

National Institute of Allergy and Infectious Diseases

NIAID Data Science (NDS) Dashboard

Office of Data Science and Emerging Technologies
Anuva Banwasi and Snipta Mallick



National Institute of
Allergy and
Infectious Diseases

Dr. Chris Marcum
Staff Scientist
Mentor

Dr. Steve Tsang
Health Scientist
Mentor

NIAID

Introductions



Anuva Banwasi

Software Engineering Fellow
Columbia University '24
Computer Science



Snipta Mallick

Software Engineering Fellow
University of Texas at Dallas '22
Computer Science & Cognitive Science

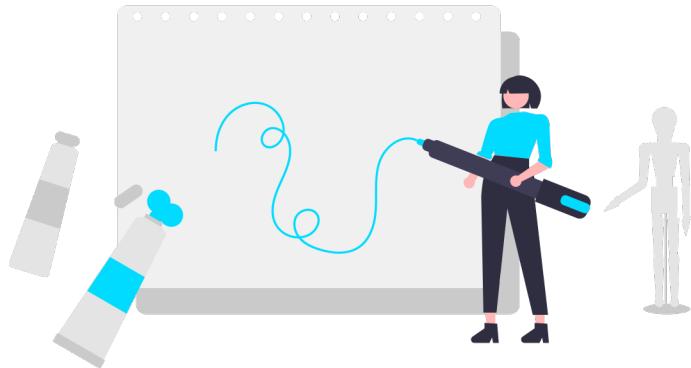


National Institute of
Allergy and
Infectious Diseases

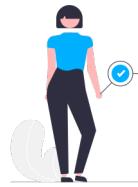


Problem Statement

Absence of an **integrated, automated** workflow to analyze, categorize, and visualize NIAID data science projects to streamline the funding application management process



Civic Digital Fellowship Project Goals



Develop an automated pipeline to characterize and analyze NIAID funded data science projects



Build a platform-independent dashboard to annotate and visualize the rich portfolio



Compile program scripts and documentation in Github repository



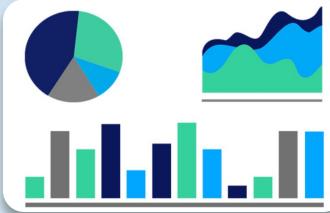
National Institute of
Allergy and
Infectious Diseases

NIAID

Project Components



DINAL most visited in the European Union ORG, and the most popular tourist destination in Italy. It is listed by UNESCO ORG as a World Heritage Site.[14] Host city for the United Nations Educational, Scientific and Cultural Organization ORG (FAO ORG), the World Food Programme, the International Fund for Agricultural Development ORG (IFAD ORG). The city also hosts the European Parliament Assembly ORG of the Union for the Mediterranean[15] (UIM ORG) and is the seat of several specialised agencies of the United Nations. Rome is also the seat of several international business companies such as Eni ORG, Enel GPE, TIM ORG, and Unicredit ORG. It is also the home of many companies involved in the oil industry, the pharmaceutical industry, and the food industry.



**NIH
RePORTER
Data
Module**

**Machine
Learning
(ML)
Classifier
Module**

**Named
Entity
Recognition
(NER)
Module**

**NIAID Data
Science
(NDS)
Dashboard**



National Institute of
Allergy and
Infectious Diseases

CIVIC
DIGITAL
FELLOWSHIP

NIAID

Data Module

- Extract publicly-available funded project data from NIH RePORTER
 - Title, Abstract, Project Number, Fiscal Year, Funding Institute, etc.
- Pipeline script that downloads data and runs classifier on project abstracts

The screenshot shows the NIH RePORTER search results page. The search criteria are set to Fiscal Year: 2020, Admin: Yes, Agency/Institute/Center: NIAID, and Text Search: \software tools\ OR (software AND develop*) OR (\R Package\ OR python OR \software bioinformatics pipeline*) OR \computational analysis tool\ OR \necessary resource\ OR github OR \social media\ OR \social media mining\ NOT \administrative core\ OR \missing data\ OR \missing at random\ OR \missing completely at random\ OR input\ OR \geospatial model\ OR \epidemic model\ OR \Mobile Health\ OR mHealth OR \Mobile Application Strategies\ OR \computational pharmacology\ OR \electronic medical record\ OR emr OR \electronic health record\ OR ehr\ OR \longitudinal pattern\ OR \community resource\ OR \Computational Biology\ AND \bioinformatic\ OR biostatistic\ OR \bioinformatics technique\ OR \database AND (develop* OR maint* OR update*) OR \data management\ OR \Bioinformatics Resource Centers\ OR \bootstrap OR \Markov Chain Monte Carlo\ OR \sequential Monte Carlo\ OR \simulation approach\ OR \text mining\ OR \network model\ OR \agent-based model\ OR \population-level model\ OR \machine learning\ OR \predictive model\ OR \deep learning\ OR \data science\ OR \big data\ OR \data analytics\ advanced) Limit to: Project Title, Project Terms, Project Abstracts

Project ID	Principal Investigator(s)/Project Leader(s)	Organization	Fiscal Year	Admin IC	Funding IC	FY Total Cost by IC	Similar Projects
SR01AT132030-05	YOUNG, SEAN C	UNIVERSITY OF CALIFORNIA-IRVINE	2020	NIAID	NIAID	\$737,082	View >
SU19AT135972-03	PACHE, LARS C	ICAHN SCHOOL OF MEDICINE AT MOUNT SINAI	2020	NIAID		\$570,130	View >
SR01AT116770-05	KENAH, EBEN C	OHIO STATE UNIVERSITY	2020	NIAID	NIAID	\$347,898	View >

