

## Project Report

### Vision:

The general purpose of the project is to perform Social License to Operate Triple-Bottom-Line topic classification on Twitter data associated with various mining companies. Social License to Operate indicates the ongoing acceptance of a company or industry's standard business practices and operating procedures by its employees, stakeholders, and the general public (Investopedia). Triple Bottom Line is a framework or theory that recommends that companies commit to focus on social and environmental concerns just as they do on profits (Investopedia). We will use supervised machine learning algorithms to perform multi-class single-label classification Tweets to predict whether their topic of discussion corresponds to social, environmental, or economic concerns.

### Background:

Our work is a revival and continuation of the work initially done at the Commonwealth Scientific and Industrial Research Organization (CSIRO) by (insert name here) on TBL topic classification. We are not directly referencing that research but instead basing our initial data pre-processing on the anonymous ACL submission titled "Classifying Stance Using Profile Text". We are however using the exact same labeled training dataset that was used in the prior research for TBL topic classification on SLO for mining companies. Our work will also involve use of the datasets available on Calvin College's Borg Supercomputer and will be uploaded to the Calvin-CS / slo-classifiers GitHub Repository. This project will be a prelude to continued research on topic, stance, and sentiment analysis utilizing machine learning for Social License to Operate of mining companies in connection with Professor VanderLinden's "Machine Learning for Social Media" research project.

As of the current status of this report, we are currently rapid prototyping using Scikit-Learn machine learning classifiers. These classifiers require minimal effort to initially setup with default hyperparameters. They train speedily and provide results in a timely manner, allowing us to adjust our hyper-parameters on-the-fly to see if there are any noticeable differences. It is also quite simple to add additional Classifiers as the Pipeline class allows literal copy/paste of a code template. All that is required is the addition of a new import statement for that Classifier and to replace the name of the old Classifier and its corresponding parameters with the new one. This design feature is one of the reasons we chose to utilize Scikit-Learn; that and it was recommended by Professor VanderLinden as the starting point.

Of note is that Scikit-Learn provides automated parameter tuning via the Grid Search and Random Search classes. Grid search methodically builds and evaluates a model for each combination of algorithm parameters specified in a grid. Random search methodically builds and evaluates a model for each combination of algorithm parameters sampled from a random distribution for a fixed number of

## Project Report

iterations. We plan to utilize one or both of these hyperparameter tuning methods in order to expedite the search for optimal hyperparameters for all of the Scikit-Learn Classifiers we are prototyping with. As we add additional Classifiers to our codebase, it becomes time-saving to automate parameter tuning as much as possible.

Once we have established which classifiers have the most potential to provide favorable metrics, we may migrate towards Keras and Tensorflow for GPU support and more versatility. Scikit-Learn does not provide GPU support for its machine learning algorithms. This does not matter at the moment as we are working with two very small datasets which in total only provide us with 330 samples. That and GPU support will primarily benefit deep neural networks while we are also using non-NN algorithms. However, if we wish to crowdsource TBL classification on significantly larger Twitter datasets and work with those, then GPU support will become necessary. We have heard it requires approximately 24 hours utilizing one Nvidia Geforce Titan on the Borg supercomputer to perform stance analysis training on the larger Twitter datasets consisting of 500k+ examples. It would be expedient to parallelize this process utilizing all 4 Nvidia Geforce Titans on Borg to cut the training time down to a quarter.

We plan to implement metric visualizations via the use of the matplotlib library and SciView in Pycharm. The Scikit-learn online documentation has a section on “Classification of text documents using sparse features” that can hopefully be modified to suit our purposes. Their codebase constructs a bar plot comparing a variety of Classifiers side-by-side visualizing the accuracy score, training time, and test time. As we are also training multiple Classifiers in the hopes of finding a suitable one(s) to further explore in the Keras and Tensorflow API, this type of visualization would be very useful. Individual charts detailing a metric summarization of the micro/macro average, weighted average and associated precision, recall, f1-score, and support values are also planned.

