

Documentación del Crawler para cuba.travel

Equipo de Desarrollo

15 de junio de 2025

Índice

1. Introducción	1
2. Configuración del Crawler	1
2.1. Archivo de configuración (crawler_config.py)	1
2.2. Configuración de Selenium	2
3. Clase CubaTravelCrawler	2
3.1. Estructura principal	2
3.2. Métodos principales	2
4. Flujo de Trabajo	3
5. Estructura de Salida	3
6. Ejemplo de Uso	3
7. Optimizaciones Implementadas	4
8. Limitaciones y Mejoras Futuras	4
8.1. Limitaciones Actuales	4
8.2. Mejoras Sugeridas	4
9. Consideraciones Éticas	5

1. Introducción

Este documento describe el funcionamiento del crawler desarrollado para extraer información de ofertas hoteleras del sitio web cuba.travel. El sistema utiliza Selenium para automatizar la navegación y extracción de datos, con configuración basada en reglas de robots.txt.

2. Configuración del Crawler

2.1. Archivo de configuración (crawler_config.py)

Contiene parámetros esenciales para el funcionamiento del crawler:

```
1 # crawler_config.py
2 CRAWLER_CONFIG = {
3     "sitemap": "https://www.cuba.travel/sitemapindex.xml",
4     "user_agents": {
5         "*": {
6             "disallow": [
```

```

7         "/admin/",
8         "/App_Browsers/",
9         # ... (lista completa de exclusiones)
10    ],
11    "crawl_delay": None,
12 },
13 # ... (configuraciones por user-agent)
14 },
15 "destinations": [
16     "La Habana",
17     "Varadero",
18     # ... (lista completa de destinos)
19 ],
20 }

```

Listing 1: Configuración principal

2.2. Configuración de Selenium

Parámetros para controlar el navegador:

```

1 SELENIUM_CONFIG = {
2     "driver": "chrome",
3     "headless": False,
4     "window_size": "1920,1080",
5     "user_agent": "Mozilla/5.0 (compatible; TourGuideCubaBot/1.0; ...)",
6     "implicit_wait": 0,
7     "page_load_timeout": 120,
8     "disable_images": True,
9     "disable_javascript": False,
10    "disable_cookies": True,
11 }

```

Listing 2: Configuración de Selenium

3. Clase CubaTravelCrawler

3.1. Estructura principal

Clase que gestiona todo el proceso de crawling:

```

1 class CubaTravelCrawler:
2     def __init__(self, base_url="https://www.cuba.travel/"):
3         self.base_url = base_url
4         self.config = CRAWLER_CONFIG
5         self.selenium_config = SELENIUM_CONFIG
6         self.driver = self._init_driver()
7         self.disallow_patterns = self._compile_disallow_patterns()
8         self.crawl_delay = self._get_crawl_delay()
9
10    # ... (funciones internos)

```

Listing 3: Clase principal del crawler

3.2. Métodos principales

- `_init_driver()`: Inicializa el controlador de Chrome con las opciones configuradas
- `_compile_disallow_patterns()`: Convierte reglas de robots.txt a patrones regex
- `is_allowed(url)`: Verifica si una URL está permitida para crawling

- `_select_destination(wait, destination)`: Selecciona un destino en la interfaz
- `extract_offers()`: Extrae datos de las ofertas en la página actual
- `crawl(urls)`: Ejecuta el proceso completo de crawling
- `close()`: Cierra el navegador y limpia recursos

4. Flujo de Trabajo

El proceso de crawling sigue estos pasos:

1. Inicializar el navegador Chrome con configuración personalizada
2. Para cada URL proporcionada:
 - a) Verificar si está permitida por robots.txt
 - b) Cargar la página principal
 - c) Para cada destino en la lista de destinos:
 - 1) Seleccionar el destino en el menú desplegable
 - 2) Hacer clic en el botón "Buscar"
 - 3) Recorrer todas las páginas de resultados
 - 4) Extraer datos de cada oferta hotelera
3. Retornar resultados estructurados por destino

5. Estructura de Salida

Los datos extraídos tienen el siguiente formato por destino:

```

1 {
2   "La Habana": [
3     {
4       "name": "Hotel Nacional de Cuba",
5       "stars": 5,
6       "address": "Calle 21 y O, Vedado",
7       "cadena": "Gran Caribe",
8       "tarifa": "Todo Incluido",
9       "price": "$150/noche",
10      "hotel_url": "https://www.cuba.travel/.../detail"
11    },
12    # ... ofertas
13  ],
14  "Varadero": [
15    # ... ofertas para Varadero
16  ],
17  # ... otros destinos
18 }
```

Listing 4: Estructura de datos de salida

6. Ejemplo de Uso

Implementación básica del crawler:

```

1 from crawler import CubaTravelCrawler
2
3 # Inicializar crawler
4 crawler = CubaTravelCrawler()
5
6 # Ejecutar crawling en URL principal
7 results = crawler.crawl(["https://www.cuba.travel/"])
8
9 # Procesar resultados
10 for destino, ofertas in results.items():
11     print(f"\nDestino: {destino} ({len(ofertas)} ofertas)")
12     for hotel in ofertas:
13         print(f"- {hotel['name']} ({hotel['stars']}*)")
14
15 # Liberar recursos
16 crawler.close()

```

Listing 5: Ejemplo de implementación

7. Optimizaciones Implementadas

- Directorio temporal para datos de usuario
- Deshabilitación de imágenes y cookies
- Patrones regex para URLs prohibidas
- Scroll automático a elementos
- Manejo robusto de paginación
- Inyección de CSS para ocultar elementos multimedia

8. Limitaciones y Mejoras Futuras

8.1. Limitaciones Actuales

- Dependencia de selectores CSS/XPATH específicos
- No manejo de CAPTCHAs o bloqueos avanzados
- Extracción de precios sensible a cambios en la estructura HTML

8.2. Mejoras Sugeridas

- Implementar sistema de proxies rotativos
- Añadir reintentos automáticos para fallos
- Integrar parámetros de búsqueda personalizables (fechas, huéspedes)
- Añadir soporte para almacenamiento en base de datos
- Implementar monitoreo de cambios en la estructura del sitio

9. Consideraciones Éticas

- Respeto estricto a robots.txt y políticas del sitio
- User-agent identificable con información de contacto
- Crawl-delay configurado para minimizar impacto
- Uso responsable de los datos extraídos