

University of Applied Sciences Würzburg-Schweinfurt Faculty of Computer Science and
Business Information Systems

Bachelor-Thesis

Evaluation and implementation of gradient descent algorithms on streaming data

**submitted to the University of Applied Sciences Würzburg-Schweinfurt in the
Faculty of Computer Science and Business Information Systems to achieve the
Bachelor of Engineering degree in 'Information Systems'**

Florian Hohn

Submitted on: 25.02.2020

First Reader: Prof. Dr. Frank-Michael Schleif
Second Reader: Moritz Heusinger

Summary

TODO

Abstract

TODO

Note of thanks

Contents

1	Introduction	1
1.1	Inspiration for this Thesis	1
1.2	Purpose of this Thesis	1
1.3	Outline of the Thesis	1
2	Basics	3
2.1	Machine Learning	3
2.1.1	Supervised Learning	5
2.1.2	Gradient Descent	7
2.2	Streaming Data	9
2.2.1	Increasing Volume of the Data	9
2.2.2	Drift in Streaming Data	9
2.3	Criteria for the Comparision	11
2.3.1	Evaluations for Data Streams	11
2.3.2	Performance Evaluation	13
2.3.3	Accuracy Evaluation	14
2.4	Framework	15
3	Used Algorithms	18
3.1	Theoretical Work	18
3.1.1	Learning Vector Quantization	18
3.1.2	Robust Soft Learning Vector Quantization	20
3.2	Momentum Based Gradient Descent	23
3.2.1	Momentum	23
3.2.2	Adagrad	24
3.2.3	Adadelata	26
3.2.4	RMSprop	27
3.2.5	Adam	28
3.3	Implementation of the RSLVQ variants	29
3.3.1	RSLVQ SGD	29
3.3.2	RSLVQ Adadelata	32
3.3.3	RSLVQ RMSprop	34
3.3.4	RSLVQ Adam	35

4	Test realisation	36
4.1	Used Streaming Datat sets	36
4.1.1	Synthetic data	36
4.1.2	Real world data	36
4.2	Implementation of the Comparision and Results	37
4.2.1	Experiment begrenzung	37
4.2.2	Experiment aufbau	37
4.2.3	Experiment durchlauf	37
5	Result Evaluation	38
6	Summary	39
	Table of Contents	41
	Literature	45
	Statutory Declaration	46

1 Introduction

The following chapter will give the reader a short explanation what inspired the work for this thesis, what its purpose is and what it wishes to achieve, as well as an explanation how this work is structured, to give the reader a better orientation for how to navigate this work.

1.1 Inspiration for this Thesis

1.2 Purpose of this Thesis

1.3 Outline of the Thesis

In the second Chapter of this thesis will give a short overview over the basics of Machine Learning that are relevant for the understanding of the premise for the thesis. It will give also a short description about what is special about Data stream in Machine learning algorithms and what is important to take into account when working with them. It will also give a short introduction about the criteria with which the comparison of the different implementations of the Algorithm take place and why these criteria are used for comparison. In the last part of the first chapter will be the framework introduced that was used for implementation and comparison.

The next chapter will then continue to give a deeper look into the deeper parts of the subsection of Machine learning and the Algorithm that is the basis for the different variation of implementation for the Algorithms. It will also give a short introduction about the version that this algorithm was based on, as this is important to later understand the variations and differences in the other implementations. In the last part of that chapter will then explain the different methods and the implementations of the variants in greater detail.

The fourth chapter will then list and describe the used synthetic and real world streaming

1 Introduction

data sets that were used to perform the experiments. Furthermore will be describe how the experiments were conducted and with which methods the results were controlled.

Following this, the next chapter will be dedicated for the result evaluation what what we informations were hopefully gained from conducting these experiments for the comparison of these implementations.

The last chapter will summarize all the work that was done and give a final evaluation to the algorithms.

2 Basics

The Goal of this Chapter is to give the Reader a basic understanding of how Machine Learning Algorithms work. Furthermore, it will describe what the difference is between unsupervised and supervised learning.

There will be also an introduction and description on how the Gradient Descent Algorithm, that is the basis for the Modified versions of the Algorithms that will be compared in this work, works and how it decides what to do. After this is done there will be also a short description on what Streaming Data is and on the framework that is used to realise the implementation of the Experiments.

Finally there will be an explanation on how the criteria for the comparison of the Algorithms were chosen and how these criteria work.

2.1 Machine Learning

The term 'Machine Learning' stands for a collection of self-learning algorithms, that learn from input Data to make fast and correct decisions and predictions. It is important to remember that it is not a real Artificial intelligence, but rather a sequence of statistical analyses on given data, with the goal to make gradual improvements to the prediction models of the algorithm. A Machine can do this optimization a lot faster and effectively than a human, especially if this optimization has to be done in real time on big data input sets. This gives also a bigger flexibility, as a machine that can learn from a changing environment needs lesser redesigns, which leads to less time that gets wasted and can be used more productively on other tasks.

