

FINAL PROJECT REPORT



Machine Learning Applications

Bachelor in Data Science and Engineering

May of 2024

Authors:

Ricardo Vázquez Álvarez

Carlos Casanova Romero

TABLE OF CONTENTS

Introduction.....	2
1. Task 1: Text Preprocessing and vectorization.....	3
1.0 Brief description of the data.....	3
1.1 SpaCy and preprocessing.....	4
1.2 Text vectorization.....	5
1.2.1 Bag of Words.....	5
1.2.3 TF-IDF.....	5
1.2.4 GloVe.....	5
1.2.5 Word2Vec.....	6
1.3 Topic Modelling.....	7
2. Task 2: Machine Learning Models.....	9
2.1 Models implemented and methods of implementation.....	9
2.2 Model Performance (Results).....	11
2.3 Results conclusions.....	12
3. Task 3: Description of the implemented dashboard.....	14
4. Conclusions.....	15
5. Acknowledgment of authorship.....	16

Introduction

It is common for people who are seeking employment to be using websites such as LinkedIn, indeed.com, and similar ones to apply to jobs. These job listings typically come with textual information regarding the company position, description of the work, requirements necessary, benefits, and many other characteristics. But because of the increase in Online Recruitment Fraud which leads to severe scams and stealing of personal information, there can be many job listings posing as authentic ones. Making the job market an attractive place for scammers to steal sensitive information on a person. These job listings may attract others through offering high wages, very flexible hours, great benefits and even low experience required.

Hence, before a person can continue to be entrapped in fraud, the purpose of our project is to tackle this very important problem with machine learning and NLP tools to extract information out of what makes a job posting fraudulent and not fraudulent.

Our objective at hand is to classify correctly whether a listing is fraudulent or not using a dataset consisting in text features and categorical ones from job listings.

Primarily, this will consist of converting our text data in different forms of NLP representations to then use classification techniques to output a binary representation (1 if it is fraudulent, 0 otherwise) of a job posting.

1. Task 1: Text Preprocessing and vectorization

1.0 Brief description of the data.

The data obtained represents different characteristics of job postings. Similar to job postings on LinkedIn, indeed.com and other websites, this dataset includes just that information. The data was downloaded from a Kaggle post created by user Shivam Bansal called [Real / Fake Job Posting Prediction](#). This data set was taken from the [University of the Aegean | Laboratory of Information & Communication Systems Security](#) and was involved in a research paper found in their website ([link to pdf](#)).

The following table is a review of the features being used in the project.

Column Name	Type of data	Description	Percentage of NAs	Importance
Title	Textual	Title of job post.	0%	NLP Feature
Location	Textual	Location of the job.	0%	NLP Feature
Company profile	Textual	Information about the company.	19%	NLP Feature
Description	Textual	Description of the job.	0%	NLP Feature
Requirements	Textual	Textual description of the requirements needed.	15%	NLP Feature
Benefits	Textual	Textual description of the benefits included as a part of the job.	40%	NLP Feature
Telecommuting	Binary	1 if one can work from home, 0 otherwise.	0%	Feature
Required experience	Categorical	Different categories corresponding to the type of experience needed.	39%	Feature
Required education	Categorical	Different levels of education required.	45%	Feature
Fraudulent	Binary	If the posting is fraudulent (1) or not (0).	0%	Target Class

Table 1: Summary of the data set used along with the type of data, description, percentage of NAs and Importance in our task.

1.1 SpaCy and preprocessing

For the preprocessing of our data, what was firstly done was to combine all the NLP features into one column. This meant making a new column called **'text'** that adds the columns title, location, company profile, description, requirements, and benefits sequentially using a space in between them. This was ultimately done in order to have all the textual information in one column so that this can all be processed later using SpaCy. If the row had an NA for these textual columns then an empty string was added (attributed nothing).

Next, the required experience and required education features (nominal) were all given numbers corresponding to their categories. Both columns included NA values that were given by just spaces. These spaces were replaced with 'Not Applicable' (for required experience) and 'Unspecified' (for required education), leaving no NA values.

The text column was preprocessed using SpaCy. Initially, two primary issues were addressed with RegEx:

- Words that were joined together, such as 'awesomeLove', were separated into 'awesome love'.
- Job descriptions that contained phrases like 'classic!)Fashion' were corrected to 'classic fashion' by removing special characters between words and separating them.

