

Deep Learning for Knowledge Graph Completion

Nick, Kevin, Cesar Bartolo-Perez, Doug

University of California - Davis

November 30th, 2017

Abstract

Contents

1 Introduction	1
2 Methods	1
3 Results	3
4 Discussion	4

1 Introduction

2 Methods

Multilayer Perceptron Model (MLP)

The dataset we considered was queried from the WikiData, a free opensource knowledge base that acts as the central storage device for its sister projects including Wikipedia, Wikivoyage, and Wikisource [1]. Wikidata offers a SPARQL type queries for obtaining data on a variety of sources. We queried data regarding the lineage of many famous people in history including: Frederik, Crown Prince of Denmark, George Windsor, Earl of St Andrews, and Shigeko Higashikuni. The relations we considered were if one person, or entity, was the Father, Mother, Spouse, Sibling, or Child of another. See Figure 1 for an example of the query we used to obtain our dataset.

The raw data is a table of values that lists each person and their Father, Mother, spouse, sibling, and child; leaving an entry empty if unknown. We processed this table into positive triplets of the form $(Q193752, P25, Q229279, 1)$ where $Q193752$ and $Q229279$ are a unique encoding of the target and related entities respectively, $P25$ represents one of the relations, or predicates, and 1 indicates that this statement is true. In english this triplet implies that $Q229279$ is the $P25$ of $Q193752$. These positive facts represent all the data we can pull directly from the knowledge graphs. However, to increase the robustness of our data, we can create negative facts as well. For example, since we know that Cleopatra VII Philopator, the last active pharaoh of Ptolemaic Egypt, was the daughter of Ptolemy XII Auletes, we can conclude she was not the daughter of Richard Nixon. Thus we can extend our data to contain these negative triplets of the form $(Q1, P1, Q2, -1)$. In all, we were able to produce 37395 positive triplets and 37391 negative triplets from our raw data.

```

SELECT ?human ?humanLabel ?spouse ?spouseLabel
?child ?childLabel ?sibling ?siblingLabel
?father ?fatherLabel
?mother ?motherLabel
WHERE {
    ?human wdt:P31 wd:Q5.
    ?human wdt:P26 ?spouse.
    ?human wdt:P40 ?child.
    ?human wdt:P3373 ?sibling.
    ?human wdt:P22 ?father.
    ?human wdt:P25 ?mother.
    SERVICE wikibase:label { bd:serviceParam wikibase:language "en". }
}
LIMIT 10000

```

Figure 1: SPARQL query used to obtain our raw data from Wikidata

In an effort to exploit the uniqueness in the names for each of our entities, we used word embeddings of each of the words in a name to generate entity embeddings. We used the Fasttext model to generate word embeddings from semantic analysis across Wikipedia. Fasttext is a library that represents words as bags of character n -grams for fast and efficient classification [2, 3]. Thus, we used Fasttext to obtain 300×1 word embeddings of each of the unique words found in all of our entities' names. For example, the name "Charles Stuart, 1st Earl of Lennox", we found an 300 dimensional embedding vector for Charles, Stuart, 1st, Earl, of, Lennox. Once each of the words had embeddings, we could construct entity embeddings for each of our entities. The method we used for this was to aggregate the word embeddings by taking the mean across each of the words within an entity name. For the MLP model, if Fasttext could not produce a valid embedding for a given word, then that word was ignored. Moreover, if an entity consisted entirely of unknown words, then that entity was thrown out. We did this to preserve the integrity of our input data because in the standard MLP model we would not be training these embeddings. However, for more complex models, another suitable technique is to randomize the embeddings for these entities, or all entities, to guarantee a more representative dataset. Overall, only about 6.4% of the entities had no suitable embedding. Once we encoded each entity with their embedding we produced a dataset where each row contains the numeric representation of a triplet. Hence, for a triplet given by (EntityA, Pred2, EntityB, 1), where the embedding for EntityA and EntityB is $[a_1, a_2, \dots, a_{300}]^T$ and $[b_1, b_2, \dots, b_{300}]^T$ respectively, the associated row in our data set is given by

$$[a_1, a_2, a_3, \dots, a_{300}, b_1, b_2, \dots, b_{300}, 0, 1, 0, 0, 0]$$

where the last 5 columns represents the truth values for each of the 5 predicates (Father, Mother, Spouse, Sibling, Child).

