# Natural Language Processing

Review of applications and opportunities for IEG

Aivin V. Solatorio

Data Scientist | DECAT

asolatorio@worldbank.org

@avsolatorio

# POLL: How familiar are you with NLP?

1. I have never heard of it
2. I have heard about it but I don't know what it is
3. I am familiar with NLP but not a practitioner
4. I have worked in projects that use NLP
5. I am an expert in NLP

# What is Natural Language Processing (NLP)?
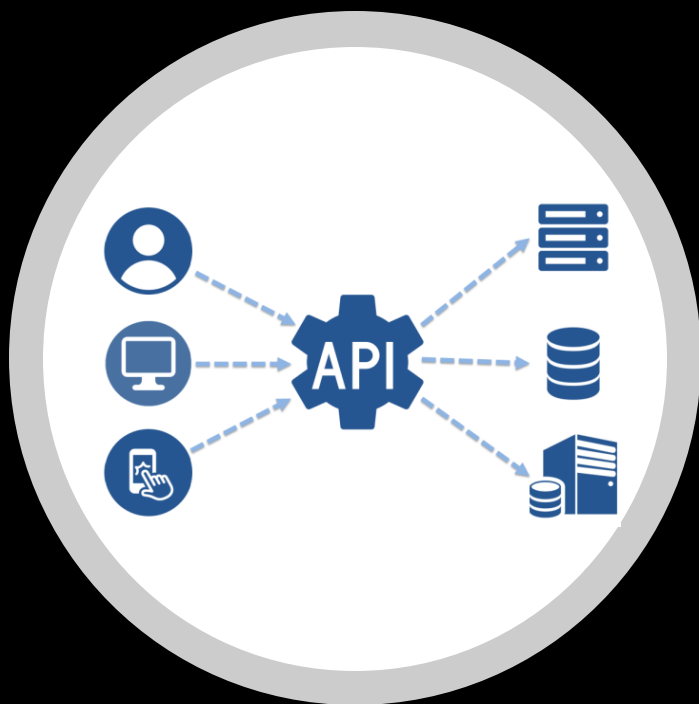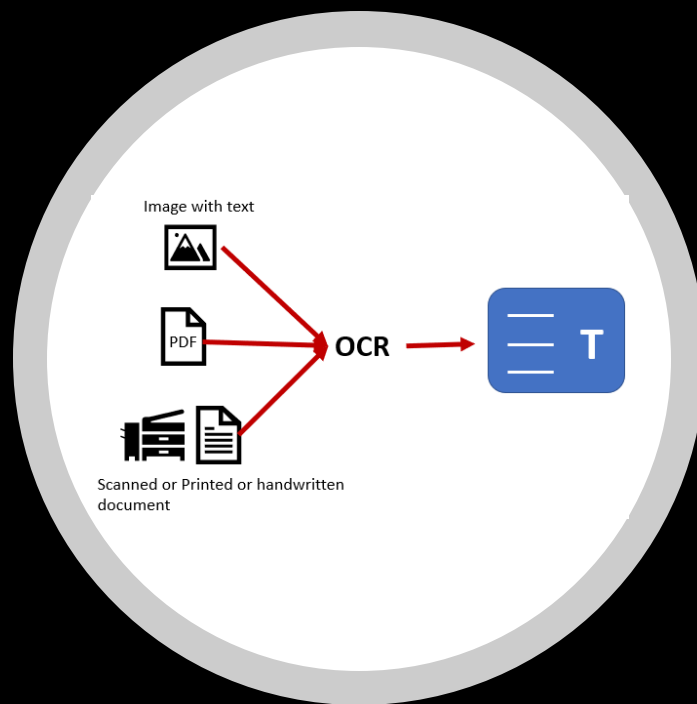
# Basic NLP Workflow
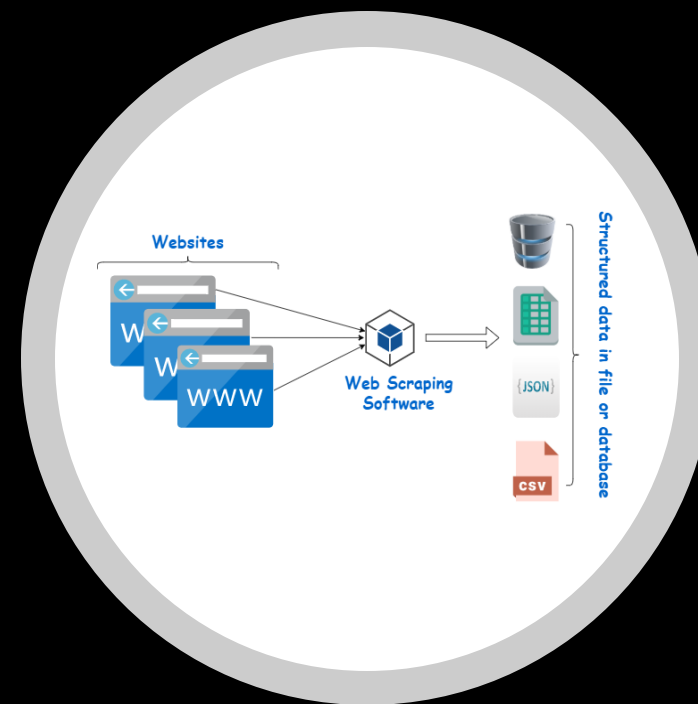
| Data acquisition | → | Text processing and cleaning | → | Model building, evaluation, and iteration |

