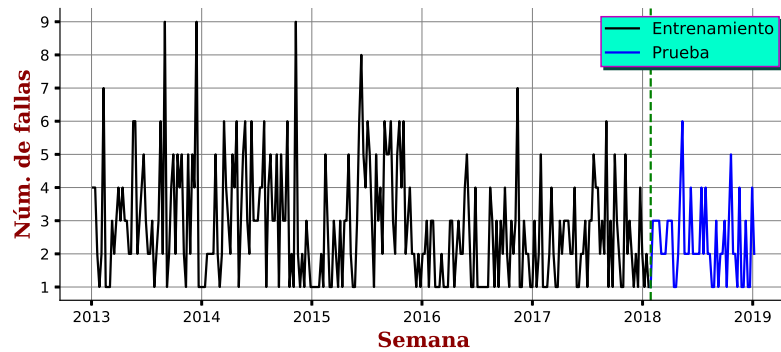


Figura 5.43: Conjuntos de entrenamiento y prueba de fallas en la Zona Metropolitana Poniente

En la figura 5.44 se muestra el MAPE obtenido variando las ventanas de pronóstico para cada familia dada la longitud del conjunto de prueba de $n = 50$. Se observa que la familia que tiene menor MAPE posible para sus ventanas de pronóstico es K vecinos más cercanos y con $V = (17, 1)$, por lo que se elige dicha ventana de pronóstico para obtener mayor ventaja de este menor MAPE.

