

USING NATURAL LANGUAGE PROCESSING TO PREDICT PATENT CLASSIFICATION

CATHERINE FRITZ – CAPSTONE



BUSINESS UNDERSTANDING

- Various types of information on patents are connected at various stages in a patent's life, one being the classification of a patent's technology (e.g. mechanical, chemistry, electrical, etc.)
- Classification is assigned by the patent Office some time after a patent is filed
- Could be useful to automate the classification process.

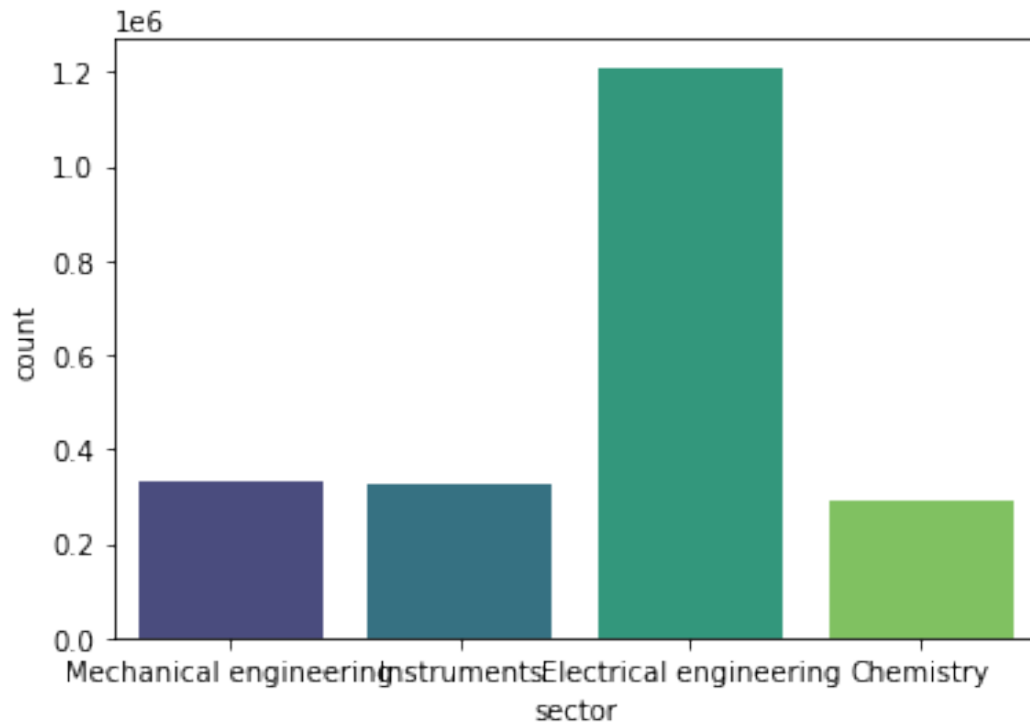
CLAIM EXAMPLE

- One sentence per claim
- Focused on invention
- Structured around features of invention

10
What is claimed is:
1. A laser detection and ranging (LADAR) system, comprising:
a two-dimensional array of detector elements, each detector element within the array including: 15
a photosensitive region configured to receive return light reflected from a target and oscillating local light from a local light source, and
local processing circuitry coupled to an output of the respective photosensitive region and configured to 20
receive an analog signal on the output and to sample the analog signal a plurality of times during each sample period clock cycle to obtain a plurality of components for a sample during each sample period clock cycle; 25
a data bus coupled to one or more outputs of each of the detector elements and configured to receive the plurality of sample components from each of the detector elements for each sample period clock cycle; and
a processor coupled to the data bus and configured to 30
receive, from the data bus, the plurality of sample components from each of the detector elements for each sample period clock cycle and to determine an amplitude and a phase for an interfering frequency corresponding to interference between the return light 35
and the oscillating local light using the plurality of sample components.
2. The system according to claim 1, wherein the two-dimensional array of detector elements comprises a large format array. 40