All these are reasons why Machine Learning has an ever increasing importance in the field of computer science, but also a growing impact on everyday life. It is thanks to Machine learning that there are intelligent assistant softwares that can recognize voice commands and questions (Apples SIRI, Amazon Alexa or Microsofts Cortana), e-mail spam filter, nearly self-driving cars, reliable search-engines, faster and more precise weather forecasts and challenging game ai's nowadays

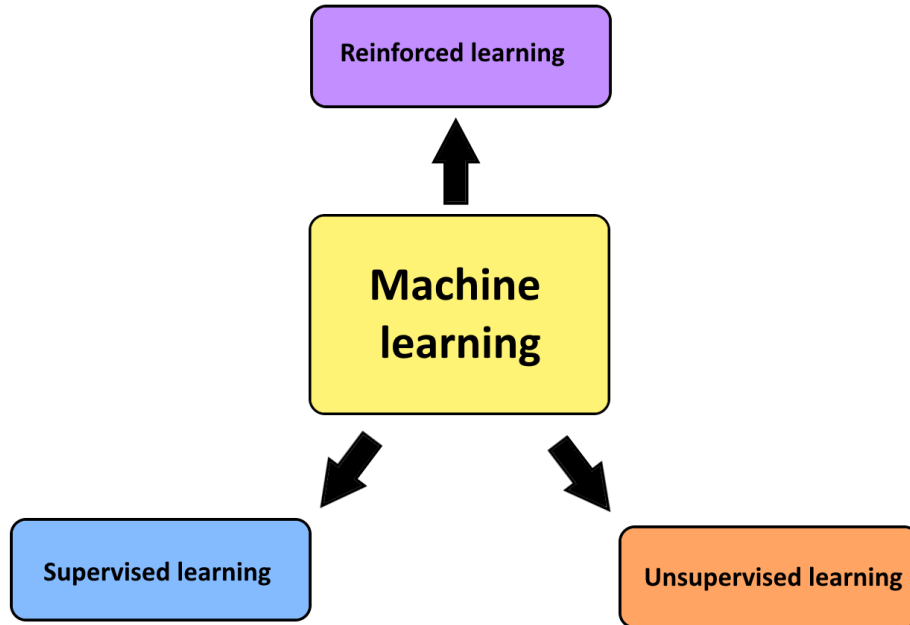


Figure 2.1: Overview of the Machine Learning Strategy's.

Machine Learning Algorithms can be categorisized depending on how they learn from the available training Data to create a model for the predictions. There are three main subdivisions, under wich these algorithms can fall. The first one is the Supervised learning strategy, then the Unsupervised learning strategy and the Reinforced learning strategy.

Only the Supervised Learning strategy will be explained in more detail on the following pages, as the focus of this thesis lays in the comparision of several different implemen-tations of classification algorthims that are based on Supervised Learning.

[19]

2.1.1 Supervised Learning

The idea behind Supervised learning is to create a model from a already labeled training data set, that then can make prediction and decision on new or unknown data based on the learned model from the training set. The following example will help to understand this easier. Lets say a fruit farmer wants to automatically sort the harvested Fruits in either categorie A or B for the market. The sorting device is equiepped with 2 sensors to measure 2 features, one for colour and another one for the size of the fruit. It then decides to wich of the two classes the fruit belongs to. As this is a System that is capable of dividing features into a discrete number of classes the system can be called a *classifier*. This is a typical exampel of an *classification task*. To configure the System in such a way that it correctly sorts the fruits into the right category, a specailist hand picks several fruits from a small test set and then sorts them depending on the already named features. In other words, the specailist classiefies them. Based on this by hand sortet training set the machine will then later decide how the other fruits should be sorted. It is used as a *model* for the classification. [10] Another type of Supervised learning is the *regression task*. The difference to the *classification task* is, that instead of a *discrete* number of classes that gets returned, a *continous* number of classes will be returned. An example of a *regressive task* is, for example, the prediction of the Stock market or the price of medications.

The Process as seen in figure 2.2 and as described as in the example of the farmer and his sorted fruits shows how the supervised learning works. The Algorithm (marked in green) learns with the help of the labeled training data. The result of this learning is then that the trained Algorithm, now called predictive model, can be used in a production enviroment. New, unlabeled, data can then be send into the predictiv model to get a predictied model as an output. [22]

To better understand how this predictiv model is formed, one has to understand how the classification task works. In the classification task the computer programm is asked to specify which of k catergories a input belongs to. To solve this task, the learning algorithm is usually asked to produce a function $f : \mathbb{R} \rightarrow 1, \dots, k$. When $y = f(\mathbf{x})$, the model assigns an input described by vector \mathbf{x} to a category identified by numeric code y . There are other variants of the classification task, for example, where f outputs a probability distribution over classes. An example of a classification task is object recognition, where the input is an image (usually described as a set of pixel brightness values), and the output is a numeric code identifying the object in the image. [12]

