# Contents

# Predicting if a Business in America will Survive COVID-19

AUTHORS: Mu Cai(mcai44), Varun Sreenivasan(vsreenivasan), Niharika Tomar(ntomar)@wisc.edu

## GO GREEN. AVOID PRINTING, OR PRINT 2-SIDED OR MULTIPAGE.

The incidence of COVID-19 this year has posed a significant challenge to businesses across various industries. Some of them are thriving whilst some are closing. Thus, more than ever before, it is important for business owners, customers, investors, and many others to know which businesses will survive. In this paper, we first construct a comprehensive dataset, including both normal business and COVID specific information, by utilizing the Yelp dataset in combination with some other external sources. Then we present various classifiers that can be used to predict whether an American business will survive the pandemic, and analyze their performance by utilizing the training and evaluation pipeline. SMOTE is used in the pipeline while training to address the issue of class imbalance in the dataset. Multiple score metrics are used to assess model performance over both classes (open and closed businesses). Our test set results show that model performance varies depending on the metric. In this paper, we also determine the features that are most influential in determining whether a business survives or not. The analysis of feature significance tells us that certain COVID-19 features of our dataset such as Grubhub, COVID Banner, and Call_To_Action_enabled play a vital role in determining if a business remains open.

## 1 Introduction

The success of a business is dependent on the axiom that it stays open. In this project, we develop several machine learning models to predict and analyze whether a business remains open or permanently closes during times of COVID-19.

We conduct research on a dataset containing records of 153,843 American businesses shown in Yelp comprising of multiple segments including COVID-19 features. The records comprise of information such as location, reviews, number of stars, business attributes, business categories, and COVID-19 features like delivery or takeout, virtual services offered, COVID Banner, etc. We then employ the technique of feature engineering to apply appropriate transformations to create a parsed dataset that can be utilized for training and evaluation. In this process, we also add data from external sources namely median income and population for the postcode the corresponding business is located in.

It is important to note that COVID-19 also indirectly affects businesses through some of the other features not cited as COVID features. For instance, the postcode median income feature. With COVID-19, increased job losses will very likely translate to lower median incomes in a locality and if median income relative to postcode is an important feature, it will have a significant impact on business survival odds.

We devise a pipeline to effectively upsample (SMOTE-NC) the minority class (closed businesses, which takes up around 20% of original dataset), train, validate, and test the performance of our classifiers, including K-Nearest Neighbors, Logistic Regression, Random Forest, and Neural Network. The results for Naive Bayes and Support Vector Machine are shown in the Appendix.

We assess feature significance using the Permutation Importance algorithm provided by `eli5` library to determine the 10 most important features in this classification task.

We believe that this paper elevates general understanding of the impact of COVID-19 on American businesses and the economy as whole. The results of our paper highlight the effect that this pandemic will have on businesses across the country. Moreover, this paper can serve as insight to people making business decisions in the time of COVID-19.

## 2 Related/Similar work

### 2.1 COVID-19 Impact on Businesses

The outbreak of the Coronavirus disease 2019 (COVID-19) has had significant ramifications for businesses of all sizes. Businesses are required to navigate through the financial and operational challenges or else face the prospect of imminent closure. For instance, an Accenture study [1] details how businesses are responding to this unique challenge. Another study [2] conducted a survey on 5,800 small businesses to study the effect of COVID.

### 2.2 Machine Learning based Classification Models

The problem of whether a business will survive is a typical classification problem. Various machine learning models have been proposed to address the general classification problem, such as K-Nearest Neighbors [3, 4], Naive Bayes [5, 6, 7], Logistic Regression [8, 9, 10], Decision Tree [11, 12, 13], and Random Forest [14, 15, 16]. Among them, Neural Networks [17, 18, 19] are recently proposed and demonstrated to be a powerful model in terms of the prediction accuracy. However, it is known to lack interpretability, *i.e.*. a black box [20], which means that it is hard to get the confidence interval or the lower bound for the final results. Although some works [21, 22] proposed to interpret the inner mechanism of neural networks, the uncertainty estimation of deep neural networks is still an open problem.

