

# ÍNDICE

---

- [1 RESUMEN](#)
- [2 ABSTRACT](#)
- [3 INTRODUCCIÓN](#)
  - [3.1 Contexto](#)
  - [3.2 Previos](#)
- [4 PLANIFICACIÓN](#)
  - [4.1 Alcance](#)
  - [4.2 Metodología](#)
  - [4.3 EDT](#)
  - [4.4 Entregables](#)
  - [4.5 Gantt](#)
  - [4.6 Comunicaciones](#)
  - [4.7 Riesgos](#)
- [5 ANÁLISIS](#)
  - [5.1 Dominio agronómico](#)
  - [5.2 Fuentes de datos](#)
  - [5.3 Automatizacion](#)
  - [5.4 Tecnologías](#)
  - [5.5 Requisitos](#)
- [6 INFRAESTRUCTURA](#)
  - [6.1 Arquitectura física](#)
  - [6.2 Vagrant](#)
  - [6.3 Entorno](#)
- [7 DISEÑO](#)
  - [7.1 Carga de datos](#)
  - [7.2 Modelo relacional](#)
  - [7.3 Servicios](#)
- [8 PIPELINE](#)
  - [8.1 Transformación de datos](#)
  - [8.2 Design pattern](#)
  - [8.3 Vista minable](#)
  - [8.4 Integración continua](#)
- [9 MODELO IA](#)
  - [9.1 Preprocesamiento](#)
  - [9.2 Regresión](#)
  - [9.3 Mejor modelo](#)
  - [9.4 Predicciones](#)
- [10 SEGUIMIENTO](#)
  - [10.1 Desviaciones](#)
  - [10.2 Memoria](#)
- [11 CONCLUSIONES](#)
- [12 BIBLIOGRAFÍA](#)

# **1 RESUMEN**

---

El trabajo consistirá en la reestructuración de una aplicación de gestión de datos agronómicos. Estos datos, que actualmente provienen de varias fuentes, van a ser tratados en una nueva arquitectura que posibilitará técnicas de integración y despliegue continuo. En su diseño pasamos por las etapas sucesivas de infraestructura, diseño de un modelo relacional, utilización de un pipeline para la carga de datos y, finalmente, la obtención de una vista minable sobre la que se ejecutarán modelos de Inteligencia Artificial. Esta vista permitirá predecir mediante regresión la cantidad de kg de cultivo que se espera cosechar en una parcela y en una fecha determinadas.

## **2 ABSTRACT**

---

The following work will handle the restructuring of an agronomic data management application. These data, which currently come from various sources, are going to be processed in a new architecture that will enable continuous integration and deployment techniques. In its design, we go through the successive stages of infrastructure, design of a relational model, use of a pipeline for data loading and, finally, obtaining a mineable view on which Artificial Intelligence models will be executed. This view will make it possible to predict, with a regression model, the amount of kilograms of crop that is expected to be harvested in a plot on a given date.

## 3 INTRODUCCIÓN

---

### 3.1 Contexto

---

El propósito de este proyecto es mejorar el flujo de trabajo del equipo de "SpectralGeo" mediante la automatización del proceso de recogida y procesado de datos agronómicos que se utilizan en la creación de modelos de inteligencia artificial para el cultivo.

AGRAI es una aplicación para la gestión de cultivo que permite al agricultor monitorizar el estado de sus parcelas. El estado actual de la aplicación se centra el despliegue de datos agronómicos y vegetativos georreferenciados a través de en una interfaz web. Los usuarios de la aplicación, normalmente agricultores o cooperativas, pueden consultar el estado de su parcelario junto con algunas predicciones, como puede ser la cantidad de Kg que se van a cosechar en una fecha determinada. El acceso a dicha aplicación puede darse desde equipo de sobremesa o un dispositivo móvil, aunque es este último lo que parece que se utiliza más.



*Imagen del frontend de la aplicación en la que se muestran, mediante un mapeo con colores, la evolución de los índices vegetativos en varias parcelas agrícolas.*

Actualmente, cómo se procesa la información que utiliza nuestra aplicación es un proceso tedioso para el equipo. En este proceso, diferentes miembros trabajan con tecnologías distintas sobre datos duplicados provenientes de fuentes comunes. Aunque se realiza una planificación y coordinación de los proyectos, se pierde bastante tiempo en la transformación de los datos que cada miembro del equipo necesita para llevar a cabo su labor.

El objetivo de este trabajo es mejorar y automatizar el flujo de trabajo del equipo, convirtiendo la aplicación AGRAI en una herramienta robusta que manipule un único repositorio de datos al cual el resto del equipo pueda acceder, utilizándolo de forma segura y -lo más importante,-sin duplicar y romper la integridad de estos datos.

Es un proyecto ambicioso debido a que no solo es importante el conocimiento técnico sino que serán necesarios cambios en la forma de trabajo del equipo, cambios que se basan en la confianza de cada miembro en la nueva forma de trabajo que se desea implementar.

## 3.2 Previos

---

SpectralGeo se ha especializado en el uso de nuevas tecnologías para sectores como el de la agricultura o el reciclaje, centrándose en proyectos con carácter reivindicativo de la sostenibilidad medioambiental. La relación con clientes como Ecoembes o el desarrollo de software para la gestión sostenible de cultivos lo demuestran.

Durante las prácticas realizadas en la empresa se identificó que era totalmente necesario dotar de una arquitectura robusta a la aplicación AGRAI. En estas prácticas se comenzó el desarrollo de un modelo de datos sobre el que dicha aplicación pudiese escalarse posteriormente.

Actualmente, AGRAI aún necesita una reestructuración para la gestión de grandes volúmenes de datos agronómicos, poniendo especial énfasis en los índices vegetativos que gestiona, los cuales provienen del procesamiento de imágenes satelitales. El estado actual de dicha aplicación consiste en la monitorización de cultivos a través del análisis de históricos de datos.

La aplicación ya está creada; el cliente obtiene los resultados que espera cuando consulta el estado de su parcelario en la interfaz de la aplicación. Como equipo, nos organizamos para que la información y los modelos predictivos lleguen al cliente a través de esta, pero son varios los puntos en los que trabajamos en exceso para finalmente mostrar una "predicción" al cliente.

El valor más importante que predecimos es el número de kg que se espera cosechar en la fecha estipulada para el agricultor. Esta predicción hace uso de valores de índices vegetativos a partir de imágenes satelitales tomadas en diferentes fechas. El histórico de datos registrado permite estudiar la evolución del cultivo y, finalmente, mediante métodos de Inteligencia Artificial IA, predecir un valor en la fecha en la que se espera cosechar.

Serán las etapas de infraestructura, diseño, implementación y modelado las que buscaremos optimizar en el contexto de todo el equipo, debido a que estas etapas se reparten entre sus miembros.

Justificamos este cambio como parte de una búsqueda de procesos esbeltos (lean), que consisten en la eliminación de los "desperdicios", o fuentes de despilfarro de tiempo y trabajo en la elaboración de productos o servicios. A través de la solución que vamos a implantar buscamos la optimización continua del proceso y la aceptación de esta cultura de optimización por el equipo.

Los puntos más importantes de los procesos "lean" son los siguientes:

- identificar los desperdicios y tratar de eliminarlos
- mejorar la comunicación interna de la organización
- reducir costes y tiempos de entrega y mejorar la calidad

## 4 PLANIFICACIÓN

---

### 4.1 Alcance

---

A continuación se expondrán los objetivos generales del proyecto para aportar una visión global del trabajo. Acotamos en la medida de lo posible nuestros objetivos con la aplicación.

- Automatizar y asilar las dependencias necesarias mediante infraestructura como código para que el equipo pueda trabajar sobre el mismo entorno.
- Identificar los conceptos del dominio agronómico y de índices vegetativos susceptibles de ser implementados.
- Diseñar un modelo de datos que permita un rápido escalado de dicha aplicación mediante el que podamos manejar grandes volúmenes de datos.
- Buscar el proceso de trasformación de los datos mediante el diseño de un pipeline que posibilite la integración continua.
- Diseñar un algoritmo que exporte la vista minable necesaria para predecir la producción de un cultivo (kg).
- Probar técnicas de Inteligencia Artificial para entrenar un modelo de regresión con el que predecir nuevos valores que no han sido usados en el entrenamiento.

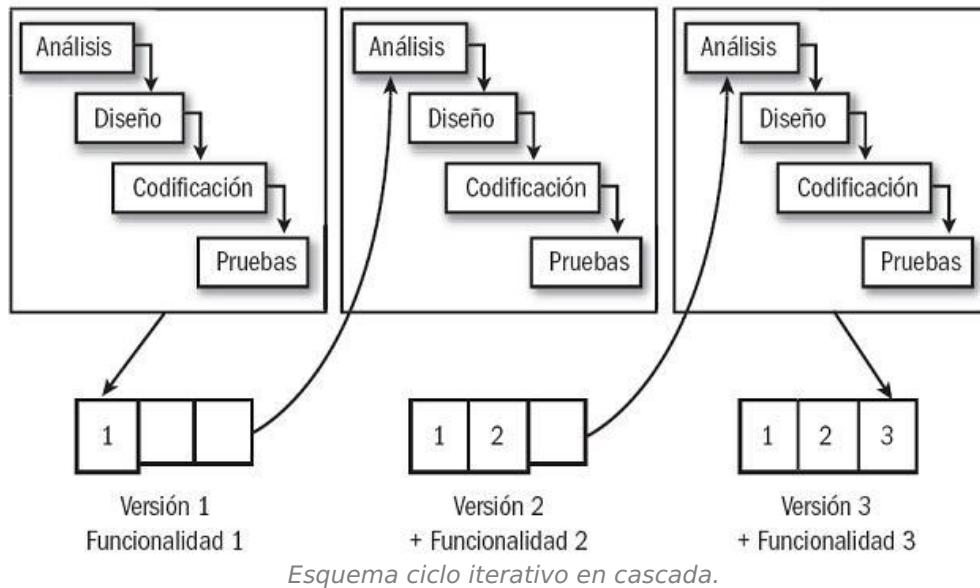
En los próximos puntos se detallan estas ideas iniciales para dar forma a los requisitos que buscaremos realizar a lo largo del trabajo.

### 4.2 Metodología

---

En el equipo hacemos uso de SCRUM para la gestión del proyecto. Dicha metodología permite un desarrollo en cascada de las diferentes tareas propuestas para la creación del pipeline.

Se usará un ciclo de vida iterativo e incremental. Dentro de las fases del proyecto (también llamadas iteraciones), se repiten de manera intencionada una o más actividades del proyecto. Con estas iteraciones el entendimiento del producto por parte del equipo va aumentando. Las fases (iteraciones) desarrollan el producto a través de una serie de ciclos repetidos, mientras que los incrementos van añadiendo sucesivamente funcionalidad al producto. Esto, en definitiva, consiste en varios ciclos de vida en cascada. Al final de cada iteración se entrega una versión mejorada. Las versiones que entregamos se corresponden con las cuatro grandes fases de la arquitectura que desarrollamos.



### 4.3 EDT

Mostramos el desglose de las tareas más importantes para la correcta realización del proyecto. Para poder realizar un posterior seguimiento y control del proyecto asignamos a cada tarea un tiempo adecuado a la complejidad que estimamos para ésta. En caso de que posteriormente nos alejemos de lo que aquí planificamos, anotaremos dichas desviaciones en la sección de 'seguimiento y control'

Tarea	Descripción	Horas
<b>PLANIFICACIÓN</b>		20
INTRODUCCIÓN	Explicamos el contexto del problema planteado y cómo vamos a implementar la solución.	2
ALCANCE	Qué pretendemos conseguir con el proyecto acontando las tareas	2
RECURSOS HUMANOS	Aquellas personas que intervienen y la explicación de sus funciones	2
COMUNICACIONES	Elementos de comunicación interna y externa, incluimos interesados y clientes.	2
METODOLOGIA	Conjunto de procedimientos que usaremos para la planificación y gestión.	2
EDT Y TIEMPOS	Estructura en árbol con el desglose en tareas del trabajo junto con su estimación en tiempos	4
ENTREGABLES	Creación de los paquetes de trabajo del proyecto	2
GANTT	Tiempo determinado en el calendario para cada uno de los paquetes de trabajo.	2
PLAN DE RIESGOS	Determinaremos los principales riesgos que pueden suceder durante el proyecto y los categorizaremos por nivel de impacto y de suceso.	2
<b>ANÁLISIS</b>		20
CONCEPTOS PREVIOS	Explicamos el contexto del problema planteado y cómo vamos a implementar la solución.	5

Tarea	Descripción	Horas
PROYECTO AGRAI	Presentamos el estado actual de la aplicación AGRAI.	5
AUTOMATIZACIÓN	Exponemos los distintos pasos a través de los cuales se transforman los datos.	5
CATÁLOGO DE REQUISITOS	Obtendremos los requisitos funcionales y no funcionales del proyecto.	5
<b>INFRAESTRUCTURA</b>		25
ENTORNO	Decidimos los entornos de desarrollo necesarios así como las tecnologías necesarias en cada punto.	10
INFRAESTRUCTURA CÓDIGO	Desarrollamos un script para dejar la máquina que contiene la aplicación en un estado estable con todas sus dependencias instaladas.	15
<b>DISEÑO</b>		100
ANÁLISIS MODELO PREVIO	Explicamos el contexto del modelo de datos previos y sus fallos.	30
REDISEÑO ENTIDADES	Explicamos qué entidades conformarán el modelo y por qué la base de datos quedará normalizada.	30
DISEÑO DE SERVICIOS	Encapsulamos las consultas más frecuentes en servicios fáciles de acceder.	10
CARGA CON DATOS	Cargamos el modelo con datos para verificar su consistencia.	20
DESPLIEGUE DE APLICACIÓN	Desplegamos la aplicación para comprobar su funcionamiento.	10
<b>PIPELINE</b>		50
FUENTES DE DATOS	Explicamos las fuentes de datos con las que trabajamos.	15
DEFINICIÓN DE ETAPAS	Acotamos y definimos los pasos en los que se transforman los datos hasta llegar al modelo de forma consistente.	15
VISTA MINABLE	Seleccionamos del modelo relacional las 'features' para el modelo IA de producción.	20
<b>MODELO IA</b>		50
DATOS	Exponemos los datos que vamos a utilizar, procedentes de la vista anterior.	10
MODELO DE PRODUCCIÓN	Usamos varias técnicas de Inteligencia Artificial para crear un modelo de producción de Kg para el cultivo.	20
PREDICCIONES	Utilizamos el modelo anterior dentro del sistema para predecir los kg de cosecha a partir de una fecha dada.	20
<b>SEGUIMIENTO</b>		10
CONTROL	Realizamos pruebas de integración para asegurar la consistencia de los datos.	5
GESTIÓN DE TIEMPOS	En una tabla registramos el tiempo estimado, el tiempo real invertido y la desviación del tiempo en cada una de las tareas del proyecto.	5

Tarea	Descripción	Horas
<b>MEMORIA</b>		25
ESTRUCTURA	Diseñamos la estructura basada en múltiples notas con la documentación y explicaciones del proyecto.	5
REDACCIÓN	Redactamos la memoria a medida que avanza el proyecto.	15
EXPORTACIÓN	La memoria se escribe en notas de "markdown" debido a la simplicidad con la que se pueden integrar diagramas y código. Al terminarla, exportaremos la memoria como página "html".	5
REVISIÓN	Al finalizar el proyecto, revisamos el escrito para asegurarnos de no cometer errores.	5
<b>TOTAL</b>		300

## 4.4 Entregables

---

Para acotar el proceso de trabajo identificamos los siguientes entregables que generaremos a lo largo del desarrollo del proyecto. Estos entregables quedarán en la memoria en sus correspondientes apartados o como anexos si tienen una extensión más larga. La siguiente tabla contiene que lleva a cada uno de estos artefactos.