Once we had a comprehensive dataset to train on, we designed our predictive model. We implemented a Multilayer Perceptron (MLP) model using the RSNNS package following the Stuttgart Neural Network Simulator[4]. Figure 2 shows a general look at the architecture for our MLP design. We tested a few variations of the MLP model to determine the suitable parameters for our dataset; including a variable number of layers, activation functions, and nodes per layer. Moreover, we tuned the hyperparameters such as the learning rate and momentum term for backpropagation. We will reference Figure 2 as we discuss the effects of the different architectures.

Where we have used standard backpropagation and backpropagation with momentum methods in our simulation models.

Different MLP architectures are built, varying number of hidden layers, nodes and iterations in order to improve the accuracy and precision of our model.

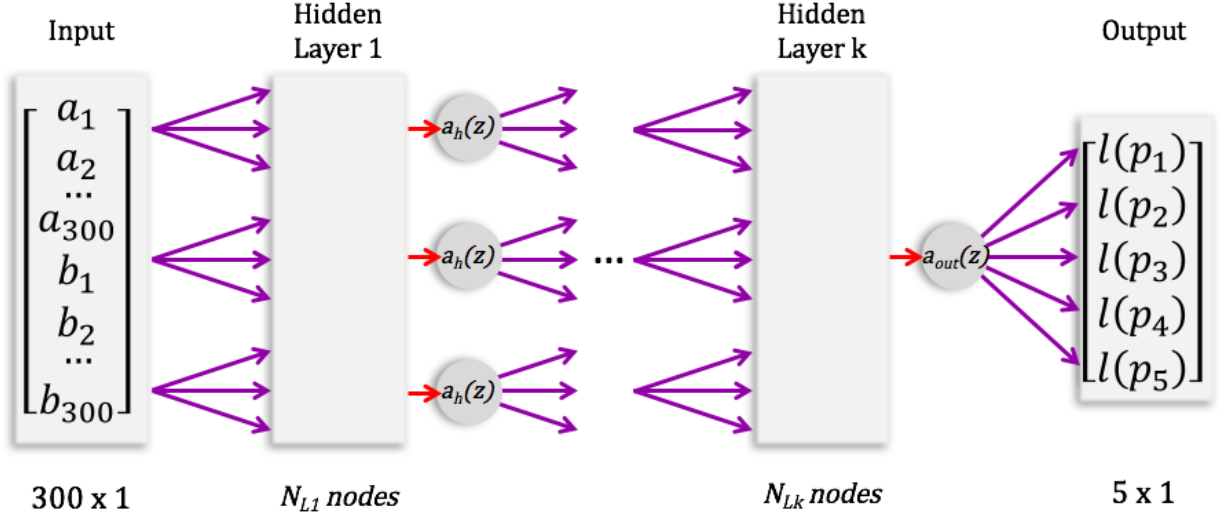


Figure 2: A general look at the Neural Network architecture used for a MLP model with k hidden layers with N_{L_i} nodes for the i th layer, a common activation function $a_h(z)$ for each layer, and a final output activation function $a_{out}(z)$. We have 600 features from our input, 1 for each dimension of both entity embeddings, and we are predicting the true/false value of our 5 output predicates. Note that $l(p_1)$ is the likelihood that the predicate p_1 is true.

Weighted SSE by number of iterations of the training and testing sets in addition to ROC curves will provide the references in our model architecture selection.

3 Results

In order to prevent overfitting in our model, we have compared our model with 100 and 50 iterations (Figure 3). We can notice that under 100 iterations, the training error drops to a weighted SSE of almost 500. However, our testing error, start increasing after 30 iterations approximately. From this results, we can infer that we have reached an overfitting in our model and a smaller number of iterations are enough. In that sense, we have run our model with 50 iterations, avoiding the increase in the SSE.

In a second running of the model, we have added a new hidden layer with different number of nodes (16,64 and 128). Before 50 iterations, adding nodes to the second hidden layer, reduce the error in our model. However, after 50 iteration, the higher number of nodes, increase rapidly the overfitting of the model.

In a following set of tests, we compared different Neural Networks structures, increasing the number of hidden layers (hidden layers/nodes: 128, 128/16 and 128/128/16). Figure 4 shows that for the model with three hidden layers the training error is not decreasing considerably before 40 iteration. Suddenly, the SSE in this NN structure drops abruptly. The training error, shows that the error keep increasing at that range of iterations. This behavior, can be attributed to an overfitting in our model.