Following text preprocessing, we proceeded with the tokenization process using SpaCy, where each word was tokenized to remove punctuation, spaces, and stopwords. Words were then lemmatized, converted to lowercase, and filtered to exclude words shorter than 2 characters or longer than 20 characters (removing elements like hashes). All these steps were conducted using the `en_core_web_sm` model as the NLP engine.

As an additional preprocessing step, we created bi-grams using the Phrases function from Gensim, with a threshold of 30 and a minimum count of 50, focusing on very common words appearing together in the corpus. Analysis of the bi-grams showed the most frequent ones: 'fast_pace', 'social_medium', 'new_york', 'ideal_candidate', and 'long_term', with a count higher than 2000. Significantly enhancing our NLP modeling.

The final output was the processed tokens, which are prepared to use in the NLP tasks.

Something to note in the data set is that there is a large disbalance between the Non Fraudulent rows and the Fraudulent rows. Resampling the data was also done using the SMOTE (Synthetic Minority Over-sampling Technique) algorithm found in [2] and also exemplified in [1] which allows one to generate synthetic samples of the minority class by using the *k* nearest neighbors and interpolating between the original sample and those neighbors. This was done after applying the text vectorization which will be shown in the following section. Something to take into consideration is that with SMOTE, this resampled data may lead to overfitting or may not make sense for NLP since we may be resampling text vectorizations that do not represent words at all.

1.2 Text vectorization

To effectively analyze our data, we first need to transform the raw text into a structured numerical format. By using Gensim's 'corpora.Dictionary' class, creating a mapping of word tokens to unique integer IDs, opening up to define models.

1.2.1 Bag of Words

We employed the Bag of Words (BoW) vectorization method to transform textual data into a numerical format.

First by using the Gensim's dictionary facilitating the conversion of text documents into BoW. Where we can also extract the vocabulary from the dictionary to examine the unique tokens in our dataset.

Secondly, utilizing the CountVectorizer from scikit-learn to implement the model, we prepared the input by ensuring it consists of a single string of concatenated tokens per document. The resultant bow_matrix represents the frequency of each word within the documents. Using sparse matrices, which store only non-zero elements, is highly efficient with memory, and therefore we use 'hstack' to integrate the other variables for classification into a single matrix.

This vectorization has prepared our dataset for the classification task, providing a strong baseline.

1.2.3 TF-IDF

Term Frequency-Inverse Document Frequency or TF-IDF is used to transform text data into numerical format. This method not only considers the frequency of terms within a document like BoW but also how unique these terms are across all documents in the corpus. With 'TfidfVectorizer' from scikit-learn it converted the tokenized string data into a matrix where each row represents a document and each column a term, where the entries of the matrix are weighted values indicating the importance of a term in a document relative to its commonness across all documents. Unlike BoW, TF-IDF down-weights terms that appear frequently across documents, emphasizing terms that are unique to certain documents. Additionally, to combine with the other variables we use the 'hstack' function.

Providing a way to vectorize text for NLP tasks, it adjusts the importance of words according to their relevance and uniqueness, with different capabilities for the classification models.

1.2.4 GloVe

Global Vectors for Word Representation, or GloVe, is a word embedding that captures not only the frequency of terms but also the contexts in which they appear. We utilized the 'glove-wiki-gigaword-50' pre-trained GloVe, which has 50-dimensional word vectors trained on Wikipedia and Gigaword data. Offering a good balance in vocabulary coverage, achieving 100% coverage with our Gensim dictionary. This ensures that no information is lost due to out-of-vocabulary terms, which is beneficial. Capturing the semantic relationships between

words based on their co-occurrence in large corpus, giving a deeper understanding of word meanings.

Using GloVe embeddings enhances our text processing capabilities by providing context-aware input for our models. However, GloVe embeddings might not be optimal for our project objectives. For instance, while GloVe offers extensive vocabulary coverage and powerful insights, it may not directly align with our goals. Additionally, the static nature of GloVe embeddings means they do not adjust based on the context in a sentence, which could be a limitation.

1.2.5 Word2Vec

Word2Vec is a prediction-based method used for generating word embeddings that capture both the context and semantic similarity of words. Unlike one-hot encoding that represents words as independent entities, Word2Vec provides a dense vector representation where semantically similar words are mapped to proximate points in the embedding space. In our project, we used the Word2Vec module from Gensim, specifically using the Skip-Gram algorithm. The Skip-Gram model is particularly effective with larger datasets and better at capturing infrequent words than its counterpart CBOW. By predicting the context given a target word, Skip-Gram deeply analyzes the relational aspects of the data.

