

Insert here your thesis' task.



**FACULTY
OF INFORMATION
TECHNOLOGY
CTU IN PRAGUE**

Master's thesis

Developing an automatic speech recognition system based on Czech spoken language

Bc. Richard Werner

Department of Applied Mathematics

Supervisor: Mgr. Alexander Kovalenko, Ph.D.

May 14, 2020

Acknowledgements

THANKS (remove entirely in case you do not wish to thank anyone)

Declaration

I hereby declare that the presented thesis is my own work and that I have cited all sources of information in accordance with the Guideline for adhering to ethical principles when elaborating an academic final thesis.

I acknowledge that my thesis is subject to the rights and obligations stipulated by the Act No.121/2000 Coll., the Copyright Act, as amended, in particular that the Czech Technical University in Prague has the right to conclude a license agreement on the utilization of this thesis as a school work under the provisions of Article 60 (1) of the Act.

In Prague on May 14, 2020

.....

Czech Technical University in Prague

Faculty of Information Technology

© 2020 Richard Werner. All rights reserved.

This thesis is school work as defined by Copyright Act of the Czech Republic. It has been submitted at Czech Technical University in Prague, Faculty of Information Technology. The thesis is protected by the Copyright Act and its usage without author's permission is prohibited (with exceptions defined by the Copyright Act).

Citation of this thesis

Werner, Richard. *Developing an automatic speech recognition system based on Czech spoken language*. Master's thesis. Czech Technical University in Prague, Faculty of Information Technology, 2020.

Abstrakt

V několika větách shrňte obsah a přínos této práce v českém jazyce.

Klíčová slova Replace with comma-separated list of keywords in Czech.

Abstract

Summarize the contents and contribution of your work in a few sentences in English language.

Keywords Replace with comma-separated list of keywords in English.

Contents

Introduction	1
1 Neural networks	3
1.1 Basic concept – Rosenblatt’s Perceptron	5
1.2 Multilayer perceptron	7
1.2.1 Backpropagation	7
1.2.2 The vanishing gradient problem	12
1.3 Recurrent neural networks	14
1.3.1 The problem of long-term dependencies	15
1.3.2 Long-short-term memory networks	15
1.3.3 Gated recurrent units	18
2 Automatic speech recognition	21
2.1 Introduction	23
2.2 Feature extraction	26
2.2.1 Wavelet vs. Fourier transform	26
2.2.2 MFCC features	32
2.2.3 Wavelets in feature extraction	37
3 Experiments	43
4 Realization	45
Conclusion	47
Bibliography	49
A Acronyms	53
B Contents of enclosed CD	55

List of Figures

1.1	Neuron anatomy	3
1.2	Perceptron schema	5
1.3	Multilayer NN example schema	7
1.4	Sigmoid function and its derivative	12
1.5	A recurrent neural network layer unfolding	14
1.6	Standard RNN chain structure	15
1.7	LSTM chain structure	16
1.8	Three LSTM gate layers	16
1.9	LSTM cell state operators	17
1.10	GRU unit	18
2.1	Text-to-Speech synthesis system	23
2.2	Generic pattern matching system	24
2.3	Haar wavelet	28
2.4	Haar wavelet with different parameters	29
2.5	Two-stage two-band analysis tree	30
2.6	Transform comparisons	31
2.7	Digital signal before and after the preemphasis	32
2.8	A frame before and after the Hamming window function application	33
2.9	Window functions	34
2.10	A frame spectrum computed with DFT	35
2.11	Mel filter banks	35
2.12	Spectrogram of the original signal	36
2.13	Cepstrum features of the original signal	36
2.14	Wavelet packet decomposition tree	39

List of Tables

2.1 MFCC features summary	37
-------------------------------------	----

Introduction

Neural networks

Inspiration to create artificial neural networks (ANN), commonly referred to as *neural networks* (NN), was taken from the entirely different way in which human brain processes and computes information, compared to conventional digital computers. The human brain is a highly complex, nonlinear, and parallel information-processing system with the capability to organize its structural components, known as *neurons*, to perform several different kinds of computation (e.g., perception, pattern recognition or motor functions). [1]

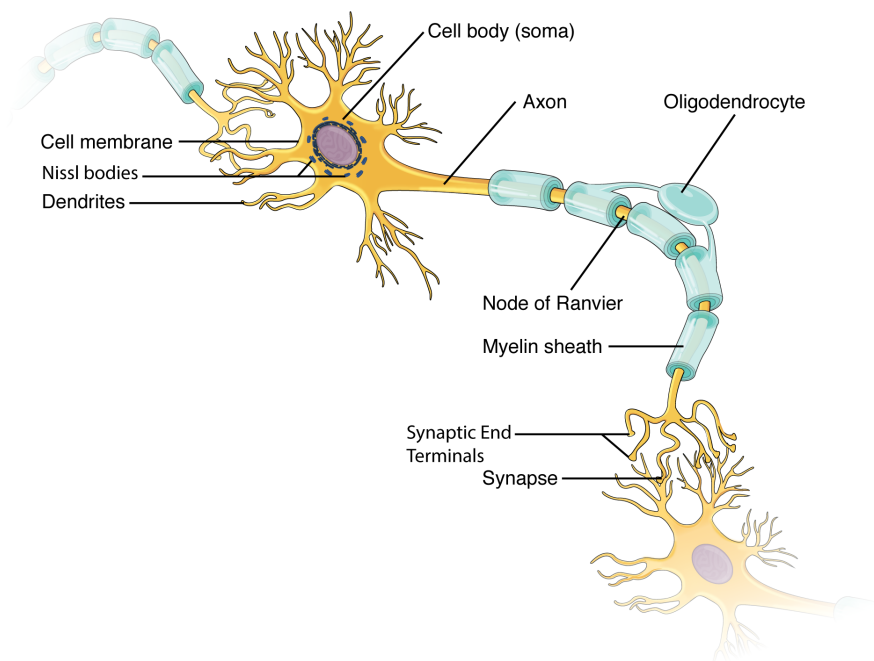


Figure 1.1: Neuron anatomy by [2]

As shown in figure 1.1, every neuron has a branch-like structure of dendrites which serve as receivers of information from connected neurons or another kind of receptor/cell. From dendrites, the information is then passed as an electric signal through the cell body (soma) to the axon. The purpose of the axon is the transportation of the signal to adjacent cells. On the end of the axon is another branch-like structure called synaptic end terminals (or axon terminals). Those are the transfer points for the information to other cells. [3]

An essential feature of the brain, which is called brain *plasticity*, is the ability to *adapt* to its surrounding environment. And just as the plasticity is the key to the functioning of the brain, the same applies to the artificial neural networks and their artificial neurons.

In its most general form, a neural network is a machine that tries to *model* how the human brain performs a particular task or function of interest. The network is usually implemented by using electronic components or is simulated by software on a digital computer. With this in mind, let us define neural networks as follows:

A neural network is a massively parallel distributed processor made up of simple processing units that has a natural propensity for storing experiential knowledge and making it available for use. It resembles the brain in two respects:

- 1. Knowledge is acquired by the network from its environment through a learning process.*
- 2. Inter-neuron connection strengths, known as synaptic weights, are used to store the acquired knowledge. [1]*

The mentioned learning process is performed by a *learning algorithm*, which is a function that modifies the synaptic weights of the network to fit it to a specific problem or objective. Apart from modifying the internal weights, some networks are capable of even modifying their topology to suit themselves to a problem better.

In the following sections, I will go through the first concepts and history of neural networks, their learning algorithms, and advanced types of NN.

1.1 Basic concept – Rosenblatt’s Perceptron

Rosenblatt’s perceptron was the first algorithmically described neural network. It was invented by psychologist Rosenblatt who inspired engineers, physicists, and mathematicians alike to devote their research effort to different aspects of neural networks in the 60s and 70s of the 20th century. Moreover, the concept of perceptron in its basic form is as valid today as it was in 1958 when Rosenblatt’s paper on the perceptron was first published. [1]

An artificial neural network, in general, can be viewed as an oriented graph, where nodes are neurons, and edges are the inter-neuron connections. Each neuron consists of three basic components [4]:

1. **weight vector**, which stores the weights of the connections between neurons,
2. **summation function**, which computes weighted sum of all inputs, and
3. **activation function** that maps the result of sum function to a number based on the specific type of the activation function.

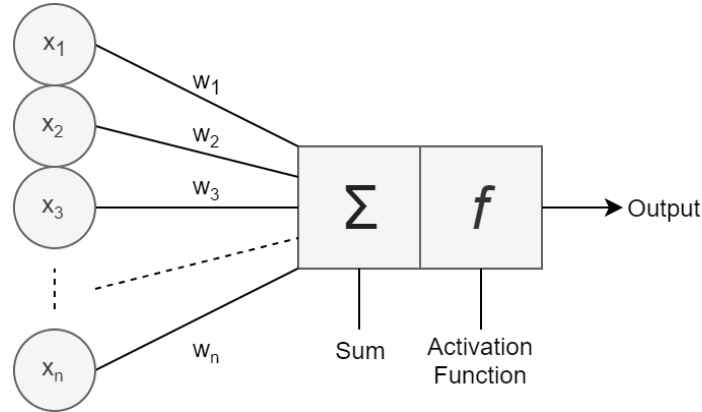


Figure 1.2: Perceptron schema

In the figure 1.2 you can see the general schema of a perceptron cell. Each of the inputs x_i has a corresponding weight w_i . The weighted sum of the inputs a would then be computed as:

$$a = \sum_{i=0}^n w_i x_i + b, \quad (1.1.0.1)$$

where b is externally applied bias.

The activation function of the original perceptron was implemented as the signum function. That means the perceptron was able to classify an input vector $X = (x_1, \dots, x_n)$ into one of two classes, because the range of this

function is a two-value set $\{-1, 1\}$. In the purest form of the perceptron, there are two decision regions separated by a *hyperplane*, which is defined by:

$$\sum_{i=0}^n w_i x_i + b = 0 \quad (1.1.0.2)$$

which would for example mean a simple line in a 2D space:

$$w_1 x_1 + x_2 x_2 + b = 0 \quad (1.1.0.3)$$

Since the basic perceptron is only able to generate a hyperplane as a decision boundary, a perceptron model is limited to linearly separable problems, thus is unable to classify any more complex tasks correctly. This leads us to a conclusion, that something more robust is needed. [1, 3]

1.2 Multilayer perceptron

The multilayer perceptron (MLP) consists of multiple layers of neurons that interact using weighted connections. After the first input layer, there are usually several more hidden layers, followed by the last output layer. There are no connections between neurons in one layer, while all neurons in one layer are fully connected to neurons in adjacent layers. While this statement tends to create an impression that all the information from one layer is copied to all the neurons of the adjacent layer, it is not necessarily so. Since the connection weights are usually real numbers, some of the connections may effectively be rendered insignificant. [5, 1] There is an example of a multilayer NN in figure 1.3.

