

BSP S3 - Speech Recognition for Luxembourgish by a Recurrent Neuronal Network

Friday 20th December, 2019 - 15:30

Le Minh Nguyen
University of Luxembourg
Email: le.nguyen.001@student.uni.lu

Vladimir Despotovic
University of Luxembourg
Email: vladimir.despotovic@uni.lu

Abstract—Recurrent Neural Networks (RNN) use their internal state to operate on data time-series. One of the domains in which RNNs are applied is speech recognition. In this paper, the focus will be on the task of end-to-end isolated word recognition. We apply RNN and feedforward Neural Networks as the baseline for the recognition and compare their performance.

1. Introduction

This paper presents the Bachelor Semester Project (BSP) made by Le Minh Nguyen together with Vladimir Despotovic as his motivated tutor. The project is divided into two parts. The first part implements the end-to-end speech word recognition model using the following concepts:

- Data collection and preprocessing.
- Deep Neural Networks architecture: Deep Feedforward Neural Networks and Recurrent Neural Networks, which are presented in the popular Deep Learning textbook [1].

In the second part, we explain the deep Neural Networks architecture and compare their performance obtained during our classification experiment.

2. Project description

2.1. Domains

- Speech Recognition
- Artificial Neural Networks
- Deep Learning
- Data preprocessing
- Training set and testing set
- Python
- Keras

2.1.1. Scientific. The scientific aspects covered by this Bachelor Semester Project are the concepts of speech recognition and Deep Learning. Different Artificial Neural Networks architectures are investigated for the end-to-end isolated word recognition task.

Speech Recognition. The objective of speech recognition is to map an audio signal which contains a set of spoken natural language expressions to the matching sequence of words produced by the speaker. In the past, Automatic Speech Recognition (ASR) was made up of different modules such as complex feature extraction, acoustic models, language and pronunciation models. [2] Sequential models, such as Hidden Markov Models (HMM), in combination with a pre-trained language model, were used to map sequences of phones to output words. [3] A different approach is to build ASR models end-to-end. With deep learning, it replaces most of the modules with a single module. This alternative method will be the main task of this report.

Artificial Neural Networks (ANN). ANNs are computing systems inspired by the biological brain. These systems are based on a set of connected units called artificial nodes. Each connection can transmit *signals* between units. A unit can process the signal and transmit it to another unit.

Deep Learning. This field deals with learning by decomposing a task's input into smaller and simpler compositions. With Deep Learning, computing systems can build complex concepts from a composition of simpler concepts.

2.1.2. Technical. The technological aspect which is covered in this project is the data collection, feature extraction and implementation of our classification model.

Data preprocessing. Data preprocessing is an important phase in machine learning. It ensures the quality of the gathered data by eliminating irrelevant and redundant information. Data preprocessing contains tasks such as cleaning, instance selection, normalization, feature extraction and selection. We will focus on feature extraction in this paper with the presentation of Mel-frequency cepstral coefficients (MFCCs) in Section 4.2.1 which are used as features extracted from the speech signal.

Training set and testing set. For a computing system to learn from and make predictions on data, a mathematical model is built from input data. This input data used to create the model consists of two datasets: The training

and testing set. The training set contains pairs of an input vector, which represent features and an output vector often called the labels. With the training set, the model learns to map the input vector to the labels. On the other hand, the testing set evaluates how well the model generalizes the prediction over the dataset previously not seen by the model.

Python. This is a programming language which is interpreted, high-level and general-purpose. [4]

Keras Library. *Keras* is a high-level Neural Networks API written in Python. It is designed for easy-to-use and efficient experimentation with Deep Learning. [5]

2.2. Targeted Deliverables

2.2.1. Scientific deliverables. One of the main deliverables is to present how we extract the features from the dataset with Mel-frequency cepstral coefficients (MFFC). Additionally, we present the notions of Deep learning. This paper gives a brief introduction to Artificial Neural Networks (ANN) and their application for classification. Further, we extend this scientific presentation by diving deeper into three different ANNs; Feedforward Neural Networks, Recurrent Neural Networks (RNN) and Long short-term memory (LSTM) as a special case of RNN. Finally, we compare their performances obtained after being trained on the given audio dataset.

2.2.2. Technical deliverables. The other main deliverable for this paper is the implementation of an end-to-end speech recognition model based on the three Neural Networks mentioned in the section above. Additionally, we collect a small dataset of Luxembourgish spoken words $w \in \{'0', \dots, '9', 'moien'\}$ to train our model to recognize these words. To the best of our knowledge, there are no publically available datasets for speech recognition of the Luxembourgish language.

