# PROJECT COMPONENTS

## 3 main tasks

- Create an interactive portfolio overview dashboard of NIGMS grant funding practices in Tableau

- Assist in current dashboarding efforts

- Help Nate and Jordan build a classifier for stem cell grants

# PORTFOLIO OVERVIEW DASHBOARD

**Goal**

- NIGMS 5-year strategic plan calls for monitoring and evaluating activities across the institute

- Very broad, exploratory purpose

- Allow NIGMS funding to be more equitable and visualize areas of improvement in the grant funding process

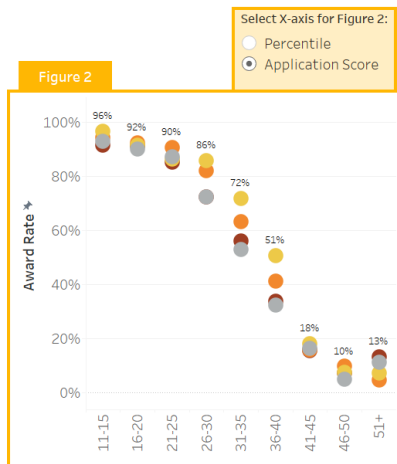# PORTFOLIO OVERVIEW DASHBOARD

**Approach**

- Challenge: maintain simplicity and intuitiveness while visualizing variables of interest for a non-specific audience

- Consistent meetings and feedback regarding changes

-  Noting places of confusion and inefficiency

- Determining variables of interest

# RESULTS

# RESULTS

# IMPROVEMENTS

# IMPROVEMENTS

Redacted Slide: Other Dashboarding Efforts

Redacted Slide: Other Dashboarding Efforts

# STEM CELL CLASSIFIER

**Goal**

- Assist Nate and Jordan in building a model that could classify the stem cell lines used within grants

- Some grants had been misclassified in the past

- Provide a tool that improves efficiency of stem cell classification

- Desire to assist POs in classifying these grants as human vs. non-human:
    - Embryonic
    - Non-embryonic
    - Induced pluripotent
    - Umbilical / placenta

# STEM CELL CLASSIFIER

**"The less data science you use, the better"**

RULE-BASED APPROACH → WORD COUNTS → TF-IDF → TRANSFORMERS → TARGET-DEPENDENT SENTIMENT ANALYSIS

# 1. RULE-BASED APPROACH

## Data

| AID | Research Strategy | Specific Aims | Label | Humans Used | Animals Used | Stem Cells Used |
|-----|-------------------|---------------|-------|-------------|--------------|-----------------|
| 999999 | lorem | dolor | Embryonic | Yes | No | Yes |
| 999998 | ipsum | sit | IPSC | No | No | No |

# 1. RULE-BASED APPROACH



| AID | Research Strategy | Specific Aims | Label | Humans Used | Animals Used | Stem Cells Used |
|---|---|---|---|---|---|---|
| 999999 | lorem | dolor | Embryonic | Yes | No | Yes |
| 999998 | ipsum | sit | IPSC | No | No | No |

Humans Used = No & Animals Used = Yes

Stem Cell Type

Have embryonic stem cells been used?

Stem Cell Type

Human
Non-Human

No
Yes

# BACK TO THE DRAWING BOARD



RULE-BASED APPROACH → WORD COUNTS

| AID | Research Strategy | Specific Aims | Label | Humans Used | Animals Used | Stem Cells Used |
|---|---|---|---|---|---|---|
| 999999 | lorem | dolor | Embryonic | Yes | | Yes |
| 999998 | ipsum | sit | IPSC | | | No |

CIVIC DIGITAL FELLOWSHIP

# 2. WORD COUNTS

| AID | Research Strategy | Specific Aims | Label | Stem Cells Used |
|---|---|---|---|---|
| 999999 | lorem | dolor | Embryonic | Yes |
| 999998 | ipsum | sit | IPSC | No |

Isolating words that modify "stem" or "cell"

"Mesenchymal stem cells (MSCs) adult stem cells that can differentiate into a variety of cell types..."