## 3 Dataset

We make use of the Yelp Dataset[1] to solve our question. We draw features from the following segments of the dataset:

`Business.json`: Data including label (`is_Open`), location, number of stars, review count, business attributes, and business categories.

`Covid_Features.json`: COVID data relevant to a business such as Grubhub_enabled, virtual services offered, COVID Banner, etc.

We make use of the `business_id` feature in the `Business` and the `Covid_Features` datasets to perform essentially a join operation to obtain joint unique records based on business id.

In this process, we also perform feature engineering on raw datasets to obtain a final parsed dataset comprising of appropriate features that can be used to develop our classifiers.

---

[1] https://www.yelp.com/dataset

## 3.1 Feature Engineering

- **COVID-19 Features**: We directly use the features' highlights, delivery or takeout, Grubhub enabled, Call To Action enabled, Request a Quote Enabled, COVID Banner, Temporarily Closed, and Virtual Services Offered from the COVID features dataset. These features are represented as binary variables.

- **Location**: The features latitude and longitude are used to represent a business location. These two features are represented as continuous variables.

- **Stars**: It refers to the rating of a business. Since float values are accepted, this feature is represented as a continuous variable.

- **Review Count**: It denotes the number of Yelp reviews for each business. This feature is represented as a continuous variable.

- **Chain**: We use a simple Natural Language Processing trick to determine whether a business (uniquely identified by its id) is a chain or not. We use the business-name feature and count the number of times each business name appears in the dataset. If it only appears one time, the corresponding business (with respect to its id) is not a chain. If it appears more than once, the corresponding business is a chain. This feature is represented as a binary variable.

- **Business Attributes**: We use the following attributes as features in our models: Business Accepts Credit Cards, Bike Parking, WiFi, Business Parking, Offers Alcohol, Has TV, Noise Level, Price Range, and Outdoor Seating. Noise Level and Price Range are ordinal categorical variables. The remaining features are represented as binary variables. When records don't have the above attributes, features represented by binary values are set to 0 (false). For Noise Level and Price Range (both range from 1-4), the value 2 is assigned.

- **Business Categories** [23]: The Yelp dataset has a feature called "categories". In spite of this name, relative to the official Yelp category list (link below), this feature provides a list of sub-categories for each business. Using the above Yelp category list, we map each sub-category to its main category in order to reduce model dimensionality and sparsity. Based on the mapping, each business is assigned categories from the following category list: Active Life, Arts and Entertainment, Automotive, Beauty and Spas, Education, Event Planning and Services, Financial Services, Food, Health and Medical, Home Services, Hotels and Travel, Local Flavor, Local Services, Mass Media, Nightlife, Pets, Professional Services, Public Services and Government, Religious Organizations, Restaurants, and Shopping.

  These categories are features represented as binary variables.

  For example, if a business in the Yelp dataset has `Aquariums and Museum` as its feature value for "categories", we utilize the Yelp category list to map Aquariums to the main category of Active Life and Museum to the main category of Arts & Entertainment. We set the corresponding binary variables to 1 and the remaining ones to 0.

- **External Sources** to boost business data:

  In the Yelp Dataset, each business has a corresponding American postal code. We make use of this to add the features of median income and population relative to the business' postcode.

  1) **Income Postcode** [24]: The median income for a postcode is obtained from the Individual Income Statistics released by the IRS. The IRS releases income data as brackets and not exact figures. This feature is an ordinal categorical variable. For instance, a value of 1 is indicative of income between \$1 and \$25,000, 2 is between \$25,000 and \$50,000, 3 is between \$50,000 and \$75,000, 4 is between \$75,000 and \$100,000, 5 is between \$100,000 and \$200,000, and 6 is income greater than \$200,000.

  2) **Population Postcode** [25]: The population for a postcode is obtained from the Zip Code Database. This feature is represented as a continuous variable.

