

Stamp detection & Extraction

Praveen Vijayan
A&I - TCS DATA Office
Email : praveen.vijayan@tcs.com

Modules/Steps in identifying stamp

Thresholding

Connected components / blobs based on min & max size

Blurring & Bounding box over probable ROI

Extract those ROI from actual image

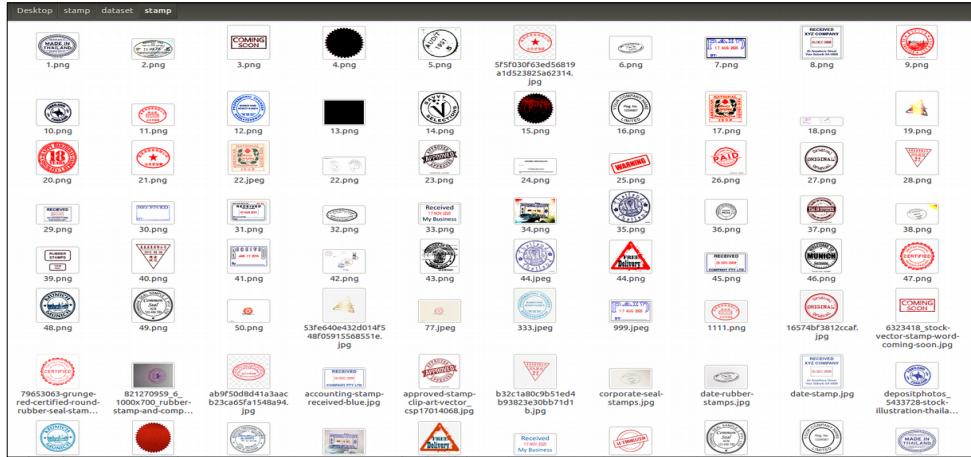
Pass each identified ROI through the **Stamp classifier****

If Classified as **1 (True)**, iterate over the image and extract text from it using tesseract

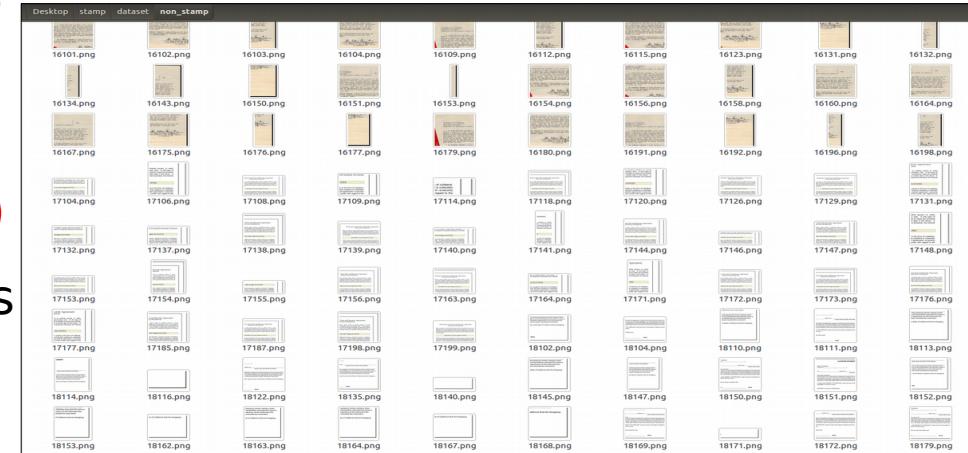
Tag the document/page to the identified text of stamp

Creating : Stamp classifier

Create Stamp dataset



Create **Non -Stamp samples**(conditioned random sampling from the documents and filter those images which doesn't have traces of stamp).



Transfer learning : Used **deep learning VGG-16** fully collected layer 2 to extract the features from images.

May use principal components to reduce dimension.

Create an optimized classifier using SVM.

