# AUTO-TAGGING GLOBAL GIVING PROJECTS

BY: MELISSA BARONA

SPRINGBOARD CAPSTONE PROJECT 3

# PROBLEM IDENTIFICATION

**Problem:**
It can be tedious, expensive, and time-consuming to manually assign each project to a theme

**Solution:**
Create an auto-tagging system to classify projects using machine learning
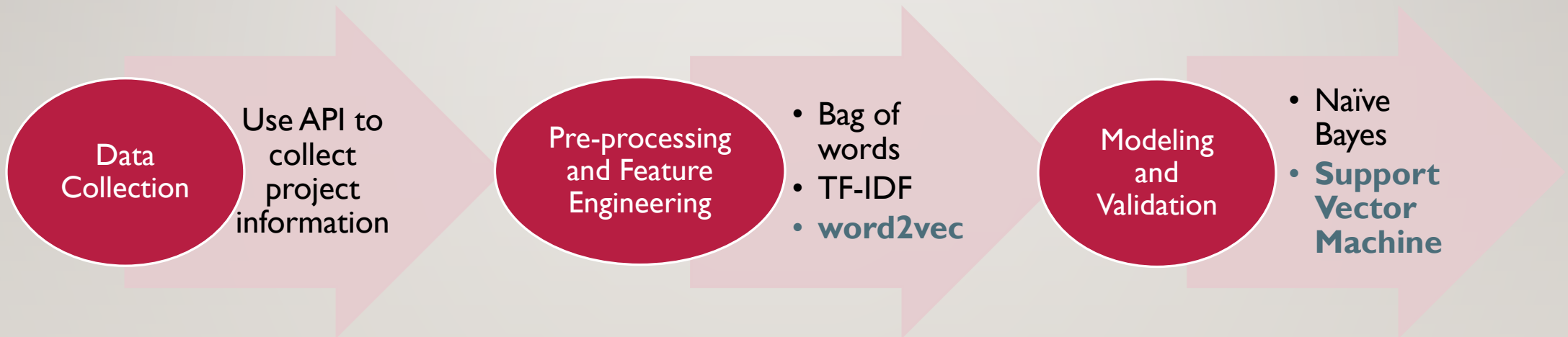
# PRODUCT VISION

- An auto-tagger that properly tags a project based on the context of its description
- An appropriate classification of a project to a theme is essential for attracting the ideal donor that will continue to support the mission of the project.