For efficient memory management and ease of use, we employed the KeyedVectors module from Gensim. This module facilitates the storage and manipulation of word vectors separately from the full model. KeyedVectors objects are more compact and can be loaded and saved in the native Word2Vec format [3]. This is crucial for frequent access to word vectors without the overhead of loading the entire model.

Despite the high-dimensional representation, our analysis revealed a vocabulary coverage of only 34.4% when matched against our Gensim dictionary. This suggests that while the model effectively captures semantic relationships for a subset of terms, a significant portion of our corpus's vocabulary remains inadequately represented. This could stem from several factors, including the specificity of domain language in our dataset or limitations in the training corpus size and diversity.

To better understand the semantic relationships captured by Word2Vec, we employed t-SNE for visualizing the high-dimensional word vectors in a 2D space. Providing a graphical representation of how words are clustered based on their similarities in the embedding space. For example, technology-related terms such as "internet", "web", "framework" are closely grouped, or "compensation" "salary" "plus" "pay" "bonus" as well, indicating effective learning of contextual relationships by our model. Showing the strengths of Word2Vec skip-gram modeling.

1.3 Topic Modelling

In this section, we explore the hidden themes or topics in our dataset by applying Latent Dirichlet Allocation (LDA). This method extracts topic distributions from the text, facilitating an understanding of the dominant themes present.

We use gensim library's LDA model to conduct our analysis, which provided us with interpretable results. Initially, we instantiated the LDA model using our Gensim corpus dictionary and set the number of topics to 20 as a benchmark. Each document in our dataset was transformed into a bag-of-words format. This is because it simplifies text data focusing on the frequency of words into numerical counts that LDA uses to determine topic distributions.

To visualize the initial results, we noted that the topics were not optimally distinct; several themes overlapped significantly. However, we were able to identify clusters such as customer service, technical job offers, teamwork, and management oriented.

To refine our model, we utilized the coherence scores for models trained with different numbers of topics. Coherence scores indicate the semantic similarity between high-scoring words within each topic, where higher scores suggest more interpretable topics. It is important to note that while coherence is a useful measure, it can be somewhat subjective. In some cases, high coherence scores may not correspond to topics that are useful for our specific task. Nonetheless, coherence scores are a valuable starting point for determining a range of topics.

The evaluation indicated that coherence was highest between 10 and 17 topics. Based on the metadata and characteristics of our dataset, we determined that a smaller number of topics would be more beneficial, ensuring that each topic is distinct and relevant to our analysis.

After setting the number of topics to 10, we obtained informative topics characterizing our data:

- Abilities: knowledge, maintenance, production, design, engineering performance, and skill.
- Team Player: Client, team experiences, opportunities, support environment, and care.
- Job Conditions: Salary, employment terms, involvement, position, benefits, vacation_hour.
- Qualifications: Career, training, quality, role.
- Customer Service Jobs: Customer, service, client, management.

These topics show an effective way of retrieving and analyzing information within our texts.

Using pyLDAvis, we visualized the distribution of topics, which allowed us to clearly observe the major themes and their distinctions.

We defined a function to examine the topics more closely by analyzing the most relevant documents, for example with 'Team Player' included descriptions from the same company

but different job positions and locations. Common phrases like "share the same values", "prioritize team's needs above your own", and "We hire for culture" show the emphasis on teamwork and cultural fit. Validating the accuracy of our topic modeling , revealing essential patterns and insights.

The LDA approach was effective in understanding the key themes in our dataset. With well-defined topics, the gensim LDA model facilitated the exploration of our text data, to make informed decisions.

2. Task 2: Machine Learning Models

2.1 Models implemented and methods of implementation

Given that our dataset includes the binary class variable (fraudulent), the objective at hand is to classify between fraudulent (1) and non-fraudulent (0) job postings. This was done using various techniques in machine learning, including Supervised Learning and indirect classification through Unsupervised Learning (clustering).

To utilize the word embedding models (Word2Vec and GloVe) we defined a function “get_review_vector” that takes the pre-trained model and a document as inputs. For each word in the document that appears in the vocabulary of the model, its vector is added to a cumulative sum, later on divided by the count of known words to produce an average vector representing the entire document, capturing the overall semantic essence of the document based on the words it contains. This enables us to utilize the embedding and to concatenate it with our variables to use them for classification.

