

ANALYSING THE EFFECTS OF DATA AUGMENTATION AND
HYPER-PARAMETERS FOR TEXT CLASSIFICATION WITH CONVOLUTIONAL
NEURAL NETWORKS

JONATHAN K. QUIJAS

Department of Computer Science

APPROVED:

Olac Fuentes, Chair, Ph.D.

Monika Akbar, Ph.D.

David Novick, Ph.D.

Charles Ambler, Ph.D.
Dean of the Graduate School

©Copyright

by

Jonathan K. Quijas

2017

to my

FAMILY

thanks for everything

ANALYSING THE EFFECTS OF DATA AUGMENTATION AND
HYPER-PARAMETERS FOR TEXT CLASSIFICATION WITH CONVOLUTIONAL
NEURAL NETWORKS

by

JONATHAN K. QUIJAS

THESIS

Presented to the Faculty of the Graduate School of

The University of Texas at El Paso

in Partial Fulfillment

of the Requirements

for the Degree of

MASTER OF SCIENCE

Department of Computer Science

THE UNIVERSITY OF TEXAS AT EL PASO

May 2017

Acknowledgements

Abstract

Convolutional neural networks have seen much success in computer vision natural language processing tasks. Because of the large number of parameters in these model's, they are prone to overfitting if not regularized appropriately or trained on a sufficiently large dataset. In this work, we study the performance of convolutional neural models when the available training datasets are of moderate size. We quantify the effects of model hyperparameters on classification performance and propose several data augmentation techniques designed for text datasets. Finally, we provide discussion on our empirical results.

Table of Contents

	Page
Acknowledgements	v
Abstract	vi
Table of Contents	vii
Chapter	
1 Introduction	1
1.1 Brief Overview of Deep Neural Networks	1
1.2 Training a Neural Network	1
1.3 Bias-Variance Tradeoff	2
2 Text Classification with Deep Neural Networks	3
2.1 Word Embeddings	3
2.2 Convolutional Neural Networks	3
2.3 Input Representation: Integer Sequences to Word Embeddings	4
3 Towards Improving Model Performance: Regularization and Data Augmentation	5
3.1 Model Regularization	5
3.1.1 Weight Regularization	5
3.1.2 Dropout	6
3.2 Dataset Augmentation: Shuffling and Noise Injection	6
3.2.1 Shuffling	6
3.2.2 Noise Injection	7
3.3 Dataset Augmentation: Padding	8
3.3.1 Wrap Padding	9
3.3.2 Keyword Padding	10
3.3.3 Proposed Padding Scheme	12
3.4 Dataset Augmentation: Padding with Similar Keywords	12

3.5	Reducing Data Granularity: From Abstract to Sentences	12
3.5.1	Voting Schemes	13
4	Model, Dataset, and Final Pipeline Description	14
4.1	Layer Descriptions	14
4.1.1	Embedding Layer	14
4.1.2	Feature Maps: Convolution + Pooling	14
4.1.3	Gated Recurrent Unit Layer	15
4.2	Model Description	15
4.3	Model Hyper-Parameters	16
4.4	Dataset Description	16
4.5	Final Pipeline	16
5	Experimental Results	17
5.1	Dataset Descriptions	17
5.2	Results	17
5.3	AstroPhysics	17
5.3.1	Baseline	18
5.3.2	Padding	18
5.3.3	Sentence Split	19
5.3.4	Noise Injection	19
5.4	Math	19
5.4.1	Sentence Split	20
5.4.2	Padding	20
5.4.3	Noise Injection	20
6	Concluding Remarks	21
6.1	Significance of the Result	21
6.2	Future Work	21
	References	22
	Curriculum Vitae	24

Chapter 1

Introduction

1.1 Brief Overview of Deep Neural Networks

Deep convolutional neural networks have seen an enormous amount of success on a wide array of application, from scene interpretation to self-driving vehicles and art generation[CITE]. Natural language processing tasks are no exception to the range of problems deep learning can solve, from sentiment analysis to language modeling. In this work, we focus our efforts to studying convolutional neural networks for text classification. Specifically, we analyse how well modern neural networks models performed using scientific abstract text data from multiple disciplines such as astro-physics and computer science. The term modern in this context refers to neural models that use popular and recently (re)discovered techniques to achieve state-of-the-art performance on most machine learning benchmark datasets[CITE].

