



UNIVERSIDAD DE BURGOS  
ESCUELA POLITÉCNICA SUPERIOR  
Grado en Ingeniería Informática



**TFG del Grado en Ingeniería  
Informática**

**Análisis de datos de  
temperatura y humedad de  
suelo procedentes de sensores  
IoT desplegados en un viñedo  
Documentación Técnica**



Presentado por Gabriel Hernández Vallejo  
en Universidad de Burgos — 26 de mayo  
de 2023

Tutores: Rubén Ruiz González, Alejandro  
Merino Gómez



---

# Índice general

---

Índice general	i
Índice de figuras	iii
Índice de tablas	iv
<b>Apéndice A Plan de Proyecto</b>	<b>1</b>
A.1. Introducción . . . . .	1
A.2. Planificación temporal . . . . .	1
A.3. Estudio de viabilidad . . . . .	7
<b>Apéndice B Comprensión del negocio</b>	<b>11</b>
B.1. Introducción . . . . .	11
B.2. Objetivos del proyecto . . . . .	11
B.3. Requisitos del proyecto . . . . .	11
<b>Apéndice C Comprensión de los datos</b>	<b>13</b>
C.1. Introducción . . . . .	13
C.2. Descripción de los datos . . . . .	13
<b>Apéndice D Preparación de los datos</b>	<b>21</b>
D.1. Introducción . . . . .	21
D.2. Valores faltantes . . . . .	21
D.3. Valores extremos . . . . .	22
D.4. Valores erróneos . . . . .	22
<b>Apéndice E Modelado</b>	<b>27</b>
E.1. Introducción . . . . .	27

<b>Apéndice F Evaluación</b>	<b>29</b>
F.1. Introducción . . . . .	29
<b>Apéndice G Documentación técnica de programación</b>	<b>31</b>
G.1. Introducción . . . . .	31
G.2. Estructura de directorios . . . . .	31
G.3. Manual del programador . . . . .	32
G.4. Pruebas del sistema . . . . .	34
<b>Apéndice H Documentación de usuario</b>	<b>35</b>
H.1. Introducción . . . . .	35
H.2. Requisitos de usuarios . . . . .	35
H.3. Instalación . . . . .	35
H.4. Manual del usuario . . . . .	36
<b>Bibliografía</b>	<b>37</b>

---

# Índice de figuras

---

C.1. Datos no procesados: sensor 1 . . . . .	15
C.2. Datos no procesados: sensor 2 . . . . .	16
C.3. Datos no procesados: sensor 3 . . . . .	16
C.4. Datos no procesados: sensor 4 . . . . .	17
C.5. Datos no procesados: sensor 5 . . . . .	17
C.6. Datos no procesados: sensor 6 . . . . .	18
C.7. Datos no procesados: sensor 7 . . . . .	18
C.8. Datos no procesados: sensor 8 . . . . .	19
C.9. Datos no procesados: pluviómetro . . . . .	19
D.1. Datos procesados: sensor 1 . . . . .	22
D.2. Datos procesados: sensor 2 . . . . .	23
D.3. Datos procesados: sensor 3 . . . . .	23
D.4. Datos procesados: sensor 4 . . . . .	24
D.5. Datos procesados: sensor 5 . . . . .	24
D.6. Datos procesados: sensor 6 . . . . .	25
D.7. Datos procesados: sensor 7 . . . . .	25
D.8. Datos procesados: sensor 8 . . . . .	26
G.1. Búsqueda del entorno virtual . . . . .	33
G.2. Selección del entorno virtual . . . . .	33

---

# Índice de tablas

---

A.1. Costes de Hardware . . . . .	7
A.2. Costes de Hardware . . . . .	8
A.3. Costes de Software . . . . .	8
A.4. Costes varios . . . . .	9
A.5. Costes totales . . . . .	9
A.6. Licencias de las dependencias utilizadas . . . . .	10
C.1. Atributos de los datos: sensores . . . . .	14
C.2. Atributos de los datos: pluviómetro . . . . .	14

## Apéndice A

---

# Plan de Proyecto

---

### A.1. Introducción

En la sección se introduce la fase de planificación del proyecto, tanto la planificación temporal, como la viabilidad económica y legal.

### A.2. Planificación temporal

La planificación del desarrollo del proyecto se realizó siguiendo la metodología ágil *Scrum*, aunque no en su definición completa, puesto que al ser una única persona en el desarrollo, no permite realizar todos los procedimientos recogidos dentro de los manuales de gestión de proyectos con metodologías ágiles [1], como las reuniones diarias “daily” para recoger diferentes cuestiones como cuáles son los problemas que ralentizan el proceso de desarrollo del proyecto.

Por otro lado, la herramienta software que apoyaba su aplicación (ZenHub [2]) dejó de estar disponible de forma abierta, de forma que no pudo seguirse con el modelo *Canvas* que aplicaba la metodología recogiendo en un tablero con varias columnas, según fuera su estado, las tareas por cada uno de los *sprints*.

Aun así, se ha continuado aplicando la filosofía ágil en líneas generales:

- Se continuó con el desarrollo iterativo incremental mediante los *sprints*.

- La duración media de estas iteraciones fueron de dos semanas al comienzo del proyecto y durante la mayor parte de su realización, reduciendo el tiempo al final de este.
- En la terminación de los *sprints* se realizaba un incremento, manteniendo reuniones con los tutores para la planificación de la iteración entrante y la revisión de posibles errores en el desarrollo.
- Las tareas surgidas de la reunión se creaba, estimaban y priorizaban, en un inicio añadiéndolas al tablero *Canvas*.

La estimación fue realizada empleando los *story points* disponibles en ZenHub, de forma que estos iban de las tareas más sencillas de implementar y rápidas, a las más complejas y que requerían de mayor tiempo de desarrollo.

## **Sprint 0 (18/01/2023 - 31/01/2023)**

En este *sprint* se presentó el proyecto, indicándose los primeros objetivos, lo que se enmarca en la fase de comprensión del negocio. Las primeras tareas se correspondían con la fase de comprensión de datos y a su vez la preparación de los mismos, puesto que el conjunto proporcionado contaba con diferentes casuísticas de las que en un comienzo se tenía constancia.

Entre los problemas encontrados previos al análisis en profundidad de los datos se encontraban lecturas incorrectas (valores erróneos) debido a problemas de *overflow* de variables en los controladores de cada sensor, de forma que cuando la batería bajaba de los 3.1V, las lecturas en el conjunto de datos se trataban de errores en su mayoría.

Además, en ciertos periodos de tiempo los sensores habían estado desconectados, por lo que se produjeron lecturas nulas (valores faltantes) que necesitaban ser tratadas.

Por otro lado, las lecturas en el pluviómetro instalado en las dependencias del viñedo en algunas de las muestras arrojaban valores inverosímiles debido a la forma en la que este se encontraba instalado, de manera que el balancín empleado para medir las diferencias de precipitaciones entre instantes se activaba artificialmente.

Durante el desarrollo de la iteración se comenzó con el tratamiento de los datos nulos, tomando la decisión de eliminar las entradas, pues se contaba con un número elevado de muestras y realizar la recuperación de estas no parecía viable debido a la naturaleza de los datos. Por otro lado, para los



problemas relacionados con el pluviómetro se solicitaron claves *API* para poder acceder a los registros de estaciones meteorológicas cercanas a la zona y poder comprobar los registros.

Por otro lado, se comenzó a realizar las gráficas de los datos para intentar encontrar correlaciones en los datos y poder hallar diferentes errores en el conjunto de datos.

