# Modular Deep Encoder-Decoders

*Group 40: Arnold, Gurunat, Huang, & Qu*

April 27, 2016

### Abstract

In this short paper we propose an approach to transfer learning using rich vector embeddings. The suggested technique can be applied to any supervised task, and it handles multiple sources and changing sources of data without the need for retraining. To verify our ideas, we apply our ideas to the task of text-classification.

## I  Introduction

Our goal is to generate rich vector embeddings from articles to classify them into predefined categories. Then, we want to extend our model with an additional source of data: the title of the videos. To handle this new knowledge source without retraining our previous model, we suggest to generate a new embedding that will be used to modify the original one. This combined embedding will then be used for the classification task.

More formally, we want to jointly learn a set of *encoder* functions $\{E_i\}_0^N$ mapping samples $x \sim \chi_i$ from a set of data distributions $\{\chi_i\}_0^N$ to a fixed-sized vector embedding $V$.

$$\forall i \in [0; N] : E_i(x) : \chi_i \to V_i \in \Re^M$$

$$V = \sum_i^N V_i$$

The embedding $V$ is then fed into a decoder function $D$ which in the case of classification learns a mapping from the vector space of $V$ to the label space $L$.

$$D(v) : \Re^M \to L \in \Re^D$$

Finally, we want to extend the set of encoders by learning a new encoder $E_{N+1}$ which handles samples from a different dataset $\chi_{N+1}$, without retraining our trained decoders and encoders.

$$E_{N+1}(x) : \chi_{N+1} \to V_{N+1} \in \Re^M$$

Although we could have used any kind of mapping, we chose to use deep learning algorithms, as they easily learn hierarchical representations and have been known to highly outperform other statistical techniques on natural language tasks. Learning embeddings has been previously done by [1], although our proposed contribution is slightly different.

[1] used deep autoencoders to obtain a better initialization of the parameters of their model. Our approach instead is much closer in spirit to the work of Sutskever [3] and Vinyals [4]. They both use encoders on a sequence of data to generate a *thought-vector* which will then be decoded in the desired terget sequence. In some sense, our proposal adds the transfer learning component to their contribution. This approach is also similar to the work of Karpathy & al [5] where they mapped images to their captions with embeddings.

# II   Database

We used the freely available dataset from `https://github.com/ParallelMazen/SaudiNewsNet` The dataset contains a total of 31,030 Arabic newspaper articles. Article objects are reprensented in json format, with UTF-8 encoding. We wrote a script to download the data from github repo, unzip each file, and read them in as json objects. Each json object contains the flollowing fields: source: A string identified of the newspaper from which the article was extracted. It can have one of the following values:

- source: A string identifier of the newspaper from which the article was extracted.

- url: The full URL from which the article was extracted.

- date_extracted: The timestamp of the date on which the article was extracted. It has the format YYYY-MM-DD hh:mm:ss.

- title: The title of the article.

- author: The author of the article.

- content: The content of the article.

# III   Procedure

### Data Pre-processing and Features

Firstly, the UTF-8 encoded data was checked for trailing spaces for removal. Next, the data was filtered by removing punctuation symbols, retaining only a subset consisting of {, . ! ( ) ?}. The data was then split into words which are used as tokens for building the vocabulary and finally obtaining a feature representation for training the neural network. A vocabulary is constructed by assigning unique integer word-ids to each unique word appearing in the corpora. The punctuations are treated as a separate token and contributes to the vocabulary. The url, date and title fields of the corpora were discarded since they were irrelevant for our purpose.

After building the vocabulary for the corpora, the original sentences of the corpora are mapped on a word-basis to obtain a series of integers representing the sentence equivalent which form the features to be input to the embedding layer of the neural network. Alternatively, the integer mapping could be replaced with one-hot sparse encoding of each word constituting the sentence. The features are shuffled and split for training and

testing. For efficient data storage and retrieval the features are compressed to HDF5 file format.

## Neural Network Architecture

### Embedding Layer

An embedding layer is used to compute a word embedding $W : words \mapsto \mathbb{R}^n$ which is a function mapping words of a language to a high dimensional vectors which are continuous, distributed representation of the vocabulary words with good semantic properties. The hope of such a representation is to project similar words to similar regions in the hyperspace from a semantic point of view.

## Baseline System

Two baseline systems are constructed, one to classify the journals based on the content of the articles and the second to classify the journals based on the authors of the articles. Identical architectures are used for both the baseline systems. The baseline system is a shallow neural network with one hidden Long Short Term Memory (LSTM) layer succeeding the embedding layer. The dimension of both the embedding layer and the LSTM was fixed to 128. The output layer is a soft-max layer with 14 dimension for each class representing each journal. Cross-Entropy loss function was used as an objective function to train the neural network using Adaptive-Gradient (Adagrad).

## Proposed System

The procedure is quite straight-forward. While part of the team worked on building our tailored dataset, the other half worked on the model definition and training procedure.

Once every pre-requisite is available, we trained the first encoder $E_1$ (implemented as a recurrent neural network) to build the embedding $V_1$. Since we only dealt with a relatively simple classification task, our decoder $D_1$ was simply a fully connected multi-layer perceptron. They were jointly trained end-to-end by propagating the gradients through the embedding from $D_1$ to $E_1$.

The second training step was to train the second encoder $E_2$. Again, we also performed training end-to-end, but specifically did not propagate the gradients through $E_1$.

**Evaluation Metrics**

# IV   Results

**Baseline**

**Proposed System**

# V   Discussion

# VI   References

Part of the relevant literature review. More literature was involved for the deep learning part.

1. Supervised Representation Learning: Transfer Learning with Deep Autoencoders, Zhuang & al, http://ijcai.org/Proceedings/15/Papers/578.pdf

2. Transfer Learning via Dimensionality Reduction, Pan & al,

   https://www.cse.ust.hk/~jamesk/papers/aaai08.pdf

3. Sequence to Sequence Learning with Neural Networks, Sutskever & al,

   http://arxiv.org/abs/1409.3215

4. Grammar as a Foreign Language, Vinyals & al, http://arxiv.org/abs/1412.7449

5. Deep Fragment Embeddings for Bidirectional Image Sentence Mapping, Karpathy & al, https://cs.stanford.edu/people/karpathy/nips2014.pdf