2.3. Constraints

For this BSP, we set the different constraints for the ANNs, for the Speech recognition for the Luxembourgish vocal dataset and the classification model's implementation.

Data set. First, we focus on training our model on a dataset containing English spoken numbers. [6] After having trained the model successfully on this dataset, we train it on the smaller Luxembourgish dataset and analyse its prediction accuracy. Since we aren't able to collect as much voice recordings as the English data set, we should not expect high accuracy.

Speech recognition. Since no datasets with continuous sentences of the Luxembourgish Language exists, we choose to work with isolated word recognition instead of continuous speech recognition. In isolated word recognition,

words are separated by pauses.

Speaker dependent. To simplify the collection of the Luxembourgish dataset, we collect the audio recordings from one person which makes the isolated word recognition speaker-dependent.

ANNs implementation. We do not implement the Artificial Neural Networks architecture since existing deep learning libraries such as Keras are already matured. Thereby, we focus on explaining how the APIs work and how to use them efficiently in our model.

3. Pre-requisites

To start on the project, certain skills in programming and mathematics are required. In particular, the preliminary requirement of the project is as follow:

- Understanding of vector and matrix algebra.
- Introductory course in Python.
- Software development.
- Knowledge of probability and statistics, is desirable although not mandatory.

3.1. Scientific pre-requisites

Linear Algebra. is a sub-field of mathematics, which works with vectors and matrices. Since the input of deep learning is data that are transformed into structures of rows and columns, linear algebra is one of the key foundations of deep learning. It is used to describe the operations of the deep learning algorithms, and implement the algorithms in code. All tasks in deep learning relate to linear algebra, from data preprocessing to the deep learning algorithms. [1]

3.2. Technical pre-requisites

Python. Python is an interpreted, high-level and general-purpose programming language, which is conceived in the late 1980s and released in 1991 by Guido van Rossum. [7] Its design philosophy accentuates readability of the code. As many high-level programming languages, Python is dynamically typed. Further, it supports multiple programming paradigms such as procedural, object-oriented and functional programming. With a proficient background in programming in Python gained during my high school and university education, I am prepared to work on this project.

Software development. Software development is the process of designing and implementation of applications and frameworks. In general, the process of software development includes writing and maintaining the source code which is often a planned and structured process. The software development contains mostly research, prototyping, modification and reuse of existing software. During the technical deliverable section, we use this structured process to design and implement our ANNs.

4. Scientific Deliverables

4.1. Requirements

- **FR01** Present the feature extraction with MFCCs
- **FR02** Investigate the notions of deep learning and Artificial Neural Networks (ANNs)
This should give an introduction to the domain of deep learning. It presents an overview of 3 different ANNs architectures, namely every Neural Network presentation will cover a small introduction and theoretical aspect.
- **NFR01** Performance evaluation and comparison
During this section, we present and discuss the results obtained using three Neural Network architectures.

4.2. Design

4.2.1. FR01: Feature extraction with MFCCs.

Before training our end-to-end ASR model, audio data have to be preprocessed. The first task is to extract features from the data which describes natural language expressions and excludes background noises and emotions.

Mel Frequency Cepstral Coefficients (MFCCs) are features widely used in ASR. MFCCs were presented by Davis and Mermelstein in the 1980s. [8]

To extract the MFCCs from the dataset, the following extraction steps have to be executed on a speech signal sampled at $16kHz$:

1. Frame the signal into $25ms$ frames. The frame length of a $16kHz$ signal:

$$0.025 * 16000 = 400samples. \quad (1)$$

With a frame step of $160 samples$ or $10ms$, the frames overlap themselves, such that the first frame containing $400 samples$ starts at sample 0 and the second frame starts at sample 160. This pattern repeats until the end of the speech signal. If the audio file is not divisible by an even number, it is padded by zeros.

2. For each frame calculate the periodogram estimate of the power spectrum. To obtain the Discrete Fourier Transform (DFT) from the frame i , we perform:

$$S_i(k) = \sum_{n=0}^{N-1} s_i(n)h(n)e^{-j2\pi kn/N} \quad (2)$$

For each speech frame $s_i(n)$, the periodogram-based power spectral estimate is calculated with:

$$P_i(k) = \frac{1}{N} |S_i(k)|^2 \quad (3)$$

This is called the Periodogram estimate of the power spectrum which identifies for every frame which frequencies are present.

