



UNIVERSIDAD DE BURGOS
ESCUELA POLITÉCNICA SUPERIOR
Grado en Ingeniería Informática



**TFG del Grado en Ingeniería
Informática**

**título del TFG
Documentación Técnica**



Presentado por nombre alumno
en Universidad de Burgos — 29 de mayo
de 2020

Tutor: nombre tutor

Índice general

Índice general	I
Índice de figuras	III
Índice de tablas	IV
Apéndice A Plan de Proyecto Software	1
A.1. Introducción	1
A.2. Planificación temporal	1
A.3. Estudio de viabilidad	6
Apéndice B Especificación de Requisitos	9
B.1. Introducción	9
B.2. Objetivos generales	9
B.3. Catalogo de requisitos	9
B.4. Especificación de requisitos	9
Apéndice C Especificación de diseño	11
C.1. Introducción	11
C.2. Diseño de datos	11
C.3. Diseño procedimental	11
C.4. Diseño arquitectónico	11
Apéndice D Documentación técnica de programación	13
D.1. Introducción	13
D.2. Estructura de directorios	13
D.3. Manual del programador	13

D.4. Compilación, instalación y ejecución del proyecto	13
D.5. Pruebas del sistema	13
Apéndice E Documentación de usuario	15
E.1. Introducción	15
E.2. Requisitos de usuarios	15
E.3. Instalación	15
E.4. Manual del usuario	15
Bibliografía	17

Índice de figuras

Índice de tablas

Apéndice A

Plan de Proyecto Software

A.1. Introducción

En este punto pasamos a detallar la planificación del proyecto y los pasos que se han dado para poder completarlo con éxito.

A.2. Planificación temporal

La planificación y el desarrollo de este proyecto se han llevado a cabo siguiendo una metodología ágil, concretamente scrum, teniendo que en cuenta que la plataforma de desarrollo usada, GitLab, presenta una nomenclatura propia que no siempre coincide con la de scrum. Así, por ejemplo, GitLab usa *milestones* o *issues* para referirse a *sprints* o *product backlog items (PBIs)*. En este documento intentaremos usar los terminos de scrum siempre que sea posible.

En relación a la plataforma, otro punto a reseñar es la no disponibilidad de gráficas de *burndown* en la versión gratuita. Se ha intentado usar otras alternativas, como **Screenful**, pero requiere acceso completo a la API de GitLab. Estando el proyecto alojado en la cuenta privada de la empresa HP SCDS, esta no es una opción viable. Así pues, la realización del proyecto se ha llevado a cabo sin disponer de estos gráficos.

El proyecto se ha llevado a cabo siguiendo un sprint de dos semanas, al final de las cuales, en reunión con los tutores, tanto por parte de la Universidad de Burgos como por parte de HP SCDS, se mostraba el trabajo realizado, se proponían cambios o mejoras, y se establecían los objetivos

para el siguiente sprint. El seguimiento de los PBIs se ha realizado a través del tablero Kanban que GitLab provee por defecto.

Pasamos a ver con más detalle el trabajo realizado en cada uno de los sprints:

01 - Arranque del proyecto

[10/02/2020 -- 17/02/2020]

El sprint de arranque del proyecto es fundamental teórico, dedicado a estudiar y decidir como vamos a hacer el proyecto, con que herramientas y que medidas usaremos para determinar la validez de un modelo. Este sprint inicial solamente duró una semana, para de esta forma, tras la primera reunión con los tutores el 17 de febrero, poder comenzar un nuevo sprint.

- Estudio comparativo de librerías de machine learning: a la hora de abordar la realización del proyecto nos encontramos con multitud de librerías disponibles. Como primera tarea seleccionamos varias de ellas como posibles candidatas y procedemos a estudiar sus características distintivas.
- Análisis estadístico: en esta tarea nos dedicamos a estudiar métricas alternativas a la precisión para valorar nuestro modelo, como la sensibilidad y la especificidad, así las curvas ROC *receiver oOperating characteristic*, o *característica operativa del receptor*.
- Estudio documentación de un proceso experimental con CRISP_DM: aquí aprendemos una metodología habitual a la hora de trabajar en ciencias de datos que usaremos durante el proyecto.
- Elección de librerías: una vez disponemos de los datos necesarios, valoramos como cada librería se adapta a las necesidades del proyecto y seleccionamos una, Pytorch.
- Estudio teórico del perceptrón multicapa: como última tarea del sprint, estudiamos el modelo básico de red neuronal profunda, el perceptrón multicapa.