### Implementation:

These sections will describe in detail (perhaps too much detail) our current implementation for SLO TBL topic classification in Python in association with the current state of the codebase. We have decided to keep all debug output statements as “log.debug()” statements that can be shown or hidden by setting “log.basicConfig(level=log.DEBUG)” to the appropriate level.

Our Tweet preprocessor file is separated into 3 individual functions that perform preprocessing specific to the datasets we are utilizing. The first is a Tweet dataset consisting of 229 labeled examples, the second is another Tweet dataset consisting of 31 labeled examples, and the third is a dataset consisting of 658983 unlabeled examples.

The first Tweet dataset we are performing text pre-processing on is the training dataset that consists of 229 Tweet examples. Not all of them are labeled with a TBL topic classification and those are

## Project Report

dropped from consideration. The data is shuffled randomly upon importation to ensure there is no biased structure to the import order. We do so by utilizing Numpy's "random. permutation" function. Then, a Pandas dataframe is constructed to store the dataset. Custom column names are added for clarity of purpose as none originally exist. The "Tweet" column stores the Tweet, "SLO1" stores the first assigned topic label, "SLO2" and "SLO3" do the same.

Pandas provide a "dropna()" method by which we drop all rows without at least 2 non-NaN values. This indicates that the example lacks any TBL classification labels and can be safely discarded. We use Boolean indexing via bitwise operations, the ".notna()" method, to construct a mask by which we isolate those examples with only a single TBL classification. These examples are placed in a new dataframe and afterward, we drop the SLO2 and SLO3 columns as they are obviously just NaN values. This procedure is effective as a preliminary analysis of the CSV file indicates that all labeled examples definitely have a label in the "SLO1" column. Our objective is to construct a dataframe consisting of a column storing the raw Tweet and another column storing a single topic classification. We rename this new dataframe to columns "Tweet" and "SLO".

Next, we construct another mask to isolate all examples with multiple SLO TBL classifications and apply the mask to construct a new dataframe containing only those examples. We then perform a "drop()" operation on the new dataframe to construct 3 separate dataframes. The first from dropping SLO2 and SLO3, second dropping SLO1 and SLO3, and third dropping SLO1 and SLO2. This inefficient but workable solution effectively create duplicates of all examples with multiple SLO TBL classifications with just a single label per example. We then name the columns "Tweet" and "SLO". This is done so that our machine learning model can take into consideration those examples that can be classified as multiple topics.

The multiple separate dataframes constructed from the above operations are then concatenated back together as a single whole Pandas dataframe. Any rows with a NaN value in any column are then dropped via "dropna()" to effectively remove all examples with multiple topic classifications that might have had a topic in SLO2 but not SLO3 or vice versa. Last, we drop all duplicated examples possessing the same TBL classification values in the "SLO" column. We do this as the initial imported dataset sometimes contained duplicate labels for the same example. We surmise this is because multiple people were manually hand-tagging the Tweets and sometimes they were in agreement.

Using the "shape()" method call, our final training dataframe contains a total of 245 Tweets with a single TBL topic classification label. ~~It should be noted that as of our current implementation the second TBL labeled dataset provided by Professor VanderLinden is not currently in use. There are 31 additional Tweets and we plan to include these in the future to help alleviate our issue of a small training and test dataset.~~ We are also using a large Twitter dataset that has already been pre-processed

## Project Report

and tokenized as the set we will make predictions on in order to test the generalization of our model(s) to new data. This set does not contain any target labels and thus we cannot use part of it to supplement our small training and test sets. There are a total of 658983 Tweets included. The CMU Tweet Tagger was used to pre-process the text but unfortunately, this is not a feasible option for us as we are working solely on Windows OS workstation(s).

