

# Computer Vision for Scraping non-searchable PDFs

Zack Boyce  
October 6, 2018

# Introduction

---

## Employment:

Markel Corporation

- Data Scientist / Jan. '17 - Present
- Business Intelligence Analyst / Oct. '13 - Dec. '17

## Education:

- Masters in Decision Analytics / '16 - '18
- Bachelors in Economics and Business / '08 - '12



[boyce.zack@gmail.com](mailto:boyce.zack@gmail.com)



[www.linkedin.com/in/zack-boyce](https://www.linkedin.com/in/zack-boyce)

# Unstructured information might account for more than 70%–80% of all data in organizations.

---

## How is it defined?

Information that either does not have a pre-defined data model or is not organized in a pre-defined manner.

## What does it look like?

- Pictures
- Emails
- PDFs
- Videos
- etc.

## Why it's important at Markel?

30k submissions = 650k pages

650k pages = minimal data captured

Minimal data captured = tough to build models for underwriters.



Computer vision is an interdisciplinary field that deals with how computers can be made to gain high-level understanding from digital images or videos.

---

Or in other words...

“This is one of the most complex processes we’ve ever attempted to comprehend – let alone recreate. Inventing a machine that sees like we do is a deceptively difficult task, not just because it’s hard to make computers do it, but because we’re not entirely sure how we do it in the first place.

What actually happens is roughly this: the image of the ball passes through your eye and strikes your retina, which does some elementary analysis and sends it along to the brain, where the visual cortex more thoroughly analyzes the image. It then sends it out to the rest of the cortex, which compares it to everything it already knows, classifies the objects and dimensions, and finally decides on something to do: raise your hand and catch the ball (having predicted its path). This takes place in a tiny fraction of a second, with almost no conscious effort, and almost never fails. So recreating human vision isn’t just a hard problem, it’s a set of them, each of which relies on the other.”

-Devin Coldewey@techcrunch / 2 years ago

But for now, let's let our friends at Stanford and MIT develop cutting edge algorithms... we're in the application game.

---



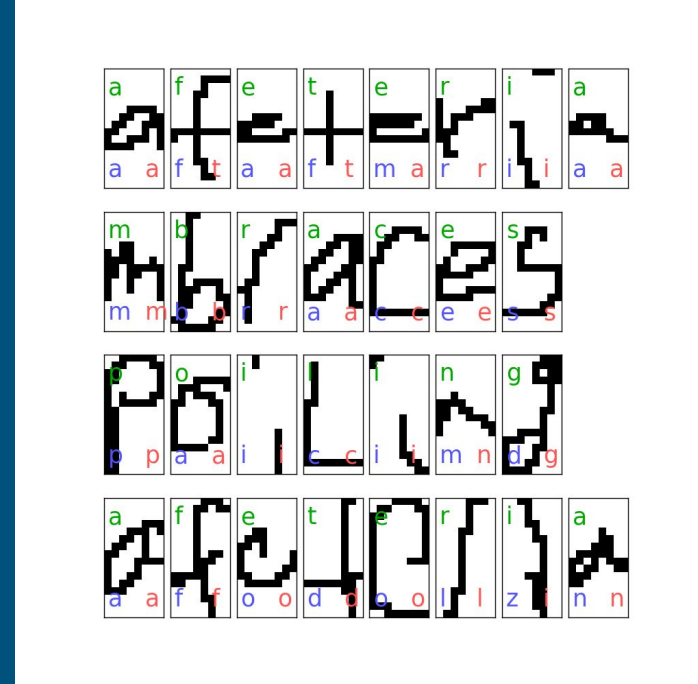
Computers can help us get value from unstructured data through process like Optical Character Recognition (OCR)

OCR is a form of image **classification**.

Common classification algorithms are used (... or used to be used) to help classify images.