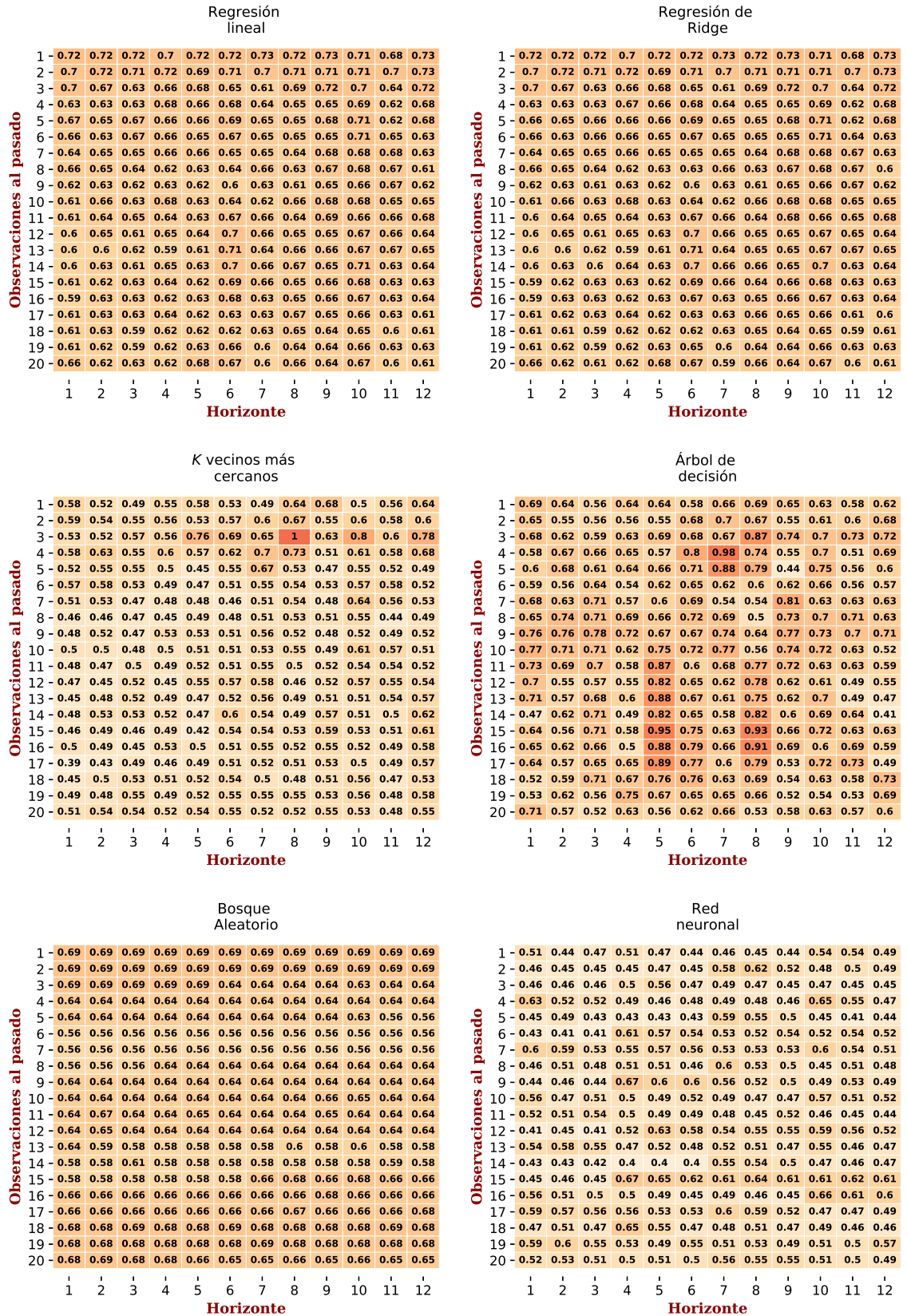


Figura 5.44: MAPE mínimo por combinación de ventana de pronóstico para cada familia en la Zona Metropolitana Poniente (N=50)

En la figura 5.45 se muestra el valor del MAPE obtenido variando la configuración de hiperparámetro(s) para cada familia de modelos en el pronóstico de fallas en la Zona Metropolitana Poniente. De esta manera, se obtiene aquel modelo que representa cada familia dada la longitud del conjunto de prueba y la ventana de pronóstico.

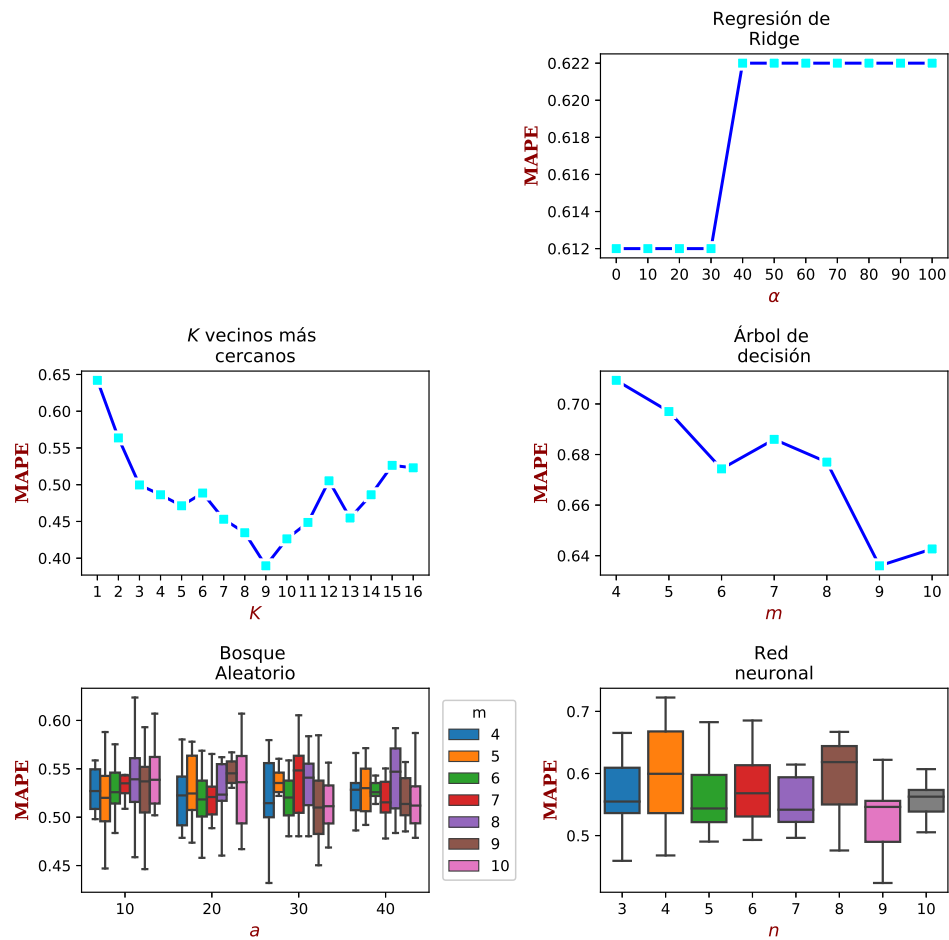


Figura 5.45: MAPE por hiperparámetro(s) de cada familia en la Zona Metropolitana Poniente ($N=50$, $V=(17,1)$)

Ya que se han elegido los hiperparámetros correspondientes de las familias de modelos para la configuración de la longitud del conjunto de prueba $n = 50$ y la ventana de pronóstico $V = (17, 1)$, en la figura 5.46 se compara el MAPE obtenido por cada uno para seleccionar aquel modelo con menor valor.

