



GRADO ...

—— TRABAJO FIN DE ESTUDIOS ——

*Introducción a
la Estadística Aplicada
con ayuda de R*

Marta García Moreno

Sevilla, Octubre de 2017

Índice general

Prólogo	III
Resumen	V
Abstract	VI
Índice de Figuras	VII
Índice de Tablas	IX
1. Capítulo 1: Introducción	1
1.1. Introducción	1
1.2. Objetivos	2
1.3. Hipótesis	3
1.4. Revisión bibliográfica	3
1.5. Análisis del problema	4
2. Preliminares	5
2.1. Datos georreferenciados	5
2.1.1. Datos Vectoriales	5
2.1.1.1. Simple features	5
2.1.2. Datos Raster	6
2.1.3. Sistemas de Referencia de Coordenadas	6
2.1.3.1. Sistemas de Coordenadas Geográficas	7
2.1.3.2. Sistemas de Coordenadas Proyectadas	7
2.2. Análisis exploratorio de datos	8
2.2.1. Depuración de los datos	8
2.3. Modelos	9
2.3.1. Regresión logística (con penalización)	9
2.3.2. Support Vector Machine	9
2.3.3. Random Forest	9
2.3.4. Redes Neuronales	9

2.3.5. Validación del ajuste	9
2.4. Evaluación modelos	9
2.4.1. Clasificación binaria	9
2.5. Herramientas	10
3. Construcción de la base de datos	11
3.1. Determinación del marco del estudio	12
3.1.1. Incendios forestales	12
3.1.2. Variables predictoras	13
3.2. Fuentes de datos	14
3.3. Procesamiento de los datos	15
4. Capítulo 4: Cuerpo	19
4.1. Obtención y organización de los datos	19
4.2. Depuración y preprocesamiento	19
4.3. Análisis exploratorio	19
4.4. Estudio de las variables	19
4.5. Modelización	19
4.5.1. Modelo 1	19
4.5.2. Modelo 2	19
4.6. Comparación	19
A. Apéndice: Título del Apéndice	21
A.1. Primera sección	21
B. Apéndice: Título del Apéndice	23
B.1. Primera sección	23
Bibliografía	25

Prólogo

Escrito colocado al comienzo de una obra en el que se hacen comentarios sobre la obra o su autor, o se introduce en su lectura; a menudo está realizado por una persona distinta del autor.

También se podrían incluir aquí los agradecimientos.

Resumen

Resumen. . .

Abstract

Abstract...

Índice de figuras

1.1. Spatial Prediction of Wildfire Susceptibility Using Field Survey GPS Data and Machine Learning Approaches, Omid Ghorbanzadeh, Khalil Valizadeh Kamran, Thomas Blaschke, Jagannath Aryal, Amin Naboureh, Jamshid Einali and Jinhu Bian.	3
---	---

Índice de tablas

3.1. Datos brutos	15
3.2. Conjunto de datos depurados	17

Capítulo 1

Capítulo 1: Introducción

1.1. Introducción

Posible estructura de la introducción obtenida de “Spatial Prediction of Wildfire Susceptibility Using Field Survey GPS Data and Machine Learning Approaches, Omid Ghorbanzadeh, Khalil Valizadeh Kamran, Thomas Blaschke, Jagannath Aryal, Amin Naboureh, Jamshid Einali and Jinhu Bian”

Superficie forestal en Andalucía. Principales usos de los bosques Biodiversidad Andalucía es la segunda comunidad de España con más terreno forestal, con 4.325.378 ha de suelo forestal que suponen el 49.37 % de su superficie.

(https://www.miteco.gob.es/es/biodiversidad/temas/inventarios-nacionales/inventario-forestal-nacional/superficie_por_uso.html)

En 2022, 15.786,64 hectarias fueron afectadas por el fuego en Andalucía, prácticamente el doble de la media anual de los 10 años anteriores, 8873,68 hectarias [Memoria Plan INFOCA 2022].

Riesgo para la biodiversidad vegetal y animal. Daño a la infraestructura y riesgo para la vida de las personas. Costes incendios forestales [Barómetro de las catástrofes Naturales en España 2022]

En el Plan Estratégico Estatal del Patrimonio Natural y de la Biodiversidad a 2030 [Cita] “identifica las principales presiones y amenazas, entre las cuales destaca el aumento de la superficie forestal afectada por incendios forestales en 2022”.

Incendios forestales (>1ha) vs conatos (<=1ha).

Incendios forestales: ocurren en monte (parte incendio).

Las causas de los incendios de pueden agregar en dos grupos: causas naturales o causas antropogénicas. Las últimas se dividen en accidentes, negligencias o intencionados. [... Explicación breve cada una]. En 2022, solo el 2.23 % de los incendios fueron debidos a causas naturales, mientras que un 35.01 % fue debido a negligencias, un 40.91 % fue intencionado y un 6.47 % fue debido a causas accidentales.

Estas cifras muestran la importancia del factor humano en el estudio de los incendios forestales. En [Human-caused wildfire risk rating for prevention planning in Spain, Jesús Martínez, Cristina Vega-Garcia, Emilio Chuvieco, 2008] se ha estudiado el factor humano

en la causa de incendios forestales en España, detectando las variables más significativas en los modelos de predicción contruidos (alcanzando una tasa de acierto del 76 % en los datos test).

En el presente trabajo se pretende estudiar el riesgo de que los distintos puntos del territorio Andaluz se vean afectados por un incendio forestal en un momento concreto a través de [26] variables explicativas (antropogénicas, hidrográficas, topográficas, meteorológicas y vegetación). A diferencia de la mayoría de estudios similares realizados, se aborda el problema desde una perspectiva dinámica. Esto significa que no se pretende obtener un valor estacionario o estático del riesgo de incendio forestal en cada punto, si no que se busca estimar el riesgo de que una zona concreta se vea afectada por un incendio forestal en un momento concreto. Para ello, se consideran los valores de las variables correspondiente al momento concreto de la observación (usando la información más desagregada temporal y especialmente de la que se ha podido disponer).