**Attributes** of final dataset that we will use to develop and evaluate our classifiers:

- N: 153843 American businesses,

- D: 45 features

- y: is_Open (1: business is open, 0: business is closed)

- **x**: all the features mentioned above

# 4 Approach

## 4.1 Pre-Processing

In Section 3.1, we address the issue of pre-processing to obtain our final, parsed dataset that we can use to train and evaluate our models. We employ some more data preparation techniques like SMOTE and Min-Max normalization in our Training & Evaluation Pipeline.

## 4.2 SMOTE

In the Yelp Dataset, only around 20 percent of American Businesses are listed as closed. To address this class imbalance, we make use of SMOTE to generate synthetic samples to obtain a balanced class distribution while training, which helps boost our models' generalization ability.

SMOTE-NC [26] stands for Synthetic Minority Over-sampling Technique for Nominal and Continuous features, which is a class balancing technique based on nearest neighbors judged by Euclidean Distance between data points in feature space. SMOTE-NC is taken from the `imbalanced-learn` library in `sci-kit learn`, which helps to create synthetic data for categorical as well as quantitative features in the dataset.

Specifically, we apply SMOTE-NC to the training set only and not the validation or test set. This is done to prevent the bleeding of information. The SMOTE algorithm creates synthetic data points by utilizing nearest neighbors of samples. Thus, if the minority class' (closed businesses) nearest neighbors end up in the validation or test set, the synthetic data points in the training set partially capture their information. Therefore, if doing so, we would get over-optimistic values for the 10-Fold cross validation and test set accuracy.

## 4.3 Min-Max Normalization

We acquire insights from an article[2] to apply Min-Max normalization appropriately to avoid data leakage in the training and evaluation pipeline for the classifiers that employ data normalization (KNN, Logistic Regression, and Neural Nets). Since our KNN model uses distance to compare feature values, it is important for features to be scaled to the same range to avoid unfairly over-weighting or under-weighting features. Normalization is applied for logistic regression to speed up solver convergence. It is also applied for Neural Networks to make the training process more stable and final classification accuracy higher. Normalization is not applied for Random Forest. The application of Min-Max normalization in the Training & Evaluation Pipeline is explained further in section 4.4.

---

[2]https://machinelearningmastery.com/data-preparation-without-data-leakage/

## 4.4  Training & Evaluation Pipeline

We adopt the following approach to train and evaluate our classifiers:

- The final parsed dataset is divided into the train and test set using the classic 80:20 split.

- The train set is then used to perform 10-Fold cross validation using balanced accuracy metric (explained in Section 4.6) to inform model selection.

- In each fold of cross validation, the train set is divided into the train fold set and the validation fold set. SMOTE-NC is applied on the train fold set to upsample the minority class (closed businesses). If required, Min-Max normalization is fit to the upsampled training fold set and the transformation is applied to the upsampled training fold set and the validation fold set. The classifier is trained on the upsampled training fold set and then evaluated on the validation fold set.

- Wherever possible and appropriate, cross validation is used to guide model parameter selection.

- After having selected model parameters through 10-fold cross validation, we then upsample the entire train set using SMOTE. If required, Min-Max normalization is fit to the upsampled train set and the transformation is applied to the train and test sets. We then train the classifier. This classifier is then evaluated on the unseen testing data.

## 4.5  Machine Learning Models

### 4.5.1  KNN

Library: `sklearn.neighbors.KNeighborsClassifier` [27]

KNN variant: Distance Weighted Nearest Neighbors

$K$ value: Given the large dataset, we use the most commonly used value for $K$, which is the $\sqrt{N}$, where $N$ is the number of samples in the train set (the upsampled train set in this case). $N$ is 200110. Thus, $K$ is set to 450.

Distance Metric: We use Euclidean based on 10-Fold cross validation results. During our experiments, we get better results with Euclidean distance than with Manhattan distance.