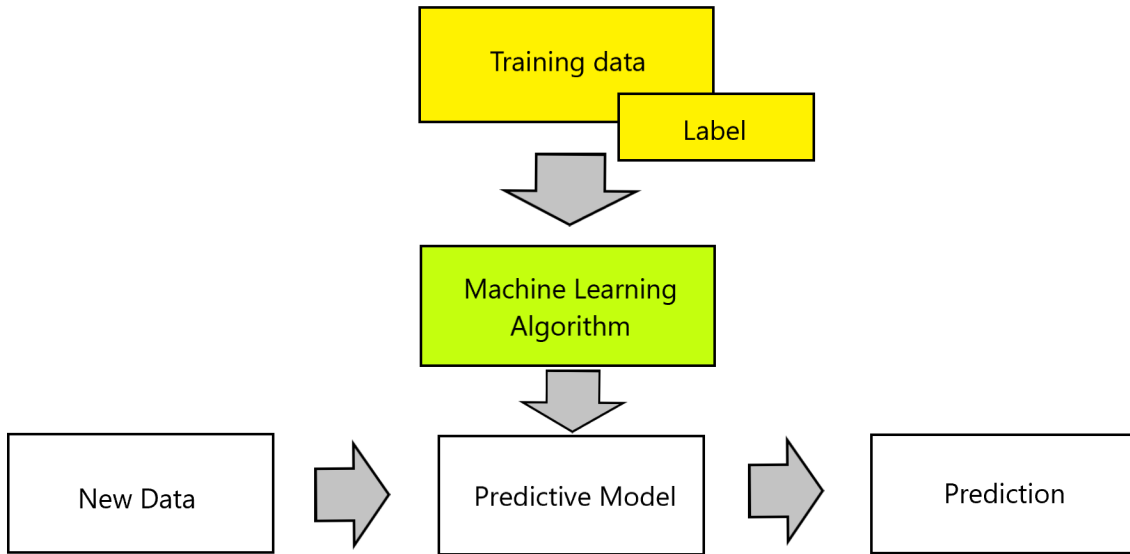


Figure 2.2: Process of the Supervised Learning. [22]

Then there must be also differentiated between the type of data that is put into the classification process. In batch or offline classification, a classifier-building algorithm is given a set of labeled examples. The algorithm creates a model, a classifier in this case. The model is then deployed, that is, used to predict the label for unlabeled instances that the classifier builder never saw. If we go into more detail, we know that it is good methodology in the first phase to split the dataset available into two parts, the training and the testing dataset, or to resort to cross-validation, to make sure that the classifier is reasonably accurate. But in any case, there is a first training phase, clearly separated in time from the prediction phase.

In the online setting, and in particular in streaming, like it is this thesis, this separation between training, evaluating, and testing is far less clear-cut, and is interleaved. We need to start making predictions before we have all the data, because the data may never end. We need to use the data whose label we predict to keep training the model, if possible. And probably we need to continuously evaluate the model in some way to decide if the model needs more or less aggressive retraining. [5]

2.1.2 Gradient Descent

The Gradient descent Algorithm is one of the most popular Algorithms to perform optimizations on and is by far the most common way to optimize neural networks. This is also the reason why nearly every state-of-the-art Deep learning library has various forms of implementations to perform optimization on the gradient descent. The Problem, however is, that these Algorithms often perform as some form of black-box, because a practical explanation of what their strengths and weaknesses is are hard to find and explain. Gradient descent is a way to minimize an objective function $J(\theta)$ parameterized by a model's parameters $\theta \in \mathbb{R}^d$ by updating the parameters in the opposite direction of the gradient of the objective function $\nabla_{\theta} J(\theta)$ with respect to the parameters. The learning rate η determines the size of the steps we take to reach a local minimum. In other words, we follow the direction of the slope of the surface created by the objective function downhill until we reach a valley. [23]

The normal gradient descent, also known as batch-gradient descent, computes the gradient of the function with respect to the parameters θ for the entire training set:

$$\theta = \theta - \eta * \nabla_{\theta} J(\theta) \quad (2.1)$$

Since we need to calculate the gradients for the complete dataset to get just one update, batch gradient descent can get very slow and is intractable for dataset that do not fit into the memory. For this reason the batch gradient descent implementation can not be used in this comparison, as it is not able to perform updates "on-the-fly", as this is a necessity for use on streaming Data.

In contrast to that, the Stochastic Gradient Descent (SGD) can perform a update for *each* training example $x^{(i)}$ and label $y^{(i)}$:

$$\theta = \theta - \eta * \nabla_{\theta} J(\theta; x^{(i)}; y^{(i)}) \quad (2.2)$$

While the batch gradient descent performed redundant computations for large datasets, as it recomputes gradients for similar examples before each update, the SGD puts away this redundancy by performing one update at a time. Because of that it performs much faster and can be used on Streaming Data. The SGD performs frequent updates with a high variance, that causes the objective function to fluctuate heavily.