### **Sprint 1 (31/01/2023 - 14/02/2023)**

En este *sprint* se presentaron los avances de la primera iteración y se acordó la continuación del proceso de comprensión de datos y preparación de los mismos.

Los esfuerzos se centraron entonces en modificar las gráficas de visualización de datos para continuar con la inspección visual de los sensores y comenzar con los datos del pluviómetro. Con estas inspecciones se consigue encontrar anomalías en ciertos sensores como en el caso del segundo y cuarto.

### **Sprint 2 (14/02/2023 - 28/02/2023)**

En la reunión de la iteración se propusieron diferentes objetivos:

- En el sensor 2 eliminar con variación excesiva en temperatura, corrigiendo los datos ruidosos.
- En el sensor 3 se encontró que el salto de humedad no era genuino.
- En el sensor 4 intentar realizar una media móvil para arreglar variables como la temperatura al emplear variaciones diarias en lugar de muestreos cada 5 minutos.

Durante el desarrollo del *sprint* se continuó con la selección de datos, encontrando que el efecto de la lluvia provocaba ciertos cambios bruscos en los datos como en el caso del sensor 5, pero estando algunos aparentemente no relacionados con este fenómeno. Se comenzó, por otro lado con la selección y limpieza de los datos mediante una columna adicional de validez en los ficheros correspondientes.

Además, se implementó la detección de *Outliers* mediante el rango intercuartílico, para eliminar el posible ruido del conjunto de datos de los sensores.

### **Sprint 3 (28/02/2023 - 14/03/2023)**

En la reunión de la iteración se acordó realizar la recuperación de los datos del sensor 8 hasta la variabilidad excesiva de los valores de las variables observadas, así como el intento de mejorar la calidad del conjunto de datos en otros sensores realizando procesos similares. Por otro lado, se propuso emplear WeatherBit como *API* meteorológica para eliminar los datos que más distaran de las lecturas del pluviómetro desplegado.

Durante el desarrollo del *sprint* se modificaron los umbrales empleados en la detección y manejo de valores extremos, para eliminar la mayoría cantidad de ruido posible en todos los sensores. Por otro lado, se fueron finalizando las selecciones de datos de los sensores, de forma que se realizó la recuperación de los datos en sensores cuyas variabilidades se debían a factores externos (fueron sacados de sus posiciones originales).

En lo referente del sensor 8 se llegó a la conclusión de que buena parte de los datos eran irrecuperables debido a las variabilidades aleatorias con las que parecía contar.

Por otro lado, se comenzó con la limpieza de los datos de los sensores empleando los datos de varias *API*, en primer lugar se trató de emplear la mencionada anteriormente, siendo el *endpoint* ideal de licencia de pago. Probando con OpenWeatherMap se presentó el caso al soporte, ampliando la licencia de usuario para poder acceder a l utilidad requerida, sin embargo no se permitía retrotraerse en más de un año en las lecturas registradas.

Finalmente se decide emplear los servicios de la AEMET, debido principalmente a que permitía obtener los datos de hacía más de un año, sin embargo, contaba con inconvenientes, y es que no permitía obtener datos por hora, sino por día y la estación meteorológica más cercana se encontraba a unos 15 Kilómetros del despliegue.

### **Sprint 4 (14/03/2023 - 28/03/2023)**

En la reunión del *sprint* se acuerda comenzar con la fase del modelado a la vista de la aparente correcta selección y limpieza de los datos proporcionados. Se plantea la creación de una matriz de correlación para observar cuáles de los atributos seleccionar para la creación del modelo.

Durante el desarrollo de la iteración los esfuerzos se centran en la creación del entorno virtual con las dependencias correspondientes que permitan realizar los modelados, además de dejar el fichero de pre-procesamiento de datos comentado con las explicaciones pertinentes.

Por otro lado, se realizan avances sustanciales en la memoria del proyecto añadiendo información a la introducción, objetivos y técnicas y herramientas, además de realizar una reestructuración de ficheros y una modificación de la forma en la que se realizan las gráficas de los datos.

Se comienza con la creación tentativa del modelo neuronal en busca de conocimiento del dominio, que permita profundizar en los módulos empleados para tal propósito, realizando regresión simple. Se investiga sobre las redes neuronales recurrentes, más concretamente sobre la aplicación de modelos como *GRU* y *LSTM* en problemas similares (predicciones de tiempo atmosférico).

### **Sprint 5 (28/03/2023 - 11/04/2023)**

En la reunión de la iteración se muestran los avances en memoria y en la creación de modelos, por otro lado, se plantea la forma de realizar la regresión empleando Keras.

Durante el *sprint* se modifica la manera en la que se crean los datos de entrada a los modelos, empleando medias diarias para tal propósito y dividiendo los datos en dos conjuntos: entrenamiento y validación. Además, se crean las gráficas de dispersión que permiten comparar visualmente las predicciones y valores reales.

Por otro lado, se continúa con el desarrollo de la memoria, añadiendo conceptos teóricos sobre el proceso de extracción de conocimiento de bases de datos (*KDD*), así como modificando algunos apartados.

### **Sprint 6 (11/04/2023 - 25/04/2023)**

En la planificación se considera entrenar a los modelos neuronales únicamente con valores adecuados, pues hasta ahora pueden existir saltos temporales que pueden afectar a la regresión, realizando medias temporales cada cierto periodo de tiempo para eliminar los posibles problemas de *offset* de los muestreos originales.

En el desarrollo del *sprint*, se modificó la selección en los datos de entrada a los modelos, de forma que se realizaban las medias diarias de los diferentes atributos y se les aplicaba el filtro de Hodrick-Presscott, un filtro de tendencias, que permite obtener las tendencias en una serie temporal, de manera que los modelos sean capaces de generalizar de forma más adecuada.

Por otro lado, se continuó con la memoria, en concreto con los conceptos teóricos relacionados con la fase de modelado.

### **Sprint 7 (25/04/2023 - 09/05/2023)**

En la reunión de la iteración se revisan los cambios realizados, llegando a la consideración de emplear varios días previos como entrada a los modelos, utilizando una ventana deslizante para seleccionar datos, de manera que se eviten los saltos temporales y estudiar la incorporación de una componente temporal como entrada en la regresión.

De esta forma, en el desarrollo se aplicó la ventana deslizante para la selección de datos de entrada en el modelo, de manera que se pudieran evitar los saltos temporales.

### **Sprint 8 (09/05/2023 - 16/05/2023)**

En la planificación se acordó intentar rebajar los periodos de tiempo, y que se emplearan datos medios por hora en lugar de diarios y probar los hiperparámetros de las redes, intentando variar las neuronas por capa, el número de estas, el ratio de aprendizaje, etc. para observar el comportamiento de los modelos y determinar tanto si continúan generalizando como si es necesario realizar cambios en el proceso de conocimiento de datos y preparación de estos.

De esta manera, se redujo la frecuencia de las medias de los datos, realizando la integración de los datos de los sensores en un único fichero, para poder entrenar los modelos empleando un único conjunto. Esto obligó a cambiar la forma en la que se graficaban los resultados, mostrando las gráficas de dispersión por sensor y atributo.

Por otro lado, se dividió el conjunto de datos en tres subconjuntos diferentes: un conjunto para entrenamiento, otro para validación y el restante para test, permitiendo la parametrización en la división de estos.