### 4.5.2  Logistic Regression

Library: `sklearn.linearmodel.LogisticRegression` [28]

Convergence Algorithm: SAGA (Stochastic Average Gradient)

We selected this convergence algorithm because it performs better than other solvers during our assessment with 10-Fold cross validation.

### 4.5.3  Random Forest

Library: `sklearn.ensemble.RandomForestClassifier` [29]

Number of Decision Trees: 100

This algorithm draws 80% of the input training set and uses bootstrap sampling to construct decision tree.

Splitting Criterion: Information Gain

Stopping Criteria: 1) Less than or equal to 500 samples at node (decided based on 10-Fold cross validation). 2) No more splits left. 3) All samples at node have same label. 3) All samples at node have same label.

### 4.5.4 Neural Network

Library: `PyTorch` [30]

We use the Multi Layer Perceptron Artificial Neural Network as our model, which is a combination of 5 fully connected layers and activation functions called ReLU. Here we choose cross entropy as the loss function, and utilize the learning rate of $1 \times 10^{-3}$, batch size 1024, 20 training epochs and Adam optimizer to conduct the experiments.

## 4.6 Metrics

Here we show some important definitions:

**True Positive Rate (Sensitivity):** Proportion of businesses classified correctly as open from the set of open businesses.

**False Positive Rate:** Proportion of businesses classified incorrectly as open from the set of closed businesses.

**True Negative Rate (Specificity):** Proportion of businesses correctly classified as closed from the set of closed businesses.

**False Negative Rate:** Proportion of businesses classified incorrectly as closed from the set of open businesses.

**Balanced Accuracy:** This is defined as the average recall obtained on each class (average of sensitivity and specificity). We use this metric across both the 10-fold cross validation and test data because the validation and test sets are imbalanced. The goal is to find a classifier that performs well across both classes (open and closed businesses) and this metric helps us in achieving this.

**Vanilla Accuracy:** The common accuracy metric is used on the unseen test data as well to add further perspective to the results.

**Class-specific metrics:** We also compute precision, recall and F1 score on the unseen test data to assess how our classifiers perform with respect to each of the two classes.

## 4.7 Library Usage

We use libraries for our classifiers instead of using the code from class for the following reasons:

- Code from libraries is cleaner and easier to organize given that only function calls are involved.

- Using libraries allows model parameters to be easily updated.

- In the case of logistic regression, the sci-kit learn model converges faster than the code used for class.

# 5 Results

## 5.1 Model Evaluation

From the ML models that we use to check for business survival including KNN, Logistic Regression, Random Forest, and Neural Net, we get the following results:

| Classifier | 10-Fold CV Balanced-Accuracy (%) | Test Balanced-Accuracy (%) | Test Vanilla Accuracy (%) |
|---|---|---|---|
| KNN | 71.20 | 70.82 | 77.45 |
| Logistic Regression | 70.34 | 70.36 | 76.87 |
| Random Forest | 70.41 | 70.56 | 82.65 |
| Neural Network | 71.32 | 71.72 | 79.31 |

Table 1: Model Evaluation

From Tab. 1 we can see that all models have relatively similar performance on the test balanced-accuracy metric with the Neural Net performing slightly better than the rest. The Random Forest classifier has the best performance according to the vanilla accuracy metric.

## 5.2 Class-specific metrics

| | Class 0 (Closed Businesses) | | | Class 1 (Open Businesses) | | |
|---|---|---|---|---|---|---|
| Model Type | Precision (%) | Recall (%) | F1 | Precision (%) | Recall (%) | F1 |
| KNN | 41.84 | 60.38 | 49.43 | 90.19 | 81.26 | 85.49 |
| Logistic Regression | 40.90 | 60.10 | 48.67 | 90.05 | 80.62 | 85.07 |
| Random Forest | 52.50 | 51.51 | 52.00 | 89.22 | 89.60 | 89.41 |
| Neural Network | 44.97 | 59.77 | 51.33 | 90.31 | 83.67 | 86.86 |

Table 2: Model Evaluation on Each Class

Across the board, the models are better at predicting open businesses than closed ones. All models have higher precision, recall and F1 scores for Open Businesses.