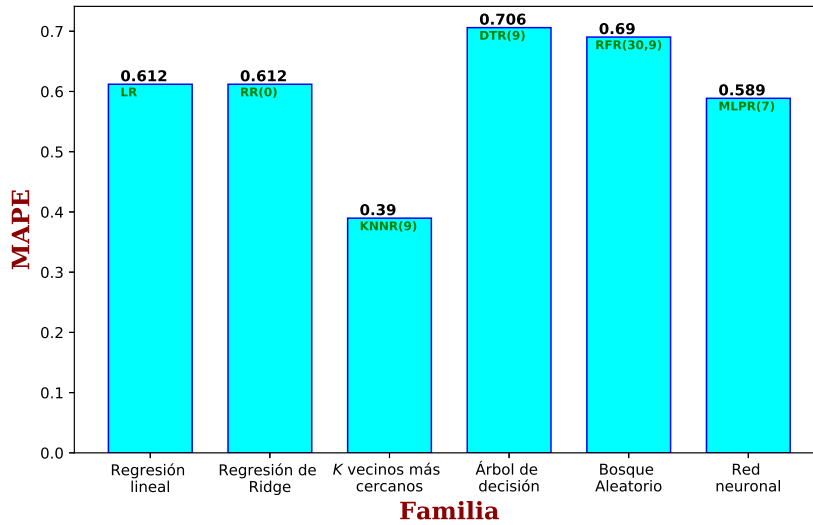


Figura 5.46: Comparación de MAPE de los modelos representativos de cada familia en la Zona Metropolitana Poniente (N=50, V=(17,1))

En la figura 5.47 se muestran las observaciones del conjunto de prueba de la serie de tiempo de fallas en la Zona Metropolitana Poniente y el pronóstico brindado por el modelo con menor MAPE para los parámetros seleccionados.

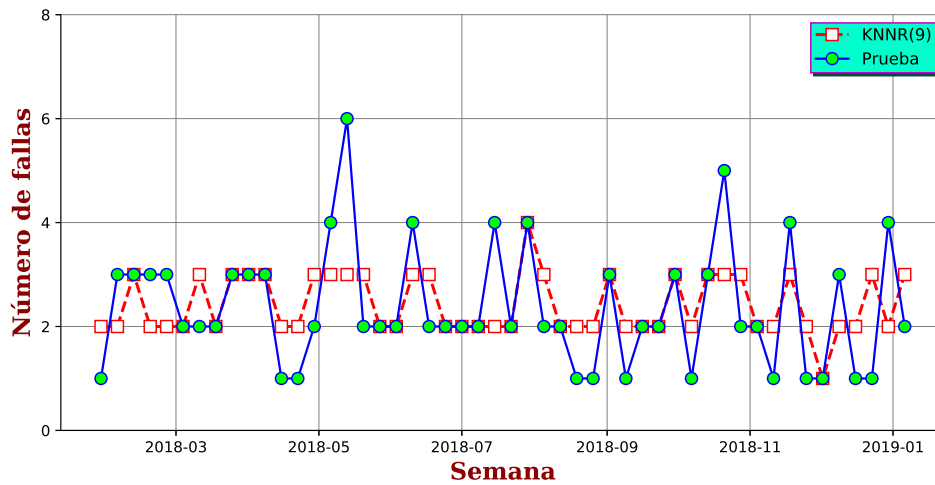


Figura 5.47: Pronóstico de fallas en la Zona Metropolitana Poniente (N=50, V=(17,1), M=KNNR(9))

En la figura 5.48 se observa la frecuencia relativa porcentual del error porcentual absoluto obtenido por el pronóstico con el respectivo modelo y las características previamente seleccionadas, en la que se resalta que el grupo con mayor frecuencia relativa porcentual es el de las observaciones pronosticadas con 0% de error absoluto porcentual.