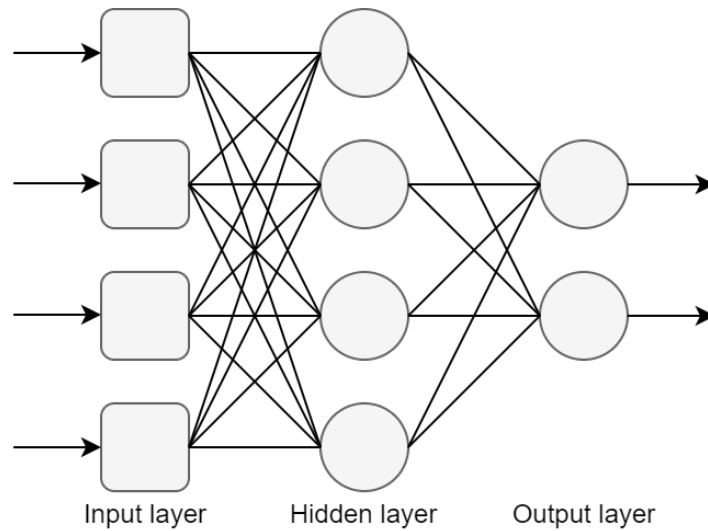


Figure 1.3: Multilayer NN example schema

Another difference over the simple perceptron is the activation function. Each model of neuron includes a nonlinear activation function that is *differentiable*.

Those characteristics, however, are responsible for the deficiencies in our understanding of the behavior of the network. The distributed non-linearity and high connectivity of the network make the analysis of the network quite hard to undertake. Moreover, the use of hidden layers makes the learning process harder to visualize.

1.2.1 Backpropagation

Backpropagation, short for “backward propagation of errors”, is an algorithm for supervised learning of ANNs, which is using the *gradient descend*. In this

section, especially in the formal definition, I've drawn from [6]. It proceeds in two phases:

1. In the **forward phase**, the connection weights of the network are fixed, and the input is propagated through the network, layer by layer, until it reaches the output. In other words, the network performs inference.
2. In the **backward phase**, the error is computed by comparing the network output with the desired output. The resulting error is then propagated through the network, but this time from the output layer to the input. In this phase, appropriate adjustments are made to the connection weights based on the amount of error introduced by the corresponding connection. [1]

Formal definition

If we were to define the backpropagation formally, there would be three things required:

1. A **dataset** consisting of input-output pairs (\vec{x}_i, \vec{y}_i) , where \vec{x}_i is the input and \vec{y}_i is the corresponding desired output. The dataset of size N is denoted $X = \{(\vec{x}_1, \vec{y}_1), \dots, (\vec{x}_N, \vec{y}_N)\}$.
2. A **feedforward NN**, e.g. our MLP, whose parameters are collectively denoted θ . Parameters of primary interest are w_{ij}^k , which is the connection weight between node j in layer l_k and node i in layer l_{k-1} , and b_i^k , the bias for node i in layer l_k .
3. An **error function** $E(X, \theta)$, which defines the error between the desired output \vec{y}_i and the calculated (inferred) output $\hat{\vec{y}}_i$ of the NN on input \vec{x}_i for a dataset X and parameters θ .

Training a NN with gradient descent requires the calculation of the gradient of the error function $E(X, \theta)$ with respect to the network parameters. Then, according to the learning rate α , each iteration of gradient descent updates the weights and biases θ as:

$$\theta^{t+1} = \theta^t - \alpha \frac{\partial E(X, \theta^t)}{\partial \theta}, \quad (1.2.1.1)$$

where θ^t denotes the parameters of the network at iteration t in gradient descent.

The following formulation is for a neural network with one output. The algorithm can be applied to a network with any number of outputs by consistent application of the chain rule and power rule. Thus, for all the examples below, input-output pair will be of the form (\vec{x}, y) , i.e. the target variable is a number, not a vector.

Preliminaries

Let us use the *mean squared error* as our error function:

$$E(X, \theta) = \frac{1}{2N} \sum_{i=1}^N (\hat{y}_i - y_i)^2, \quad (1.2.1.2)$$

where y_i is the target output value and \hat{y}_i is the computed output value for the input \vec{x}_i . This function is used for numeric (not a vector) output. If we were to predict vectors, the use of a vector similarity function would be necessary.

Lets us lay down a few preliminaries:

- the sum of weighted inputs for the i -th node in the k -th layer a_i^k is:

$$a_i^k = b_i^k + \sum_{j=1}^{r_{k-1}} w_{ij}^k o_j^{k-1} = \sum_{j=0}^{r_{k-1}} w_{ij}^k o_j^{k-1} \quad (1.2.1.3)$$

where b_i^k is bias for node i in layer k , r_k is number of nodes in layer l_k , w_{ij}^k is weight for node j in layer l_k from incoming node i and o_j^k is output for node i in layer l_k . After the standard version with the explicit bias follows the one with the bias incorporated into the weights:

$$w_{i0}^k = b_i^k. \quad (1.2.1.4)$$

- Backpropagation attempts to minimize the error function with respect to the NN weights by calculating the value of $\frac{\partial E}{\partial w_{ij}^k}$ for each weight w_{ij}^k . We can decompose the derivative with respect to each input-output pair as follows:

$$\frac{\partial E(X, \theta)}{\partial w_{ij}^k} = \frac{1}{N} \sum_{d=1}^N \frac{\partial}{\partial w_{ij}^k} \left(\frac{1}{2} (\hat{y}_d - y_d)^2 \right) = \frac{1}{N} \sum_{d=1}^N \frac{\partial E_d}{\partial w_{ij}^k}. \quad (1.2.1.5)$$

Error function derivatives

The derivation of the backpropagation algorithm begins by applying the *chain rule* to the error function partial derivative:

$$\frac{\partial E}{\partial w_{ij}^k} = \frac{\partial E}{\partial a_j^k} \frac{\partial a_j^k}{\partial w_{ij}^k}, \quad (1.2.1.6)$$

where a_j^k is the weighted sum (activation) of node j in layer k before it is passed to the nonlinear activation function. This decomposition simply says the change in the error function due to a weight is a product of the change in

the error function E due to the sum a_j^k times the change in the sum a_j^k due to the weight w_{ij}^k .

Mostly, the whole equation is simplified by calling the first term the *error* and denoting it

$$\delta_j^k = \frac{\partial E}{\partial a_j^k}. \quad (1.2.1.7)$$

The second term is calculated from the equation (1.2.1.3) for a_i^k above:

$$\frac{\partial a_j^k}{\partial w_{ij}^k} = \frac{\partial}{\partial w_{ij}^k} \left(\sum_{l=0}^{r_{k-1}} w_{lj}^k o_l^{k-1} \right) = o_i^{k-1}. \quad (1.2.1.8)$$

And finally the partial derivative of the error function E with respect to a weight w_{ij}^k is

$$\frac{\partial E}{\partial w_{ij}^k} = \delta_j^k o_i^{k-1}. \quad (1.2.1.9)$$

This is quite intuitive, since the weight w_{ij}^k connects the output of node i in layer $k-1$ to the input of node j in layer k in the computation graph.

The output layer

Backpropagation starts from the final layer, thus it attempts to define the value δ_1^m , where m is the index of the final layer and the subscript is constant 1, because we are concerned with only one-output NN. We are able to express the output value by a neuron from the weighted sum from the previous layer and the activation function of the current neuron. Therefore, the error function can be expressed as

$$E = \frac{1}{2}(\hat{y} - y)^2 = \frac{1}{2}(g_0(a_1^m) - y)^2, \quad (1.2.1.10)$$

where $g_0(x)$ is the activation function for the single neuron in the output layer. Applying the partial derivative and using the chain rule gives

$$\delta_1^m = (g_0(a_1^m) - y)g_0'(a_1^m) = (\hat{y} - y)g_0'(a_1^m). \quad (1.2.1.11)$$

Putting this together with the partial derivative of the error function E with respect to a weight in the final layer w_{i1}^m gives

$$\frac{\partial E}{\partial w_{i1}^m} = \delta_1^m o_i^{m-1} = (\hat{y} - y)g_0'(a_1^m)o_i^{m-1}. \quad (1.2.1.12)$$

The hidden layers

Firstly, we need an equation for the error term δ_j^k for layers $1 \leq k < m$. Again, with the help from the chain rule, we get:

$$\delta_j^k = \frac{\partial E}{\partial a_j^k} = \sum_{l=1}^{r_{k+1}} \frac{\partial E}{\partial a_l^{k+1}} \frac{\partial a_l^{k+1}}{\partial a_j^k}, \quad (1.2.1.13)$$

where l ranges from 1 to r^{k+1} . l does not take the value of zero because the bias input o_0^k corresponding to w_{0j}^{k+1} is fixed, and its value is not dependent on the outputs of previous layers. Using the error term δ_l^{k+1} gives the following equation:

$$\delta_j^k = \sum_{l=1}^{r_{k+1}} \delta_l^{k+1} \frac{\partial a_l^{k+1}}{\partial a_j^k}. \quad (1.2.1.14)$$

With definition of a_l^{k+1} as

$$a_l^{k+1} = \sum_{j=1}^{r_k} w_{jl}^{k+1} g(a_j^k), \quad (1.2.1.15)$$

where $g(x)$ is the activation function for the hidden layers, we get:

$$\frac{\partial a_l^{k+1}}{\partial a_j^k} = w_{jl}^{k+1} g'(a_j^k). \quad (1.2.1.16)$$

Using this in the above equation yields a final equation for the error term δ_j^k for the hidden layers called the *backpropagation formula*:

$$\delta_j^k = \sum_{l=1}^{r_{k+1}} \delta_l^{k+1} w_{jl}^{k+1} g'(a_j^k) = g'(a_j^k) \sum_{l=1}^{r_{k+1}} w_{jl}^{k+1} \delta_l^{k+1}. \quad (1.2.1.17)$$

And finally, the partial derivative of the error function E with respect to a weight in the hidden layers w_{ij}^k for $1 \leq k < m$ is

$$\frac{\partial E}{\partial w_{ij}^k} = \delta_j^k o_i^{k-1} = g'(a_j^k) o_i^{k-1} \sum_{l=1}^{r_{k+1}} w_{jl}^{k+1} \delta_l^{k+1}. \quad (1.2.1.18)$$

The algorithm itself

Assuming a suitable value for learning rate α and a random initialization of the parameters w_{ij}^k , the backpropagation algorithm proceeds in four steps:

1. **Calculate the forward phase** for each input-output pair (\vec{x}_d, y_d) and store the results and store the NN output \hat{y}_d and inner neuron outputs a_j^k , and o_j^k for each node j in layer k by proceeding from layer 0, the input layer, to layer m , the output layer.
2. **Calculate the backward phase** for each input-output pair and store the results $\frac{\partial E_d}{\partial w_{ij}^k}$ for each weight w_{ij}^k connecting node i in layer $k-1$ to node j in layer k by proceeding from layer m to layer 1. There are three steps to compute the partial derivatives:
 - a) Evaluate the error term for the output layer δ_1^m .