### **Sprint 9 (16/05/2023 - 23/05/2023)**

En la reunión de la iteración se acordó tratar de predecir tiempos más largos en lugar de un único instante como realizaban los modelos hasta el momento; previsiblemente el error sería mayor, por lo que se mencionó permitir parametrizar este valor.

Ante el objetivo principal planteado hubo que realizar diferentes modificaciones en la selección de los datos de salida y los parámetros, penalizando el rendimiento general de los modelos.

Se avanzó en la memoria en los conceptos teóricos y se realizaron algunas correcciones.

### **Sprint 10 (23/05/2023 - 30/05/2023)**

En la planificación básicamente se indicó la continuación en los esfuerzos de realización de memoria y anexos. De esta manera, este *sprint* se centró en la finalización de los mismos.

## **A.3. Estudio de viabilidad**

En esta sección se desglosarán la viabilidad económica y legal del proyecto, en cuanto a la primera indicará los costes derivados del desarrollo en un entorno real. En lo referente la segunda, se presentarán las licencias empleadas en el proyecto y su implicación con librerías de terceros.

### **Viabilidad económica**

En términos de viabilidad económica es necesario hacer una diferenciación entre los costes y los beneficios que lleva realizar el proyecto.

#### **Costes**

Los costes que pueden surgir del proyecto en un entorno empresarial pueden desglosarse en los siguientes:

##### **Costes de personal:**

El desarrollo se ha realizado con un única persona empleada a tiempo completo en aproximadamente 4 meses, de forma que se consideran los siguientes costes:

Concepto	Coste
Salario mensual neto	1.080€
Retención I.R.P.F(12 %)	185,46€
Seguridad Social(31,1 %)	480,64€
Salario mensual bruto	1.545€
<b>Total 4 meses</b>	<b>6.181,92€</b>

Tabla A.1: Costes de Hardware

El porcentaje mensual de aporte a la Seguridad Social se calcula como el 0,2 % al Fondo de Garantía Salarial (FOGASA), más el 0,6 %, el 6,7 % de la prestación por desempleo y el 23,6 % de contingencias comunes [3].

#### Costes de hardware:

El hardware empleado tiene un tiempo de amortización aproximado de 4 años.

Concepto	Coste	Coste amortizado
Ordenador portátil	900€	45€
<b>Total</b>	900€	45€

Tabla A.2: Costes de Hardware

#### Costes de software:

En cuanto a los costes de software, hay ciertos programas o herramientas empleadas durante el desarrollo que requieren de licencia de pago, estos contarán con un tiempo de amortización estimado de 2 años.

Concepto	Coste	Coste amortizado
Windows 10 Home	148€	18,5€
<b>Total</b>	148€	18,5€

Tabla A.3: Costes de Software

#### Costes varios:

Al igual que en otras situaciones empresariales, en los desarrollos de software surgen tanto costes inesperados, como fijos. Se reflejarán aquellos costes variados que aparentemente se tratarían de costes fijos.

Concepto	Coste
Memoria y anexos	50€
Alquiler de espacio de trabajo	592€
Internet	120€
<b>Total</b>	<b>762€</b>

Tabla A.4: Costes varios

**Costes totales:**

Los costes del desarrollo del proyecto son los siguientes:

Concepto	Coste
Personal	6.181,92€
Hardware	900€
Software	148€
Varios	762€
<b>Total</b>	<b>7.991,92€</b>

Tabla A.5: Costes totales

**Beneficios**

En cuanto a los beneficios generados, en el *scope* del proyecto no había cabida a un despliegue en los meses de desarrollo de los modelos, por lo que a corto plazo no habría ningún beneficio aparente.

Sin embargo, la forma de obtener los beneficios entraría dentro de una segunda fase, cuando se realice el despliegue del resultados en algún producto.

**Viabilidad legal**

En esta sección se desglosarán los temas relacionados con las licencias de uso de los productos software.

En primer lugar se analizará la licencia más conveniente en el desarrollo del proyecto, teniendo en cuenta las dependencias empleadas y sus respectivas licencias, puesto que la selección estará limitada por las restricciones de estas principalmente.

Dependencia	Versión	Descripción	Licencia
TensorFlow	2.11.0	Biblioteca de aprendizaje automático	Apache v2.0
Pandas	1.5.3	Biblioteca especializada en manipulación y análisis de datos	BSD
Matplotlib	2.11.0	Biblioteca para generación de gráficos	BSD
IPython Kernel for Jupyter	6.21.3	Biblioteca para manipulación de Python Notebooks	BSD
Statsmodel	0.13.5	Biblioteca con modelos y funciones estadísticas	BSD

Tabla A.6: Licencias de las dependencias utilizadas

De esta forma, se debe seleccionar una licencia compatible con Apache v2.0 y BSD, siendo la primera la más restrictiva de las licencias, que requiere la conservación del aviso de derecho de autor y un descargo de responsabilidad, sin embargo, no es *copyleft* [4].

Por su versatilidad se escoge una licencia MIT, de forma que el se trata de una licencia permisiva de software libre, imponiendo muy pocas restricciones y otorgando una buena compatibilidad con otras licencias.



## *Apéndice B*

---

# **Comprensión del negocio**

---

### **B.1. Introducción**

En la fase de comprensión del negocio se pretende analizar los objetivos y requisitos del proyecto.

Estos quedarán bien marcados por las necesidades establecidas en las descripciones del TFG, quedando reflejadas en la sección de introducción de la memoria.

A continuación se resumen tanto los requisitos como los objetivos.

### **B.2. Objetivos del proyecto**

Entre los principales objetivos del proyecto se encuentran la realización de un análisis de los datos proporcionados de forma que se pueda facilitar la comprensión de los datos recogidos mediante representaciones gráficas, buscando, de esta forma correlaciones y para poder encontrar modelos predictivos para anticipar la evolución futura de las diferentes variables en un plazo de tiempo determinado.

### **B.3. Requisitos del proyecto**

Entre los requisitos del proyecto se encontraría la creación de modelos adecuados a los datos que predigan correctamente las diferentes variables objetivo.



## *Apéndice C*

---

# Comprensión de los datos

---

## C.1. Introducción

En la etapa de comprensión de los datos se creará el conjunto inicial y se comprobará si este es adecuado, para, en caso contrario, poder seguir recopilando.

Los plazos del proyecto impiden realizar esta fase en toda su extensión, puesto que a pesar de ser el proceso de análisis de datos un proceso iterativo, se encuentra sujeto a unos plazos inamovibles.

## C.2. Descripción de los datos

Se proporcionan un conjunto de datos compuesto por varios sensores (en concreto 8 sensores *IoT*) y un pluviómetro desplegados en un viñedo.

Los datos de los sensores proporcionados cuentan con los siguientes atributos:

Atributo	Descripción	Medida
ts	<i>Timestamp</i> de la muestra	ms
fecha	Fecha de la muestra	Fecha en formato yy/m-m/dd, hh:mm:ss
batería	Nivel de batería del sensor	V
t_ext	Temperatura ambiental	°C
h_ext	Humedad ambiental	%
t_C_cal	Temperatura superficial a 20cm.	°C
h_C_cal	Humedad superficial a 20cm.	%
t_L_cal	Temperatura a profundidad mayor a t_C_cal	°C
h_L_cal	Humedad superficial a profundiad mayor a h_C_cal	%
h_C	Humedad superficial sin calibrar	Valor relacionado con capacitancias
h_L	Humedad profunda sin calibrar	Valor relacionado con capacitancias

Tabla C.1: Atributos de los datos: sensores

Los datos del pluviómetro cuentan con los siguientes atributos:

Atributo	Descripción	Medida
ts	<i>Timestamp</i> de la muestra	ms
fecha	Fecha de la muestra	Fecha en formato yy/m-m/dd, hh:mm:ss
batería	Nivel de batería del sensor	V
pluv	Contador del incremento de activación del balancín	N/A
pluv_delta	Incremento de activación del balancín	N/A
pluv_deltaMM	Precipitaciones registradas desde la última lectura	litros/ $m^2$

Tabla C.2: Atributos de los datos: pluviómetro

## Visualización de los datos

Para dar una idea del estado de los datos se proporcionan unas gráficas de dispersión con las lecturas realizadas en cada uno de los sensores y el pluviómetro. En estas podrán observarse los valores **sin procesar** de cada atributo en función del tiempo.