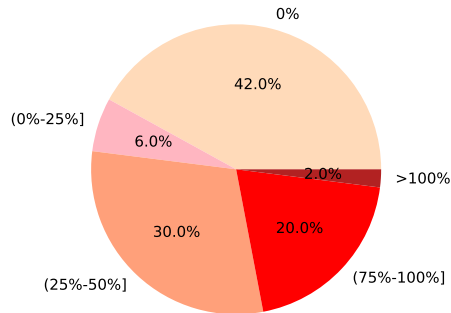


Figura 5.48: Distribución del error porcentual absoluto por intervalos en la Zona Metropolitana Poniente ($N=50$, $V=(17,1)$, $M=KNNR(9)$)

En la tabla 5.6 se muestran los valores de distintas métricas, las cuales representan el desempeño del pronóstico obtenido para la serie de tiempo de las fallas en la Zona Metropolitana Poniente con $n = 50$, $V = (17, 1)$ y empleando el modelo de regresión de los k vecinos más cercanos con hiperparámetro $k = 9$.

Tabla 5.6: Desempeño del pronóstico en la Zona Metropolitana Poniente

MAPE	RMSE	MAE	R-cuad
0.3900	0.1800	0.7000	0.8626

CAPÍTULO 6

CONCLUSIONES

Por último, en este capítulo se discuten las conclusiones obtenidas por los resultados que proporciona la experimentación con las características anteriormente descritas. En la sección 6.1 se muestran las conclusiones generales de este trabajo, en la sección 6.2 se presentan las contribuciones que este aporta y finalmente en la sección 6.3 se hace mención del trabajo a futuro.

6.1 CONCLUSIONES

En este trabajo se determinó un modelo matemático de pronóstico para cada una de las zonas que comprenden la división de distribución Golfo Norte. Además, se utilizó la estrategia directa de pronóstico y la métrica MAPE para establecer las configuraciones adecuadas para determinar dicho modelo, entre ellas la longitud del conjunto de prueba y la ventana de pronóstico. En la tabla 6.2 se muestran los valores adecuados obtenidos para cada una de estas variables.

Tabla 6.1: Parámetros obtenidos para cada zona

Zona	Longitud de prueba	Ventana de pronóstico	Modelo
Montemorelos	30	(1,3)	KNNR(10)
Norte	50	(1,1)	KNNR(1)
Oriente	40	(14,8)	KNNR(3)
Poniente	50	(17,1)	KNNR(9)

Tabla 6.2: Desempeño del pronóstico para cada zona

Zona	MAPE	RMSE	MAE	R-cuad
Montemorelos	0.3640	3.1300	2.6666	0.7876
Norte	0.2552	1.8800	0.7600	-0.0979
Oriente	0.4730	1.4250	1.7250	0.7060
Poniente	0.3900	0.1800	0.7000	0.8626

Así mismo se calculó que el 70 %, 78 %, 67.5 % y 78 % de las observaciones pronosticadas tienen un error porcentual absoluto menor o igual que 50 % para las fallas de Montemorelos, Zona Metropolitana Norte, Zona Metropolitana Oriente y Zona Metropolitana Poniente, respectivamente. Además, de acuerdo a la métrica MAE se obtienen valores cercanos a cero para la Zona Metropolitana Norte y Zona Metropolitana Poniente, mientras que en Montemorelos se obtiene el mayor valor del MAE, lo cual es aceptable ya que es la zona en que hay mayor cantidad de fallas y dicho valor es relativamente bajo, la zona con menor precisión con esta métrica fue la Zona Metropolitana Oriente, debido a que obtiene un valor cercano a dos y el máximo valor de fallas es bajo. Se obtienen conclusiones similares en el caso de la métrica RMSE. Por otra parte, para el R-cuad la Zona Metropolitana Norte obtiene un valor cercano a cero, ya que las predicciones son aproximadamente iguales entre sí, mientras que para el resto de las zonas se obtiene un valor mayor a 0.70.

De esta manera, se puede concluir que los modelos tienen un desempeño aceptable, confirmando la hipótesis de investigación.

6.2 CONTRIBUCIONES

La principal contribución de este trabajo es que ahora se tiene un modelo matemático que permite pronosticar el número de fallas que habrá en un futuro a corto plazo en cada una de las zonas de la división de distribución Golfo Norte. Además, cada modelo fue seleccionado comparándolo metodológicamente con otros, también se evaluó el desempeño para dicho modelo, por lo que estadísticamente, y en base a los resultados obtenidos, se puede decir que los modelos de pronóstico son aceptables.

