

Documento de Principios Éticos y Legales

Proyecto: Chatbot de Recomendación de Videojuegos (DataMind)

1. Introducción y Contexto

El presente documento detalla el marco de cumplimiento normativo, ético y legal bajo el cual se ha desarrollado el asistente virtual de recomendación de videojuegos. El desarrollo se ha regido por los principios de **Privacy by Design** (Privacidad desde el Diseño) y **Security by Design** (Seguridad desde el Diseño), garantizando un uso responsable de la Inteligencia Artificial y la gestión de datos.

2. Privacidad y Protección de Datos

2.1. Origen de los Datos (Ingesta Ética)

Para la obtención de la base de conocimientos, se ha implementado un sistema híbrido que prioriza el uso de canales oficiales y el respeto a la infraestructura del proveedor:

- **Fuentes:**
 - **Steam Web API:** Para la obtención de detalles técnicos, precios y descripciones (canal oficial).
 - **Steam Search Results:** Para la indexación de IDs y popularidad.
- **Mecanismos de Scraping Responsable:**
 - **Limitación de Velocidad (Rate Limiting):** El script de extracción (sacar-datos-games.py) implementa pausas programadas (`time.sleep`) entre peticiones para no saturar los servidores de Valve.
 - **Gestión de Errores 429:** El sistema detecta automáticamente códigos de estado HTTP 429 (Too Many Requests) y detiene la ejecución durante 60 segundos antes de reintentar, cumpliendo estrictamente con las políticas de tráfico del servidor.
- **Privacidad:** Los datos extraídos son exclusivamente información pública del producto. No se extraen perfiles de usuarios, reseñas con nombres reales ni listas de amigos.

2.2. Tratamiento de Datos del Usuario (Logs e Interacciones)

El sistema interactúa con usuarios humanos a través de lenguaje natural. Para garantizar la privacidad:

- **Minimización de Datos:** El sistema **no requiere registro** ni solicita datos personales para funcionar.
- **Gestión de Logs (JSON):** Se almacenan los *prompts* de forma anónima en un archivo JSON local únicamente para depuración técnica.

- **Advertencia:** Se asume el principio de responsabilidad compartida, advirtiendo al usuario de no introducir datos sensibles.

3. Propiedad Intelectual y Licencias

3.1. Inventario de Licencias de Software

El proyecto utiliza software de código abierto y librerías de terceros, respetando sus licencias:

Herramienta / Librería	Licencia	Uso en el proyecto
Python	PSFL	Lenguaje base.
Elasticsearch	SSPL / Elastic License	Motor de búsqueda y almacenamiento.
Sentence-Transformers	Apache 2.0	Generación de embeddings locales (vectorizador.py).
BeautifulSoup4	MIT License	Parseo de HTML en la ingesta.
Pandas / NumPy	BSD-3-Clause	Procesamiento de datos.
Steam Web API	Términos de Uso de Valve	Fuente de datos.

3.2. Modelos de IA Utilizados

El sistema utiliza dos modelos de Inteligencia Artificial distintos:

1. **Modelo de Embeddings (Local):**
 - **Nombre:** paraphrase-multilingual-mpnet-base-v2 (HuggingFace).
 - **Licencia:** Apache 2.0.
 - **Uso:** Vectorización de descripciones para búsqueda semántica.
2. **Modelo de Generación (LLM - Nube):**
 - **Nombre:** Google Gemini 2.0 Flash Lite (vía OpenRouter).
 - **Licencia:** Propietaria (Términos de Servicio de Google Generative AI).
 - **Uso:** Generación de respuestas en lenguaje natural.

4. Protección frente a Errores y Seguridad (Security by Design)

4.1. Filtrado de Contenidos en Origen

Se ha implementado una capa de seguridad en la fase de ingestión de datos (filter-games.py) para evitar que el sistema procese contenido inapropiado:

- **Listas Negras (Blacklisting):** El scraper elimina automáticamente cualquier videojuego

que contenga etiquetas o palabras clave como adult, sexual, xxx o hentai en su título o metadatos.

- **Limpieza de Ruido:** El script clean-tags.py elimina etiquetas irrelevantes ("hardware", "soundtrack") para evitar alucinaciones del modelo basadas en datos técnicos irrelevantes.

4.2. Control en Tiempo de Ejecución

- **Gestión de Alucinaciones:** El sistema utiliza RAG (Retrieval-Augmented Generation). Si Elasticsearch no devuelve videojuegos relevantes para una consulta, el *System Prompt* instruye al modelo para admitir que no tiene información, en lugar de inventar títulos inexistentes.
- **Protección de Credenciales:** Las API Keys no se almacenan en el código fuente.

5. Consideraciones Éticas y Gestión de Sesgos

5.1. Sesgos de Género y Representación

Somos conscientes de los sesgos inherentes en los datos históricos de videojuegos.

- **Mitigación:** Se ha instruido al modelo para utilizar lenguaje inclusivo y neutro en sus recomendaciones.
- **Filtrado Ético:** La eliminación proactiva de contenido explícito o sexualmente violento durante la fase de scraping asegura que el chatbot sea seguro para audiencias generales (Safe for Work).

5.2. Transparencia

El sistema se presenta claramente como una herramienta automatizada ("Chatbot") y no intenta suplantar a un operador humano.

Descargo de Responsabilidad:

Esta herramienta ha sido desarrollada como parte de un reto académico (Reto DataMind). Los desarrolladores no se hacen responsables de las decisiones de compra tomadas en base a estas recomendaciones.