1.2 Training a Neural Network

A neural network is a function $f(\mathbf{x}; \boldsymbol{\theta})$ that maps its input \mathbf{x} to some response variable y . When we *train* a neural network, we *learn* the model parameters, or weights, $\boldsymbol{\theta}$ that minimize some cost function $J(\boldsymbol{\theta})$. For a regression task, where the model's output is a continuous variable, a common cost function is the **Mean Square Error**:

$$J(\theta) = \frac{1}{m} \sum_{i=1}^m (y_i - f(\mathbf{x}_i; \boldsymbol{\theta}))^2$$

For categorical or discrete output variables found in e.g. classification tasks, we use the

Categorical Cross-Entropy:

$$J(\boldsymbol{\theta}) = -\mathbb{E}_{\mathbf{x}, y \sim \hat{p}_{data}} \log p(y|\mathbf{x}; \boldsymbol{\theta})$$

Given a *training* set of observations \mathbf{x}_i and their true labels y_i , we compute weights that minimize the cost, or error, via maximum likelihood (ML) estimation:

$$\boldsymbol{\theta}_{ML} = \underset{\boldsymbol{\theta}}{\operatorname{argmax}} \sum_i^m \log P(y_i|\mathbf{x}_i; \boldsymbol{\theta})$$

,

which one can see is the equivalent of computing the weights that *minimize* the cross-entropy cost function.

1.3 Bias-Variance Tradeoff

When learning a neural network's weights, we use a *training set* so that we can later generalize previously unseen data with high accuracy (or some other determined metric). That means that during training, we obtain $\boldsymbol{\theta}_{ML}$ by minimizing $J_{train}(\boldsymbol{\theta})$, but we care about having low $J_{test}(\boldsymbol{\theta})$ i.e. low cost on test data points.

Overfitting occurs when a network is able to predict its training set extremely well i.e. very close to zero error, but fails to predict unseen data points. This is because the network's weights have been extremely fine-tuned to *fit* its training data, but do not fit or represent data points outside of its training sample. An overfitted model is said to have large **variance** and small **bias**. Conversely, underfitting occurs when the model fails to predict the training set because it generalizes too harshly. This model is said to have large bias and small variance.

Because of the commonly large amount of weights in deep convolutional networks, it is easy to overfit even a moderate size training set[CITE]. In this work, we study the effects of multiple regularization techniques used to avoid overfitting in a neural network[REFER TO CHAPTER].

Chapter 2

Text Classification with Deep Neural Networks

2.1 Word Embeddings

A very common and simple vector representation for words is the one-hot representation. The length of a one-hot vector is the size of the data vocabulary i.e. how many distinct words are found in our data set. For any word, its one-hot vector is zeros everywhere except for a 1 at the word's index. This representation, although simple, fails to capture any meaning other than an identifier. Neural language models i.e. a language model learned using a network, learn to represent words as continuous, dense vectors. These dense, continuous vector representations are commonly called word embeddings. Due to the of the underlying algorithm used to learn these word embeddings, similar words tend to lie closer to each other on embedding space. Because of this, word embeddings are said to capture semantic relations, and thus encode more information than just a word identifier.

[SHOW EMBEDDING PROJECTION]

2.2 Convolutional Neural Networks

Convolutional neural networks are known for their abilities to learn high-level features from raw data. As input signals advance forward through the network, they produce latent signals as linear combinations with learned parameters, have non-linearities applied to them, and have a pooling or selection mechanism based on simple functions such as the average

or maximum operations.

When dealing with image data, images are convolved with multiple filters, each convolution applied to overlapping subimages called as receptive fields. This localized convolution process leads to discovery of low level features of images in the training set such as edges. As data flows forward through the model, higher level features are discovered e.g. wheels or headlights in a vehicle image dataset.

These networks are comprised of *feature maps*. A feature map is a **convolution** layer paired with a **pooling** layer afterwards. The convolution stage creates *activations*, whereas the pooling stage reduces dimensionality and creates translation invariance[CITE].

[INCLUDE IMAGE]

We can take advantage of word embeddings to apply convolutions to text in a fashion similar to convolutions with image data. We apply the convolutions to overlapping subregions of the input text i.e. bi-grams, tri-grams, etc. After convolving these individual and overlapping subregions, we apply a non-linear function and reduce data dimensionality by pooling.

