

Grado en Ingeniería Informática

Asignatura: Visión por Computador

TRABAJO FINAL DE ASIGNATURA

Simon Dice

Autores:

Laura Herrera Negrín
Dunia Suárez Rodríguez
Ayman Asbai Ghoudan

31 de diciembre de 2025

Índice

1. Introducción	2
1.1. Descripción del proyecto	2
1.2. Objetivos	2
2. Tecnologías y herramientas	3
2.1. Lenguaje y entorno	3
2.2. Librerías principales	3
3. Estructura del repositorio	4
4. Metodología y desarrollo	5
4.1. Fase 1: Generación del dataset	5
4.2. Fase 2: Procesamiento y entrenamiento	5
4.2.1. Extracción de características	5
4.2.2. Preprocesamiento y normalización	5
4.2.3. Data augmentation (aumento de datos)	6
4.2.4. Modelo de clasificación	6
4.3. Fase 3: Lógica del juego (Simon Says)	6
5. Manual de usuario	7
5.1. Requisitos previos	7
5.2. Ejecución	7
5.3. Cómo jugar	7
6. Conclusiones	8
7. Bibliografía	9

1. Introducción

1.1. Descripción del proyecto

El presente proyecto, titulado «Simon Dice», consiste en el desarrollo de una versión moderna e interactiva del clásico juego de memoria «Simón Dice». La innovación principal radica en la interfaz de usuario: en lugar de utilizar botones físicos o periféricos convencionales (teclado/ratón), el juego se controla enteramente mediante **visión artificial (Computer Vision)**.

El sistema es capaz de reconocer gestos faciales y manuales del jugador en tiempo real a través de una webcam, interpretándolos como comandos de juego. Esto permite una experiencia inmersiva donde el jugador debe replicar secuencias de movimientos corporales.

1.2. Objetivos

- Implementar un sistema de reconocimiento de gestos robusto utilizando técnicas de aprendizaje automático (Machine Learning).
- Desarrollar una lógica de juego interactiva que responda en tiempo real a las acciones del usuario.
- Integrar librerías de visión por computador como OpenCV y MediaPipe para la extracción de características.
- Entrenar clasificadores neuronales (MLP) para distinguir entre diferentes gestos (cara y manos).

2. Tecnologías y herramientas

Para el desarrollo de este proyecto se ha utilizado un stack tecnológico basado en Python, aprovechando su extenso ecosistema para ciencia de datos y visión artificial.

2.1. Lenguaje y entorno

- **Python 3.11**: Lenguaje principal del proyecto.
- **Conda**: Gestor de entornos utilizado para aislar las dependencias (entorno `VC_Trabajo` o `memo_env`).
- **Jupyter Notebooks**: Utilizados para el prototipado, generación de datasets y entrenamiento de modelos.

2.2. Librerías principales

- **OpenCV (cv2)**: Herramienta fundamental para la captura de vídeo desde la webcam, procesamiento de imágenes (cambios de espacio de color BGR a RGB) y visualización en tiempo real.
- **MediaPipe**: Framework desarrollado por Google, utilizado para la detección de hitos (landmarks) faciales y manuales. Específicamente se usan los módulos `FaceMesh` (malla facial) y `Hands` (detección de manos).
- **Scikit-Learn**: Utilizada para la implementación de algoritmos de Machine Learning. En concreto, se ha empleado el `MLPClassifier` (Perceptrón Multicapa) para la clasificación de los gestos basándose en los landmarks extraídos.
- **NumPy**: Para operaciones matemáticas eficientes y manipulación de arrays, crucial para la normalización de coordenadas.
- **Pickle**: Para la serialización y guardado de los modelos entrenados (`.pk1`).

3. Estructura del repositorio

El proyecto se organiza en la siguiente estructura de directorios, separando el código fuente (`src`), los datos (`dataset`) y la documentación:

```
/  
|-- dataset/          # Imágenes para entrenamiento  
|   |-- gestos_cara/  # Dataset de gestos faciales  
|   |-- gestos_manos/ # Dataset de gestos manuales  
|  
|-- src/              # Código fuente del proyecto  
|   |-- game/          # Lógica del juego (Simon Says)  
|   |-- generate_images/ # Scripts para captura de datos  
|   |-- memory/         # Documentación LaTeX  
|   |-- scripts/        # Utilidades auxiliares  
|   |-- train/          # Notebooks de entrenamiento (MLP)  
|  
|-- README.md          # Documentación general y setup
```

4. Metodología y desarrollo

El proyecto se divide en tres fases principales: generación del dataset, entrenamiento de modelos y lógica del juego.

4.1. Fase 1: Generación del dataset

Para la construcción del *dataset*, inicialmente se realizó una búsqueda de imágenes en repositorios públicos como Kaggle. Se identificaron conjuntos de datos relevantes, tales como:

- CEW Dataset (Closed Eyes in the Wild), útil para la detección de ojos cerrados.
- HaGRID (HAnd Gesture Recognition Image Dataset), para el reconocimiento de posturas de la mano.