This fluctuation enables the SGD to jump to new and potentially better local minima, while the batch gradient only converges to the minimum of the basin the parameters are placed in. But this also ultimately complicates the convergence to the exact minimum, as

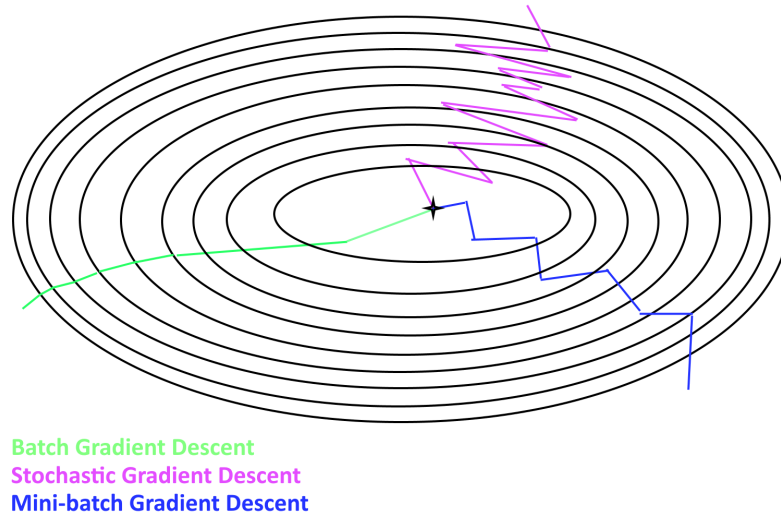


Figure 2.3: Visualisation differences of the three Gradient Descent Algorithms

the SGD keeps overshooting. This can be fixed by slowly decreasing the learning rate, as the SGD then shows the same convergence behaviour as the batch gradient descent. It then almost certainly converges to a local or global minimum for non-convex and convex optimization respectively. [23] These advantages lead to the fact that the SGD is the most common Gradient descent algorithm.

Then there is also the Mini-batch gradient descent algorithm. It takes the best of both worlds and performs an update for every mini-batch on n training examples:

$$\theta = \theta - \eta * \nabla_{\theta} J(\theta; x^{(i:i+n)}; y^{(i:i+n)}) \quad (2.3)$$

This way it reduces the variance of the parameter updates, which can lead to a more stable convergence. It can also make use of highly optimized matrix optimizations common for state-of-the-art deep learning libraries, that makes computing the gradient with respect to a mini-batch very efficient. Common mini-batch sizes range between 50 and 256, but can vary for different applications. Typically the Mini-batch gradient descent is the algorithm of choice when there is a neural network to train and the term SGD is usually also employed even when a mini-batches are used. [23]

2.2 Streaming Data

Data streams are an algorithmic abstraction to support real-time analytics. They are sequences of items, possibly infinite, each item having a timestamp, and so a temporal order. Data items arrive one by one, and we would like to build and maintain models, such as patterns or predictors, of these items in real time. There are two main algorithmic challenges when dealing with streaming data: the stream is large and fast, and we need to extract information in real time from it. That means that usually we need to accept approximate solutions in order to use less time and memory. Also, the data may be evolving, so our models have to adapt when there are changes in the data. [5]

2.2.1 Increasing Volume of the Data

When the Data items keep arriving fast and in large quantities an algorithm is no longer able to process the data efficiently by using multiple passes. Instead, it has to be designed in such a way that it is able to process the data with one pass [1]. [5] describes five requirements that an algorithm has to fulfill, if it should be able to process streaming data successfully:

- Process an instance at a time, and inspect it (at most) once.
- Use a limited amount of time to process each instance.
- Use a limited amount of memory.
- Be ready to give an answer (prediction, clustering, patterns) at any time.
- Adapt to temporal changes (Concept drift).

The last point, "Adapt to temporal changes", has to do with evolving data streams and will be described further in the following section.

2.2.2 Drift in Streaming Data

Often there is an inherent temporal component to the stream mining process. This is because the data may evolve over time. This behaviour of a data stream is known as temporal locality[1] (also known as Concept drift)[5]. Therefore, a straightforward adaptation of one-pass mining algorithms may not be an effective solution to the task. Thus, Stream mining algorithms need to be carefully designed with a clear focus on the

evolution of the underlying data, which leads to the fifth requirement (see enumeration 2.2.1) described in [5].

Now, to be more specific, the term Concept drift will be described further. It does not mean that the items that are observed today are not exactly the same as those that were observed yesterday. A more reasonable notion is that statistical properties of the data change more than what can be attributed to chance fluctuations. To better understand this, it helps to assume that the data is in fact the result of a random process that at each time generates an item according to a probability distribution that is used at that exact time, and that may or may not be the same that is used at any other given time. There is no change (drift) when this underlying generating distribution remains the same. Drift occurs whenever it varies from one time step to the next [5].

Often, an additional independence assumption is implicitly or explicitly used: that the item generated at time t is independent of those generated at previous time steps. This assumption is often incorrect or false, because in many situations the stream has memory, and experiences bursts of highly correlated events. For example, a fault in a component of a system is likely to generate a burst of faults in related components [5].