API scraping

Document parsing

Web scraping

Data acquisition

| Raw | Lowercased |
|---|---|
| Canada<br>CanadA<br>CANADA | canada |
| TOMCAT<br>Tomcat<br>toMcat | tomcat |

| Raw | Normalized |
|---|---|
| 2moro<br>2mrrw<br>2morrow<br>2mrw<br>tomrw | tomorrow |
| b4 | before |
| otw | on the way |
| :)<br>:-)<br>;-) | smile |

| text-word | to | lemma |
|---|---|---|
| help | | help (v) |
| helps | | help (v) |
| helping | | help (v) |
| helped | | help (v) |

Stop Words

Common Words, Determiners, Conjunctions ...

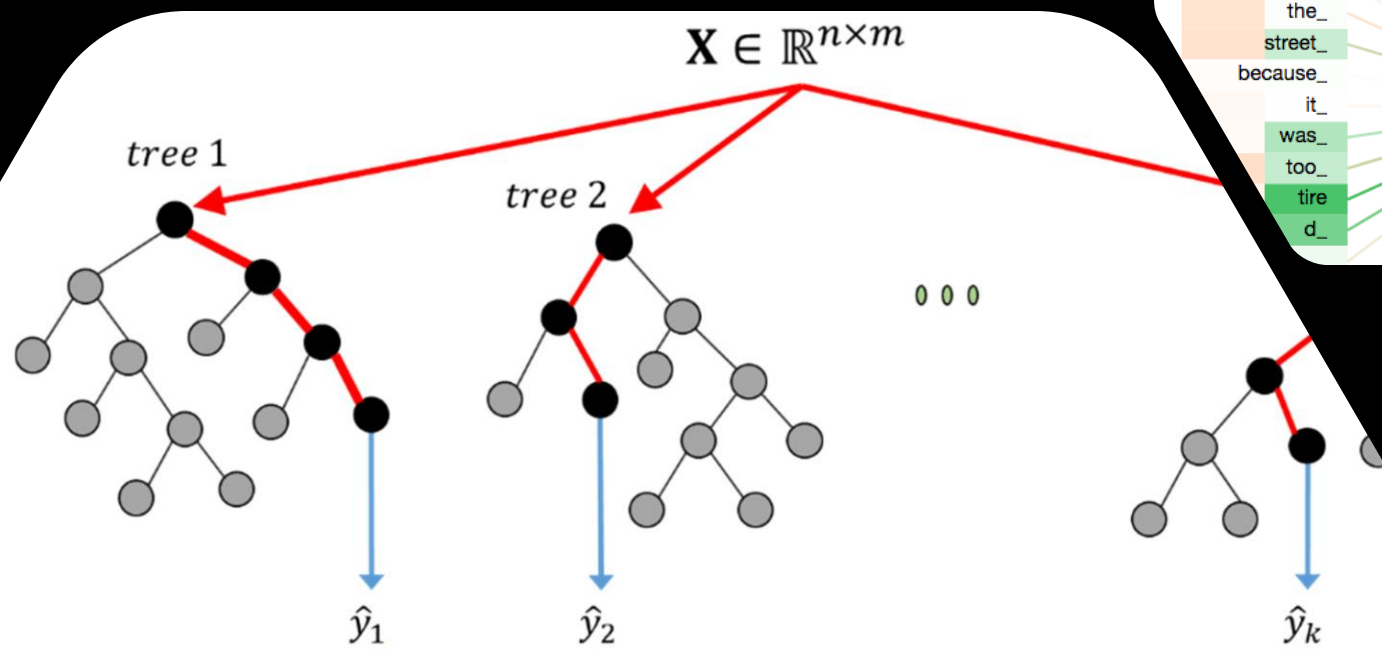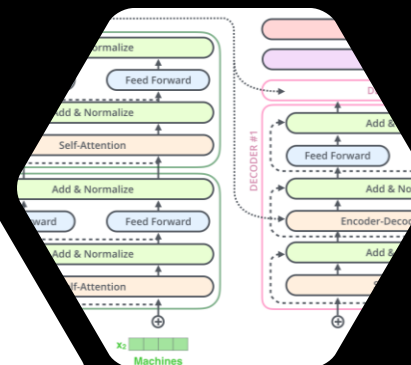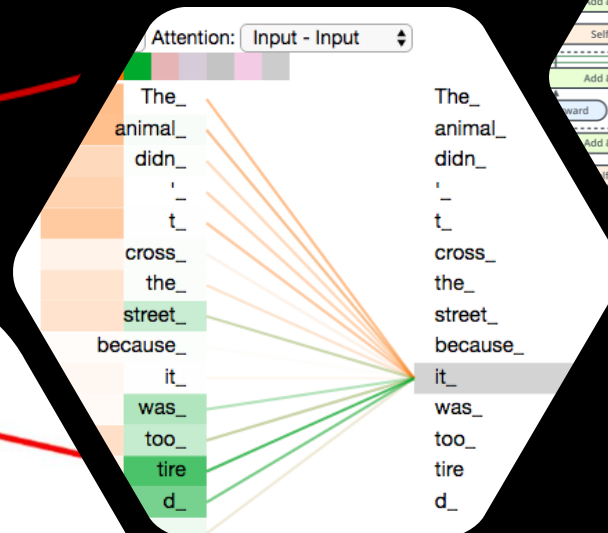Lowercase text  Token normalization  Word lemmatization  Stopwords removal