Random Forest is the best according to the sensitivity metric, KNN is the best according to the specificity metric, and as mentioned above, Neural Network is the best when balancing the sensitivity and specificity metrics (balanced accuracy).

## 5.3 Feature Importance

Algorithm: Permutation Importance, Library: `eli5`

This method determines feature importance by assessing how much the score for the respective metric (we use the balanced accuracy metric) decreases when the values for a feature are randomly shuffled.

We use this method because it is model agnostic and allows us to determine the important features for various models (didn't use neural networks and KNN since they take too long). In this paper, we present the important features for the Random Forest model, which achieves the third highest balanced accuracy on the test set. The following are the 10 most important features in determining business survival:

1) Restaurants 2) Call_To_Action_enabled 3) Grubhub_enabled 4) Home_Services 5) Covid_Banner

```
6) Health_&_Medical 7) Review_Count 8) delivery_or_takeout 9) Stars 10) Local_Services
```
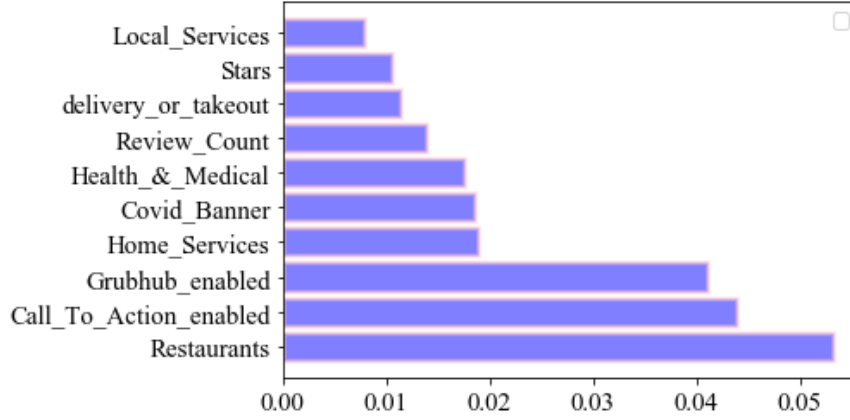


Figure 1: The Top 10 most important features. For example, a weight of 0.05 means balanced-accuracy decreases by 5% without the corresponding feature.

Based on Fig. 1, COVID features like `Call_To_Action_enabled, Grubhub_enabled, Covid_Banner`, and `delivery_or_takeout` are important in this classification task. It is also worth noting that Health_&_Medical is an important feature.

# 6 Conclusions and Future Work

## 6.1 Conclusions

In this paper, we build a comprehensive business dataset with rich features, collected from diverse information sources. Then we formulate the problem of predicting business survival and employed several machine learning methods to solve it. In this process, we develop an exhaustive training and evaluation pipeline combining techniques like SMOTE, cross validation and make use of multiple performance metrics. Results show that all models are significantly better at classifying open businesses than closed ones. Finally, we determine the 10 most important features that influence business survival. The presence of four COVID features suggests that a business' survival prospects are reliant on providing services that ensure customer health and safety. For eg, restaurants having Grubhub enabled. We hope that our work can help save the economy by informing and guiding business owners, investors and other people in power to make wise business decisions.

## 6.2 Future Work

The Yelp dataset has a file called `review.json` that contains reviews written by users for businesses. Users and Businesses are identified through `user_id` and `business_id` respectively. Natural Language Processing can be leveraged and this data can be used to perform sentiment analysis. This will give us a sense of what customers think about the business and this sentiment can be modelled as a feature to help us possibly obtain better classifiers. In addition to this we would also like to gather a larger dataset with more countries in order to build models that generalize better and aren't limited to specific countries.

# References

[1] Accenture. Outmaneuver uncertainty: Navigating the human and business impact of covid-19: https://www.accenture.com/us-en/about/company/coronavirus-business-economic-impact, 2020.