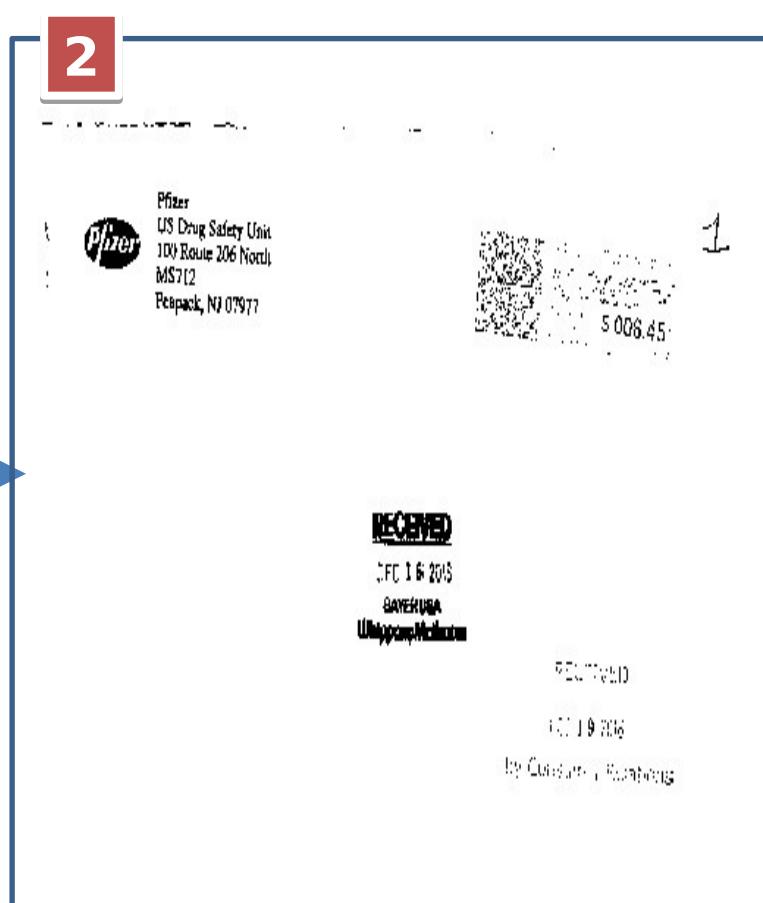
Save the classifier using Joblib to use in the main pipeline.

Given example:

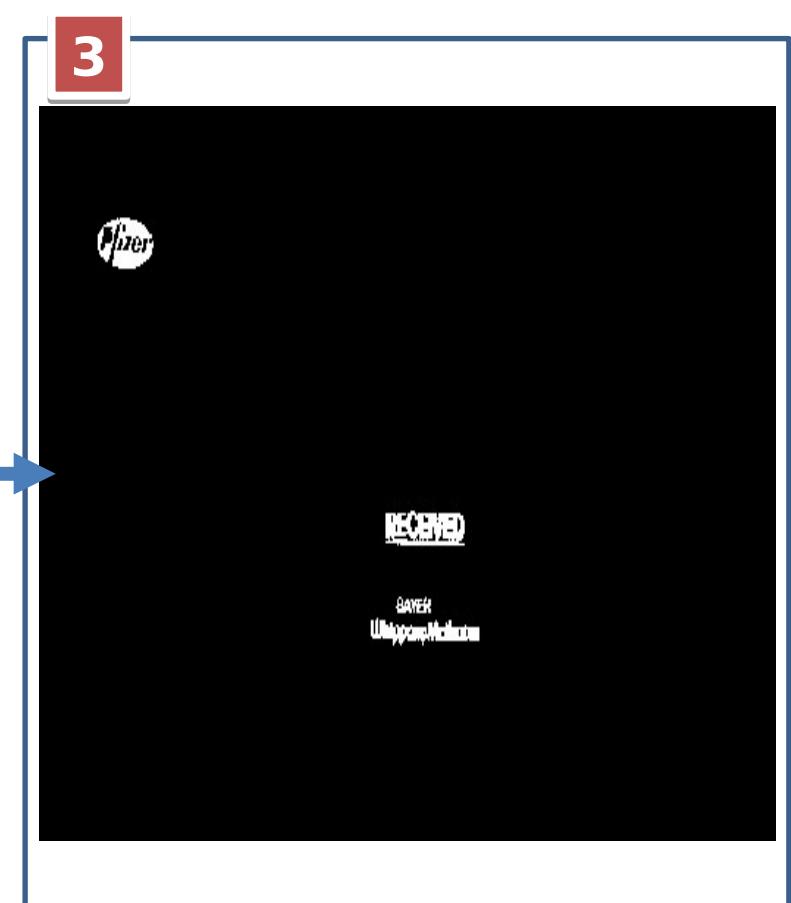
Original image



Thresholded image



Connected component/ blobs in the image



Given example:

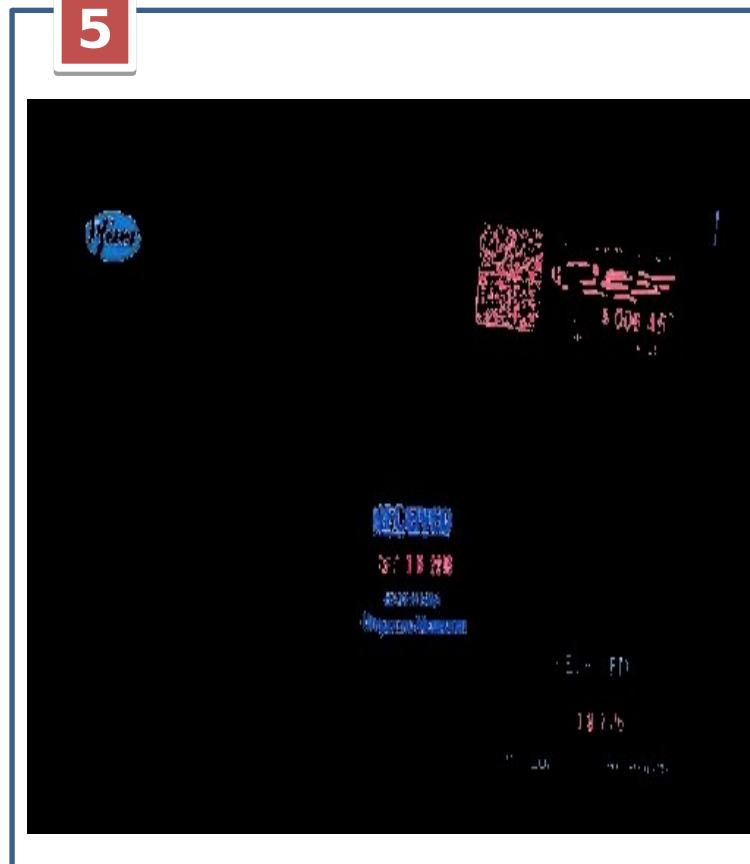
Blurring & Bounding box over big
Connected components