- b) Backpropagate the error terms to the hidden layers δ_j^k , working from the final hidden layer.
 - c) Evaluate the partial derivatives of the individual error E_d with respect to w_{ij}^k .
3. **Combine the individual gradients** $\frac{\partial E_d}{\partial w_{ij}^k}$ for each pair to get the total gradient $\frac{\partial E(X, \theta)}{\partial w_{ij}^k}$ for the entire set of pairs $X = (\vec{x}_1, y_1), \dots, (\vec{x}_N, y_M)$ by using the above mentioned (section 1.2.1) simple average of the individual gradients.
4. **Update the weights** according to the learning rate α and the total gradient by using

$$\Delta w_{ij}^k = -\alpha \frac{\partial E(X, \theta)}{\partial w_{ij}^k}, \quad (1.2.1.19)$$

which is moving the weights in the direction of the negative gradient – gradient descend.

1.2.2 The vanishing gradient problem

The problem lies in many layers using certain activation functions, a sigmoid function, for example, which squishes a large input space into a small input space between 0 and 1. That means a significant change in the input leads to a small change on the output; hence its derivative becomes small.

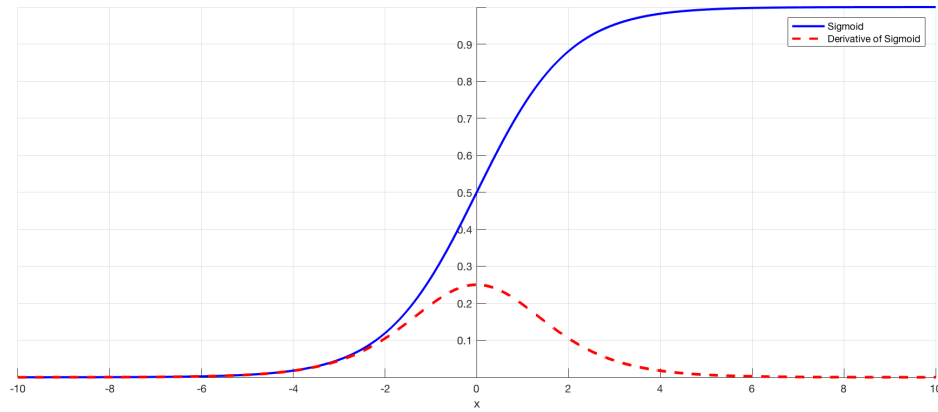


Figure 1.4: Sigmoid function and its derivative by [7]

There is an example in figure 1.4 with the sigmoid function and its derivative. When the input of the sigmoid becomes either larger or smaller, the derivative becomes close to zero. For “shallow” networks with only a few layers with this type of activation function, it is not a big problem. However,

when more layers are added, it can cause the gradient to be too small to have any significant impact on the learning process.

This is caused by the chain rule used in the backpropagation, which is described for feedforward networks in section 1.2.1. It computes the gradients by moving layer by layer from the final one to the initial one. By the chain rule, each subsequent derivative is multiplied by the already computed value. Therefore, when there are n hidden layers using sigmoid-like activation function, n small derivatives are multiplied together. Thus, the gradient value decreases exponentially as the backpropagation algorithm advances to the initial layers. [7]

There are a few well-known solutions [7]:

- Probably the simplest solution is to use a different activation function, such as ReLU (Rectified Linear Unit), which does not cause a small derivative.
- Another solution is *residual networks* (ResNets). They provide residual connections (without weight parameters) straight to the next layers, effectively skipping the activation functions. That results in overall higher derivatives, and therefore the ability to train much deeper networks. [8]
- The last one is *batch normalization* layers. As mentioned before, the problem appears when a large input is mapped to small output, causing the derivatives to disappear. The batch normalization method normalizes the input on a predefined scale, where the sigmoid derivative is not that small.

1.3 Recurrent neural networks

A standard neural network takes in a fixed size vector as input, which limits its usage in situations that involve a *series* type input with no predetermined size. In other words, vanilla neural networks are unable to take into account any context or history, even if the current input is in a way based on the previous one. That is the reason why the *recurrent neural networks* (RNN) were designed and implemented.

The RNNs are made to take a series of input with no predetermined size. A single input item from the series is related to others, and likely influences its neighbors. So while they remember information from the training phase, they also learn things from prior inputs while generating outputs.

A recurrent network can take one or more input vectors and produce one or more output vectors. The outputs are influenced not just by weights applied on inputs like a regular NN, but also by a “hidden” state vector representing the context based on prior inputs/outputs. That means the same input could produce a different output depending on previous inputs in the series. [9]

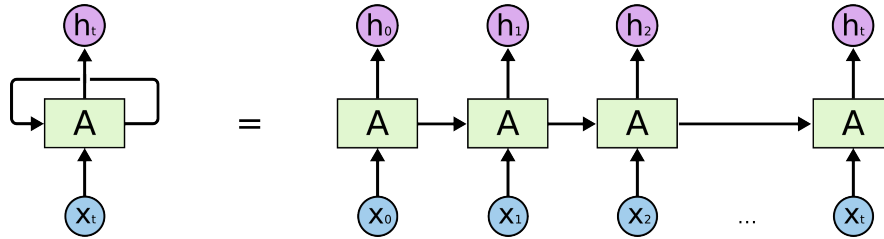


Figure 1.5: A recurrent neural network layer unfolding by [10]

The above-described concept is pictured in figure 1.5. On the unfolded neural network (the right side of the picture), we can see that part of the information from the NN in time 0 is passed to the network in time 1. That continues throughout the whole series of inputs, and in theory, the whole history of inputs and calculations will influence every other computation of the network.

There are many varieties, solutions, and constructive elements of recurrent neural networks. The difficulty of a recurrent network is that if we take into account each time step, each one of them must create its layer of neurons. That causes severe computational complexity. Besides, multilayer implementations are computationally unstable, since they tend to disappear or go off the scale. Restricting the computation to a fixed time window resolves those problems, but the model will not reflect the long-term trends. Various approaches are trying to improve the model of RNN memory and the mechanisms of remembering and forgetting. [11]

1.3.1 The problem of long-term dependencies

Sometimes, for a task, we have enough recent information to be able to compute the output with certainty. For example, imagine a language model that predicts the next word based on the previous inputs/outputs. If it tries to predict the last word in the sentence “Clouds in the sky,” we do not need any context other than the rest of the sentence – it is quite apparent what would the last word regard (the sky/heaven). In such cases, where there is only a small gap between the necessary information, the standard RNNs can learn to use the recent information.

However, there are also times when we need a broader context. If a model were to predict the last words in the sentence “I grew up in France . . . I speak French fluently,” the recent context would probably suggest a name of a language. If it were to clarify which one, it would need the previous context, up to the beginning of the sentence. It is not rare that the gap between the necessary information and the place where it is needed becomes very large. Unfortunately, as the gap grows, the standard recurrent networks become unable to learn how to connect these pieces of information. As was mentioned above, in theory, the RNNs should handle such long-term dependencies. However, in practice, RNNs are not able to learn this.

1.3.2 Long-short-term memory networks

Long-short-term memory (LSTM) networks are a special type of recurrent neural networks capable of learning long-term dependencies as well as short-term ones. They work very well on a large variety of problems and are currently widely used. LSTMs are specifically designed to avoid the problem mentioned above of long-term dependencies, and they also manage to avoid the vanishing gradient problem.

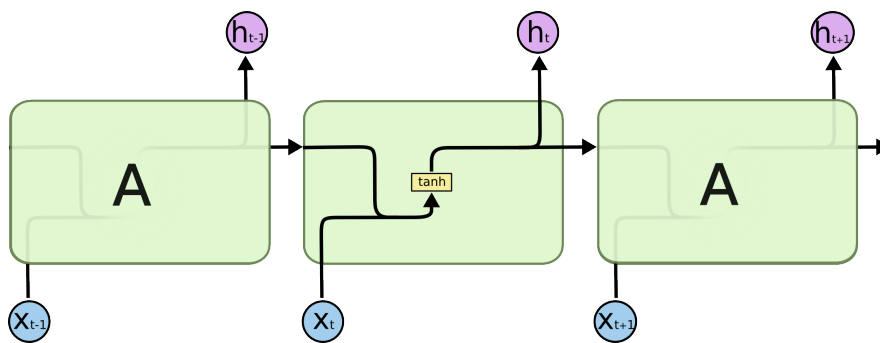


Figure 1.6: Standard RNN chain structure by [10]

Recurrent neural networks are in the form of a chain of repeating elements (neurons). In standard RNNs, those elements are usually of a straightforward structure – for example, one layer of a hyperbolic tangent (tanh). LSTM

networks have that same chain structure, but each of the elements implements a much more complex structure. Instead of one layer, they implement four layers that uniquely interact with each other.

The difference between the inner structures of standard and LSTM recurrent networks can be seen in figures 1.6 for the standard approach and 1.7 for the LSTMs. In the second diagram, each black line transmits a whole vector from one node to another. Merging lines mean vector concatenation, and branching lines are copying content without any alteration. Pink circles represent pointwise operators – vector addition or multiplication. Lastly, the yellow rectangles are the trained layers of the neural network. [12]

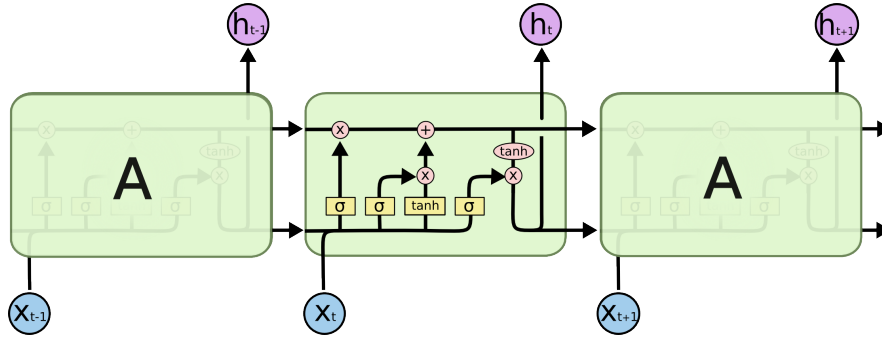


Figure 1.7: LSTM chain structure by [10]

The main feature of LSTMs is the *cell/cellular state*. That is the top horizontal line in the network cell. The neurons can remove or add information to a cellular state, while so-called *gates* carefully regulate the changes. They are composed of a sigmoid function and a pointwise multiplication operation. The sigmoid function returns values between 0 and 1, and it determines the magnitude of the change to the cellular state, whether it is information removal or addition.

To explain the information flow process through the LSTM cell, let us go step by step through it [10]:

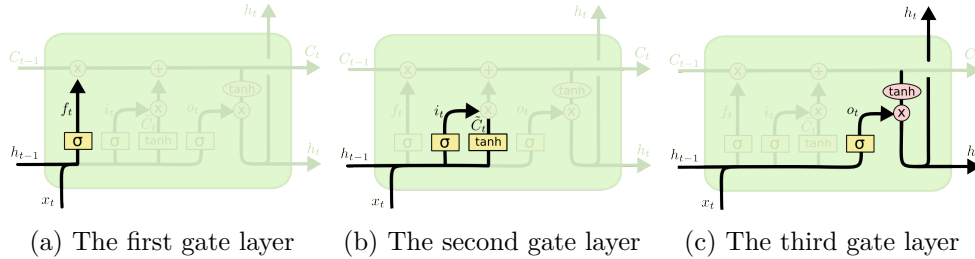


Figure 1.8: Three LSTM gate layers by [10]

1. The first step is to decide what information we are going to remove or keep in the cell state. It is decided by the first gate called the *forget gate layer* displayed in figure 1.8a. It combines values h_{t-1} and x_t , sends it to the sigmoid function, and the function returns a number between 0 and 1 for each number in the cell state C_{t-1} . A 1 means “keep as is,” while a 0 represents “completely remove.” Mathematically, this step would look like this:

$$f_t = \sigma(W_f[h_{t-1}, x_t] + b_f), \quad (1.3.2.1)$$

where W_f are the weights of forget gate layer and b_f is the the bias.