6.3 TRABAJO A FUTURO

Existen diversas áreas de interés que se generan después de los resultados obtenidos en este trabajo, entre ellas:

1. **Utilizar los resultados obtenidos como instancias a un problema de planeación para la realización de las podas que se realizan de manera frecuente.**

En [51] se menciona que dentro de las aplicaciones de planificación industrial, se lleva a cabo inicialmente un proceso de pronóstico para así, obtener una estimación de las demandas futuras esperadas, y dado que este trabajo permite estimar el número de fallas al futuro para las zonas que comprende la división de distribución Golfo Norte es posible utilizar dichas estimaciones como demandas.

El problema de programación matemática reportado en la literatura que presenta mayor similitud con este problema es el *Single-Item Capacitated Lot Sizing Problem (SICLSP)*. La formulación es como sigue

Descripción

Dados los conjuntos N y S de zonas y semanas, respectivamente, se desea planear la producción e inventario de podas conociendo el pronóstico de fallas $D_{i,j}$ que habrá en la zona i durante la semana j , donde dicho pronóstico tiene un error de desviación d_i para la zona i . También se tiene un límite de podas Cap_j a realizar en la semana j , el cual no puede ser superado y finalmente, se cuenta con costos de producción $C_{1,i,j}$ e inventario $C_{2,i,j}$ de realizar y almacenar una poda en la zona i durante la semana j .

El objetivo es determinar las unidades a producir y almacenar en el inventario para minimizar los costos de producción e inventario satisfaciendo restricciones de capacidad y de demanda.

Es importante observar que en este problema el término producción se refiere a llevar a cabo una poda, mientras que el término inventario está relacionado con pagar el servicio de producción para que se lleve a cabo la poda pero no en la semana que se produjo.

Parámetros

Z : Conjunto de zonas.

S : Conjunto de semanas.

$D_{i,j} \in \mathbb{N} \cup \{0\}$: Pronóstico de fallas en la zona i durante la semana j .

$d_i \in \mathbb{R}_+ \cup \{0\}$: Error de desviación del pronóstico en la zona i .

$Cap_j \in \mathbb{N} \cup \{0\}$: Límite de podas a realizar en la semana j .

$C_{1,i,j} \in \mathbb{R}_+$: Costo de realizar una poda en la zona i durante la semana j .

$C_{2,i,j} \in \mathbb{R}_+$: Costo de almacenar una poda durante la semana j para realizarse en la zona i .

$I_{i,0} \in \mathbb{N} \cup \{0\}$: Podas almacenadas inicialmente en la zona i .

Variables de decisión

$P_{i,j} \in \mathbb{N} \cup \{0\}$: Número de podas a producir en la zona i durante la semana j .

$I_{i,j} \in \mathbb{N} \cup \{0\}$: Número de podas a almacenar en la zona i durante la semana j .

Función objetivo

$$\min \sum_{i \in Z} \sum_{j \in S} [C_{1,i,j} P_{i,j} + C_{2,i,j} I_{i,j}] \quad (6.1)$$

Restricciones

$$I_{i,j} = \sum_{k=1}^j [P_{i,k} - D_{i,k} - \alpha d_i] + I_{i,0} \quad \forall i \in Z, j \in S \quad (6.2)$$

$$\sum_{i \in Z} P_{i,j} \leq Cap_j \quad \forall j \in S \quad (6.3)$$

Naturaleza de las variables de decisión

$$P_{i,j} \in \mathbb{N} \cup \{0\} \quad \forall i \in Z, j \in S \quad (6.4)$$

$$I_{i,j} \in \mathbb{N} \cup \{0\} \quad \forall i \in Z, j \in S \quad (6.5)$$

donde la expresión 6.1 hace referencia a minimizar el costo total durante el horizonte de planeación, mientras que la 6.2 indica la relación entre las unidades producidas y aquellas que se almacenan en donde $\alpha \in [0, 1]$ representa el porcentaje de cobertura en el peor caso del pronóstico de las fallas, por su parte la 6.3 señala que las podas realizadas semanalmente deben respetar el límite de podas establecido. Finalmente las expresiones 6.4 y 6.5 muestran la naturaleza de las variables las cuales son enteras y representan la producción y el inventario, respectivamente.

2. Implementar otros modelos matemáticos de pronóstico para verificar su desempeño en comparación con los utilizados en este trabajo.

