

# Práctica 4 - Reto Text-based Personality Computing: Prediciendo Big-5 a partir de textos en Reddit

Marcos Fernández Pichel

Curso 2025-2026 Grado en Inteligencia Artificial - USC

El **Text-based Personality Computing (TPC)** es un área de la inteligencia artificial y la psicología computacional que se centra en analizar y predecir rasgos de personalidad a partir del lenguaje escrito. Utiliza algoritmos de NLP para examinar textos, mensajes o publicaciones en redes sociales, identificando patrones lingüísticos que se correlacionan con dimensiones de la personalidad, como las del modelo Big Five (apertura, responsabilidad, extraversion, amabilidad y neuroticismo). Su objetivo es inferir características psicológicas de manera automática.

En lo que resta de la parte práctica de la asignatura, llevaremos a cabo una competición por equipos<sup>1</sup> con 3 integrantes por equipo (excepcionalmente habrá equipos de 2), cuyo objetivo será desarrollar un predictor de **rasgos de personalidad** a nivel de **usuario** a partir de posts publicados en Reddit.

Para ello, os proporcionamos lo siguiente:

- Un fichero **posts.csv** que contiene los posts de los/as usuarios/as tanto de train como de test.
- Un fichero llamado **authors\_train.csv** que contiene una serie de usuarios para los que tenemos sus valores (del 0 al 100 y con un incremento de 10 en 10) de las 5 dimensiones de personalidad recogidas en el Big-5. Por ejemplo, un 0 de extraversion, significaría que la persona es muy introvertida.
- Un fichero llamado **authors\_test.csv** que contiene una serie de usuarios para los que no tenemos predicciones.

La práctica es completamente libre, el objetivo es estimar las dimensiones Big-5 de los/as usuarios/as de test. Habrá varias entregas de variantes, en la que cada equipo puede enviar entre 2-3 soluciones y se actualizará la leaderboard.

## 1 Posibles estrategias

Cada equipo puede decidir qué estrategia seguir. La idea es aplicar los conocimientos obtenidos en la asignatura e ir probando cómo de bien funciona cada solución en las entregas de variantes.

A continuación se listan una serie de ideas que podríais utilizar, pero no son las únicas:

---

<sup>1</sup>Habrá una *leaderboard*, pero el puesto final no tiene influencia en la nota. Es importante que le pongáis un **nombre** al equipo

- Entrenar un modelo supervisado con los datos de entrenamiento, bien sea basado en representaciones vectoriales clásicas, *embeddings* o haciendo *fine-tuning* de un Transformer completo. A continuación, se os deja un enlace de cómo fine-tunear un modelo de BERT: <https://colab.research.google.com/drive/1WerT1ApNY3Sqv8QpSK9q1R25wCqRn11?usp=sharing>.
  - Utilizar estrategias no supervisadas de alineamiento con cuestionarios psicológicos. Muchas veces, se utilizan cuestionarios de preguntas como el NEO-FFI<sup>2</sup> para estimar la personalidad de pacientes psicológicos, pero no es el único y podéis buscar más.
- Una posible estrategia no supervisada podría ser utilizar *embeddings* de SBERT y aplicar similaridad coseno para estimar las diferentes dimensiones. **Ojo:** Tened en cuenta que en la mayoría de cuestionarios tienen oraciones para un polo positivo y otro negativo, en nuestro caso, el polo negativo de un rasgo se refleja como una puntuación cercana a 0 en el mismo.
- Estrategias zero- o few-shot con LLMs.
  - etc.

También puede ser útil que apliquéis una fase de **filtrado** de los posts u oraciones más relevantes para cada rasgo de personalidad antes de aplicar alguna de las estrategias previas, ya que no todos los textos serán igual de importantes para predecir cada rasgo.

## 2 Entregables

Existen una serie de fechas intermedias de **envíos de variantes que no requieren Notebook**:

- **Variantes iniciales el 12 de noviembre a las 23:59** y el primer **leaderboard** se publicará el **17 de noviembre**.
- El **26 de noviembre a las 23:59** y el segundo **leaderboard** se publicará el **1 de diciembre**.
- El **7 de diciembre a las 23:59** será el envío de variantes finales. El leaderboard final no se revelará hasta el día 11 de diciembre en la última sesión de presentación de soluciones.
- Sesiones de presentación de soluciones en las clases de teoría del **9** y el **11 de diciembre**. De esta presentación, saldrá el punto de trabajos para la evaluación de la asignatura.

En cada envío, cada equipo puede enviar entre **2-3 variantes** y la única entrega obligatoria de variantes es la del **7 de diciembre**. Se evaluarán con un script automático que calcula el **MSE medio de las 5 dimensiones**. Ese será el criterio para rankear los equipos. Por tanto, debéis entregar para los usuarios de test un fichero con el siguiente formato de columnas:

*team\_name, variant\_name, username, agreeableness, openness, conscientiousness, extraversion, neuroticism*

Nombrad las variantes como **nombre\_equipo\_nombre\_variante.csv**. En el siguiente form podéis registrar oficialmente a vuestro equipo: <https://forms.office.com/Pages/ResponsePage.aspx?id=LEUNj6S3ZE4EIw5c3RHe1gLXxZ3Wk1KrKJaRzDcyqVUM1RLTzJEREZYMDkyWDIxRTRFQkRKUkFXSi4u>.

En la entrega final, debéis entregar un **Python Notebook** con vuestra práctica y que se llame **nombre\_equipo\_reto.ipynb**. Es fundamental que el Notebook sea autoexplicativo de todos los pasos (con celdas textuales acompañando a celdas con código y que contenga explícitamente los resultados -sin tener que ejecutar las celdas de nuevo-). **No** es necesario entregar la presentación.

---

<sup>2</sup>Enlace: [https://docs.google.com/document/d/1pIYif6b2XVwoKe\\_zhVMLW9hHMhLb3n1-19-u8RC3KZE/edit?usp=sharing](https://docs.google.com/document/d/1pIYif6b2XVwoKe_zhVMLW9hHMhLb3n1-19-u8RC3KZE/edit?usp=sharing)

### 3 Valoración y Fecha de Entrega

Esta práctica tiene una valoración de 6 puntos (sobre el total de 10 puntos de la parte práctica de la materia). Fecha límite entrega: **8 de diciembre a las 23:59**. En esta entrega, sí que se debe depositar el **Notebook** que explique todo el procedimiento seguido por los equipos.

Se permiten entregas retrasadas pero se reducirá la puntuación del siguiente modo:

- Cada día tarde reduce en un 10% la máxima nota alcanzable.