If we were to use the language model example again, the LSTM network could be deciding whether to keep information about the gender currently being spoken of or information about the subject.

2. The next step is deciding which information the network will add to the cellular state. That consists of two parts. Firstly, a sigmoid layer called the *input gate layer* decides which values will be updated. The next step is performed by the *tanh* layer, which creates a vector of new candidate values \tilde{C}_t . The equations for this step are:

$$i_t = \sigma(W_i[h_{t-1}, x_t] + b_i), \quad (1.3.2.2)$$

and

$$\tilde{C}_t = \tanh(W_C[h_{t-1}, x_t] + b_C), \quad (1.3.2.3)$$

where W_i/b_i and W_C/b_C are weights/biases for the input gate layer and the *tanh* layer respectively.

In the language model example, the network would add new information about gender or a new subject. This step, or components active during this step, is displayed in figure 1.8b.

3. Now it is time to update the old cellular state C_{t-1} into the new cellular state C_t . This step is only the execution of the two steps described above.

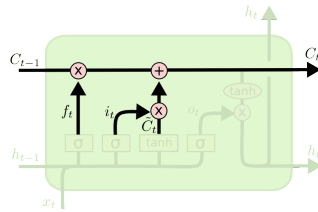


Figure 1.9: LSTM cell state operators by [10]

The old state is multiplied by f_t , forgetting, or keeping the information we decided in step 1. Then, as described in step 2, new information

$i_t \tilde{C}_t$ is added, where \tilde{C}_t is the information and i_t the scale deciding the magnitude of an update for each cell state value. The equation is then:

$$C_t = f_t \otimes C_{t-1} + i_t \otimes \tilde{C}_t. \quad (1.3.2.4)$$

4. Finally, we decide what the output will be. The output is a filtered version of the cell state. It is again a combination of sigmoid values that decide which parts of information go through, and cell state values “activated” by a \tanh layer. This step is displayed in figure 1.8c. The final two equations for this process are:

$$o_t = \sigma(W_o[h_{t-1}, x_t] + b_o), \quad (1.3.2.5)$$

and

$$h_t = o_t \otimes \tanh(C_t), \quad (1.3.2.6)$$

where W_o are the weights for the output gate layer and b_o is the bias.

As for the language model example, since the network just saw a subject, it might output information regarding a verb. For example, information about singularity or plurality so that we know what form of a verb should be used.

1.3.3 Gated recurrent units

Gated recurrent unit (GRU) is a simplification of LSTM-like units. It contains only two gates [13]:

- an **update gate** z_t , whose role is similar to the LSTM forget gate, and
- a **reset gate** r_t , whose role is somewhat similar to the input gate.

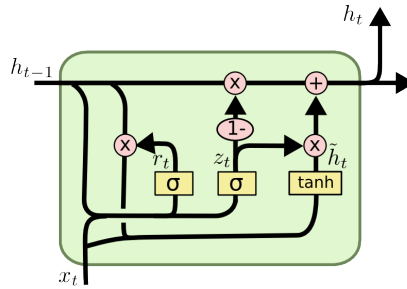


Figure 1.10: GRU unit by [10]

Other than removing the output gate from LSTM, GRU also removed the cellular state C_t . The schema of this unit is displayed in figure 1.10. To

represent the data flow through the GRU unit mathematically, we will start with the update gate:

$$z_t = \sigma (W_z [h_{t-1}, x_t] + b_z) . \quad (1.3.3.1)$$

The reset gate is quite similar:

$$r_t = \sigma (W_r [h_{t-1}, x_t] + b_r) . \quad (1.3.3.2)$$

As for the computation of the preliminary output \tilde{h}_t :

$$\tilde{h}_t = \tanh (W_h [r_t \otimes h_{t-1}, x_t] + b_h) . \quad (1.3.3.3)$$

Finally, the output:

$$h_t = (1 - z_t) \otimes h_{t-1} + z_t \otimes \tilde{h}_t . \quad (1.3.3.4)$$

Evaluations found that when LSTM and GRU have the same amount of parameters, GRU slightly outperforms LSTM. [13]

Automatic speech recognition

Automatic speech recognition (ASR) has been undergoing active research for more than fifty years. It is an essential milestone in both human-human and human-machine communication. Due to the insufficient performance of technologies in the past, ASR has not become a desirable part of human-machine communication. It was because the lack of computing power did not allow passing the usability bar for real users, and other means of communication, such as keyboard and mouse, significantly outperformed speech in most aspects of communication efficiency with computers. [14]

This all changed in recent years. Speech technologies started to change how we live and work and, for some devices, became the primary means of interaction with them. By [14], there are several key areas of which progress allowed this trend:

- the first area is *Moore's law*. The law states that approximately every two years, the number of transistors in a dense integrated circuit doubles roughly every two years. [15] That leads to the computational capabilities of CPU/GPU clusters also being doubled every two years. That makes training of more complex and powerful models possible, and thus the error rates of ASR systems lower.
- The second area is the access to more data due to the continued advance of the Internet and cloud computing. By using and training models on big data collections, it is possible to build much more robust and assumption-less models than before.
- The third is that mobile, wearable, and intelligent living room devices and in-vehicle infotainment [16] became quite popular. Since the use of alternative interaction means, such as keyboard and mouse, is mostly not possible in these cases, speech communication, which is natural for humans, becomes more convenient.

2. AUTOMATIC SPEECH RECOGNITION

There are several approaches to the ASR with different models, such as Gaussian mixture models or hidden Markov models. Since this thesis deals with neural networks, we are going to go through deep neural network (DNN) models in ASR.

Aside from diverse models, there are several ways to preprocess the raw audio sound. In the following sections, I am going to introduce some of the methods of preprocessing.

2.1 Introduction

Both in human and electronic communication, the speech information is encoded in the form of a continuously varying (analog) waveform that can be transmitted, recorded, manipulated, and ultimately decoded by a human listener. The primary analog form of the message is an acoustic waveform, which we call the *speech signal*. Those can be converted to an electrical waveform by a microphone, further manipulated by analog and digital signal processing, and then converted back to acoustic form by a loudspeaker or another electronic device.

Before we can apply any (digital) processing techniques, we must convert to the acoustic waveform, analog signal, to a sequence of numbers – digital signal. System or tool able to do such converting is called A-to-D converter, which creates digital representation by sampling the waveform in a very high rate, applies filters to preserve a prescribed bandwidth, and then reduces the sampling rate to the desired sampling rate. This discrete-time representation is the starting point for most digital signal processing applications.

Several areas that work with the digital representation of the speech signal follow. [17]

Speech coding

Perhaps the most widespread applications of digital speech are the A-to-D and D-to-A converting mentioned briefly above. It is commonly referred to as *speech coding* or *speech compression*, and its goal is to compress the digital waveform representation of speech into a lower bit-rate representation. The D-to-A decoder part of the system is often called the *synthesizer* because it must reconstruct the speech waveform from the discrete digital data.

Text-to-Speech

A Text-to-Speech analysis is another area that uses the digital signal. The goal of these systems is to start with text and automatically produce speech. The input is an ordinary text such as a newspaper article or an e-mail message. The system is depicted in figure 2.1. The first block named *Linguistic Rules*

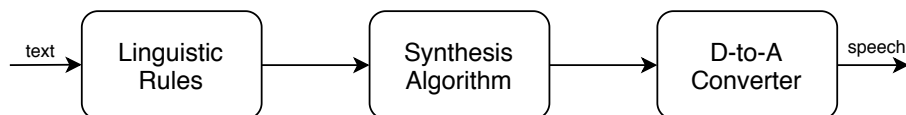


Figure 2.1: Text-to-Speech synthesis system inspired by [17]

has the job of converting the ordinary text into a set of sounds, which must be synthesized afterward. The system should include a set of linguistic rules that must determine appropriate sounds such as emphasis, pauses, or rates

of speaking. Moreover, pronounce acronyms and ambiguous words like *read*, *bass*, or *object*, how to pronounce abbreviations like *St.* (street or Saint?), *Dr.* (Doctor or drive?), and adequately pronounce specialized terms, and names. After the system builds the pronunciation set of sounds, it is time to synthesize the speech – to create the appropriate sound sequence to represent the text message in the form of speech.

Text-to-Speech synthesis systems are used to do things like to provide voice output from GPS systems, handle call center help desks, or providing information from handheld devices such as foreign language phrasebooks and dictionaries. They are also being used in announcement machines that provide information – stock quotes and airline schedules. Finally, perhaps the most important application is in reading machines for the blind, where an optical character recognition system provides the text input.

Speech recognition

Quite the opposite of the Text-to-Speech problem described above is the *speech recognition* or *Speech-to-Text* problem. It is concerned with the extraction of information from the speech signal. Figure 2.2 shows a diagram of a generic ap-

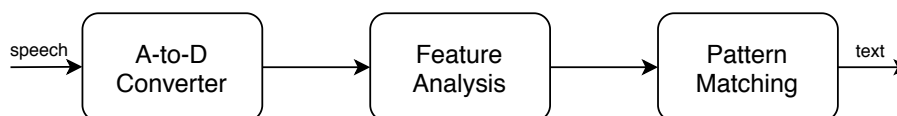


Figure 2.2: Generic pattern matching system inspired by [17]

proach to pattern matching problems in ASR. There are several sub-problems in this class, such as

- **speech recognition**, where the goal is to extract the message from the speech signal,
- **speaker recognition**, where the goal is to identify the speaker,
- **speaker verification**, which verifies a speakers claimed identity from analysis of their speech signal,
- **word spotting**, which involves monitoring a speech signal for the occurrence of a specified word or phrase, and
- **automatic indexing** of speech recordings based on the recognition of spoken keywords.

The first block in the pattern matching system diagram 2.2 converts the analog speech waveform to the digital signal using an A-to-D converter. The *feature analysis* module, the second block, converts the sampled speech signal to a

set of feature vectors. The third and final block in the system called *pattern matching* block dynamically time aligns the set of feature vectors representing the speech signal with a concatenated set of stored patterns. Then it chooses the identity associated with the pattern, which is the closest match to the time-aligned set of feature vectors of the speech signal. The symbolic output consists of a set of recognized words in the case of speech recognition. In the case of speaker recognition, it would be the identity of the best matching talker.

