# Business
# News Classification Engine

**Springboard Capstone Project 2**

**James Flint | mail@jamesflint.Net | 2018-03-30**

# REQUIREMENT

- **Take CurationCorp's labelled news database of 43,502 articles**

- **Train a neural net-based topic classification engine**

- **Make this functionality available via a cloud-based API**

# APPROACH

- **Data wrangling**

- **Compare Classifiers**

  - A multi-layer neural net (NN)

  - A convolutional neural net (CNN)

  - A long/short term memory neural net (LSTM)

  - A very deep convolutional neural net (VDCNN)

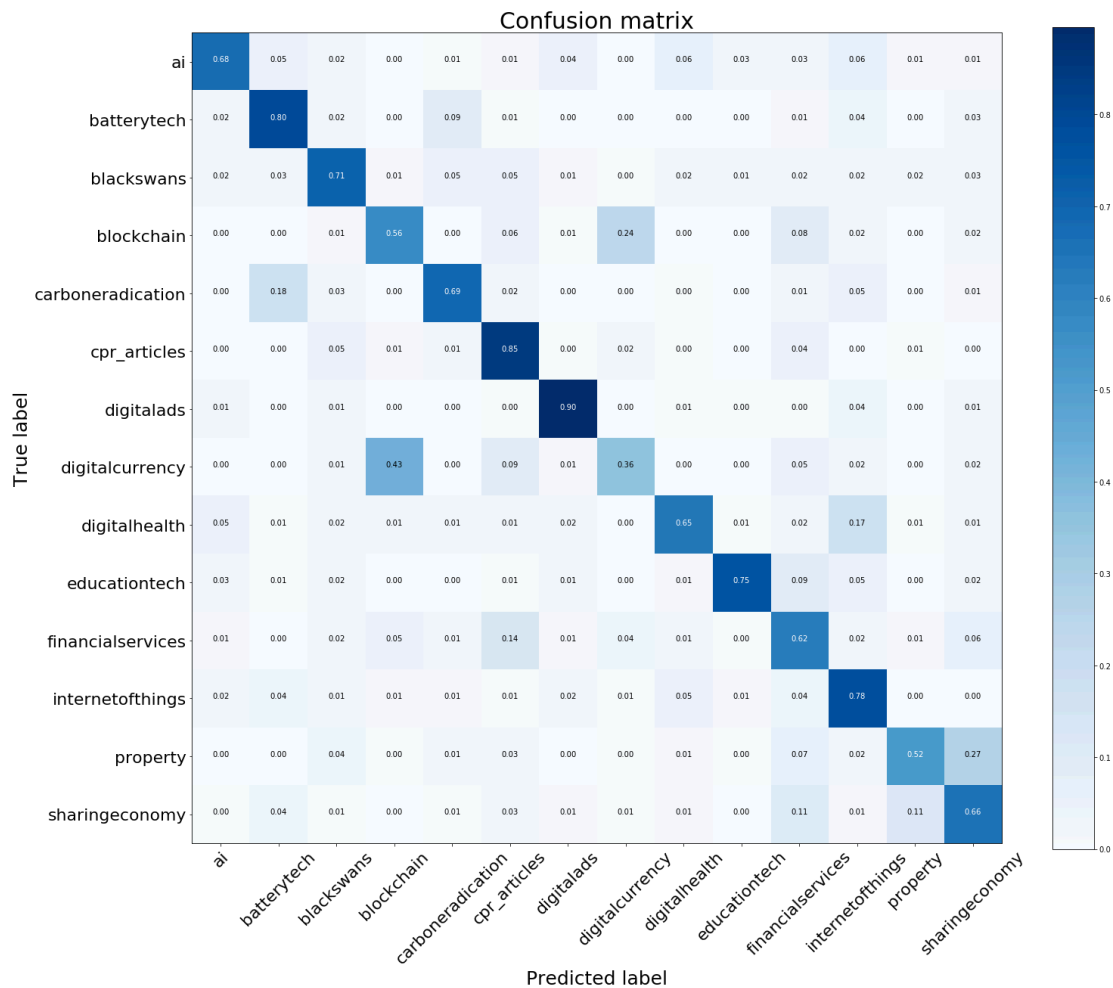- **Build a prediction API**

- **Future research**

# UNBALANCED DATASET

| | |
|---|---|
| cpr_articles | 6846 |
| blackswans | 4837 |
| batterytech | 4092 |
| financialservices | 3986 |
| carboneradication | 3690 |
| sharingeconomy | 3574 |
| digitalads | 2920 |
| internetofthings | 2627 |
| property | 2132 |
| digitalhealth | 1943 |
| digitalcurrency | 1914 |
| ai | 1722 |
| blockchain | 1650 |
| educationtech | 1555 |