Although Concept drift in the item distribution may be arbitrary, it helps to name a few generic types, which are not exclusive within a stream:

- *Sudden* drift occurs when the distribution has remained unchanged for a long time, then changes in a few steps to a significantly different one. It is often called *shift*.
- *Gradual* or *incremental* drift occurs when, for a long time, the distribution experiences at each time step a tiny, barely noticeable change, but these accumulated changes become significant over time.
- Drift may be *global* or *partial* depending on whether it affects all of the item space or just a part of it. In ML terminology, partial change might affect only instances of certain forms, or only some of the instance attributes.
- *Recurrent* concepts occur when distributions that have appeared in the past tend to reappear later. An example is seasonality, where summer distributions are similar among themselves and different from winter distributions. A different example is the distortions in city traffic and public transportation due to mass events or accidents, which happen at irregular, unpredictable times.
- In prediction scenarios, we are expected to predict some outcome feature Y of an item given the values of input features X observed in the item. *Real* drift occurs when $Pr[Y | X]$ changes, with or without changes in $Pr[X]$. *Virtual* drift occurs when $Pr[X]$ changes but $Pr[Y | X]$ remains unchanged. In other words, in

real drift the rule used to label instances changes, while in *virtual* drift the input distribution changes [5].

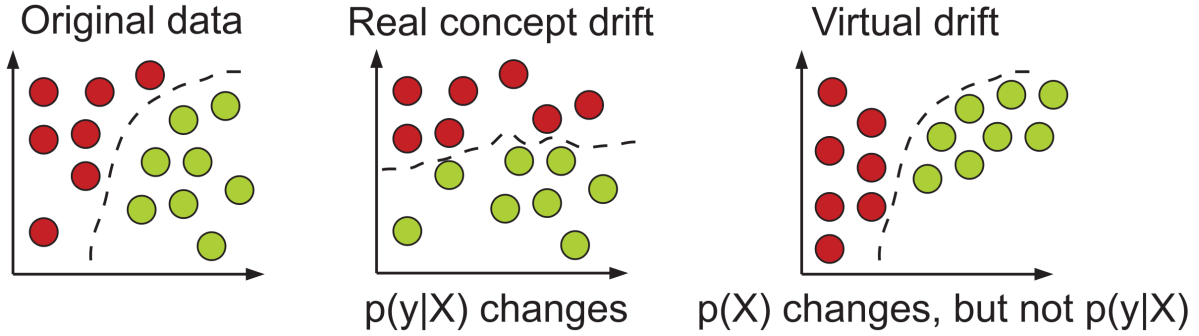


Figure 2.4: Visualization of types of Concept drift [11]

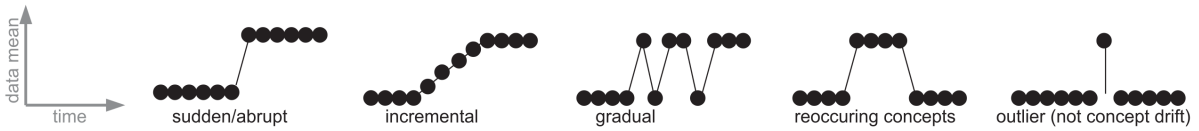


Figure 2.5: Visualization of real and virtual drift [11]

[11]

2.3 Criteria for the Comparision

This section will explain on what criteriuy this thesis will base the comparision of the different implementations. It will also descirbe what the difference between the cireria is and why this is importend to know. It will define what this thesis understands under the terms performance and accuracy, as well as what the differences between the used error-estimations are.

2.3.1 Evaluations for Data Streams

To evaluate a learning algorithmus, it must be know which examples were used to train the algorithmus and which are used to test the model output by the algorithms. Which procedure is used in batch learning has partly depended on data size. For small data sets with ≤ 1000 examples the procedure of cross-validation was well suited, as it made

maximum use of the data that it could use. A Problem arose when the size of the data kept increasing, as this practically created a time limit for procedures that kept repeating the training process to many times. For this reason it is common practice to try to reduce the number of folds and repetitions for big data sets to allow experiments to conclude in reasonable time. One way to achieve this is to take several hundreds of thousand samples in a single batch and take them for a single holdout run, as this requires the least computational effort. The rationalization for this, beside the practical time issues, is that the reliability we may lose by not having repeated runs is compensated by the reliability gained by the sheer number of samples that are used. There are two main approaches that we have to consider for the evaluation process, the first one being Holdout and the second being Prequential. [2]

These two approaches will now be described in more detail:

- **Holdout:** It is often acceptable to measure the performance with a single holdout set when the batch learning would otherwise reach a scale where cross-validation would be too time consuming. This can get very helpful, as the results of different studies can be directly compared to each other. To do this, the division between the train and test set have to be predefined. To track model performance over time, the model can be evaluated periodically, for example, after every one million training examples. Testing the model too often has potential to significantly slow the evaluation process, depending on the size of the test set. A source for holdout data sets are for example sets of data of the stream that have not yet been used to train the learning algorithm. These sets could be taken from "further ahead" in the stream to be used as test sets and could then later be used again to train the learning algorithm again after the testing is complete. This procedure would be recommendable in scenarios with concept-drift, as it would measure a model's ability to adapt to the latest trends in the data. In scenarios without drift, a single static holdout set should be sufficient, as this would avoid varying estimates between test sets. If it is assumed that the test set is independent and sufficiently large in relation to the complexity of the target concept, it can be concluded that it then will provide an accurate measurement of generalization accuracy. [2]
- **Prequential:** An alternative approach to evaluate data stream algorithm is to interleave testing with training. Each individual set can be used to test the model before it is used for training and then the accuracy can be incrementally updated from this. When it is performed in this order the model will always be tested by sets that it has not seen yet. This has the advantage that no holdout set is needed to perform testing. This also results in a smooth plot of the accuracy over time, as each individual set will become increasingly more insignificant to the overall average. This approach also has its disadvantages, as it makes it difficult to accurately separate and measure training and testing times. It also obscures the true accuracy that the algorithm can achieve at a given point, as the algorithm gets punished for mistakes that it made earlier regardless of the level of accuracy they

are eventually capable of. It also has to be said that this effect will diminish over time. With this approach the statistics are updated with every set from the stream and can be recorded at this level of detail if desired. To increase the efficiency a sampling parameter can be used to reduce the required storage for the results, by recording only at periodic intervals like the holdout method. [2]

2.3.2 Performance Evaluation

In real data streams, the number of instances for each class may be evolving and changing. It may be argued that the prequential accuracy measure is only appropriate when all classes are balanced and have approximately the same number of examples. The Kappa statistic is a more sensitive measure for quantifying the predictive performance of streaming classifiers. [5] There are 3 different versions of the Kappa statistic that are relevant for the analysis of streaming data. The first Kappa statistic, Kappa κ , was introduced by Cohen in his work [7] and was defined as follows:

$$\kappa = \frac{p_0 - p_c}{1 - p_c} \quad (2.4)$$

The quantity p_0 is the classifier's prequential accuracy, and p_c is the probability that a chance classifier—one that randomly assigns to each class the same number of examples as the classifier under consideration—makes a correct prediction. If the classifier is always correct, then $\kappa = 1$. If its predictions coincide with the correct ones as often as those of a chance classifier, then $\kappa = 0$ [5]. The second Kappa statistic, κ_m [4], is a measure that compares against a majority class classifier p_m , that is a classifier that stores a count for each of the class labels and predicts as the class of a new instance the most frequent class label, instead of a chance classifier:

$$\kappa_m = \frac{p_0 - p_m}{1 - p_m} \quad (2.5)$$

In cases where the distribution of predicted classes is substantially different from the distribution of the actual classes, the majority class classifier p_m can perform better than a given classifier while the classifier has a positive κ statistic [5]. The third Kappa statistic that exist, is the Kappa temporal statistic [6, 28]. It considers the presence of temporal dependencies in data streams and is defined as:

$$\kappa_{temp} = \frac{p_0 - p'_e}{1 - p'_e}, \quad (2.6)$$

where p'_e is the accuracy of the No-change classifier, that is a classifier that predict the last class in the data stream. It does this by exploiting autocorrelation in the label assignments. It is a simple and useful classifier when the same labels appear together in bursts. Statistic κ_{temp} takes values from 0 to 1. The interpretation is similar to that of κ : if the classifier is perfectly correct, then $\kappa_{temp} = 1$. If the classifier is achieving the same accuracy as the No-change classifier, then $\kappa_{temp} = 0$. Classifiers that outperform the No-change classifier fall between 0 and 1. Sometimes $\kappa_{temp} < 0$, which means that the classifier is performing worse than the No-change baseline [5].

Using κ_{temp} instead of κ_m , we can detect misleading classifier performance for data that is not $\prod D$. For highly imbalanced but independently distributed data, the majority class classifier may beat the No-change classifier. The κ_{temp} and κ_m measures can be seen as orthogonal, since they measure different aspects of the performance [5].

2.3.3 Accuracy Evaluation

Accuracy describes the ability of the classifier to correctly predict the data points in relation to the whole training set, which can be defined as:

$$Accuracy = \frac{Number of correct predictions}{Total number of predictions} \quad (2.7)$$

For a binary classification, accuracy can also be calculated in terms of positive and negatives as follows:

$$Accuracy = \frac{TP + TN}{TP + TN + FP + FN}, \quad (2.8)$$

with TP = True Positives, TN = True Negatives, FP = False Positives and FN = False Negatives. The accuracy is often used to get an measurement of how good the classifier is at returning correct predictions. It should be noted that the accuracy can give a misleading result for cases where there is a high class imbalance. In this case the classifier may achieve a high accuracy, but because the model is too simple the result is too crude to be of use. An example would be, if the Label for A would occur in 99% of the cases, then the prediction that every case would be A would have an accuracy of 99%. This is called the Accuracy Paradox. For this reason measurements like Kappa, $Kappa_m$ and $Kappa_{temp}$ are also required to get a useful result of the performance and accuracy of a classifier [5].