Si bien las variables topográficas, hidrográfica y las ligadas a la infraestructura se han considerado constantes en los 20 años del periodo en estudio (por considerarse estructurales), las variables ligadas a las condiciones climáticas, al número de habitantes del municipio y al estado de la vegetación se consideran relativas al momento concreto de cada observación.

Población: variaciones en las tendencias demográficas tienen un impacto importante sobre los regímenes de incendios que se producen [Añadir fuente]

Variabes climáticas y del estado de la vegetación: Variables con una gran variabilidad temporal que tienen un impacto directo en la aparición y propagación de los incendios forestales.

Este enfoque hace posible que se tengan en cuenta las variaciones que se producen en las variables a lo largo del tiempo (cambios en el número de habitantes, en las condiciones climáticas y en el estado de la vegetación), haciendo así posible añadir una nueva dimensión al estudio.

Spain could be considered as a key area for wildfire modeling since it is, by far, the most fire-affected territory within the European Union -> Alguna característica relevante de Andalucía ¿Biodiversidad?

INFOCA

1.2. Objetivos

El objetivo de esta investigación será construir modelos que permitan predecir el riesgo de incendio forestal en la Comunidad Autónoma de Andalucía.

Subobjetivos:

1. Construir un conjunto de datos que permita la realización de análisis y la posterior construcción de modelos de Machine Learning para la predicción del riesgo de incendio forestal en Andalucía a partir de un estudio previo del problema.
2. Modelizar el riesgo de incendio forestal usando distintos algoritmos de ML y comparar sus resultados

Table 6. *Cont.*

No	Factors	Impacts	References
3	Altitude (m)	Altitude is an essential feature of fire danger distribution that should be considered. The wildfires that occur at higher altitudes are less severe because of the increase in moisture.	Koutsias et al. 2002, [30]; Canteaume, et al. 2013, [31] Jaafari et al. 2019, [26]
4	Annual temperature (°C)	There is a direct relationship between temperature increase and wildfires.	Baltar et al. 2015, [32]; Oulad Sayad et al. 2019, [10]
5	Annual rainfall (mm)	The annual rainfall parameter is one of the most significant variables of wildfires; rainfall moisture influences the speed of wildfires, which makes more extension of the burned area.	Vasilakos et al. 2009, [33]; Tanskanen et al. 2005, [34]
6	Wind effect	Wind can affect the extension and direction of the wildfires immediately after their ignition.	Darvishsefat et al. 2018, [11]; Sakellariou et al. 2016, [3]; Fovell and Gallagher et al. 2018, [35]
7	Plan curvature (100/m)	The positive curvature can be considered convex, such as the top of the hills, while negative curvature is concave, which refers to features like valleys. These criteria have different effects on the dynamics of wildfires.	Hilton et al. 2016, [36]; Pourtaghi et al. 2015, [4]
8	Topographic wetness index (TWI)	Fuel moisture is directly related to the required heat of ignition occurs. The actual relationship between the TWI and wildfires differs from other ground conditions and features.	Porensky et al. 2018, [37]; Ghorbanzadeh and Blaschke, 2018, [12]
9	Landform	Areas with steep slopes usually present the highest percentage of wildfires	Cantarello et al. 2011, [38];
10	Land use	Land use patterns based on shape and type have different impacts on wildfire risk.	Pourghasemi et al. 2016, [29]
11	NDVI	Reduction of the NDVI can cause an increase in water stress and the risk of fire.	Verbesselt et al. 2006, [39]; Pourtaghi et al. 2015, [4]
12	Distance to stream (m)	There is an indirect relationship between the distance from water sources and wildfire risk.	Razali and Sheriza 2010, [40]; Lee et al. 2010
13	Distance to road (m)	Roads provide access to forest areas; as a result, the risk of wildfire increases.	Syphard et al. 2008 Lee et al. 2010, [9]
14	Recreation area (m)	Recreation areas are places for human gatherings; humans, intentional or unintentional, can increase the risk of wildfire.	Stephens, 2005, [41]; Keeley and Fotheringham, 2003, [42]
15	Potential solar radiation	Increasing solar radiation can cause a reduction in the soil moisture and an increase in temperature and, consequently, wildfire risk.	Peters et al. 2013, [43]; Oulad Sayad et al. 2019, [10]
16	Distance to villages (m)	Expansion of residential area can increase the risk of wildfires, mostly because of human activities.	Canu et al. 2017, [44]; Lee et al. 2010, [9]

Figura 1.1: Spatial Prediction of Wildfire Susceptibility Using Field Survey GPS Data and Machine Learning Approaches, Omid Ghorbanzadeh, Khalil Valizadeh Kamran, Thomas Blaschke, Jagannath Aryal, Amin Naboureh, Jamshid Einali and Jinhu Bian.

3. Analizar potenciales casos de interés.

1.3. Hipótesis

“Spatial Prediction of Wildfire Susceptibility Using Field Survey GPS Data and Machine Learning Approaches, Omid Ghorbanzadeh, Khalil Valizadeh Kamran, Thomas Blaschke, Jagannath Aryal, Amin Naboureh, Jamshid Einali and Jinhu Bian”

1.4. Revisión bibliográfica

“Spain on fire: A novel wildfire risk assessment model based on image satellite processing and atmospheric information, Helena Liz-López , Javier Huertas-Tato a , Jorge Pérez-

Aracil, Carlos Casanova-Mateo, Julia Sanz-Justo, David Camacho”

“A review of machine learning applications in wildfire science and management Piyush Jain, Sean C.P. Coogan, Sriram Ganapathi Subramanian, Mark Crowley, Steve Taylor, and Mike D. Flannigan”

“Los incendios forestales en Andalucía: investigación exploratoria y modelos explicativos” Oliver Gutiérrez-Hernández (1*), José María Senciales-González (2), Luis V. García (1)

1.5. Análisis del problema

El área de estudio abarca el conjunto de la comunidad autónoma de Andalucía (España), la región más meridional de la península Ibérica, un territorio de 87.268 km², donde el 50,8 % de la superficie está ocupada por usos y cubiertas forestales.