Text cleaning and lemmatization

**Word Bank**

{0:"mesenchymal",
1:"adult",
2:"mutated"}

| AID | Cleaned Text |
|---|---|
| 999999 | [mesenchymal stem cell msc adult stem cell] |

| AID | Cleaned Text |
|---|---|
| 999999 | [mesenchymal stem cell msc adult stem cell] |

Run clean text through word bank to obtain frequencies per application

**Word Bank**
{0:"mesenchymal",
1:"adult",
2:"mutated"}

**Word Frequency by Stem Cell Type**

| | Embryonic | Non-Embryonic | Induced Pluripotent | Placenta |
|---|---|---|---|---|
| mesenchymal | 300 | 30000 | 40 | 1 |
| adult | 10000 | 1 | 50 | 0 |
| mutated | 30 | 2000 | 20 | 10 |

Aggregate frequencies by stem cell label

*Assumption: "mesenchymal" in text suggests non-embryonic application*

Create counts of words that appear according to category

**Count of Words Belonging to Stem Cell Type**

| AID | Embryonic | Non-Embryonic | Induced Pluripotent | Placenta |
|---|---|---|---|---|
| 999999 | 200 | 30 | 20 | 0 |
| 999998 | 10 | 20 | 100 | 1 |

# SMOTE

## Count of Words Belonging to Stem Cell Type

| AID | Embryonic | Non-Embryonic | Induced Pluripotent | Placenta |
|-----|-----------|---------------|---------------------|----------|
| 999999 | 200 | 30 | 20 | 0 |
| 999998 | 10 | 20 | 100 | 1 |



Stem Cell Type

SMOTE:
Synthetic
Oversampling

# Synthetic Minority Oversampling Technique



Original Dataset      Generating Samples      Resampled Dataset

| AID | Placenta |
|---|---|
| 999999 | 0 |
| 999998 | 1 |

| AID | Placenta |
|---|---|
| 999999 | 0 |
| ---------- | 0.5 |
| 999998 | 1 |

# 3. TFIDF

RULE-BASED APPROACH → WORD COUNTS → TF-IDF

CIVIC
DIGITAL
FELLOWSHIP

# SEGWAY: BALANCED ACCURACY

## Confusion Matrix

|  | Actually Positive (1) | Actually Negative (0) |
|---|---|---|
| Predicted Positive (1) | True Positives (TPs) | False Positives (FPs) |
| Predicted Negative (0) | False Negatives (FNs) | True Negatives (TNs) |

**Accuracy – Fraction correctly identified:**
(True Positive + True Negative) / Total

**Sensitivity – Fraction of Positives correctly identified:**
True Positive / (True Positive + False Negative)

**Specificity – Fraction of Negatives Correctly Identified:**
True Negative / (True Negative + False Positive)

**Balanced Accuracy = Average of Sensitivity & Specificity**

# 3. TFIDF

$$w_{i,j} = tf_{i,j} \times \log\left(\frac{N}{df_i}\right)$$

$tf_{i,j}$ = number of occurrences of $i$ in $j$

$df_i$ = number of documents containing $i$

$N$ = total number of documents

**Word Bank**

{0:"mesenchymal",
1:"adult",
2:"mutated"}

**+**

**RCDC CATEGORIES**

["stemness",
"retinal",
"progenitor"]

| AID | TF-IDF |
|---|---|
| 999999 | [["stemness", .30], ["adult", .01], ["progenitor", .15]] |
| 999998 | [["stemness", .01], ["adult", .70], ["progenitor", .30]] |