[DEVELOP]

2.3 Input Representation: Integer Sequences to Word Embeddings

Given a set of texts $\mathbf{s}_1, \dots, \mathbf{s}_m$, we build a vocabulary \mathbb{V} and a bag-of-words model from \mathbb{V} to create a mapping $BoW : \mathbb{V} \mapsto \{1, \dots, |\mathbb{V}|\}$. We represent a training text \mathbf{s}_i as a sequence of integers $BoW(\mathbf{s}) = \mathbf{x} = x_1, \dots, x_k$, each integer being simply a word index in \mathbb{V} .

Chapter 3

Towards Improving Model Performance: Regularization and Data Augmentation

3.1 Model Regularization

Regularization is any approach used to reduce test error but not training error. [ELABORATE]

3.1.1 Weight Regularization

One way to regularize a model to impose a constraint on its weights. By adding a penalty term to the cost function, we can shrink weights to avoid them from blowing up in magnitude and overfitting to the training set. L2 regularization adds the L2 norm of the weights to the cost function:

$$J(\boldsymbol{\theta}) = -\mathbb{E}_{\mathbf{x}, y \sim \hat{p}_{data}} \log p(y|\mathbf{x}; \boldsymbol{\theta}) + \lambda \|\boldsymbol{\theta}\|_2$$

where λ is the regularization control parameter. Higher values of λ penalize more and a value of 0 leads to no regularization.

If the embedding layer's weights are allowed to be fine-tuned during training, it is important to regularize them. Because of the large amount of parameters in our model's embedding layer (number of words times embedding size), overfitting occurs quite rapidly

if the weights are allowed to grow without constraint.

3.1.2 Dropout

A more recent and highly effective way to regularize a neural network is via **dropout**. Dropout "deactivates" a neuron or unit with probability p_{drop} . We deactivate a unit by setting its output to 0. This forces the network to not rely on certain activations, since any unit will only be *active* with probability $1 - p_{drop}$ [CITE].

[IMAGE]

3.2 Dataset Augmentation: Shuffling and Noise Injection

Another way to improve model performance is to **augment** the dataset. Data augmentation refers to any transformation of the input data in a way that the label value does not change. Data augmentation is ubiquitous in computer vision tasks. Example transformations include translations, rotations, and even color intensity jitters via Principal Component Analysis[CITE]. All these transformations should be subtle enough that the overall structure is preserved, but allow for the model to process more distinct training data points.

In order to introduce more variance into our dataset, we propose to **shuffle** and **add noise** to the input sequences. We shuffle by randomly permuting non-overlapping neighborhoods in the input text sequence i.e. **context windows**. We further propose to add small amounts of noise to the input sequence by randomly replacing words with words taken from the training vocabulary.

3.2.1 Shuffling

We propose to augment our dataset by making small *context* changes that don't change the global structure of the input sequences. Concretely, we move a non-overlapping *context*

window along the input sequence, randomly shuffling the words indexes inside it. For example using a context window of size 2 i.e. *bi-gram*, an input sequence

$$\mathbf{x} = x_1, x_2, x_3, x_4, x_5, x_6$$

could be shuffled to:

$$\mathbf{x}' = x_2, x_1, x_3, x_4, \underbrace{x_6, x_5}_{\text{bi-gram}}$$

In the first pass, the first two indexes get shuffled. The second pass led to no changes in ordering. The third and final pass shuffles the last two indexes.

The motivation for this dataset augmentation technique is that small changes in the ordering of the words will simulate a larger training sample; a training input sequence has a low chance of being repeated as its length increases. Smaller context windows preserve the most structure. This method will preserve most of the original context structure in the sequence, as long as the shuffling is not too harsh e.g. a complete random permutation of the sequence.

3.2.2 Noise Injection

One common augmentation technique for image datasets is to add small amount of noise. By replacing a relatively small amount of pixel values with noise, the model gets to process and train on a different but similar instance. The noise should be subtle enough as to not distort the image too much, otherwise we could potentially train the model using mostly noise. We propose to augment our dataset by injecting small amounts of noise to each training sample. Again, we aim to simulate a larger training sample while avoiding harsh changes to the original inputs. Concretely, each word in a text sequence gets replaced with a specified probability with a word randomly chosen from our training vocabulary.