The input to our machine learning models were the vector representations of the dataset from Bag of Words, TF-IDF, Word2Vec, and GloVe. Where we tested each of the models to be explained against the above representations in order to compare which resulted in better results. We also compared the performance of each model with the data set including the other features included in Table 1 (categorical variables). These categorical variables were changed to numerical factors using LabelEncoder() as explained previously.

Because our data was highly imbalanced, the performance of the models could not simply be explained by the accuracy of the predicted values. Therefore, to correctly determine how well a model is doing in classifying the data, they were evaluated according to the **recall** and **macro F1-score** which shows the clear proportion of the True negatives, False positives and False negatives. Indicating exactly how well each class is being represented by the models used. The following equations show the metrics of recall, F1-Score and Macro F1-Score where TP is for True positives, FN is for False negatives and FP is for False positives.

$$\text{Recall} = \frac{TP}{TP+FN}$$
$$\text{F1-Score} = \frac{TP}{TP + 0.5 \cdot (FP + FN)}$$
$$\text{Macro F1-Score} = \frac{1}{n} \sum_{i=1}^n F1\ Score_i$$

Supervised Learning:

For supervised learning, the models implemented were Logistic Regression, Random Forest Classifier, and Support Vector Machine. Logistic regression aims to map a logit function to the data set using the columns and output you give. In the end, the logit function gives a probability of a data point being in a class, if the probability is greater than 0.5 then it belongs to class 1 (fraudulent), else it belongs to class 0. The Random Forest model for classification

works by creating decision trees on random subsets of the model to then combine the outputs of the decision trees to create a robust classifier. Support vector machines work by projecting the data onto a kernel space that induces a nonlinear mapping. It then tries to maximize the margin between the support vectors.

Hyperparameter Tuning:

The hyperparameters were found using a Pipeline of the model itself along with a Grid Search Cross Validation scheme to find the best parameters that resulted in the best performance. This was done for all vector representations. The following are the parameters that were considered in the Cross validation method across all the vector representations. Three folds were used for validating.

Random Forest:

- **n_estimators** : [100, 200, and 300]
- **max_depth**: [None, 10, and 20]

Logistic Regression:

- **C**: [0.1, 1, and 10]
- **penalty** = ['L1', 'L2']

SVM:

- **RBF kernel** only.
- **'C'**: [0.1, 1, and 10].
- **'gamma'**: ['scale', 'auto'].

Unsupervised Learning:

Although this is a classification task, which is done with supervised learning techniques, unsupervised learning may also provide indirect classification to our data. For this reason the following unsupervised learning methods were used. For novelty detection (outlier detection) a One class SVM was trained on the non fraudulent data and then tested on the fraudulent data. This way we can detect the fraudulent data as an outlier. PCA was also applied in order to then project the data onto a lower dimensional space. This was done to obtain more insights on the data while reducing dimension and maintaining most of the variance explained by the data. Although for PCA, the data should be scaled to provide better results, since our input to the PCA model are word embeddings we should not scale since this would provide very distinct interpretations of the data. Gaussian mixture models were also computed to try and cluster between two classes in the PCA projected data. The results of these models will be shown in the next section. In this case, the hyperparameters were not chosen using the Grid Search Cross Validation algorithm but rather the parameters chosen were according to their use. Because we have that the fraudulent class and non fraudulent classes are very similar, it was apparent that a small nu value of 0.1 should be chosen as well as a kernel that can consider complex relationships such as the RBF kernel. For PCA, only three components were considered for the sake of simplicity since we were only going to visualize the first two components. Finally, for the Gaussian Mixture Model, the hyperparameter for **'n_components'** was chosen as 2 since we wanted to be able to classify between our classes which are binary.

Class imbalance problems:

Given the class imbalance problem, the SMOTE resampling technique explained in section 1 was used prioritizing the minority class (fraudulent rows) so as to obtain a more balanced training set. Also including the extra features in the resampling. To avoid flawed conclusions, the data set was first split using `train_test_split`, and the training set was used as input to the SMOTE algorithm. Then this resampled data was used as a part of the training to our machine learning models. The testing data for the machine learning models were then the testing set obtained separately from the SMOTE algorithm, which was a part of the original data set. This was only applied to the Bag of Words representation in an attempt to improve accuracy. In the end, it did not make much sense to use this technique for training a model since what is being resampled is really just numbers and these numbers do not actually represent coherent text. Not only that but unsupervised learning techniques such as PCA and One-Class SVM were not good techniques to use in the SMOTE algorithm since we want to find a representation of our realistic data not the resampled data that is not so indicative of the text of job listings. Nonetheless, the PCA technique was used for the SMOTE data to see how well it performed in obtaining the principal components that maximize the variance. The performance of the SMOTE output data is not shown for the unsupervised techniques though only, the PCA graphs are shown. Finally, the same optimal hyperparameters found for the Bag of Words inputted models were used in the SMOTE models.