Sin embargo, para adaptar el sistema a las necesidades específicas del juego (gestos concretos como «guiño» o combinaciones particulares) y mejorar la robustez ante el entorno real de uso, se desarrolló un código específico para enriquecer el conjunto de datos. Mediante los scripts `GuardarGestosCara.ipynb` y `GuardarGestosManoCuerpo.ipynb`, se capturaron imágenes propias utilizando la webcam, almacenándolas en carpetas etiquetadas (ej. `dataset/gestos_cara/0_Neutro`) para complementar las muestras obtenidas de fuentes externas.

4.2. Fase 2: Procesamiento y entrenamiento

El entrenamiento se realiza en los notebooks `EntrenoCara.ipynb` y `EntrenoMano.ipynb`. El flujo de trabajo es el siguiente:

4.2.1. Extracción de características

En lugar de procesar la imagen cruda (píxeles), lo cual sería costoso y propenso a errores por iluminación, se utilizan los **landmarks** de MediaPipe.

- Para la cara, se extrae la malla facial.
- Para las manos, se extraen los 21 puntos clave de la mano.

4.2.2. Preprocesamiento y normalización

Los puntos extraídos (x, y, z) se normalizan para que el modelo sea invariante a la posición del usuario en la pantalla y a la distancia a la cámara.

```

1 def normalizar_puntos(landmarks):
2     coords = np.array([[lm.x, lm.y] for lm in landmarks])
3     centroid = np.mean(coords, axis=0) # Centrar
4     centered = coords - centroid
5     max_dist = np.max(np.abs(centered)) # Escalar
6     return (centered / max_dist).flatten()

```

Listing 1: Función de normalización

4.2.3. Data augmentation (aumento de datos)

Para mejorar la robustez del modelo con pocos datos, se generan variaciones sintéticas de cada muestra original:

- **Ruido Gaussiano:** Simula imperfecciones en la detección.
- **Escalado:** Simula ligeros acercamientos o alejamientos.
- **Rotación:** Simula inclinaciones de la cabeza o mano.

4.2.4. Modelo de clasificación

Se utiliza una red neuronal artificial (MLP).

- **Arquitectura cara:** Capas ocultas de (128, 64) neuronas.
- **Arquitectura manos:** Capas ocultas de (64, 32) neuronas.
- **Validación:** Se utiliza Cross-Validation (K-Fold estratificado) para asegurar que el modelo generaliza correctamente, obteniendo altas tasas de precisión en las pruebas.

4.3. Fase 3: Lógica del juego (Simon Says)

El núcleo del juego se encuentra en `src/game/game.ipynb`. El sistema funciona como una máquina de estados:

1. **Generación de secuencia:** El sistema elige una serie de gestos aleatorios.
2. **Muestra:** Se indica al usuario qué gestos debe realizar (mediante texto o iconos en pantalla).
3. **Escucha (inferencia):** Se activa la cámara y el modelo predice en tiempo real qué gesto está haciendo el usuario.
4. **Validación:** Si el gesto coincide con el esperado en la secuencia, se avanza. Si falla o tarda demasiado, pierde.

Mecánicas adicionales:

- **Trampa «Simón Dice»:** El juego puede intentar engañar al usuario mostrando instrucciones con un nombre diferente (ej. «Modesto dice»). Si el usuario realiza el gesto cuando no lo dijo «Simón», pierde.
- **Dificultad incremental:** Aumenta la longitud de la secuencia y disminuye el tiempo permitido por ronda.

5. Manual de usuario

5.1. Requisitos previos

Es necesario disponer de una webcam funcional y un entorno con las librerías instaladas:

```
pip install opencv-python numpy mediapipe scikit-learn
```

5.2. Ejecución

1. Abrir el proyecto en un editor compatible con Jupyter Notebooks (VS Code o Jupyter Lab).
2. Ejecutar el archivo principal del juego: `src/game/game.ipynb`.
3. Seguir las instrucciones en pantalla.

5.3. Cómo jugar

1. El juego mostrará una secuencia de acciones (ej. «Cerrar Ojos», «Mano Arriba»).
2. Memoriza la secuencia.
3. Cuando sea tu turno, repite los gestos frente a la cámara en el mismo orden.
4. ¡Cuidado! Si la instrucción no empieza por «Simón dice» (o la indicación visual correspondiente), no debes moverte.

6. Conclusiones

Este proyecto demuestra la viabilidad de crear interfaces hombre-máquina naturales y accesibles utilizando hardware de consumo estándar (webcam). La combinación de la extracción de características geométricas (mediante MediaPipe) con clasificadores ligeros (MLP) permite una inferencia en tiempo real extremadamente rápida, adecuada para videojuegos.

El uso de técnicas como la normalización de coordenadas y el aumento de datos ha sido crucial para obtener un sistema robusto que funcione con diferentes usuarios y condiciones de iluminación.

7. Bibliografía

- Documentación de OpenCV: <https://docs.opencv.org/>
- MediaPipe Solutions: <https://google.github.io/mediapipe/>
- Scikit-Learn Documentation: <https://scikit-learn.org/stable/>
- Ultralytics YOLO (referencia conceptual): <https://github.com/ultralytics/ultralytics>