3. Apply the Mel filterbank to the power spectral and sum the energy in each filter. The Mel filterbank is a set of 26 triangular filters. The filterbank energies are calculated by multiplying every filterbank with the power spectrum.
4. Take the logarithm of the 26 filterbank energies.
5. Take the Discrete Cosine Transform (DCT) of the 26 log filterbank energies resulting in 26 cepstral coefficients. We only used the lower 12-13 of the 26 coefficients for ASR.

In Section 4.3.1 we will use the Python library *Librosa* to extract the MFCCs features from the dataset.

4.2.2. FR02: Deep Learning and Artificial Neural Networks.

The notions of Deep Learning and examples of ANNs presented in this section are based on the textbook *Deep Learning* [1]. These presentations should only give a high-level introduction to these notions.

4.2.3. Deep Learning.

Artificial Intelligence (AI) is a field with many practical applications such as understanding speech, ... AI deals with challenging problems which are easy to perform but hard to describe formally by humans. Deep Learning is a solution to these intuitive problems. This approach to AI allows computer systems to study from experience and understand the world through a hierarchal stack of concepts. The computing system gathers knowledge from experience which eliminates the need to describe formal rules. The computer understands complex concepts with this stack of concepts by decomposing them with simpler ones. The representation of this hierarchy of concepts shows a deep graph with many layers, hence the name for Deep Learning. [1]

4.2.4. Feedforward Neural Network.

Feedforward neural networks or multilayer perceptrons (MLPs) are very important deep learning models. The objective of MLPs is to approximate a given function f^* . A feedforward network defines a mapping $\hat{y} = f(x; \theta)$ and learns the values of the parameter θ with a given input vector x which yields the best approximation of f^* . This function can be for example a classifier, $y = f^*(x)$ which maps an input x to a class y .

These structures are called **feedforward** since there is a progress of information. In one direction: From input to

output.

MLPs are called **networks** because they are made up of composed functions. Let f_1 , f_2 and f_3 be three function forming a chain $f(x) = f_3(f_2(f_1(x)))$. The chain structures are commonly used in ANNs. In the example, f_1 is called the first layer of the network, f_2 is the second layer, etc. The length of the chain represents the depth of the network. The final layer is also called the output layer.

During the training of the neural network, we try to find $f(x)$ to approximate $f^*(x)$. The training data contains examples of $f^*(x)$ obtained with different training points where every example x is referring to a label $y \approx f^*(x)$. These examples determine how the output layer should approximate a value close to y with each example x . The other layers are called the **hidden layers** because the training examples are not determining their behaviour and it doesn't contain the desired output of each layer. Therefore, a learning algorithm has to determine how to utilize the hidden layers to obtain a good approximation of $f^*(x)$.

The **width** of the feedforward model is defined by the dimensionality of each hidden layer which is usually a vector containing values. These vector layers contain many **units** which process data in parallel. Each unit refers to a vector-to-scalar function, by taking input from other units and calculating its proper activation value.

We will present a simple MLP with one hidden layer which contains hidden units. The hidden layer represents a vector of hidden units h which are calculated by a function $f_1(x; \theta)$. θ consists of W_{input} and b_{input} which are the *weights* and *biases* respectively such that:

$$\begin{aligned} h &= f_1(x; W_{input}, b_{input}) \\ \Leftrightarrow h &= W_{input}^T x + b_{input} \end{aligned} \quad (4)$$

h is then used as input for the second layer which is the output layer; $f_2(h; W_{output}, b_{output})$.

$$\begin{aligned} y &= f_2(h; W_{output}, b_{output}) \\ \Leftrightarrow y &= W_{output}^T h + b_{output} \end{aligned} \quad (5)$$

The network can be illustrated as a chain of two functions:

$$\begin{aligned} f(x; W_{input}, b_{input}, W_{output}, b_{output}) \\ = f_2(f_1(x)) \end{aligned} \quad (6)$$