-> MAPA

2002-2022 Razones: - Disponibilidad datos - Cambios en los regímenes de incendios

Capítulo 2

Preliminares

2.1. Datos georreferenciados

Todos los datos empleados en este trabajo son georreferenciados, lo que significa que están asociados a ubicaciones geográficas específicas. Por ello, resulta esencial introducir, aunque sea de forma general, los tipos de datos más utilizados para trabajar con esta información, sus características y las herramientas disponibles para manipularlos. Se tratarán los datos vectoriales y los datos rasters, al ser los tipos de datos fundamentales en este contexto, con características bien diferenciadas entre ellos.

2.1.1. Datos Vectoriales

El modelo de datos vectoriales geográficos se basa en puntos ubicados dentro de un sistema de referencia de coordenadas (CRS, por sus siglas en inglés). Estos puntos pueden representar características independientes o pueden estar conectados para formar geometrías más complejas como líneas y polígonos.

2.1.1.1. Simple features

Las “Simple features” son un estándar abierto ampliamente usado para la representación de datos vectoriales, desarrollado y respaldado por el Open Geospatial Consortium (OGC, por sus siglas en inglés), una organización sin ánimo de lucro dedicada a la creación de estándares abiertos e interoperables a nivel global dentro del marco de los sistemas geográficos de información (GIS, por sus siglas en inglés) y de la World Wide Web.

El paquete `sf` proporciona clases para datos vectoriales geográficos y una interfaz de línea de comandos consistente para importantes bibliotecas de bajo nivel para geoprocesamiento (GDAL, PROJ, GEOS, S2, ...).

Los objetos `sf` son fáciles de manipular ya que son `dataframes` o `tibbles` con dos características fundamentales. En primer lugar, contienen metadatos geográficos adicionales: tipo de geometría, dimensión, “Bounding Box” (límites o extensión geográfica) e información sobre el Sistema de referencia de coordenadas. Y además, presentan una columna de geometrías que tiene el nombre de “geom”. Algunas ventajas del uso del modelo de “simple features” en R son que en la mayoría de operaciones los objetos `sf` se pueden tratar como

data frames, los nombres de las funciones son consistentes (todos empiezan por `st_`), las funciones se pueden combinar con el operador tubería y además funcionan bien con el ecosistema de paquetes tidyverse.

El paquete `sf` de R soporta 18 tipos de geometrías para las simple features, de las cuales las más utilizadas son: POINT, LINESTRING, POLYGON, MULTIPOINT, MULTILINESTRING, MULTIPOLYGON and GEOMETRYCOLLECTION.

2.1.2. Datos Raster

El modelo de datos raster representa el espacio con una cuadrícula de celdas (también llamadas píxeles), que generalmente es regular, es decir, con todas las celdas de igual tamaño. Aunque no se tratarán en el presente trabajo, cabe mencionar que existen otros modelos de raster más complejos en los que se usan cuadrículas irregulares (rotadas, truncadas, rectilíneas o curvilíneas) y que pueden manipularse con el paquete de R (stars)[<https://cran.r-project.org/web/packages/stars/index.html>]. A cada una de estas celdas se le asocia uno (rasters de una sola capa) o varios (rasters multicapa).

Los datos en formato raster constan de una cabecera y una matriz cuyos elementos representan celdas equiespaciadas. En la cabecera del raster se definen el Sistema de referencia de coordenadas, la extensión (o límites espaciales del área cubierta por el ráster), la resolución y el origen. El origen son las coordenadas de uno de los píxeles del ráster, que sirve de referencia para los demás, siendo generalmente utilizado el de la esquina inferior izquierda (aunque el paquete TERRA usado en este trabajo usa por defecto el de la esquina superior izquierda). La resolución se calcula como:

$$resolution = \frac{x_{max} - x_{min}}{ncol}, \frac{y_{max} - y_{min}}{nrow}$$

La representación en forma de matriz evita tener que almacenar explícitamente las coordenadas de cada una de las cuatro esquinas de cada píxel, debiendo almacenar solamente las coordenadas de un punto (el origen). Esto, unido a las operaciones del álgebra de mapas hacen que el procesamiento de datos raster sea mucho más eficiente que el de datos vectoriales.

Se usará el paquete TERRA para tratar los datos en formato ráster. Este paquete permite tratar el modelo de rásters regulares con una o varias capas a través de la clase de objetos `SpatRaster`. Sin embargo, existen otras alternativas, como el paquete (stars)[<https://cran.r-project.org/web/packages/stars/index.html>], que además de ser más potente, permite trabajar con rásters no regulares y ofrece una mejor integración con el paquete `sf` y el entorno tidyverse.

2.1.3. Sistemas de Referencia de Coordenadas

Intrínseco a cualquier modelo de datos espaciales está el concepto de Sistema de referencia de coordenadas (CRS), que establece cómo la geometría de los datos se relaciona con la superficie terrestre. Es decir, es el nexo de unión entre el modelo de datos y la realidad, por lo que juega un papel fundamental. Los CRS pueden ser de dos tipos: geográficos o proyectados.

2.1.3.1. Sistemas de Coordenadas Geográficas

Los sistemas de coordenadas geográficas (GCS por sus siglas en inglés) identifican cada punto de la superficie terrestre utilizando la longitud y la latitud. La longitud es la distancia angular al Meridiano de Greenwich medida en la dirección Este-Oeste. La latitud es la distancia angular al Ecuador medida en la dirección Sur-Norte.

Cualquier sistema de coordenadas geográficas se compone de tres elementos: el elipsoide, el geoide y el datum. El primero es el elipsoide (o esfera) utilizado para representar de forma simplificada la superficie terrestre, sobre el que se supone que se encuentran los datos y el que permitirá realizar mediciones. El segundo, el geoide, es el modelo matemático que representa la verdadera forma de la Tierra, que no es suave sino que presenta ondulaciones debidas a las fluctuaciones del campo gravitatorio a lo largo de la superficie terrestre, que además cambian a una amplia escala temporal. Y el tercero, el datum, indica cómo se alinean el elipsoide y el geoide, es decir, cómo el modelo matemático se ajusta a la realidad. Este puede ser local o geocéntrico, en función de si el elipsoide se ajusta al geoide en un punto concreto de la superficie terrestre o de si el centro del elipsoide es el que se alinea con el centro de la Tierra. Ejemplos de datums geocéntricos usados en este trabajo son:

- European Terrestrial Reference System 1989 (ETRS89), usado ampliamente en la Europa Occidental.
- World Geodetic System 1984 (WGS84), usado a nivel global.