# Basic text processing and cleaning
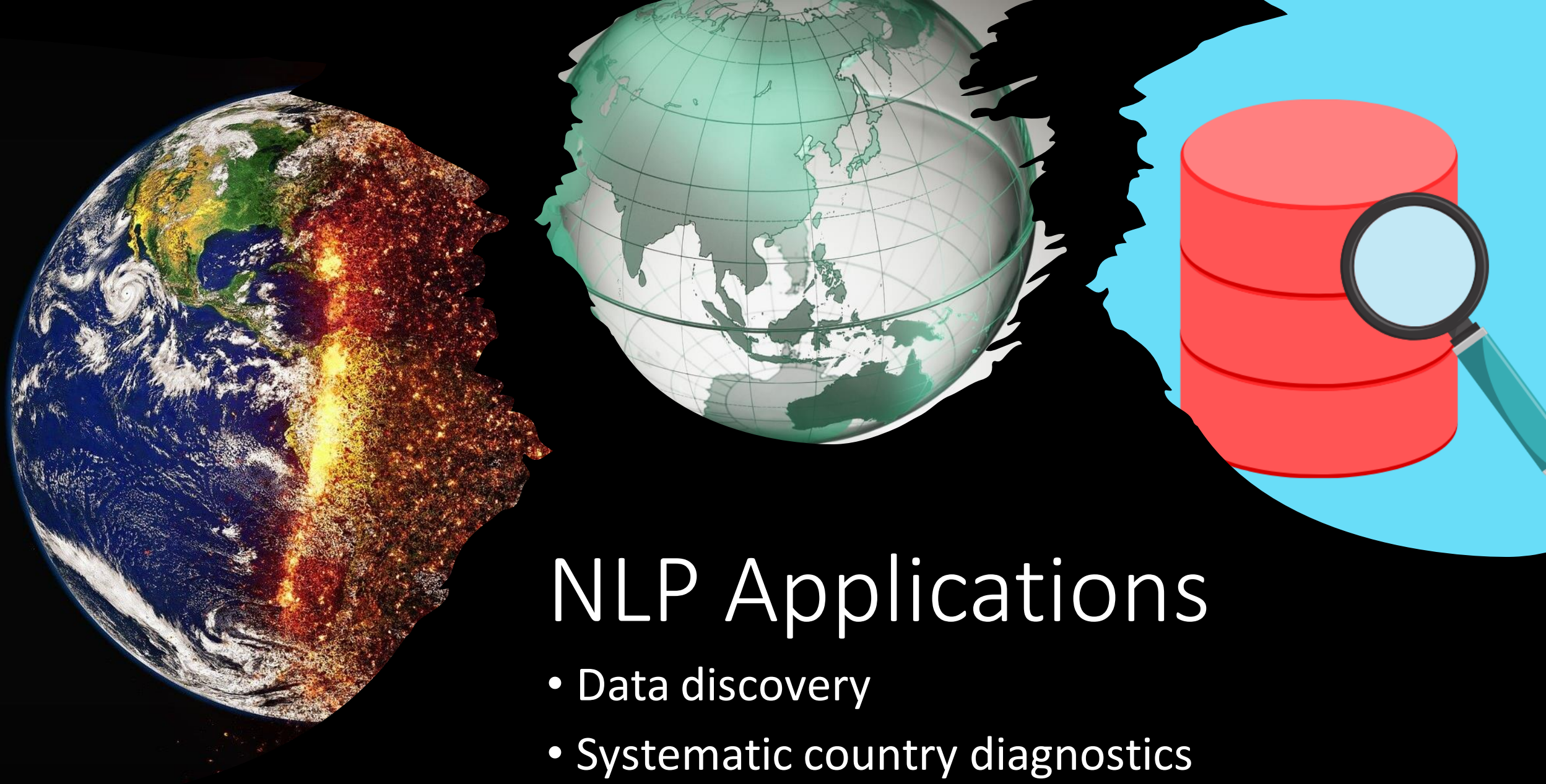
Training Data

Validation Data

Attention: Input - Input

The_      The_
animal_     animal_
didn_      didn_
'_       '_
t_       t_
cross_     cross_
the_      the_
street_    street_
because_   because_
it_       it_
was_     was_
too_      too_
tire      tire
d_       d_

$\mathbf{X} \in \mathbb{R}^{n \times m}$

tree 1      tree 2     o o o

$\hat{y}_1$      $\hat{y}_2$      $\hat{y}_k$

$$\hat{y} = \frac{1}{K} \sum_{k=1}^{K} \hat{y}_k$$

Model
building, evaluation,
and iteration

# NLP Applications

- Data discovery
- Systematic country diagnostics
- Climate change co-benefits prediction

Data Discovery

# Why is data discovery important?

- Different types of data are available across our institution and beyond.

- This is an asset when we can find the best data or document for our needs out of all the available data.

- However, it **proves to be a liability** when these data only reside in silos and databases that are **not equiped with discovery tools** to allow users to use them.

# Semantic understanding

# How could data discovery be relevant for IEG?

- May be able to provide a tool for more comprehensive review of reports.

- Discover related projects based on previous reports.

- Curate audit materials based on project topics.

# Systematic Country Diagnostics

# Systematic Country Diagnostics

- Use topic modeling to understand changes in legislative agenda across change in government.