After using an affine transformation, the hidden units are computed by a nonlinear function g called an activation function:

$$h = g(W_{input}^T x + b_{input}) \quad (7)$$

The common activation function in deep learning is the **rectified linear unit** or ReLU defined as:

$$g(z) = \max\{0, z\} \quad (8)$$

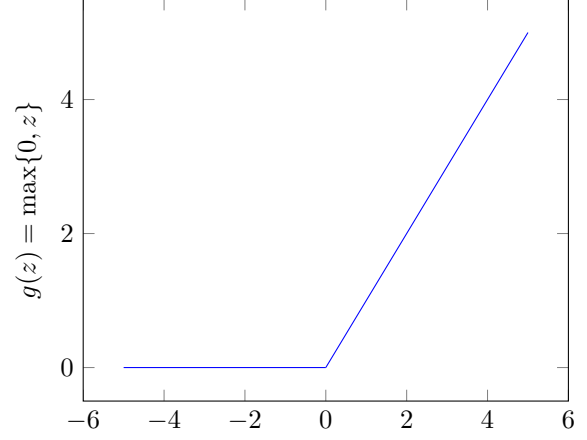


Figure 1. The rectified linear activation function.

The complete feedforward network looks as follows:

$$\begin{aligned} f(x; W_{input}, b_{input}, W_{output}, b_{output}) \\ = W_{output}^T \max\{0, W_{input}^T x + b_{input}\} + b_{output} \end{aligned} \quad (9)$$

A representation of a single perceptron and its position in an MLP is given in Figure 2.

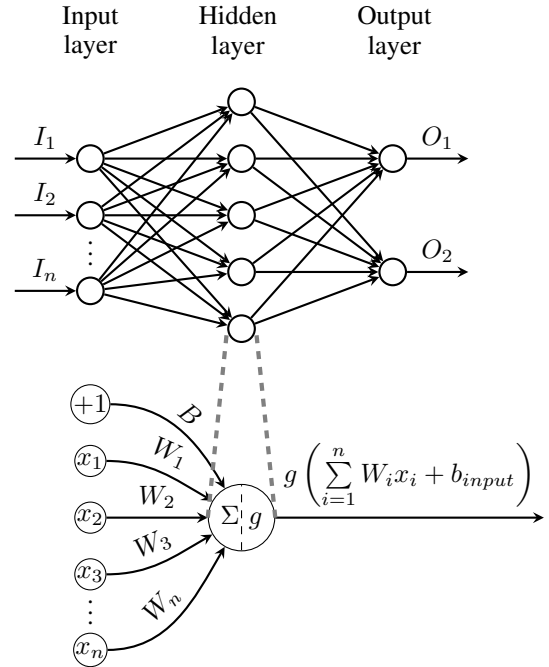


Figure 2. A figure of a single perceptron with its position in a large-scale MLP.

Gradient-Based Learning. During the training of the model, deep learning algorithms solve an optimization problem. The problem is to minimize the error represented by an objective function (loss). With the loss function the deep learning model optimizes the approximation of f^* . An example of a loss function for our MLP can be the mean squared error (MSE):

$$\frac{1}{n} \sum_{i=1}^n (f^*(x_i) - f(x_i; \theta))^2 \quad (10)$$

ANNs usually uses gradient-based optimizers to obtain a minimal loss function.

Back-propagation. We call forward propagation when an MLP takes x as input and computes \hat{y} as output since the information passes through the network. The back-propagation algorithm passes information from the loss function backwards through the network. This algorithm computes the gradient while passing back the network. With this gradient, another algorithm called the *stochastic gradient descent* performs the learning of the *weight* = $\{W_{input}, W_{output}\}$ and the *biases* = $\{b_{input}, b_{output}\}$ and adjusts them to minimize the loss function.

4.2.5. Recurrent Neural Network.

The idea behind RNNs is to make use of sequential information. In a traditional neural network we assume that all inputs (and outputs) are independent of each other. But for many tasks that's a very bad idea. If you want to predict the next word in a sentence you better know which words came before it. RNNs are called recurrent because they perform the same task for every element of a sequence, with the output being depended on the previous computations and you already know that they have a "memory" which captures information about what has been calculated so far.

Recurrent Neural Networks (RNNs) are a particular architecture of ANNs which can process sequential data $S = x_1, \dots, x_\tau$. To explain how RNNs are implemented, we have to introduce the notion of computational graphs. A computational graph is a formal structure representing a set of computations. We obtain a chain of events by unfolding a recurrent computation into a computational graph with a repetitive structure.