Education

Child
Protection

Gender
Equality

School
Protect children in schools
**girls in schools**

# OUR BEST CLASSIFIER:

## A WORD2VEC TRANSFORMATION AND SUPPORT VECTOR MACHINE MODEL

**Data Collection** → Use API to collect project information

**Pre-processing and Feature Engineering** →
- Bag of words
- TF-IDF
- **word2vec**

**Modeling and Validation** →
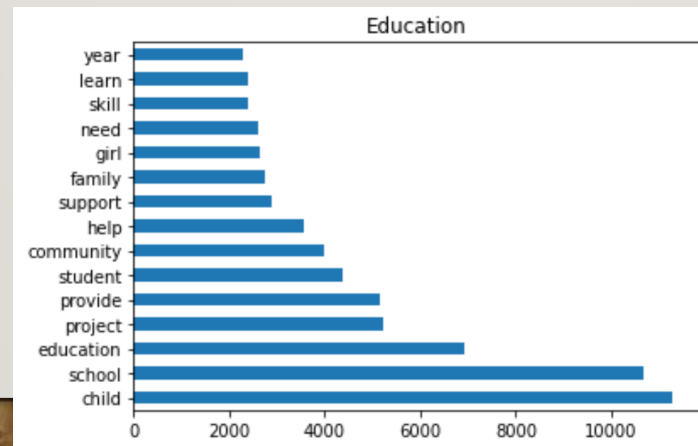- Naïve Bayes
- **Support Vector Machine**

# OUR BEST CLASSIFIER:

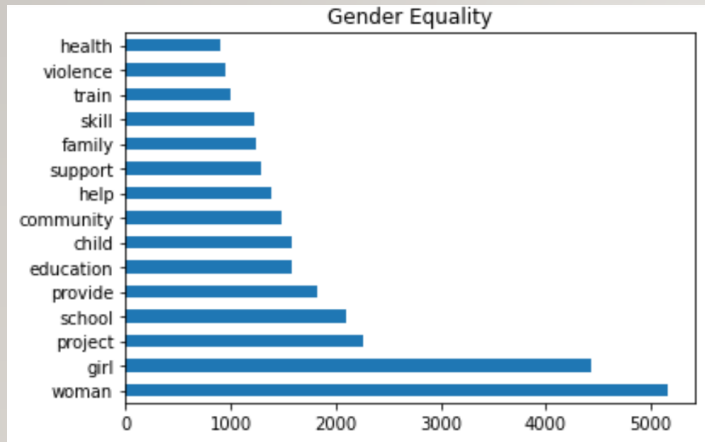## A WORD2VEC TRANSFORMATION AND SUPPORT VECTOR MACHINE MODEL

| Model | Overall Accuracy | |
|---|---|---|
| | Training Set | Testing Set |
| **TF-IDF/Naïve Bayes** | 0.94 | 0.92 |
| **TF-IDF/Linear SVM** | 0.94 | 0.84 |
| **Word2vec/Naïve Bayes** | 0.70 | 0.71 |
| **Word2vec/RDF SVM** | **0.99** | **0.84** |

# BUILDING THE AUTO-TAGGING CLASSIFIER:

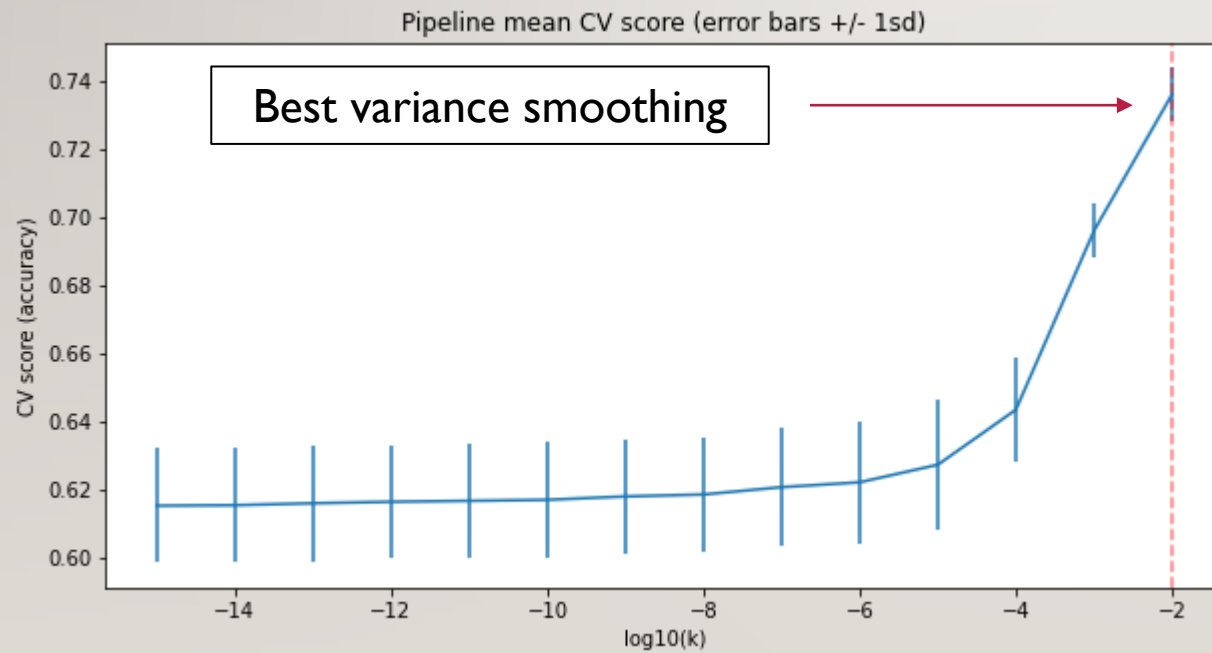## WHAT WORDS WILL MATTER THE MOST WHEN CLASSIFYING?



