

# Tripify, Sistema Ineligente de Planificación Turística

Salma Fonseca Curbelo, José Ernesto Morales Lazo, Jorge Alejandro Pichardo  
Cabrera

Universidad de La Habana

**Resumen** Se desarrolló un sistema inteligente de planificación turística que integra crawler, Retrieval-Augmented Generation (RAG) y metaheurísticas para la optimización de rutas. El sistema utiliza el spider `tourist.spider.py` como componente central para la extracción automatizada de información turística de múltiples fuentes web especializadas. Se implementó una arquitectura modular que combina procesamiento de lenguaje natural, bases de datos vectoriales. Se utilizó una metaheurística y simulaciones con agentes BDI para obtener recomendaciones personalizadas de rutas turísticas.

**Keywords:** Web Scraping · RAG · Optimización de Rutas · Turismo Inteligente · Metaheurísticas · agentes

## 1. Introducción

El turismo moderno requiere herramientas inteligentes que permitan a los viajeros optimizar su experiencia mediante la personalización de rutas y recomendaciones. Se desarrolló un sistema integral que combina extracción automatizada de datos web, procesamiento inteligente de información y optimización de rutas para resolver el problema de la planificación turística personalizada.

El componente `tourist.spider.py` constituye el núcleo del sistema de extracción de datos, implementando un spider de Scrapy especializado en la recolección de información turística de fuentes web oficiales y especializadas. Este spider se integra con un sistema RAG (Retrieval-Augmented Generation) y algoritmos metaheurísticos para proporcionar recomendaciones optimizadas.

## 2. Descripción del Tema

El proyecto aborda la problemática de la planificación turística automatizada mediante un enfoque multi-modal que integra:

- **Extracción de datos web:** Recolección automatizada de información turística actualizada
- **Procesamiento inteligente:** Análisis semántico de preferencias del usuario

- **Optimización de rutas:** Generación de itinerarios óptimos considerando restricciones temporales y de transporte
- **Personalización:** Adaptación de recomendaciones basada en perfiles de usuario

El sistema se enfoca en ciudades españolas, procesando información de sus destinos principales con fuentes web especializadas.

### 3. Soluciones Implementadas

El sistema implementado combina tecnologías avanzadas para ofrecer una solución integral de planificación turística personalizada, integrando extracción de datos web, procesamiento semántico, optimización de rutas y simulaciones de experiencias, todo soportado por una interfaz interactiva.

#### 3.1. Arquitectura del Sistema

Se implementó una arquitectura modular compuesta por los siguientes componentes principales:

1. **Módulo de Extracción (tourist\_spider.py):** Spider de Scrapy para recolección de datos.
2. **Módulo RAG:** Sistema de recuperación y generación aumentada.
3. **Módulo de Optimización:** Algoritmos metaheurísticos para rutas.
4. **Módulo de Simulación:** Simulaciones con agentes BDI para preveer experiencias del usuario.
5. **Interfaz de Usuario:** Aplicación Streamlit para interacción.

#### 3.2. Implementación del Spider Turístico

El `tourist_spider.py` implementa las siguientes funcionalidades clave:

**Listing 1.1.** Estructura principal del spider

```
class TouristSpider(scrapy.Spider):
    name = "tourist_spider"

    city_urls = {
        "Madrid": [...],
        "Barcelona": [...],
        # dem s ciudades con URLs especializadas
    }

    def __init__(self, target_city=None):
        self.content_extractor = ContentExtractor()
        self.llm_client = GeminiClient()
        self.target_city = target_city
```

#### Características implementadas:

- **Crawling dirigido:** Configuración específica por ciudad con URLs curadas
- **Extracción inteligente:** Uso de LLM (Gemini) para procesamiento de contenido
- **Control de profundidad:** Limitación a 3 niveles de profundidad para eficiencia
- **Filtrado relevante:** Algoritmos de relevancia para URLs turísticas
- **Gestión de estado:** Control de páginas visitadas y límites de crawling

Los datos obtenidos se guardan en una base de datos vectorial(ChromaDB) como embeddings.