2.4 Framework

The following section will shortly introduce the python based scikit-multiflow framework, with which the the comperssion of the different implementations of the RSLVQ algorithm where done.

Scikit-multiflow is a framework that is inspired by the Massive Online Analysis (MOA), the most popular open source framework for machine learning on data streams, and MEKA, an open source implementation of methods for multi labeling. Scikit-multiflow is also inpired by scikit-learn, the most popular framework for machine learning in Python. Following the Scikit philosophy, scikit-multiflow is an open source machine-learning framework for multi-output/multi-label and stream data.

The Scikit-multiflow does not provide a graphical user interface like the MOA framework that it was inspired by. On the other hand can it be used within Jupyter notebooks, a popular interface in the data scientist community that is based on python. It is mainly a library that contains stream generators, learning methods, change detectors and evaluation methods for Machien learning algorithms on streaming data. Scilit-multiflow lays a special focus on its design, so that it is user friendly to new user and familiar for more experienced ones.[18]

It achieves this by having a simple workflow that can be easily describe in 4 steps:

1. Create a stream:
A stream can be generate by use of either one of the predefiend Stream generators or by importing a Stream from a file.
2. Instantiate a classifier:
In this step the user has to choose by which type of algorithm the classifikation process should be accomplished.
3. Setup the evaluator:
Next the user has to choose if the evaluation process should be done presequential or by Hold-out method. Prequential-evaluation or interleaved-testthen-train evaluation, is a popular performance evaluation method for the stream setting only, where tests are performed on new data before using it to train the model. Hold-out evaluation is a popular performance evaluation method for batch and stream settings, where tests are performed in a separate test set. [18]
4. Run the evaluation:
The evaluator will perform the following subtask in this final step, first, it will check if there are still samples in the stream and then it will pass the next sample

to the classifier, which will either test the sample with the `predict()` method or it will update the classifier with the `partial_fit()` method of the classifier.

To give better understanding how this exactly works, 2.1 will provide a small example where a synthetic waveform gets generated with the `WaveformGenerator`. This generator will create a synthetic stream that generates default samples with 21 numeric attributes and 3 target values, based on a random differentiation of some base waveforms. For the classifier is the Hoeffding Tree chosen with default parameters. After this we chose a sequential evaluator with the following parameters:

- `show_plot = true`, this will generate us a dynamic plot that will update as the classifier is trained.
- `pretrain_size = 200`, this parameter sets the number of samples that is used for the first training set.
- `max_sample = 20000`, this parameter set the maximum number of samples that should be used.

In the last step we then call evaluator method, that will then update the results and the plot. It is to note that at the time this thesis was written, the accuracy score stand for the performance in the framework. [18]

Listing 2.1: Simple example of the scikit-multiflow workflow

```

1  """1. Create a stream"""
2  stream = WaveformGenerator()
3  stream.prepare_for_use()
4
5
6  """2. Instantiate the HoeffdingTree classifier"""
7  ht = HoeffdingTree()
8
9  """3. Setup the evaluator"""
10 evaluator = EvaluatePrequential(show_plot=True,
11                                pretrain_size=200,
12                                max_samples=20000)
13
14 """4. Run evaluation"""
15 evaluator.evaluate(stream=stream, model=ht)