4



Extracted color segments

5



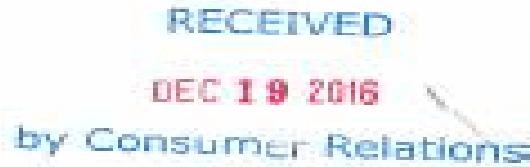
Probable stamps extracted

6



7

Classifier output of stamps



1

1

1

1

```
['RECEIVED', '',
'DEC 19 2016 x',
'by Consumu Relations']
```

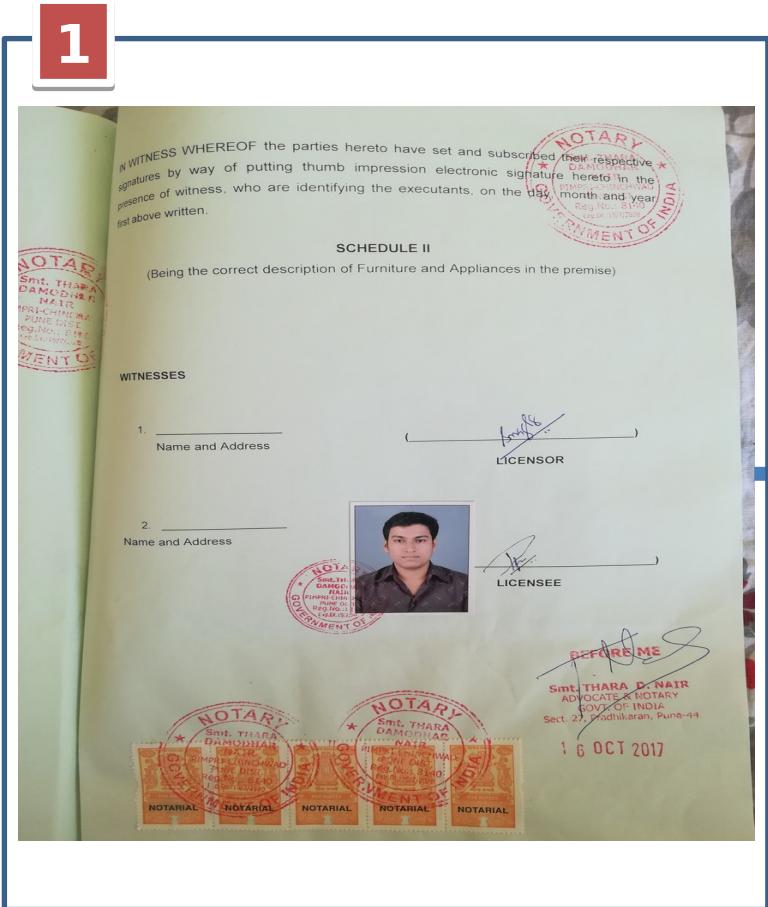
```
[["CHM"]]
```

```
["]
```

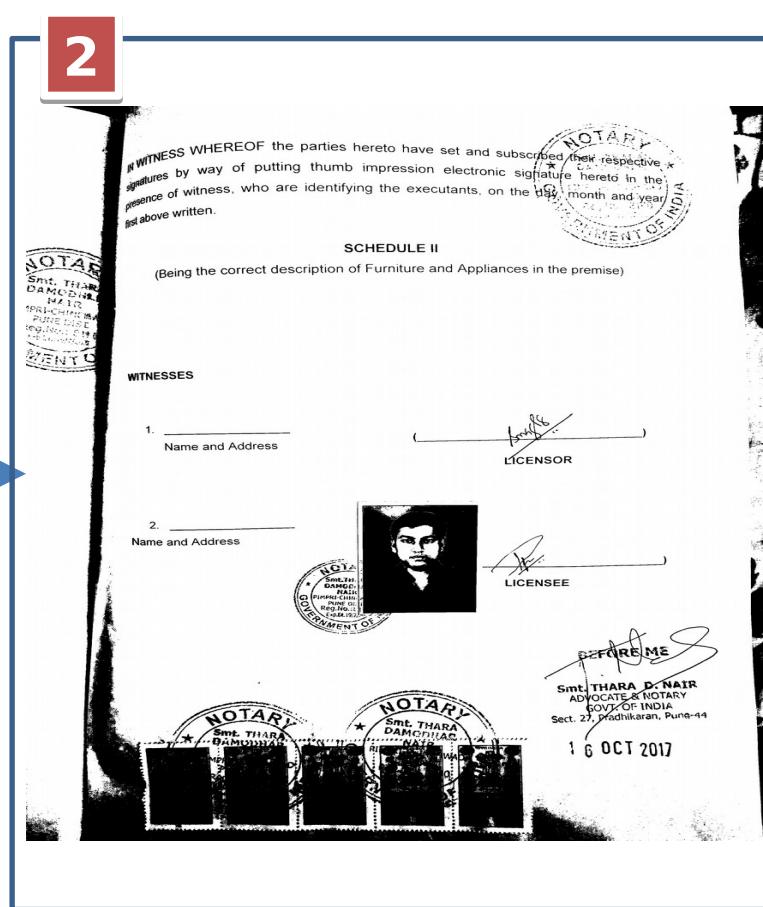
```
['US D', 'le', 'M871']
```

Real example :