IDENT	ENTREGABLE	DESCRIPCIÓN
E01	Módulo de Inicio	Análisis de viabilidad que permitirá determinar si es posible desarrollar el proyecto.
E02	Planificación del proyecto	Determinaremos los requisitos del proyecto, crearemos un enunciado para el alcance, realizaremos una descripción detallada de tareas junto a una EDT, estimaremos el tiempo y coste de los paquetes de trabajo, desarrollaremos un cronograma, estableceremos los estándares, procesos y métricas de calidad, determinaremos un plan de identificación de riesgos y crearemos el plan de gestión de cambios.
E03	Infraestructura como código	Desarrollo de un 'script' con las instrucciones bash para dejar la máquina que despliega la aplicación en un estado estable.
E04	Documento de Diseño	Documento en el que se explican las decisiones tomadas para la creación de las entidades del modelo y sus relaciones.
E05	Esquema modelo de Datos	Esquema UML que representa las entidades del modelo y sus relaciones. Modelo simplificado de datos extraído del ORM Django
E06	Arquitectura Pipeline	Arquitectura y diagramas del proceso de transformación de datos, junto con sus scripts. Memoria + Código
E07	Vista Minable	Método algorítmico que genera la estructura tabular de datos necesaria para el modelo de producción
E08	Datos Vista	Estructura tabulada con las 'features' necesarias y los datos requeridos. (Estudio con 25 parcelas)
E09	Modelo de producción	Cuaderno jupyter en el que buscamos el mejor modelo posible

## 4.5 Gantt

El siguiente diagrama propone un desglose de tareas para la planificación del proyecto. Este diagrama es una estipulación de cómo se desarrolla el proyecto, en el eje horizontal, medimos el tiempo en días. Considerando 4 h de trabajo diarias, obtenemos una duración de 75 días para repartir las 300 h que dedicamos al trabajo. En la sección final de seguimiento y control analizamos las desviaciones encontradas y cómo han sido solucionados estos contratiempos.

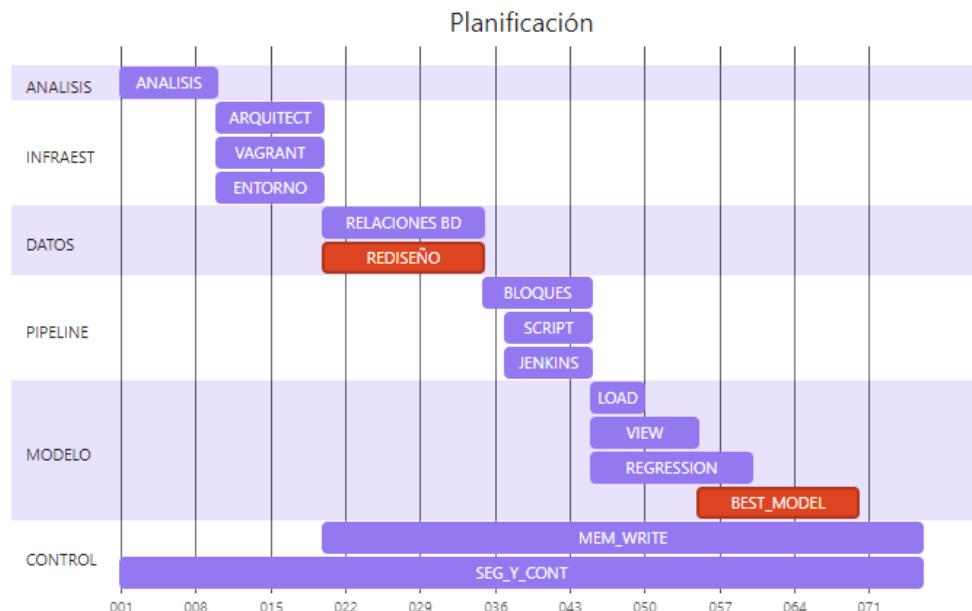


Diagrama Gantt con el desglose de tareas a lo largo de los días estipulados para la realización del proyecto.

Destacamos que las tareas de escritura de la memoria y seguimiento y control se realizan durante casi todo el proyecto. Las pruebas que realizamos durante el seguimiento al final de cada bloque, nos permiten asegurarnos de que el proyecto funciona adecuadamente: así no pasamos a una nueva tarea si no hemos dejado la aplicación estable al terminar la anterior.

## 4.6 Comunicaciones

Para mantener informados tanto al tutor académico como a la tutora de la empresa, utilizaremos los siguientes canales:

CANAL	DESCRIPCIÓN
Reuniones	Presenciales, tanto con el tutor académico como con el tutor de la empresa.
Email	Forma de comunicación electrónica para concretar determinadas pautas o establecimiento de citas.
Discord	Aplicación de comunicación interna para mensajes directos con los miembros del equipo.

## 4.7 Riesgos

El objetivo de esta tabla es aumentar la probabilidad de eventos positivos y disminuir la de los negativos.

FUENTE	RIESGO	SI SUCEDE	MINIMIZAR
Ausencia de interesados	Médio	Utilizar menos reuniones y ser independiente	Planificar reuniones con antelación
Tecnología equivocada	Alto	Estudiar nuevas posibilidades	Uso de patrones abstractos elegantes que se puedan implementar en diferentes frameworks
Pérdida código fuente	Alto	Empezar casi de cero	Backups y control de versiones
Fallo de dependencias o librerías	Alto	Volver a crear los entornos	Utilizar técnicas de infraestructura como código donde sea posible
Sobreajuste en los modelos	Medio	Probar otros algoritmos	Automatizar el proceso de búsqueda de los mejores algoritmos probando hiperparámetros.
Falta de horas para terminar el proyecto	Medio	Dedicarse primero a los requisitos más importantes.	Realizar una buena planificación y ajustar el alcance

De la tabla anterior se prestará especial atención a los riesgos valorados como 'Alto', ya que por su incidencia en el proyecto requieren un seguimiento exhaustivo, para los cuales se desarrollarán sus correspondientes planes de contingencia. Los riesgos de valoración Media y Baja tendrán un control basado en un plan de mitigación simple de seguimientos del cronograma, con intervalos cortos de tiempo.

## 5 ANÁLISIS

### 5.1 Dominio agronómico

Una vez en contexto, procedemos a justificar las decisiones de diseño que la aplicación AGRAI necesita para llegar al rediseño que hemos identificado como clave.

AGRAI es una aplicación para la gestión de cultivo que permite al agricultor monitorizar el estado de sus parcelas. El estado actual de la aplicación se centra el despliegue de datos georreferenciados por medio de una interfaz web. Los resultados para el usuario de la aplicación son satisfactorios, permitiendo a éste consultar el estado de su parcelario junto con algunas predicciones en cualquier momento.

Tenemos que situarnos en el uso que se le está dando a la aplicación en relación con el diseño de ésta. La monitorización efectuada es válida para cualquier tipo de cultivo susceptible de ser fotografiado desde una imagen satelital y/o de dron. El modelo de datos que utiliza, aunque pobre al comienzo del trabajo, permite registrar tantos cultivos como un potencial cliente deseé. Por otro lado, el caso real de uso de la aplicación proviene de bodegas o cooperativas vitivinícolas que quieren predecir la cantidad de kg de cosecha en una campaña determinada (2020-2021).

Como punto de partida buscaremos identificar las entidades del dominio necesarias para que nuestro nuevo modelo de datos represente y relacione toda la información con precisión.

En la interfaz web, el cliente observa un mapa con sus parcelas resaltadas en colores. Este color reflejado en los píxeles hace referencia a los índices vegetativos provenientes de las imágenes satelitales. Los clientes utilizan dicha interfaz para ver el estado de sus cultivos y observar la producción que predicen los modelos de Inteligencia Artificial generados.



Captura del visor web en la que se aprecia cómo se monitorizan varias parcelas

A nivel interno, el equipo trabaja con datos provenientes de imágenes descargadas de satélites como SENTINEL, o de grabaciones realizadas por dron cuando se requiere una mayor calidad. Estos datos se mezclan con información proporcionada por estaciones meteorológicas como el SIAR y con la información que los clientes pueden proporcionar sobre campañas anteriores.

Por otro lado, el proceso actual de esta información es tedioso para el equipo. Los diferentes miembros trabajan con tecnologías distintas sobre datos, muchas veces duplicados, que provienen de fuentes comunes. Los requisitos que identificamos en el siguiente punto consideran la idea de optimización continua necesaria.

## 5.2 Fuentes de datos

---

Mencionábamos la búsqueda de optimización del proceso como una de las principales ideas para mejorar el flujo de datos que utiliza la aplicación. El rediseño que se quiere implantar y la arquitectura resultante tienen que cumplir con los principios de dicha metodología. Para que este trabajo resulte satisfactorio, el software tiene que poder utilizarse por el equipo con agilidad, consiguiendo que quede como una herramienta física que coordine bien todas las tareas de los participantes que la usan.

Para poder realizar este proyecto decidimos trabajar con un número pequeño representativo de datos de la aplicación, actualmente en producción. Un menor volumen de información permite realizar pruebas y da pie a fijarnos en las relaciones y los esquemas estructurales con más precisión.

Como estamos buscando una optimización continua, a medida que voy haciendo pruebas con mi entorno local, un compañero se encarga de utilizar el mismo proyecto con volúmenes mayores de datos, que la aplicación utiliza en producción. La reestructuración del modelo de BD que realizamos contempla un posterior escalado de la aplicación que permitirá la integración continua con más datos y nuevas tecnologías.

La siguiente tabla muestra las principales fuentes de datos de las que se obtienen y enlaza la información. El volumen de datos que se puede llegar a manejar es grande ya que por cada parcela se almacenan varios índices vegetativos en cada uno de sus píxeles.

FUENTE DATOS	DESCRIPCIÓN
QGIS	Información geométrica de parcelas y sus pixéles
Índices Vegetativos	Información provenientes de imágenes satelitales descargadas en diferentes fechas
Cultivos / Variedad	Información tabulada en excels sobre tipo de cultivos y sus variedades

Para poder realizar el trabajo utilizaremos, por tanto, una muestra representativa de los datos debido a que el proceso de descarga de índices y persistencia de datos es largo para realizar las pruebas. El sistema está pensado para trabajar con muchas parcelas; en las pruebas que yo voy a realizar escogemos una muestra de 25 parcelas y pensaremos una descarga de índices para no más de 4 fechas diferentes.

Para hacernos una idea, contemplando solo 25 parcelas, podemos almacenar 3000 píxeles. Por cada pixel vamos a registrar varios índices vegetativos (ndvi, ndre) y para generar el histórico de datos esta información se multiplica por el número de fechas contempladas. Es decir, que aunque trabajamos sobre un volumen reducido para probar la automatización, sigue siendo mucha la información que tiene que almacenarse en la base de datos.

Es el histórico de índices en diferentes fechas lo que permite al sistema realizar modelos predictivos. A mayor volumen de datos, más precisión podremos obtener en los modelos posteriormente. El punto importante de este trabajo es la mejora del

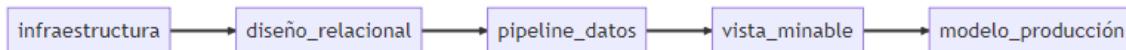
proceso y la automatización del flujo de datos; ello nos llevará a que finalmente podamos obtener un modelo de regresión para predecir los kg de cultivo, tal y como el agricultor necesita. La muestra que utilizamos, al ser representativa con 25 parcelas, predecirá con un cierto sobreajuste a los datos. Por tanto nuestro esfuerzo se centra en que el proceso que implantamos sea reproducible para cualquier volumen de datos, y en que sin desarrollo adicional podamos utilizar el mismo proceso para predecir en el volumen que queramos.

### 5.3 Automatizacion

---

El proceso que vamos a implantar requiere de un diseño que posibilite estructurar y manejar grandes volúmenes de datos. Actualmente, existen varias tecnologías que abordan el concepto de *Pipeline* desde un punto de vista global. No vamos a utilizar ninguna de estas tecnologías; nuestro proceso se va a desarrollar a medida, modelando un Pipeline como patrón de diseño, encargado de la transformación de los datos en bloques reutilizables para su persistencia en el modelo.

El primer paso para realizar una carga de datos, es el diseño de un modelo relacional que soporte esta información de la forma más estructurada y organizada posible. Hará falta la comprensión de conceptos del dominio agronómico para dar sentido a esta información y obtener así un esquema sólido sobre el que poder preguntar sin limitaciones.



Con un modelo relacional para los datos, los bloques atómicos de carga se podrán ejecutar periódicamente utilizando las fuentes de datos necesarias. El proceso de automatización termina, una vez que están almacenados los datos, con la extracción de una vista minable que permite obtener el mejor modelo de regresión para los kg de cultivo que se van a cosechar en una fecha determinada.

### 5.4 Tecnologías

---

El código de la aplicación que heredamos está escrito en Python. El lenguaje es una decisión adecuada debido a la necesidad de integrar técnicas de procesamiento de datos e inteligencia artificial. El framework de Django para Python permite construir un proyecto robusto y desplegarlo en un servidor con una interfaz web.

