Radius Collider Submission - The Approximators

Myles Scolnick, Shane Barratt and Alexander Danilychev Jr

I. INTRODUCTION

The goal of this project was to determine the best North American Industry Classification System (NAICS) code for a business based on its name, address, description, and website. To solve this problem, we incorporated several techniques from a variety of fields: text cleaning, data augmentation, tf-idf [2], wordnet [3], neural word embeddings [4] and cross-validation. A pure supervised learner would be tough to model/train with less training data points than classes, so we went with what one could call a 'topic modeling' approach.

II. DATA

To increase the amount of data on businesses and NAICS codes, we aggregated some data from public APIs:

- Google Places business type
- NAICS code titles
- NAICS code descriptions
- Frequency of NAICS 2-long industry codes

The respective scripts are listed in List 1 in the appendix.

We also created a function to remove stop words, lemmatize, lowercase, tokenize. We did this to normalize the text being compared. We also enhanced the richness of the text by adding synonyms obtained through the WordNet model.

III. WEB APP

A. Assisted Hand Classifier

To ease the painful process of hand-classification, we built a web application in Flask/Python that displays information about the business to be classified and a list of NAICS codes sorted by a preliminary algorithm's scoring function. This allowed more interactive classification - that is automatically written to a file. This greatly sped up hand-classifying 1,000 businesses. Refer to Figure 1 of the appendix.

B. Database View and Code Comparison

To further understand the structure of the data and the quality of our classifier, we also built an additional endpoint to read the hand-classification and algorithm-classification databases. It displays a table of businesses with our hand-classification and algorithm-classification so that we could discern what we got it wrong and why. Refer to Figure 2 of the Appendix.

IV. TECHNICAL MODELING DETAILS

Let the set of businesses be denoted by $\{b_1, b_2, \ldots, b_n\}$ and the set of NAICS codes by $\{c_1, c_2, \ldots, c_m\}$. The algorithm attempts to construct a similarity matrix \mathbf{S} of dimension $n\mathbf{x}m$ where $\mathbf{S}_{i,j}$ denotes the similarity between business i and NAICS code j. The classification for business (row) i is $\arg\max_{1 \leq j \leq m} \mathbf{S}_{i,j}$, the most similar NAICS code to that

business. To construct this matrix, we first construct 9 similarity matrices, $\{S_1, S_2, \ldots, S_8\}$ with the same dimension. There are 4 ways to pair the business[.name .description] and NAICS[.title .description]. So the *tf-idf* cosine similarity and *word2vec* cosine similarity between these 4 combinations as well as a similarity matrix of 2-digit industry priors (independent of the business) constitute the 9 similarity matrices. The final similarity matrix is a convex combination of the preliminaries:

$$\mathbf{S} = \sum_{i=1}^{8} w_i \mathbf{S}_i,$$

where $\sum_{i=1}^{9} w_i = 1$. The weights were selected through random hyper-parameter search [1].

V. CONCLUSION

A. Results

After training, our algorithmic classifier scores 18.5 %. We hand-classified 1000 businesses.

B. Confidence in our Predictions

• We are not particularly confident in the predictions our algorithm made, as it has lots of false-positives.

C. Limitations

 The algorithm did not take into account the hierarchical structure of the NAICS codes.

D. Challenges

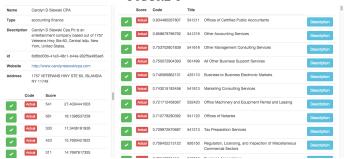
- The data did not come perfectly formatted and sometimes the title, description, website did not actually match and was a mix of 2 businesses.
- Some of the businesses had names and descriptions in other languages.
- Some businesses have multiple obvious classifications making it difficult to assign a single NAICS code (e.g. a bakery which also is a sit-down restaurant).

REFERENCES

- Bergstra, James, and Yoshua Bengio. "Random Search for Hyperparameter Optimization." The Journal of Machine Learning Research 13.1 (2012): 281-305.
- [2] Brown, Peter F., et al. "A Statistical Approach to Machine Translation." Computational Linguistics 16.2 (1990): 79-85.
- [3] Miller, George A., et al. "Introduction to Wordnet: An On-line Lexical Database*." *International Journal of Lexicography* 3.4 (1990): 235-244.
- [4] Mikolov, Tomas, et al. "Distributed Representations of Words and Phrases and their Compositionality." Advances in Neural Information Processing Systems. 2013.