Consider a recurrent system,

$$s_t = f(s_{t-1}; \theta) \quad (11)$$

with s being the state of the system. The graph can be unfolded considering the time step $\tau = 3$, we have,

$$\begin{aligned} s_3 &= f(s_2; \theta) \\ &= f(f(s_1; \theta); \theta) \end{aligned} \quad (12)$$

This expression can be illustrated by a directed acyclic computational graph. See figure 3

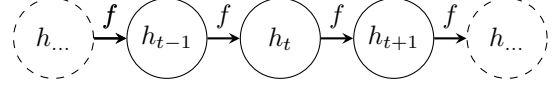


Figure 3. The recurrent system described by equation 12, illustrated as an unfolded computational graph.

Now we consider recurrent system accepting an external signal x_t ,

$$s^t = f(s^{t-1}, x_t; \theta) \quad (13)$$

where the state s is containing information about the past sequence S . The hidden units can be described with equation 13 using h to denote the state,

$$h_t = f(h_{t-1}, x_t; \theta) \quad (14)$$

represented in figure 4.

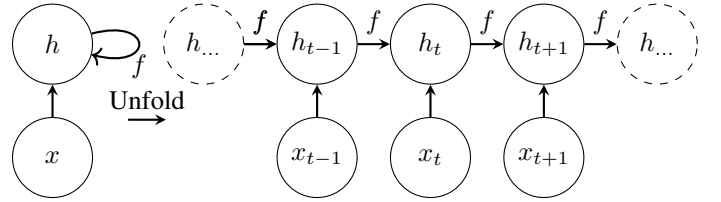


Figure 4. This is an RNN with no outputs, it processes the input x and passes it into the state h which is then passed through time. On the right graph, we have, the same network unfolded as a computational graph. [1]

With graph unfolding, we can now present some common examples of an RNN. Here are some important designs of RNNs:

- RNNs which produce an output at every time step τ and are connected recurrently between hidden units.
- RNNs which produce an output at every time step τ and are connected recurrently only from the output from one time step between hidden units at another time step.
- RNNs which are connected recurrently between hidden units, read an entire data sequence S and produce only one single output.

We pick the RNN represented in figure 5 to develop forward propagation equations. This RNN takes in a data sequence S and outputs for every time step τ an output. In this figure, we didn't include an activation function for the hidden units h . However, for the equations, we assume the activation function \mathcal{H} to be the hyperbolic tangent.

Let h_0 be the initial state. For every time step from $t = 1$ to $t = \tau$, we apply the following equations:

$$a_t = b + W h_{t-1} + U x_t \quad (15)$$

$$h_t = \tanh(a_t) \quad (16)$$

$$o_t = c + V h_t \quad (17)$$

$$\hat{y}_t = \text{softmax}(o_t) \quad (18)$$

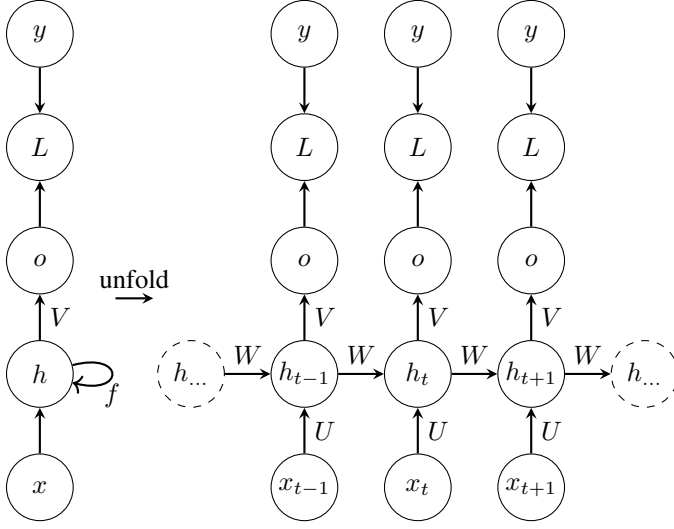


Figure 5. this computational graph is for calculating the training loss of an RNN which labels a sequence of x values to corresponding output o values. A loss function L is calculating the distance from o to the corresponding y target. On the left, we have the network as a recurrent graph. On the right, we have its unfolded computational graph. [1]

where b and c are the bias vectors and U , V and W are the weight matrices.

The Challenge of Long-Term Dependencies.

