

Propuesta de proyecto de Fin de Carrera

Análisis automático de estructura de documentos

Resumen

El presente proyecto trata del estudio, implementación y evaluación sistemática de métodos de análisis de documentos (en inglés, *Document Layout Analysis*). El cometido principal de dichos métodos es, dado un documento, identificar sus partes y la relación lógica entre ellos: títulos, columnas, párrafos, encabezados, pie, gráficos/figuras, imágenes, tablas, etc., como puede verse en el siguiente ejemplo:

The screenshot displays a technical manual page with various elements highlighted and labeled. The labels are color-coded and numbered, corresponding to a legend at the bottom. The legend includes: Text (yellow circle), Caption (orange circle), List-Item (blue circle), Formula (grey circle), Table (white circle), Picture (white circle), Section-Header (red circle), Page-Header (green circle), Page-Footer (green circle), and Title (red circle).

Key elements identified in the document include:

- Page-Header:** 28 Vehicle Care
- Section-Header:** 2 LED Lighting
- Text:** 3 This vehicle has several LED lamps. For replacement of any LED lighting assembly, contact your dealer.
- Section-Header:** 4 Headlamps, Front Turn Signal, Sidemarker, and Parking Lamps
- Section-Header:** 5 Standard Headlamp
- Image:** 26 Diagram of a headlamp assembly with numbered callouts 1, 2, and 3.
- Section-Header:** 6 Sidemarker Lamp
- Section-Header:** 7 Low/High-Beam Headlamp
- Section-Header:** 8 Turn Signal Lamp
- Section-Header:** 9 High/Low-Beam Headlamp
- Text:** 10 replace a headlamp bulb
- Text:** 11 Open the hood. See Hood > 161.
- Image:** 25 Photograph of the engine compartment showing the headlamp assembly area.
- Text:** 12 Remove the washer funnel by turning it counterclockwise and pulling it straight out.
- Image:** 24 Photograph of the headlamp assembly with a callout 27 pointing to the bulb socket.
- Text:** 27 Remove the headlamp bulb socket
- Text:** 14 Turn the bulb counterclockwise and pull straight back.
- Text:** 15 Disconnect the wiring harness connector from the bulb.
- Text:** 16 Install the new bulb in the headlamp assembly by turning clockwise.
- Text:** 17 Reconnect the wiring harness connector.
- Text:** 18 Install the bulb socket into the headlamp assembly by turning clockwise.
- Section-Header:** 19 Sidemarker Lamps
- Text:** 20 replace a sidemarker bulb
- Text:** 21 Open the hood. See Hood > 161.
- Image:** 23 Photograph of the headlamp assembly with a callout 22 pointing to the sidemarker lamp bulb socket.
- Text:** 22 Remove the sidemarker lamp bulb socket from the headlamp assembly by turning counterclockwise.

Objetivos

Se estudiarán los métodos existentes en la literatura, se implementará aquellos relevantes que no tengan implementación disponible, y se evaluarán sobre dos conjuntos de datos: a) conjuntos de evaluación de uso académico, público, como [PubLayNet](#) y b) los documentos del Proyecto Cruzar.uy, que incluyen el Archivo Berruti y otros archivos digitalizados de organismos represivos que operaron durante la dictadura militar en Uruguay (1973-1984).

Antecedentes

A nivel global, el problema de *Document Layout Analysis* (DLA) es tan antiguo como complejo, lejos de estar cerrado a nivel académico e industrial, y paso fundamental en la transcripción e interpretación automática de documentos. Existe muchísima literatura al respecto, e involucra varios sub problemas al punto de que su evaluación suele ser restringida a ciertas familias o tipos de documentos.

Desde el punto de vista local, la motivación surge del proyecto Cruzar.uy, que consiste en el análisis de varios millones de documentos que deben ser procesados de manera automática, rápida y eficaz. La dificultad adicional que dicho conjunto de datos plantea con respecto a bases de datos públicas es el estado de degradación de los documentos, a veces muy avanzado, lo cual dificulta enormemente el análisis. Además de ello, existen ciertos tipos de estructuras y documentos únicos o muy particulares en este conjunto de datos como ser: fichas personales, tablas, formularios complejos.

Debido a lo anterior, los objetivos de este proyecto serán apropiadamente delimitados de manera de poder generar una comparación útil y fiable, y en base a ella seleccionar e implementar un método que pueda ser aplicado en los archivos del proyecto Cruzar.uy.

Resultados esperados

La finalización del proyecto implica lograr dos objetivos bien concretos:

1. Una evaluación sistemática de los métodos de DLA más relevantes de la literatura
2. La implementación correcta, eficiente y con interfaz de línea de comandos de un método seleccionado para su aplicación en los archivos de Cruzar.uy

Interesados dirigirse a Ignacio Ramírez, nacho@fing.edu.uy