2.1.3.2. Sistemas de Coordenadas Proyectadas

Un Sistema de Coordenadas Proyectadas (PCS por sus siglas en inglés) es un sistema de referencia que permite identificar localizaciones terrestres y realizar mediciones en una superficie plana, es decir, en un mapa. Estos sistemas de coordenadas se basan en las coordenadas cartesianas, por lo que tienen un origen, un eje X y un eje Y y usan una unidad lineal de medida (en este trabajo, metro). Pasar de una superficie elíptica (GCR) a una superficie plana (PCS) requiere de transformaciones matemáticas apropiadas y siempre induce deformaciones en los datos.

Al proyectar la superficie terrestre en una superficie plana siempre se modifican algunas propiedades de los objetos, como el área, la dirección, la distancia o la forma. Un PCS solo puede conservar alguna de estas propiedades, por lo que es habitual clasificar los PCS en función de la propiedad que mantienen: las proyecciones de igual área preservan el área, las azimutales preservan la dirección, las equidistantes preservan la distancia y las conformales preservan la forma local. La mayoría de las proyecciones también se pueden clasificar en planas, cilíndricas o cónicas en función de cómo se realiza la proyección.

Un caso particular y ampliamente usado de PCS cilíndrico son los Universe Transverse Mercator (UTM), en el los que se proyecta el elipsoide sobre un cilindro tangente a este por las líneas de longitud (los meridianos). De esta forma, se divide el globo en 60 zonas de 6° de longitud, para cada una de las cuales existe un PCS UTM correspondiente que está asociado al meridiano central. Se trata de proyecciones conformales, por lo que preservan ángulos y formas en pequeñas regiones, pero distorsionan distancias y áreas.

A lo largo de este trabajo se utilizará ampliamente el Sistema de coordenadas proyectadas UTM30N (es habitual especificar el hemisferio para evitar confusión en los valores del eje Y, ya que miden distancia al ecuador, de ahí la N de hemisferio norte).

2.2. Análisis exploratorio de datos

El análisis exploratorio de datos (EDA, por sus siglas en inglés), es una parte fundamental de todo proyecto de Machine Learning y en general de cualquier proyecto en el que se deba trabajar con datos de cualquier procedencia para extraer de ellos conclusiones. Antes del procesamiento de los datos es siempre necesario explorar, entender y evaluar la calidad de estos, pues como indica la expresión inglesa *garbage in, garbage out*, si trabajamos con datos pobres, no podemos esperar obtener buenos resultados con ellos.

El EDA hace referencia al conjunto de técnicas estadísticas con las que se pretende explorar, describir y resumir la naturaleza de los datos, comprender las relaciones existentes entre las distintas variables presentes, identificar posibles errores o revelar posibles valores atípicos, todo esto con el objetivo de maximizar nuestra comprensión sobre el conjunto de datos.

2.2.1. Depuración de los datos

La depuración de los datos o *data cleaning* es el proceso de detectar y corregir o eliminar datos incorrectos, corruptos, con formato incorrecto, duplicados o incompletos dentro de un conjunto de datos. Puede considerarse una fase dentro del EDA (como se sugiere en R4DS, Wickman) o una fase previa a este.

Puede entenderse que el *data cleaning* es el proceso de pasar de *raw data* o datos en bruto a datos técnicamente correctos y finalmente a datos consistentes.

Entendemos por datos técnicamente correcto cuando cada valor pertenece a una variable y está almacenado en el tipo que le corresponde en base al conocimiento del dominio del problema. Para ello se debe reajustar el tipo de cada variable al que le corresponda en base al conocimiento que se tenga sobre esta, codificando los valores en las clases adecuadas si fuese necesario.

Decimos que un conjunto de datos es consistente cuando es técnicamente correcto y adecuado para el análisis estadístico. Se trata, por tanto, de datos que han eliminado, corregido o imputado los valores faltantes, los valores especiales, los valores atípicos y los errores.

2.3. Modelos

2.3.1. Regresión logística (con penalización)

2.3.2. Support Vector Machine

2.3.3. Random Forest

2.3.4. Redes Neuronales

2.3.5. Validación del ajuste

Partición entrenamiento/ validación / test

2.4. Evaluación modelos

Una vez construido un modelo predictivo es necesario conocer el rendimiento de este sobre nuevos datos, con el objetivo de estimar su capacidad de generalización. Esto es fundamental de cara a determinar si el modelo es adecuado para el propósito previsto o si necesita ajustes o mejoras. Además, la evaluación del rendimiento permite comparar entre diferentes modelos y seleccionar el que mejor se adapte a las necesidades específicas del problema en cuestión. Para ello, se recurre a distintas métricas, en función de las características propias de cada problema.

2.4.1. Clasificación binaria

En el presente trabajo el problema que se aborda es un problema de clasificación binaria, pues tenemos solo dos clases que son la clase positiva y la clase negativa. A la hora de clasificar una nueva instancia pueden darse 4 situaciones:

- Que se clasifique como positiva siendo realmente positiva, en cuyo caso se dirá que forma parte de las *True Positives (TP)*
- Que se clasifique como negativa siendo realmente negativa, en cuyo caso se dirá que forma parte de las *True Negatives (TN)*
- Que se clasifique como positiva siendo realmente negativa, en cuyo caso se dirá que forma parte de las *False Positives (FP)*
- Que se clasifique como negativa siendo realmente positiva, en cuyo caso se dirá que forma parte de las *False Negatives (FN)*

Se definen las siguientes métricas de rendimiento de un modelo de clasificación binaria:

Tasa de acierto o exactitud. Mide la proporción de casos que han sido correctamente clasificados.

$$Exactitud = \frac{TP + TN}{TP + FP + TN + FN}$$

Precisión. Mide la proporción de casos clasificados como positivos que realmente lo son.

$$Precisión = \frac{TP}{TP + FP}$$

Especificidad. Mide la proporción de casos negativos que han sido correctamente clasificados por el modelo.

$$Especificidad = \frac{TN}{TN + FP}$$

Sensibilidad o recall. Mide la proporción de casos positivos que han sido correctamente clasificados por el modelo.

$$Recall = \frac{TP}{TP + FN}$$

AUC-ROC. Mide el área bajo la curva ROC (*Receiver Operating Characteristic* o Característica Operativa del Receptor en castellano). Esta curva es una representación gráfica del rendimiento de un modelo de clasificación binaria para todos los umbrales de clasificación.