- Analyze the general sentiment in the country from global news data.
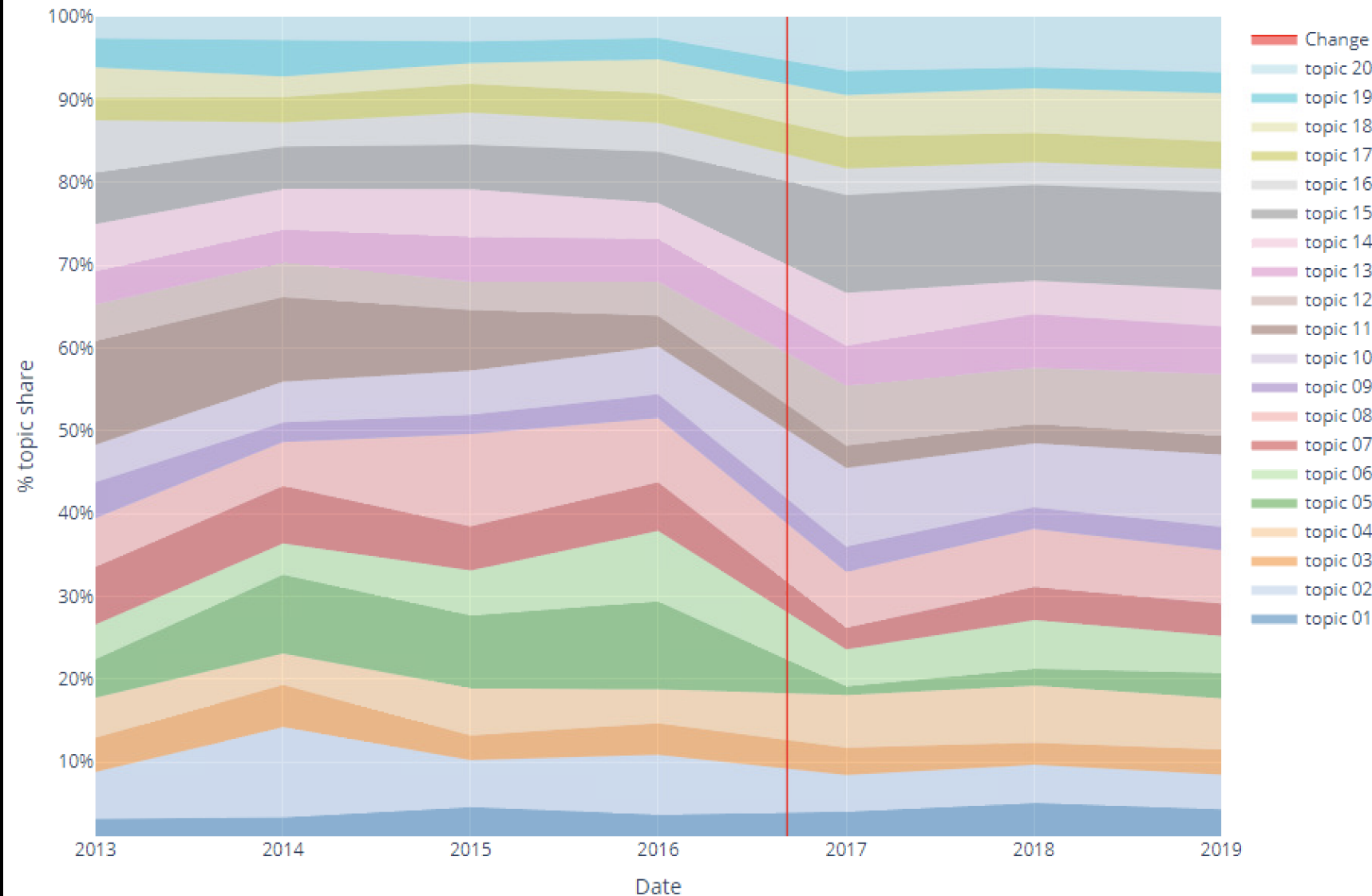
- Demonstrate the feasibility of proxy methods