

This is my title

Bachelorarbeit

Christopher Schütz | 2462248

Bachelor Wirtschaftsinformatik



TECHNISCHE
UNIVERSITÄT
DARMSTADT

WIRTSCHAFTS
INFORMATIK



Christopher Schütz
Matrikelnummer: 2462248
Studiengang: Bachelor Wirtschaftsinformatik

Bachelorarbeit
Thema: "This is my title"

Eingereicht: 18. März 2020

Betreuer: M.Sc. Felix Peters

Prof. Dr. Peter Buxmann
Fachgebiet Wirtschaftsinformatik | Software & Digital Business
Fachbereich Rechts- und Wirtschaftswissenschaften
Technische Universität Darmstadt
Hochschulstraße 1
64289 Darmstadt

Ehrenwörtliche Erklärung gemäß §22 Abs. 7 APB der TU Darmstadt

Hiermit versichere ich, Chrisopher Schütz, die vorliegende Master-Thesis gemäß §22 Abs. 7 APB der TU Darmstadt ohne Hilfe Dritter und nur mit den angegebenen Quellen und Hilfsmitteln angefertigt zu haben. Alle Stellen, die Quellen entnommen wurden, sind als solche kenntlich gemacht worden. Diese Arbeit hat in gleicher oder ähnlicher Form noch keiner Prüfungsbehörde vorgelegen.

Mir ist bekannt, dass im Falle eines Plagiats (§38 Abs. 2 APB) ein Täuschungsversuch vorliegt, der dazu führt, dass die Arbeit mit 5,0 bewertet und damit ein Prüfungsversuch verbraucht wird. Abschlussarbeiten dürfen nur einmal wiederholt werden.

Bei der abgegebenen Thesis stimmen die schriftliche und die zur Archivierung eingereichte elektronische Fassung gemäß §23 Abs. 7 APB überein.

Darmstadt, den 18. März 2020



Abstract

This is my abstract.

Table of contents

List of figures	vi
List of tables	vii
1 Background	1
1.1 History of language modeling	1
1.1.1 N-Gram models	2
1.1.2 Neural language models	3
1.1.2.1 Neural network	3
1.1.2.2 Recurrent neural networks	3
1.1.2.3 Convolutional neural networks	3
1.1.2.4 Attention and transformers	3
1.2 Deep Learning	3
2 Methodology	4
2.1 Dataset	4
2.1.1 Extraction Process	5
2.1.2 Generation Parameter Combination	5
2.1.3 Data Building	7
2.2 Approach	8
2.3 Infrastructure	8
References	I
A Appendix	II
A.1 User groups	II
A.2 Confusion matrices	II



List of figures

Figure 1: Confusion matrices for linear retweet classification models II

Figure 2: Confusion matrices for deep feedforward classification models III

Figure 3: Confusion matrices for multi-input deep neural networks IV

List of tables

Table 1: Value ranges for all modified parameters. The sampling method indices refer to nucleus sampling (1), top k (2) and beam search (3).	6
Table 2: Exemplary word and character split inputs and according outputs	6

1 Background

This chapter serves to explain the foundations of natural language processing (NLP), especially the subpart of language modeling, needed to understand the problem and methodology used in this thesis. The two main building blocks of this thesis are the examination of a state-of-the-art language model (LM) and its ability to generate high quality, human-like text on the one hand and methods to distinguish such generations from human written text on the other hand. In order to understand the problem and methodology used in this thesis, this chapter explains the foundations of NLP (ch. 1.1), especially the subpart of language modeling, and the foundations of deep learning (ch. 1.2) especially under the aspect of sequence classification.

1.1 History of language modeling

Natural language processing (NLP) is “an area of research and application that explores how computers can be used to understand and manipulate natural language text or speech to do useful things”¹. The applications of NLP such as speech recognition, sentiment analysis, question answering and others are numerous². Moreover, these applications are already being heavily used by industry and consumers alike e.g. in the forms of digital voice assistants, sentiment analysis for recommender systems and browser search bars³. The subcomponent of NLP needed when it comes to tasks like machine translation, predictive typing or summarization that involve either generating text or estimating the probability of text is called language modeling. The following notation mostly follows the one from the CS224 Stanford Natural Language Processing with Deep Learning lecture by Chris Manning.