Es posible aplicar modelos de otras familias de regresión que no fueron contempladas en este estudio, por ejemplo modelos de máquinas de soporte vectorial, aprendizaje profundo, entre otros.

3. Utilizar datos de variables climatológicas y ambientales para aumentar la precisión y el horizonte de pronóstico.

Como se explica en la sección 2 la precisión y el horizonte de pronóstico pueden aumentar si se tiene mayor cantidad de variables.

El Sistema Integral de Monitoreo Ambiental (SIMA) es una organización que se encarga de medir los niveles de distintas variables climatológicas y contaminantes en el Área Metropolitana de Monterrey (AMM). Para poder brindar su servicio cuenta con 13 estaciones de monitoreo en el AMM. Entre las principales variables que se monitorean se encuentran: lluvia, viento (dirección y velocidad), radiación solar y algunos contaminantes como CO_2 , $PM_{2,5}$, etc.

Dado que las zonas que conforman la división de distribución Golfo Norte se encuentran dentro o cerca del AMM, se pretende considerar variables que puedan aportar información para mejorar el pronóstico tal como la velocidad del viento y la lluvia desde el punto de vista climatológico.

En la figura 6.1 se muestra el caso en el que se contemplan las fallas en la zona de Montemorelos y la velocidad del viento en la estación Suroeste en el año 2013.

Dado que la velocidad del viento es monitoreada cada hora de cada día, para poder considerar frecuencias del mismo paso que las fallas (semanal) se considera el promedio de la velocidad del viento para todas las observaciones por semana.

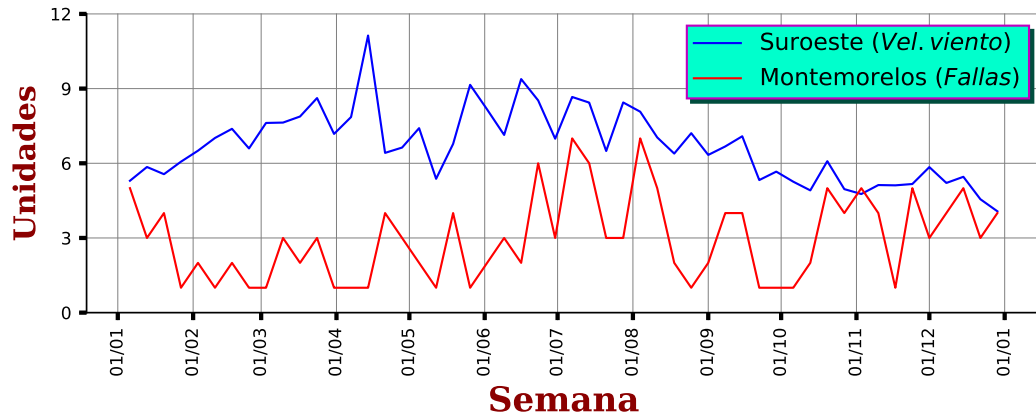


Figura 6.1: Serie de tiempo de las fallas en Montemorelos y la velocidad del viento en la estación Suroeste en el 2013

De manera general, se tiene el problema de pronosticar la serie de tiempo de fallas por semana para alguna zona de de las cuatro que comprenden la división de distribución Golfo Norte utilizando dicha serie de tiempo y también la serie de tiempo de la velocidad del viento, es decir se trata del pronóstico de una serie de tiempo con múltiples series de tiempo como entradas.

BIBLIOGRAFÍA

- [1] CFE (2020). Comisión federal de electricidad. <https://www.cfe.mx/Pages/Index.aspx>.
- [2] Cantilever (2017). Divisiones de distribución cfe. <https://www.cantilever.com.mx/>.
- [3] Comisión Federal de Electricidad (1992). *Terminología de CFE*.
- [4] Astudillo M. (2012). *Fundamentos de economía*. Probooks, 1era edición, ISBN 978-607-02-2974-9.
- [5] Pacheco J. (2019). Estructura económica (estructura económica mundial). <https://www.webyempresas.com/estructura-economica/>.
- [6] Sy Corvo H. (2020). Estructura económica: características, elementos y ejemplos. <https://www.lifeder.com/estructura-economica/>.
- [7] Pérez M., Orlandoni G., and Ramoni J. (2013). Aplicación de la metodología de series de tiempo en la estimación de los niveles de exportaciones de café de colombia periodo 1958-2011. *Revista Innovaciencia*, 1(1):11–16.
- [8] Contreras A., Atziray C., Martínez José, and Sánchez D. (2016). Análisis de series de tiempo en el pronóstico de la demanda de almacenamiento de productos perecederos. *Estudios Gerenciales*, 32(141):387–396.
- [9] Gallego-Nicasio Moraleda J., Rodríguez A., Mínguez J., and Jiménez F. (2018).