Probably the most common use of speech recognition is in portable communication devices. Spoken name speech recognition in cell phones enables voice dialing, and today basically every smart phone has its own voice assistant that is capable of multitude functions including speech-to-text when writing messages, the device options management, or information retrieval from the internet. [18]

The long-time dream of speech researchers are *automatic language translation* systems. The goal of such systems is to convert spoken words in one language to spoken words in another language to facilitate natural language voice dialogues between people speaking different languages.

2.2 Feature extraction

In this section, we are going to introduce two popular methods of speech signal feature extraction – namely, the Mel Frequency Cepstral Coefficients (MFCC) and the Wavelet approach.

The MFCC approach is computing the short-term *power spectrum*, which is based on an *inverse Fourier transform* of a log *power spectrum* on a nonlinear *mel scale* of frequency. [19] There are several definitions of the cepstrum, and they sometimes deviate from the steps used in this thesis. The most notable difference could be the use of a discrete cosine transform (DCT) or the lack of a mel scale filter bank. [20]

The wavelet approach uses the *wavelet transform* in the process of feature extraction. This transform does not expect a stationary signal and is well localized both in time and frequency using a *multiresolution* approach. The following section introduces both Fourier and wavelet transforms, their differences, and their advantages.

2.2.1 Wavelet vs. Fourier transform

Let us start by stating that the Fourier transform can be viewed as a particular case of wavelet transforms. They are both expanding the input signal onto a set of basis functions and transforms that will give an efficient and informative description of the signal.

We are going to start by defining generic transforms and pinpointing the main differences between wavelet and Fourier transforms. After that, we will build an intuitive idea around the transforms to better understand what is hiding behind the mathematics. This section draws from [21, 22, 23].

A generic discrete linear decomposition of a signal $f(t)$ can be expressed as:

$$f(t) = \sum_l a_l \psi_l(t), \quad (2.2.1.1)$$

where a_l are the real-valued expansion coefficients, and $\psi_l(t)$ is a set of real-valued functions of t called the expansion set. If the expansion set is unique, it is called a *basis* for the class of functions that can be so expressed.

This is how a discrete Fourier transform looks like:

$$\hat{f}_n = \sum_{k=0}^{N-1} f_k e^{-\frac{2\pi i n k}{N}}, \quad (2.2.1.2)$$

with its inverse counterpart defined as:

$$f_k = \frac{1}{N} \sum_{n=0}^{N-1} \hat{f}_n e^{\frac{2\pi i n k}{N}}. \quad (2.2.1.3)$$

For a Fourier transform, the basis functions are sine and cosine functions – generally the complex exponential, which we can see in the definitions above.

The problem with the Fourier transform is that it expects a stationary signal in time. This presumption of a stationary signal comes from the basis this transform uses – sine and cosine. Those are periodic, infinite functions, and any linear combination of such functions creates, again, a periodic, infinite one. In other words, this transformation discards all the information about time but greatly captures the frequency information. We can say that the Fourier transform has great *frequency localization* and weak (none) *time localization*. That, however, does not go along with the properties of human speech.

Human speech is very far from being stationary. To be able to recognize words, syllables, and phonemes, we need to have both time and frequency localization. The first approach to tackle this problem is the *Gabor* transform named after Dennis Gabor.

The Gabor transform

Gabor transform is a special case of short-time Fourier transform (STFT). STFT splits the signal into short segments of equal length and computes the Fourier transform for each segment separately. Gabor transform is additionally using a Gaussian window function to separate the desired part of the signal from the rest, squishing undesired parts of the signal (close) to zero.

MFCC approach to the feature extraction is using a very similar process with the signal splitting and windowing. The whole process is described in section 2.2.2.

As hinted above, the Gabor transform is still not ideal. Increasing the window size lowers the time resolution for better frequency information. Smaller window, on the other hand, increases the time resolution, but we lose the low-frequency content available in longer time intervals. How to avoid such trade-off and keep both frequency and time resolution as high as possible? Let us first introduce mother wavelets.

A mother wavelet

A *mother wavelet* is a fundamental parametric function defined as:

$$\psi_{a,b}(t) = \frac{1}{\sqrt{a}} \psi\left(\frac{t-b}{a}\right) \quad (2.2.1.4)$$

This is a generic description of a wavelet, without a specific wavelet function given. It has two real parameters:

- parameter a , which is a non-zero scaling parameter, and
- parameter b – a translation parameter.

Now, we have a function that is both scalable and translatable in time. In other words, we have got a sliding window that has a dynamic size. Now,

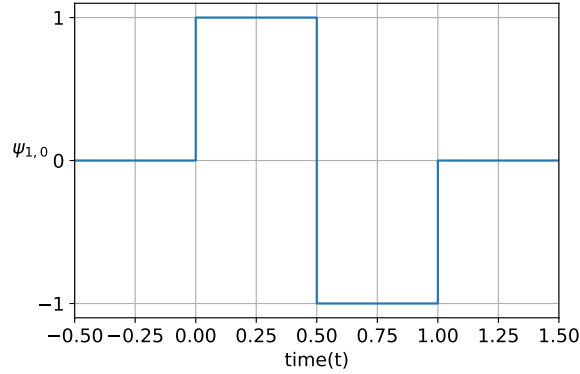


Figure 2.3: Haar wavelet

the only thing that remains is to define the wavelet function itself. There is a wide variety of wavelet functions that can be defined. Depending on the type of problem, we choose the one with the best properties for us. In the end, the wavelet is simply another expansion basis for the signal representation.

The Haar wavelet

The first recognized wavelet function was the Haar wavelet in 1910, meaning the idea of wavelets is more than a century old. It can be described as:

$$\psi(t) = \begin{cases} 1 & 0 \leq t < \frac{1}{2} \\ -1 & \frac{1}{2} \leq t < 1 \\ 0 & \text{else} \end{cases} \quad (2.2.1.5)$$

Figure 2.3 shows the Haar wavelet step function. It is an ideal wavelet for describing localized functions because it has *compact support*. Compact support means that a function returns values on a defined interval and returns zero outside the interval. That leads to Haar function having a strong time localization. However, its frequency localization is weak, because it decays like a sine function in powers of $\frac{1}{\omega}$ in the frequency domain (Fourier transform of the wavelet). This wavelet has two more properties, that are important for us. The sum of the “area under the curve” equals zero:

$$\int_{-\infty}^{\infty} \psi(t) dt = 0, \quad (2.2.1.6)$$

and the area of the squared absolute value of the curve equals one:

$$\int_{-\infty}^{\infty} |\psi(t)|^2 dt = 1. \quad (2.2.1.7)$$

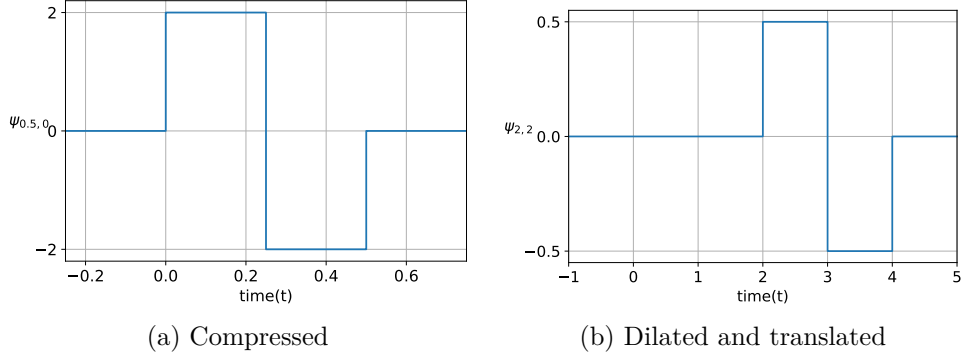


Figure 2.4: Haar wavelet with different parameters

The wavelet transform

To be precise, figure 2.3, and the wavelet description (2.2.1.5) are both the wavelet $\psi_{1,0}(t)$. Meaning its translation is zero, and its scaling is unity. Figure 2.4 shows other parameter configurations for Haar wavelets. Plot 2.4a depicts wavelet $\psi_{\frac{1}{2},0}$, which is a compressed function without translation, and plot 2.4b depicts $\psi_{2,2}$, which is a dilated wavelet with positive translation.

That allows for large scale structures in time in the signal to be captured by wide time-domain Haar wavelets. On this broad scale, the time resolution is pretty bad. However, by rescaling the wavelet in time, we get a more delicate time resolution along with high-frequency information. Therefore all the information at both low and high scales is preserved. That allows a complete reconstruction of the original signal. In the end, the only limit in this process is the number of rescaling levels. That is also where the *multiresolution* term comes from. The wavelet transform captures the information on multiple time-frequency resolutions.

The wavelet basis is computed as

$$(Tf)(\omega) = \int_t K(t, \omega) f(t) dt, \quad (2.2.1.8)$$

where $K(t, \omega)$ is the kernel of the transform. The same applies to the Fourier transform, where the kernel is the oscillations given by $K(t, \omega) = \exp(-i\omega t) = e^{-i\omega t}$. Now, we need to define a transform that incorporates the mother wavelet as the kernel. The following equation is the definition of the *continuous wavelet transform* (CWT):

$$\mathcal{W}_\psi[f](a, b) = (f, \psi_{a,b}) = \int_{-\infty}^{\infty} f(t) \bar{\psi}_{a,b}(t) dt. \quad (2.2.1.9)$$

CWT is a functions of the dilatation parameter a and the translation parameter b . A more simple yet equivalent formulation of the CWT by [21] is:

$$F(a, b) = \int_t f(t) \psi\left(\frac{t-a}{b}\right) dt. \quad (2.2.1.10)$$

2. AUTOMATIC SPEECH RECOGNITION

Since the *discrete wavelet transform* is somewhat complicated to define, we will first define the inverse wavelet transform. It is even a more intuitive idea – it says how the signal is reconstructed from the set of wavelets. It goes as:

$$f(t) = \sum_k \sum_j c_{a,b} \psi_{a,b}(t), \quad (2.2.1.11)$$

where the set of expansion coefficients $a_{j,k}$ is called the *discrete wavelet transform*.

The result of a DWT are coefficients of a so-called two-band analysis tree. It is a binary tree, where only one of the branches is expanded into the next stages. Figure 2.5 depicts such a tree, which visualizes two stages of the DWT algorithm.

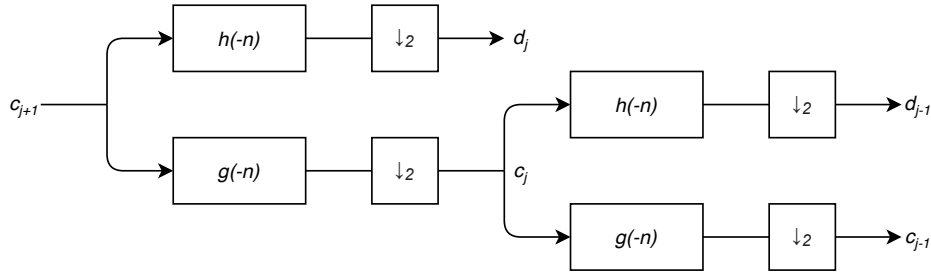


Figure 2.5: Two-stage two-band analysis tree, inspired by [21]

Each stage implements two digital filters, and two down-samplers (sometimes also called *decimators*). One of the filters, denoted as $h(\dots)$, is a high-pass filter, which is letting only higher frequencies of a signal through. The other one, denoted as $g(\dots)$, is a low-pass filter doing the opposite. The filters are where the wavelets play their role, and the filtering is achieved by convolving the signal sequence with the filter (wavelet) coefficients. The down-samplers discard half of the data that come through them; in this case, taking every other (even) value as $y(n) = x(2n)$. This decimation causes rougher resolution with every stage of the algorithm because we apply the same filter to fewer data. Also, it causes that the output from the tree has, on average, the same length as the input, meaning that there is a possibility of no information loss, and therefore complete signal recovery.