Algorithm 1 Add noise to input sequence

```
1: procedure NOISEINJECTION( $\mathbf{x} = [x_1, \dots, x_n], \mathbb{V}, p_{noise}$ )  
2:   for  $k = 1$  to  $n$  do  
3:     if  $p_k \sim U(0, 1) \leq p_{noise}$  then  
4:        $x_k \leftarrow x'_k \in \mathbb{V}$   
5:     end if  
6:   end for  
   return  $\mathbf{x}$   
7: end procedure
```

3.3 Dataset Augmentation: Padding

Although more complex neural models are designed to cope with variable length input, in practice a more common and simple approach is to pad data to be of some specified length as described in chapter [CITE CHAPTER]. In order to enforce uniform input size for our neural networks, we apply **zero-padding**. For any arbitrary training instance $BoW(\mathbf{s}) = \mathbf{x} = x_1, \dots, x_k$, we enforce that $k = n$, for the specified input size n . Thus, if $k < n$, we transform it into $\mathbf{x}_{pad} = x_1, \dots, x_k, 0_{k+1}, \dots, 0_n$. Conversely, if $k > n$, we simply truncate \mathbf{x} to be of size n . The input length introduces another model hyper-parameter that should be fine-tuned, but a reasonable approach is to pad enough to fully accomodate the length of most input sequences i.e. its preferable to pad than truncate and lose information.

Our network's input is a sequence of integers, each integer being a word index: a number representing a word in our vocabulary \mathbb{V} . Word indexes range from 1 to $|\mathbb{V}|$, and the index 0 being left for out-of-vocabulary words. When we pad our input sequence with 0s, we don't add any additional information; we simply create a constant input length. We propose that if instead we add values that characterize or help describe the input sequence more thoroughly, we may increase the amount of useful information available during neural network training.

Consider a very simple input sequence, and assume its true label can be determined from a single word:

”This paper is about computer graphics”

where the label is ”Graphics” i.e. a publication about computer graphics. In this simple case, the label is determined by the word ”graphics”. A single, **very informative** word is enough for the classifier to predict the label correctly. This very informative word however is only 1/6 of the entire text.

Now consider a the padded version, where we enforce an input length of 10:

”This paper is about computer graphics PAD PAD PAD PAD”

In the padded version, the word graphics is now only 1/10 of the entire text. If we could find this very informative word from within the text and pad using it instead of some meaningless token e.g. 0’s, we could increase the information of this word over the entire input sequence by making it comprise a larger amount of the input sequence. In other words, make important words be present more frequently.

We propose to pad input sequences with **meaningful** values found already within the input text sequence instead of 0’s only. In the following section, we will develop on how to extract these meaningful words via a cosine similarity-based metric and apply them to pad our data. In the following subsections, we propose several padding schemes in order reduce the amount of non-informative when training.

3.3.1 Wrap Padding

Our first padding scheme is to *wrap around* the text, repeating words once we reach the padding portion. Referring back to the first example, the input sequence:

”This paper is about computer graphics PAD PAD PAD PAD”

would then be

"This paper is about computer graphics This paper is about"

One thing to observe is that using this simple scheme, we remove all 0's i.e. non-informative padding indexes from the text sequence, but we also don't have selection mechanism and can miss useful words such as missing the word "*graphics*" in this example. Nonetheless, it is a simple first approach to pad our data.

3.3.2 Keyword Padding

One characteristic of the wrap padding scheme described above is that there is no selection mechanism. Words are not selected based on some metric, but rather just *wrapped*. We now propose a simple algorithm to select *keywords* from within the input text sequence based on a cosine similarity based metric, as well as a padding scheme using those keywords.

Using word embeddings, word tokens i.e. strings can now be represented as dense, continuous vectors of dimensionality much smaller than vocabulary size. Another useful property of embeddings is that semantic relationships are captured in this embedding space and can be measured between any two words using the cosine similarity metric[CITE]. This is an attractive property that allows for term comparison by semantic relationships. These relationships can range from complete similarity i.e. the same term, to term orthogonality i.e. unrelated words, and even capture the direction of the relation e.g. negative values of cosine similarity.