*Django* nos ofrece las herramientas necesarias para trabajar desde un alto nivel de abstracción y poder diseñar una aplicación sólida y estable. Algunas de estas herramientas son:

- ORM, Object-Relational-Mapping: permite crear un modelo de datos y gestiona automáticamente la BD (base de datos) subyacente. Abstacta las consultas SQL y evita tener que realizar migraciones manuales de los esquemas.
- Static File Generator: podemos diseñar la interfaz de la aplicación en formato web y desplegar en un servidor.
- Commands-System: gestión de comandos internos mediante los que se puede automatizar tareas; utilizaremos esta arquitectura para diseñar el pipeline de datos y almacenar la información proveniente de diferentes fuentes.

```
Stage Logs (sentinel)
@ Shell Script – python3.10 manage.py 4-download_img -a "20220601" -b "20220701" -p "data/parcelas/25/pixeles.shp" -i ndvi ndre (self time 11min 21s)

+ python3.10 manage.py 4-download_img -a 20220601 -b 20220701 -p data/parcelas/25/pixeles.shp -i ndvi ndre

 0%|          | 0/2 [00:00<?, ?it/s]
 50%|██████| 1/2 [05:32<05:32, 332.20s/it]
100%|██████| 2/2 [11:06<00:00, 333.27s/it]
100%|██████| 2/2 [11:06<00:00, 333.11s/it]
```

Ejemplo de uso de un comando en Django para ejecutar la descarga de los índices (salida del comando mostrada en Jenkins).

Aunque el framework es muy potente, harán falta otras herramientas y entornos para completar con éxito la automatización que buscamos para así conseguir un proceso de optimización continua. En el siguiente punto hablaremos de la infraestructura que la aplicación requiere y de cómo podemos solventar algunos de los problemas de integración más importantes. Mencionamos algunas de estas tecnologías:

- Vagrant: virtualización de una máquina Unix en la que podamos provisionar el proyecto y reproducir los pasos necesarios de nuestro pipeline.
- Jenkins: herramienta de integración continua que permite declarar un pipeline a modo de infraestructura como código.

La utilización de *Python*, aparte de incluir el framework de Django, nos da la posibilidad de utilizar las librerías de inteligencia artificial y ciencia de datos que son necesarias para la creación del modelo de producción de cultivo. Algunas de librerías principales que vamos a utilizar son las siguientes:

LIBRERÍA	DESCRIPCIÓN
scikit-learn	Creación de modelos de inteligencia artificial
numpy	Tratamiento de datos multidimensionales
pandas	Uso de datos tabulados con herramientas para su procesamiento
jupyter	Integración de cuadernos procesables con el resto del proyecto

## 5.5 Requisitos

La siguiente tabla muestra los requisitos funcionales críticos para completar con éxito el proyecto, terminando el trabajo con un producto estable que el equipo pueda utilizar.

REQUISITO	DESCRIPCIÓN
RF1	Automatizar el despliegue de la aplicación identificando las librerías y dependencias necesarias.
RF2	Rediseñar el modelo de datos para poder escalar la aplicación y su gestión de datos.
RF3	Documentar el proceso de desarrollo generando los documentos de diseño pertinentes.
RF4	Crear un entorno común que permita trabajar al equipo sobre las mismas fuentes de datos.
RF5	Obtener una vista minable con la selección de <i>features</i> provenientes del modelo.
RF6	

<b>REQUISITO</b>	<b>DESCRIPCIÓN</b>
	Creación de un modelo de producción probando varios algoritmos de inteligencia artificial.
RF7	Despliegue de los datos del modelo en la interfaz de la aplicación.

Como requisitos no funcionales para la aplicación identificamos los siguientes, entendiendo que se completarán a lo largo de todo el proyecto.

<b>REQUISITO</b>	<b>DESCRIPCIÓN</b>
RNF1	Crear un entorno Linux local con la misma configuración necesaria en producción.
RNF2	Preparar las fuentes de datos con una muestra pequeña representativa de la aplicación en producción.
RNF3	Estudiar las posibles features para la vista minable que utilizará el modelo de producción de cultivo
RNF4	Implantar un flujo de trabajo en el equipo para utilizar las mismas fuentes de datos
RNF5	Utilizar software libre

# 6 INFRAESTRUCTURA

---

## 6.1 Arquitectura física

---

El entorno local, sobre el cual desarrollamos, y el de producción, donde se despliega la aplicación, son actualmente diferentes. Todos los miembros del equipo utilizamos máquinas Windows, mientras que en el servidor de producción encontramos una máquina Ubuntu. Esto crea un problema importante en el proceso de instalación de las librerías necesarias para trabajar con los datos; se pierde gran cantidad de tiempo en preparar el entorno de cada persona que va a trabajar con el repositorio de la aplicación; además, las librerías necesarias necesitan de una gran cantidad de dependencias que varían en cuanto a las versiones. Como solución propondremos una nueva forma de trabajo que nos asegure un entorno común con las mismas librerías y paquetes.

Parte del trabajo consiste en la automatización del despliegue de la aplicación. En este punto veremos cómo la infraestructura como código permite aislar las librerías necesarias creando un script que deja una máquina en estado estable para ejecutar la aplicación. Para poder utilizar tecnologías que posibiliten la integración continua, transformamos el entorno de desarrollo desplegando la aplicación en UNIX mediante una máquina virtual Vagrant en la que podamos reproducir varias veces el proceso. Identificamos las siguientes dependencias que necesitan ser instaladas:

- Postgres 14 + Postgis
- Python3.10
- GDAL 3.3.2
- Django
- Jenkins

Suele ser complicado encontrar todas las dependencias con sus versiones correctas; en este caso el punto más complicado ha sido la instalación de Python con su versión correspondiente de GDAL, librería que permite tratar con los datos geoespaciales. Además, el sistema gestor de BD, 'postgres', necesita una extensión especial, 'postgis', para poder guardar los datos georreferenciados. Como este tipo de trabajo es casi de prueba y error, la utilidad del script que obtenemos es de gran valor.

Jenkins es una herramienta de integración continua muy relevante para el proceso de automatización y mejora continua que queremos implantar. Aquí mencionamos cómo se incluye en la máquina Vagrant que contiene el proyecto. En los siguientes puntos veremos cómo permite integrar los distintos scripts de carga de datos junto con la generación de la vista minable y posterior modelo que conformar el pipeline que queremos conseguir.

## 6.2 Vagrant

---

Hacemos uso de Vagrant para crear, en una máquina virtual Unix (Ubuntu), el entorno de desarrollo necesario para el proyecto. De esta forma conseguimos que los entornos de pruebas y producción sean muy similares, permitiendo automatizar el despliegue de la aplicación mediante técnicas de infraestructura como código. Una de las ventajas de tener una máquina Ubuntu, es que, mediante scripts .sh, instalamos todas las dependencias necesarias para dejar dicha máquina lista para el despliegue de la aplicación.

Para el entorno de desarrollo virtualizamos con Vagrant una máquina Ubuntu 20.04. El archivo *Vagrantfile* nos permite provisionarla con el código del proyecto y el script que deja a dicha máquina con las librerías necesarias para el correcto funcionamiento de este.

```
Vagrant.configure("2") do |config|
  config.vm.box = "bento/ubuntu-20.04"

  config.vm.provision "file", source: "~/Desktop/agrai", destination: "$HOME/"
  config.vm.provision :shell, :path => "start.sh"
  config.vm.provision :shell, :path => "jenkins.sh"

  config.vm.define "server" do |server|
    server.vm.network "private_network", ip: "192.168.56.19"
    server.vm.provider "virtualbox" do |vb|
      vb.memory = "4096"
      vb.cpus   = "2"
    end
  end
end
```

Se enlazan como entregables los dos scripts con los que provisionamos dicha máquina. El primero, 'start.sh', hace referencia a la instalación de las librerías necesarias, el segundo, instala la herramienta Jenkins para la posterior ejecución de tareas dentro de un pipeline. A continuación mostramos algunas de las dependencias que más ha costado que funcionen conjuntamente.

```
#!/bin/bash
sudo apt install -y postgresql-14
sudo apt install -y postgresql-14-postgis-scripts

# gdal native:
sudo add-apt-repository ppa:ubuntugis/ppa && sudo apt-get update

# python concrete installation
sudo apt install -y python3.10 # version concreta
sudo apt install -y python3.10-dev python3.10-venv
sudo apt install -y virtualenv
sudo apt install -y build-essential
sudo apt install -y python3.10-tk

# Tkinter
sudo apt install -y python3-tk

# GDAL
sudo apt install -y libgdal-dev
sudo apt install -y gdal-bin

# env gdal lib variables
export CPLUS_INCLUDE_PATH=/usr/include/gdal
export C_INCLUDE_PATH=/usr/include/gdal
```

## 6.3 Entorno

---

Por último, es importante preparar un entorno común para el equipo. Django es un framework muy útil por la integración de las herramientas que hemos visto (ORM, integración continua, etc.), pero para conseguir un *workflow* adecuado necesitamos incluir otras herramientas propias de ciencia de datos como pueden ser cuadernos de Jupyter y sus entornos, con las librerías necesarias para ejecutar modelos predictivos.

Dedicamos varios días a preparar la integración de Django con cuadernos Jupyter y a ejecutar diferentes entornos virtuales con las librerías adecuadas en cada momento. Los siguientes ejemplos muestran cómo se ha conseguido integrar las librerías necesarias para el despliegue de los cuadernos y su integración con los datos y los modelos de Django. El siguiente snippet es relevante debido a que se necesita ejecutar en la primera celda del cuaderno para tener acceso a los modelos y servicios definidos en la arquitectura de Django.

```
# base code for jupyter integration

import sys, os, django

BASE_DIR = os.path.dirname(os.path.abspath('../..../agrai/'))

sys.path.insert(0, BASE_DIR)

os.environ.setdefault("DJANGO_SETTINGS_MODULE", "agrai.settings")
os.environ["DJANGO_ALLOW_ASYNC_UNSAFE"] = "true"

django.setup()
```

Los cuadernos Jupyter que se encuentran en esta memoria se han exportado desde la propia aplicación. Vemos la facilidad con la que el equipo puede redactar informes para clientes concretos, utilizando los datos de la aplicación directamente. El código anterior es el 'snippet' necesario para poder integrar los cuadernos con las llamadas concretas a los modelos de Django. Estos cuadernos son un punto tan importante de la arquitectura como puede ser el módulo de servicios.

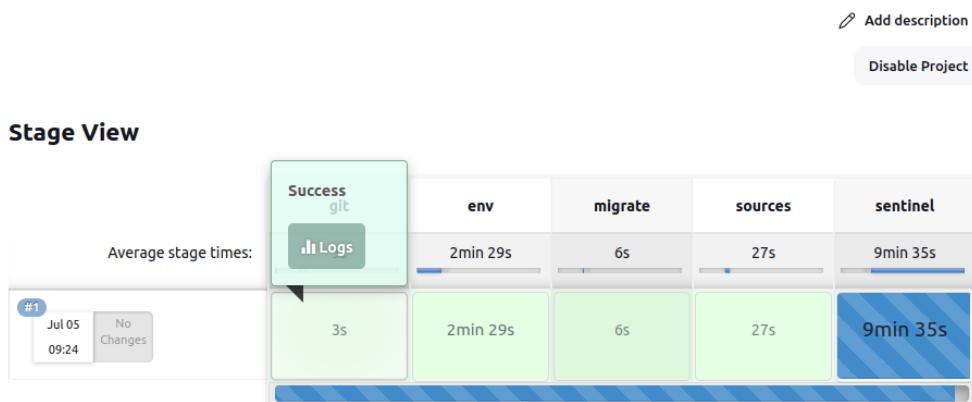
# 7 DISEÑO

## 7.1 Carga de datos

Un modelo claro y robusto va a permitir cargar datos en la aplicación de forma ordenada con posibilidad de ser escalados en el futuro. En la implementación del pipeline que veremos en siguientes apartados haremos hincapié en los estados de carga por los que pasan los datos; de momento solo importa destacar que el modelo actual permite mantener un histórico de datos mucho más sólido que el que estaba implantado cuando comenzamos el desarrollo.

Es importante explicar por qué vamos a poner tanto énfasis en que estos datos se gestionen de manera fluida y eficaz. La interfaz web de la aplicación muestra la punta de un iceberg en la que el usuario observa las parcelas coloreadas basándose en los índices vegetativos, también se muestra la producción estimada.

### Pipeline tfg



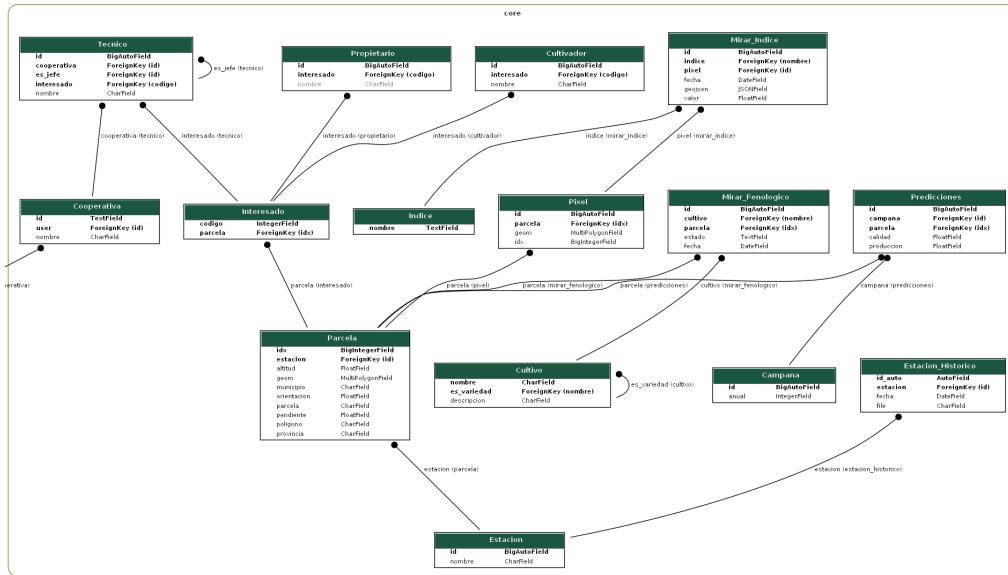
Ejecución del pipeline mientras realiza la descarga de datos del satélite Copernicus.

El proceso de carga de datos hace uso de un módulo de descarga de imágenes satelitales que obtiene el valor de estos índices vegetativos. Este script de descarga ha sido desarrollado por el equipo y ha sido colocado dentro del módulo de automatización junto con otras tareas similares. Utilizaremos el módulo como una caja negra y aseguraremos el correcto uso de este para que los índices vegetativos terminen en sus tablas correspondientes dentro de la BD.