[2] Alexander W Bartik, Marianne Bertrand, Zoë B Cullen, Edward L Glaeser, Michael Luca, and Christopher T Stanton. How are small businesses adjusting to covid-19? early evidence from a survey. Technical report, National Bureau of Economic Research, 2020.

[3] Paul Horton and Kenta Nakai. Better prediction of protein cellular localization sites with the it k nearest neighbors classifier. In *Ismb*, volume 5, pages 147–152, 1997.

[4] AnalyticsVidhya. Introduction to k-nearest neighbors: A powerful machine learning algorithm: https://www.analyticsvidhya.com/blog/2018/03/introduction-k-neighbours-algorithm-clustering/, 2018.

[5] Irina Rish et al. An empirical study of the naive bayes classifier. In *IJCAI 2001 workshop on empirical methods in artificial intelligence*, volume 3, pages 41–46, 2001.

[6] Andrew McCallum, Kamal Nigam, et al. A comparison of event models for naive bayes text classification. In *AAAI-98 workshop on learning for text categorization*, volume 752, pages 41–48. Citeseer, 1998.

[7] AnalyticsVidhya. naive-bayes-for-machine-learning:https://machinelearningmastery.com/naive-bayes-for-machine-learning/, 2020.

[8] Scott Menard. *Applied logistic regression analysis*, volume 106. Sage, 2002.

[9] David G Kleinbaum, K Dietz, M Gail, Mitchel Klein, and Mitchell Klein. *Logistic regression*. Springer, 2002.

[10] TowardsDataScience. Logistic regression:detailed overview:https:https://towardsdatascience.com/logistic-regression-detailed-overview-46c4da4303bc, 2018.

[11] S Rasoul Safavian and David Landgrebe. A survey of decision tree classifier methodology. *IEEE transactions on systems, man, and cybernetics*, 21(3):660–674, 1991.

[12] Yoav Freund and Llew Mason. The alternating decision tree learning algorithm. In *icml*, volume 99, pages 124–133, 1999.

[13] TowardsDataScience. Decison tree in ml:https://towardsdatascience.com/decision-trees-in-machine-learning-641b9c4e8052, 2017.

[14] Andy Liaw, Matthew Wiener, et al. Classification and regression by randomforest. *R news*, 2(3):18–22, 2002.

[15] Vladimir Svetnik, Andy Liaw, Christopher Tong, J Christopher Culberson, Robert P Sheridan, and Bradley P Feuston. Random forest: a classification and regression tool for compound classification and qsar modeling. *Journal of chemical information and computer sciences*, 43(6):1947–1958, 2003.

[16] TowardsDataScience. Understanding random forest:https://towardsdatascience.com/understanding-random-forest-58381e0602d2, 2019.

[17] Andrzej Cichocki, Rolf Unbehauen, and Roman W Swiniarski. *Neural networks for optimization and signal processing*, volume 253. wiley New York, 1993.

[18] Sepp Hochreiter and Jürgen Schmidhuber. Long short-term memory. *Neural computation*, 9(8):1735–1780, 1997.

[19] Bayya Yegnanarayana. *Artificial neural networks*. PHI Learning Pvt. Ltd., 2009.

[20] Julian D Olden and Donald A Jackson. Illuminating the "black box": a randomization approach for understanding variable contributions in artificial neural networks. *Ecological modelling*, 154(1-2):135–150, 2002.

[21] David Alvarez Melis and Tommi Jaakkola. Towards robust interpretability with self-explaining neural networks. *Advances in Neural Information Processing Systems*, 31:7775–7784, 2018.

[22] Quanshi Zhang, Ying Nian Wu, and Song-Chun Zhu. Interpretable convolutional neural networks. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 8827–8836, 2018.

[23] Yelp. categories yelp:https://blog.yelp.com/2018/01/yelp_category_list#section20, 2020.