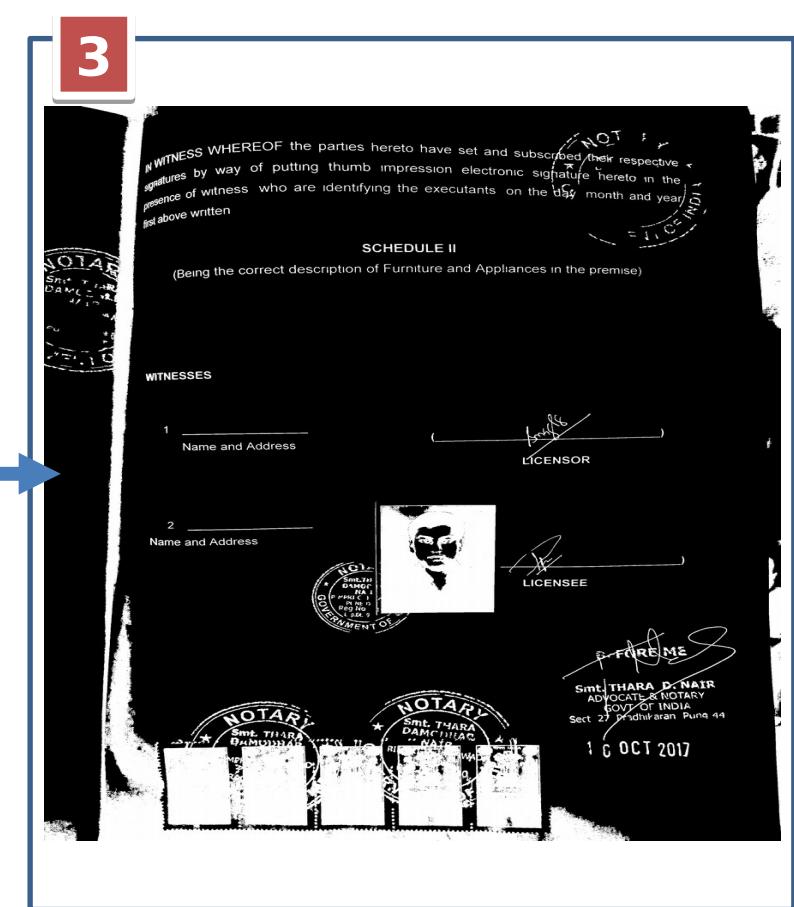
Original image



Thresholded image



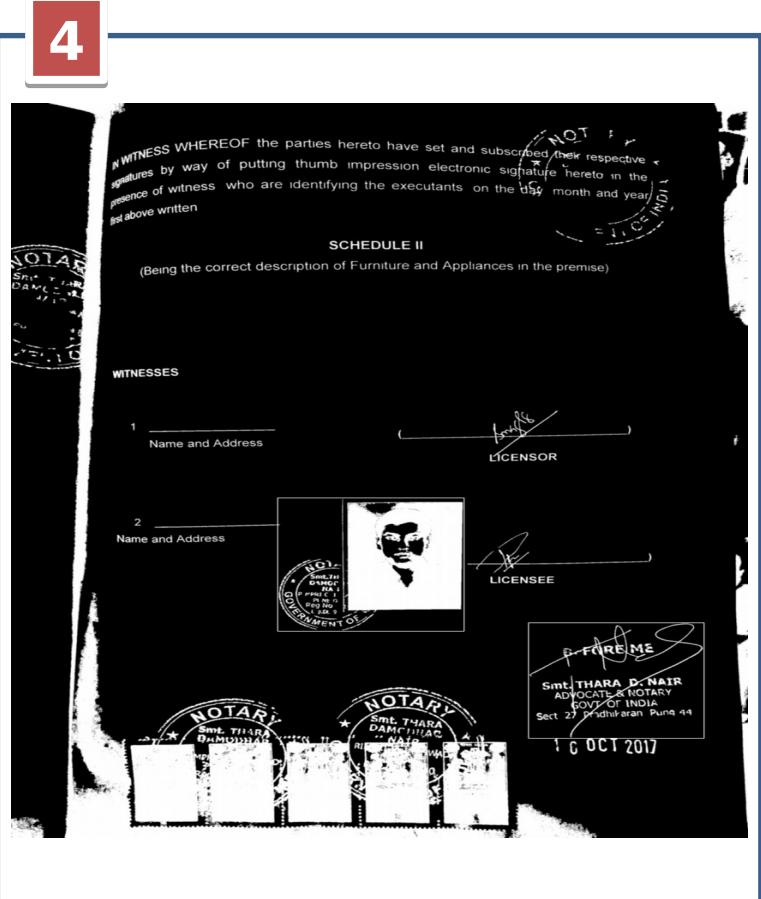
Connected component/ blobs in the image



Real example :