This trend described above is similar for two hidden layers. For the case of just one hidden layer with 128 nodes, the error does not drop abruptly for the training set but we can see an increase in the error for after 30 iterations approximately.

These results, make us suggest that only one hidden layer is enough for our MLP model, since there is not a big improvement in the reduction of error with more layers added and with the benefit of less computational resources are required. This decision is confirmed comparing their ROC curves for each NN structure.

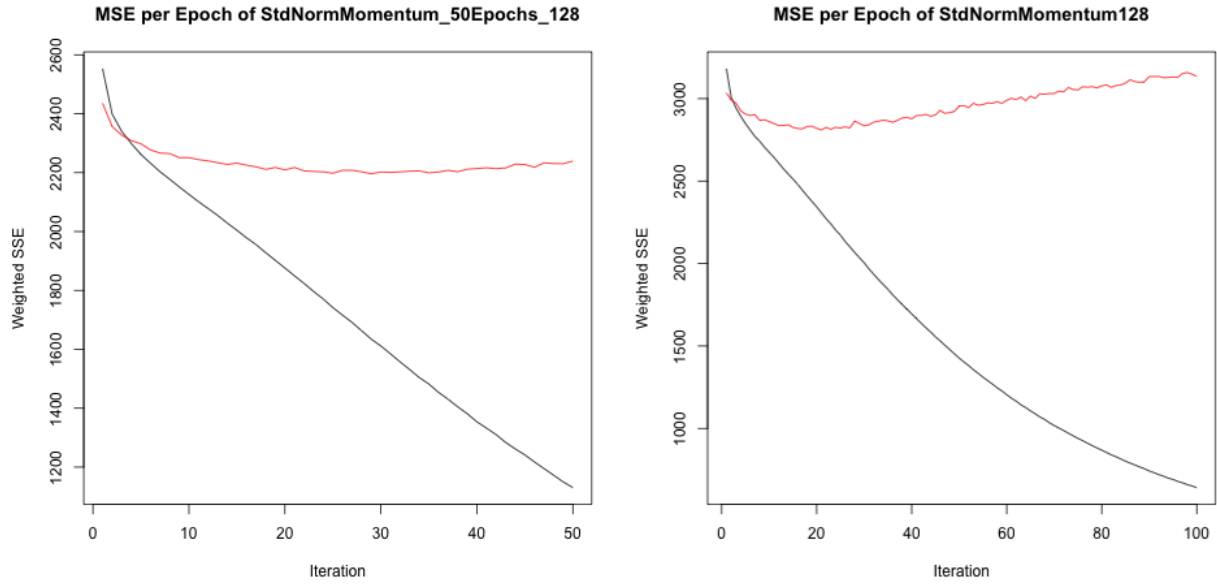


Figure 3: MSE at different iterations. The structure presents a single hidden layer with 128 nodes

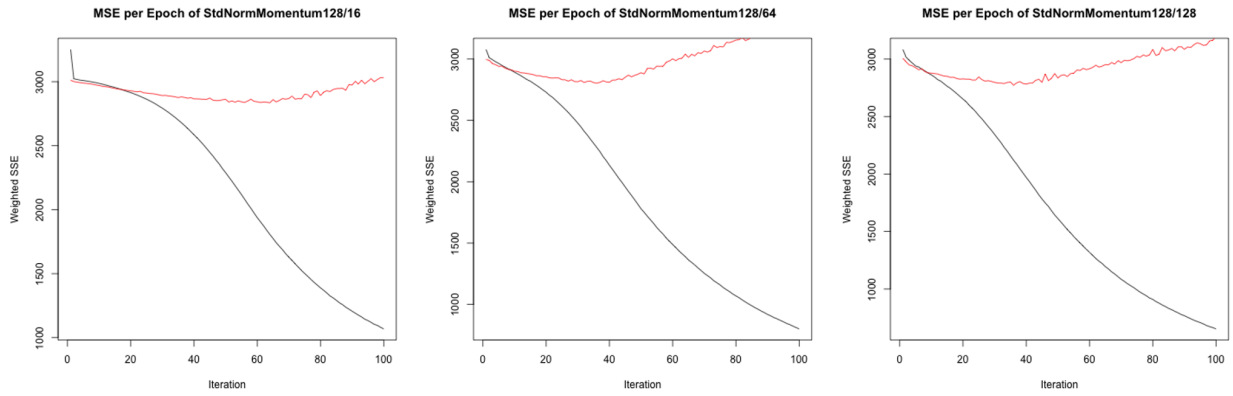


Figure 4: MSE for A NN with two hidden layers with different nodes at second hidden layer.