# BALANCED DATASET

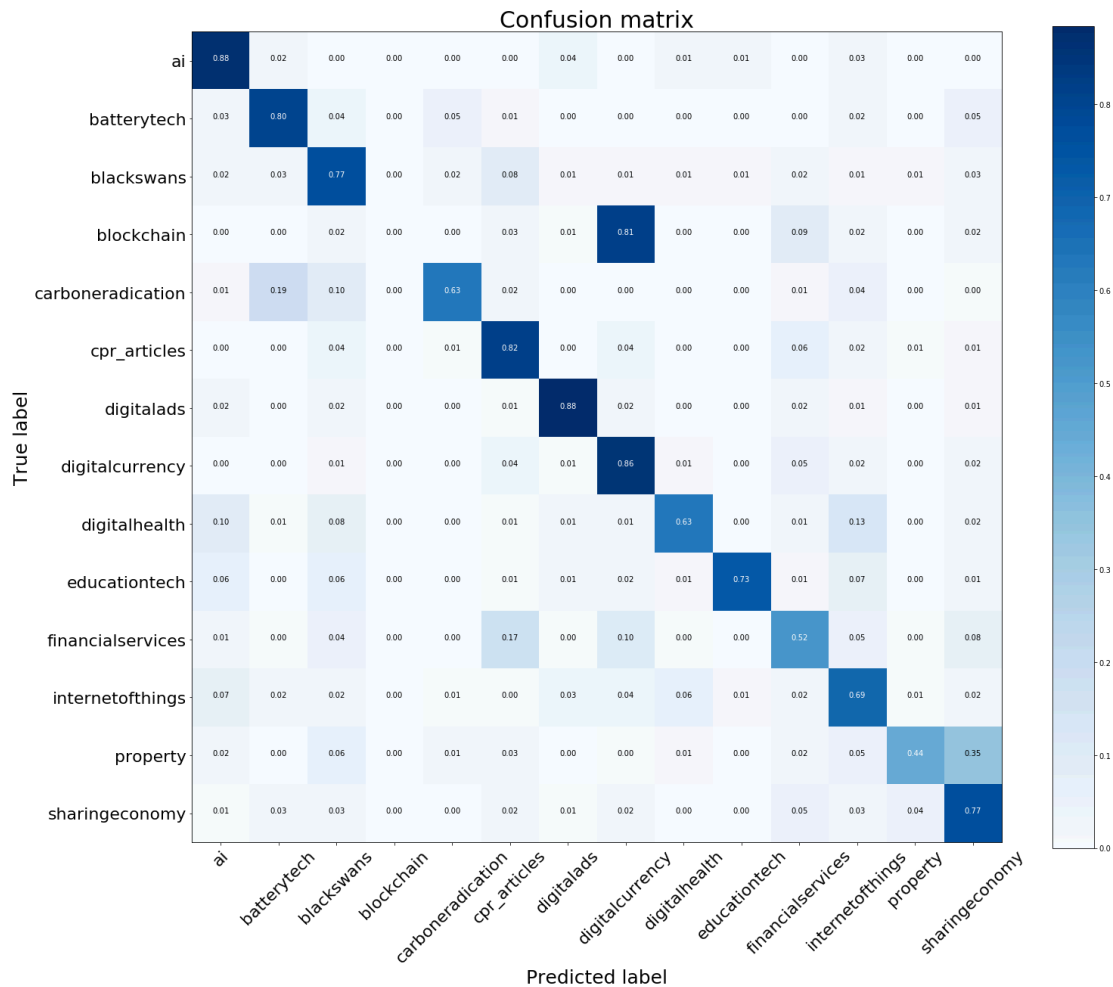| | |
|---|---|
| batterytech | 16368 |
| financialservices | 15944 |
| internetofthings | 15762 |
| digitalhealth | 15544 |
| digitalcurrency | 15312 |
| property | 14924 |
| carboneradication | 14760 |
| digitalads | 14600 |
| blackswans | 14511 |
| sharingeconomy | 14296 |
| educationtech | 13995 |
| ai | 13776 |
| cpr_articles | 13692 |
| blockchain | 13200 |