- Modelos arima para la predicción del gasto conjunto de oxígeno de vuelo y otros gases en el ejército del aire. *Sanidad Militar*, 74(4):223–229.
- [10] Spyros M., Evangelos S., and Vassilios A.(2018). Statistical and machine learning forecasting methods: Concerns and ways forward. *PloS one*, 13(3):e0194889.
- [11] Bohdan P. (2019). Machine-learning models for sales time series forecasting. *Data*, 4(1):15.
- [12] Juárez P. (2016). El sector energético, factor clave para el desarrollo económico. <https://www.eluniversal.com.mx/articulo/cartera/finanzas/2016/01/22/el-sector-energetico-factor-clave-para-el-desarrollo-economico>.
- [13] D.Econosignal (2019). Tendencias de industrias. Deloitte. Pág. 6.
- [14] Azadeh A., Asadzadeh SM., Saberi M., Nadimi V., Tajvidi A., and Sheikalis-hahi M. (2011). A neuro-fuzzy-stochastic frontier analysis approach for long-term natural gas consumption forecasting and behavior analysis: the cases of bahrain, saudi arabia, syria, and uae. *Applied Energy*, 88(11):3850–3859.
- [15] Pappas S., Ekonomou L., Karamousantas D., Chatzarakis G., Katsikas S., and Liatsis P. (2008). Electricity demand loads modeling using autoregressive moving average (arma) models. *Energy*, 33(9):1353–1360.
- [16] Kadir K., Halim C., Harun O., and Olcay C. (2009). Modeling and prediction of turkey’s electricity consumption using artificial neural networks. *Energy Conversion and Management*, 50(11):2719–2727, 2009.
- [17] Rahmani R., Yusof R., Seyedmahmoudian M., and Mekhilef S. (2013). Hybrid technique of ant colony and particle swarm optimization for short term wind energy forecasting. *Journal of Wind Engineering and Industrial Aerodynamics*, 123:163–170.

-
- [18] Osório G., Matias J., and Catalão J. (2015). Short-term wind power forecasting using adaptive neuro-fuzzy inference system combined with evolutionary particle swarm optimization, wavelet transform and mutual information. *Renewable Energy*, 75:301–307.
- [19] Azadeh A., Asadzadeh S., Mirseraji G., and Saberi M. (2015). An emotional learning-neuro-fuzzy inference approach for optimum training and forecasting of gas consumption estimation models with cognitive data. *Technological Forecasting and Social Change*, 91:47–63.
- [20] Vázquez R. (2010). Energía eléctrica, estratégica para la economía. <https://realestatemarket.com.mx/infraestructura-y-construccion/11285-energia-electrica-estrategica-para-la-economia>.
- [21] Mishra A. (2018). Metrics to evaluate your machine learning algorithm. recuperado de: <https://towardsdatascience.com/metrics-to-evaluate-your-machine-learning-algorithm-f10ba6e38234>.
- [22] Brownle J. (2016). Metrics to evaluate machine learning algorithms in python. recuperado de: <https://machinelearningmastery.com/metrics-evaluate-machine-learning-algorithms-python/>.
- [23] Berrones M. (2019). *Clasificación de mamografías mediante redes neuronales convolucionales*. PhD thesis, Universidad Autónoma de Nuevo León.
- [24] Peralta J. (2018). Modelo de elaboración de pronóstico de ventas mediante el uso de redes neuronales artificiales y svr.
- [25] Coba L. (2017). Pronóstico de ventas de la farmacéutica sanofi usando series temporales.
- [26] Chiou-Jye H. and Ping-Huan K. (2018). A deep cnn-lstm model for particulate matter (pm_{2.5}) forecasting in smart cities. *Sensors*, 18(7):2220.