The challenge of learning long-term dependencies is that gradients, propagated through a very deep RNN, can vanish or explode. Gradients can vanish when the gradients are converging to 0 and therefore the weights won't be updated while learning. Exploding gradients happen if the gradients are too large and the loss function will never reach the optimized minimum. There is a solution to this problem. We will introduce Long Short-term Memory (LSTM) and Gated Recurrent Unit (GRU) RNNs which are an RNN architecture dealing with this challenge. [9]

4.2.6. Long short-term memory.

As a solution to the challenges of RNNs, effective sequence models used in practice are called gated RNNs. These include the long short-term memory (LSTM) and the gated recurrent unit (GRU) based RNNs. Gated RNNs are based on the idea of creating paths through time that have gradients that neither vanish nor explode.

An LSTM unit is composed of a memory cell, an input gate, an output gate and a forget gate. The cell stores information over a long period and the three gates control the flow of information through the cell. RNNs using LSTM units deal with the vanishing gradient problem because these units allow gradients to flow unchanged. [9]

The standard formulation of a single LSTM cell can be given by the following equations:

$$f_t = \sigma(W_f h_{t-1} + V_f x_t + b_f) \quad (19)$$

$$i_t = \sigma(W_i h_{t-1} + V_i x_t + b_i) \quad (20)$$

$$C'_t = \tanh(W_C h_{t-1} + V_C x_t + b_C) \quad (21)$$

$$C_t = f_t C_{t-1} + i_t C'_t \quad (22)$$

$$o_t = \sigma(W_o h_{t-1} + V_o x_t + b_o) \quad (23)$$

$$h_t = o_t \tanh(C_t) \quad (24)$$

where σ is the sigmoid function, \tanh is the hyperbolic tangent function, f, i, C', C, o are the forget gate, input gate, new memory cell content, and memory cell content, output gate respectively. See representation of an LSTM cell in figure 6.

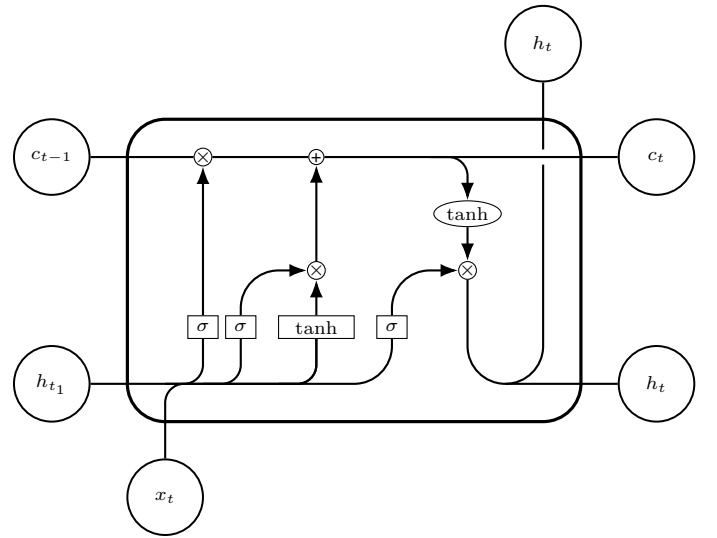


Figure 6. An LSTM unit

4.2.7. Gated Recurrent Unit.

Identical to LSTM units, the gated recurrent unit (GRU) consists of gated units which control the flow of information through the unit without having a dedicated memory cell. Compared to LSTM units, GRUs don't have an output gate. The output gate is responsible to control the amount of memory content seen by the other units in the network. Thus GRU exposes its entire content without control. [10]

The activation of a GRU unit is computed with:

$$h_t = (1 - z_t)h_{t-1} + z_t \tilde{h}_t \quad (25)$$

with an update gate z_t which decides how much the unit updates its content. Its content consists of the previous activation h_{t-1} and the candidate activation \tilde{h}_t . The update gate is given by the equation:

$$z_t = \sigma(W_z x_t + U_z h_{t-1}) \quad (26)$$

This update gate is similar to the one used by LSTM. However, it exposes its state each time step since it does not have a procedure to control how the state is exposed. The candidate activation is given by:

$$\tilde{h}_t = \tanh(Wx_t + U(r_t \odot h_{t-1})) \quad (27)$$

Where r_t is a reset gate and is computed comparably to the update gate:

$$r_t = \sigma(W_r x_t + U_r h_{t-1}) \quad (28)$$

See the illustration of a GRU uni in Figure 7.