2.5. Herramientas

Toda la parte práctica del presente trabajo se ha llevado a cabo empleado el lenguaje de programación R a través del entorno de desarrollo integrado (IDE) que ofrece RStudio. R es un lenguaje y entorno de programación de código abierto desarrollado dentro del proyecto GNU y orientado a la computación estadística. R puede extender sus funcionalidades fácilmente a través de la gran cantidad de paquetes disponibles dentro del repositorio de paquetes de CRAN (The Comprehensive R Archive Network), siendo este uno de sus puntos fuertes, dada la gran comunidad de usuarios y desarrolladores con las que cuenta este lenguaje.

Los paquetes que se han utilizado han sido:

- tidyverse:
 - ggplot2, para la visualización.
 - dplyr, para la manipulación.
 - tidyr, para la ordenación.
 - readr, para la importación.
 - purrr, para la programación funcional
- tidymodels
- sf
- terra
- nasapower: obtención de información climática satelital
- mapSpain:

... (podemos seguir hasta el infinito)

Capítulo 3

Construcción de la base de datos

El primer paso a la hora de construir cualquier modelo de predicción es disponer de datos adecuados que permitan explicar correctamente el fenómeno en estudio, en este caso los incendios forestales en Andalucía. Con este fin, se ha llevado a cabo un extenso estudio previo del dominio del problema para conocer qué variables pueden ser relevantes de cara a la predicción de incendios forestales, analizando estudios similares realizados anteriormente así como otras fuentes relativas a la ecología del fuego, que nos permitiesen conocer el efecto que cabría esperar de estas variables.

Se ha querido adoptar un enfoque dinámico, es decir, el objetivo no es construir un modelo estacionario que nos indique si una determinada zona se verá afectada por un incendio forestal a lo largo de un amplio periodo temporal, si no que se pretende ser capaz de predecir si un determinado punto del territorio andaluz se verá afectado por un incendio forestal en un momento concreto, en base a las covariables correspondientes a ese lugar en ese momento. Es decir, se considera no solo la dimensión espacial de los datos si no también la temporal, al mayor nivel de desagregación disponible. Este es un enfoque mucho menos explorado, debido fundamentalmente a dos factores:

1. La dificultad de disponer de información fiable y de calidad desagregada espacio-temporalmente
2. La dificultad de trabajar con datos de estas características de cara al análisis y principalmente a la modelización, ya que son datos correlados en el tiempo y en el espacio.

Queda claro, por tanto, que se trata de un problema complejo que requiere de simplificaciones para poder ser abordado, más aun dadas las limitaciones en los recursos computacionales disponibles y la enorme cantidad de de datos que se están considerando y que requieren de un procesamiento sumamente costoso desde un punto de vista computacional.

Por todo ello, esta sección es probablemente la de mayor importancia y dificultad de todo el trabajo, ya que implica la toma de decisiones que serán determinantes de cara al correcto desempeño de los modelos que se construirán más adelante, requiere de un vasto conocimiento del problema que permita un enfoque adecuado que haga posible la consecución de los objetivos que se esperan conseguir, necesita del uso de técnicas específicas de procesamiento de datos espaciales que no han sido tratadas durante el grado y se ve fuertemente limitada por los escasos recursos computacionales disponibles.

3.1. Determinación del marco del estudio

El primer paso ha sido limitar el área y la franja temporal que abarcará el estudio. Para ello, ha sido necesario basarse principalmente en la disponibilidad y consistencia de la información requerida para el proyecto y en las limitaciones computacionales impuestas por el equipo disponible.

En cuanto a la disponibilidad de información, hay que diferenciar entre la información de incendios forestales y la información de variables que permitan explicar este fenómeno considerando la mayor desagregación espacial y temporal posible.

3.1.1. Incendios forestales

En lo referente a los datos sobre incendios forestales cabe mencionar que España cuenta con una de las mayores y más completas bases de datos sobre incendios forestales a nivel europeo. Se trata de la Estadística General de Incendios Forestales (EGIF), que en su versión definitiva actualmente contiene toda la información que se recoge en cada parte de incendio forestal que ha tenido lugar en España desde 1983 hasta 2015, incluyendo su información espacial con sus coordenadas de origen. Se ha explorado extensamente el uso de esta base de datos para el proyecto, dada su exhaustividad y completitud. Sin embargo, lamentablemente no ha sido posible en este caso incorporarla al trabajo por diversas razones.

La principal de ellas fue que hasta marzo de 2024 la base de datos de la EGIF solo se encontraba disponible en el Catálogo de Datos del Gobierno de España en formato TURTLE y esto conllevó numerosas dificultades. Se exploraron distintas librerías de R (y alguna de Python) para el manejo de datos en este formato como RDFlib. Sin embargo, al tratarse de una base de datos de un tamaño considerable (aproximadamente 1GB y con más de una decena de millones de tripletas), esta librería no era suficientemente eficiente para poder realizar consultas en un tiempo razonable al conjunto de datos. Tras explorar otras alternativas, se valoró la posibilidad de usar un triplestore, es decir, una base de datos especialmente diseñada para el almacenamiento y recuperación de tripletas a través de consultas semánticas. En este caso se usó Apache Jena Fuseki, ya que cuenta con una interfaz que facilita su uso. Sin embargo, aunque esto supuso una mejora considerable en la eficiencia y permitió realizar consultas sencillas a la base de datos, en este caso fue la complejidad del gráfico de datos (ontología) y la escasa documentación disponible sobre esta, la impidió que pudiese realizar las consultas más complejas que requería para llevar a cabo el proyecto. Además, se debe tener en cuenta que se trata de una base de datos muy heterogénea y con numerosos datos faltantes debida su naturaleza, por lo que requiere de un preprocesamiento que probablemente será complicado y costoso en tiempo y en recursos computacionales. Al no disponer de ninguno de estos, finalmente se optó por buscar una alternativa más abaricable dada las limitaciones con las que cuenta un Trabajo de Fin de Estudios, aunque queda abierta la posibilidad de explorar esta base de datos en futuros estudios, la cual aportar nuevas dimensiones al estudio de los incendios forestales en España gracias a la enorme cantidad de información que ofrece.