The cosine similarity between two vectors can be computed as follows:

$$\frac{\mathbf{x} \cdot \mathbf{y}}{\|\mathbf{x}\| \|\mathbf{y}\|}$$

We can therefore quantify similarity between words in a sentence as a function of their cosine similarity with respect to each other. We refer to this measure as the *within-sentence* cosine similarity. Similarly, we can measure this similarity for words within a noun phrase in a sentence. We call this the *within-noun-phrase* cosine similarity. We represent a noun phrase

\mathbf{np} as a matrix composed of horizontally stacked vectors of size d , each vector a word in the noun phrase:

$$\mathbf{np} = \begin{bmatrix} \mathbf{np}_1 \\ \vdots \\ \mathbf{np}_n \end{bmatrix}, \mathbf{np}_i \in \mathbb{R}^{1 \times d}$$

Given a noun phrase \mathbf{np} comprised of n unit norm embeddings \mathbf{np}_i , we can compute its *within noun phrase* cosine similarity by:

$$c = \frac{1}{n^2} \sum_{i=1}^n \sum_{j=1}^n \begin{bmatrix} \mathbf{np}_1 \\ \vdots \\ \mathbf{np}_n \end{bmatrix} \times \begin{bmatrix} \mathbf{np}_1 \\ \vdots \\ \mathbf{np}_n \end{bmatrix}^T$$

Thus, given a set of noun phrases $\{\mathbf{np}^{(1)}, \dots, \mathbf{np}^{(m)}\}$ from a sentence, we compute the corresponding set of *within noun phrase* cosine similarities $\mathbf{c} = \{c_1 \dots c_m\}$, and we select the noun phrase which corresponds to minimum

$$\mathbf{np}^{(j)}, j = \operatorname{argmin} \mathbf{c}$$

The reasoning behind the cosine similarity minimization is as follows: An informative noun phrase within a sentence will contain a word which will stand out from the rest of the words in the noun phrase. We quantify this notion of *standing out from the rest* using the mean within-noun-phrase cosine similarity. We therefore compute, for all noun phrases in a sentence, the mean *within-noun-phrase* average cosine similarity, as described in [REFER TO ALGORITHM]. We select the noun phrase that minimizes this measure.

In practice, we find both maximization and minimization of the mean *within-nounphrase* average cosine similarity leads to interesting results i.e. reasonable words to be considered keywords, so we include them both for padding. [INCLUDE SNAPSHOT OF CSO OUTPUT]

3.3.3 Proposed Padding Scheme

We propose to pad an input sequence $\mathbf{x}_{pad} = x_1, \dots, x_k, 0_{k+1}, \dots, 0_n$ using the indexes of extracted keywords $x_{key^1}, \dots, x_{key^j}$ by concatenating the non-zero entries x_1, \dots, x_k with the extracted keyword indexes:

$$\mathbf{x}_{cso} = x_1, \dots, x_k, x_{key^1}, \dots, x_{key^j}, 0_{k+j+1}, \dots, 0_n$$

3.4 Dataset Augmentation: Padding with Similar Keywords

Our CSO padding scheme can only pad with words already present in the input text sequence. It could prove beneficial to further pad the input with words that are *similar* to its keywords i.e. words extracted via CSO. We therefore further pad the input sequence with the nearest neighbor of each keyword. In this context, we consider the nearest neighbor of a word embeddings \mathbf{w} as:

$$\min_{\mathbf{w}_i} 1 - \frac{\mathbf{w} \cdot \mathbf{w}_i}{\|\mathbf{w}\| \|\mathbf{w}_i\|}, \mathbf{w} \neq \mathbf{w}_i$$

We use a Locality Sensitive Hashing forest for our approximate nearest neighbor search[CITE]. After having computed the nearest neighbor of each keyword, we pad as follows:

$$\mathbf{x}_{lsh} = x_1, \dots, x_k, x_{key^1}, \dots, x_{key^j}, x_{lsh^1}, \dots, x_{lsh^j}, 0_{k+j+1}, \dots, 0_n$$

3.5 Reducing Data Granularity: From Abstract to Sentences

Our last proposed dataset augmentation scheme is to consider the data as a set of sentences rather than abstracts.

3.5.1 Voting Schemes

Here we introduce our proposed voting schemes

Chapter 4

Model, Dataset, and Final Pipeline Description

In this chapter we describe our model architecture. We start by describing our choice of layers, followed by our network’s final architecture. We conclude the chapter by describing our dataset.

4.1 Layer Descriptions

4.1.1 Embedding Layer

This layer maps an word index, or integer, into its corresponding embedding. This is a layer in the neural network because the embeddings can be further fine-tuned during training. Because of the large number of parameters in this layer (number of words allowed times embedding size), we add a L_2 regularization. Given input text s , this layer receives as input a sequence of word indexes $BoW(\mathbf{s}) = \mathbf{x} = x_1, \dots, x_k$ and maps each word into an embedding. Thus, for an input text \mathbf{s} , we transform it into a sequence of integers \mathbf{x} , and from there into $\mathbf{E} \in \mathbb{R}^{n \times d}$, where d is the embedding size.