# 4. TRANSFORMERS

| RULE-BASED APPROACH | → | WORD COUNTS | → | TF-IDF | → | TRANSFORMERS |

**Balanced Accuracies:**

Embryonic – 85.35%

Non-Embryonic – 73.15%

Induced Pluripotent – 69.5%

Umbilical – 56.1%

# 4. TRANSFORMERS

# 4. TRANSFORMERS

```
┌─────────────────┐      ┌─────────────────┐      ┌─────────────────┐
│   RULE-BASED    │  →   │   WORD COUNTS   │  →   │     TF-IDF      │
│    APPROACH     │      │                 │      │                 │
└─────────────────┘      └─────────────────┘      └─────────────────┘
                                                           ↓
┌─────────────────┐      ┌─────────────────┐
│    TARGET-      │  ←   │  TRANSFORMERS   │
│   DEPENDENT     │      │                 │
│   SENTIMENT     │      │                 │
│   ANALYSIS      │      │                 │
└─────────────────┘      └─────────────────┘
```
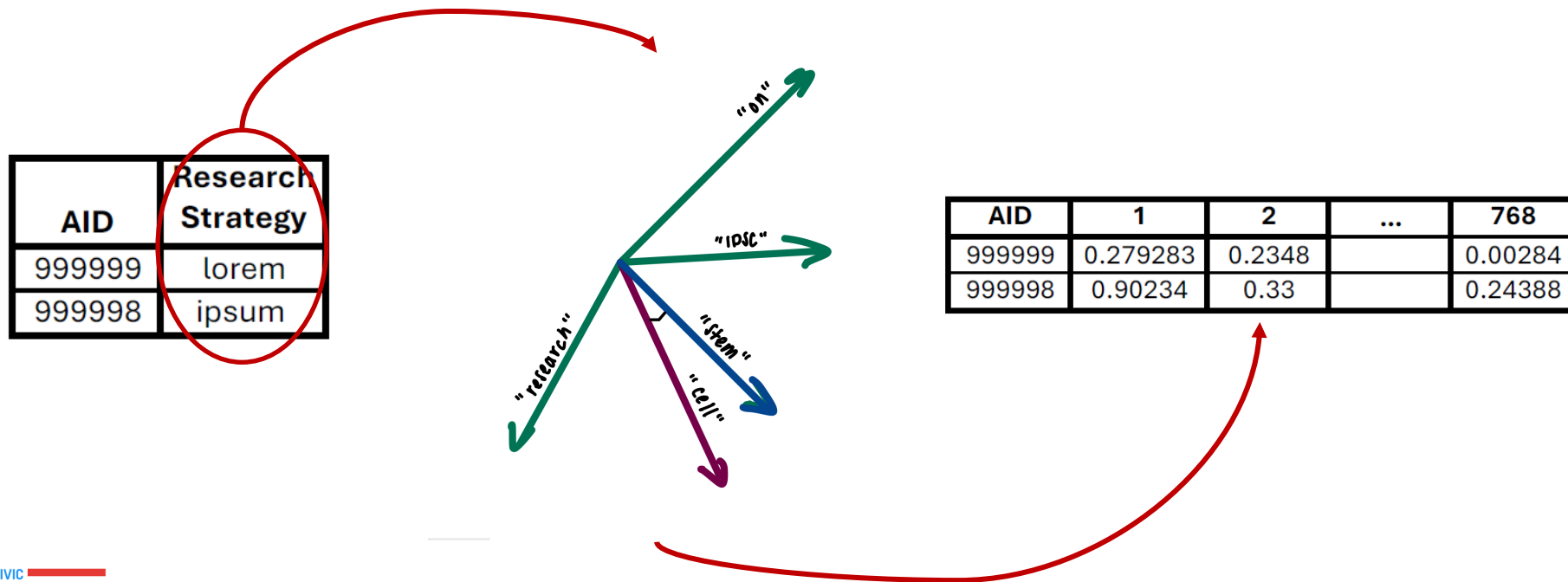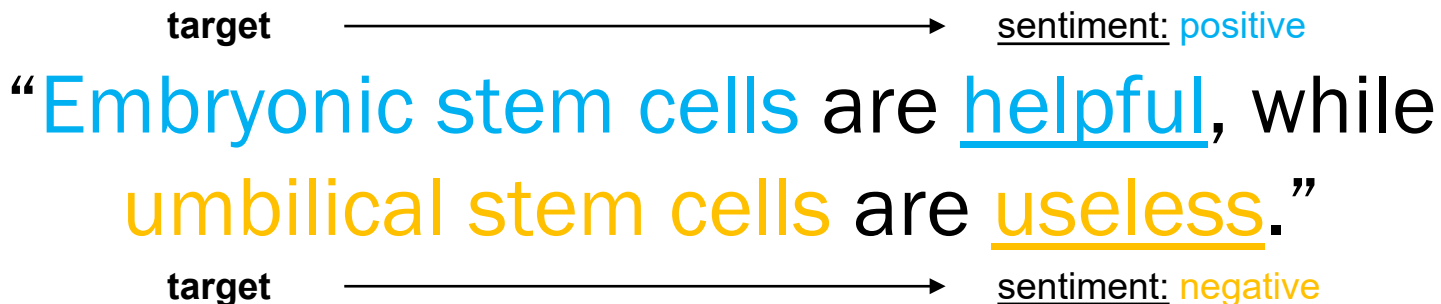
**Balanced Accuracies:**

Embryonic – 88%

Non-Embryonic – 70.2%

Induced Pluripotent – 77.2%

# 5. SENTIMENT ANALYSIS

target ⟶ sentiment: positive

"Embryonic stem cells are helpful, while umbilical stem cells are useless."

target ⟶ sentiment: negative

CIVIC
DIGITAL
FELLOWSHIP