02 - Desarrollo del perceptrón

[02/03/2020 – 16/03/2020]

En este sprint nos enfocamos en la parte formal del TFG, estudiando LaTeX y las plantillas de la memoria así como la estructura de archivos, y en los primeros pasos con Pytorch y el dataset.

- Estudio de LaTeX y plantillas de memoria: estudiamos los documentos a entregar y sus diferentes secciones a la vez que preparamos el entorno para trabajar con documentos LaTeX.
- Añadir a la memoria LaTeX la documentación generada en markdown: pasamos la documentación generada en el primer sprint, escrita en markdown en la wiki del proyecto, a LaTeX y la incorporamos a la memoria.
- Estructurar las carpetas del repositorio y actualizar fichero Readme.md: procedemos a estructurar los archivos de nuestra solución siguiendo ejemplos de proyectos de éxito y otros TFG.
- Aprendizaje de Pytorch: en este punto nos centramos en aprender la librería que vamos a usar.
- Preprocesado de los datos: finalmente estudiamos los datos con los que vamos a trabajar y realizamos una implementación básica de red.

03 - Modelo de perceptrón

[16/03/2020 - 06/04/2020]

Por circunstancias personales imprevistas, este sprint se alarga una semana más. En este tiempo apenas se puede avanzar en el proyecto; solo se completa una tarea y se comienza a trabajar en el modelo básico del perceptrón aunque no se completa.

- Gestión de la configuración del repositorio GitLab, donde limpiamos el repositorio de ficheros temporales de Python y Pytorch

04 - Finalizar perceptrón

[07/04/2020 - 20/04/2020]

En este sprint nos centramos en el desarrollo de los diversos modelos del perceptrón.

- Modelo simple de perceptron: terminamos la tarea comenzada en el sprint anterior donde generamos nuestro modelo base de perceptrón probando con los algoritmos SGD y Adam.
- Entrenar perceptrón con BCEWithLogitsLoss: ampliamos nuestro modelo inicial para usar la función de pérdida BCEWithLogitsLoss.
- Mejorar modelo de perceptrón: terminado nuestro modelo base, pasamos ahora a incluir el filtro gaussiano.
- Añadir SMOTE en el modelo de perceptrón: procedemos a entrenar nuevos modelos de perceptrón usando SMOTE.
- Hacer análisis de frecuencia con la transformada de Fourier: para terminar con los modelos basados en perceptrón, cambiamos el análisis de la intensidad de luz por el de la frecuencia.
- Caso de estudio de documentación TFG en LaTeX relacionado: comprobamos otros TFG para estudiar diferentes enfoques a la hora de encarar la escritura de la memoria.
- Crear notebook para testear los modelos: para finalizar el sprint, creamos un notebook donde poder ejecutar nuestros modelos contra el dataset de testing.

05 - Probar LSTM y crear web

[21/04/2020 - 04/05/2020]

En este sprint nos marcamos como objetivo desarrollar modelos basados en redes LSTM y crear una primera versión de aplicación web para poder ejecutar los modelos.

- Formación en conjuntos de datos desequilibrados: estudiamos algunos métodos adicionales para trabajar con dataset desbalanceados.
- Integración documentación Wiki a memoria LaTeX: migramos el resto de información que teníamos en la wiki a la memoria.
- Crear modelo de red LSTM: desarrollamos varios modelos usando redes LSTMs.
- Crear web para testear modelos: creamos una primera versión donde sea posible subir un archivo para probar la ejecución del modelo.

- Escribir capítulos de Introducción y Objetivos de la memoria: pasamos a completar los capítulos iniciales de la memoria.

06 - Preparar memoria

[05/05/2020 - 18/05/2020]

Estaba previsto que el foco de este sprint fuese trabajar en la memoria pero algunos problemas en las tareas relacionadas con la aplicación web provocan que no haya mucho avance en la documentación.

- Desplegar web en Heroku: desplegamos de forma manual la web en Heroku para comprobar que todo funciona correctamente.
- Crear pipeline de release: después del despliegue manual procedemos a automatizar el proceso generando un pipeline en GitLab.
- Completar web: añadimos algunas páginas de información sobre la misión Kepler, sobre los datos y sobre el proyecto. También la localizamos soportando el idioma inglés.
- Ampliar descripción de las herramientas usadas: añadimos más detalles sobre las herramientas usadas en el proyecto.