[24] IRS. Irs:https://www.irs.gov/statistics/soi-tax-stats-individual-income-tax-statistics-2016-zip-code-data-soi, 2020.

[25] ZIP Code List. Zip code database:https://www.unitedstateszipcodes.org/zip-code-database/, 2020.

[26] Nitesh V Chawla, Kevin W Bowyer, Lawrence O Hall, and W Philip Kegelmeyer. Smote: synthetic minority over-sampling technique. *Journal of artificial intelligence research*, 16:321–357, 2002.

[27] SciKitLearn. knn:https://scikit-learn.org/stable/modules/generated/sklearn.neighbors.KNeighborsClassifier.html, 2020.

[28] SciKitLearn. lr:https://scikit-learn.org/stable/modules/generated/sklearn.linear_model.LogisticRegression.html, 2020.

[29] SciKitLearn. Dt:https://scikit-learn.org/stable/modules/generated/sklearn.ensemble.RandomForestClassifier.html, 2020.

[30] Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, Alban Desmaison, Andreas Kopf, Edward Yang, Zachary DeVito, Martin Raison, Alykhan Tejani, Sasank Chilamkurthy, Benoit Steiner, Lu Fang, Junjie Bai, and Soumith Chintala. Pytorch: An imperative style, high-performance deep learning library. In H. Wallach, H. Larochelle, A. Beygelzimer, F. d'Alché-Buc, E. Fox, and R. Garnett, editors, *Advances in Neural Information Processing Systems 32*, pages 8024–8035. Curran Associates, Inc., 2019.

[31] WordStream. How the coronavirus (covid-19) pandemic is affecting small businesses & marketers, 2020.

[32] SciKitLearn. knn:https://scikit-learn.org/stable/modules/generated/sklearn.svm.LinearSVC.html, 2020.

[33] SciKitLearn. Dt:https://scikit-learn.org/stable/modules/generated/sklearn.tree.DecisionTreeClassifier.html, 2020.

[34] PyPi. nb:https://pypi.org/project/mixed-naive-bayes/, 2019.

[35] Alexander W Bartik, Marianne Bertrand, Zoe Cullen, Edward L Glaeser, Michael Luca, and Christopher Stanton. The impact of covid-19 on small business outcomes and expectations. *Proceedings of the National Academy of Sciences*, 117(30):17656–17666, 2020.

[36] Marti A. Hearst. Support vector machines. *IEEE Intelligent Systems*, 13(4):18–28, July 1998.

# 7  Appendix

## 7.1  Additional Machine Learning Models

### 7.1.1  Support Vector Machine Model

Library: `sklearn.svm.LinearSVC()` [36, 32]

We use `sklearn.svm.LinearSVC()` because `sklearn.svm.SVC()` is not suitable for the large-scale dataset because of it time complexity. Here we choose $l_2$ norm to represent the distance between different features, and use squared hinge loss as our loss function.

### 7.1.2  Naive Bayes

Library: `mixed_naive_bayes.MixedNB` [34]

## 7.2  Additional Results

| Classifier | 10-Fold CV Balanced-Accuracy (%) | Test Balanced-Accuracy (%) | Test Vanilla Accuracy (%) |
|---|---|---|---|
| Support Vector Machines | 70.45 | 70.53 | 76.66 |
| Naive Bayes | 68.65 | 68.31 | 67.81 |

Table 3: Model Evaluation for Support Vector Machines and Naive Bayes

| | Class 0 (Closed Businesses) | | | Class 1 (Open Businesses) | | |
|---|---|---|---|---|---|---|
| Model Type | Precision (%) | Recall (%) | F1 | Precision (%) | Recall (%) | F1 |
| Support Vector Machines | 40.67 | 60.88 | 48.76 | 90.18 | 80.17 | 84.88 |
| Naive Bayes | 32.19 | 69.10 | 43.92 | 90.73 | 67.52 | 77.42 |

Table 4: Model Evaluation on Each Class for Support Vector Machines and Naive Bayes