En el caso de los sensores, los atributos  $h\_C$  y  $h\_L$  se obviarán debido a que se tratan de datos redundantes con las respectivas humedades calibradas. Además, los datos de batería y fecha tampoco se emplearán en la creación de las gráficas.

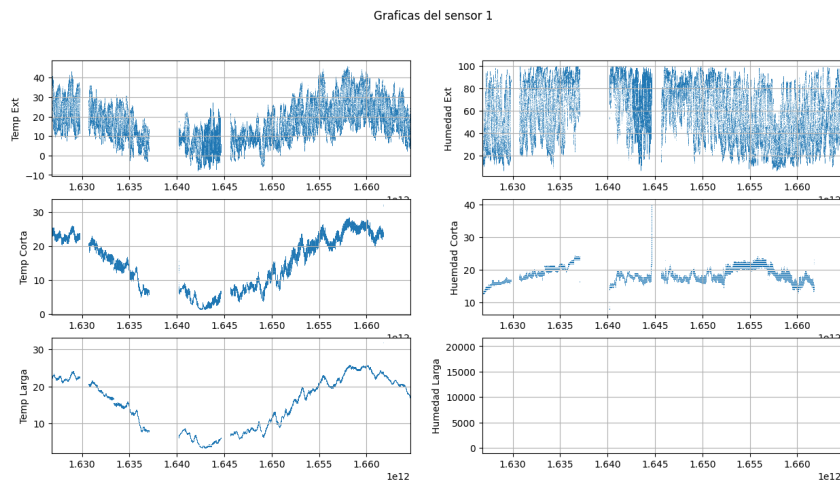


Figura C.1: Datos no procesados: sensor 1

Gráficas del sensor 2

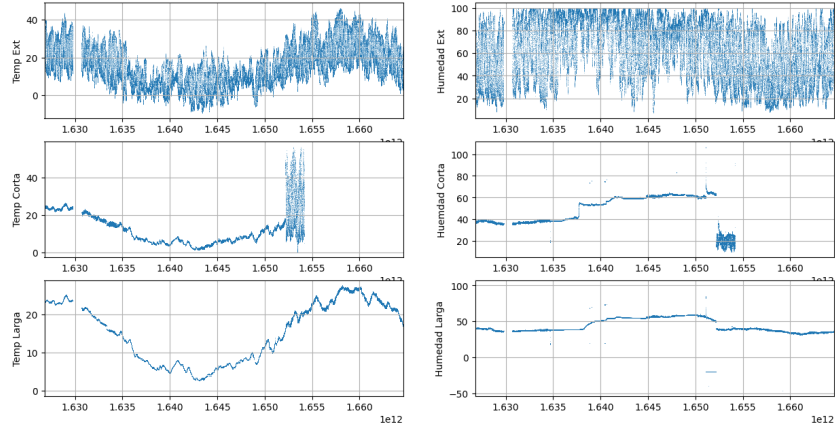


Figura C.2: Datos no procesados: sensor 2

Gráficas del sensor 3

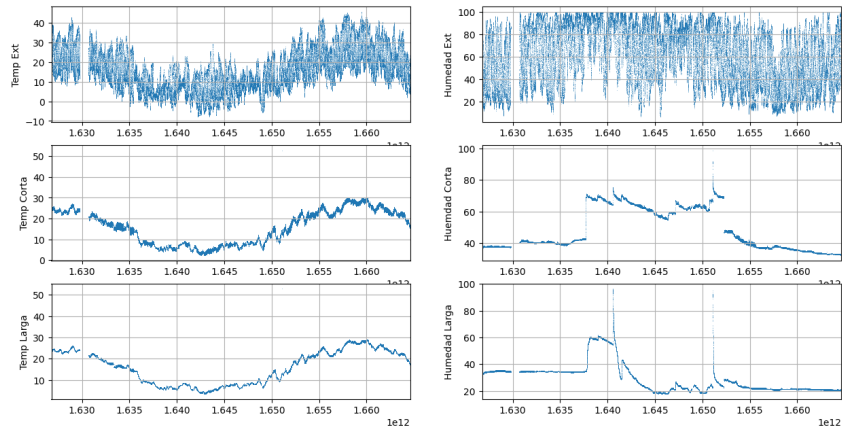


Figura C.3: Datos no procesados: sensor 3

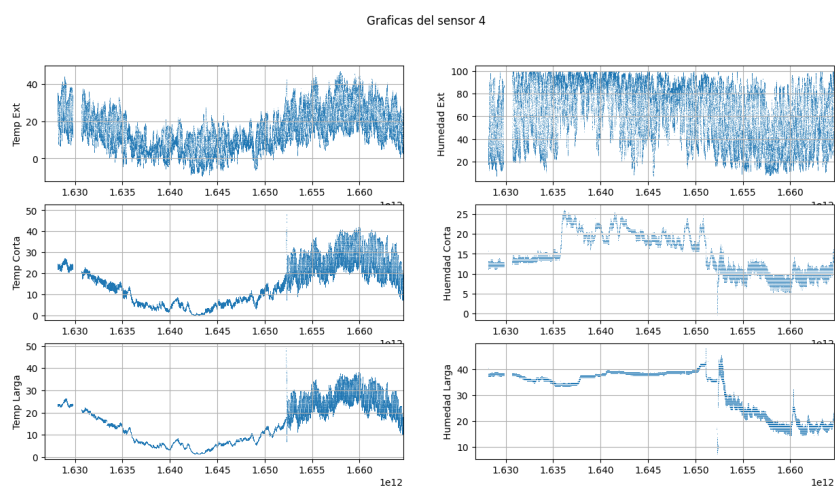


Figura C.4: Datos no procesados: sensor 4

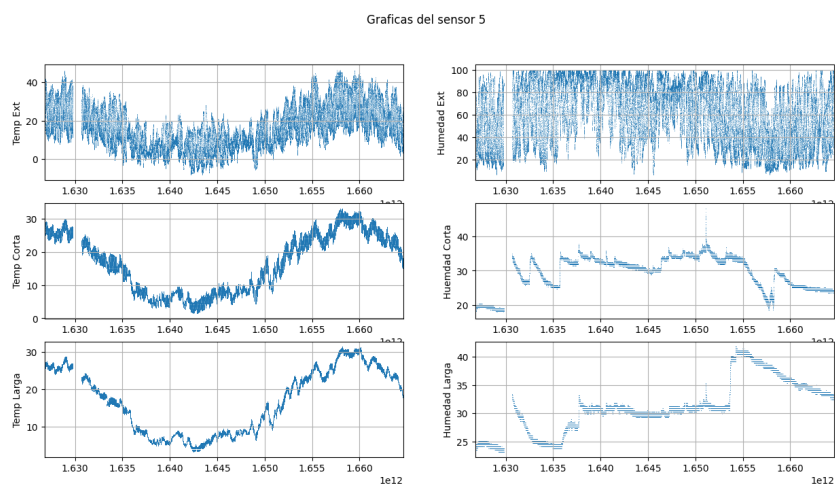


Figura C.5: Datos no procesados: sensor 5

Gráficas del sensor 6

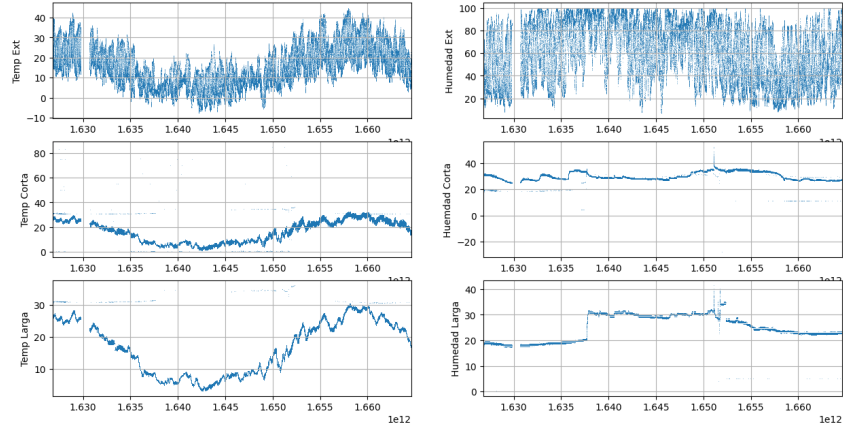


Figura C.6: Datos no procesados: sensor 6

Gráficas del sensor 7

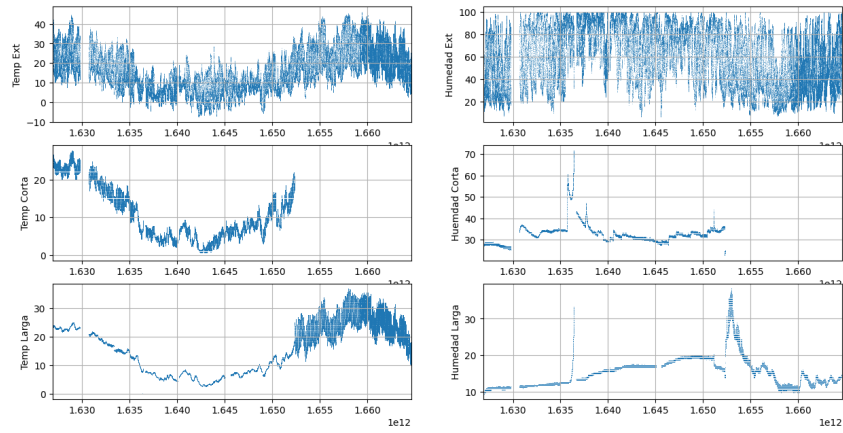


Figura C.7: Datos no procesados: sensor 7



Gráficas del sensor 8

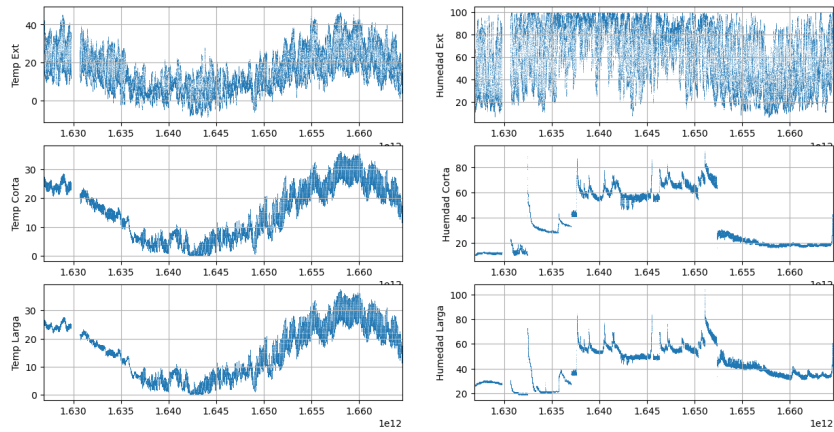


Figura C.8: Datos no procesados: sensor 8

Gráficas del pluviometro

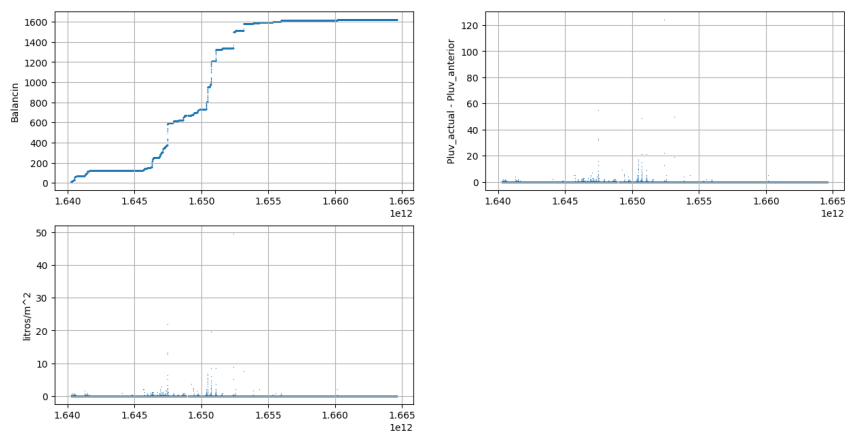


Figura C.9: Datos no procesados: pluviómetro



## *Apéndice D*

---

# **Preparación de los datos**

---

## **D.1. Introducción**

En la fase de preparación de los datos se realiza la limpieza de los mismos para que posteriormente puedan ser empleados por los modelos y obtener resultados adecuados para ser empleados en un posterior despliegue.