As we are incapable of using the Linux/Mac only CMU Tweet Tagger for pre-processing, our decision was to manually clean the raw Tweet using Python regular expressions and other libraries. The Natural Language Toolkit was considered as an alternative but ultimately we chose to just use built-in Python libraries and functions. A for loop is used to send each Tweet to a preprocessing function that does the following:

- a) Removes “RT” tags indicating retweets.
- b) Removes URL. (e.g. `https://...`) and replace with `slo_url`.
- c) Removes Tweet mentions (e.g. `@mention`) and replaces with `slo_mention`.
- d) Removes Tweet hashtags (e.g. `#hashtag`) and replaces with `slo_hashtag`.
- e) Removes all punctuation from the Tweet.

We also down-case all text from upper to lower case letters. On our TODO list is to implement regular expressions or other methods in order to:

- a) Shrink character elongations (e.g. “yeees” → “yes”)
- b) Remove non-English tweets
- c) Remove non-company associated Tweets.
- d) Remove year and time.

~~For our current two datasets,~~ the yet-to-be-implemented preprocessing features do not seem to be an issue as the preliminary analysis indicates those elements are not present or have already been considered. We save the processed dataframe to a comma-delimited CSV file to be used in training our Scikit-Learn Classifiers.

The second Tweet dataset we are performing text preprocessing on consists of 31 hand-labeled examples provided by Professor VanderLinden. We follow a similar path as above with our first dataset of 229 labeled examples. We noticed that there was a spelling error present in one of the examples where “environmental” was misspelled to “environmental”, resulting in the erroneous creation of a 4<sup>th</sup> target label later on when we were training our Classifiers. This was corrected manually by editing the original CSV file before re-preprocessing and saving out to a comma-delimited CSV file. Of note, is that

## Project Report

each example only possesses up to two different TBL classifications as opposed to up to three with the first dataset. The Tweet itself was in the same format as the other and thus we could trust that preprocessing, in the same manner, would yield similar processed data.

The third Tweet dataset we are performing text preprocessing on consists of 658983 unlabeled examples with 11 columns of different data including Tweet ID#, language of the Tweet, whether it is a re-Tweet, associated hashtags, associated mining company, Tweet text with mentions, user screen name, user description, Tweet text without mentions (replaced with slo\_mention), and Tweet author profile description. While this dataset has technically been previously preprocessed by earlier research on SLO stance classification, we noticed some discrepancies between these processed Tweets and ours.

- a) The Tweets still had “#” hashtags, whereas we replaced with slo\_hashtag in ours.
- b) The Tweets still had punctuation, whereas we removed them in ours.

Consequently, we decided to run the entire set through our custom preprocessor in order to normalize the Tweets to be consistent with ours. Python’s timer class records that it took approximately 11412.2 seconds to process the entire dataset of 658,983 Tweets. This was done overnight and the results were again saved to a comma-delimited CSV file.

Please refer to “SLO\_TBL\_Tweet\_Preprocessor\_Specialized.py” for the codebase. It has also been included in our “proposal.ipynb” Jupyter Notebook file.

Our “slo\_topic\_classification\_clean.py” program implements Scikit-Learn Classifier training, prediction, and parameter tuning via the Pipeline and GridSearchCV classes. We import our processed datasets, re-index and shuffle the data, and generate a Pandas dataframe for each. We then concatenate the individual datasets together into one cohesive dataframe and again re-index to ensure our range starts from 0. The total number of useable labeled examples is at 277, each with a single TBL topic classification of economic, environmental, or social.

The next step was the input feature set created using the “Tweet” column and a target label set created using the “SLO” column. We chose to refactor our code for this into a separate function so that we can run multiple iterations for training our Classifiers on randomized Tweet and target label test and training sets each iteration. Scikit-Learn included a handy function “train\_test\_split()” which allowed us to easily split our input feature and target labels into a training and test set for each.