DATA PREPARATION

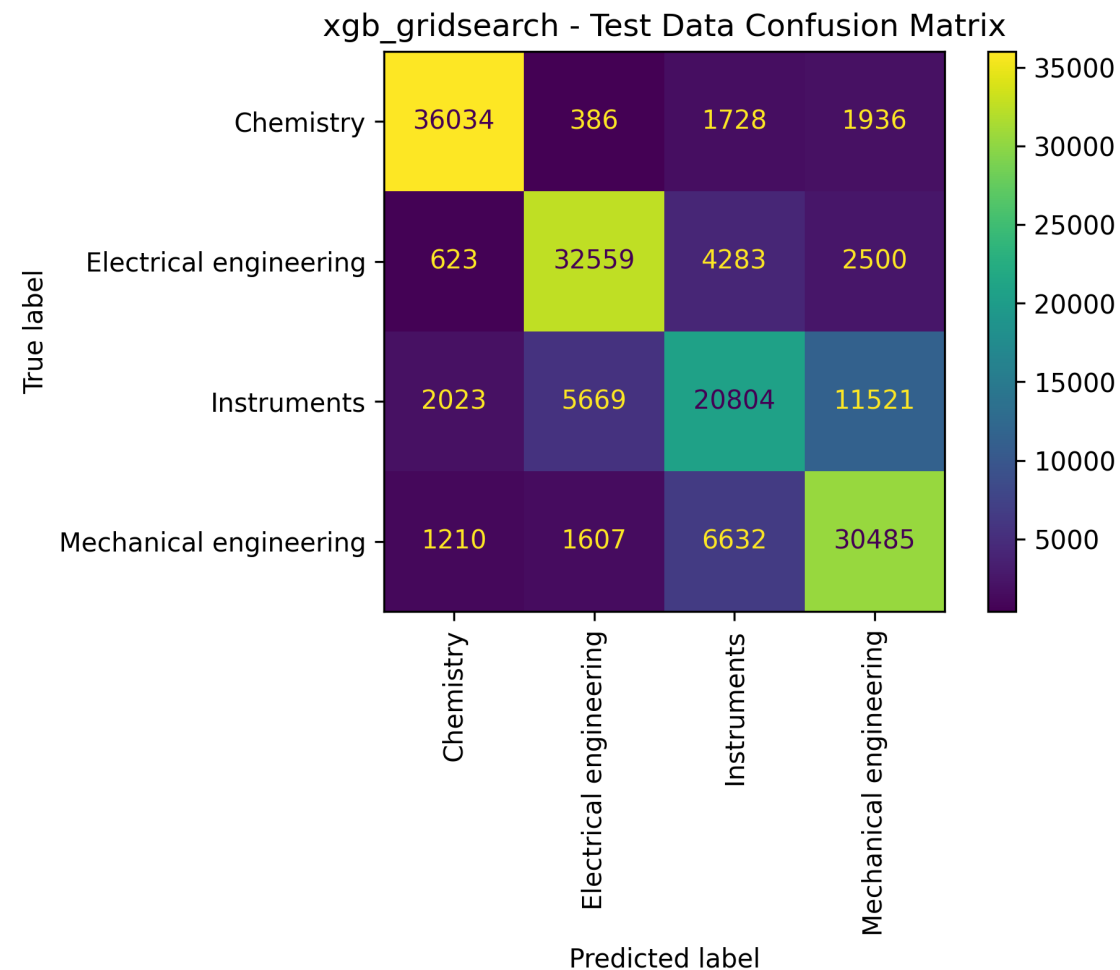
Distribution of Classes for Patents granted since 2011