to understand certain country-level information using NLP.
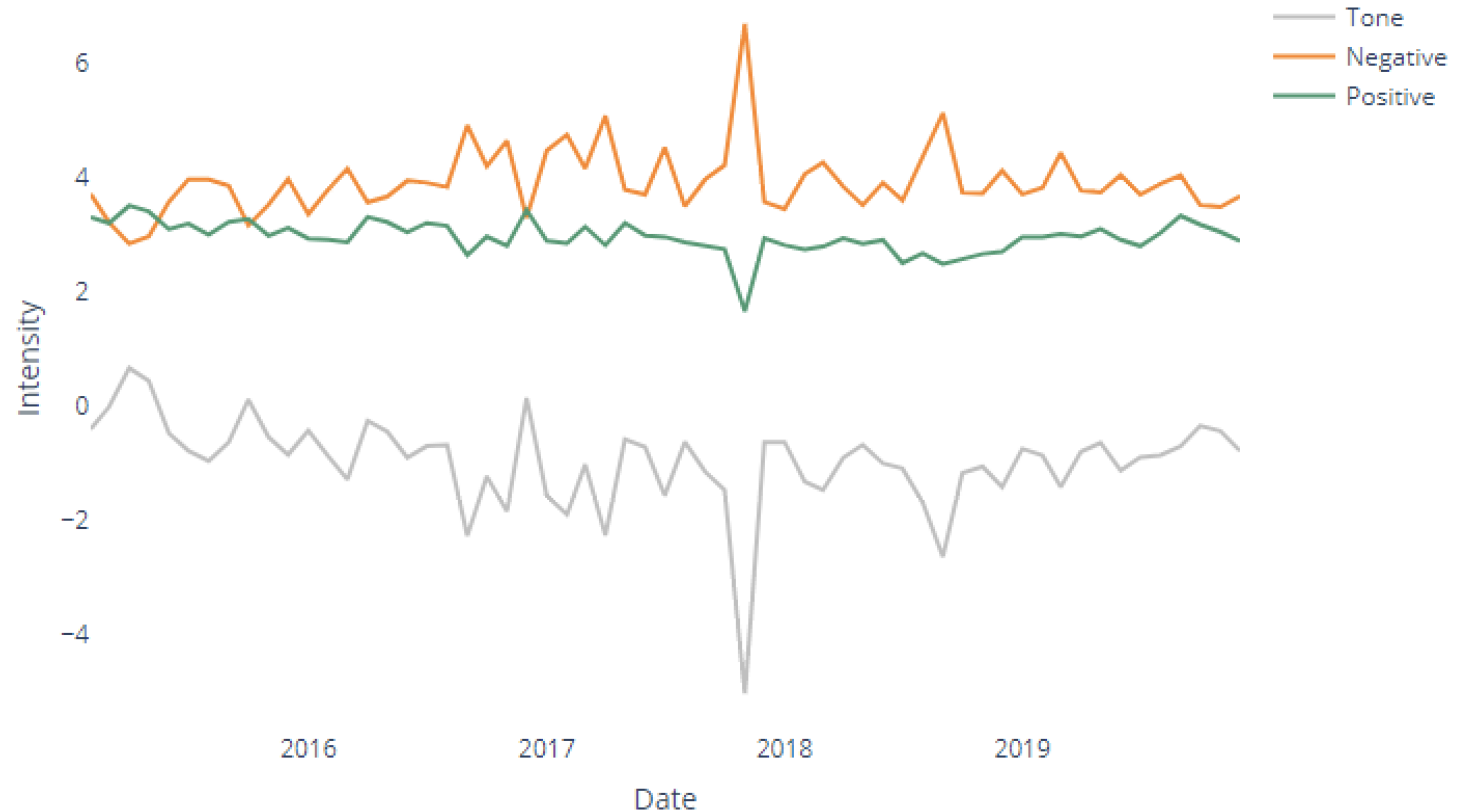
# Topic evolution and changes in legislation priorities