- 'Project', 'provide', 'help', 'child', 'community', and 'family' appear in all categories
- 'School', 'education', 'woman', and 'girl' appear often in both 'Gender Equality' and 'Education' texts which can lead to some texts being misclassified under these themes
- 'Water' and 'health' are most frequent in 'Physical Health' and 'development'
- 'income' are most frequent in 'Economic Growth' texts

# BUILDING THE AUTO-TAGGING CLASSIFIER:

## THE TF-IDF/NAÏVE BAYES MODEL MISCLASSIFIES "EDUCATION" PROJECTS THE MOST

# BUILDING THE AUTO-TAGGING CLASSIFIER:

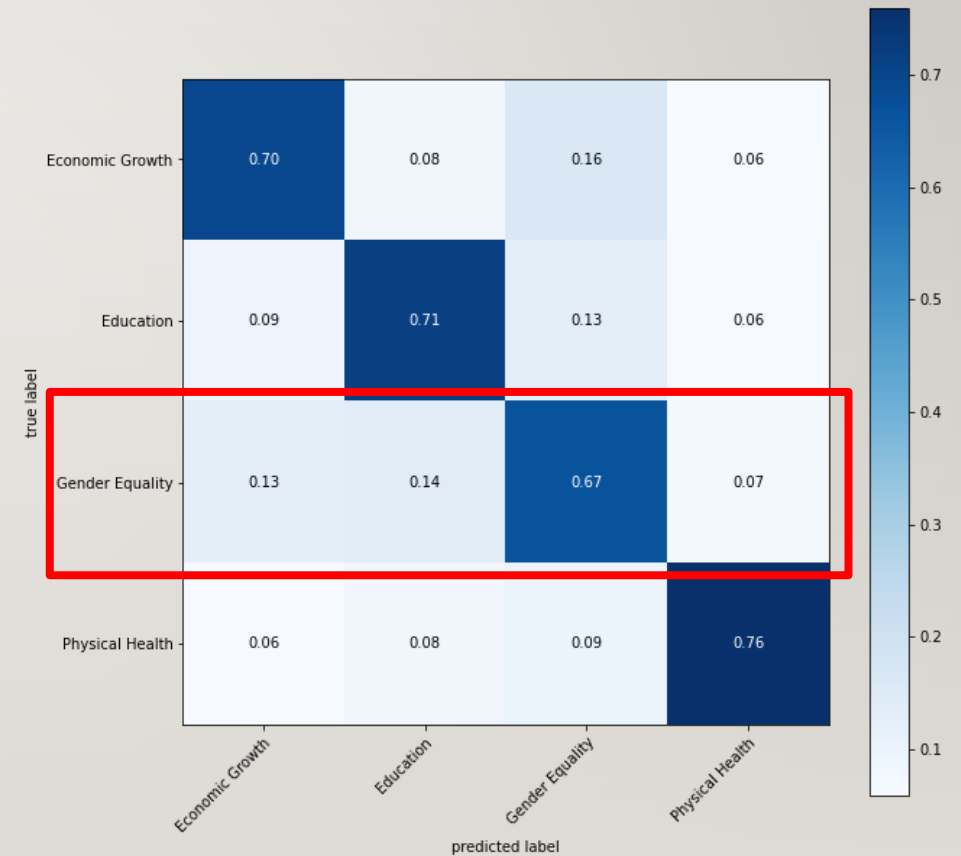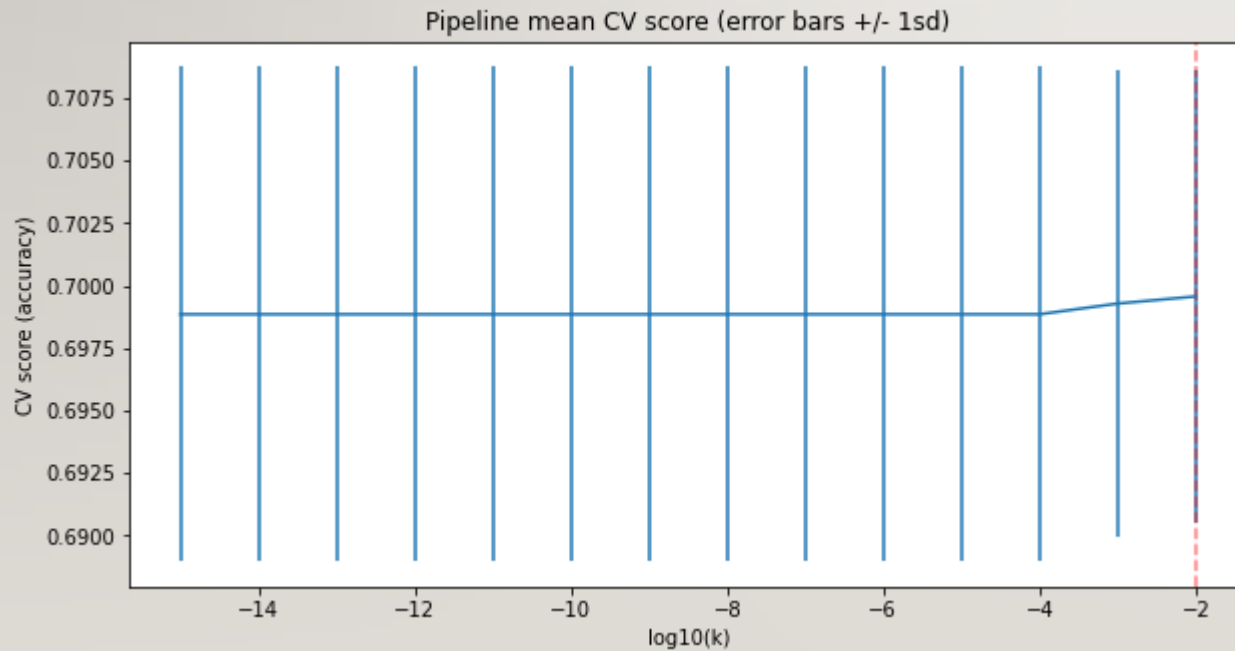## THE TF-IDF/SVM MODEL IS MORE ACCURATE, BUT MUCH SLOWER TO TRAIN