Machine Learning Classifier

Determine probability that a project is data-science related

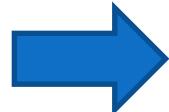
1

Load and Prepare Dataset

Bag of words model

2020 RePORTER
projects within NIAID

Manually coded:
200 positives
200 negatives



word	weight
lorem	17
ipsum	16
amet	13
sit	13
vitae	13
egestas	12
quis	11
sem	11
massa	10
nec	10
sed	10
semper	9
egestas	8
libero	8
nisl	8
odio	8
scelerisque	8

Figures Designed By Dr. Chris Marcum



National Institute of
Allergy and
Infectious Diseases



Machine Learning Classifier

2

Train and test supervised ML models

2020 RePORTER NIAID manually coded data

- Linear SVC, Ridge classifier, Perceptron, etc.
- Outputs probability that project is data science

3

Second test on 2019 RCDC data

Research, Condition, and Disease Categorization



National Institute of
Allergy and
Infectious Diseases



NIAID

Classifier Validation

PREDICTED

		TRUE	
		Data Science	Not Data Science
Data Science	Data Science	25	4
	Not Data Science	0	21

2019 RCDC Data Science

PREDICTED

		TRUE	
		Data Science	Not Data Science
Data Science	Data Science	192	2484
	Not Data Science	3	5933

2020 NIAID RePORTER

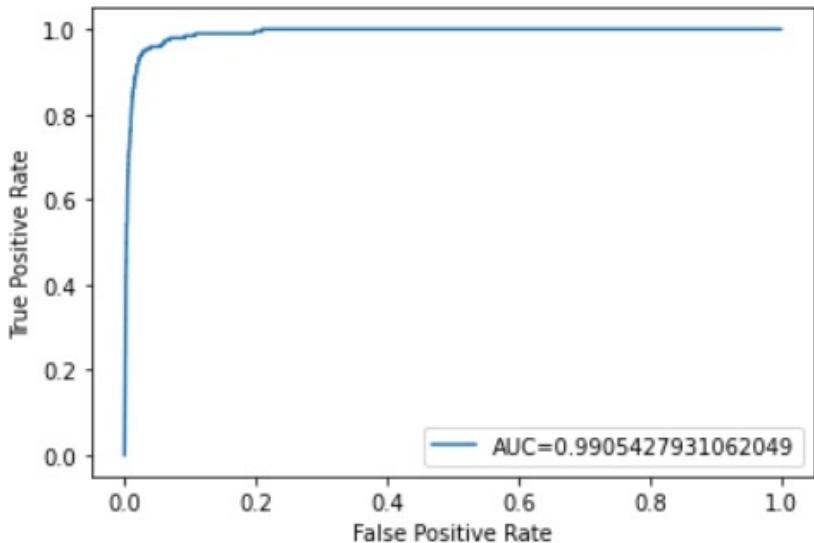


National Institute of
Allergy and
Infectious Diseases



ML Classifier Results and Insights

ROC Curve for 2020 Grants



Ridge classifier

F1 score: 0.94

Weighted average: 0.92



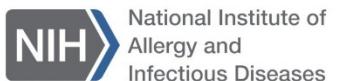
National Institute of
Allergy and
Infectious Diseases



Named Entity Recognition Module

" Biostatistics ✕ / Epidemiology ✕ Training ✕ Grants in AIDS. Summary/Abstract This is a request for an extension of the Biostatistics ✕ and Epidemiology ✕ in AIDS Training ✕ Program at the Harvard T. H. Chan School of Public Health (HSPH). The program prepares pre-doctoral and postdoctoral fellows ✕ in the application of biostatistics ✕ and epidemiology ✕ to HIV research. The program is an active collaboration between the Departments of Biostatistics ✕ and Epidemiology ✕. Trainees ✕ receive high-quality instruction in basic biostatistical ✕ theory and methods, such as probability, statistical inference, computing and data analysis ✕. The program also provides training ✕ in specialized topics of particular relevance for HIV/AIDS applications, such as survival, longitudinal and multivariate data analysis ✕, causal inference and statistical genetics. Options are available for training in

Sample NER Annotation



Career/Workforce Development

Data Management Centers

Informatics Research

Repository/Knowledgebases

Software/Tools

ODSET Data Science Subcategories



Named Entity Recognition Model

1

Create Training Set using Positively Labeled Data

Doccano: open-source text annotation tool

2

Build Natural Language Processing Model

spaCy: open-source NLP python library

3

Display Named Entity Recognition Results

spacy-streamlit: library to visualize spacy models



National Institute of
Allergy and
Infectious Diseases



Potential Users



National Institute of
Allergy and
Infectious Diseases

CIVIC
DIGITAL
FELLOWSHIP

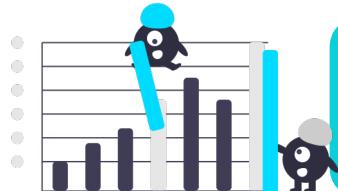
Future Directions



Support data analysis initiatives for additional institutes and projects across the NIH



Visualize the downstream impacts of grant funding by displaying network analysis models



Add containerization to improve the FAIRness and reproducibility



National Institute of
Allergy and
Infectious Diseases



Acknowledgements

- **NIAID/ODSET**

- Dr. Chris Marcum
- Dr. Steve Tsang
- Sydney Foote
- Sara Jones

- **Coding it Forward**

- Rachel Dodell
- Ariana Soto



National Institute of
Allergy and
Infectious Diseases

CIVIC DIGITAL FELLOWSHIP

coding it forward >

Thank you!

Any questions?

Slide Image Sources

- Figures on Slide 6
 - Google Images
- Figures on Slides 4, 5, and 15
 - Katerina Limpitsouni, Undraw.co
- Figures on Slide 8
 - Dr. Chris Marcum



National Institute of
Allergy and
Infectious Diseases

NIAID