# Capturing sentiment from news

## Sentiment for all data on GDELT

|  | Tone | Negative | Positive | Polarity |
|---|---|---|---|---|
| 5th quantile | -2.1412 | 3.2086 | 2.5757 | 6.3134 |
| mean | -0.9680 | 3.9267 | 2.9587 | 6.8854 |
| 95th quantile | 0.0753 | 4.8696 | 3.3518 | 7.6016 |

# How might IEG benefit from sentiment and topic models?

- Monitor impact of projects based on social media sentiment.

- Improve the management of projects by using topic models.

- Understand trends between successful and unsuccessful projects.

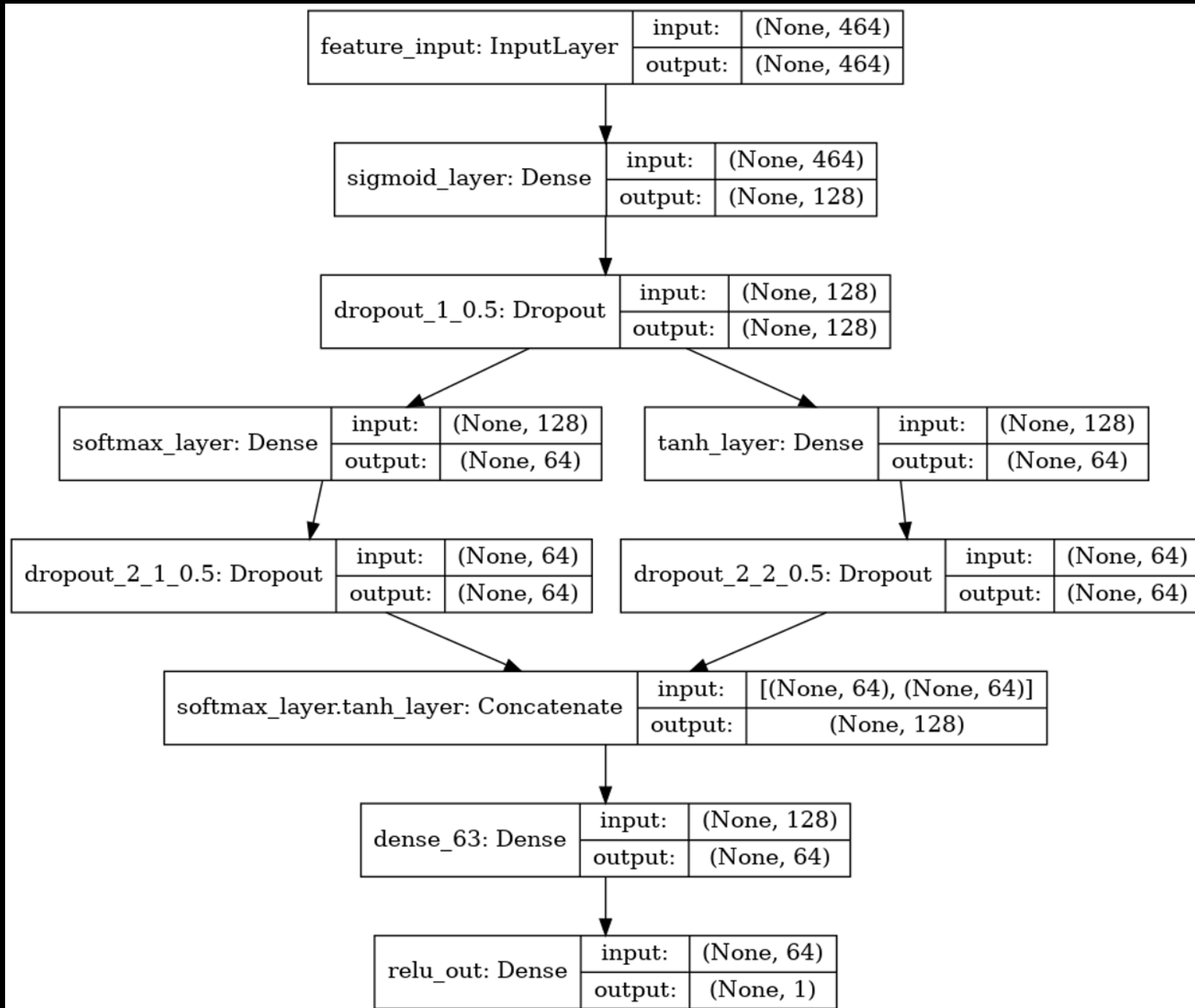# Climate change co-benefits prediction
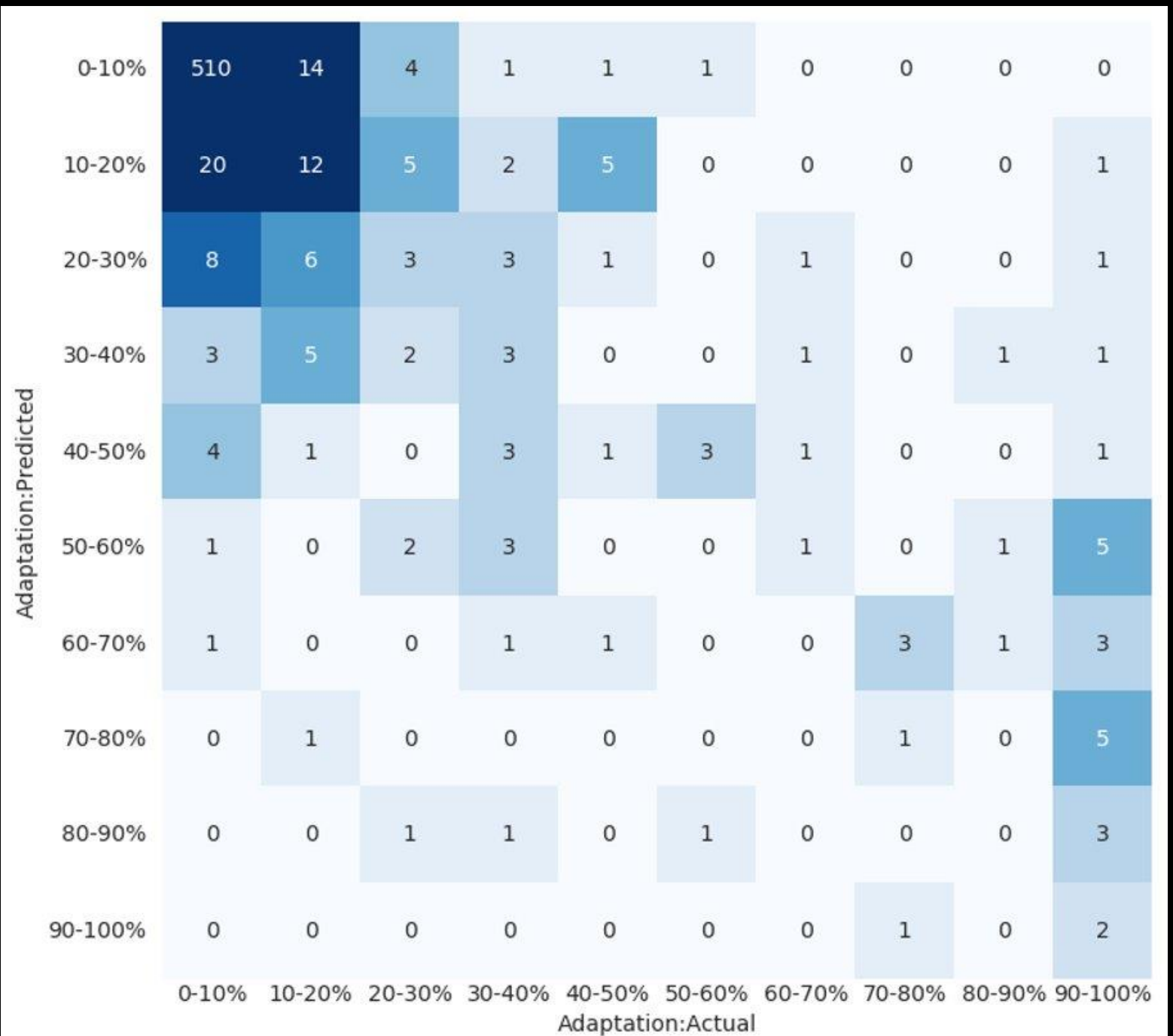
# Climate change co-benefits prediction

- Process and perform text mining on project documents and enrich the available data using other sources.

- Use pretrained NLP models and machine learning to predict climate change co-benefits.

- The project was exploratory in nature. A baseline exists and our model has outperformed this baseline.

# Neural network model

# Predictive performance

# Climate change co-benefits prediction

- While we have shown feasibility of improved performance, the degree of accuracy was still not at the desired level for automation.

- The primary reason for poor performance was the limited volume of available data.

# How machine learning with NLP could be beneficial in IEG?

- Automate data extraction from text.

- Detect fraud or potential issues in projects given historical data.

# POLL: Which NLP application will be immediately useful for IEG?

1. May be able to provide a tool for more comprehensive review of reports.

2. Discover related projects based on previous reports.

3. Curate audit materials based on project topics.

4. Monitor impact of projects based on social media sentiment.

5. Improve the management of projects by using topic models.

6. Understand trends between successful and unsuccessful projects.

7. Automate data extraction from text.

8. Detect fraud or potential issues in projects given historical data.

# Natural Language Processing

Review of applications and opportunities for IEG

Aivin V. Solatorio

Data Scientist | DECAT

asolatorio@worldbank.org

@avsolatorio