Language modeling is the task of predicting what word comes next. More formally, this means: Given a sequence of words $x^{(1)}, x^{(2)}, \dots, x^{(t)}$, compute the probability of the next word $x^{(t+1)}$:

$$P(x^{(t+1)} | x^{(1)}, \dots, x^{(t)}) \quad (1)$$

where $x^{(t+1)}$ can be any word in the vocabulary $V = \{w_1, \dots, w_{|V|}\}$

Having a system that does allows us to assign a probability to a snippet of text of length T :

$$\begin{aligned} P(x^{(1)}, \dots, x^{(T)}) &= P(x^{(1)}) \times P(x^{(2)} | x^{(1)}) \times \dots \times P(x^{(T)} | x^{(1)}, \dots, x^{(T-1)}) \\ &= \prod_{t=1}^T P(x^{(t)} | x^{(1)}, \dots, x^{(t-1)}) \end{aligned} \quad (2)$$

Developing and improving language models is a task central to language understanding by which we can measure how well machine learning systems actually comprehend natural language [cite either stanford or ice paper 3 or 4]. This is demonstrated by the fact that “often (although not

¹ Chowdhury (2003).

² Gatt & Krahmer (2017).

³ Valdivia; Luzón; Herrera (2017); Klopfenstein et al. (2017); Pandu (2019).

always), training better language models improves the underlying metrics of the downstream task (such as word error rate for speech recognition, or BLEU score for translation), which makes the task of training better LMs valuable by itself”⁴.

Since the first significant language model was proposed back in 1980⁵, language models and their architectures have gone through many changes. Especially the rise of Deep Learning and new network models such as RNNs or Transformers have fueled language modeling research in the past few years. The following chapters will cover the most common model types, heuristics and architectures (ch. 2.1.1 - 2.1.3 [STILL HARDCODED]).

1.1.1 N-Gram models

One solution in dealing with the problem of predicting a word after a sequence of $(n-1)$ words in the form of a Markov model, i.e. the probability of each event depends only on the state attained through the previous event, is called an n-gram model. An “n-gram” hereby denotes a chunk of n consecutive words. The core idea is that the probability of a word w_i occurring in the i^{th} instance after a sequence of $(i-1)$ preceding words can be approximated by observing only the preceding context of $(n-1)$ words.

Following this insight we can compute the probability of all n-grams in a corpus of text by simply counting their occurrences. Doing so allows us to calculate these conditional probabilities like so:

$$\begin{aligned} P(x^{(t+1)}|x^{(t)}, \dots, x^{(1)}) &= P(x^{(t+1)}|x^{(t)}, \dots, x^{(t-2+2)}) \\ &= \frac{P(x^{(t+1)}, x^{(t)}, \dots, x^{(t-2+2)})}{P(x^{(t)}, \dots, x^{(t-2+2)})} \\ &\approx \frac{\text{count}(x^{(t+1)}, x^{(t)}, \dots, x^{(t-n+2)})}{\text{count}(x^{(t)}, \dots, x^{(t-n+2)})} \text{ (statistical approximation)} \end{aligned} \tag{3}$$

N-Gram Models where $n=1$ are called Unigram, $n=2$ Bigram and $n=3$ respectively Trigram.

[Talk about the text generation process]

Even though n-gram models have been widely used especially due to their simplicity and scalability they do face certain limitations that have led to a decrease in their popularity. One problem that can arise when computing n-gram frequencies/probabilities is that n-grams encountered in a test setting do not appear in the corpus that the model was trained on. This leads to the probability of that n-gram being 0 [insert formula]. In order to counter this different smoothing techniques can be applied.

⁴ Jozefowicz et al. (2016).

⁵ Rosenfeld (2000).

Sparsity: what if “students opened their w” or “students opened their” never occurred? solutions
-> smoothing, backoff [cite stanford textbook chapter]

Storage: Need to store counts for all n-grams you saw in the corpus. NO SOLUTION.

Lack of understanding: Perhaps the biggest drawback of n-gram models, though, is their very limited context size. In practice no $n > 5$ usually. While produced output is often grammatically correct, there is an evident [lack] of coherence.