TURTLE es una sintaxis para RDF con una sintaxis compatible con SPARQL. RDF (Resource Description Framework) es un estándar de semántica web utilizado para el intercambiar de datos en la Web.

Ante esta situación, la solución planteada fue limitar el área en estudio a la Comunidad Autónoma de Andalucía, aprovechando la enorme disponibilidad de información medioambiental que ofrece la Red de Información Ambiental de Andalucía (REDIAM). En particular, se emplea la cartografía generada por la REDIAM sobre las áreas recorridas por los incendios forestales entre 1975 y 2022. Esta contiene los perímetros de incendios forestales mayores de 100 ha en Andalucía obtenidos a partir de imágenes de satélite y datos de campo. Se trata por tanto de una información que no es exhaustiva, pues los incendios con una extensión inferior a 100ha no han sido considerados. Sin embargo, frente a no disponer de otra información operativa de mayor calidad, se utilizará esta teniendo en cuenta que tendrá un efecto sobre las conclusiones que se puedan sacar de los modelos que se construyan.

3.1.2. Variables predictoras

Una vez limitada la extensión territorial del estudio el siguiente paso era acotar la franja temporal que abarcaría el estudio en base a la disponibilidad de datos adecuados para explicar el fenómeno en cuestión desagregados espacial y temporalmente.

Los incendios forestales son un proceso sumamente complejo, en el que actúan numerosos factores de muy distinta índole (...). Además, dentro de un incendio forestal se pueden distinguir distintas fases que presentan características muy diversas y sobre las que actúan distintos agentes: ignición, propagación y extinción. Dada la información sobre incendios forestales disponible, se está obligado a adoptar un enfoque global, pues no se dispone de los puntos de ignición u origen de los incendios forestales. El enfoque será, por tanto, intentar predecir si una determinada localización se verá afectada por un incendio forestal (de más de 100 ha) en un momento concreto.

Además, es importante tener en cuenta que existen factores estructurales que tienen una influencia directa sobre los regímenes de incendios forestales como son las tendencias de uso y explotación de los bosques, la presencia de interfaz urbano forestal, los tipos y técnicas de agricultura que se llevan a cabo, la presencia e intensidad del pastoreo, los cambios en los usos de suelo e incluso conductas sociales y tendencias demográficas diversas. Se trata de variables que cambian a lo largo de periodos relativamente largos de tiempo y que muy difícilmente pueden ser incluidos en los modelos, dada la falta de datos sobre ellas, así como su carácter transversal. Por ello, se ha considerado conveniente no extender en exceso el periodo de estudio, reconocida la imposibilidad de incluir en el modelo todas las variables que tienen un impacto relevante en la aparición de incendios y que son cambiantes en el tiempo.

Todo ello hace necesario que el conjunto de datos utilizado contenga información sobre todas las dimensiones (o al menos las principales) que influyen en cualquiera de las fases de un incendio forestal. Es decir, se deben incluir la dimensión antropogénica, la demográfica, la hidrográficas, la topográfica, la meteorológica y la vegetación. Es importante recalcar que siempre se hace referencia a datos geoespaciales pues debe ser la información relativa al lugar (y al momento) del incendio, con la dificultad posterior que esto supondrá.

Por último, es importante diferenciar entre características que se considerarán estructurales (y por tanto invariantes a lo largo del periodo de estudio) y aquellas que se considerarán variables en el tiempo. Dentro de las primeras se encuentran todas las características relacionadas con la topografía del terreno, las infraestructuras y los usos del suelo, como

por ejemplo el modelo de elevaciones, la distribución de asentamientos de población, la red de carreteras y el uso de suelo. Todas las demás variables de carácter demográfico, meteorológico o de vegetación se considerarán, por tanto, desagregadas temporalmente.

En base a todo lo mencionado y a la disponibilidad de información de calidad de las categorías comentadas, se ha decidido limitar la franja temporal del estudio a 20 años que van de 2002 a 2022, ambos inclusive.

3.2. Fuentes de datos

Como se ha comentado en la sección anterior, los datos sobre los incendios forestales se han obtenido de los perímetros de incendios forestales mayores de 100 ha en Andalucía entre 1975 y 2020 disponibles la REDIAM. De cada incendio registrado se dispone de su fecha de inicio, del área recorrida por el fuego y del municipio en el que originó, así como de otras variables que dependen del año de la campaña y que no son relevantes de cara a nuestro estudio.

Tomando como base estudios similares (...) y partiendo de las 6 categorías ya mencionadas se han recopilado 23 conjuntos de datos de distinto tipo que se usarán para explicar y predecir los incendios forestales en Andalucía. Estos conjuntos se recogen en la Tabla 3.1, donde también se indica la fuente de la que ha sido obtenido cada uno de ellos, el tipo de datos que contiene (indicando su resolución en el caso de los datos ráster) y la frecuencia de las observaciones (o resolución temporal) en el caso de las variables temporales.

Es relevante la heterogeneidad de los datos recopilados, pues se dispone tanto de datos tabulares como de datos espaciales y dentro de estos últimos de datos vectoriales y datos ráster, con distintas resoluciones, distintas frecuencias y distintos sistemas de referencia de coordenadas. Esto hará que el procesamiento de estos datos hasta obtener datos adecuados para el análisis estadístico sea costoso y que deban utilizarse técnicas específicas de geocumulación.