Esta fase está compuesta por la detección y el tratamiento de varias casuísticas en los datos, como los valores faltantes, que se tratan de valores de los que no tenemos registros por diferentes motivos y los valores extremos, que son aquellos estadísticamente lejanos al resto, lo que no quiere decir que sean erróneos.

Por otro lado, también se pueden encontrar valores erróneos, que pueden ser estadísticamente correctos, pero son datos que no tienen sentido contextual, en este caso concreto se tratará de subidas o bajadas inverosímiles de temperaturas y/o humedades del suelo, que se encuentran dentro de los valores naturales, pero que no se pueden producir a la velocidad a la que tienen lugar. En este caso, se empleará una exploración visual para realizar su tratamiento.

## **D.2. Valores faltantes**

En cuanto a los valores faltantes, la decisión fue eliminar las muestras que tuvieran indeterminados en alguna de sus columnas, puesto que la recuperación empleando técnicas como la imputación de valores era prácticamente inviable y por otro lado, se contaban con ejemplos suficientes como para poder suprimir parte de estos.

### D.3. Valores extremos

En cuanto a la detección y el tratamiento de los valores faltantes, se empleó el rango intercuartílico para realizar la primera de las tareas, sustituyendo las variables del atributo concreto de los ejemplos detectados con la mediana diaria.

De esta manera se eliminaba el ruido presente en alguno de los sensores, a la par que no se producían cambios bruscos en el conjunto de datos gracias a emplear la mediana del grupo.