07 - Anexos

[19/05/2020 – 31/05/2020]

Entrando ya en las etapas finales, este sprint se dedica a continuar el trabajo en la memoria y a ir puliendo detalles del proyecto en general.

- Documentar proceso experimental: completamos el apartado quinta de la memoria detallando las lecciones aprendidas durante el proyecto.
- Incluir dataset y datos del modelo en la web: ampliamos la web incluyendo un dataset de pruebas, información sobre el modelo desplegado y enlazando al repositorio original de los datos en Kaggle.
- Evaluar la calidad del código: usamos una herramienta para controlar la calidad de nuestro código, Codebeat.
- Añadir badges al repositorio: incluimos algunos badges en el repositorio y en el Readme.md para conocer la calidad del código, como fue la ejecución del pipeline de CD o como se encuentra la web en Heroku.

A.3. Estudio de viabilidad

Abordamos a continuación un breve estudio de la viabilidad tanto legal como económica de este proyecto.

Viabilidad económica

Este proyecto nunca se ha enfocado a una viabilidad económica y sería difícil venderlo *tal cual*. Quizá la razón más importante es que apenas hay demanda; el conocimiento de que una estrella es orbitada o no por exoplanetas es, a día de hoy, de poca utilidad comercial, por lo que hay pocas, por no decir ninguna, empresas dispuestas a pagar por ello.

A parte de las empresas privadas, tenemos a diferentes organismos públicos interesados en la exploración espacial, como la NASA, la ESA, la CNSA o similares. Estas agencias se enfrentan a un grave problema común: hay muchas estrellas con posibles exoplanetas y disponen de pocos recursos para explorarlos todos. Dado el carácter probabilístico de este tipo de modelos, puede tener sentido comercial disponer de diferentes aplicaciones y enfocar sus recursos en estudiar aquellas estrellas seleccionadas como más probables por todos, o la mayoría, de las aplicaciones de las que dispongan. A nivel de modelos, en aprendizaje automático, tenemos los métodos de *ensemble* que, básicamente, consisten en combinar varios modelos base para producir otro modelo óptimo. Así pues, no sería ilógico realizar un proceso similar a nivel de aplicaciones.

Por otro lado, otro colectivo que podría estar interesado en nuestra aplicación sería la comunidad científica, más allá de aquellos que trabajan en las agencias espaciales. La confirmación de la presencia o no de exoplanetas en torno a diferentes estrellas puede ayudar a generar nuevo conocimiento sobre este tipo de sistemas, sirviendo de base para confirmar hipótesis o para descartar teorías erróneas.

Finalmente, tenemos otro colectivo que podría estar interesado en nuestro proyecto, los aficionados a la astronomía. En este caso, se podría complementar la aplicación con alguna funcionalidad más como, por ejemplo, proporcionar las coordenadas de la estrella en la que se han identificado exoplanetas de forma que el aficionado pueda apuntar su telescopio directamente a ella.

Viabilidad legal

En lo referente a la viabilidad legal del proyecto debemos separarla en dos puntos: las herramientas y los datos. Respecto a las primeras, no

se encuentra ningún problema. Todas las herramientas pueden usarse en aplicaciones comerciales, presentando licencias permisivas (BSD o MIT en su mayoría).

En el apartado de los datos si podemos encontrarnos con mayor problema. Los usados para el entrenamiento de estos modelos, pertenecientes a la campaña 3, son de dominio público, esperando solamente ser citados por investigadores. Sin embargo, si deseamos mejorar nuestros modelos con datos más recientes de otras campañas si podemos encontrarnos con datos restringidos o con copyright de forma temporal. Además, en el caso de querer implementar la versión para aficionados a la astronomía, el catálogo de estrellas si tiene copyright.

Apéndice B

Especificación de Requisitos

- B.1. Introducción
- B.2. Objetivos generales
- B.3. Catalogo de requisitos
- B.4. Especificación de requisitos

Apéndice C

Especificación de diseño

- C.1. Introducción
- C.2. Diseño de datos
- C.3. Diseño procedimental
- C.4. Diseño arquitectónico

Apéndice D

Documentación técnica de programación

- D.1. Introducción
- D.2. Estructura de directorios
- D.3. Manual del programador
- D.4. Compilación, instalación y ejecución del proyecto
- D.5. Pruebas del sistema

Apéndice E

Documentación de usuario

- E.1. Introducción
- E.2. Requisitos de usuarios
- E.3. Instalación
- E.4. Manual del usuario

Bibliografía