Los índices que renderiza el visor de la aplicación son el resultado de todo el proceso que estamos exponiendo. Es el punto más delicado que requiere del diseño de varias entidades para guardar un histórico de datos. En la visualización del histórico de índices vegetativos es el mayor valor que obtiene el cliente cuando accede a la aplicación.

## 7.2 Modelo relacional

Como inicio de la posterior automatización, se reestructura la arquitectura de la aplicación junto con su BD para soportar el almacenamiento de nuevos datos. El siguiente esquema de BD es el resultado de la implementación del análisis descrita en el planteamiento.



Esquema relacional en uso del componente principal de la aplicación

Para la creación de este modelo he trabajado con el equipo en la identificación de los conceptos que necesitaban ser representados. El diseño en papel da la posibilidad de pensar abiertamente sobre las relaciones entre entidades, además de permitir la transmisión de ideas de forma sencilla durante las reuniones.

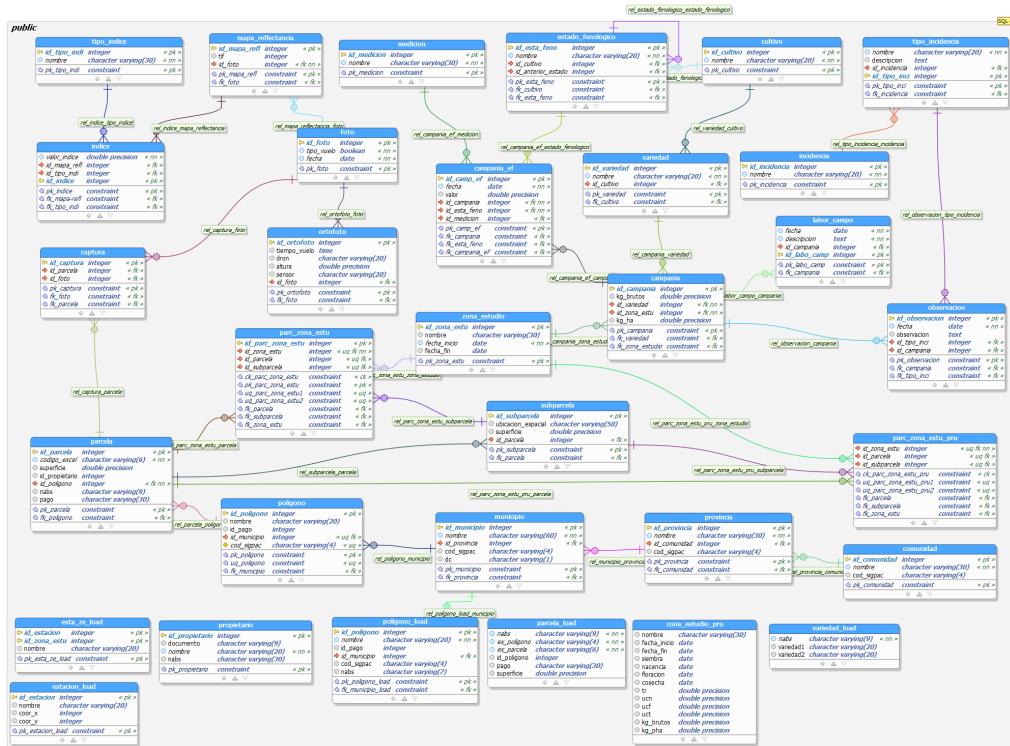
Dentro del entorno de trabajo de *Django*, el ORM proporcionado aísla la base de datos y nos permite diseñar directamente en *Python* dicho modelo. Las entidades se diseñan como clases y las relaciones entre ellas se especifican mediante el lenguaje de mapeo proporcionado. De dicha forma creamos las claves foráneas que físicamente contiene la base de datos.

La utilización de este ORM es una ventaja que nos evita usar SQL directamente y permite poblar la base de datos mediante comandos *Python* desde la terminal. De todas formas, la localización física de la BD necesita ser enlazada correctamente. Este aislamiento posibilita utilizar diferentes bases de datos en los entornos locales y de producción.

Varias etapas han sido necesarias hasta llegar a un punto más o menos estable. La herramienta *Graphviz* ha permitido obtener diagramas UML que visualizan la estructura de clases del modelo a partir del código. Esta representación gráfica ha sido realmente útil para poder pensar sobre el diseño a medida que avanzaba.

## Entidades Principales

De la situación inicial, nos quedamos solo con las entidades principales. Esta selección de entidades se realiza cuidadosamente para que la aplicación pueda contemplar la mayor información posible haciendo uso de un esquema sencillo pero robusto. Analizamos el diseño detalladamente en los siguientes puntos.



*Esquema relacional inicial de partida, con entidades definidas, pero sin implementar, se toma como idea para el comienzo del diseño.*

Este rediseño contempla la mayoría de casos posibles para el posterior tratamiento y procesamiento de datos, con posibilidad de crear buenos modelos predictivos. A continuación se describen las entidades principales contempladas en la base de datos.

ENTIDAD	DESCRIPCIÓN
Parcela	Como toda la información que se va a almacenar es referente a las parcelas y sus datos, toda la BD gira en torno a la tabla “Parcela”.
Cultivo	Esta tabla está directamente unida a la anterior, ya que un cultivo puede tener muchas variedades diferentes se modela de forma recursiva con una clave foránea a sí mismo (leer el siguiente punto, Ej: Guisante es un tipo de cultivo, pero tiene varias variedades: tirabeque, snap peas...).
Indice	Esta tabla es de vital importancia para el funcionamiento de la aplicación, toda nuestra información lleva a estos datos vegetativos. Serán datos calculados en un dominio, esto se hace así ya que existen cientos de formas de denominar un mismo índice de vegetación y podría ser confuso a la hora de realizar las consultas. Para ello se establecerán algunos campos determinados en el campo “tipo_indice” y se les dará un valor en “valor_indice”.
Mirar_Indice	Contempla un histórico de datos y permite dar valor concreto a un índice vegetativo en una fecha única.
Fenologico	Permite registrar diferentes estados fenológicos por los que puede pasar un cultivo. La relación entre el estado fenológico y el cultivo se realiza a través de la tabla Mirar_Fenologico.
Mirar_Fenologico	Permite dar carácter temporal al estado fenológico de un cultivo. Registramos un estado para un cultivo en una fecha concreta.
Campaña	Es una tabla muy importante, sirve para unir distintos tipos de registros (desde variedades hasta labores de campo pasando por unir los datos de parcelas y subparcelas).

ENTIDAD	DESCRIPCIÓN
Interesado	Permite realizar el diseño de roles mediante su clave foránea a una parcela. Un interesado (stakeholder) es una persona que tiene relación con una o varias parcelas (Ej, cultivador, propietario, etc).

## Diseño del Cultivo y Fenología

Una decisión complicada sobre el posible histórico de datos es el registro de la evolución de un cultivo en una parcela concreta. Mediante la tabla MIRAR\_FENOLÓGICO damos carácter temporal al cultivo concreto de una parcela. De esta forma podemos registrar cuáles son los estados por los que ha pasado un cultivo, desde su "siembra" hasta su "recolección".

El modelado de la entidad FENOLÓGICO ha sido una decisión complicada debido a que no se enlaza directamente con un cultivo. Entendemos que un estado como puede ser el de "siembra" tiene un carácter general y tiene sentido como entidad en sí misma (todos los cultivos pasan por siembra). Es su relación con CULTIVO mediante MIRAR\_FENOLÓGICO lo que dice que dicho cultivo está en un estado fenológico concreto en un momento determinado. Por ejemplo, podemos decir que el cultivo "vid" estaba en estado de "siembra" el día "2022-01-23".

También es importante destacar que CULTIVO es simplemente el tipo que se ha registrado, por lo tanto, para que la información sea completa, un avistamiento fenológico sucede sobre un tipo de cultivo físicamente sembrado en una PARCELA. En la próxima figura podemos ver cómo estas tres entidades se relacionan para dar carácter temporal y espacial a un tipo de cultivo.

Otro punto importante de la aplicación es la contemplación de variedades. Para poder mantener una jerarquía con las posibles entidades registradas en el sistema, enlazamos de forma recursiva el cultivo con una clave foránea a su misma tabla. Esta estructura permite el desglose de una jerarquía de cultivos en la que sabemos qué entidad es una subvariedad de un cultivo, dando la posibilidad de almacenar múltiples niveles.

agrai=# select * from core_mirar_fenologico;				
id	estado	fecha	cultivo_id	parcela_id
126	SIEMBRA	2021-03-01	TEMPRANILLO	3502
756	SIEMBRA	2021-03-01	TEMPRANILLO	3553
740	SIEMBRA	2021-03-01	TEMPRANILLO	3582
945	SIEMBRA	2021-03-01	TEMPRANILLO	3804
634	SIEMBRA	2021-03-01	TEMPRANILLO	3882
633	SIEMBRA	2021-03-01	VIURA	3914
313	SIEMBRA	2021-03-01	TEMPRANILLO	3942
299	SIEMBRA	2021-03-01	TEMPRANILLO	3998

Ejemplo de una consulta SQL sobre la tabla MIRAR\_FENOLÓGICO

## Diseño de Índices Vegetativos

Para nuestro sistema de información una parcela es una agrupación de varios píxeles. Entendemos como PIXEL a la imagen satelital más pequeña que se puede obtener sobre el terreno, a partir de la cual obtendremos los índices vegetativos. (*El valor del índice es una caja negra para este trabajo, proviene de scripts desarrollados por el resto del equipo*)

El trabajo con índices vegetativos por parte del equipo es uno de los puntos más importantes para que el rediseño sea satisfactorio. La entidad MIRAR\_ÍNDICE es de vital importancia en el modelo, permite gestionar el histórico de valores de los índices para todos los píxeles. Esta tabla es una de las candidatas para realizar

optimización en las consultas, algo que desde el ORM que utilizamos todavía queda un poco lejos.

El volumen de datos que la tabla MIRAR\_INDICE puede contener es grande. Solo con la muestra representativa de las 25 parcelas con las que estamos trabajando (3000 píxeles), dos índices registrados, y contemplando sus valores en 6 fechas distintas, obtenemos un conjunto de 36.000 valores. Cuando la aplicación escala a un mayor número de parcelas, esta tabla será susceptible de utilizar índices (BD).

Este diseño de índices permite mantener un histórico de datos preciso. Gestiona la evolución de los cultivos a través de las diferentes fechas en las que se persisten los valores de los índices. Destacamos que cada índice se corresponde con una capa de visualización sobre el visor GIS que contiene la presentación de la aplicación, requisito que hace que tenga mucha importancia el poder añadir índices a los datos a medida que se vayan contemplando.

## Diseño de Roles

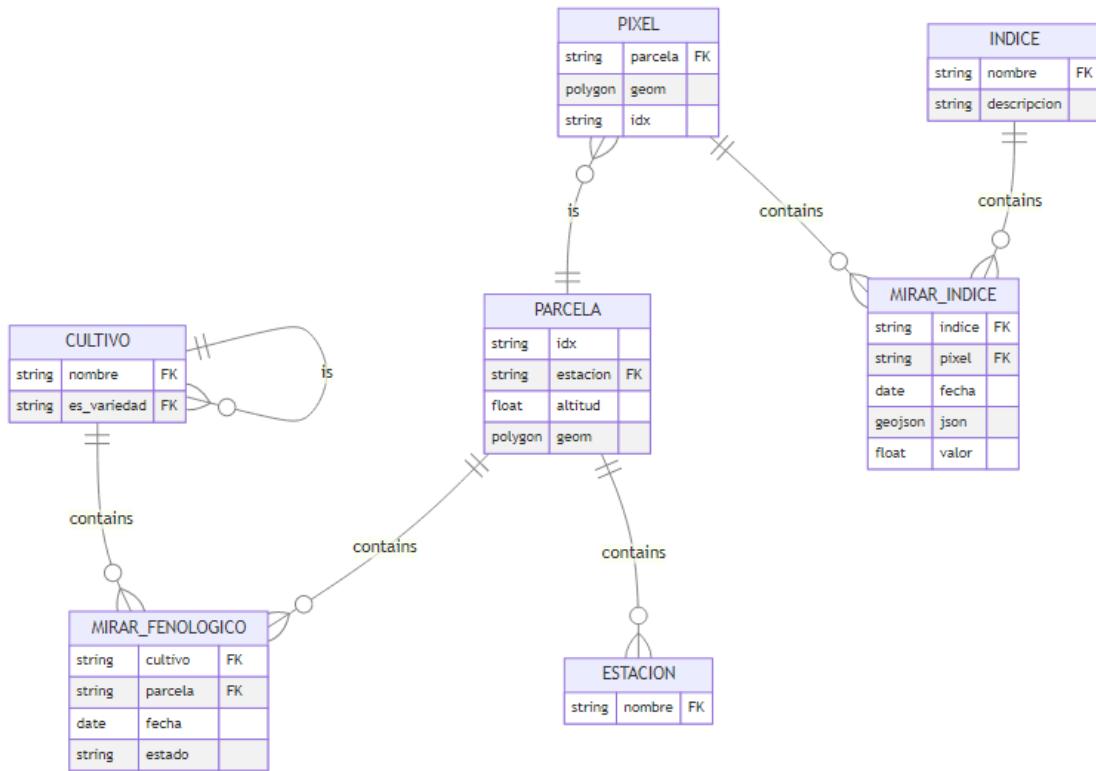
Aunque es un punto secundario, la aplicación contempla que diferentes usuarios puedan realizar diferentes acciones. Recordamos que habíamos separado en componentes la aplicación, y los roles se encuentran en un módulo secundario independiente de las entidades anteriores, que representaban los principales conceptos agronómicos y vegetativos. El diseño modular del modelo de datos permite al equipo añadir diferentes roles sin que estos estén *hardcoded* en el código de la aplicación. La siguiente tabla describe los roles hasta ahora contemplados por el equipo y su relación con el parcelario.

ROL	DESCRIPCIÓN
Cultivador	persona encargada de realizar los mantenimientos en campo
Propietario	propietario catastral de las parcelas
Cooperativa	entidad que agrupa a técnicos con diferentes cargos sobre un parcelario
Técnico	personal asignado a un número de parcelas.

Aunque la aplicación hace uso de estas entidades, no serán tan importantes para el desarrollo posterior del trabajo, debido a que la parte que necesita automatizarse es aquella relacionada con los índices vegetativos y la geometría de las parcelas anteriores. Como decíamos, estas entidades quedan separadas en un módulo de la aplicación con la única funcionalidad de controlar el acceso a la aplicación y dar permiso a las funcionalidades que cada cliente puede realizar dependiendo del rol que tenga asignado.