2.2 Model Performance (Results)

The model performance is shown with two types of testing, with just the NLP representations (**green columns**) and also with additional features ('telecommuting', 'required_education', 'required_experience') (**blue columns**). Only the results using the best hyperparameters and best resulting models are shown in this section. Best performing models are highlighted in light green.

Model	C0, recall		C1 recall		Macro-F1 Score	
<i>Logistic Regression (BOW)</i>	1.00	1.00	0.71	0.72	0.90	0.90
<i>SVM (BOW)</i>	1.00	1.00	0.71	0.72	0.90	0.90
<i>SVM (BOW + SMOTE)</i>	0.99	0.99	0.73	0.73	0.86	0.86
Random Forest (TFIDF)	1.00	1.00	0.58	0.59	0.86	0.87
Random Forest (W2V)	1.00	1.00	0.42	0.44	0.79	0.80
Random Forest (GloVe)	1.00	1.00	0.40	0.38	0.78	0.77

Table 2: Results of the best machine learning models using all representations. Macro F1-score and recall for each class.

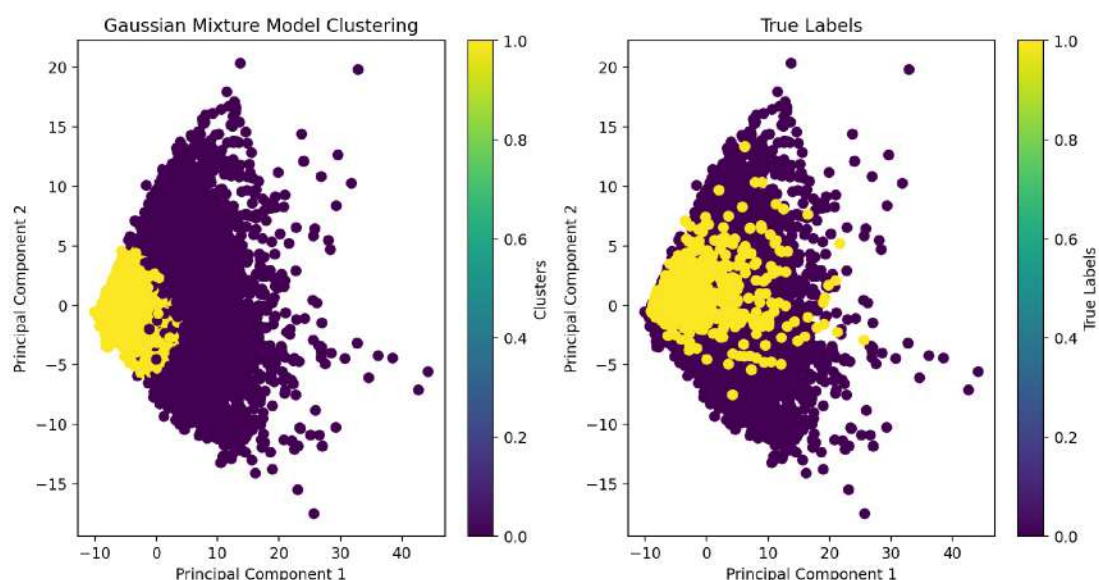
2.3 Results conclusions

Out of all of these methods, the Logistic Regression model and SVM presented the best and nearly identical performance using the Bag of Words representation. With almost perfect results for the non fraudulent case, then 0.71 recall and 0.90 macro- F1 score for the fraudulent class. The better model would be in this case the logistic regression model since it provides the most interpretability on the coefficients of the features (vocabulary of our documents) and also being the most simple model (not computationally expensive and easy to understand).

The best performance found in Unsupervised Learning was the PCA technique followed by a Gaussian Mixture Model. The GMM model gave high recall scores for the fraud class but worsened in every aspect for the other scores. Showing that there was a large amount of correctly classified observations in the class one but too many false positives and negatives (hence the low macro-f1 score). Nonetheless, these techniques did not provide such good results for classification but indeed did present good insights as to how our data's variance is distributed across multiple components when comparing the different representations.