~~It is at this point in the code base that we also import a very large Tweet dataset consisting of some 600k+ Tweets that are unlabeled to be used as the input feature for making predictions and seeing how well our model generalizes to new data. These Tweets have already been preprocessed and tokenized by the CMU Tweet Tagger. We simply have to run the entire dataset into a Pandas dataframe,~~

## Project Report

isolate the “tweet\_t” column that contains the Tweet, and use the CountVectorizer and TfidfTransformer class to normalize from categorical to numerical data. The details are similar to what is described above. For the future, we plan to do further post-processing on these Tweets in order to minimize the discrepancies between how we pre-processed and post-processed our training and test datasets and how it was done on this Twitter dataset. Two things we have noticed is that those Tweets still seem to contain hashtag items and some punctuation. These should be removed as we removed both in our training and test sets. The predictive ability of our trained model may otherwise be compromised when using these Tweets.

With the training, test, and generalization set properly prepared, we utilized Scikit-Learn’s Pipeline class in order to set up various Classifiers. Each Classifier is contained in its own module and provided log output is set to “debug” or lower, will display accuracy metrics and a classification report summary. The summary includes statistics on precision, recall, f1-score, as well as the micro, macro, and weighted averages for each. A for loop is used to generate metrics over N iterations and a mean accuracy metric is provided. The trained Classifier is also passed to our “make\_predictions” method afterward which attempts topic classification using the large unlabeled 658,983 Tweet processed dataset. The currently implemented Classifiers include:

- a) Multinomial Naïve Bayes’
- b) Stochastic Gradient Descent (SGD)
- c) Support Vector Machine – Support Vector Classifier.
- d) Support Vector Machine – Linear Support Vector Classifier.
- e) Nearest Neighbor KNeighbors Classifier.
- f) Decision Tree Classifier.
- g) Multi-layer Perceptron Neural Network Classifier.
- h) Logistic Regression Classifier.

These are many of the Classifiers capable of multi-class single-label topic classification. As such, we have decided to implement as many as we can to see which one will be the most performant and worthy of further consideration in the Keras and Tensorflow API, provided those API’s support or can be made to support that Classifier.

For each Scikit-Learn Classifier Pipeline, we implement a CountVectorizer(), TfidfTransformer(), and the relevant Classifier Class(). The following 2 sections describe in some detail the reason we utilize these three classes:

The target label train and test sets were encoded using the Scikit-Learn LabelEncoder class. This converted our categorical labels of “economic”, “environmental”, and “social”, into associated integer values of 0, 1, and 2, respectively. A necessary step as most machine learning algorithms we are interested in prototyping with require and support only numerical data. (Note: this is deprecated – may or may not use in the future)

## Project Report

The Scikit-Learn CountVectorizer class was used to convert the processed Tweet training and test set into feature vectors with binary values of 0 and 1. Documentation indicates that the class converts a collection of text documents to a matrix of token counts and produces a sparse representation of the counts. As we did not provide an a-priori dictionary and analyzer for feature selection, the total number of features is equal to the vocabulary size of the analyzed data. Hence, we have a very high dimensionality in our feature vectors compared to our small number of samples. This effectively creates the bag-of-words that we used to represent our categorical Tweet data. The occurrences of each word are stored in the feature vector. Console output shows that we are dealing with a vocabulary size of 809 in comparison to 164 examples for the training set and 81 examples for the test set (*Note: deprecated numbers, TODO - update for new training set*).

The Scikit-Learn TfidfTransformer class was used to convert the vectorized categorical Tweet data into term-frequency \* inverse document-frequency. The purpose of this is to scale down the impact of tokens that occur very frequently and are therefore empirically less informative than features that occur in a small fraction of the training set. Term frequencies, in general, are better than raw occurrences as larger corpuses will have higher average word occurrence values than smaller corpuses. So, normalization of this kind provides better input feature vectors for training our model.

Each Classifier is also paired with a Grid Search Function utilizing Scikit-Learn's GridSearchCV() class that provides automated parameter tuning. The grid search requires the setup of a classifier (which we did via Pipeline) and the specification of a dictionary storing all the keys (parameters) and values (parameter values) to tune with. The dictionary is passed as an argument to the GridSearchCV() class along with the Classifier Pipeline. We also passed along optional arguments specifying it should run in parallel using all available cores and perform 5-fold cross-validation splitting. This class provides an exhaustive search of all possibilities, meaning it tries all possible combinations of the parameters and associated values you provide it with. Hence, the time to find optimal parameters using our grid searches varied drastically from a few minutes to a few hours.