4.1.2 Feature Maps: Convolution + Pooling

We refer to a pair of convolutional layer followed by a pooling layer as a feature map[CITE].

4.1.3 Gated Recurrent Unit Layer

4.2 Model Description

The network's general architecture is as follows:

- Input layer: word index vector
- Embedding layer: maps word index vector to embedding matrix
- Dropout layer
- Feature Map 1:
 - Convolution layer: number of kernels:32, activation: rectified linear
 - MaxPooling layer
- Dropout layer
- Feature Map 2:
 - Convolution layer: number of kernels:32, activation: rectified linear
 - MaxPooling layer
- Dropout layer
- Gated Recurrent Unit layer
- Dropout layer
- Dense layer: output size: number of classes, activation: softmax

[IMAGE OF ARCHITECTURE]

4.3 Model Hyper-Parameters

- number of training words
- number of kernels
- kernel size
- pool size
- learning rate
- regularization rate
- dropout rate
- maximum sequence length

4.4 Dataset Description

[DESCRIBE AND CITE EMBEDDINGS]

We gathered scientific paper abstracts from the online repository Arxiv.org. We scraped papers for 5 departments: computer science, mathematics, astrophysics, physics, quantitative biology, and quantitative finance.

4.5 Final Pipeline

[INCLUDE PIPELINE DIAGRAM]

Chapter 5

Experimental Results

In this chapter we will show our results.

5.1 Dataset Descriptions

We scraped a scientific publication abstract dataset from Arxiv.org. Concretely, we gathered publications from the physics, mathematics, computer science, quantitative biology, and quantitative finance departments. For each department, we considered each **topic** as a class. For example, when considering computer science publications, we considered the "computational complexity" topic as one class, "artificial intelligence" as another class, and so on. We apply simple preprocessing to the texts by removing non-alphanumeric characters, and converting to lower case characters. [Overlap Figure] [Class Distribution Figures]

5.2 Results

In this section we show our experimental results. [Add tables and most relevant plots]

5.3 AstroPhysics

[Distribution]

5.3.1 Baseline

NWords	Batch	MaxLen	p_{drop}	Kern	Pool	α	λ	Acc
30000	64	200	0.442	3	5	1e-3	1e-4	0.773
80000	64	250	0.059	3	2	1e-4	1e-5	0.760
80000	64	250	0.056	2	2	1e-3	1e-4	0.759

5.3.2 Padding

Wrap Padding

NWords	Batch	MaxLen	p_{drop}	Kern	Pool	α	λ	Acc
30000	64	200	0.442	3	5	1e-3	1e-4	0.773
80000	64	250	0.059	3	2	1e-4	1e-5	0.761
80000	64	250	0.056	2	2	1e-3	1e-4	0.760

CSO Padding

NWords	Batch	MaxLen	p_{drop}	Kern	Pool	α	λ	Acc
30000	64	200	0.442	3	5	1e-3	1e-4	0.774
80000	64	250	0.059	3	2	1e-4	1e-5	0.752
80000	64	250	0.056	2	2	1e-3	1e-4	0.761

LSH Padding

NWords	Batch	MaxLen	p_{drop}	Kern	Pool	α	λ	Acc
30000	64	200	0.442	3	5	1e-3	1e-4	0.768
80000	64	250	0.059	3	2	1e-4	1e-5	0.752
80000	64	250	0.056	2	2	1e-3	1e-4	0.755

5.3.3 Sentence Split

Mean Vote

NWords	Batch	MaxLen	p_{drop}	Kern	Pool	α	λ	Acc
30000	512	50	0.395	3	3	1e-5	1e-4	0.730
60000	512	40	0.051	4	4	1e-5	1e-6	0.7152
20000	512	40	0.062	5	5	1e-4	1e-6	0.685
20000	512	60	0.247	3	5	1e-3	1e-5	0.710
10000	512	60	0.121	5	3	1e-3	1e-5	0.712
70000	512	40	0.223	4	5	1e-5	1e-5	0.690
40000	128	50	0.225	5	3	1e-5	1e-4	0.734
30000	64	40	0.132	3	4	1e-3	1e-4	0.711
40000	1024	50	0.225	5	3	1e-4	1e-4	0.689