Cabe también mencionar que se ha optado por el uso de datos meteorológicos basados en modelos y en observaciones satelitares, en lugar del uso de datos provenientes de estaciones meteorológicas. Si bien la información de estaciones meteorológica puede ser más precisa, la dificultad de disponer de datos consistentes y continuos en el tiempo a lo largo del periodo de estudio de las variables meteorológicas seleccionadas ha hecho que este enfoque no sea viable. En esta dirección se ha explorado la API de la AEMET y algunos paquetes de R como `climate`, sin llegar a resultados satisfactorios. Por otro lado, el paquete `nasapower` permite la descarga de una gran cantidad de variables meteorológicas con frecuencia diaria y con una resolución de aproximadamente 0.5×0.625 grados de latitud y longitud (unos 50km). Si bien es cierto que no es lo ideal, es la única opción que se ha considerado viable y de cara a la construcción de unos primeros modelos aproximativos podría ser suficiente. Si quisiese extenderse el estudio, sería conveniente profundizar en la búsqueda de alternativas que permitan obtener información meteorológica de una mayor calidad.

Categoría	Dato	Fuente	Tipo de dato	Frecuencia
Topográficas	Altitud	DERA ^a	TIFF (100m)	-
	Orientación	REDIAM ^b	TIFF (100m)	-
	Pendiente	REDIAM	TIFF (100m)	-
	Curvatura	REDIAM	TIFF (100m)	-
Vegetación	NDVI	REDIAM	TIFF (250m)	Mensual
Antropogénicas	Uso de suelo	DERA	Shapefile	-
	Red de carreteras	DERA	Shapefile	-
	Red de ferrocarril	DERA	Shapefile	-
	Línea eléctrica	DERA	Shapefile	-
	Espacio protegido	DERA	Shapefile	-
	Senderos / Vías Verde / Carriles Bici	DERA	Shapefile	-
	Camino / Vías Pecuarias	DERA	Shapefile	-
Demográficas	Población del municipio	IECA ^c	csv	Anual
Hidrográficas	Principales Ríos	MAGRAMA ^d	Shapefile	-
Meteorológicas	Precipitación (mm/day)	NASA POWER ^e	df (0.5° x 0.625°)	Diaria
	Temperatura a 2m sobre la superficie (°)	NASA POWER	df (0.5° x 0.625°)	Diaria
	Humedad del suelo (%)	NASA POWER	df (0.5° x 0.625°)	Diaria
	Dirección del viento a 10 metros sobre la superficie terrestre(°)	NASA POWER	df (0.5° x 0.625°)	Diaria
	Humedad relativa a 2m sobre la superficie (%)	NASA POWER	df (0.5° x 0.625°)	Diaria
	Cantidad de precipitaciones (mm/day)	NASA POWER	dfdf (0.5° x 0.625°)	Diaria

Fuente: Elaboración propia

^a Datos Espaciales de Referencia de Andalucía (DERA)^b Descargas Rediam^c Instituto de Estadística y Cartografía de Andalucía (IECA)^d Ministerio de Agricultura, Alimentación y Medio Ambiente (MAGRAMA)^e NASA Prediction Of Worldwide Energy Resources (NASA POWER)

Tabla 3.1: Datos brutos

3.3. Procesamiento de los datos

Una vez se dispone de todos los conjuntos de datos que se usarán en el estudio, el siguiente paso será combinarlos de manera adecuada y transformarlos a un formato apto para el análisis estadístico y la construcción de modelos predictivos, es decir, a un data frame. Dado que el objetivo que se persigue es predecir si, dada unas condiciones meteorológicas concretas en un momento dado, un punto del territorio andaluz se verá afectado o no por un incendio forestal, será necesario disponer de una cantidad suficiente de muestras negativas y positivas distribuidas espacial y temporalmente que tengan asociadas las variables explicativas correspondientes.

Intuitivamente, las muestras positivas serán aquellas observaciones (puntos definidos en el tiempo y en el espacio) dentro del marco espacio-temporal del estudio en las que se ha detectado un incendio forestal en el día de la observación. Es decir, son observaciones dentro de los polígonos de incendios el día que estos se han producido. Por tanto, las muestras negativas serán observaciones dentro del marco espacio-temporal definido en las que no se ha detectado un incendio forestal. Es importante tener en cuenta que dado que solo se dispone de los incendios con una extensión mayor a 100 ha, la muestra cuenta con un importante sesgo, ya que los casos positivos están infrarepresentados. Por ello, no podremos hacer inferencia a todos los incendios forestales, si no solo a los de una extensión superior a 100ha.

A continuación se detalla el proceso seguido para generar el conjunto de datos depurado sobre el que desarrollar el estudio a partir de los distintos conjuntos de datos en bruto:

1. Generación de una muestra balanceada de casos positivos y negativos.

Para poder construir cualquier modelo de clasificación binaria se necesita disponer de

una muestra que cuente con un número suficiente de casos positivos y negativo. Además, es aconsejable trabajar con conjuntos de datos balanceados para evitar sesgos en los modelos de clasificación [ARTICULO].

Como ya se ha comentado, se considerarán observaciones positivas aquellas que se hayan visto afectadas por un incendio forestal en el día y lugar de la observación. En cambio, serán observaciones negativas aquellas que no se hayan visto afectadas por un incendio forestal en el día y lugar de la observación. Estas observaciones deberán generarse a partir de los polígonos de incendios disponibles. Para ello, se usará un enfoque similar al utilizado en [https://www.researchgate.net/publication/228527438_Learning_to_predict_forest_fires_with_different_data_mining_techniques], con una diferencia fundamental. En [cita al paper] usan los puntos de ignición como muestras positivas y su objetivo es predecir los puntos de origen de los incendios forestales. En cambio, en el presente trabajo no se disponen de los puntos de ignición de los incendios, por lo que el enfoque adoptado es ligeramente diferente; el objetivo es predecir las zonas que pueden verse afectadas por un incendio forestal (superior a 100ha) bajo unas circunstancias concretas. De esta forma, los casos positivos serán puntos aleatorios tomados dentro de los polígonos de incendio en los días que estos ocurrieron. Los casos negativos se generarán igual que en [cita paper]: se toman fechas aleatorias dentro del periodo de estudio y a cada una de ellas se le asocia una localización aleatoria dentro del área de estudio satisfaciendo que deben estar a al menos 15km de cualquier incendio detectado en un margen de ± 3 días. Esta forma de tomar los casos negativos asegura que estén lo suficientemente alejados de los incendios forestales para representar condiciones no influidas por estos, dando prioridad así a las áreas con una menor prioridad de ocurrencia de incendio en un período definido. Las ubicaciones de los ejemplos positivos y negativos de ocurrencia de incendios estaban vinculadas espacial y temporalmente a los datos descriptivos.

