

# Combinación de texto e información visual

# Texto plano

Texto plano extraído con pypdf

## Información visual

Texto extraído con OCR  
y reconstruido

## Separar texto en palabras (txt words)

## Separar texto en palabras (ocr words)

## Obtener lista de diferencias con difflib

No diferencia

Capítulo  
Capítulo

# Capítulo

## Palabras cortadas

```

- Constituciona
+ Constituciona-
?                +
- -

```

# Constituciona-

### Palabra diferente

$$\begin{array}{rcl} - & \text{II.} & - 5. \\ ? & ^ & \\ + & \text{IV.} & + 8. \\ ? & ^ & \end{array}$$

II. 5.

## Múltiples palabras diferentes

- Orgánica
- de
- la
- + ORGÁNICA
- + DE
- + LA

Orgánica  
de  
la

## Palabra separada

- DÍA  
+ D  
+ Í  
+ A

D  
Í  
A

## Una palabra sin correspondencia

- 1.  
La

## 1. La