“today the price of gold per ton , while production of shoe lasts and shoe industry , the bank intervened just after it considered and rejected an imf demand to rebuild depleted european stocks , sept 30 end primary 76 cts a share”

Even though n-gram language models are still widely used in speech recognition due to their high efficiency in inference, their limitations caused by poor generalization to unobserved n-grams and inability to capture long range dependencies led to the rise of neural language models [ice paper 8].

1.1.2 Neural language models

1.1.2.1 Neural network

This is some text.

1.1.2.2 Recurrent neural networks

This is some text.

1.1.2.3 Convolutional neural networks

This is some text.

1.1.2.4 Attention and transformers

This is some text.

1.2 Deep Learning

This is deep learning.

2 Methodology

This chapter explains the methodology.

2.1 Dataset

In order to create a dataset suited for the classification task a few aspects had to be considered. As deep learning methods with many parameters that need to be optimized were going to be used for the classification task, many labeled training samples were needed. The dataset should ideally be comprised of equal amounts of synthetically and human generated texts in order to improve the accuracy of training.

The first idea was to look for already built datasets that are freely available as these would not only reduce the time and amount of work needed for the creation of the dataset, but also provide a benchmark against other models used on them. Because the examined classification task is fairly novel and powerful language models have just started to emerge in recent years, there is a lack of standardized data sets - prior research often focused on detection of artificially generated academic papers instead of short texts [reference to papers that use academic papers]. Furthermore, the incentive of using metadata related to text snippets and inspecting the changes in detection accuracy through it led to the motivation of building a new dataset.

For this purpose, the following sources of human created text were inspected - taking different aspects like data availability, extensibility by metadata, potential for use of text generation and minimal overlap with the pretrained GPT2 Model into account: Wikipedia, Twitter, reviews (e.g. Amazon, IMDb), Reddit comments, Reuters Corpus.

With everything taken into consideration, Wikipedia articles were chosen as the best fit for this work. Reasons for this choice were the ease of access, the vast amount of data entries (currently there are more than 6 million articles in the english Wikipedia⁶), the extensibility by metadata such as pageviews or categories and most importantly the fact that the ready-to-use GPT-2 model was explicitly not trained on any Wikipedia article since "it is a common data source for other datasets and could complicate analysis due to overlapping training data with test evaluation tasks."⁷

The reviews and reddit comments datasets were not chosen, because the metadata was not seen as decisive in improving detection quality. The disadvantage of the Reuters Corpus was that the training dataset used for the ready-to-use GPT2 model is comprised of many newspaper sources and thus is more likely to generate results that are rather similar to their human written counterparts. Twitter was seen as a promising data source, but the access to its public api was shut down which is why it could not be further considered.

⁶ https://en.wikipedia.org/wiki/Wikipedia:Size_comparisons#Wikipedia

⁷ Radford et al. (2019).

2.1.1 Extraction Process

The two usual ways of scraping Wikipedia articles and other data like article metadata or media files from the Wikimedia Foundation⁸ are through the APIs or the database dumps provided by the MediaWiki platform⁹. Both of these channels were used for different purposes:

- **Database Dumps:** Titles, clear text and article ids were extracted using a python library called WikiExtractor¹⁰. This tool parses the compressed and xml-formatted dumps and extracts the aforementioned fields. The output is stored in “.jsonl” (json lines) files, where each line denotes a complete JSON object.
- **Api:** The MediaWiki Api was used to retrieve metadata such as pageviews, edits, the latest edit timestamp and namespaces (categories) linked to each article via the extracted page ids, but also to access the latest news listed in WikiNews to generate exemplary news messages.

The extension of text samples by metadata was made in order to examine the following hypotheses:

1. The detection accuracy varies (significantly) across different categories.
2. The higher the edits and/or views on a page are, the lower the detection accuracy will be as the human text will be more sophisticated and better worded.
3. The more recent the last edit timestamp of an article is the lower the detection rate will be as newer information will be less likely to be present in the training data used for GPT-2.

It should be noted that only articles with a minimum length of 1000 characters were taken into consideration in order to filter out many entries that would have diminished the quality of the dataset (e.g. entries that only have a redirect notice to another article). For each article entry the last edit timestamp, the aggregated amount of page views over the last 30 days and the namespace was retrieved.