APPENDIX LIST 1

FIGURE 1



- $\bullet \hspace{0.2cm} get_business_latlot.py: \hspace{0.2cm} generates \hspace{0.1cm} id_to_loc.pickle \\$
- get_business_types.py: generates business_types.pickle
- get_naics_data.py: generates naics_list.json

FIGURE 2

							0.700023030700	200110	LAGRA		423	16 9666223369
3	11111	811198	e4e95937- cb14-43b6-	Midas	store car repair	Midas is one of the world's largest providers of auto repair services, including brakes, oil change, tires, maintenance, stearing, and exhaust services. Visit your local Midas for additional services.	1.09538479258		All Other Automotive Repair and Maintenance	0	333	26.575619758
			aaf5- 3ead191a04f6									24.2358276003
											311	22.9722858433
							1.07409094769	811118	Other	0	392	22 4224605487
6	41310	441310	b3b4b350- 4ab5-45c7- b245- 5fec410df635	CARQUEST Auto Parts	store car repair	GARQUEST Auto Parts is the premier supplier of replacement car parts, truck parts, products, accessories, supplies and equipment for virtually all makes and models	1.02984011627		Parts and Accessories Stores	Θ	333	25.3351418132
											541	23.3779691387
											332	22.4708595204
							1.01911408181	441110	New Car	0	423	21 7787447743
	24460	114112	ef5d5c98- 389a-4870-	Langsford Road Lobster		We are a retail fish market and have live or cooked lobsters. We have a huge selection of fresh fish.	0.758370702619	114112		0	311	22.275019893
-2			af17- 365194972106	& Fish		Located on picturesque Cape Porpoise Harbor, We are open from May to the end of Sept. We ship lobsters overnight by UPS all year long.	0.726748048268	721214	Pishing Recreational and Vacation Camps (except	0	423	15.7205413163
											424	15.5128382911
											111	15 1711607889
	22320	813930	11906302- 8d73-4e70- b04e- 3004e81f9ab4	western union	finance	Western Union - When you need to send money in person, Western Union ' Fremont agent clerks will help you send money quickly and easily.	1.17889035472	813930	and Similar Labor Organizations	0	541	14.6252252593
-2											561	13.5959069048
											423	10.0651817087
							1.11611869498	522130	Credit Unions	0	333	9.41583712179
6	22513	722513	47745e51- f9f1-4aa1- befc- c0e31a3c7ea2	BOB EVANS RESTAURANT 2069	meal takeaway restaurant	Goodness of the Farm' to you through quality food delivered in a friendly atmosphere at fair prices. Open for breakfast, funch & dinner, As of 34/12, Bob Evans owns & operates 564 family restaurants in 18 states. For more info, places visit www.bobevans.com	0.777894635189	722513	Limited- Service Restaurants	0	311	18.5327977263
											424	13.4640352215
							0.664684038522	424330	Women's, Children's, and	0	111	13.3050203376
											423	10 9035281207
0 713110	13110	0	0e16ta6e- 13e9-49f5- 8f5c- 0e743c0778f3	Laugh Out Loud Stations		celebration, and entertainment center where families, friends, and co-workers can run, play, eat, compee, and justLaugh Out Loudl Laser Tag, Beltadium, Spider Mountain, Land of Laugh-A-Lot Playground, Toddier World, Big Kid Bouncers,	0.594100656206 0.577185289578		and Garages Other Building Material	0	311	10.1857762552
											423	9.72452905755
											332	9.54718061495
						Climbino Wall, Game Arena, Kiddle Rides, &			Dealers		333	9.33353436105
621498	21498	3	d178ea2c- 6c0b-4a2a- bc78-	Jing		At Jing, the premier holistic healthcare center in Louisville, Kentucky, we care about your health and well being. Whether you're suffering from	0.352484785101	621498	All Other Outpatient Care Centers	0	541	8.34780868699
											333	5.80382187685
			0856979566d7			chronic pain, stress, or just want to improve your				-	601	5.78234314258