# MULTILAYER NEURAL NET



Confusion matrix

**Best performance:**

**72.5% accuracy (unbalanced data)**

**Support vector machine benchmark:**

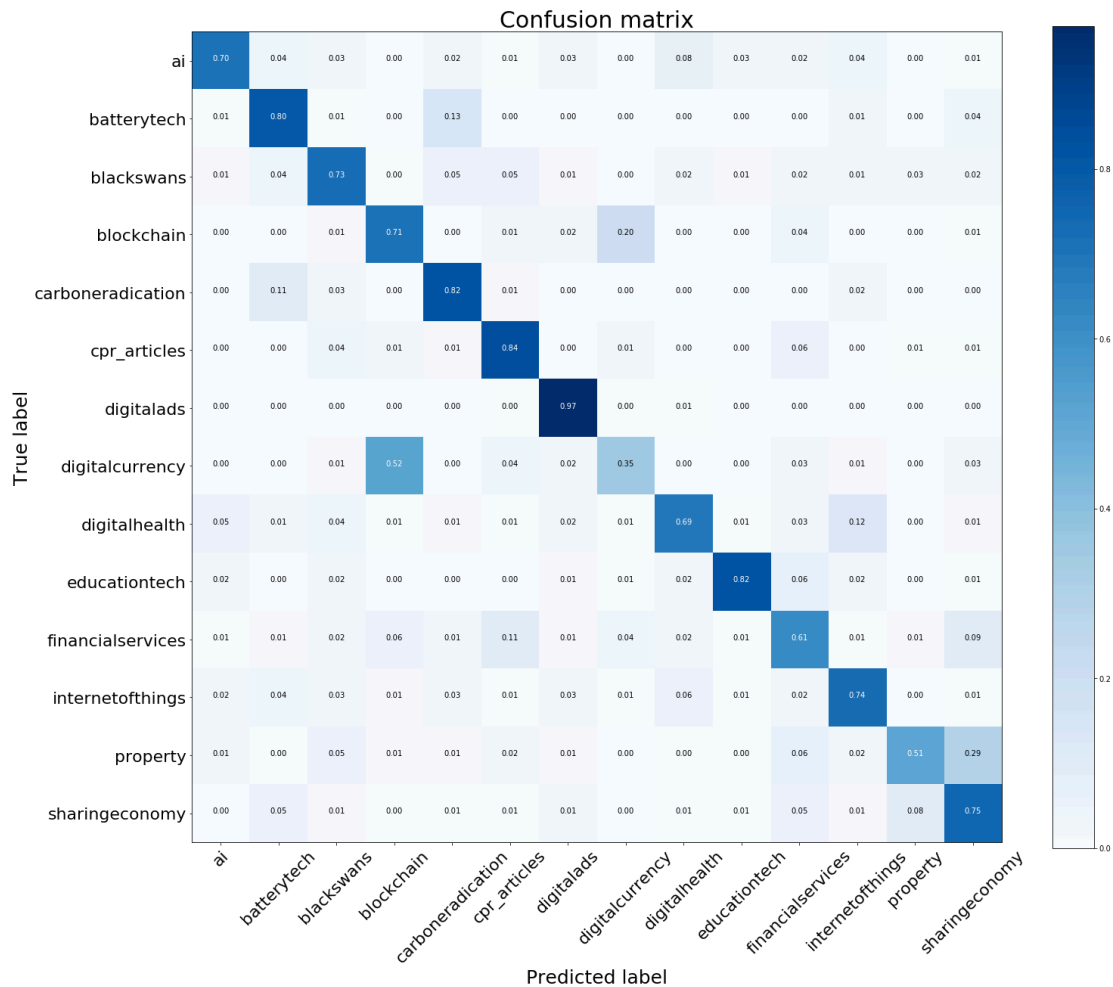**77% accuracy**

# CONVOLUTIONAL NEURAL NET



Confusion matrix

**Best performance:**

**72.6% accuracy
(unbalanced data)**

**Support vector
machine benchmark:**
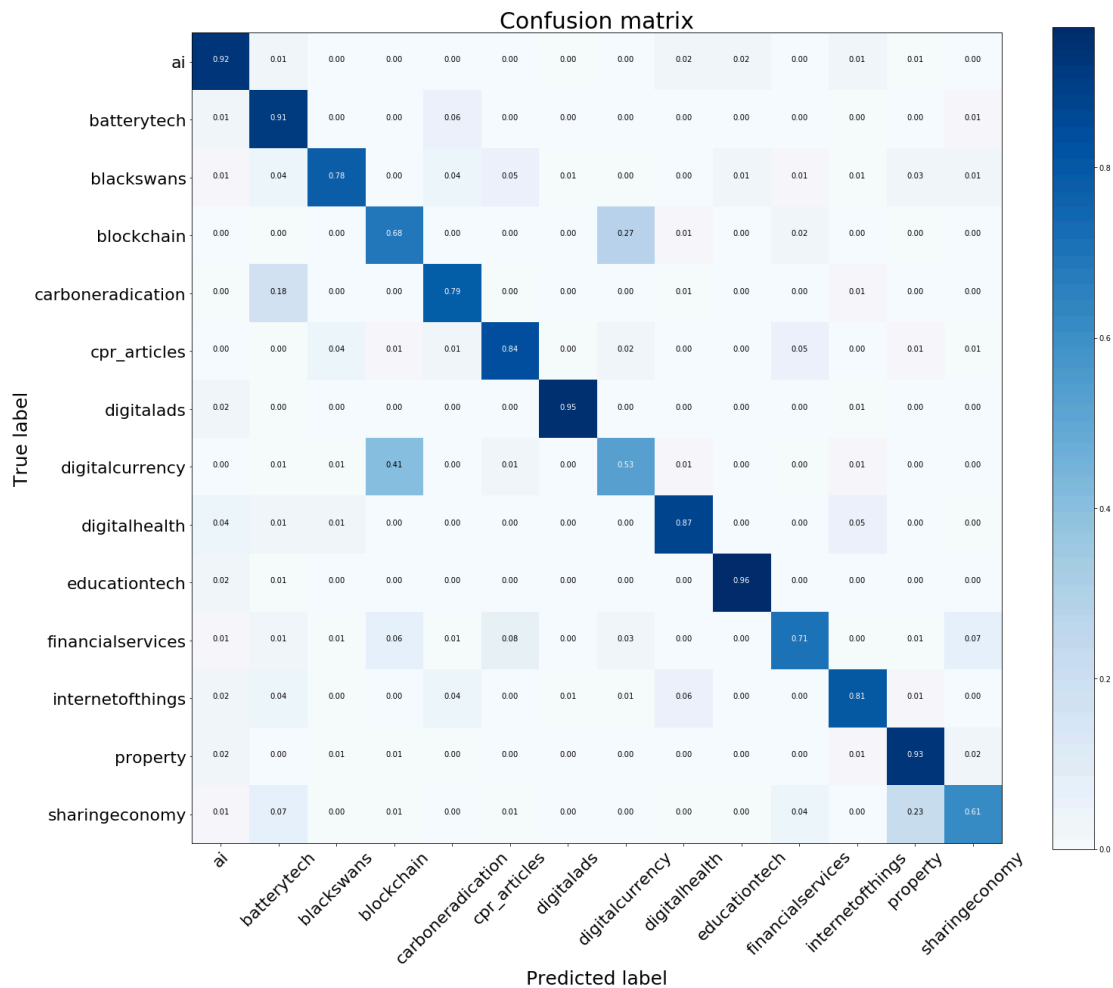
**77% accuracy**

# LONG SHORT-TERM MEMORY


Confusion matrix

**Best performance:**

**81.3% accuracy
(balanced data)**

**Support vector
machine benchmark:**

**77% accuracy**

# VERY DEEP CNN



Confusion matrix

**Best performance:**

**80.6% accuracy
(balanced data)**

**Support vector
machine benchmark:**

**77% accuracy**

# PREDICTION WEB FORM

**Please classify an article using the form below**

Title (max 15 words):

[_____]

Article (max 135 words) :

[_____
_____
_____
_____
_____]

[CLASSIFY!]

**https://afternoon-shelf-15457.herokuapp.com/form**

# PREDICTION API

```python
import requests
import json

url = 'https://afternoon-shelf-15457.herokuapp.com/predict'
s = {"title":"foo", "body":"bar"}
s_json = json.dumps(s)
headers = {'Content-Type': 'application/json'}
r = requests.post(url, data=s_json, headers=headers)
print(r.text)
```

**https://afternoon-shelf-15457.herokuapp.com/predict**

# NEXT STEPS

- **Text Summarisation**

- **Tag generation**

- **Retrain using user feedback**

- **Enable batch processing**

*See REPORT document for details on each of these*

# THANKS & GOOD LUCK