2.1.2 Generation Parameter Combination

The output style of GPT2 differs depending on the chosen parameter combination. The parameters that influence the produced output the most are input length, maximum output length, repetition penalty and temperature.

In order to determine a parameter combination that generates convincing text all possible permutations between the parameters given the values listed in table 1 were used to generate samples for 50 articles. As a metric for evaluation a fine-tuned large (1.5GB weights size as a .pt file) RoBERTa-based model with a mixture of temperature 1 and nucleus sampling outputs was chosen. This configuration was elected as it generalizes well to outputs generated using different

⁸ <https://wikimediafoundation.org/>

⁹ <https://www.mediawiki.org/wiki/MediaWiki>

¹⁰ <https://github.com/attardi/wikiextractor>

Table 1: Value ranges for all modified parameters. The sampling method indices refer to nucleus sampling (1), top k (2) and beam search (3).

Parameter name	Value range		
max. input length (in chars)	20	40	60
max. output length (in tokens)	50	100	200
temperature	0.7	1.0	1.3
repetition penalty		1.0	1.3
sampling method	1	2	3
number of beams	5	7	

Table 2: Exemplary word and character split inputs and according outputs

word split input	
Input	tennis is a sport
Output	tennis is a sport which requires a lot of energy and runs at a very fast pace. Getting tired can really impair your performance.

character split input	
Input	tennis is a spo
Output	tennis is a spoilt brat by all accounts, but can he work with the likes of Krajicek? Can he become a player on the big stage who can keep his cool with the big boys?

sampling methods [citation gpt2 paper]. On top of that, human evaluation was performed while reading through samples created by the best performing parameter combinations according to the RoBERTa evaluation.

The main findings were that a higher repetition penalty as well as shorter output length were key factors for better text generation. In addition to that feeding the LM word split input sentences as opposed to character split input sequences also improved quality (table 2).

The finally selected parameter values were:

Maximum Input Length - 60 characters

The input that was fed into GPT-2 XL was the first sequence of plain text in a Wikipedia article, i.e. no infobox or content table text was considered. Additionally, before feeding the 60 character String into the LM it was split at the last word. This was done because the quality of the generated output tended to increase when input was given in full words rather than in characters and thus having many times split words.

Maximum Output Length - 50 Tokens

This corresponds to an average of about 240 - 260 characters per text (in comparison: the max tweet length is 280 characters). As this is a size that occurs a lot especially in social media or breaking news headlines (with the subtitle), the focus was placed on shorter text snippets [cite techcrunch and socialreport]. [insert figure that shows average character length per platform that achieves the best ‘virality’]

Temperature - 1.0

Both lowering and increasing (to 0.7 and 1.3) the temperature led to an increased detection of synthetically generated text, which is why the temperature was left at its original value.

Repetition Penalty - 1.3

As repetitions were not desired in the generated output the repetition penalty was increased from its default value of 1.0 to 1.3. Under the assumption that most Wikipedia articles do not contain repetitions (especially in their abstracts) unlike for instance dramatic text, where text repetition can be used as a stylistic device.

Sampling method - beam search (3)

Number of beams - 5

Initially thought of as a parameter that would have a big impact on the quality of the generated text, it was found out that altering the number of beams used in the beam sampling strategy when predicting the next most probable tokens had only little impact on the LM text generation. The values chosen ranged between 5 and 10 as this is currently the de facto standard [cite Stanford lectures] in research.

2.1.3 Data Building

After downloading the Wikipedia dumps with a compressed size of 17GB the natural text was extracted via the *wikiextractor*¹¹ library (filtering out texts with a character length less than 1000) and the metadata was extracted and parsed using *SAX Parser*¹² ("Simple Api for XML"). In order to finalize the data set creation, synthetic text had to be generated for each article's first 60 characters twice. The double generation was performed in order to select the sample that GPT-2 felt more confident on and improve the data set quality.