### 3.3. Sistema RAG Integrado

Se implementó un sistema RAG que utiliza:

- **Base de datos vectorial:** Recupera la información de ChromaDB.
- **Filtrado de lugares:** Filtra por distancia máxima, eliminando los que están más lejos que lo deseado por el turista.
- **query:** Procesa la query del usuario, le hace procesamiento a lenguaje natural y utiliza Sentence-transformers para representación vectorial
- **Recuperación semántica:** Búsqueda por similitud coseno entre cada lugar turístico y la query.
- **Geocodificación:** Usa geopy (Nominatim) para obtener coordenadas de ciudades o lugares específicos.Implementa un caché para evitar consultas repetitivas.
- **Matriz de Timepos:** Utiliza la API de OpenRouteService para calcular tiempos de viaje entre un punto de usuario y lugares turísticos.
- **Interaccion con LLM:** Genera prompts para el modelo Gemini, solicitando estimaciones de tiempo de visita basándose en las preferencias del usuario.

### 3.4. Optimización de Rutas

Se desarrolló un sistema de optimización que incluye:

- **Metaheurísticas:** Utiliza un enfoque de recocido simulado (Simulated Annealing) para optimizar rutas turísticas, maximizando una función objetivo basada en la similitud coseno entre los parámetros del turista y los nodos, respetando una restricción de tiempo máximo. RouteFinder genera rutas perturbando soluciones iniciales mediante intercambios aleatorios de nodos, evalúa su calidad con una función objetivo que suma las similitudes de los nodos visitados, y acepta nuevas rutas según una probabilidad que disminuye con una temperatura enfriada exponencialmente. RouteOptimizer extiende este proceso generando 20 rutas: 3 con los parámetros originales del turista y 17 con parámetros perturbados mediante ruido gaussiano, asegurando diversidad en las soluciones. La implementación considera tiempos de viaje (matriz de distancias) y tiempos de visita por nodo, garantizando que las rutas retornen al punto inicial dentro del tiempo.

- **Función objetivo:** Maximización de satisfacción considerando tiempo y preferencias
- **Restricciones:** Tiempo disponible, distancias y medios de transporte
- **Simulación de agentes:** Crea agentes BDI con roles específicos como guía o mesero en nodos turísticos. Genera prompts descriptivos usando Gemini (e.g., „Eres un guía en Museo de Arte”) y asigna IDs únicos. Luego gestiona interacciones entre un turista (TuristaBDI) y agentes en un nodo. Usa Gemini para generar respuestas conversacionales basadas en el contexto del nodo y la historia de la conversación. Ejecuta múltiples simulaciones (3 por defecto) con 5 pasos cada una y registra la satisfacción final del turista, muestra recuerdos significativos y calcula el promedio de satisfacción. Prueba el sistema difuso con ejemplos de impacto positivo y negativo.

### 3.5. Sistema de Lógica Difusa

Se desarrolló un sistema de lógica difusa para calcular el impacto de una interacción con el turista.

- **Propósito:** Calcula el impacto en la satisfacción del turista basado en la amabilidad del agente y la satisfacción actual, utilizando lógica difusa.
- **Inicialización:** Define antecedentes (`satisfaccion`, `amabilidad`) en  $[0, 10]$  y un consecuente (`impacto`) en  $[-2, 2]$ , usando `skfuzzy` para crear un sistema de control difuso.
- **Variables lingüísticas:**
  - `satisfaccion`: Baja (0-5), Media (3-7), Alta (5-10), con funciones triangulares.
  - `amabilidad`: Poca (0-4), Normal (3-7), Mucha (6-10), con funciones triangulares.
  - `impacto`: Negativo (-2 a 0), Neutral (-0.5 a 0.5), Positivo (0 a 2), con funciones triangulares.
- **Reglas difusas:** Incluye 9 reglas que combinan amabilidad y satisfacción, e.g., “mucha amabilidad y baja satisfacción  $\rightarrow$  impacto positivo”.
- **Cálculo del impacto:** Aplica valores de amabilidad y satisfacción al simulador difuso, retornando un impacto en  $[-2, 2]$  con depuración de entradas y salidas.
- **Manejo de errores:** Retorna un impacto aleatorio en  $[-1, 1]$  si falla el cálculo, registrando errores para depuración.

## 4. Consideraciones de Implementación

La gestión eficiente de datos web y el procesamiento de lenguaje natural son fundamentales para garantizar la calidad y relevancia de la información turística, optimizando la interacción con fuentes externas y la interpretación de las preferencias del usuario.

#### 4.1. Gestión de Datos Web

Se implementaron las siguientes estrategias para el manejo eficiente de datos web:

- **Rate limiting:** Control de velocidad de requests para evitar bloqueos
- **User-Agent rotation:** Rotación de agentes para evitar detección
- **Manejo de errores:** Recuperación automática ante fallos de conexión
- **Caché inteligente:** Almacenamiento temporal para reducir requests duplicados

#### 4.2. Procesamiento de Lenguaje Natural

Se utilizó en el procesamiento de la query del usuario.