### D.4. Valores erróneos

En lo referente a los valores erróneos se trata de una de las tareas más arduas del proyecto, puesto que parte de este se desarrolló de forma visual directamente sobre el conjunto de datos de los diferentes sensores y el pluviómetro, siendo, de esta manera, en su mayoría un proceso manual.

Para tal propósito se estableció una columna adicional al conjunto de datos original que indicaba la validez de la muestra concreta.

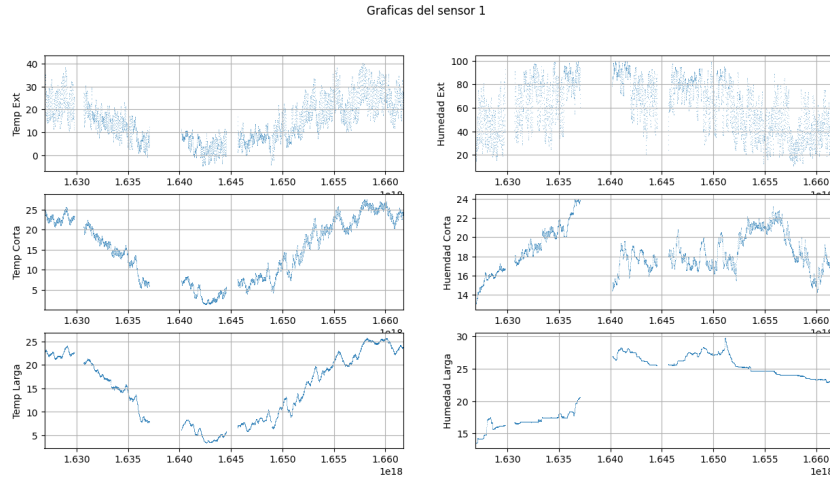


Figura D.1: Datos procesados: sensor 1

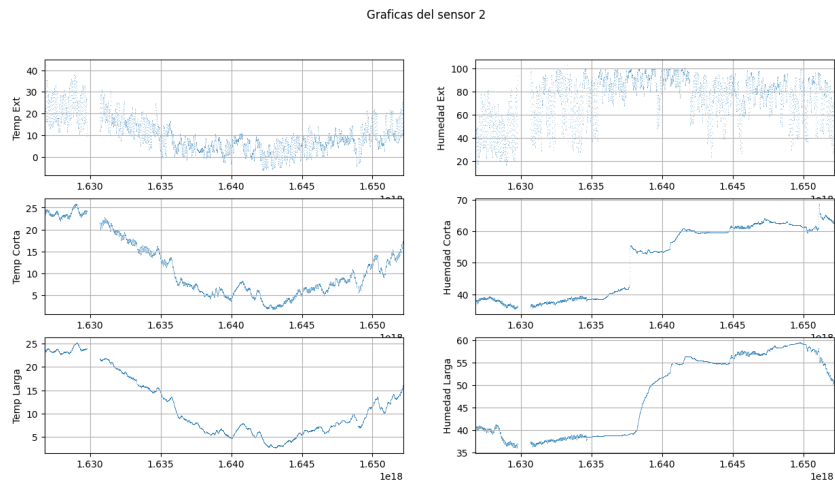


Figura D.2: Datos procesados: sensor 2

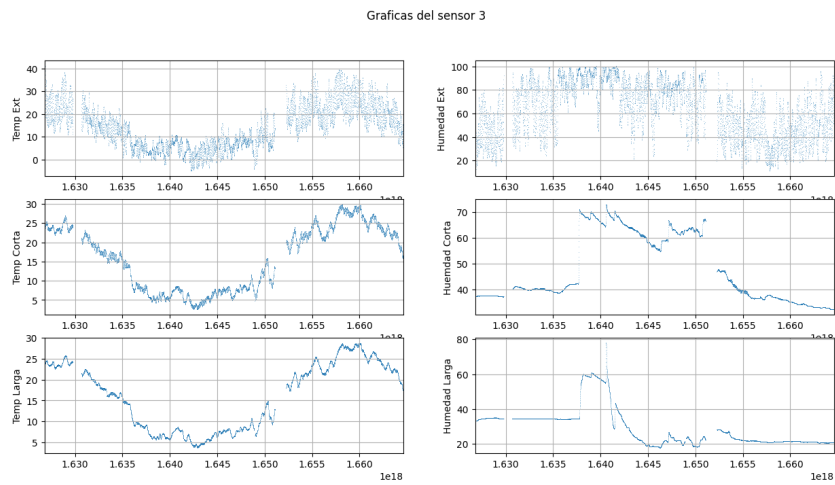


Figura D.3: Datos procesados: sensor 3

Gráficas del sensor 4

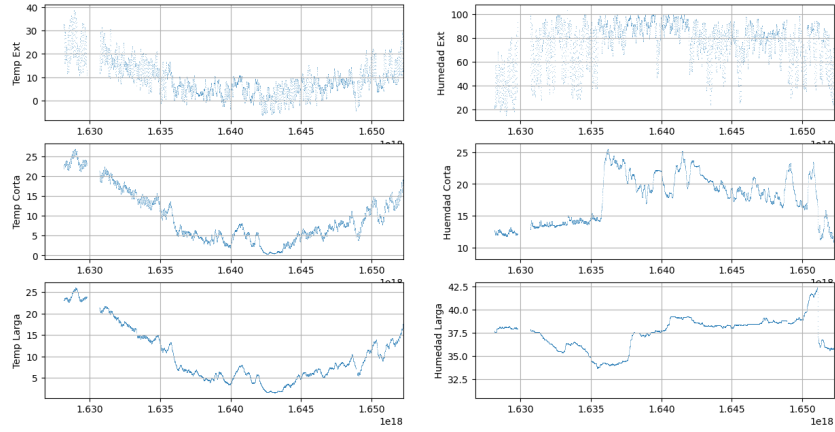


Figura D.4: Datos procesados: sensor 4

Gráficas del sensor 5

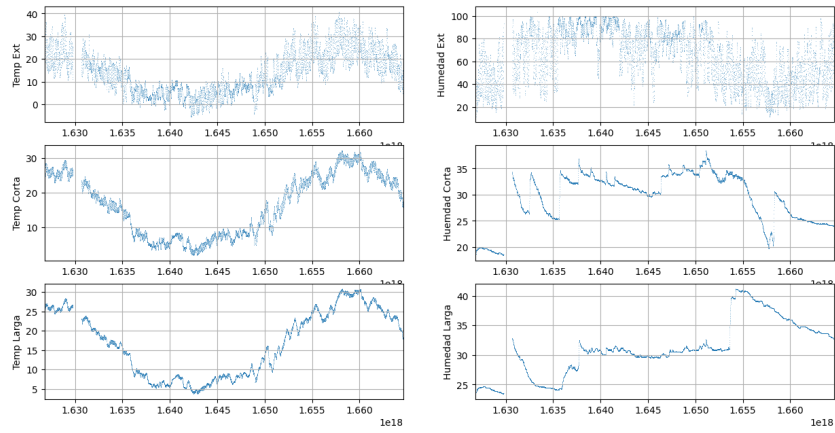


Figura D.5: Datos procesados: sensor 5

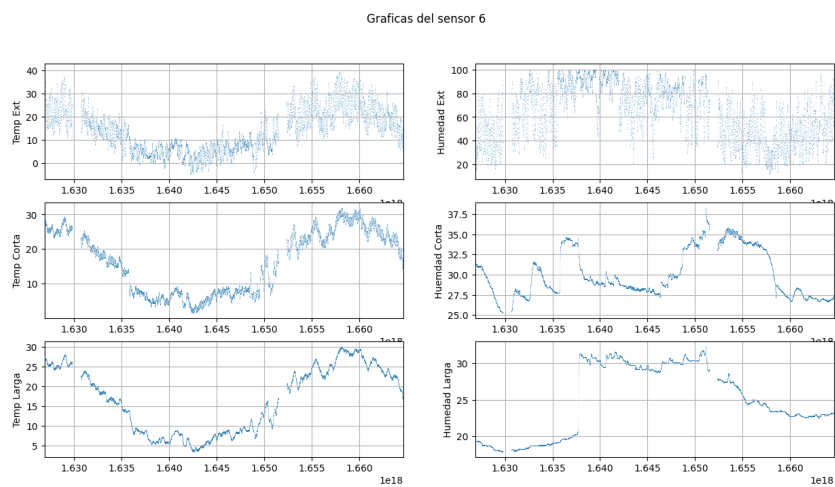


Figura D.6: Datos procesados: sensor 6

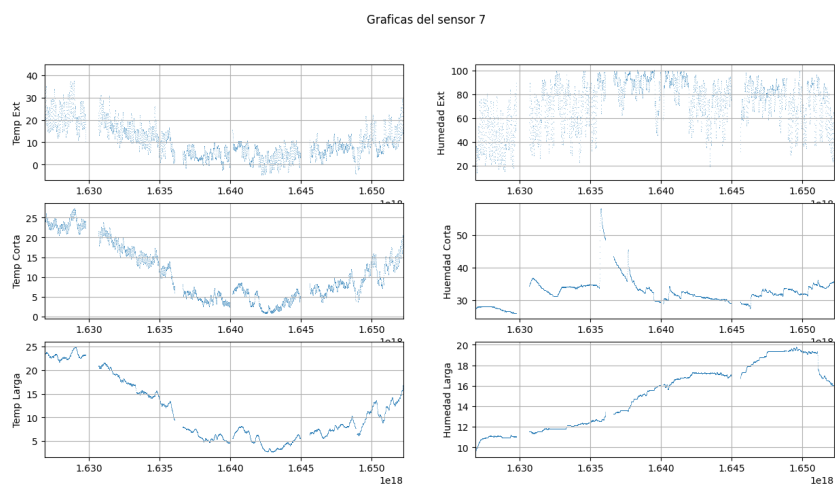


Figura D.7: Datos procesados: sensor 7

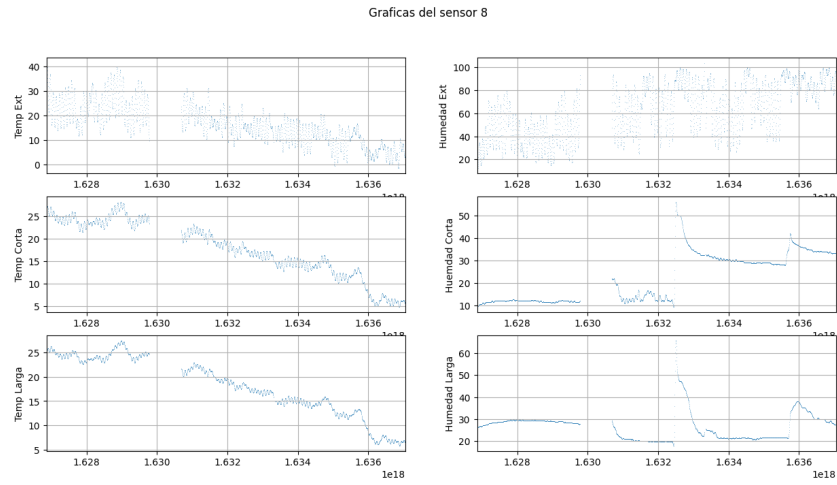


Figura D.8: Datos procesados: sensor 8

Hablándo en términos de sesgos, es posible que al realizar la recuperación de estos valores en los sensores mencionados, puedan haberse sesgado en cierta medida los conjuntos de datos, sin embargo, los datos no dejarán de estar más cercanos a los valores reales y, por tanto, contener un error más reducido que las muestras originales que contenían esta casuística.



## *Apéndice E*

---

# **Modelado**

---

### **E.1. Introducción**



## *Apéndice $F$*

---

# **Evaluación**

---

### **F.1. Introducción**



## *Apéndice G*

---

# Documentación técnica de programación

---

## G.1. Introducción

En esta sección se incluyen la documentación técnica del programador, incluyendo la estructura de directorios del proyecto, junto con el manual para realizar la correcta instalación y ejecución del mismo.

## G.2. Estructura de directorios

El proyecto cuenta con la siguiente estructura de directorios:

- **/data/**: directorio que contiene los diferentes datos del proyecto, tanto procesados, sin procesar y los datos integrados que se emplearán en el modelado.
  - /data/raw/**: directorio con los datos sin procesar (únicamente con la selección previa de validez).
  - /data/processed/**: directorio con los datos procesados.
  - /data/integrated/**: directorio con los datos integrados en un único fichero.
- **/img/graphics/**: directorio con las diferentes gráficas resultado de la ejecución de los scripts de graficado.

- **/scripts/**: directorio con los scripts para la instalación de los entornos virtuales de Python junto con los requerimientos para ejecutar todos los ficheros fuente del proyecto.