Given the aforementioned parameter configurations, the first pipeline runs took place in Google Colaboratory¹³, a free jupyter notebook cloud hosting platform which provides users with a GPU, in this case an Nvidia Tesla K80¹⁴. For the generation of synthetic text, the large GPT-2 model (774M parameters) was used. The generation of each text snippet took about 11 seconds. As soon as the pipeline was fully functioning the jupyter notebooks were converted into python modules and the code was transferred to a Google Cloud¹⁵ deep learning virtual machine¹⁶ (vm) with 13GB of memory, a 100GB standard persistent disk, and 2 Nvidia K80 GPUs. Google Cloud was elected for the reason of providing users with a free credit of 300\$, while other cloud computing

¹¹ <https://github.com/attardi/wikiextractor>

¹² <https://docs.python.org/3/library/xml.sax.reader.html>

¹³ <https://colab.research.google.com/notebooks/intro.ipynb>

¹⁴ <https://www.nvidia.com/en-gb/data-center/tesla-k80/>

¹⁵ <https://cloud.google.com/>

¹⁶ https://console.cloud.google.com/marketplace/details/click-to-deploy-images/deeplearning?_ga=2.99839453.-2121832178.1549923708

providers such as AWS¹⁷ or Azure¹⁸ only provide a free credit of 100\$. Furthermore, the specialized deep learning vm was chosen as it is configured to support common GPU workloads out of the box and removes the task of setting up a high-performance computing environment by having all the needed libraries (e.g. PyTorch, CUDA) and corresponding drivers preinstalled and with a cost of 295.20\$ per month also being less expensive than choosing the same hardware configuration oneself with a cost of about 500\$ per month.

The creation of the data points was parallelized and the total time needed for the dataset creation was approximately 14 days. The data points and log files were stored daily in a Google Storage Bucket¹⁹ and could then be downloaded into a local machine for manual inspection.

The finalized JSON format of a datapoint is shown in Listing 1. The key "label" indicates whether the text entry was created by a human (0) or by GPT-2 (1).

Listing 1: Example of a data point

```
1 {
2     "meta": {
3         "id": "565318",
4         "input": "Lost in Space is a 1998 American science-
5             fiction adventure",
6         "pageviews": 57062,
7         "touched": "2020-02-26T06:53:05Z",
8     },
9     "label": 0,
10    "title": "Lost in Space (film)",
11    "text": "Lost in Space is a 1998 American science-fiction
            adventure film directed by Stephen Hopkins, and starring
            William Hurt, Matt LeBlanc, and Gary Oldman. The plot is
            adapted from the 1965-1968 CBS television series \"of the
            same name\". Seve"
```

2.2 Approach

2.3 Infrastructure

¹⁷ <https://aws.amazon.com/>

¹⁸ <https://azure.microsoft.com/en-us/features/azure-portal/>

¹⁹ https://cloud.google.com/storage/docs/json_api/v1/buckets

References

- Chowdhury, Gobinda G. (2003): *Natural language processing*, In: Annual Review of Information Science and Technology, 37 (1), pp. 51–89 <https://asistdl.onlinelibrary.wiley.com/doi/abs/10.1002/aris.1440370103>.
- Gatt, Albert & Krahmer, Emiel (2017): *Survey of the State of the Art in Natural Language Generation: Core tasks, applications and evaluation*, In: CoRR, abs/1703.09902 <http://arxiv.org/abs/1703.09902>.
- Jozefowicz, Rafal et al. (2016): *Exploring the Limits of Language Modeling*,.
- Klopfenstein, Lorenz Cuno et al. (2017): *The Rise of Bots: A Survey of Conversational Interfaces, Patterns, and Paradigms*, In: Proceedings of the 2017 Conference on Designing Interactive Systems, New York, NY, USA Association for Computing Machinery, DIS '17 <https://doi.org/10.1145/3064663.3064672>, ISBN 9781450349222, p. 555–565.
- Pandu, Nayuk (2019): *Understanding searches better than ever before*, No address in <https://www.blog.google/products/search/search-language-understanding-bert/>, Last accessed: 20.02.2020.
- Radford, Alec et al. (2019): *Language models are unsupervised multitask learners*, In: OpenAI Blog, 1 (8), pp. 2–3.
- Rosenfeld, R. (2000): *Two decades of statistical language modeling: where do we go from here?* In: Proceedings of the IEEE, 88 (8), pp. 1270–1278, ISSN 1558–2256.
- Valdivia, A.; Luzón, M. V. & Herrera, F. (2017): *Sentiment Analysis in TripAdvisor*, In: IEEE Intelligent Systems, 32 (4), pp. 72–77, ISSN 1941–1294.