The output from high-pass filters d_i is called *detail coefficients*, and the output from low-pass filters c_i is called *approximation coefficients*.

One stage of the tree in figure 2.5 is implemented by two recursive equations. The equation for the high-pass filter coefficients:

$$d_a(b) = \sum_m h(m - 2b) c_{a+1}(m), \quad (2.2.1.12)$$

and the one for the low-pass filter coefficients:

$$c_a(b) = \sum_m g(m - 2b) c_{a+1}(m), \quad (2.2.1.13)$$

where $n = m - 2b$ to correspond with the figure 2.5. Also note, that the tree uses time-reversed recursion coefficients $h(-n)$, and $g(-n)$.

Summary

We have explained differences and (dis)advantages of Fourier, Gabor, and wavelet transforms. Figure 2.6 is quite a famous picture to visualize these differences. Understanding this cartoon is critical for understanding the wavelets and the wavelet transform.

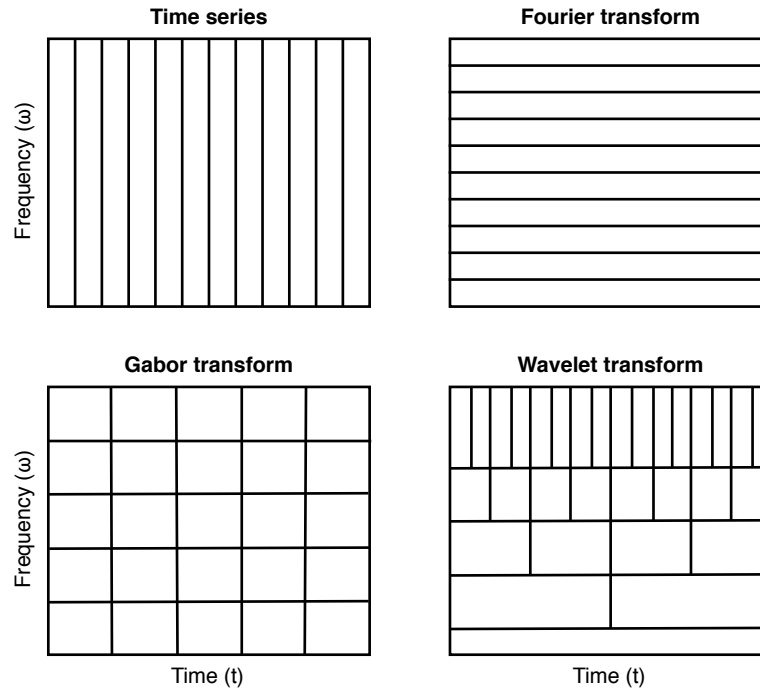


Figure 2.6: Transform comparisons, inspired by [22]

It shows us that a time series, the top left box, has all the time information we can get. We know what exactly happens in every point of time in the signal, but it tells us nothing about the magnitude of the signal at different frequencies.

The second box, the top right one, visualizes the Fourier transform of the time series. We cast the signal from the time domain into the frequency domain and therefore lose all the time information. So, as described above, we have perfect frequency resolution, but no way to say or reconstruct when any event in the signal happened. That is why Fourier transform expects a stationary signal in time.

The third box, the bottom left one, is the Gabor transform. We are splitting the signal into windows and applying Fourier onto those windows.

We trade some time resolution for frequency resolution. It is a trade-off, a compromise where the window length parameter says, how much of each we want to trade.

The last box, the bottom right one, is the wavelet transform. It starts at the bottom of the box. The first thing the algorithm does is to take the whole signal and (to some extent) perform a transform on it (with the mother wavelet scaled to the entire signal). Now, we take the mother wavelet, scale it to half its size and slide it across the time series. Repeat this until we have enough granular time resolution for the analysis we need to perform. After this, we have high frequency (bottom of the box) resolution and also the time (top of the box) resolution.

2.2.2 MFCC features

The goal of this section is to describe the transformation process of the digital signal into a sequence of acoustic feature vectors. Each of the vectors represents the information in a small time window of the signal. This section draws from [24]. Most of the figures in this section are my own, created with help from [25]. The original wave file, which is the whole process of feature extraction demonstrated on, is my recording saying “yellow dandelion.”

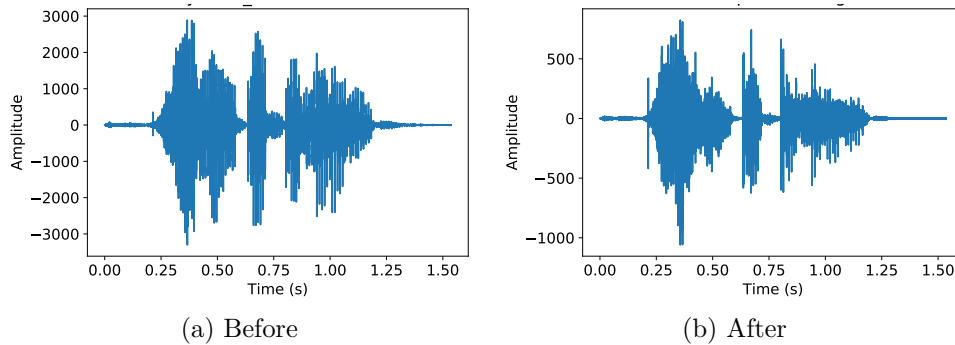


Figure 2.7: Digital signal before and after the preemphasis

Preemphasis

The first step in this process is to boost the amount of energy in the high frequencies. If we look at the frequency spectrum of voiced segments like vowels, there is more energy at the lower frequencies than the higher frequencies. This drop of energy is called the *spectral tilt*.

This preemphasis is done by applying a filter which goes as

$$y_n = x_n - \alpha x_{n-1}, \quad (2.2.2.1)$$

where signal x and parameter $0.9 \leq \alpha \leq 1.0$ are the inputs and y is the signal after preemphasis. Figure 2.7 shows a digital signal before and after the preemphasis is applied, where the value of α is set to 0.97.

Windowing

Since the goal of this feature extraction is to get spectral features that help us build a phone or sub-phone classifier, we do not want to extract our features from an entire recording of a sentence or whole conversation. In such long pieces, the spectrum changes very quickly. Technically, speech in a *non-stationary* signal, meaning that its statistical properties are not constant across time. Therefore, we need to extract spectral features from a small *window* of speech that characterizes a single sub-phone. For this window, we can make a rough assumption that the signal is stationary (i.e., its statistical properties are constant within the window).

The windowing process takes three parameters:

- the **width** of the window (in milliseconds),
- the **offset** between successive windows, and
- the **shape** of the window.

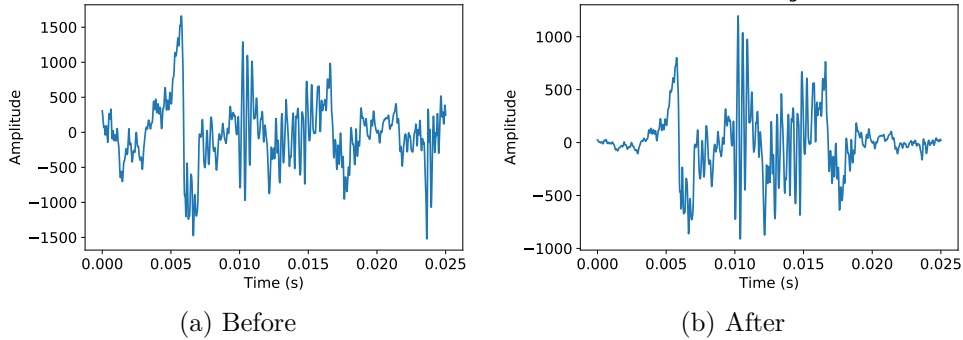


Figure 2.8: A frame before and after the Hamming window function application

Each piece of speech extracted by the window function is called a *frame*, and its length in milliseconds the *frame size*. A number of milliseconds between the left edges of two consecutive frames (their overlap) is *frame shift*. Figure 2.8 shows the 40th frame of the original signal before and after the application of the Hamming function. We can see that the length of each frame is 25 ms.

The third parameter, the shape of the window, can provide additional transformation to the frame. The most basic type of window is the *rectangular* window, which keeps the signal as is. That can cause problems, however,

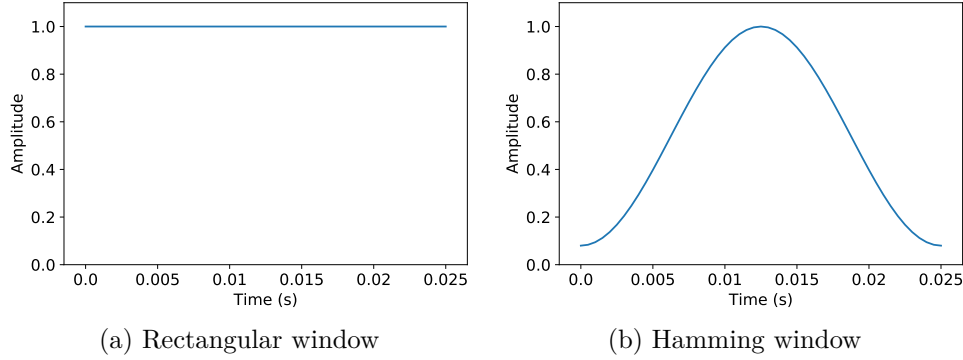


Figure 2.9: Window functions

because it abruptly cuts the signal at the window edges. Such discontinuities create issues when we use the Fourier transform. That is the reason a more typical window used in MFCC extraction is the *Hamming* window. It shrinks the values of the signal near the boundaries toward zero, avoiding discontinuities. Figure 2.9 shows both the rectangular window function and the Hamming window function. We can see how the left and right edges are being squished close to the zero when the function is applied in the case of the Hamming window.

Discrete Fourier Transform

The next step is extracting the spectral information from our windowed signal – we need to know the magnitude of energy the signal contains at different frequency bands. The tool used to extract spectral information for discrete frequency bands for a discrete-time (sampled) signal is the *Discrete Fourier Transform* (DFT).

The input for the DFT is signal $x_m \dots x_n$, and the output is a complex number X_k which represents the magnitude of that frequency in the original signal (frame). Plotting the magnitude against the frequency visualizes the *spectrum*. Figure 2.10 shows such spectrum computed for the same frame, on which was the hamming window function demonstrated.

The definition and the intuitive idea of the discrete Fourier transform is introduced in section 2.2.1. Popular algorithm for the computation of the DFT is the Fast Fourier Transform (FFT). This implementation is very efficient with one constraint – N must equal to values which are powers of two.