2. Asignación a cada observación los valores correspondientes a ese día y a esa localización concreta de todas las variables predictoras a partir de los conjuntos de datos que se han recopilado (recogidos en la Tabla 3.1). Para ello se ha hecho uso de las funciones disponibles en los paquetes `terra` y `sf`.
3. Depuración de la muestra generada, se eliminan los valores perdidos y se ajustan adecuadamente los tipos de las variables. La variable `WD10M` se codifica mediante los 4 puntos cardinales y sus bisectrices, generando así 8 clases. En el caso de la variable `orientacion` se procede de manera idéntica pero se incluye también la clase “plano”, si la pendiente del punto es 0.

El resultado es un conjunto de datos con 20998 observaciones de las variables que se muestran en la Tabla 3.2.

Categoría	Nombre	Descripción	Tipo
Topográficas	elevation	Elevación sobre el nivel del mar (m)	numérica
	orientacion	Orientación de la pendiente descendiente	categoría
	pendiente	Pendiente del terreno ($^{\circ}$)	numérica
	curvatura	Curvatura de la superficie	numérica
Vegetación	NDVI	Índice de vegetación de diferencia normalizada	numérica
Antropogénicas	uso_suelo	Clasificación del uso del suelo	categoría
	dist_carretera	Distancia a la carretera más cercana (m)	numérica
	dist_ferrocarril	Distancia a la vía de ferrocarril más cercana (m)	numérica
	dist_electr	Distancia a la línea eléctrica más cercana (m)	numérica
	enp	Espacio Natural Protegido	categoría
	dist_sendero	Distancia a la vía verde, al carril bici o al sendero más cercano (m)	numérica
	dist_camino	Distancia al camino o a la vía pecuaria más cercano (m)	numérica
Demográficas	poblacion	Número de habitantes del municipio	numérica
Hidrográficas	dist_rios	Distancia al río más próximo (m)	numérica
Meteorológicas	PRECTRORCORR	Promedio corregido del total de precipitaciones en la superficie de la tierra en masa de agua (incluye el contenido de agua en la nieve) (mm/día)	numérica
	T2M	Temperatura promedio del aire a 2 metros sobre la superficie de la tierra ($^{\circ}\text{C}$)	numérica
	GWETTOP	Porcentaje de humedad del suelo	numérica
	WD10M	Promedio de la dirección del viento a 10 metros sobre la superficie de la tierra	categoría
	WS10M	Promedio de la velocidad del viento a 10 metros sobre la superficie de la tierra (m/s)	numérica
	RH2M	Humedad relativa a 2 metros sobre la superficie de la tierra	numérica
Variable Objetivo	fire	Incendio forestal	categoría
Identificadoras	date	Fecha de la observación	fecha
	municipio	Nombre del municipio	texto
	cod_municipio	Código del municipio	texto
	geometry	Geometría de los puntos	sfc

Tabla 3.2: Conjunto de datos depurados

Capítulo 4

Capítulo 4: Cuerpo

4.1. Obtención y organización de los datos

4.2. Depuración y preprocesamiento

En este caso, se ha realizado una fase previa de depuración del conjunto de datos con el objetivo de tener un conjunto de datos correcto técnicamente. Para ello, lo primero ha sido detectar, comprender y eliminar o imputar los valores faltantes, así como reajustar los tipos de las variables, codificando para ello las que así lo requiriesen en base al conocimiento previo del dominio del problema. Posteriormente, se ha estudiado la consistencia de los datos, mediante la identificación de valores atípicos durante el análisis descriptivo del conjunto de datos.

4.3. Análisis exploratorio

4.4. Estudio de las variables

4.5. Modelización

4.5.1. Modelo 1

4.5.2. Modelo 2

4.6. Comparación

Apéndice A

Apéndice: Título del Apéndice

A.1. Primera sección

Apéndice B

Apéndice: Título del Apéndice

B.1. Primera sección

Bibliografía

- [1] (Página web). «Universidad de Sevilla». Disponible en <https://www.us.es>.
- [2] ALLAIRE, JJ; XIE, YIHUI; DERVIEUX, CHRISTOPHE; MCPHERSON, JONATHAN; LURASCHI, JAVIER; USHEY, KEVIN; ATKINS, ARON; WICKHAM, HADLEY; CHENG, JOE; CHANG, WINSTON y IANNONE, RICHARD (2024). *rmarkdown: Dynamic Documents for R*.
<https://github.com/rstudio/rmarkdown>. R package version 2.26,
<https://pkgs.rstudio.com/rmarkdown/>.
- [3] FACULTAD DE MATEMÁTICAS (UNIV. SEVILLA) (s.f.).
<https://www.matematicas.us.es>.
- [4] LOPEZ, JUAN FERNANDO; FERNÁNDEZ HENAO, SERGIO y MORALES, MARCELA MARÍA (2007). «Aplicación de la programación por metas en la distribución de servicios entre empresas operadoras del sistema de transporte masivo». *Scientia et technica*, **13**(37), pp. 339–343.
- [5] LUQUE CALVO, PEDRO L. (2017). *Escribir un Trabajo Fin de Estudios con R Markdown*.
<http://destio.us.es/calvo>.
- [6] — (2019). *Cómo crear Tablas de información en R Markdown*.
<http://destio.us.es/calvo>.
- [7] LUQUE CALVO, PEDRO L. (2021). «Página personal de Pedro L. Luque».
<http://destio.us.es/calvo>.
- [8] R CORE TEAM (2016). *R: A Language and Environment for Statistical Computing*. R Foundation for Statistical Computing, Vienna, Austria.
<https://www.R-project.org/>.
- [9] RSTUDIO TEAM (2015). *RStudio: Integrated Development Environment for R*. RStudio, Inc., Boston, MA.
<http://www.rstudio.com/>.
- [10] XIE, YIHUI (2023). *knitr: A General-Purpose Package for Dynamic Report Generation in R*.
<https://yihui.org/knitr/>. R package version 1.45.