A Appendix

A.1 User groups

A detailed list of all examined Twitter accounts is omitted here for the sake of brevity. User groups were assembled as Twitter lists in the author's profile. Find links to all lists in the following:

1. Celebrity user group: https://twitter.com/_fpeters/lists/celebrities
2. Politician user group: https://twitter.com/_fpeters/lists/us-politicians
3. Company user group: https://twitter.com/_fpeters/lists/fortune-500

A.2 Confusion matrices

Confusion matrices for favorite classification were omitted in the results chapter. They are displayed in the following.

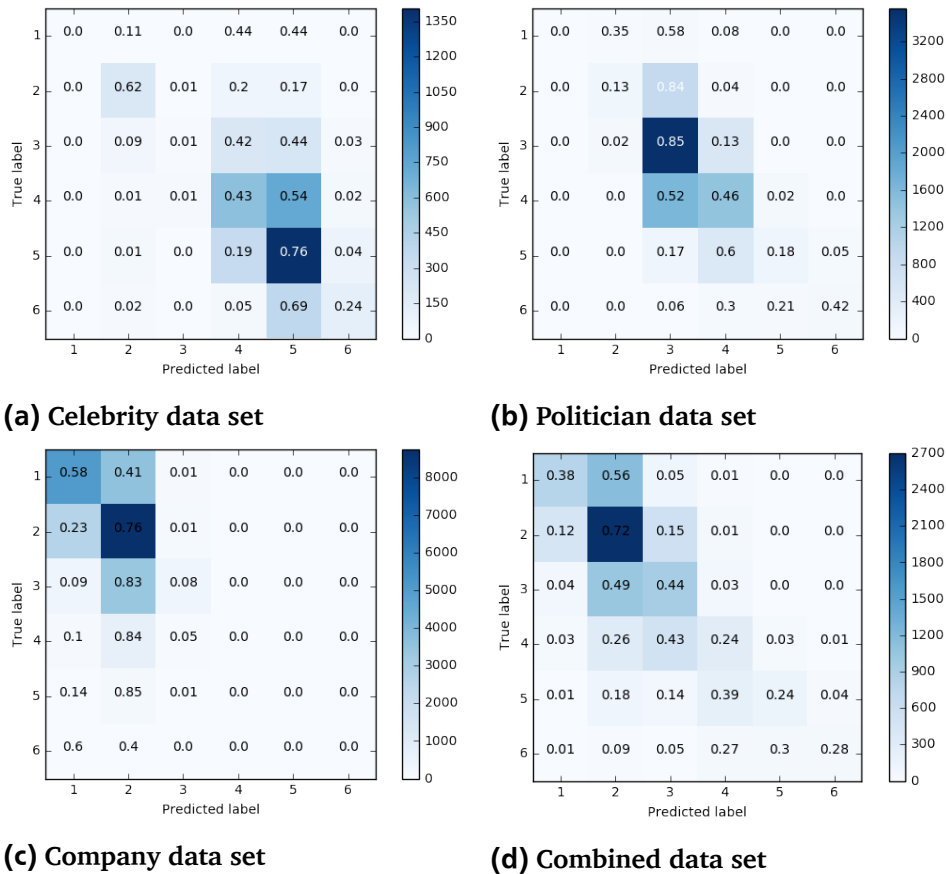
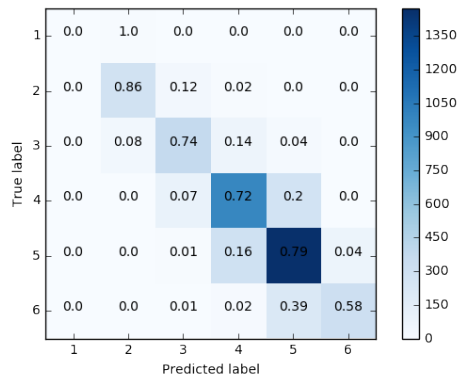
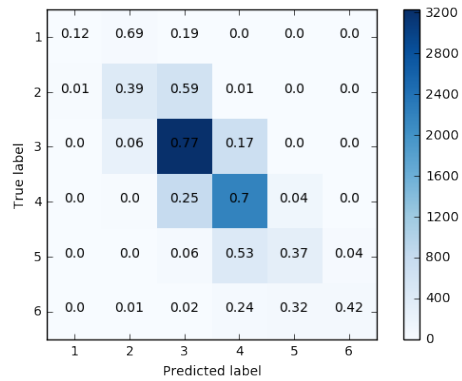