- [27] Hong T. (2014). Very short, short, medium and long term load forecasting. recuperado de: <http://blog.drhongtao.com/2014/10/very-short-short-medium-long-term-load-forecasting.html>.
- [28] Kolla A. (2008). Short, medium and long-term forecasting - marketing management. recuperado de: <https://www.wisdomjobs.com/e-university/marketing-management-tutorial-294/short-medium-and-long-term-forecasting-9586.html>.
- [29] Antti S. and Amaury L. (2006). Time series prediction using dirrec strategy. In *Esann*, volume 6, pages 143–148.
- [30] Souhaib B. (2014). Machine learning strategies for multi-step-ahead time series forecasting. *Universit Libre de Bruxelles, Belgium*, pages 75–86.
- [31] Brownle J. (2017). 4 strategies for multi-step time series forecasting. recuperado de: <https://machinelearningmastery.com/multi-step-time-series-forecasting/>.
- [32] Adil A. and Muhammad K. (2017). Multi-step ahead wind forecasting using nonlinear autoregressive neural networks. *Energy Procedia*, 134:192–204.
- [33] Yongnan J., Jin H., Nima R, and Amaury L. (2005). Direct and recursive prediction of time series using mutual information selection. In *International Work-Conference on Artificial Neural Networks*, pages 1010–1017. Springer.
- [34] Brockwell P. and Davis R. (2016). *Introduction to time series and forecasting*. springer.
- [35] Ríos G. (2008). *Series de Tiempo*. FCFM, Universidad de Chile, semestre primavera 2008.
- [36] Aneiros G. (2009). *Series de Tiempo*. Universidad de la Coruña. Curso 2008-09.

- [37] Faraway J. and Chatfield C. (1998). Time series forecasting with neural networks: A comparative study using the airline data. *Applied Statistics*, 47(2):231–250.
- [38] Velásquez J. and Franco C. (2010). Predicción de series temporales usando máquinas de vectores de soporte. *Revista chilena de ingeniería*, 18(1):64–75.
- [39] Marulanda C. (2007). *Análisis de series temporales con R (II): Estacionariedad y raíces unitarias*. FinanzasZone.
- [40] Nesbitt J. (2016). Adfstest - prueba estacionaria de dickey fuller aumentada. <https://support.numxl.com/hc/es/articles/215571923-ADFTest-Prueba-estacionaria-de-Dickey-Fuller-Aumentada>.
- [41] Du Boisberranger J. (2020). Sklearn versión 0.22. <https://scikit-learn.org/stable/index.html>.
- [42] Du Boisberranger J. and Van den Bossche J. (2019a). Linear models. https://scikit-learn.org/stable/modules/linear_model.html.
- [43] Du Boisberranger J. and Van den Bossche J. (2019b). Nearest neighbors. <https://scikit-learn.org/stable/modules/neighbors.html>.
- [44] Du Boisberranger J. and Van den Bossche J. (2019c). Decision trees. <https://scikit-learn.org/stable/modules/tree.html>.
- [45] Du Boisberranger J. and Van den Bossche J. (2019d). Ensemble methods. <https://scikit-learn.org/stable/modules/ensemble.html>.
- [46] Du Boisberranger J. and Van den Bossche J. (2019e). Neural network models (supervised). https://scikit-learn.org/stable/modules/neural_networks_supervised.html#regression.
- [47] Desarrolladores de python (2018). Python versión 3.7.2. <https://www.python.org/>.

-
- [48] Hunter J. and Firing E. (2020). Matplotlib versión 3.0.2. <https://matplotlib.org/>.
- [49] Augspurger T. (2020). Pandas versión 0.24.2. <https://pandas.pydata.org/>.
- [50] Perktold J. (2020). Statsmodels versión 0.9.0. <https://www.statsmodels.org/stable/index.html>.
- [51] Horst T. (2013). Stochastic lot sizing problems. In *Handbook of stochastic models and analysis of manufacturing system operations*, pages 313–344. Springer.

RESUMEN AUTOBIOGRÁFICO

Mario Alberto Gutiérrez Carrales

Candidato para obtener el grado de
Maestría en Ciencias de la Ingeniería
con Orientación en Sistemas

Universidad Autónoma de Nuevo León
Facultad de Ingeniería Mecánica y Eléctrica

Tesis:

PRONÓSTICO DE FALLAS EN LA DISTRIBUCIÓN DE ENERGÍA
ELÉCTRICA

Nací el 15 de agosto de 1995 en la ciudad de Monterrey, Nuevo León, México. Mis padres son José Alfredo Gutiérrez Ríos y María de los Ángeles Carrales Osoria. En 2017 egresé de la Licenciatura en Matemáticas en la Facultad de Ciencias Físico Matemáticas de la Universidad Autónoma de Nuevo León.