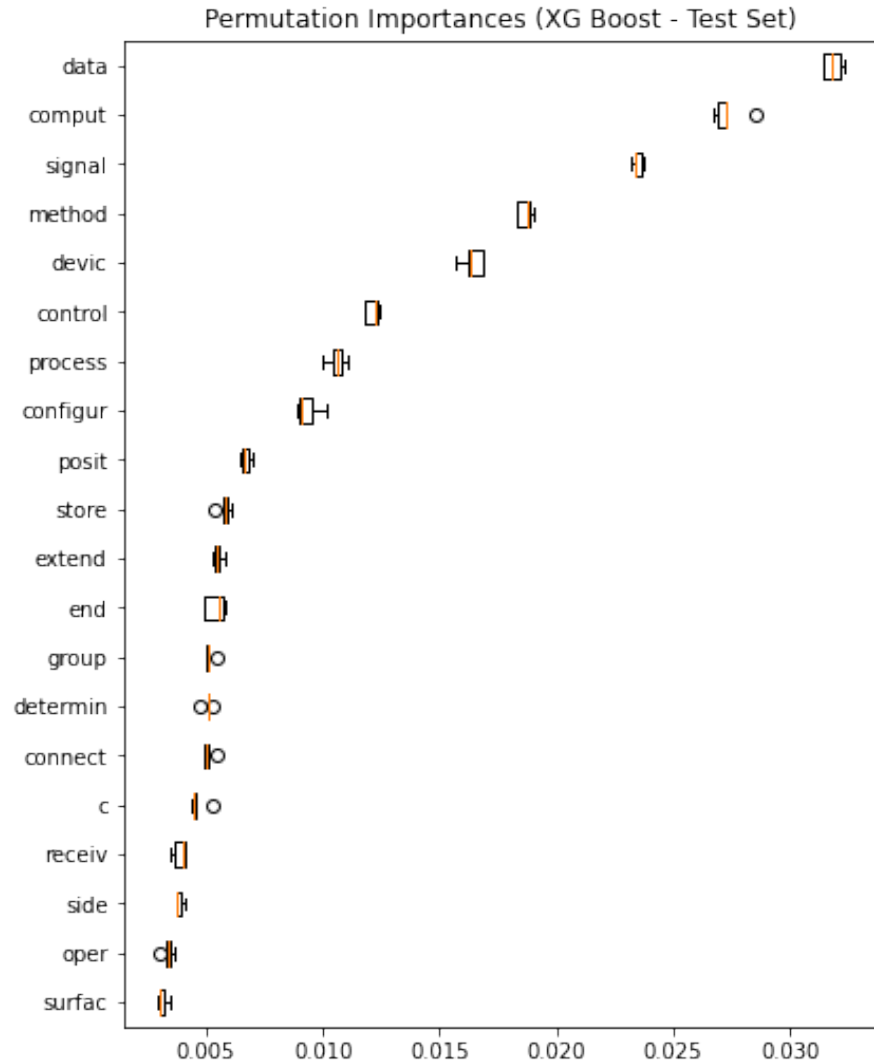
- Classifications:
 - Mechanical
 - Electrical
 - Chemical
 - Instruments
- Data imbalance

MODEL – XG BOOST

- Best Model: XG Boost
- Captures classifications
- Training Accuracy is: 72.1%
- Validation Accuracy is: 72.0%



EVALUATION



- From these importances, we can see the top 5 words are:
 - data
 - comput (stem of computer, computing, computation, etc.)
 - signal
 - method
 - devic (stem of device, devices, etc.)
- Ramifications for instruments class

RECOMMENDATIONS

- Use machine learning to help automate the classification process for patent Offices or third parties.
- Proof of concept for usefulness of machine learning



FUTURE WORK

1. Apply model to more complex classification systems
2. Finalize deep learning model



THANK YOU

NAME: CATHERINE FRITZ

EMAIL: CMFRITZO@GMAIL.COM

GITHUB: [@CMFRITZ](#)

LINKEDIN: [LINKEDIN.COM/IN/CATFRITZ](https://www.linkedin.com/in/catfritz)

ADDITIONAL INFORMATION CAN BE FOUND AT
[HTTPS://GITHUB.COM/CMFRITZ/PROJECT_4_NLP](https://github.com/CMFRITZ/PROJECT_4_NLP)