- TF-IDF/Linear SVM model correctly classifies 20% more "Education" projects than the TF-IDF/Naïve Bayes model

# BUILDING THE AUTO-TAGGING CLASSIFIER:

## THE WORD2VEC/NAÏVE BAYES MODEL MISCLASSIFIES "GENDER EQUALITY" PROJECTS THE MOST

Tuning the variance smoothing parameter has no effect on the performance of the model
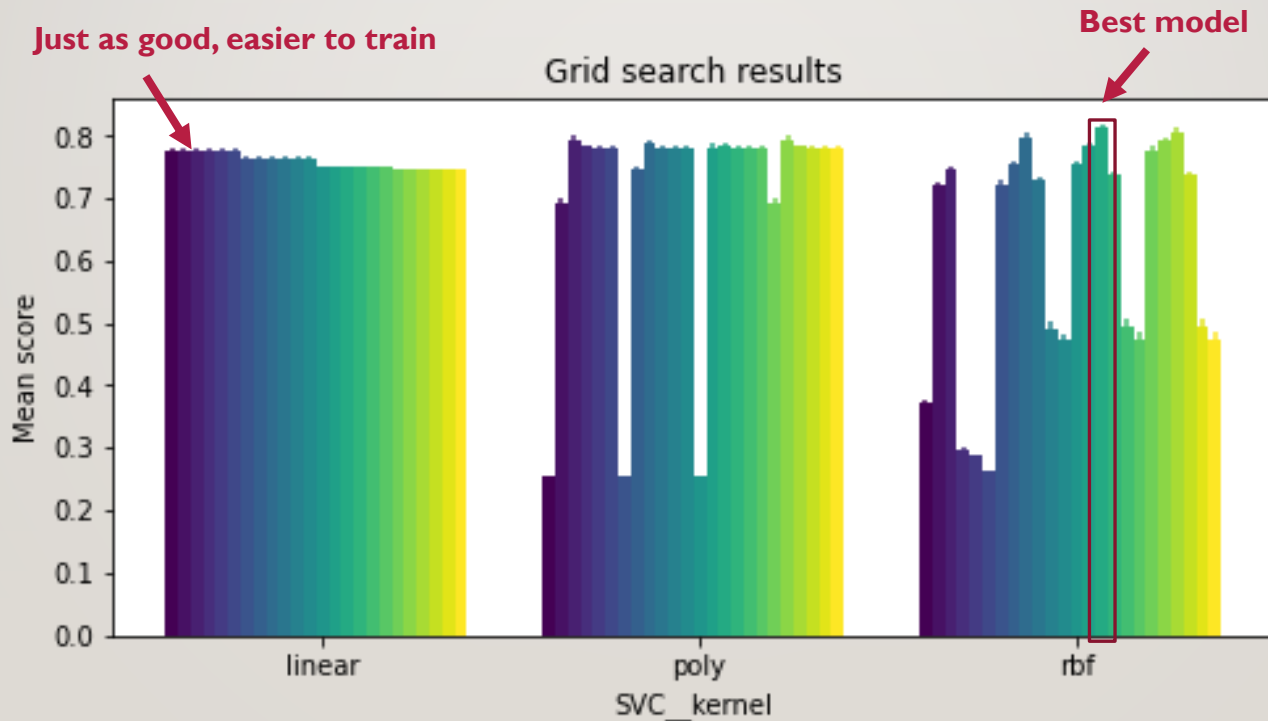
# THE MOST OPTIMAL CLASSIFIER:

## A WORD2VEC/SVM MODEL

The SVM model was tuned using different kernels and different values for the parameters C and gamma
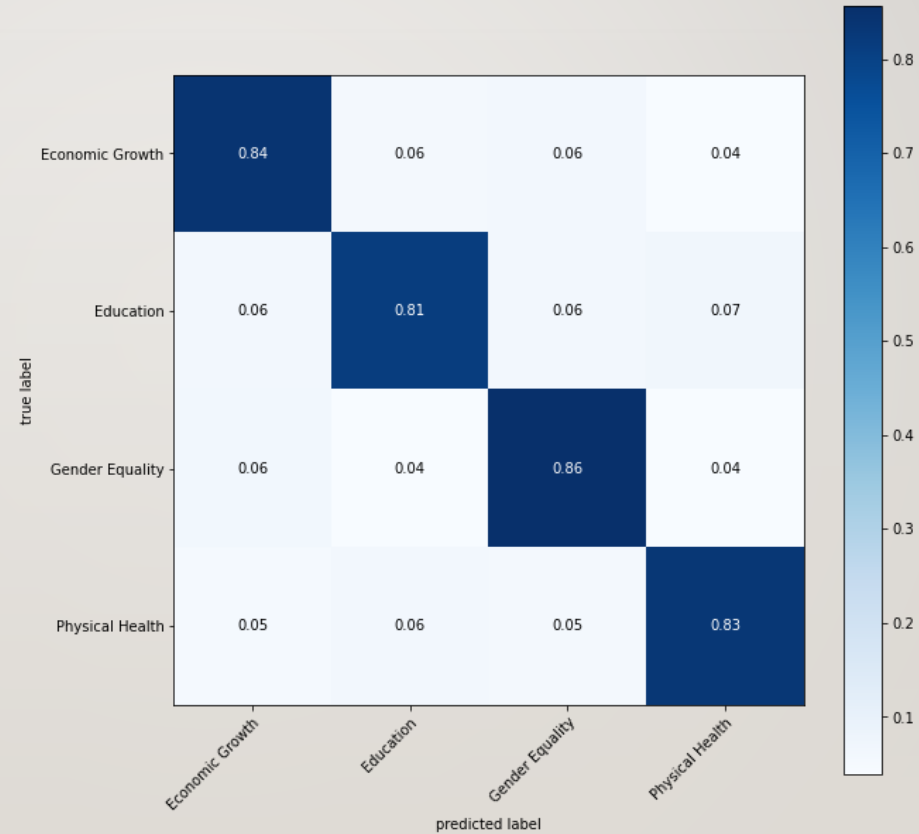


- The best performing kernel is RBF with parameter values C = 10 and gamma = 0.01
- A linear SVM model with C = 0.1 performs well with an accuracy of nearly 0.80

# THE MOST OPTIMAL CLASSIFIER:

## A WORD2VEC/SVM MODEL

- Word2vec/RBF SVM model correctly classifies 19% more "Gender Equality" projects than the word2vec/Naïve Bayes model

# FUTURE IMPROVEMENTS

- Apply the word2vec word-embedding method to the entire dataset of containing all 38,811 projects and fit a radial SVM model

- Depending on how accurate this model is, the challenges of this next step will continue to be distinguishing among projects with similar language and identifying a metric that can assess how correctly the model classifies these projects, but also gives merit to partially correct classification