Mel filter banks

This part of the extraction process tries to simulate the human hearing, which is not equally sensitive at all frequencies. It is, in fact, less sensitive at higher frequencies. Modeling this property of human hearing in the feature extraction

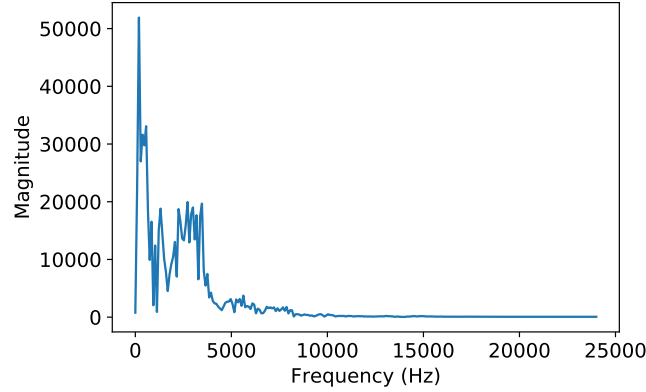


Figure 2.10: A frame spectrum computed with DFT

process turned out to improve the speech recognition performance. It is done by warping the output from DFT onto the *mel* scale. A *mel* is a unit of pitch defined so that an equal number of mels separates a pair of sounds that are perceptually equidistant in tone. In essence, two pairs of adjacent notes in different octaves would have different distances in hertz, yet the same distance in mels. Mels can be computed from raw frequency as follows:

$$mel(f) = 1127 \ln \left(1 + \frac{f}{700} \right), \quad (2.2.2.2)$$

and therefore computing frequency from mels:

$$f(m) = 700 \left(e^{\frac{m}{1127}} - 1 \right). \quad (2.2.2.3)$$

There are ten filters spaced linearly below 1000 Hz, and the remaining filters logarithmically above 1000 Hz. Figure 2.11 depicts such a bank of triangular filters.

In figure 2.12, there is visualized a spectrogram of the original signal. It is 152 feature vectors (one vector for each frame) of length 26 (number of triangular filters in the bank). The spectrogram is normalized by mean normalization to make the image clearer.

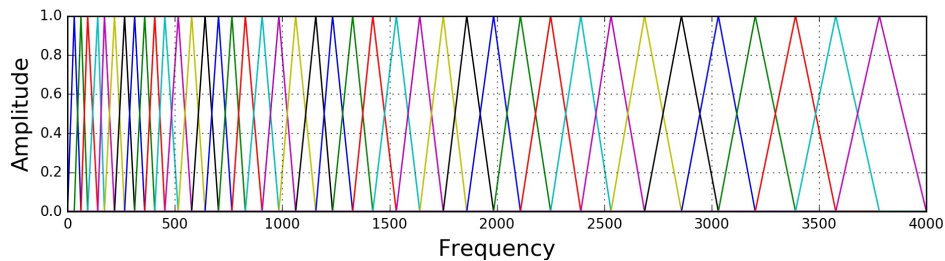


Figure 2.11: Mel filter banks by [25]

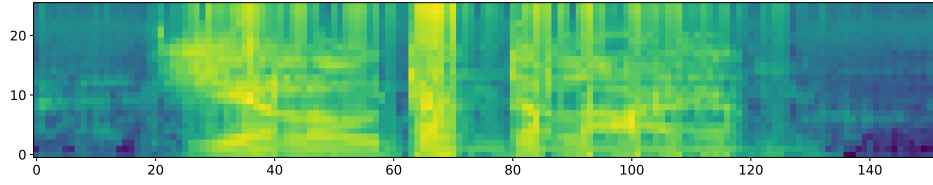


Figure 2.12: Spectrogram of the original signal

Cepstrum

It would be possible to use the mel spectrum by itself as a feature representation for phone detection. Yet, there is still a lot of unimportant information in the so-far extracted features. One way to further refine the features is the computation of the so-called *cepstrum*.

The cepstrum can be thought of as a *spectrum of the log spectrum*. We take the computed spectrum and apply a logarithmic function to it. Now, we visualize the logarithmic spectrum as if it were an ordinary waveform – its domain is still frequency, but we pretend the domain is time. And the last step is computing the spectrum of this pseudo-waveform we have just created. Generally, we only take the first 12 cepstral features, that solely represent the information of the human vocal tract without unnecessary information. Figure 2.13 visualizes the extracted cepstral features from the original signal.

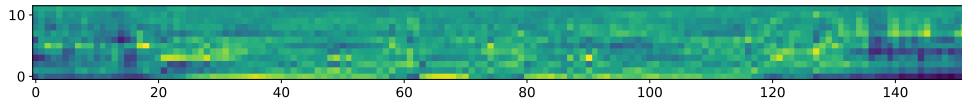


Figure 2.13: Cepstrum features of the original signal

The formal definition for the cepstrum is *inverse DFT of the log magnitude of the DFT of a signal*. Mathematically:

$$c_n = \sum_{n=0}^{N-1} \log \left(\left| \sum_{k=0}^{N-1} x_k e^{-\frac{i2\pi}{N} kn} \right| \right) e^{\frac{i2\pi}{N} kn}. \quad (2.2.2.4)$$

It turns out that cepstral coefficients have the handy property that the variance of the different coefficients at different frequency tends to be uncorrelated. That is not true for the spectrum. Spectral coefficients at various frequency bands are correlated. That significantly eases the job for the acoustic model.

Deltas and energy

We have got 12 features extracted from each frame. Now, we next add the thirteenth feature: the energy from the frame. The energy correlates with

phone identity and therefore is useful for phone detection (vowels and sibilants have more energy than stops, etc.). It is the power sum over time in the frame. Thus, for a signal x in a window from time t_1 to time t_2 :

$$E = \sum_{t=t_1}^{t_2} x_t^2. \quad (2.2.2.5)$$

In the next step, since the speech signal is not constant from frame to frame, we are adding features related to the change in cepstral features over time. We do this by adding for each of the current 13 features a delta or velocity feature, and a double delta or acceleration feature.

Each of the 13 delta features represents the change between frames in the corresponding cepstral/energy feature. Each of the 13 double delta features represents the change between frames in the corresponding delta features. A simple way to compute the deltas is by calculating the difference between frames. Thus, the delta value d_t for a particular cepstral value c_t at time t is:

$$d_t = \frac{c_{t+1} - c_{t-1}}{2}. \quad (2.2.2.6)$$

Usually, a much more sophisticated estimate of the slope, which uses a broader context of frames, is used.

Summary

Now, when we have all the features computed, we are at 39 MFCC features for each frame. In table 2.1, we can see the summary of the features with counts for each category.

Table 2.1: MFCC features summary

12	cepstral coefficients
12	delta cepstral coefficients
12	double delta cepstral coefficients
1	energy coefficient
1	delta energy coefficient
1	double delta energy coefficient
39	MFCC features

Again, the most useful fact about MFCC features is that the features are uncorrelated, which turns out to make the acoustic model much simpler.

2.2.3 Wavelets in feature extraction

Since the usage of wavelets in feature extraction is still considered relatively new, and most of the models still use MFCC vectors in practice, we are going to

introduce a few of those approaches in this section superficially. Some research papers propose to substitute the inverse Fourier transform (or the discrete cosine transform) in MFCC vectors with the discrete wavelet transform (DWT) [26]. Those features are called MFDWC vectors. Another research suggests first applying the DWT to the signal to extract the approximation coefficients and then computing MFCCs from those coefficients [27]. The third extraction method is the *wavelet packet decomposition* (WPD), or sometimes called the *wavelet packet transform* (WPT), which is one of the implementations for the DWT [28, 29].

MFDWC features

MFDWC stands for *mel-frequency discrete wavelet coefficients*, and it was proposed by [26]. Those features are obtained by applying the DWT to the mel-scaled log filterbank energies of a speech frame, thus not using the discrete cosine transform (DCT) as in MFCCs. Using the DWT has the benefit of its localization properties, which are described in section 2.2.1. This approach is trying to deal with two drawbacks of the DCT. Again, this has already been described above, and we will only recap the main points. These are the drawbacks:

1. the DCT basis covers all frequency bands. That means that a corruption of the speech signal on one band affects all the MFCCs. However, if we use a localized transform instead of DCT, corruption could be reduced to only a few coefficients.
2. One frame of speech may (and very probably will) contain information about two adjacent phonemes. If one of them is voiced and the other is not, then the voiced phoneme might dominate the low-frequency spectrum and vice versa. Since MFCCs assumes a stationary signal in one frame, it also assumes there can only be one phoneme present in one frame.

MFDWCs consist of fifteen static feature coefficients for each frame. There are eight coefficients at scale four, four coefficients at scale eight, two at scale sixteen, and one at scale thirty-two. The feature vector also includes delta coefficients and delta energy to include dynamic features. The delta extraction is described in section 2.2.2 for the MFCCs.

These proposed features performed better than other speech features tested in the paper. Those included SUB (subband-based), MULT (multi-resolution), and MFCC features. MFDWCs had relatively better performance on a noisy signal, and the difference from other approaches on a clean signal was lower. The experiments were conducted on the TIMIT database on a phoneme recognition task.

DWT based MFCCs

This approach, proposed in [27], computes the MFCC feature vectors from the approximation coefficients of the DWT of a speech signal. Again, the reason to use the wavelet transform is to reduce the impact of noise in a signal.

The extraction consists of three phases:

1. Use a voice activity detection (VAD) model to separate speech and non-speech (utterances and silence between words) segments,
2. compute the discrete wavelet transform of the speech signal, and get the approximation coefficients, and finally
3. compute the MFCCs from the approximation coefficients.

This model was tested on the PIEAS speech database. It contains 2000 short samples of both telephonic (noisy) and non-telephonic (clean) speech of both males and females. Its goal was speaker recognition, and these DWT based MFCCs outperformed LPC (linear predictive coding), especially in noisy speech. Comparison with standard MFCCs was unfortunately not present.

Wavelet packet decomposition

Wavelet packet decomposition (WPD) is an implementation of the DWT, which decomposes both sides of approximation and detail coefficients. Therefore, it is creating a whole binary tree as opposed to the classic DWT, which decomposes only one branch in each stage. Such a tree is depicted in figure 2.14.

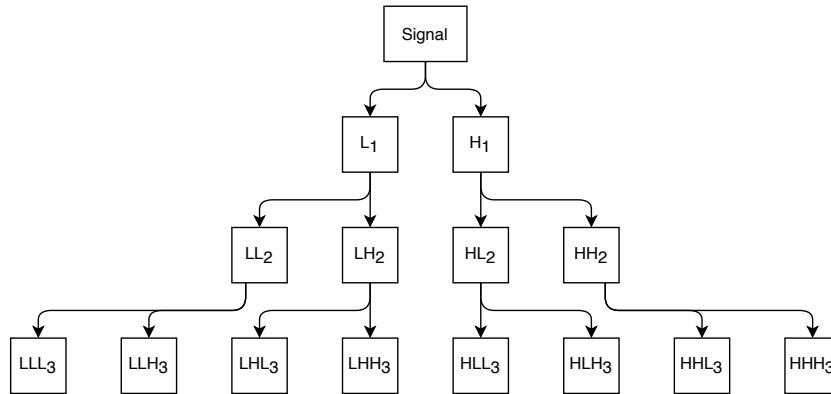


Figure 2.14: Wavelet packet decomposition tree, inspired by [30]

The WPD produces 2^n sets of coefficients (for each leaf node) as opposed to the standard DWT, where n is the level of decomposition. The standard DWT provides only $n + 1$ coefficients – one for each level from the high-pass filters, and +1 from the last level as the output from the low-pass filter. That

lets us choose from a wider variety of coefficients and experiment with different combinations.