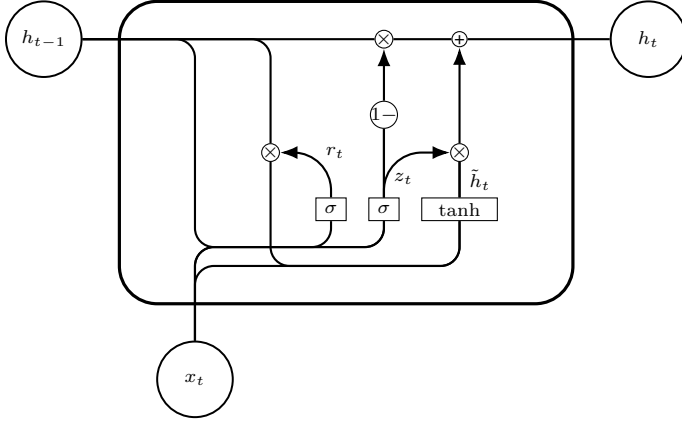


Figure 7. A GRU unit.

4.3. Production

4.3.1. Feature extraction with MFCC in Python.

For this experiment, we use an audio dataset of English spoken words [6]. From this dataset we choose the words $w \in \{'zero', 'one', 'two', \dots, 'eight', 'nine'\}$. This dataset contains many persons pronouncing these different words. Every file is roughly 1s long and is sampled at 16kHz.

We use the Python library *librosa* [11] to extract the MFCCs from the dataset. Instead of processing the whole dataset, the following steps are applied to one single audio file f .

1. Let f be an audio file sampled at 16kHz, *Librosa* loads the song as follows:

```
1 import librosa
2
3 y, sr = librosa.load(f, sr=16000, mono=True,
    duration=1)
```

with y being the audio time series and sr the sample rate of y .

2. The next step is to pad the time series y if it is smaller than 1s:

```
1 import numpy as np
2
3 if y.size < 16000:
4     rest = 16000 - y.size
5     left = rest // 2
6     right = rest - left
7     y = np.pad(y, (left, right), 'reflect')
```

The time series is reflected on both sides to get the correct size.

3. After matching the time series to the correct size, the MFCCs can be extracted with:

```
1 mfccs = librosa.feature.mfcc(y=y, sr=sr)
```

This function *.mfcc* returns an *np.ndarray* with the dimension:

$$shape = (n_mfcc, t)$$

with n_mfcc being the number of MFCCs and t the number of frames.

4. The last step is to flatten this matrix to a vector which can be stored in a data frame for later training.

```
1 row = mfcc.T.flatten()
```

The matrix *mfcc* is transposed and flattened. This concatenates the columns together forming a vector.

Theses steps are then applied to the entire dataset. Each computed vector is appended to the data frame which is going to be stored as a *.csv* file. The entire code snippet is in the Appendix in Figure 8.

4.4. Assessment

4.4.1. NFR01: Performance evaluation.

In this section, the performance of the four ANNs mentioned in Section 4.2.2 are evaluated. We use the English and Luxembourgish dataset for the final performance analysis. Before starting to train and evaluate the four models, the loss optimizer and number of epochs have to chosen. The english dataset is used to designate an optimal optimizer and number of epochs.

The first step is to choose an optimal optimizer. The choice is between using the stochastic gradient descent (SGD) or the Adam optimizer. The second is a combination of two extensions of SGD, Adaptive Gradient Algorithm (AdaGrad) and Root Mean Square Propagation (RMSProp) [12].

After testing both optimizers, the results are the following:

The results using Adam as optimizer yield better accuracies compared to SGD. For the following performance analysis, the Adam optimizer is considered.

The next step is to choose the number of epochs which will be used for the performance analysis. The following

results were made with the number of epochs being between 1 – 60:

The graphs show the loss converges around n epochs. Therefore, we choose $epochs = n$.

Let $epochs = n$ we evaluate each model's performance.

5. A Technical Deliverable

5.1. Requirements ($\pm 15\%$ of section's words)

- **FR01** Implementation of the three classification models
We use the Keras library to implement our models with three different Neural Network architecture
- **FR02** Collect a small dataset of Luxembourgish spoken words
We collect a small dataset containing Luxembourgish audio samples containing the words $w \in \{'0', \dots, '9', 'moien'\}$

5.2. Design ($\pm 30\%$ of section's words)

We explain how to use the Keras Library.

5.3. Production ($\pm 40\%$ of section's words)

5.3.1. Implementation of classification models. We present the implementation of our classification models

5.3.2. FR02: Luxembourgish dataset collection. We present how we collect our Luxembourgish dataset.

5.4. Assessment ($\pm 15\%$ of section's words)

Acknowledgment

I would like to thank my tutor Vladimir Despotovic for his constructive feedback and mentorship. His introduction and explanation of neural networks were outstanding. I would recommend fellow BiCS Students interested in this field to work with Vladimir Despotovic. Additionally, I thank him for supervising my paper.