As mentioned above, we also utilize a large 658,983 Tweet dataset upon which we make predictions using each of our trained Classifiers. The prediction set, so to speak, is prepared in its own function. We drop all columns except the "tweet\_t" column containing the processed Tweet to create an input feature dataframe. Our prediction function is then called by each Classifier's module, passing in the Classifier itself. The Classifier makes predictions on all Tweets and we use counter variables to calculate what percentage of Tweets were classified as economic, social, or environmental among the entire dataset.

Results:

## Project Report

~~As we have just begun initial implementation of our machine learning system, most of these classifiers have been using default hyperparameters and thus our results have been pretty dismal, at best. The highest accuracy metric obtained was almost 56% with the lowest dipping in the 20<sup>th</sup> percentile. It is our plan to use parameter tuning via Grid Search or Random Search to assist in finding hyperparameters that will improve our metrics. As of the moment, the predictive ability of our Scikit-Learn trained models is less useful than flipping a three-way coin for some Classifiers.~~

Grid Search was the essential component to obtaining the best possible results with our limited training and test sets consisting of a total of 277 labeled examples. With default and manual parameter tuning, our accuracy metrics were abysmally low and inconsistent. The inconsistency was due in part to initially not running 1000 iterations and then taking a mean of the accuracy metric to find a consistent percentile. Utilizing the suggested optimal parameters from exhaustive grid search, we were able to raise our accuracy metrics for each Classifier to around 50%. The lowest was the Multi-Layer Perceptron Classifier at 0.490, Stochastic Gradient Descent Classifier at 0.492, and Multinomial Bayes Classifier at 0.493, approximately. The highest was the Support Vector Classification Classifier at 0.535 and the Decision Tree Classifier at 0.532. The remainder fell somewhere in between.

Our prediction results for each Classifier indicates that they will not generalize well to new Tweets. At least, not the processed Tweets we are utilizing. “Social” was the favored classification for our trained models, with the Stochastic Gradient Descent Classifier predicting all the Tweets as social in nature (100%). On the flip side, the Decision Tree Classifier was the most balanced in identifying 40% as social, 54% as environmental, and 6% as economic. The rest of the trained models overwhelmingly predicted almost all Tweets as “social”. This obviously means that we are improperly utilizing machine learning methodologies, almost all the Tweets are actually “social” in nature in that dataset, or we simply do not have enough relevant Twitter data in order to train decent models for TBL topic classification, let alone any deep neural networks. Refer to the code output in the notebook for further details.

Of particular concern to us is performing the proper and necessary pre-processing and post-processing of the Twitter data into useable sparse feature vectors. Regretfully, we will need to obtain the assistance of other researchers with a Linux/Mac workstation and the proper set up in order to use the CMU Tweet Tagger on the labeled TBL datasets. ~~Otherwise, we can only find other alternatives.~~

It is also within our planned schedule to implement matplotlib visualizations of our metric summaries to display the results of training our models and their predictive abilities in generalizing to new data. As of the current writing of this report, this is where we are at in our research efforts. Please refer to the code modules included in this Jupyter Notebook for further details.



## Project Report

Placeholder – discuss comparison with similar works. (not really possible since the similar work was internal at CSIRO and Professor VanderLinden is unsure he can retrieve the relevant materials from years ago; otherwise we are using the prior summer's stance classification research material as a reference for our own work)

## Implications:

Social and ethical implications would be that a machine learning algorithm would be the substitute for the voice of the local population and stakeholders concerning the project. Perhaps the future holds a system where the Social License to Operate could be maintained simply by plugging in a Tweet dataset and if above a certain metric threshold, the company or organization would keep that SLO. There is the danger of the company or organization using a trained model to predict SLO levels and assuming that the results are reliable when reality could be different. These are hypothetical situations that may or may not (probably not) ever occur as we are currently just performing stance, sentiment, and topic classification on Twitter data purely for the sake of research.