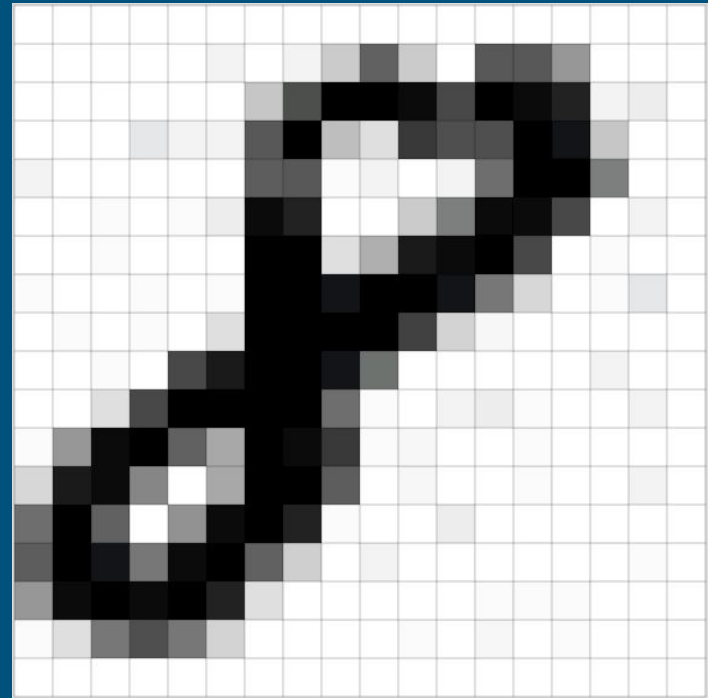
- Naive Bayes
- SVM
- KNN

But now.... Convolutional Neural Networks (CNN)



# But... what exactly are you classifying? I mean thought you needed data for classification models?

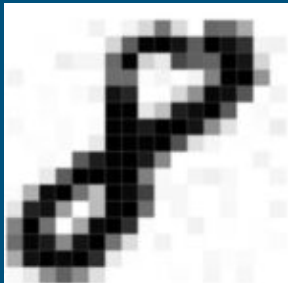
- A computer sees an image as an array of values.
  - Grayscale images take the form of a 2d array with values ranging from 0 - 255.
  - Color images take the form of a 3d array with r, g, b values ranging from 0 - 255.





And voilà, flatten those values and now you have a single observation and n-features.

---

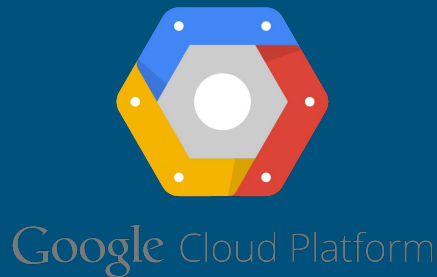


=

```
[0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 1, 12, 0, 11, 39, 137, 37, 0, 152, 147, 84, 0, 0, 0, 0, 0, 1, 0, 0, 0, 41, 160, 250, 255, 235, 162, 255, 238, 206, 11, 13, 0, 0, 0, 0, 16, 9, 9, 150, 251, 45, 21, 184, 159, 154, 255, 233, 40, 0, 0, 10, 0, 0, 0, 0, 0, 145, 146, 3, 10, 0, 11, 124, 253, 255, 107, 0, 0, 0, 0, 3, 0, 4, 15, 236, 216, 0, 0, 38, 109, 247, 240, 169, 0, 11, 0, 1, 0, 2, 0, 0, 0, 253, 253, 23, 62, 224, 241, 255, 164, 0, 5, 0, 0, 6, 0, 0, 4, 0, 3, 252, 250, 228, 255, 255, 234, 112, 28, 0, 2, 17, 0, 0, 2, 1, 4, 0, 21, 255, 253, 251, 255, 172, 31, 8, 0, 1, 0, 0, 0, 0, 0, 4, 0, 163, 225, 251, 255, 229, 120, 0, 0, 0, 0, 0, 11, 0, 0, 0, 0, 21, 162, 255, 255, 254, 255, 126, 6, 0, 10, 14, 6, 0, 0, 9, 0, 3, 79, 242, 255, 141, 66, 255, 245, 189, 7, 8, 0, 0, 5, 0, 0, 0, 0, 26, 221, 237, 98, 0, 67, 251, 255, 144, 0, 8, 0, 0, 7, 0, 0, 11, 0, 125, 255, 141, 0, 87, 244, 255, 208, 3, 0, 0, 13, 0, 1, 0, 1, 0, 0, 145, 248, 228, 116, 235, 255, 141, 34, 0, 11, 0, 1, 0, 0, 0, 1, 3, 0, 85, 237, 253, 246, 255, 210, 21, 1, 0, 1, 0, 0, 6, 2, 4, 0, 0, 0, 6, 23, 112, 157, 114, 32, 0, 0, 0, 0, 2, 0, 8, 0, 7, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0]
```

While understanding the technology is important, you likely will never need to build your own OCR engine.

---



**ABBYY®**



# Resources

---

## Modules/Utilities:

- XpdfReader
  - Pdftopng
- OpenCV
- Numpy
- Keras (for deep learning)

## Resources:

- <http://cs231n.stanford.edu/>
- <https://becominghuman.ai/>
- <https://hackernoon.com/>

So let's get to it