## Modelo Final

El siguiente diagrama recoge todas las relaciones descritas anteriormente y muestra el esquema completo que usa la aplicación. En los siguientes puntos hablaremos sobre los métodos definidos sobre este esquema de datos, dónde se localizan y cómo han de ser utilizados.



## 7.3 Servicios

Ahora tenemos un modelo relacional que soporta la gestión de grandes volúmenes de datos. En este punto tenemos que diseñar las diferentes operaciones que necesitamos hacer sobre los datos almacenados. Discutimos cómo tiene que ser la implementación de las operaciones más frecuentes y por qué en un modelo de datos sin estado, con información histórica, es más importante diseñar estas operaciones sobre el modelo relacional como servicios.

El término servicio está sobrecargado y su significado adquiere diferentes matices según el contexto en que estemos. Como resultado, existe una nube de confusión en torno a la noción de servicios cuando se trata de distinguir entre servicios de aplicación, servicios de dominio, servicios de infraestructura, servicios SOA, etc. Las funciones de estos son diferentes y pueden abarcar todas las capas de una aplicación.

De hecho, un servicio es un título un tanto genérico para un bloque de creación de una aplicación porque implica muy poco. En primer lugar, un servicio implica un cliente para cuyas solicitudes está diseñado. Otra característica de una operación de servicio es la de entrada y salida: se proporcionan argumentos como entrada a una operación y se devuelve un resultado. Más allá de esta implicación suelen estar los supuestos de "statelessness" y la idea de "pure fabrication", según GRASP:

**When a significant process or transformation in the domain is not a natural responsibility of an ENTITY or VALUE OBJECT, add an operation to the model as standalone interface declared as a SERVICE. Define the interface in terms of the language of the model and make sure the operation name is part of the UBIQUITOUS LANGUAGE. Make the SERVICE stateless.**

**Eric Evans** Domain-Driven Design

El tipo de servicios que estamos diseñando e implementando para nuestra aplicación forman parte de la capa de dominio. Estos servicios de dominio a menudo se pasan por alto como bloques de construcción clave, confundidos por el enfoque de las entidades del modelo (o value objects).

Cumpliendo con los principios mencionados, los servicios conforman la siguiente capa de abstracción del modelo de datos implementado. Colocaremos aquellas operaciones que dependan o relacionen más de una entidad en su módulo correspondiente de servicios y no como un método de la clase del modelo. Este tipo de diseño en el que dejamos el modelo casi sin métodos propios puede llegar a entenderse como un anti-patrón, [anemic domain model](#), pero que en nuestro caso, tras buscar la forma más sencilla de acceder a la información, será de gran utilidad.

Justificamos la implementación de gran parte de los métodos sobre la capa de servicios por el tipo de información histórica sobre la que necesitamos hacer las consultas. La mayoría de entidades necesitan de una relación con otra segunda o tercera entidad para devolver la información pertinente. **Como el modelo está orientado a manejar información histórica y sin estado**, optamos por colocar casi todos los métodos de acceso a los datos en un módulo de servicios aparte. Miremos el siguiente ejemplo sobre los históricos fenológicos en el parcelario registrado.

```
class ParcelaCultivos_Service():

    def get_cultivo(parcela_p):
        """
        Último cultivo que se está cultivando en una parcela """
        mirar_feno = Mirar_Fenologico.objects.filter(parcela = parcela_p.idx)
        if (mirar_feno.count() > 0):
            return mirar_feno.order_by('fecha')[0].cultivo
        else:
            return None

    def get_historico_fenologicos(parcela_p):
        """
        todos los estados fenológicos por los que ha pasado una parcela """
        return Mirar_Fenologico.objects.filter(parcela = parcela_p.referencia)

    def get_historico_range(fecha_inicio, fecha_end):
        """
        histórico de todas las parcelas en un rango de fechas """
        return Mirar_Fenologico.objects.filter(fecha__range=[fecha_inicio, fecha_end])

    def get_historico_parcela(parcela1):
        """
        todo el histórico de una parcela """
        return Mirar_Fenologico.objects.filter(parcela = parcela1)

    def get_historico_parcela_range(parcela1, fecha_inicio, fecha_end):
        """
        histórico de una parcela en un rango de fechas """
        historico = ParcelaCultivos_Service.get_historico_range(fecha_inicio, fecha_end)
```

```

    return historico.filter(parcela = parcela1)

def get_historico_mismo_fenologico(valor_fenologico):
    """ todo los cultivos que están en un mismo estado fenológico """
    return Mirar_Fenologico.objects.filter(estado = valor_fenologico)

```

Adjuntamos como entregables del proyecto algunos de los servicios diseñados y usados más frecuentemente. El código de estos servicios es utilizado por el resto de miembros del equipo para otras tareas de la aplicación, como puede ser la presentación de los datos en la interfaz web.

El pipeline que vamos a diseñar es posible gracias a la factorización en módulos de la aplicación y a la sencillez y versatilidad del modelo relacional de datos. Los servicios que ahora implementamos abstraen al resto del equipo de la funcionalidad subyacente y me permiten crear diferentes interfaces con propósitos distintos. Son muy importantes, ya que nos acercan al flujo de trabajo de integración continua que estamos buscando.

Por otro lado, tenemos métodos más clásicos en los que hacemos uso de las entidades registradas para obtener valores concretos. La siguiente función es una de las más usadas debido a la integración directa con el proceso de descarga de imágenes satelitales del cultivo. Fijémonos en cómo utilizamos las entidades creadas anteriormente para acceder a la información de la BD, el ORM subyacente nos abstracta de otras consultas complejas que podríamos hacer en lenguajes como SQL.

```

def get_indice_func(parcela, indice_p, fecha_p, func, filtro):
    """Método que encapsula el uso de una función agregada junto con un filtro"""

    indice_p = Indice.objects.get(nombre=indice_p)

    if ParcelaIndices_Service.is_parcela_filter(parcela, filtro):
        pixels = parcela.get_pixeles()

        indices_pixels = Mirar_Indice.objects.filter(
            pixel_in = pixels,
            indice = indice_p,
            fecha = fecha_p
        )

        if ParcelaIndices_Service.get_indice_filter(indices_pixels, filtro):
            array_valores = list(indices_pixels.values_list('valor', flat=True))

            return func(array_valores)

```

## 8 PIPELINE

---

### 8.1 Transformación de datos

---

Una vez tenemos un modelo que soporta las principales entidades del dominio agronómico que necesitamos, procedemos a cargar la información proveniente de distintas fuentes en la BD bajo sus correspondientes entidades y con las relaciones necesarias para las consultas. Para esta carga vamos a utilizar un pipeline de datos; veamos a qué nos referimos.

Es importante diferenciar qué tipo de estructura estamos construyendo. Existen actualmente varias arquitecturas para 'pipelines', como pueden ser ETLs u otras. Para la correcta carga y transformación de datos no utilizaremos ninguna de estas arquitecturas, sino que diseñaremos los bloques necesarios en cada paso para poblar nuestra base de datos. Este proceso puede entenderse como una canalización de datos, en la que recibimos información en diferentes fuentes y la dotamos de contexto dentro de la BD.

#### Canalización de datos

Una canalización de datos hace referencia a los pasos necesarios para mover datos del sistema de origen al sistema de destino. Estos pasos incluyen copiar datos, transferirlos desde una ubicación a otra y combinarlos con otras fuentes de datos. El objetivo principal de una canalización de datos es garantizar que todos estos pasos se produzcan de forma coherente con todos los datos.

En el apartado anterior hemos visto cómo se ha diseñado el modelo de datos. Ahora nos centramos en los pequeños pasos de carga que vamos a dar para que los datos agronómicos de clientes e imágenes satelitales persistan en dicho modelo. Valoramos la identificación de unidades atómicas de persistencia de información que puedan ejecutarse reiteradamente. Es decir, buscaremos acotar pequeños procesos de carga que puedan ejecutarse en varios puntos dependiendo del volumen de datos de clientes que maneje el equipo en un momento dado.

#### Procesos de una canalización

Identificamos tres conceptos que definen la canalización de datos que vamos a llevar a cabo y exponemos qué hace el *pipeline* de nuestra aplicación en cada uno de ellos.

- Data Ingestion: diseñamos bloques atómicos de carga que almacenan datos de fuentes diferentes en la base de datos de la aplicación.
- Data Transformation: utilizamos los bloques de carga para almacenar los datos en nuestro modelo de Django.
- Data Storage: almacenamos los datos en el modelo a través de Django, manteniendo una base de datos relacional en la máquina en la que se encuentra la aplicación.

Como hemos dicho, hay varias formas de realizar este proceso de canalización; en nuestro caso identificaremos los bloques y diseñaremos el proceso casi de forma manual, para que se ajuste con exactitud a lo que queremos hacer. En el siguiente punto exponemos la tecnología subyacente que utilizaremos para ello.

## 8.2 Design pattern

---

En ingeniería de software, un **patrón de diseño** es una solución general y reutilizable para un problema común dentro de un contexto dado. El término "pipeline" puede hacer referencia a diferentes soluciones y contextos. De momento vamos a centrarnos en su aplicación como patrón de diseño ya que en nuestra aplicación se está utilizando como tal.

**Podemos definir un pipeline como una cadena de elementos de procesamiento (procesos, subprocessos, rutinas, funciones, etc.) dispuestos de modo que la salida de cada elemento sea la entrada del siguiente; el nombre hace referencia al flujo de una o varias tuberías.**

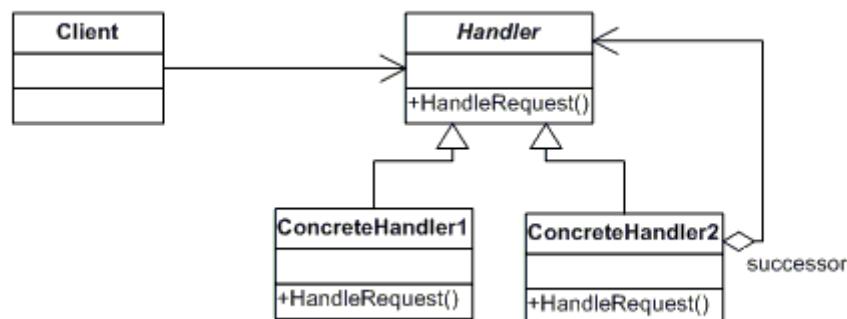
En nuestro caso vamos a utilizar la arquitectura de comandos que proporciona Django para diseñar los procesos de transformación de datos que van a poblar la base de datos. La definición y reutilización del proceso de población es muy importante debido a la forma de trabajo que tiene la empresa con diferentes clientes. Dependiendo del cliente, se prepara una copia de la aplicación para trabajar con sus datos; por ello, una vez definidas las unidades atómicas del proceso de carga de datos, se podrán diseñar diferentes pipelines para cada cliente. Para este trabajo nos conformaremos con el diseño y la implementación de un proceso general, pero manteniendo la idea de que pueda ser escalado más adelante.

```
from django.core.management.base import BaseCommand

class Command(BaseCommand):

    def handle(self, **options):
        # now do the things that you want with your models here
```

Las líneas de código anteriores hacen referencia a una implementación concreta del patrón mencionado. Fijándonos en la siguiente figura, la clase **Command** sería el handler; el método **handle** corresponde con *handleRequest* del esquema. Los próximos scripts que creemos serán los clientes de estas clases. Destacaremos que la implementación de Django de este patrón es un poco más compleja porque los comandos que encapsulan la funcionalidad permiten tomar diferentes parámetros.



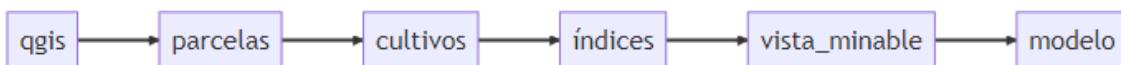
*Pipeline design pattern, diagrama de clases para exponer cómo vamos a utilizar los bloques de carga creamos.*

Con este patrón desarrollamos la infraestructura necesaria para cargar los datos en el modelo. Cada comando representa un proceso de carga de datos, el cual se puede componer posteriormente dentro de un proceso más complejo. Cada unidad de carga se comporta como un filtro que añade de forma ordenada la información al modelo de datos, enlazando las entidades con sus datos correspondientes.

Utilizaremos estas unidades para generar diferentes tuberías que permitan dejar varias instancias de la aplicación en estados diferentes.

Para AGRAI hace falta transformar los datos de índices vegetativos que maneja el equipo en información ordenada que llegue al modelo de datos desarrollado anteriormente. A continuación explicamos el flujo de procesos de carga más utilizado.

Los datos principalmente provienen de *Qgis*, software para trabajar con información georreferenciada. El equipo pre-procesa la información sobre los índices vegetativos y da valores a los píxeles mediante dicho software. La principal función de este trabajo es dar sentido a esa información, guardando los datos tomados junto con sus relaciones con el parcelario físico.



La información de los píxeles tiene que enlazarse con los datos de parcelas y cultivos. Cuando obtenemos los índices vegetativos solo contemplamos la información sobre la geometría de la parcela, no sabemos nada más. Esta información tiene que almacenarse junto con el cultivo, el estado fenológico, los datos físicos de la parcela, etc. En la siguiente tabla explicamos cada uno de los [scripts](#) que conforman los procesos del pipeline.

SCRIPT	TAREA
load_qgis	Extrae la información geométrica de cada pixel proveniente de "qgis" e inserta dicha información en las tablas de "parcela" y "pixel" de la BD.
load_parcela	Lee de un excel información fisca sobre las parcelas e inserta los datos en la BD.
load_cultivos	Lee de un excel información sobre cultivos con sus variedades e inserta los datos en la BD.
download_img	Usa uno de los módulos desarrollado por el equipo para descargar índices vegetativos a partir de imágenes satelitales. No almacena nada directamente, devuelve un dataframe con los datos descargados.
load_indice	Integra la información de índices descargada en el paso anterior en la tabla de "pixel" de la BD.
load_range	Utiliza el script anterior para cargar varios índices.
create_view	Construye la vista minable de la que hablaremos en el siguiente punto.
create_model	Busca y entrena el mejor modelo IA posible sobre la vista minable anterior.
predict_model	Utiliza el modelo entrenado para predecir sobre los datos de prueba que contiene la vista minable.

Para entender detalladamente qué hacen los scripts de carga de índices tenemos que explicar algunos de los procesos de automatización de los que hacemos uso.