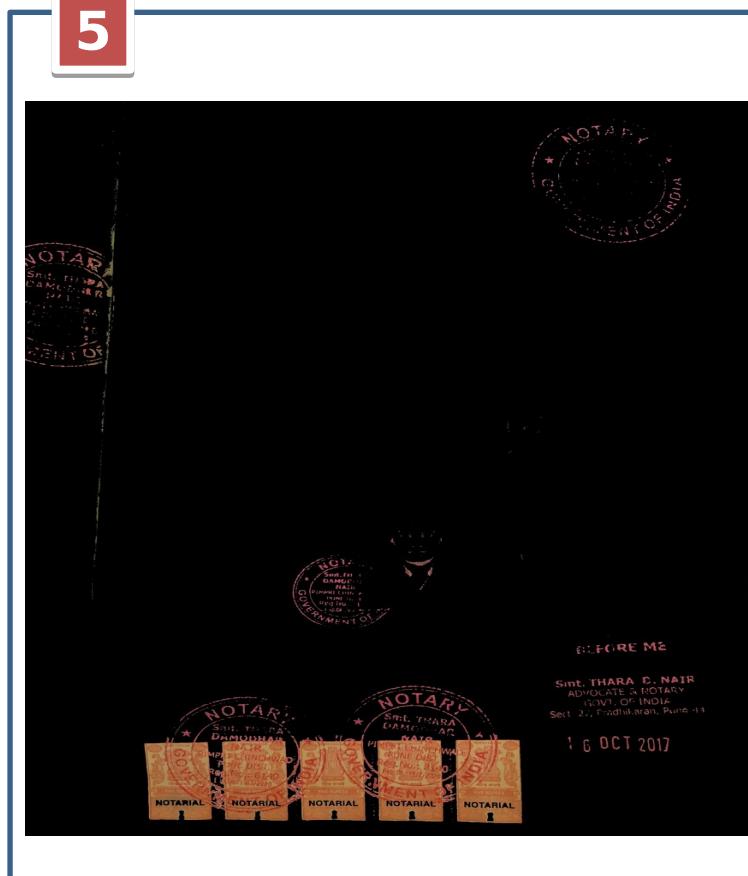
Blurring & Bounding box over big
Connected components

4



Extracted color segments

5

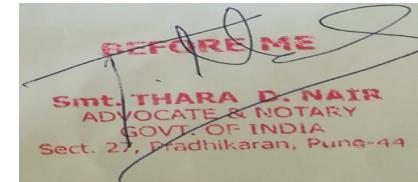
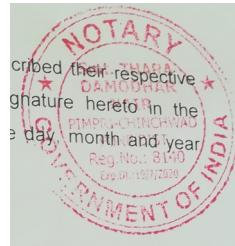


Probable stamps extracted

6



Classifier output of stamps



1

1

0

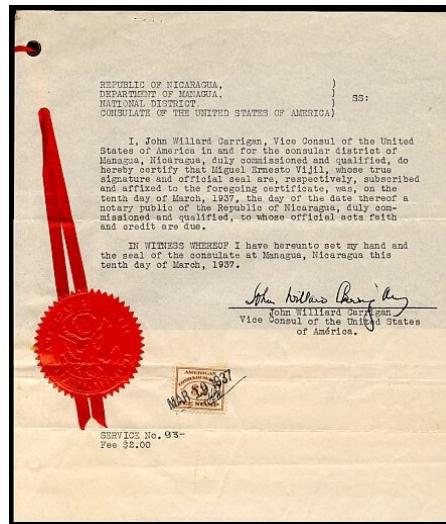
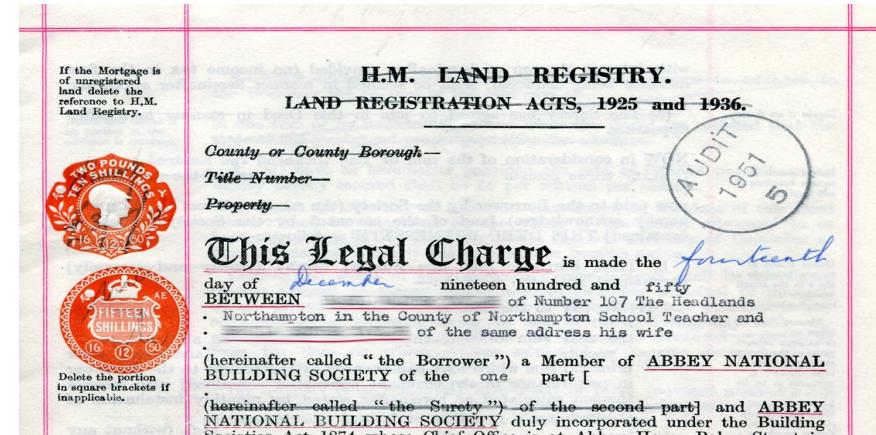
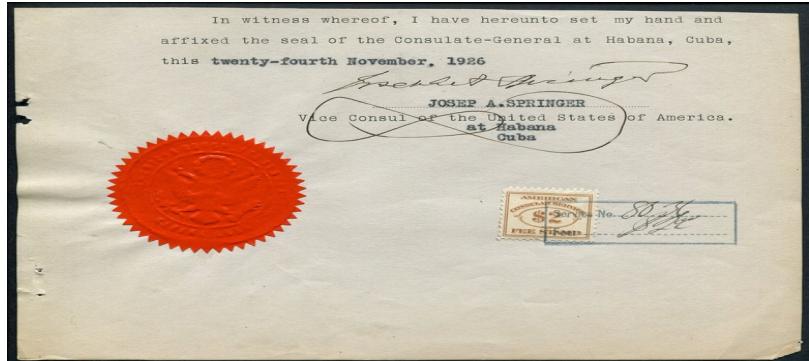
1