- **/src/**: directorio con los diferentes ficheros fuente y variables de entorno y globales.
- **/models/**: directorio con los diferentes modelos obtenidos en el proceso final. Uno subdirectorio para cada diferente modelo neuronal implementado.

### G.3. Manual del programador

En esta subsección se explicará cómo realizar una correcta descarga e instalación de los entornos necesarios para llevar a cabo la ejecución del proyecto.

Para descargar todo el contenido es necesario tener instalado en el sistema **Git**. Es posible clonar el repositorio introduciendo en la consola de git: **git clone <https://github.com/GabiHV/TFG22-23>**

De igual forma, para poder llevar a cabo la ejecución e instalación del resto de las dependencias es necesario tener instalado Python 3.9.13.

Para instalar el intérprete del lenguaje empleado en el proyecto es necesario acudir a la página web oficial de los desarrolladores e instalar el ejecutable de instalación oficial. La instalación puede realizarse en el siguiente enlace [5]. En la fuente mencionada se pueden escoger diferentes formas de instalación. Dependiendo del sistema operativo instalado en la máquina en la que se ejecutará el proyecto se debe seleccionar una u otra y seguir los pasos establecidos.

Durante el desarrollo del proyecto se empleó como entorno de programación Visual Studio Code [6], sin embargo para su ejecución podemos emplear otros entornos como Anaconda Navigator [7]. Se explicará la ejecución con el editor mencionado, puesto que simplifica el trabajo al disponer de scripts que realizan de forma automática la instalación de las dependencias. Los ficheros mencionados se encuentran en el directorio **/scripts/**.

Para ejecutar el script correspondiente al entorno de PowerShell de Windows se necesita establecer la política que permita ejecutarlo. Para ello se debe abrir la terminal mencionada como administradores del sistema e introducir:

```
Set-ExecutionPolicy Unrestricted
```

Tras esto, se puede introducir para iniciar el proceso:

```
./Virtual_env.ps1
```

Para ejecutar el script en el CMD de Windows se introduce:

```
virtual_env.bat
```

De forma similar en Linux Bash:

```
chmod +x virtual_env.sh && ./virtual_env.sh
```

Una vez finalice el proceso de instalación de todas las dependencias se podrá ejecutar los diferentes ficheros fuente de Python Notebook abriendo el proyecto en Visual Studio Code y estableciendo el Kernel de ejecución al entorno configurado. Está definido que el entorno virtual se denomine **.venv**, por lo que será necesario buscar entre los diferentes instalados haciendo click en la parte superior derecha del notebook (en el botón para la selección del intérprete de ejecución).