4 Discussion

References

- [1] “Wikidata knowledge base (main page).” https://www.wikidata.org/wiki/Wikidata:Main_Page. Accessed: 2017-11-26.
- [2] P. Bojanowski, E. Grave, A. Joulin, and T. Mikolov, “Enriching word vectors with subword information,” *CoRR*, vol. abs/1607.04606, 2016.
- [3] A. Joulin, E. Grave, P. Bojanowski, and T. Mikolov, “Bag of tricks for efficient text classification,” *CoRR*, vol. abs/1607.01759, 2016.

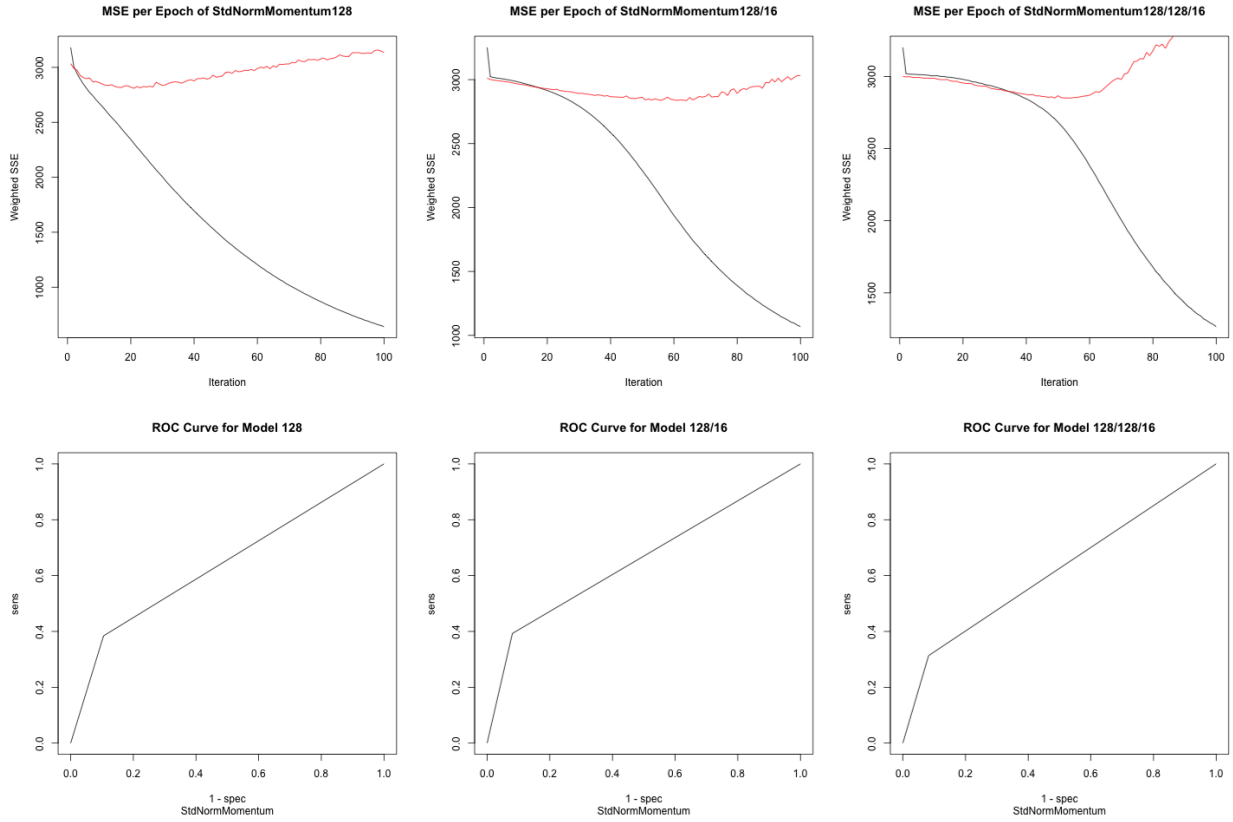


Figure 5: MSE and ROC curves for different levels of hidden layers

- [4] C. Bergmeir and J. M. Benítez, “Neural networks in R using the stuttgart neural network simulator: RSNNs,” *Journal of Statistical Software*, vol. 46, no. 7, pp. 1–26, 2012.