1

1

['the", ", 'KNAU S§_']

Tested on random documents



PATIENT PRIVACY AND CONSENT (Read and sign below.)

The information you provide will be used by Pfizer, the Pfizer Patient Assistance Foundation and parties acting on their behalf to determine eligibility, to evaluate and improve the Pfizer RxPathways program, products and services, to communicate with you about your experience with the Pfizer RxPathways program, and/or to send you materials and other helpful information and updates related to the Pfizer RxPathways Program.

By signing below, I affirm that my answers and my proof-of-income documents are complete, true and accurate to the best of my knowledge. I understand that:

- Completing this enrollment form does not guarantee that I will qualify for Pfizer RxPathways.
- Pfizer may verify the accuracy of the information I have provided and may ask for more financial and insurance information at any time.
- Any medications supplied by the Pfizer RxPathways program shall not be sold, traded, bartered or resold.
- Pfizer reserves the right to change or cancel the Pfizer RxPathways program, or terminate my enrollment at any time.
- The support provided in this program is not contingent on any future purchases.

I certify and attest that if I receive medications provided by Pfizer through the Pfizer RxPathways program:

- I will pay copays and coinsurance for my medications based on my financial status or insurance coverage changes.
- I will not seek any rebates or credits for medications I counted in my Medicare Part D out-of-pocket expenses, including my Medicare Part D plans for any costs of medications.
- I will notify my insurance provider of the receipt of any medications through Pfizer RxPathways.
- I have a signed copy of a current and completed HIPAA Authorization Form on record with my Prescriber so that my Prescriber may share health information about me with the Pfizer RxPathways program, Pfizer Inc., and the Pfizer Patient Assistance Foundation Inc.

Signature of Patient
 (Print or e-signature, if under 18 years of age) *[Signature]*

Witnessed by: *[Signature]*

Date: *01/05/2016*

Signature of Patient
(Parent or guardian, if under 18 years of age) X Virginia L Smith Date: 11/28/2016

DEC 19 2016

2016570667 DEC 19 2016

Central Standard Time • SVR:BLMMMS50144 • DNS:10107 • CSID:1 • DURATION mm:ss:02:59

Pros of this method

Instead of going for a cascading scanning of probable stamps with a Classifier, preferred **image processing method**, which is way faster to find ROIs.

Created **dataset for non-stamp examples** by conditional random sampling from the documents and removing those containing traces of stamps.

Used **Deep learning(VGG-16)** to capture better features that can distinguish stamps and non stamps.

At present system is designed in such a way that **False positive is better than False negative**. Because missing a valid stamp may end up not getting valuable info. Even if its false negative, the last stage of stamp extraction would filter the wrong ones.

Further improvements possible

Adaptive thresholding.

Better filter for identifying blobs in a much **generalized** way.

More examples of Stamps & non-stamps images.

Train for OCR extraction, as tesseract not giving good results on multiple lines or curved areas.

Finding angle of the text is hard, have to automate orientation of texts.

Thank you