## Project Report

## Works Referenced

- 1) "1. Supervised Learning¶." *Scikit*, [scikit-learn.org/stable/supervised\\_learning.html#supervised-learning](https://scikit-learn.org/stable/supervised_learning.html#supervised-learning).
- 2) "A Gentle Introduction to the Bag-of-Words Model." *Machine Learning Mastery*, 12 Mar. 2019, [machinelearningmastery.com/gentle-introduction-bag-words-model/](https://machinelearningmastery.com/gentle-introduction-bag-words-model/).
- 3) "A Gentle Introduction to k-Fold Cross-Validation." *Machine Learning Mastery*, 21 May 2018, [machinelearningmastery.com/k-fold-cross-validation/](https://machinelearningmastery.com/k-fold-cross-validation/).
- 4) "Classification of Text Documents Using Sparse Features¶." *Scikit*, [scikit-learn.org/stable/auto\\_examples/text/plot\\_document\\_classification\\_20newsgroups.html#sphx-gl-auto-examples-text-plot-document-classification-20newsgroups-py](https://scikit-learn.org/stable/auto_examples/text/plot_document_classification_20newsgroups.html#sphx-gl-auto-examples-text-plot-document-classification-20newsgroups-py).
- 5) "Introduction to Machine Learning | Machine Learning Crash Course | Google Developers." *Google*, Google, [developers.google.com/machine-learning/crash-course/ml-intro](https://developers.google.com/machine-learning/crash-course/ml-intro).
- 6) "How to Tune Algorithm Parameters with Scikit-Learn." *Machine Learning Mastery*, 1 Nov. 2018, [machinelearningmastery.com/how-to-tune-algorithm-parameters-with-scikit-learn/](https://machinelearningmastery.com/how-to-tune-algorithm-parameters-with-scikit-learn/).
- 7) Kenton, Will. "How Can There Be Three Bottom Lines?" *Investopedia*, Investopedia, 9 Apr. 2019, [www.investopedia.com/terms/t/triple-bottom-line.asp](https://www.investopedia.com/terms/t/triple-bottom-line.asp).
- 8) Littman, Justin. "Where to Get Twitter Data for Academic Research." *Social Feed Manager*, 14 Sept. 2017, [gwu-libraries.github.io/sfm-ui/posts/2017-09-14-twitter-data](https://gwu-libraries.github.io/sfm-ui/posts/2017-09-14-twitter-data).
- 9) Mohammad, Saif, et al. "SemEval-2016 Task 6: Detecting Stance in Tweets." *Proceedings of the 10th International Workshop on Semantic Evaluation (SemEval-2016)*, 2016, doi:10.18653/v1/s16-1003.
- 10) "Multiclass Classification." *Wikipedia*, Wikimedia Foundation, 18 Apr. 2019, [en.wikipedia.org/wiki/Multiclass\\_classification](https://en.wikipedia.org/wiki/Multiclass_classification).

Project Report

- 11) "Symbolic Reasoning (Symbolic AI) and Machine Learning." *Skymind*, [skymind.ai/wiki/symbolic-reasoning](https://skymind.ai/wiki/symbolic-reasoning).
- 12) Walker, Leslie. "Learn Tweeting Slang: A Twitter Dictionary." *Lifewire*, Lifewire, 8 Nov. 2017, [www.lifewire.com/twitter-slang-and-key-terms-explained-2655399](https://www.lifewire.com/twitter-slang-and-key-terms-explained-2655399).
- 13) "What Is the Social License?" *The Social License To Operate*, [sociallicense.com/definition.html](https://sociallicense.com/definition.html).
- 14) "Working With Text Data¶." *Scikit*, [scikit-learn.org/stable/tutorial/text\\_analytics/working\\_with\\_text\\_data.html](https://scikit-learn.org/stable/tutorial/text_analytics/working_with_text_data.html).