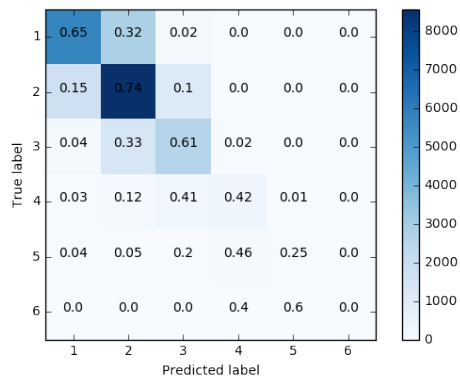
Figure 1: Confusion matrices for linear retweet classification models



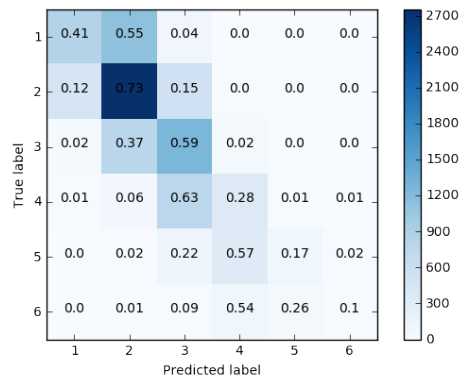
(a) Celebrity data set



(b) Politician data set

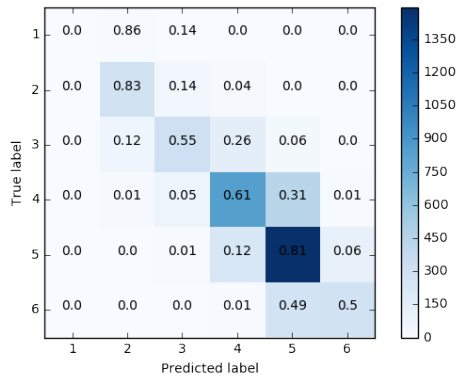


(c) Company data set

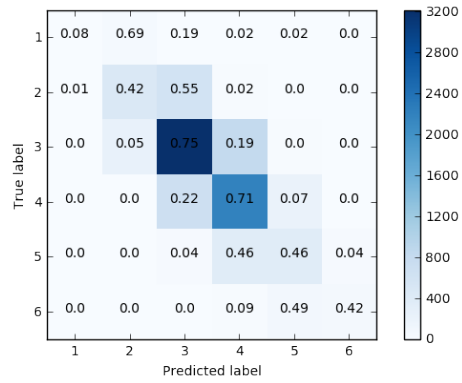


(d) Combined data set

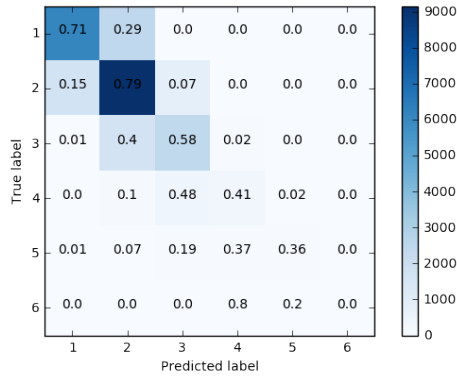
Figure 2: Confusion matrices for deep feedforward classification models



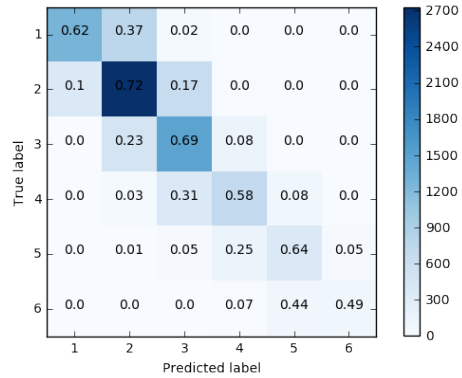
(a) Celebrity data set



(b) Politician data set



(c) Company data set



(d) Combined data set

Figure 3: Confusion matrices for multi-input deep neural networks