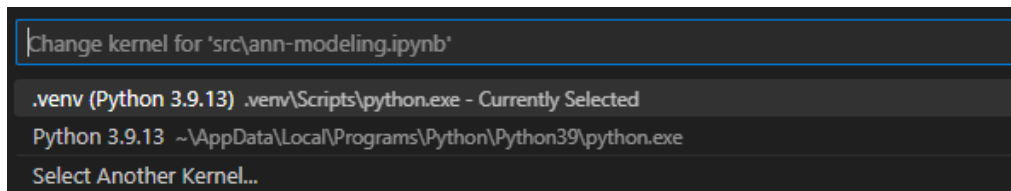


Figura G.1: Búsqueda del entorno virtual

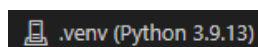


Figura G.2: Selección del entorno virtual

Posteriormente podrá ejecutarse cualquier fichero Python Notebook en el botón de “Execute All”.

En lo referente a los ficheros fuente de Python, con el entorno virtual instalado abriendo una consola en el sistema operativo en **/.venv/bin/** en Linux y **/.venv/Scripts/** en Windows, podremos ejecutar la activación del entorno virtual con los scripts incluidos en el directorio (ejecutando

**activate.bat** y **Activate.ps1** en Windows dependiendo del terminal empleado). Esta operación se realizará automáticamente al realizar la instalación de los módulos de Python incluidos en el fichero de requerimientos, por lo que se puede aprovechar la terminal en ejecución para este propósito.

## G.4. Pruebas del sistema

En esta subsección se presentará la forma de realizar las modificaciones en los hiperparámetros de los modelos, de forma que estos puedan variar de acuerdo a los nuevos requerimientos introducidos.

Los modelos contienen los siguientes hiperparámetros:

- **learning\_rate:** ratio de aprendizaje empleado en la variación de los pesos en los modelos neuronales.
- **batch\_size:** tamaño del conjunto de datos que se emplea en una única iteración en el proceso de aprendizaje.
- **epochs:** cantidad de épocas que se entrenará cada modelo.
- **window\_size\_inputs:** tamaño de la ventana de datos que se introduce como datos de entrada al modelo (se corresponde con el número de horas previas para realizar X predicciones).
- **window\_size\_targets:** tamaño de la ventana de datos que se emplean como datos a predecir.
- **train\_frac:** fracción del conjunto total de datos que se empleará para entrenar los modelos.
- **val\_frac:** fracción del conjunto total de datos que se empleará para validar los modelos, siendo el  $1 - \text{train\_frac} - \text{val\_frac}$  la fracción del conjunto de test.

Por cada uno de los diferentes modelos se proporcionará una gráfica del error de entrenamiento y validación durante el proceso de entrenamiento de la red correspondiente, además de las gráficas comparativas de los valores predichos y los reales por sensor y atributo, así como el error máximo total para todos los sensores en cada uno de estos, para dar una idea de los valores en los que ronda el error en cada una de las variables.



## Apéndice *H*

---

# Documentación de usuario

---

### H.1. Introducción

En esta sección se presentará la forma de cargar los modelos resultantes en un fichero Python para poder ser desplegados en un producto software, así como la forma que deberá tener el tensor de entrada a la red y la que tendrán los datos de salida.

### H.2. Requisitos de usuarios

En cuanto a los requisitos del usuario, se deberá tener instalado **Python 3.9.13** [5], junto con la versión 2.11.0 de **TensorFlow** [8], de forma que se puedan cargar los modelos almacenados empleando la función de la *API Keras* correspondiente.

### H.3. Instalación

En cuanto a la instalación de la versión concreta del intérprete de Python puede realizarse en [5], mientras que para instalar la dependencia de la *API*, se puede introducir el comando:

```
pip3 install tensorflow=2.11.0
```

Para cargar un modelo con la librería mencionada, se debe emplear [9]:

```
from tensorflow import keras
keras.models.load_model('<path_del_modelo>')
```

## H.4. Manual del usuario

Para realizar predicciones con el modelo pertinente, las entradas deben tener una estructura concreta que dependerá de cómo se haya entrenado a cada una de las diferentes redes neuronales. Es decir, el tamaño del número de muestras de entrada dependerá de la cantidad de “*backtracking*” que se haya establecido en el entrenamiento.

En este caso se debe introducir un tensor bidimensional de 6 muestras (se han obtenido modelos que aceptan 6 horas) con 7 variables de entrada cada una que se corresponde a:

- **t\_ext:** temperatura exterior media en una hora (-50, 50).
- **h\_ext:** humedad exterior media en una hora (0, 100).
- **t\_C\_cal:** temperatura media de la sonda de temperatura más superficial en una hora (-50, 50).
- **h\_C\_cal:** humedad media de la sonda de humedad más superficial en una hora (0, 100).
- **t\_L\_cal:** temperatura media de la sonda de temperatura interna en una hora (-50, 50).
- **h\_L\_cal:** humedad media de la sonda de humedad interna en una hora (0, 100).
- **sensor:** sensor al que se corresponde los datos (0, 7).

Por otro lado, las variables deberán estar normalizadas en el rango 0-1, con los máximos y mínimos especificados anteriormente.

En el caso del modelo *MLP*, en lugar de un tensor tridimensional, se deberá modificar para que se corresponda con un vector del número de variables por el de muestras.

La salida será igualmente un tensor bidimensional del número de predicciones establecidas en el entrenamiento del modelo con 6 variables cada una que se corresponden a las mencionadas anteriormente a excepción del número de sensor. De forma inversa, el modelo proporcionará datos normalizados, por lo que para obtener cada atributo en un rango correcto deberá de denormalizarse.

---

## Bibliografía

---

- [1] M. Palacio, *Scrum Master*. Iubaris Info 4 Media SL, 2021, ch. 1, p. 11.
- [2] “Zenhub.” [Online]. Available: <https://www.zenhub.com/>
- [3] “Seguridad social: Cotización / recaudación de trabajadores.” [Online]. Available: <https://www.seg-social.es/wps/portal/wss/internet/Trabajadores/CotizacionRecaudacionTrabajadores/36537#36538>
- [4] “Welcome to the apache software foundation!” [Online]. Available: <https://www.apache.org/licenses/LICENSE-2.0>
- [5] “Python.” [Online]. Available: <https://www.python.org/downloads/release/python-3913/>
- [6] Microsoft, “Visual studio code - code editing. redefined,” Nov 2021. [Online]. Available: <https://code.visualstudio.com/>
- [7] “Anaconda navigator.” [Online]. Available: <https://anaconda.org/anaconda/anaconda-navigator>
- [8] “Tensorflow.” [Online]. Available: <https://www.tensorflow.org/?hl=es-419>
- [9] “Guardando y serializando modelos con tensorflow keras | tensorflow core.” [Online]. Available: [https://www.tensorflow.org/guide/keras/save\\_and\\_serialize?hl=es-419](https://www.tensorflow.org/guide/keras/save_and_serialize?hl=es-419)