When looking at the PCA outputted data for the different representations it was clear to see which vector representation provided the better conclusions about our data. Using PCA with TF-IDF and Word2Vec it was clear that no information related to the separation of our classes could be found, in general the classes were mixed well with each other and no conclusions can be made about the results. Meanwhile adding new features to the data actually created more noise in the PCA results making it impossible to use the clustering task for classification.

For the Bag of Words representation PCA did help in better showing the distinction between classes and after adding the new features, a new cluster appeared in the results (therefore it did not provide more noise to the algorithm). Either way, for indirect classification using clustering it became difficult to find distinction between the classes. Nonetheless, having class imbalance was the bottleneck for the classification and clustering tasks that were used since there was insufficient information about the fraudulent class.



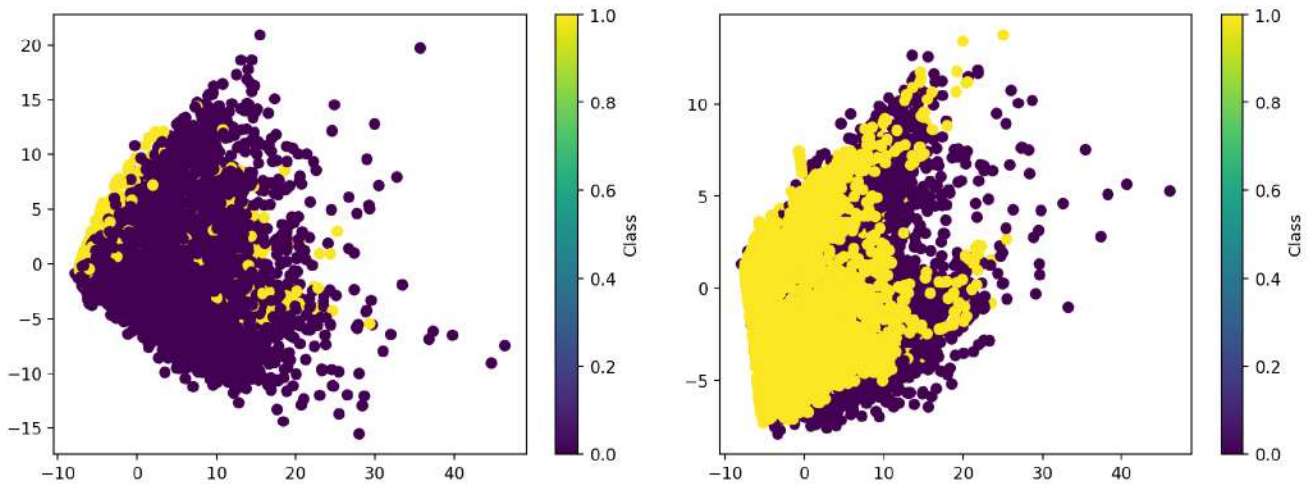


Figure 1: PCA representation of the Bag of Words data using original data set (above) and SMOTE resampled data set (below). Gaussian Mixture model results are shown on the left, meanwhile the true labels on the PCA are shown on the right.

When using SMOTE on the Bag of Words data and PCA, it was clear that all SMOTE really did was provide more overlapping between the classes. This was expected and also one of the reasons why SMOTE was not ultimately used for classification since the resampled data does not actually help in obtaining more insight on the fraudulent class. This SMOTE data was really just numbers being resampled, which did not pertain to real words from the documents. This can be seen by the comparison of both the PCA applied to SMOTE and also PCA applied to Bag of words in the above Figure 1.

In general the Bag of Words representation was the best representation to use for classification. The TF-IDF representation provided worse results meanwhile not providing any additional information on the data set. Finally, for Word2Vec, the classification results were by far the worst in all aspects. This was better used to provide better insights on the text data finding most similar words to a word of interest or using T-SNE to find a low dimensional representation of the text column instead of the classification task.

Ultimately, the extra features being used in the models did not provide a notable change in results and in most cases ended up worsening the performance. This shows that it mainly provided more noise than anything as shown by the PCA representations.

To conclude, the SMOTE resampling technique although in the logistic regression and SVM model improved on the Class 1 accuracy, the F1-score decreased meaning there is a relatively high amount of False positives when compared to the original dataset. Hence, showing that the resampled dataset provided a model that was more sensitive.