Estos scripts nos permiten crear el flujo necesario para que la aplicación procese la información y termine en un estado consistente. Parte del proceso necesita ser automatizado y ejecutado reiteradamente, añadiendo los datos a la BD y relacionándolos adecuadamente. En nuestro contexto agronómico, es la descarga de

índices la parte del proceso que necesita ser automatizada, mientras que los datos de las parcelas y cultivos y la información geoespacial se pueden cargar directamente en el despliegue de la aplicación.

Con esta estructura de tuberías lista, podemos reutilizar cualquier parte del proceso. Un planificador de tareas nos ayuda a establecer el tiempo que tiene que transcurrir entre las descargas. Se ha utilizado el scheduler por defecto del framework de Django que estamos utilizando para esta gestión de tiempos y descarga,

Finalmente, para realizar pruebas y exponer el proceso en esta memoria, Jenkins nos permite ejecutar de forma ordenada el flujo completo del pipeline mediante infraestructura como código. El resto de tuberías que se utilizan en el proyecto tienen una forma similar y hacen el mismo uso de los comandos de Django que hemos explicado.

Por último, este patrón nos permite ejecutar comandos de forma asíncrona y gestionar los procesos con un planificador. Actualmente, aunque el tiempo de carga es alto, podemos asumir una carga lenta. Pero es este mismo patrón de diseño el que nos puede permitir a futuro la carga asíncrona de los datos en cada una de los componentes que lo forman.

```
def run (self, DATA, INDICES):

    for indice in INDICES:

        df = DATA.copy()

        drop_cols = []

        for col_name in list(df):

            if indice.upper() not in col_name:

                if (col_name == 'ID'):

                    continue

                drop_cols.append(col_name)

        df1 = df.drop(drop_cols, axis=1)

        call_command('4-load_indice', data=df1, indice=indice)
```

Adjuntamos en los entregables los [scripts](#) que forman este proceso de almacenamiento de datos.

## 8.3 Vista minable

---

Hasta este punto habíamos insertado en la base de datos la información que provenía de diferentes fuentes. Hemos dotado de consistencia y relación a los datos y se ha creado un sistema con capacidad para hacer consultas sobre estos de forma ordenada y sencilla.

Ahora, fuera de la estructura de comandos de *Django* y utilizando el entorno virtual creado y los cuadernos integrados en la aplicación, escribimos el algoritmo para extraer los datos y dar la forma necesaria para los siguientes pasos.

El código que mostramos en este apartado corresponde con la generación de la vista tabulada con la que vamos a crear el modelo de producción de cultivo. Ahora realizamos el proceso inverso a la persistencia de datos que hemos realizado hasta

este punto. Utilizamos la información y sus relaciones perisitidas en la base de datos de nuestro sistema para generar la vista minable que necesita el modelo predictivo.

```
def statistic_indices(indices = ['ndvi', 'ndre'], func=np.mean):

    df = pd.DataFrame()

    col_ids = []

    for p in Parcela.objects.all():

        col_ids.append(p.idx)

    df['IDX'] = col_ids

    for ind in indices:

        indice_p = Indice.objects.get(nombre=ind)

        for fecha in fechas:

            col_data = []

            # por cada iteración creamos una fila en el df:

            for p in Parcela.objects.all():

                # la media de todos sus índices

                p_indices = Pixel.objects.filter(parcela=p)

                list_values = []

                for p_itr in p_indices:

                    qs = Mirar_Indice.objects.get(
                        pixel = p_itr,
                        indice = indice_p,
                        fecha = fecha # fecha para la columna:
                    )

                    import math

                    if ( not math.isnan(qs.valor) ):

                        list_values.append(qs.valor)

                # estadistico:añadir columna

                res = func(list_values)

                col_data.append(res)

            df[func.__name__ + '_' + ind + '_' + str(fecha)] = col_data

    return df
```

## Selección de features

El diseño de la vista minable es complicado; para ello trabajo con el equipo seleccionando los campos de la BD más importantes que formarán la tabla. Nos centramos principalmente en los índices vegetativos registrados en diferentes fechas.

Para poder predecir la producción de cultivo, incluimos en la vista varios estadísticos, como la media y el sumatorio para cada índice en cuatro fechas representativas de la evolución de este. El diseño de esta estructura proviene de un estudio previo realizado por el equipo en el que se ha concluido que ciertos estadísticos en unas fechas calibradas funcionan bien para predecir la producción.

Hablábamos anteriormente de entornos virtuales, por el momento la aplicación utiliza un solo entorno virtual Python, pero este punto del trabajo podría hacer uso de un entorno separado con las librerías de ciencia de datos instaladas, de tal forma que separaríamos dos procesos importantes. Hay que destacar que en el trabajo me estoy encargando del proceso para hacer pruebas, pero serán diferentes miembros del equipo los que usen estas tecnologías: para que podamos seguir escalando a medida que evolucione la aplicación, estos miembros tienen que poder utilizar los entornos creados.

## **Variable Objetivo**

Junto con los datos extraídos añadiremos la variable objetivo, "producción". Esta variable, aún sin almacenar en el modelo relacional, hace referencia a los kilos de producción de cada parcela en su cosecha, se obtiene a partir de históricos de campañas anteriores. A partir de ella, junto con las features anteriores, vamos a poder predecir la evolución del cultivo con precisión en estados concretos de su maduración. Se explica con detalle en el siguiente punto.

## **8.4 Integración continua**

---

Una vez hemos identificado los pasos atómicos que queremos ejecutar de forma secuencial, utilizamos la herramienta de integración continua, Jenkins, para ejecutar estos pasos de forma ordenada y ver poco a poco el proceso.

El siguiente script reproduce el flujo de uso común de la aplicación desde su creación, pasando por todos los pasos de carga de datos en el modelo, hasta la vista minable sobre la cual se obtiene el modelo de producción de kg de cultivo.

La ejecución del siguiente pipeline tiene lugar tras lanzar la máquina Vagrant sobre la que corre el proyecto. En el comienzo de la memoria habíamos hablado de que dicha máquina se provisiona con los scripts necesarios para instalar las dependencias del proyecto. En este punto con la máquina en dicho estado, los primeros pasos del pipeline descargan la última versión estable del código del proyecto y crean un entorno virtual que hace uso de las librerías correctas.

Los siguientes pasos ejecutan uno a uno los comandos definidos anteriormente, cargando así en la base de datos los datos provenientes de las distintas fuentes (en este caso excels) y procesando su vista minable. El proceso se realiza de forma ordenada para las 25 parcelas de muestra con las que estamos trabajando, remarcamos que el proceso tiene una alta escalabilidad debido a que podemos reproducir aquellos 'stages' necesarios en los momentos que nos interese. Es decir, podemos comenzar con un pipeline para insertar la información de 25 parcelas, pero posteriormente podríamos añadir más pasos a medida que dispongamos de más información. Por otra parte, la infraestructura como código que conseguimos posibilita definir el flujo concreto que un cliente puede tener con la aplicación.

```
pipeline {
    agent any

    stages {
        stage('git') {
            steps {
```

```

        git url: '/home/vagrant/agrai/.git'
    }
}
stage('env') {
    steps {
        sh('python3.10 -m venv ./tfq_venv')
        sh('./tfq_venv/bin/activate')
        sh('pip install -r ./unix-dep/requirements.txt')
        sh('pip install --upgrade numpy')
    }
}
stage('migrate') {
    steps {
        sh('python3.10 manage.py makemigrations core')
        sh('python3.10 manage.py migrate')
        sh('python3.10 manage.py flush --no-input')
    }
}
stage('sources') {
    steps {
        sh("""python3.10 manage.py 1-load-qgis
            --parcelas "data/parcelas/25/parcelas.shp"
            --pixels "data/parcelas/25/pixels.shp" """)
        sh("""python3.10 manage.py 2-load-parcela-data
            -xls "data/excels/datos-test-2.xls" """)
        sh("""python3.10 manage.py 3-load-cultivos
            -xls "data/excels/datos-test-2.xls" """)
    }
}
stage('sentinel') {
    steps {
        sh("""python3.10 manage.py 4-download_img
            -a "20220601"
            -b "20220701"
            -p "data/parcelas/25/pixels.shp"
            -i ndvi ndre""")
    }
}
stage('indices') { steps { sh('python3.10 manage.py 5-load_range
            -i ndvi ndre') } }

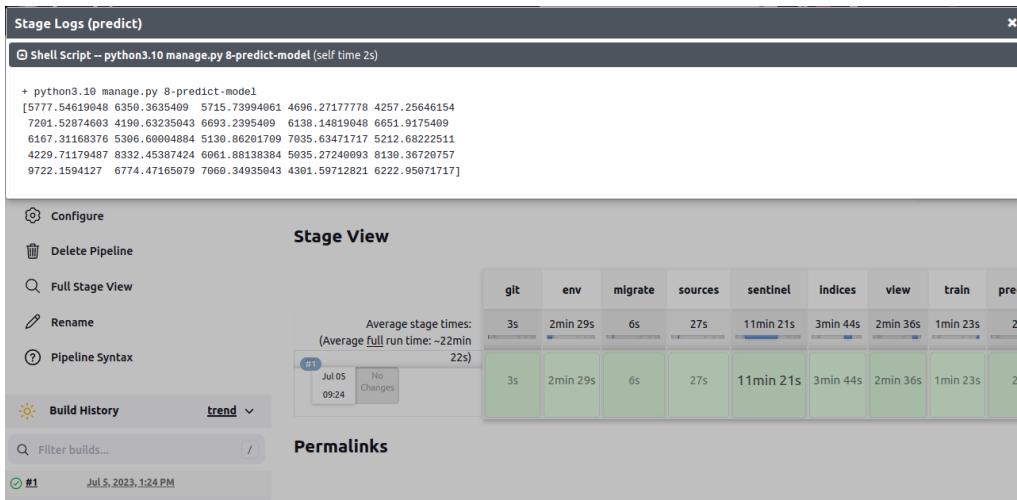
stage('view') { steps { sh('python3.10 manage.py 6-create-view') } }

stage('train') { steps { sh('python3.10 manage.py 7-create-mode') } }

stage('predict') { steps { sh('python3.10 manage.py 8-predict-model') } }
}
}

```

Además de este último, podemos crear varios entornos con tuberías diferentes que den como resultado bases de datos con estados distintos. Como mencionábamos anteriormente, esto es muy útil debido al funcionamiento del equipo en relación con el procesamiento de datos con varios clientes. La aplicación puede trabajar con instancias diferentes de la misma base de datos dependiendo del proyecto en el que se encuentre; por ejemplo, la misma instancia de la BD sirve tanto para una bodega con variedades de vino como para una cooperativa que contempla varios cultivos como guisantes, olivas, peras, etc. Por ello, definir scripts como el anterior utilizando los mismos comandos con distintas fuentes de datos es una forma clara y ordenada de automatizar los procesos.



Ejecución del pipeline implementado desde la herramienta de integración continua Jenkins

Por último, mencionar que en los últimos pasos del pipeline se entrena el mejor modelo y se predice a modo de ejemplo la cantidad de kg de cultivo que van a ser cosechados para las 25 parcelas con las que estamos trabajando.

## 9 MODELO IA

### 9.1 Preprocesamiento

Como los datos con los que trabajamos en los cuadernos provienen de una base de datos estructurada y previamente estudiada, no va a hacer falta un paso de preprocesamiento previo. De todas formas, sí que intentaremos que los datos de producción con los que vamos a predecir estén normalizados y que se haya hecho una búsqueda previa de 'outliers'.

Buscamos un modelo de regresión para el número de Kg de cultivo en cada parcela. Aunque los datos son representativos de un parcelario pequeño, diseñaremos los cuadernos para automatizar la búsqueda del mejor modelo con cualquier volumen de datos que podamos necesitar más adelante.

El rango de valores para la mayoría de los índices que obtenemos va de -1 a 1. Los datos de índices negativos que llegan a la vista se pueden considerar como 'outliers'. Esto es debido a que, si los índices provienen de una imagen con nubes, no se diferencian los colores del terreno y el índice acaba teniendo un valor muy malo. Para evitar valores corruptos eliminaremos los índices que provengan de la toma de la imagen en una fecha en la que había nubes.

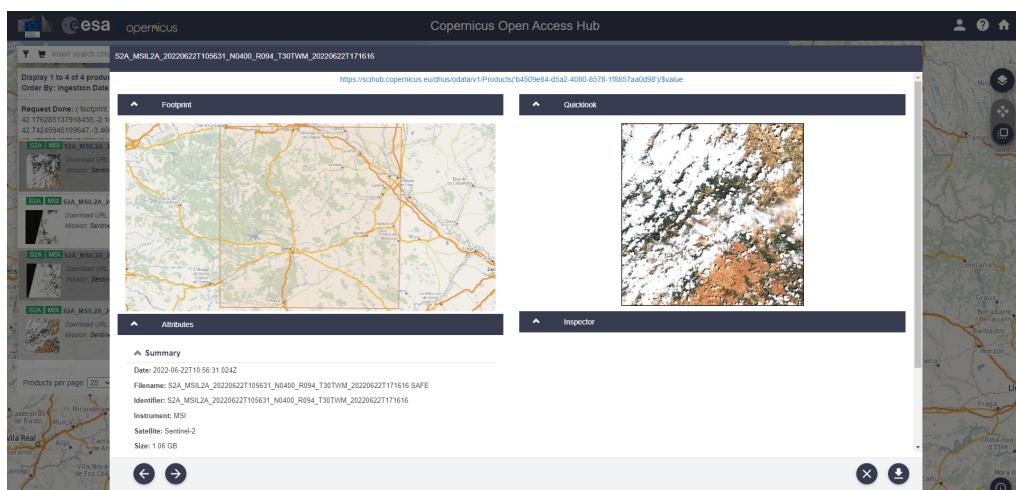


Imagen satelital con nubes para indicar cómo afectan a los índices vegetativos calculados en la vista minable.

El [siguiente repositorio](#) contiene los cuadernos necesarios para la gestión y automatización de los modelos de producción generados a partir los datos de las 25 parcelas estudiadas. Recordamos que estamos pensando en escribir el código necesario para automatizar en la medida de lo posible la búsqueda del mejor modelo de producción, de tal forma que cuando el sistema contemple un mayor número de parcelas, la ejecución de los cuadernos nos siga generando el mejor modelo de producción posible.

### 9.2 Regresión

La vista minable que hemos obtenido representa los datos necesarios para crear un modelo de producción para el cultivo. El [próximo entregable](#) es el desarrollo de uno o varios cuadernos de jupyter con modelos de inteligencia artificial para la tabla anterior.