6. Conclusion

References

- [1] I. Goodfellow, Y. Bengio, and A. Courville, *Deep Learning*. MIT Press, 2016, <http://www.deeplearningbook.org>.
- [2] D. Amodei, R. Anubhai, E. Battenberg, C. Case, J. Casper, B. Catanzaro, J. Chen, M. Chrzanowski, A. Coates, G. Diamos, E. Elsen, J. H. Engel, L. Fan, C. Fougner, T. Han, A. Y. Hannun, B. Jun, P. LeGresley, L. Lin, S. Narang, A. Y. Ng, S. Ozair, R. Prenger, J. Raiman, S. Satheesh, D. Seetapun, S. Sengupta, Y. Wang, Z. Wang, C. Wang, B. Xiao, D. Yogatama, J. Zhan, and Z. Zhu, "Deep speech 2: End-to-end speech recognition in english and mandarin," *CoRR*, vol. abs/1512.02595, 2015. [Online]. Available: <http://arxiv.org/abs/1512.02595>
- [3] W. S. J. Cai, "End-to-end deep neural network for automatic speech recognition," 2015. [Online]. Available: <https://cs224d.stanford.edu/reports/SongWilliam.pdf>
- [4] "Python software foundation, python language reference, version 3.7. available at," <http://www.python.org/>.
- [5] F. Chollet *et al.*, "Keras," <https://keras.io>, 2015.
- [6] P. Warden, "Speech commands: A dataset for limited-vocabulary speech recognition," *CoRR*, vol. abs/1804.03209, 2018. [Online]. Available: <http://arxiv.org/abs/1804.03209>
- [7] "General python faq," <https://docs.python.org/3/faq/general.html#what-is-python>, accessed 19/05/19.
- [8] "Mel frequency cepstral coefficient (mfcc) tutorial," <http://practicalcryptography.com/miscellaneous/machine-learning/guide-mel-frequency-cepstral-coefficients-mfccs/#references>, accessed 02/12/19.
- [9] S. Hochreiter and J. Schmidhuber, "Long short-term memory," *Neural Computation*, vol. 9, no. 8, pp. 1735–1780, 1997. [Online]. Available: <https://doi.org/10.1162/neco.1997.9.8.1735>
- [10] J. Chung, Ç. Gülçehre, K. Cho, and Y. Bengio, "Empirical evaluation of gated recurrent neural networks on sequence modeling," *CoRR*, vol. abs/1412.3555, 2014. [Online]. Available: <http://arxiv.org/abs/1412.3555>
- [11] [Online]. Available: <http://doi.org/10.5281/zenodo.3478579>
- [12] D. P. Kingma and J. Ba, "Adam: A method for stochastic optimization," 2014, cite arxiv:1412.6980Comment: Published as a conference paper at the 3rd International Conference for Learning Representations, San Diego, 2015. [Online]. Available: <http://arxiv.org/abs/1412.6980>

7. Appendix

```
1 import librosa
2 import numpy as np
3 import pandas as pd
4 import os
5 import csv
6
7 numbers = '0 1 2 3 4 5 6 7 8 9'.split()
8
9 header = ''
10 for i in range(1, 641):
11     header += f'mfcc{i} '
12 header += ' label'
13 header = header.split()
14
15 file = open('largetrainingset.csv', 'w', newline='')
16 with file:
17     writer = csv.writer(file)
18     writer.writerow(header)
19
20 for number in numbers:
21     for filename in os.listdir(f'./recordings/{number}'):
22         nfile = f'./recordings/{number}/{filename}'
23
24         y, sr = librosa.load(nfile, sr=None, mono=True, duration=1)
25
26         if y.size < 16000:
27             rest = 16000 - y.size
28             left = rest // 2
29             right = rest - left
30
31             y = np.pad(y, (left, right), 'reflect')
32
33         mfccs = librosa.feature.mfcc(y=y, sr=sr)
34
35         row = mfccs.T.flatten()
36         row = np.append(row, number)
37
38         file = open('largetrainingset.csv', 'a', newline='')
39         with file:
40             writer = csv.writer(file)
41             writer.writerow(row)
```

Figure 8. Code snippet to extract MFCCs from the dataset.