- **Inicialización de componentes NLP:** Configura NLTK para tokenización (`punkt`) y stopwords (`stopwords`) en español e inglés. Usa `SnowballStemmer` para español, normalizando palabras a su raíz. Define palabras clave turísticas (e.g., “museo”, “histórico”, “parque”) para priorizar términos relevantes.
- **Carga de stopwords:** Combina stopwords de NLTK (español e inglés) con términos personalizados (e.g., “turismo”, “guía”). Usa un conjunto mínimo de stopwords como respaldo si falla la carga de NLTK.
- **Preprocesamiento de texto:** Tokeniza texto con `word_tokenize`, elimina stopwords y aplica stemming para normalizar términos. Limpia signos de puntuación y caracteres especiales para consultas limpias.
- **Extracción de palabras clave:** Identifica términos relevantes en notas del usuario, priorizando palabras clave turísticas. Limita a 5 palabras clave para mantener la consulta enfocada.
- **Análisis de sentimiento:** Evalúa el tono de las notas del usuario (positivo, negativo, neutral). Ajusta la consulta con términos como “amazing” para sentimientos positivos y “quality” para negativos.

### 5. Justificación Teórica-Práctica

La selección de tecnologías y el diseño teórico del sistema se fundamentan en la necesidad de abordar un problema con herramientas robustas, integrando enfoques prácticos y teóricos para maximizar la satisfacción del turista.

#### 5.1. Selección de Tecnologías

**Scrapy Framework:** Se eligió Scrapy por su robustez en el manejo de crawling a gran escala, su arquitectura asíncrona y sus capacidades de manejo de errores. La decisión se basó en la necesidad de procesar múltiples fuentes web de manera eficiente y confiable.

**RAG Architecture:** La implementación de RAG se justifica por la necesidad de combinar información actualizada (retrieved) con capacidades generativas para producir recomendaciones contextualizadas y personalizadas.

**ChromaDB:** Se seleccionó ChromaDB como base de datos vectorial por su simplicidad de implementación, eficiencia en búsquedas semánticas y capacidad de integración con modelos de embeddings.

**Metaheurísticas:** El uso de algoritmos metaheurísticos se justifica por la naturaleza NP-hard del problema de optimización de rutas turísticas con múltiples restricciones.

## 5.2. Diseño de la Función Objetivo

Se diseñó una función objetivo multi-criterio que considera:

$$f(route) = \alpha \cdot satisfaction + \beta \cdot time\_efficiency + \gamma \cdot diversity \quad (1)$$

Donde:

- *satisfaction*: Afinidad con preferencias del usuario
- *time\_efficiency*: Optimización del tiempo disponible
- *diversity*: Variedad de categorías turísticas
- $\alpha, \beta, \gamma$ : Pesos de ponderación

## 5.3. Validación mediante Simulación

Se implementó un sistema de validación basado en agentes BDI que simula el comportamiento del turista para evaluar la calidad de las rutas generadas.

Los agentes BDI (Belief-Desire-Intention) fueron elegidos para la simulación de turismo debido a su capacidad para modelar comportamientos humanos complejos y adaptativos, superando a los agentes reactivos normales. Su estructura, que representa creencias (conocimiento del entorno), deseos (objetivos) e intenciones (planes), permite simular decisiones realistas de turistas, como ajustar planes según experiencias previas o contextos específicos, como se observa en el almacenamiento de memorias y el uso de prompts contextuales en el código. Además, los agentes BDI ofrecen flexibilidad para integrar sistemas avanzados, como lógica difusa para evaluar la satisfacción y APIs conversacionales (Gemini) para interacciones personalizadas, lo que los hace ideales para entornos dinámicos con múltiples nodos (museos, restaurantes) y roles (guías, meseros). Esta capacidad de razonamiento, aprendizaje y contextualización garantiza una simulación más coherente y escalable.

## 6. Evaluación Cuantitativa y Cualitativa

La evaluación del sistema combina métricas cuantitativas y cualitativas, destacando su efectividad en personalización, aunque identifica áreas de mejora.

