

PDF Classifier

By: Matthew Wong

Motivation

- ❑ Businesses get hundreds of PDFs
- ❑ PDFs are sorted by employees
- ❑ Takes a lot of time and costs money
- ❑ Goal: Classify PDFs



Data

- ❏ 14 different classifications
 - ❏ Appraisal
 - ❏ Income
 - ❏ Insurance
- ❏ ~700 documents and uneven classes
- ❏ PDFs in each classification vary
 - ❏ Forms
 - ❏ Scans
 - ❏ Faxes

Text Extraction

- ❑ PDFminer
 - ❑ Works well for PDFs with text layers
 - ❑ Much quicker than PyOCR
- ❑ PyOCR (python wrapper for Tesseract)
 - ❑ Works well for most PDFs
 - ❑ Takes a long time

Different Models

Naive Bayes

Simpler NLP model

Worked well for some categories and not so great with others

TF-IDF

Average Cosine Similarity

Tried to capture signal from variation in documents of the same category

TF-IDF

kNN Cosine Similarity

Only looked at the n most similar documents

Cross Validated Accuracy

- ❑ Used Stratified Cross Validation due to uneven classes

Naive Bayes

PDFminer:

~ 0.58

PDFminer and PyOCR:

~0.62

PyOCR:

~0.48

Average Cosine Similarity

PDFminer:

~ 0.67

PDFminer and PyOCR:

~0.75

PyOCR:

~0.67

kNN Cosine Similarity

PDFminer:

~ 0.71

PDFminer and PyOCR:

~0.82

PyOCR:

~0.81

Future Work

- ❑ Renaming/Filing documents automatically
- ❑ Extracting specific text strings
 - ❑ Name
 - ❑ Address
 - ❑ \$\$\$ Values

Thank You Questions?

Matthew Wong

mjwong1991@gmail.com

<https://www.linkedin.com/in/matthew-j-wong/>

<https://github.com/Unrelenting>