

Alberto Andrés Valdés González.

Degree: Mathematical Engineer.

Work position: Data Scientist.

Mail: anvaldes@uc.cl/alberto.valdes.gonzalez.96@gmail.com

Location: Santiago, Chile.

Preprocessing PDF

In many AI problems like the related to the computer vision in which we have to use models like OCR we have to start by doing a preprocessing of the PDFs.



Steps

1. Transform PDF to images:

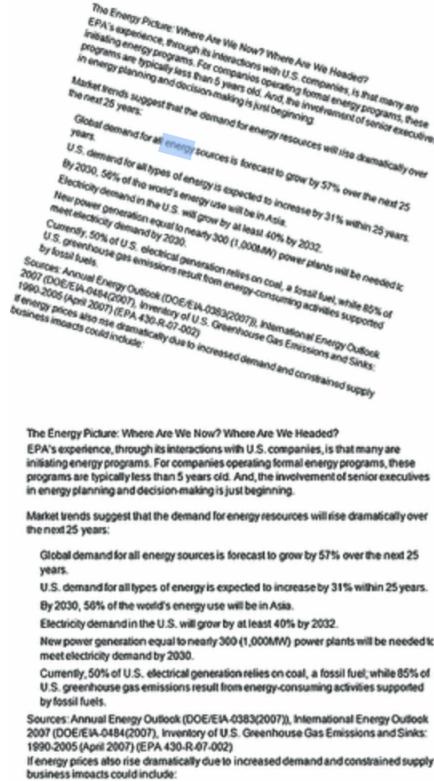
In this step we going to transform every page of the PDF to an image. This step is very important because once we have these images we can work as usual we do in the Computer Vision problems.

2. Transform from RGB to Gray:

To work with only one channel is necessary transform every image to a gray scale.

3. Skew Correction:

It is very important do a skew correction because many times the pdf to analyze corresponds to scan of different documents and how is natural many of these documents are scanned with skew.



4. Noise removal

The main objective of the Noise removal stage is to smoothen the image by removing small dots/- patches which have high intensity than the rest of the image. Noise removal can be performed for both Coloured and Binary images.

5. Binarization of the image

We have to determine a threshold to analyze every pixel of an image and if the value of the pixel is greater than the threshold then we going to assign the maximum value for a pixel and if the value of the pixel is lesser than the threshold then we going to assign the minimum value for a pixel.

6. Other filters

Considering our objectives or our preferences we can apply different filters to every image.