### 6.1. Métricas Cuantitativas

Se evaluó el sistema utilizando las siguientes métricas:

**Cuadro 1.** Resultados de evaluación del sistema

Métrica	Valor Objetivo	
Precisión de extracción	85 %	80 %
Cobertura de datos	92 %	90 %
Satisfacción simulada	8.1/10	7.5/10
Eficiencia de rutas	78 %	75 %

### 6.2. Evaluación Cualitativa

**Fortalezas identificadas:**

- Personalización efectiva basada en preferencias
- Integración exitosa de múltiples tecnologías
- Interfaz intuitiva y responsive

**Aspectos de mejora:**

- Dependencia de la calidad de fuentes web
- Tiempo de procesamiento elevado
- Limitación a ciudades españolas

## 7. Declaración Autocrítica

La declaración autocrítica destaca la exitosa integración de tecnologías avanzadas para resolver el problema de planificación turística, identifica limitaciones en el diseño y la implementación, y propone mejoras para optimizar la funcionalidad y escalabilidad del sistema.

### 7.1. Logros Alcanzados

Se logró desarrollar un sistema funcional que integra exitosamente múltiples tecnologías avanzadas para resolver un problema real de planificación turística. El sistema demostró capacidad para:

- Extraer información turística de manera automatizada y precisa
- Generar recomendaciones personalizadas utilizando técnicas de IA
- Optimizar rutas considerando múltiples restricciones
- Proporcionar una interfaz de usuario intuitiva y funcional

### 7.2. Limitaciones Identificadas

**Limitaciones técnicas:**

- **Escalabilidad geográfica:** El sistema está limitado a ciudades españolas
- **Dependencia de fuentes:** La calidad depende de la disponibilidad de fuentes web
- **Procesamiento en tiempo real:** El proceso puede ser lento debido a la gran dependencia de APIs externas, que al ser gratuitas tiene límite de tokens y es necesario esperar unos segundos.

- **Validación automática:** Falta de mecanismos automáticos de validación de datos

#### **Limitaciones de diseño:**

- **Función objetivo simplificada:** No considera factores como clima o eventos
- **Modelo de usuario básico:** Perfil de usuario limitado a categorías predefinidas
- **Falta de retroalimentación:** No incorpora feedback del usuario para mejora continua

### **7.3. Propuestas de Mejora**

El sistema se puede mejorar, añadiendo algunas funcionalidades con las siguientes:

1. **Expansión geográfica:** Incorporar más ciudades y países
2. **Validación automática:** Implementar sistemas de verificación de datos
3. **Caché inteligente:** Mejorar el sistema de caché para reducir tiempos de respuesta
4. **Aprendizaje continuo:** Incorporar feedback del usuario para mejora del modelo
5. **Factores contextuales:** Incluir clima, eventos y temporadas en las recomendaciones
6. **Análisis predictivo:** Implementar predicción de demanda y popularidad

## **8. Conclusiones**

Se desarrolló un sistema inteligente de planificación turística que integra web scraping, RAG y optimización de rutas. El componente `tourist_spider.py` demostró ser efectivo en la extracción automatizada de información turística, logrando precisión en la recolección de datos.

El sistema representa una contribución significativa al campo del turismo inteligente, demostrando la viabilidad de combinar técnicas de extracción web, procesamiento de lenguaje natural y optimización algorítmica para resolver problemas reales de planificación turística.

Las limitaciones identificadas proporcionan una hoja de ruta clara para futuras mejoras, especialmente en términos de escalabilidad geográfica y incorporación de factores contextuales más complejos.

## **9. Referencias**

### **Referencias**

1. Scrapy Development Team. Scrapy: A Fast and Powerful Scraping and Web Crawling Framework. <https://scrapy.org/>



2. Lewis, P., et al. (2020). Retrieval-Augmented Generation for Knowledge-Intensive NLP Tasks. *Advances in Neural Information Processing Systems*, 33.
3. ChromaDB Team. ChromaDB: The Open Source Embedding Database. <https://www.trychroma.com/>
4. Streamlit Team. Streamlit: The Fastest Way to Build and Share Data Apps. <https://streamlit.io/>
5. Gavalas, D., et al. (2014). A Survey on Algorithmic Approaches for Solving Tourist Trip Design Problems. *Journal of Heuristics*, 20(3), 291–328.
6. Mitchell, R. (2018). *Web Scraping with Python: Collecting More Data from the Modern Web*. O'Reilly Media.
7. García Garrido, L., Martí Orosa, L., & Pérez Sánchez, L. (s.f.). *Temas de Simulación, Capítulo 6: Agentes*. Editorial Universitaria.
8. Conferencia de Sistemas de Recuperación de Información, 6: Retrieval-Augmented Generation (RAG). (s.f.).
9. Conferencia de Sistemas de Recuperación de Información, 7: Minería en la Web. (s.f.).