5.3.4 Noise Injection

NWords	Batch	MaxLen	p_{drop}	Kern	Pool	α	λ	p_{aug}	p_{noise}	Acc
80000	64	250	0.548	4	5	1e-4	1e-5	0.653	0.195	0.537
70000	64	300	0.423	5	4	1e-5	1e-3	0.107	.004	0.665
70000	32	200	0.512	3	4	1e-3	1e-5	0.179	.373	0.764
80000	128	200	0.185	4	4	1e-5	1e-6	0.973	0.128	0.680
50000	32	250	0.167	4	4	1e-4	1e-4	0.866	0.364	0.761

5.4 Math

[Distribution] [Distribution]

5.4.1 Sentence Split

Mean Vote

5.4.2 Padding

Wrap Padding

NWords	Batch	MaxLen	p_{drop}	Kern	Pool	α	λ	Loss	Acc
80000	64	250	0.548	4	5	1e-4	1e-5	1.818	0.537
2	7	78	5415						

CSO Padding

NWords	Batch	MaxLen	p_{drop}	Kern	Pool	α	λ	Loss	Acc
80000	64	250	0.548	4	5	1e-4	1e-5	1.818	0.537
2	7	78	5415						

LSH Padding

NWords	Batch	MaxLen	p_{drop}	Kern	Pool	α	λ	Loss	Acc
80000	64	250	0.548	4	5	1e-4	1e-5	1.818	0.537

5.4.3 Noise Injection

NWords	Batch	MaxLen	p_{drop}	Kern	Pool	α	λ	p_{aug}	p_{noise}	Acc
80000	64	250	0.548	4	5	1e-4	1e-5	0.653	0.195	1.818 0.537

Chapter 6

Concluding Remarks

6.1 Significance of the Result

6.2 Future Work

References

- [1] S. Cook, “The complexity of theorem-proving procedures,” *Proceedings of the 3rd ACM Symposium on Theory of Computing*, Shaker Heights, Ohio, 1971, pp. 151–158.
- [2] M. R. Garey and D. S. Johnson, *Computers and Intractability: A Guide to the Theory of NP-Completeness*, W. H. Freeman, San Francisco, 1979.
- [3] R. Karp, “Reducibility among combinatorial problems,” in: R. Miller and J. Thatcher (eds.), *Complexity of Computer Computations*, Plenum Press, New York, 1972, pp. 85–103.
- [4] R. B. Kearfott and V. Kreinovich (eds.), *Applications of Interval Computations*, Kluwer Academic Publishers, Norwell, MA, 1996.
- [5] V. Kreinovich, A. V. Lakeyev and S. I. Noskov, “Optimal solution of interval linear systems is intractable (NP-hard),” *Interval Computations*, 1993, No. 1, pp. 6–14.
- [6] V. Kreinovich, A. V. Lakeyev and J. Rohn , “Computational complexity of interval algebraic problems: some are feasible and some are computationally intractable: a survey,” in: G. Alefeld and A. Frommer (eds.), *Scientific Computing and Validated Numerics*, Akademie-Verlag, Berlin, 1996, pp. 293–306.
- [7] V. Kreinovich, A. V. Lakeyev, J. Rohn and P. Kahl, *Feasible? Intractable? On Computational Complexity of Data Processing and Interval Computations*, Kluwer Academic Publishers, Norwell, MA, 1996 (to appear).
- [8] U. Kulisch and W. L. Miranker, *Computer Arithmetic in Theory and Practice*, Academic Press, NY, 1981.
- [9] A. V. Lakeyev and V. Kreinovich, “If input intervals are small enough, then interval computations are almost always easy,” *Reliable Computing*, Supplement (Extended

Abstracts of APIC'95: International Workshop on Applications of Interval Computations), 1995, pp. 134–139.

- [10] L. Levin, “Universal sequential search problems,” *Problems of Information Transmission*, 1973, Vol. 9, No. 3, pp. 265–266.
- [11] R. E. Moore, “Automatic error analysis in digital computation,” *Technical Report LMSD-48421*, Lockheed Missiles and Space Co., Palo Alto, CA, January 1959.
- [12] R. E. Moore, *Interval Analysis*, Prentice Hall, Englewood Cliffs, NJ, 1966.
- [13] A. Neumaier, *Interval Methods for Systems of Equations*, Cambridge University Press, Cambridge, 1990.
- [14] S. G. Rabinovich, *Measurement Errors: Theory and Practice*, American Institute of Physics, NY, 1995.

Curriculum Vitae