IDX	mean_ndvi_2022-06-02	mean_ndvi_2022-06-17	mean_ndvi_2022-06-07	mean_ndvi_2022-06-27	mean_ndvi_2022-06-22	mean_ndvi_2022-06-12	mean_ndre_2022-06-02	mean_ndre_2022-06-17	mean_ndre_2022-06-07	mean_ndre_2022-06-27	mean_ndre_2022-06-22	mean_ndre_2022-06-12	mean_ndre_2022-06-02
0	3942	0.007476	0.212770	0.180844	-0.007705	0.038478	0.219896	-0.017073					
1	4053	0.008709	0.202058	0.165046	-0.005460	-0.019584	0.198812	-0.007738					
2	4196	0.204576	0.192939	0.156745	0.001065	-0.011448	0.186642	0.226347					
3	4422	-0.001838	0.198717	0.181343	-0.009038	-0.002894	0.215428	-0.036878					
4	4424	-0.046762	0.205482	0.185708	-0.010250	-0.005377	0.218237	-0.094255					
5	3882	0.000592	0.224865	0.189626	-0.004974	0.062246	0.233515	-0.018172					
6	4732	0.256245	0.177747	0.141133	-0.000639	-0.054281	0.170076	0.252110					
7	3553	0.036759	0.119855	0.116412	-0.006666	-0.002644	0.106559	0.028544					
8	3582	0.115627	0.216081	0.176353	0.004865	-0.002795	0.209227	0.091589					
9	3914	0.172558	0.315941	0.291676	0.005314	-0.028690	0.364501	0.127867					
10	3998	-0.048620	0.227125	0.203715	-0.006980	0.011323	0.253654	-0.082753					
11	4156	-0.022582	0.206010	0.177684	-0.004524	0.016386	0.215706	-0.059363					

Tabla con la vista minable calculada para un rango de fechas concretas.

En el punto en el que estamos, podemos pensar en las múltiples vistas que se pueden generar a partir de los datos persistidos en la BD. Desde aquí desarrollaremos la automatización de la búsqueda del mejor modelo posible para la vista anterior, pero siempre teniendo en cuenta que a partir de los datos almacenados podemos predecir muchos otros valores, no solo la producción de un cultivo.

Resaltamos esta idea debido a que es el punto más fuerte del trabajo realizado. La arquitectura obtenida para la aplicación nos permite generar diferentes vistas con el objetivo de utilizar la información en diversos estudios. Los índices vegetativos y la información almacenada sobre los cultivos se pueden presentar en el formato necesario que el estudio requiera.

Por último, deberíamos poder utilizar este modelo obtenido dentro de la arquitectura de la aplicación para predecir, cómo evolucionan los cultivos que se están monitorizando. El punto en el que se encuentra la arquitectura soporta casi de forma directa la inclusión de las predicciones de la aplicación.

## **9.3 Mejor modelo**

Tras conseguir los datos con la forma correspondiente, exponemos los modelos de inteligencia artificial que vamos a utilizar para buscar el mejor modelo de producción de cultivo (kg) posible.

El problema que estamos buscando resolver es un problema de regresión, en el que conocemos la variable objetivo por el estudio de campañas anteriores. Los clientes, en este caso agricultores, quieren poder predecir cuántos kg de cultivo van a poder cosechar al final de la temporada (septiembre), o en estados previos de maduración (enero).

Los cuadernos que conforman el entregable contienen código que busca el mejor modelo posible para los datos importados con los índices y sus estadísticos. Los índices vegetativos registrados en diferentes fechas, junto con la variable de producción de campañas anteriores, tienen que poder predecir con precisión la campaña actual.

```
# drop columns with negative mean.  
  
for col in df.columns:  
    df2 = df[col].mean()  
  
    if (df2 < 0):  
        df = df.drop(columns=[col], axis=1)
```

Antes de ejecutar código para los modelos, analizamos las relaciones entre las columnas y hacemos estudios de correlación y análisis de componentes principales, para ver si algunos datos de fechas que han llegado a este punto no añaden valor a la información que queremos predecir. Recordemos que puede que lleguen índices provenientes de imágenes tomadas entre nubes y entonces la información empeora el modelo.

Un análisis de componentes principales y la observación de las mejores columnas son el siguiente paso a la limpieza de columnas que se ha realizado en el preprocesamiento.

```
rfe = RFE(estimator=DecisionTreeRegressor(), n_features_to_select=5)
model = DecisionTreeRegressor()
pipeline = Pipeline(steps=[('rfe', rfe), ('m', model)])

cv = RepeatedKFold(n_splits=10, n_repeats=3, random_state=1)
n_scores = cross_val_score(pipeline, X, y, scoring='neg_mean_absolute_error', cv=cv,
n_jobs=-1)

print('MAE: %.3f (%.3f)' % (mean(n_scores), std(n_scores)))
```

Para la variable objetivo de producción también hacemos un pequeño análisis. Queremos ver si esta contiene outliers, para ello analizamos la distribución y seleccionamos aquellos valores que se alejan mucho de la media y desviación típica de la muestra. Es importante que los datos con los que vamos a predecir estén limpios para que el modelo obtenga una buena métrica, como puede ser el "mean absolute error", MAE, puntuación para problemas de regresión en el que el resultado está contemplado en las unidades de la variable analizada.

Una vez hemos analizado los datos y hemos limpiado los valores erróneos de los índices, procedemos a probar aquellos modelos de inteligencia artificial que mejor puedan predecir este valor de kg de producción. Los modelos que probamos son los siguientes.

	<b>MODEL</b>	<b>PARAMS</b>	<b>MAE</b>	<b>MSE</b>	<b>RMSE</b>
<b>2</b>	RandomForest	{'rf_max_depth': 7, 'rf_max_features': 'sqrt', 'rf_n_estimators': 5, 'rf_random_state': 42}	2282.013333	1.464411e+07	3826.762121
<b>0</b>	ElasticNet	{'en_alpha': 10, 'en_l1_ratio': 0.7000000000000001, 'en_max_iter': 5}	4668.784757	4.217480e+07	6494.212821
<b>4</b>	KNeighbors	{'knn_metric': 'manhattan', 'knn_n_neighbors': 11, 'knn_weights': 'uniform'}	5109.534545	4.310770e+07	6565.645200
<b>3</b>	MLP	{'mlpr_activation': 'tanh', 'mlpr_alpha': 0.05, 'mlpr_hidden_layer_sizes': (100, 50, 30), 'mlpr_learning_rate': 'adaptive', 'mlpr_max_iter': 50, 'mlpr_solver': 'sgd'}	4599.475956	4.323998e+07	6575.711553
<b>1</b>	SVR	{'svr_C': 5.0, 'svr_gamma': 0.01, 'svr_kernel': 'linear'}	4129.801479	4.749271e+07	6891.495313

La tabla anterior muestra ciertas métricas para los algoritmos probados. Al ordenar la tabla por el RMSE vemos qué algoritmo predice un poco mejor y observamos aquellos hiperparámetros que mejor score obtienen.

```
df_sorted = df.sort_values(by='RMSE', ascending=True, na_position='first');

best_models = []

for index, row in df_sorted.iterrows():

    model_str, model_prm= row["MODEL"], row["PARAMS"]
    model_type = getattr(sys.modules[__name__], model_str)
    rm_pre_parms = rm_dict_pre(model_prm) # modify dict keys
    best_model = model_type(**rm_pre_parms)
    best_models.append(best_model)

best_estimators = []

for idx, est in enumerate(best_models):

    best_estimators.append((est.__class__.__name__, est))
```

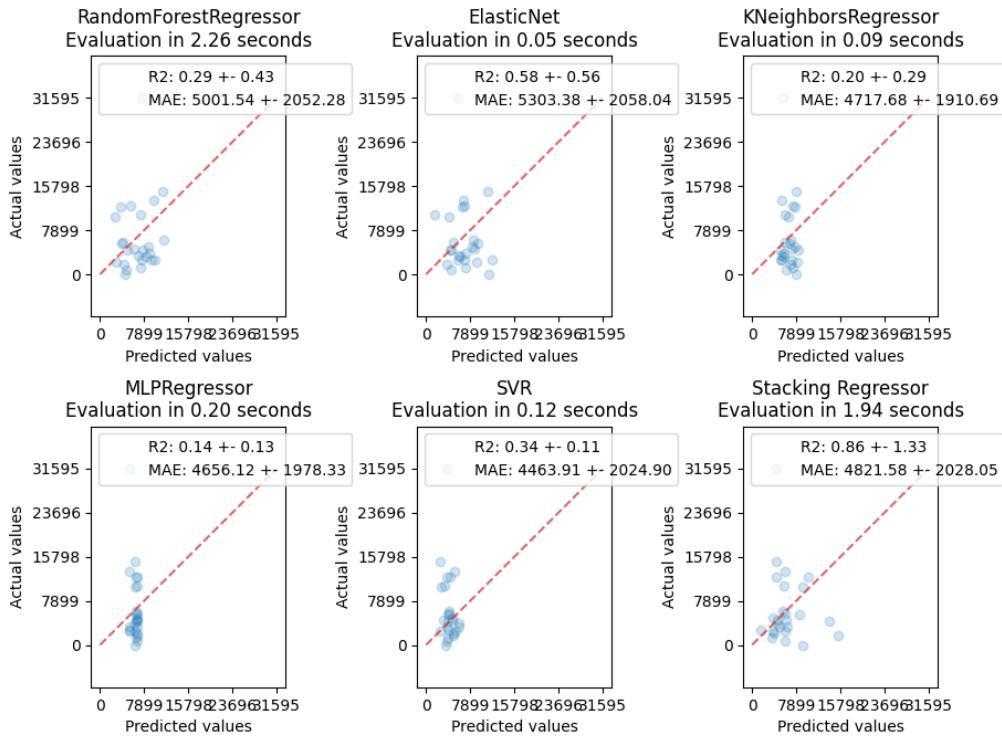
Como cada modelo aporta un punto de vista diferente, vamos a utilizar los resultados obtenidos para generar un modelo ensemble que generalice las predicciones de los anteriores. Utilizaremos las métricas obtenidas para establecer el orden en el que construir un modelo 'stacking'. Como estimador final nos quedamos como el mejor modelo obtenido al ordenar la tabla anterior por el valor del RMSE obtenido para cada uno de ellos, es decir, el RandomForestRegressor.

```
reg = StackingRegressor(
    estimators=estimators,
    final_estimator=estimators[0][1])

reg.fit(x, y)

y_pred_e = reg.predict(x)
```

A continuación mostramos una gráfica para cada modelo en la que vemos la relación entre los valores actuales y los valores predichos por el modelo. Esta pequeña visualización nos permite ver si cada algoritmo sobreajusta a los datos de entrenamiento o realmente tiene el rango que necesitamos. El punto más interesante es el análisis del modelo 'stacking' en el que vemos que es capaz de inferir los estilos de los algoritmos anteriores.



Comparación de los algoritmos. Valores reales vs. los predichos por cada uno de los modelos utilizados.

Hemos analizado cada algoritmo con la intención de ver si modeliza correctamente nuestro problema de regresión, pero no es nuestro objetivo final que el modelo prediga con exactitud: solo queremos dejar automatizada una pequeña visualización de las relaciones entre los índices para así poder intuir qué modelo nos puede funcionar mejor. De hecho, los datos representativos con los que trabajamos nos impiden buscar un modelo real preciso.

## 9.4 Predicciones

---

Nuestro último paso en el pipeline mostrado anteriormente es la capacidad para utilizar el modelo y predecir qué cantidad de kg se pueden cosechar cuando se pasa una nueva parcela al sistema. A continuación mostramos cómo el último 'stage' del pipeline utiliza el modelo que acabamos de entregar para predecir sobre la vista minable.

```
filename = 'script/finalized_model.sav'

df_x = pd.read_csv('script/view.csv')

loaded_model = pickle.load(open(filename, 'rb'))

y_pred_e = loaded_model.predict(df_x)

print(y_pred_e)
```

En este punto, deberíamos poder utilizar este modelo obtenido dentro de la arquitectura de la aplicación para predecir, cómo evolucionan los cultivos que se están monitorizando. El punto en el que se encuentra la arquitectura soporta casi de forma directa la inclusión de las predicciones de la aplicación.

```
Stage Logs (predict)
Shell Script -- python3.10 manage.py 8-predict-model (self time 2s)

+ python3.10 manage.py 8-predict-model
[5777.54619048 6350.3635409 5715.73994061 4696.27177778 4257.25646154
 7201.52874603 4190.63235043 6693.2395409 6138.14819048 6651.9175409
 6167.31168376 5306.60004884 5130.86201709 7035.63471717 5212.68222511
 4229.71179487 8332.45387424 6061.88138384 5035.27240093 8130.36720757
 9722.1594127 6774.47165079 7060.34935043 4301.59712821 6222.95071717]
```

*Predicción de la cantidad de Kg de cultivo que se va a cosechar. Resultado tras la ejecución del último 'stage' en el pipeline.*

# 10 SEGUIMIENTO

---

## 10.1 Desviaciones

---

El proyecto se ha desarrollado sin grandes cambios relativos a la planificación estipulada, pero a pesar de realizar la planificación del proyecto con su desglose en entregables y horas para los paquetes de trabajo, el tiempo de desarrollo ha diferido de la misma en algunas ocasiones, especialmente en la memoria.

El número de horas planificadas era de 300, las cuales se repartían en varias etapas marcadas por los diversos puntos de control: planificación, análisis, infraestructura, diseño, implementación, modelos y seguimiento y control.

No ha habido grandes desviaciones en cuanto a la planificación estipulada en el desglose de tareas. Sí que podemos remarcar que se han dedicado más horas a la redacción de la memoria y a algunas tareas secundarias relacionadas con esta, como pueden ser la integración de los diversos diagramas que aparecen la automatización en Github Actions. Podemos estimar entre 15 o 20 horas más para esta tarea sobre las 20 estimadas al principio.

## 10.2 Memoria

---

Esta memoria se incluye como un entregable del proyecto debido a las explicaciones y esquemas que contiene (principalmente para el modelo relacional). Para poder mejorar el flujo de trabajo, ciertos conceptos necesitan ser entendidos por los desarrolladores. La documentación aquí expuesta tiene un gran valor por el lenguaje ubicuo que implanta en el proyecto.

El documento se ha escrito como notas en *Markdown* por la facilidad que da para incluir fragmentos de código relativos al proyecto y poder explicar decisiones y conceptos. Para integrar los diagramas expuestos utilizamos tecnologías como *Mermaid*, que posibilita editar los diagramas directamente en el documento, algo totalmente necesario cuando estos modelos se modifican a medida que avanzamos.

Finalmente, exportamos las notas en el orden necesario para la memoria y renderizamos el documento en *HTML*. El fácil acceso al documento tiene bastante relevancia por el uso que el equipo puede darle al consultar la documentación, especialmente la del modelo relacional y los servicios.