Research [29] experiments with the recognition rate of Indonesian syllables, and uses three datasets: a training syllable dataset DLSK, a testing syllable dataset DUSK, and a testing word dataset DUK. It compares MFCC feature vectors and WPD (with DB3 and DB7 wavelets) as a feature extraction method on a *hidden Markov model* (HMM) with Euclidean distance as a classifier. The best results for the MFCC method are 75/38/47 % for the DLSK/DUSK/DUK datasets, respectively. As for the WPD approach, the best results are 70.8/37.5/53.1 %. Therefore, it depends on a specific problem, which one of the strategies to use.

Research [28] deals with speaker recognition over two datasets. The first one is a custom dataset created in laboratory conditions, and it consists of 783 recorded samples of both male and female speech. These samples are several minutes long, and there are ten speakers. The second dataset consists of human dialogues of fifty speakers created for Czech Railways.

There are four conventional methods compared with the WPD method: MFCCs, *perceptual linear prediction coefficients* (PLPC), *linear prediction reflection coefficients* (LPREFC), and *linear prediction cepstral coefficients* (LP-CEPSTRA). All of these were tested on two models: a *Gaussian mixture model* (GMM) and a *multilayer perceptron* (MLP).

The WPD, in this case, uses three discrete wavelets: *Daubechies*, *Symlets*, and *Coiflets*. It computes the *fifth* level of decomposition, which therefore outputs 2^5 sub-signals. Those sub-signals are then used to calculate the final feature vectors as follows:

$$f_i = \log_{10} \left(\frac{1}{N_i} \sum_{k=1}^{N_i-1} |(w_i(k))^2 - w_i(k-1) \cdot w_i(k+1)| \right), \quad (2.2.3.1)$$

where f_i are values of the feature vector (one value for each leaf node, thus $2^5 = 32$ values), $w_i(k)$ is the k -th coefficient in leaf node i , and N_i is the number of coefficients in leaf node i .

The results are:

- MLP requires less training data per speaker than GMM to classify speakers correctly. In this case, it needed two times fewer data.
- The minimal utterance duration for correct classification is two seconds for both models and all feature extraction methods.
- The accuracy results are dependent on the model:
 - GMM: all three wavelet families had very similar results (with the Symlets in the lead), and they were all slightly outperformed by other conventional methods.

- MLP: all methods were clustered together with LPREFC and LP-CEPSTRA in the lead. The only exception were MFCCs, which were far behind every other approach.

The choice of feature extraction method again seems to be dependent on the task we want to perform, and on the dataset the model needs to work with. All of the presented approaches had very similar results, except for the significant negative deviation of MFCC features in the MLP model.

2. AUTOMATIC SPEECH RECOGNITION

Test [31, 32]

Experiments

Realization

Conclusion

Bibliography

- [1] Haykin, S. *Neural Networks and Learning Machines (3rd Edition)*. Pearson Education India, 2010, ISBN 0131471392.
- [2] Biga, L. M.; Dawson, S.; et al. Anatomy & Physiology, Nervous tissue. Online; accessed 18-February-2020. Available from: <https://open.oregonstate.edu/aandp/chapter/12-2-nervous-tissue/>
- [3] Kuželá, O. *Neural networks with memory*. Master's thesis, CTU in Prague, FIT, Thákurova 9, May 2016.
- [4] Patel, N. Data Mining – Artificial Neural Networks. Online; accessed 25-February-2020. Available from: <https://ocw.mit.edu/courses/sloan-school-of-management/15-062-data-mining-spring-2003/lecture-notes/NeuralNet2002.pdf>
- [5] Pal, S. K.; Mitra, S. Multilayer perceptron, fuzzy sets, classification. 1992.
- [6] Brilliant.org. Backpropagation. Online; accessed 10-March-2020. Available from: <https://brilliant.org/wiki/backpropagation/>
- [7] Shorten, C. Introduction to ResNets. Online; accessed 17-March-2020. Available from: <https://towardsdatascience.com/introduction-to-resnets-c0a830a288a4>
- [8] Wang, C.-F. The Vanishing Gradient Problem. Online; accessed 17-March-2020. Available from: <https://towardsdatascience.com/the-vanishing-gradient-problem-69bf08b15484>
- [9] Venkatachalam, M. Recurrent neural networks. Online; accessed 16-March-2020. Available from: <https://towardsdatascience.com/recurrent-neural-networks-d4642c9bc7ce>

- [10] Olah, C. Understanding LSTM Networks. Online; accessed 18-March-2020. Available from: <https://colah.github.io/posts/2015-08-Understanding-LSTMs/>
- [11] learn neural networks. Recurrent neural networks. Online; accessed 16-March-2020. Available from: <https://learn-neural-networks.com/recurrent-neural-networks/>
- [12] learn neural networks. Long-short-term memory (LSTM) networks. Online; accessed 17-March-2020. Available from: <https://learn-neural-networks.com/lstm-networks/>
- [13] Zhou, G.-B.; Wu, J.; et al. Minimal gated unit for recurrent neural networks. *International Journal of Automation and Computing*, volume 13, no. 3, 2016: pp. 226–234.
- [14] Yu, D.; Deng, L. *AUTOMATIC SPEECH RECOGNITION*. Springer, 2016.
- [15] Schaller, R. R. Moore’s law: past, present and future. *IEEE spectrum*, volume 34, no. 6, 1997: pp. 52–59.
- [16] Dictionary.com. Definition of Infotainment. Online; accessed 25-March-2020. Available from: <https://www.dictionary.com/browse/infotainment>
- [17] Rabiner, L. R.; Schafer, R. W.; et al. Introduction to digital speech processing. *Foundations and Trends® in Signal Processing*, volume 1, no. 1–2, 2007: pp. 1–194.
- [18] Hoy, M. B. Alexa, Siri, Cortana, and more: an introduction to voice assistants. *Medical reference services quarterly*, volume 37, no. 1, 2018: pp. 81–88.
- [19] Oppenheim, A. V. Speech analysis-synthesis system based on homomorphic filtering. *The Journal of the Acoustical Society of America*, volume 45, no. 2, 1969: pp. 458–465.
- [20] Davis, S.; Mermelstein, P. Comparison of parametric representations for monosyllabic word recognition in continuously spoken sentences. *IEEE transactions on acoustics, speech, and signal processing*, volume 28, no. 4, 1980: pp. 357–366.
- [21] Burrus, C.; Gopinath, R.; et al. Introduction to Wavelets and Wavelet Transform—A Primer. *Recherche*, volume 67, 01 1998.

-
- [22] Kutz, J. N. *Data-driven modeling & scientific computation: methods for complex systems & big data, PART III, 13.: Time-Frequency Analysis: Fourier Transform and Wavelets*. Oxford University Press, 2013, ISBN 978-0199660346.
- [23] Kutz, N. Time Frequency Analysis & Wavelets. Available from: <https://www.youtube.com/watch?v=ViZYXxuxUKA>
- [24] Pearson Education, Inc. *Chapter 9, Automatic Speech Recognition*. Pearson Prentice Hall, 2008, available from: [http://www.cs.columbia.edu/~julia/courses/CS6998-2019/\[09\]AutomaticSpeechRecognition.pdf](http://www.cs.columbia.edu/~julia/courses/CS6998-2019/[09]AutomaticSpeechRecognition.pdf).
- [25] Fayek, H. Speech Processing for Machine Learning: Filter banks, Mel-Frequency Cepstral Coefficients (MFCCs) and What's In-Between. 2016, online; accessed 6-April-2020. Available from: <https://haythamfayek.com/2016/04/21/speech-processing-for-machine-learning.html>
- [26] Gowdy, J. N.; Tufekci, Z. Mel-scaled discrete wavelet coefficients for speech recognition. In *2000 IEEE International Conference on Acoustics, Speech, and Signal Processing. Proceedings (Cat. No. 00CH37100)*, volume 3, IEEE, 2000, pp. 1351–1354.
- [27] Malik, S.; Afsar, F. A. Wavelet transform based automatic speaker recognition. In *2009 IEEE 13th International Multitopic Conference*, IEEE, 2009, pp. 1–4.
- [28] Král, P. Discrete Wavelet Transform for automatic speaker recognition. In *2010 3rd International Congress on Image and Signal Processing*, volume 7, IEEE, 2010, pp. 3514–3518.
- [29] Hidayat, S. Speech recognition of KV-patterned Indonesian syllable using MFCC, wavelet and HMM. *Kursor*, volume 8, no. 2, 2016: pp. 67–78.
- [30] Muo, U.; Madamedon, M.; et al. Wavelet packet analysis and empirical mode decomposition for the fault diagnosis of reciprocating compressors. 09 2017, pp. 1–6, doi:10.23919/IConAC.2017.8082065.
- [31] Hannun, A.; Case, C.; et al. Deep Speech: Scaling up end-to-end speech recognition. 2014, 1412.5567.
- [32] Abadi, M.; Agarwal, A.; et al. TensorFlow: Large-Scale Machine Learning on Heterogeneous Systems. 2015. Available from: <https://www.tensorflow.org/>

Acronyms

ANN	Artificial neural network
ASR	Automatic speech recognition
CPU	Central processing unit
CWT	Continuous wavelet transform
DCT	Discrete cosine transform
DFT	Discrete Fourier transform
DNN	Deep neural network
DWT	Discrete wavelet transform
FFT	Fast Fourier transform
GMM	Gaussian mixture model
GPS	Global positioning system
GPU	Graphics processing unit
GRU	Gated recurrent unit
LPC	Linear predictive coding
LPCEPSTRA	Linear prediction cepstral coefficients
LPREFC	Linear prediction reflection coefficients
LSTM	Long-short-term memory (network)
MLP	Multilayer perceptron
MFCC	Mel-frequency cepstral coefficients

A. ACRONYMS

MFDWC Mel-frequency discrete wavelet coefficients

MULT Multi-resolution (features)

NN Neural network

PLPC Perceptual linear prediction coefficients

ResNet Residual network

RNN Recurrent neural network

STFT Short-time Fourier transform

SUB Subband-based (features)

VAD Voice activity detection

WPT Wavelet packet transform

XML Extensible markup language

Contents of enclosed CD

	readme.txt	the file with CD contents description
	exe	the directory with executables
	src	the directory of source codes
	wbdcm	implementation sources
	thesis	the directory of \LaTeX source codes of the thesis
	text	the thesis text directory
	thesis.pdf	the thesis text in PDF format
	thesis.ps	the thesis text in PS format