3. Task 3: Description of the implemented dashboard

The implemented dashboard has three main windows which the user can use to navigate through the different visualizations implemented.

The first window, "Word2Vec Visualization", shows a scatter plot of the t-SNE for visualizing the high-dimensional word vectors in a 2D space. It shows how words are clustered based on their similarities in the embedding space. The user can click a point/word from the scatter plot in order to display under the graph the 5 most similar words to the one that the user has selected.

The second window, "Document Topic Distribution Viewer", provides mainly a pie chart with the top 5 predominant topics distribution in a document. Here, the user has a lot of freedom, since at the top of the window there is a search bar, where the user can choose one or more documents among all the documents of the corpus. More than one document can be chosen in order to see the predominant topics in more than one document. In addition, the documents can be filtered in 'All Documents', 'Fraudulent' and 'Not Fraudulent', so that when the user chooses the documents to be displayed in the graph, only the documents belonging to the class selected in the filtering will be displayed. Also, under the graph, the documents are shown, so the user can see the whole document which topic distribution is being plotted.

Finally, the third window "Model Performance Dashboard", shows the performance of the different models obtained during the second task of this project. As with the previous windows, the user can use a search bar in order to filter the model results that wants to visualize. In this case, there are two search bars: one to choose among the different vectorization techniques; and another one to choose the elaborated models of that vectorization technique. The result is a vertical bar chart with the performance of the different models following three measurements: recall of not-fraudulent, recall of fraudulent and macro f1 score.

4. Conclusions

This project's purpose is to tackle a very known and increasingly dangerous problem in today's world. Not only is Online Recruitment Fraud (ORF) consistently growing but the difficulty in recognizing whether something is truly fraudulent or not is worsening. Most importantly, through the use of job websites such as linkedin and indeed, stealing of personal information can be easily done by attracting unemployed people to a job posting that aims to provide exactly what they are looking for with their pleasing descriptions and benefits. Masking themselves as real job postings, using similar sentence structures and words to real postings, it is incredibly important to use advanced measures to protect others from this danger. That is why the purpose of this project was to use machine learning techniques to provide a classification of the listings so as to avoid this problem.

Ultimately, this project was done by converting the textual data of job postings (both real and fake) into different word embeddings using text vectorization techniques such as Bag of Words, TF-IDF, Word2Vec and GloVe as well as exploring hidden themes and topics using Topic modeling. These text representations were then used into our machine learning models for use of classification and clustering techniques. The classification models used were logistic regression, SVMs, and Random Forest Classifier. Meanwhile for clustering, One class SVM, PCA, and GMM were used to provide better insights as to how our data is represented. In the end, after tuning the best models used were found by the best macro f1 score from the Logistic Regression Model and also the SVM achieving a 0.90 macro f1 score. Also, extra features were tested but ended up providing more noise or complexity to the model, not improving the macro f1 score.

SMOTE resampling was also tested with the various models, but it was unfortunately not appropriate for the use of classification since it did not resample similar words but words that were not coherent or noisy.

Although this project suffered from class imbalancing, using Logistic regression and Support Vector Machines we were still able to provide good enough results for determining whether a job posting is fraudulent or not.

Finally, a dashboard was implemented, where three main windows were included. Firstly, a scatter plot of the t-SNE for visualizing the high-dimensional word vectors in a 2D space, providing clusters based on the similarities in the embedding space. Secondly, a pie chart with the top 5 predominant topics distribution in a document. Thirdly, the performance of the different models obtained during the second task of this project for better visualization and comparison (with bar plots).

In conclusion, this project has made use of the NLP techniques for text representation and also machine learning tools to tackle the ORF problem existing in our current world.

5. Acknowledgment of authorship

- [1] J. Brownlee, "Smote for imbalanced classification with python," MachineLearningMastery.com, <https://machinelearningmastery.com/smote-oversampling-for-imbalanced-classification/> (accessed Apr. 28, 2024).
- [2] "Smote Documentation," SMOTE - Version 0.12.2, https://imbalanced-learn.org/stable/references/generated/imblearn.over_sampling.SMOTE.html (accessed Apr. 28, 2024).
- [3] "Store and query word vectors", GENSIM, [models.keyedvectors – Store and query word vectors — gensim \(radimrehurek.com\)](https://radimrehurek.com/gensim/models/keyedvectors.html) (accessed May 03, 2024).