Utilizamos algunas herramientas de integración continua como *Github Actions* para automatizar este proceso. En un primer lugar convertimos el grafo de notas *Markdown* a un documento *HTML* mediante un proceso recursivo que ordena las notas y genera así la memoria del proyecto. Además, el índice se crea a partir de los nombres de las notas y sus relaciones con las que están enlazadas.

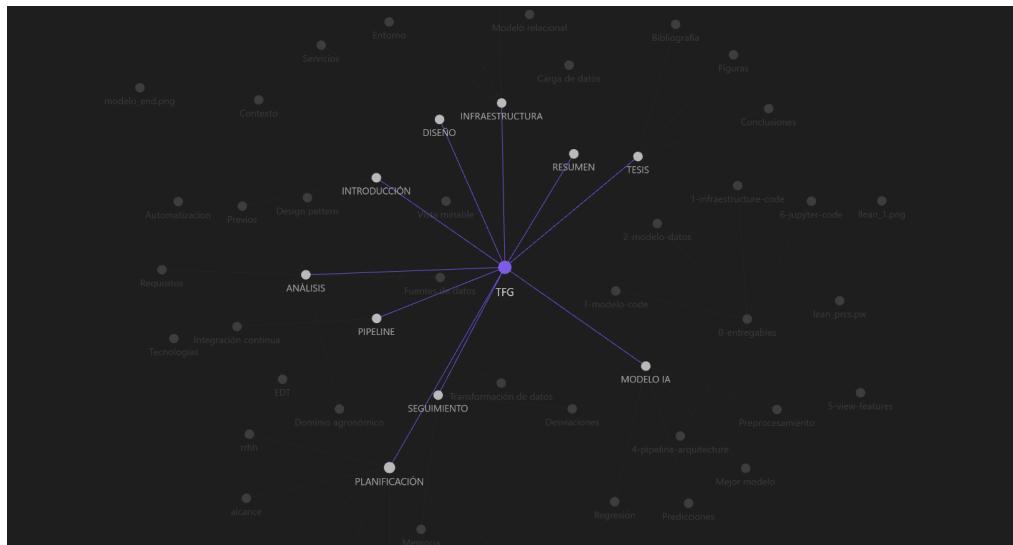


Imagen del grafo asociado a las notas en las que se ha escrito esta memoria.

El principal interés de que la redacción se realice de esta forma es que al comienzo del proyecto no había una idea clara de cómo se tenían que desarrollar los contenidos. Escribir un capítulo o tema en una nota es algo más sencillo (por lo tanto, se puede terminar) que pensar desde el principio dónde colocar esta sección que quiero escribir en la memoria. Utilizando Obsidian como gestor de las notas, podemos definir relaciones entre ellas simplemente conectando palabras clave, de esta forma conseguimos un esquema mental del proyecto visualizando el grafo generado. Finalmente, a medida que nos acercamos al cierre del proyecto, escribimos en un notebook una serie de algoritmos que nos permiten presentar el resultado de la memoria de la forma que deseamos. En ellos implementamos los pasos necesarios para transformar cada nota a HTML, integrar aquellos *snippets* de código que pueda contener, representar tablas y diagramas, etc. Muy importante, el grafo permite generar los niveles de indentación asociados a cada capítulo o sección del documento. A continuación, mostramos algunas de las funciones utilizadas en el cuaderno para obtener el resultado que está leyendo.

```
def rec_note_prcs(G, start='TFG', ident='1'):

    """ recursive note building """

    ident = ident + '.0'

    neigh = list(G.neighbors(start))

    if (len(neigh) == 0):

        return note_decorate(start,ident[:-2])

    end_doc = note_decorate(start,ident)

    for con in neigh:

        ident = ident[:-1] + str(int(ident[-1])+1)
        end_doc = end_doc + rec_note_prcs(G, con, ident)

    return end_doc
```

Estas celdas de código son parte del proceso de transformación de Markdown a HTML. Aunque existen librerías con funciones para realizar dicha transformación, necesitamos realizar un gran trabajo de procesado para que el documento enlace las notas en el orden correcto y se tengan en cuenta conceptos como la indentación en

los distintos apartados. El grafo que conforman las notas pasa a una estructura de árbol en la que es fácil crear enlaces entre los puntos a los que pertenece cada nota.

```
def mermaid_process(text):

    """ transforms any mermaid diagrams to images in png format """

    global parent_directory

    result = re.findall(```mermad(?s:.*?)```', text, re.M)

    if len(result) <= 0:

        return text

    sub_text = result[0][len(```mermad'):-3]

    mermaid_dir = os.path.join(parent_directory, "mermad")

    text = re.sub(
        ```mermad(?s:.*?)```',
        mmd_to_img(sub_text, mermaid_dir) ,
        text, 1
    )

    return mermaid_process(text)
```

Estas notas se encuentran en un repositorio separado al proyecto para poder así construir la memoria. Utilizamos *Github Actions* para definir las acciones necesarias cuando se realicen cambios en el repositorio, es decir, a medida que escribimos en las notas. En primer lugar, queremos conseguir que el documento se renderice al hacer un *commit* en el repositorio y se publique posteriormente dicha memoria en una pequeña web bajo el dominio asociado a nuestro usuario de *Github Pages* <https://alesteba.github.io/tfg/>. En segundo lugar, una vez tenemos publicada la memoria, una segunda acción convertirá el código HTML asociado al documento en el archivo PDF que está leyendo actualmente.

```
build_notes:

  name: Building web notes from graph
  runs-on: ubuntu-latest

  steps:

    - name: Checkout
      uses: actions/checkout@v3

    - name: Setup Python
      uses: actions/setup-python@v3
      with:
        python-version: '3.9'
        cache: 'pip'

    - name: Install Dependencies
      run: pip install -r requirements.txt

    - name: Install Jupyter
      run: sudo -H pip install jupyter

    - name: Install Mermaid Converter
      run: |
        npm install -g @mermaid-js/mermaid-cli

    - name: Run Script and Build Notes
```

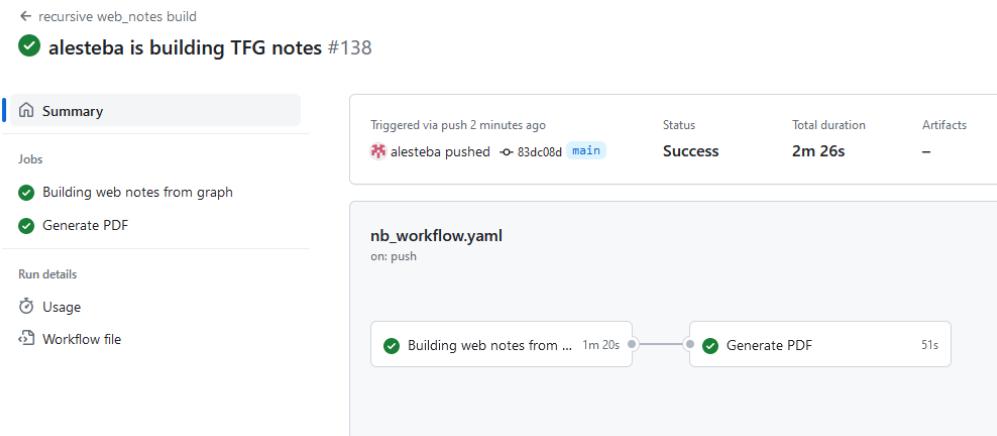
```

run: |
  cd _k.pcs.zen/
  jupyter nbconvert --to script ./automate_md.ipynb
  python ./automate_md.py

- name: Pushes to another repository
  uses: cpina/github-action-push-to-another-repository@main
  env:
    API_TOKEN_GITHUB: ${{ secrets.API_TOKEN_GITHUB }}
  with:
    source-directory: '_k.pcs.zen/publish'
    destination-github-username: 'alesteba'
    destination-repository-name: 'tfg'
    user-email: alesteba@unirioja.es
    target-branch: main

```

El código anterior se corresponde con la infraestructura como código necesaria para transformar las notas de Markdown al documento en HTML. Como podemos observar en los pasos se utiliza una máquina Ubuntu en la que se instalan librerías necesarias para transformar los diagramas en imágenes y poder trabajar con cuadernos Jupyter. Posteriormente, se transforma el cuaderno con el código de conversión de las notas en un único script, se ejecuta y se publican los resultados en otro repositorio diferente, el cual mediante *Github Pages* crea la página HTML correspondiente con la memoria final.



Ejemplo construcción de la memoria desde Github Actions.

## 11 CONCLUSIONES

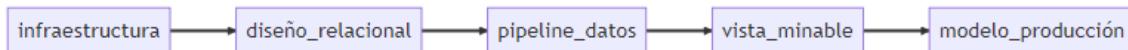
---

Este trabajo se ha centrado en el desarrollo de un proceso de automatización que permite la integración continua del flujo de trabajo del equipo para una aplicación de gestión agronómica. Hemos pasado por casi todas las etapas de diseño e implementación, además de la final inclusión de los modelos predictivos para la producción de cultivo.

Durante el proceso hemos influido en los principales bloques de diseño e implementación relativos a la arquitectura de la aplicación, pasando por la infraestructura, los modelos de datos, la carga de estos datos, y la búsqueda de modelos predictivos a partir de ellos. Acorde al alcance planteado, hemos llegado satisfactoriamente a todos los puntos estipulados; recordamos estos puntos:

- Obtenemos un script de creación para el entorno necesario.
- Diseñamos un modelo relacional de datos que permita un rápido escalado de la aplicación.
- Automatizamos el proceso de transformación de los datos mediante el diseño de un pipeline.
- Diseñamos un algoritmo que exporte la vista minable necesaria para predecir la producción.
- Utilizamos técnicas de IA para entrenar el mejor modelo de regresión para la producción anterior.
- Implementamos el modelo para predecir nuevos valores que no han sido usados en el entrenamiento.

Para finalizar, volvemos a los conceptos de los procesos 'lean' que mencionábamos al principio. La fácil comprensión que permiten los esquemas de datos aquí planteados dotan al equipo de trabajo del lenguaje ubicuo necesario para que el software pueda crecer y ser escalado, evitando bloqueos por falta de comprensión. El lenguaje que maneja el equipo está directamente relacionado con su buen funcionamiento y garantiza el proceso de optimización continua. Los artefactos obtenidos y la explicación de los modelos quedan así como documentación interna que debe ser mantenida al mismo nivel de importancia que el software que la implementa.



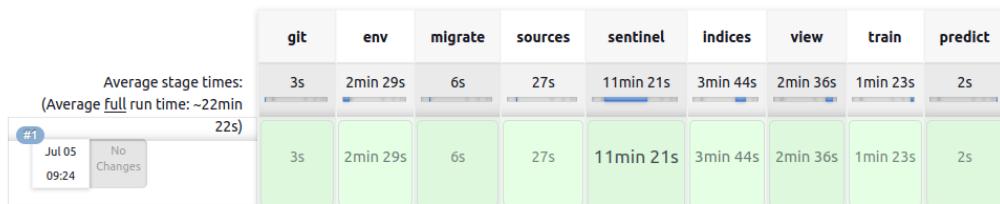
En cuanto a la aplicación, se ha conseguido terminar con un producto funcional que cumple con los requisitos planteados. El proceso de arquitectura conseguido facilita el uso de múltiples datos provenientes de las fuentes estudiadas. El esquema de datos relacional implementado soporta la carga de los índices vegetativos, así como el resto de información relativa al parcelario. Los bloques atómicos de carga diseñados permiten planificar y automatizar la ingestión de los datos provenientes de las fuentes descritas. Por último, la capacidad para utilizar modelos de regresión sobre los datos dota a la arquitectura de la capa necesaria para realizar predicciones y configurar los estudios necesarios. Hemos conseguido crear las bases de un proceso que tiene un gran potencial para monitorizar y proporcionar una adecuada toma de decisiones sobre conceptos agronómicos.

## Pipeline tfg

[Add description](#)

[Disable Project](#)

### Stage View



## **12 BIBLIOGRAFÍA**

---

- Planificación

<https://www.pmi.org/pmbok-guide-standards/foundational/pmbok#>

<https://mermaid.js.org/syntax/gantt.html>

<https://www.pmi.org/learning/library/applying-work-breakdown-structure-project-lifecycle-6979>

- Infraestructura

[https://en.wikipedia.org/wiki/Infrastructure\\_as\\_code](https://en.wikipedia.org/wiki/Infrastructure_as_code)

<https://www.jenkins.io/>

<https://www.vagrantup.com/>

- Diseño modelo relacional

<https://mermaid.js.org/syntax/entityRelationshipDiagram.html>

<https://docs.djangoproject.com/en/4.1/intro/tutorial01/>

<https://docs.djangoproject.com/en/4.1/topics/db/queries/>

<http://gorodinski.com/blog/2012/04/14/services-in-domain-driven-design-ddd/>

Domain Driven Design - Eric Evans

- Automatización Pipeline

<https://simpleisbetterthancomplex.com/tutorial/2018/08/27/how-to-create-custom-django-management-commands.html>

<https://medium.com/@bonnotguillaume/software-architecture-the-pipeline-design-pattern-from-zero-to-hero-b5c43d8a4e60>

<https://www.astera.com/es/type/blog/etl-pipeline-vs-data-pipeline/>

<https://www.jenkins.io/doc/book/pipeline/syntax/>

- Modelo IA

<https://stackabuse.com/random-forest-algorithm-with-python-and-scikit-learn/>

<https://towardsdatascience.com/random-forest-regression-5f605132d19d>

[https://scikit-learn.org/stable/auto\\_examples/miscellaneous/plot\\_pipeline\\_display.html#sphx-glr-auto-examples-miscellaneous-plot-pipeline-display-py](https://scikit-learn.org/stable/auto_examples/miscellaneous/plot_pipeline_display.html#sphx-glr-auto-examples-miscellaneous-plot-pipeline-display-py)

[https://scikit-learn.org/stable/modules/cross\\_validation.html](https://scikit-learn.org/stable/modules/cross_validation.html)

[https://scikit-learn.org/stable/auto\\_examples/ensemble/plot\\_stack\\_predictors.html#sphx-glr-auto-examples-ensemble-plot-stack-predictors-py](https://scikit-learn.org/stable/auto_examples/ensemble/plot_stack_predictors.html#sphx-glr-auto-examples-ensemble-plot-stack-predictors-py)

- Redacción memoria

<https://networkx.org/>

<https://github.com/obsidianmd/obsidian-api>

<https://docs.github.com/en/actions>