```

2 Basics

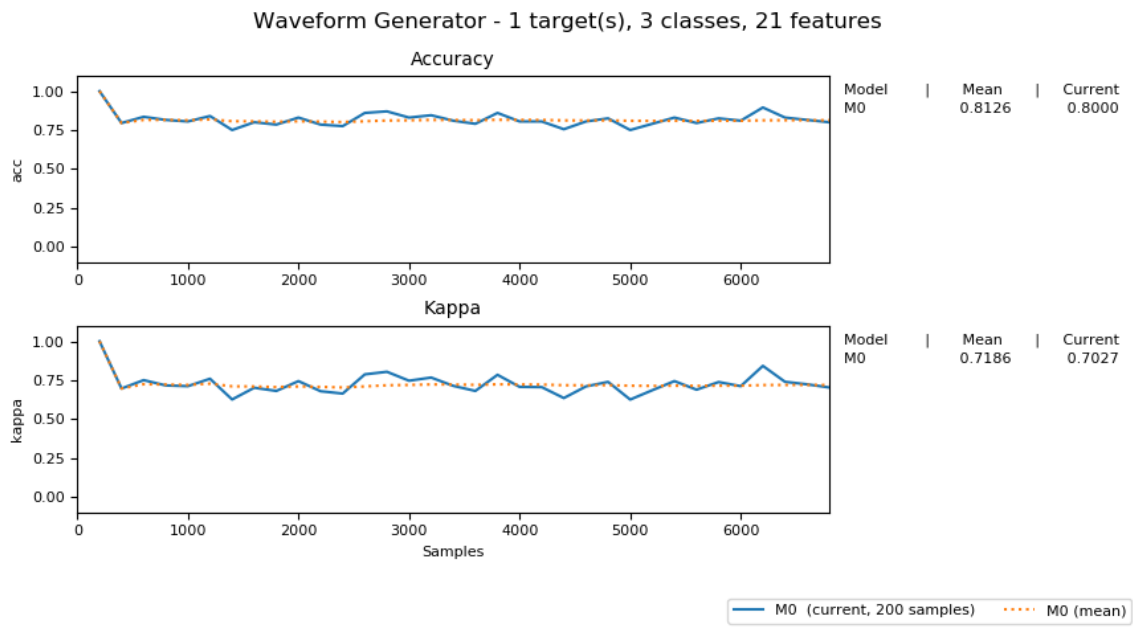


Figure 2.6: Result Table of the example in 2.1

3 Used Algorithms

3.1 Theoretical Work

This chapter will try to give a short explanation about how the Robust Soft Learning Vector Quantization (RSLVQ) works and will also explain the algorithm that it is based on. Furthermore will it explain what momentum based gradient descent algorithms there are and how they work. These are used in the different optimization implementation of the RSLVQ that will be compared in the later chapters.

3.1.1 Learning Vector Quantization

The first part of the theoretical work will give an introduction to the Learning Vector Quantization (LVQ), from that later the RSLVQ was derived from. The LVQ is a class of learning algorithms for nearest prototype classification (NPC). It stores a set of appropriately chosen prototype vectors instead of storing all the data points of a training set. This change makes the method computationally more efficient, as it does not have to store a high number of items to which new data must be compared to for classification anymore. Instead, the current data point only has to be compared to the prototypes to achieve classification. This led to a high popularity of the LVQ for real time applications, like i.e. speech recognition [16, 17], where it was used with good success.

A NPC consists of a set $\tau = (\theta_j, c_j)_{j=1}^M$ of labeled prototype vectors. The parameter $\theta_j \in \chi \equiv \mathbb{R}^D$ are vectors in the data space and $c_j \in \iota$, $i = 1, \dots, N_y$ are their corresponding labels. Typically, the number of prototypes is larger than the number of classes, such that every classification boundary is determined by more than one prototype. The class y of a new datapoint x is determined

- by selecting the prototype θ_i which is closest to x ,

$$y = \min_i d(x, \theta_i), \quad (3.1)$$

3 Used Algorithms

where $d(.,.)$ is the distance between data point and prototype, and

- by assigning the label y of this prototype to the datapoint x .

A popular choice for d is the Euclidian distance

$$d(q, p) = \sqrt{\sum_{i=1}^n (q_i - p_i)^2}. \quad (3.2)$$

LVQ is a class of algorithm for the determination of prototype vectors for NPC, which make use of the distribution of the training data as well as their class labels when selecting useful set of prototype vectors. The LVQ 2.1, for example, selects for every pair (x, y) , where x is the data point and y is the label, from the training set $S = (x_i, y_i)_{i=1}^N$, the two nearest prototypes θ_l, θ_m (see equation 3.1 for selection of one prototype) according to the Euclidian distance (see equation: 3.2). If the labels c_l and c_m are different and if one of them is equal to the label y of the data point, then the two nearest prototypes are adjusted according to:

$$\begin{aligned} \theta_l(t+1) &= \theta_l(t) + \alpha(t)(x - \theta_l), & c_l &= y, \\ \theta_m(t+1) &= \theta_m(t) - \alpha(t)(x - \theta_m), & c_m &\neq y. \end{aligned} \quad (3.3)$$

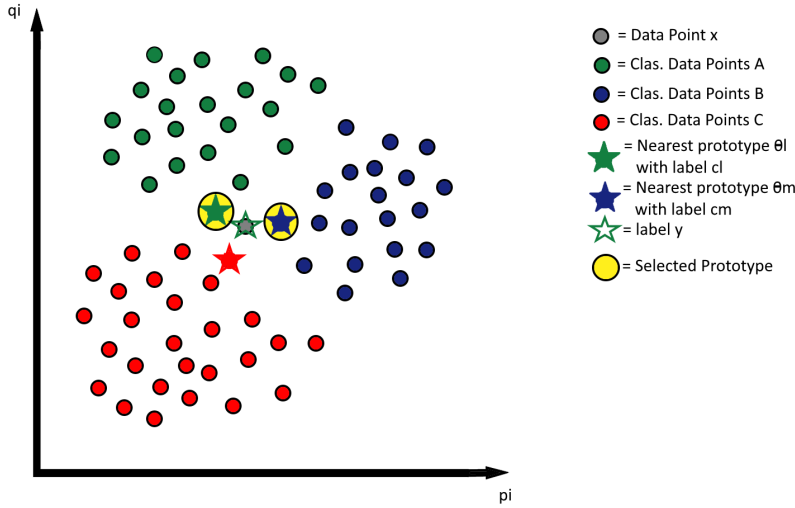


Figure 3.1: Graphic depiction of the LVQ classification task.

3 Used Algorithms

If the labels c_l and c_m are equal or both labels differ from the label y of the data point, no parameter update is being performed. The prototype, however, are changed only if the data point x is close to the classification boundry, i.e. if it falls into a *window*

$$\min\left(\frac{d(x, \theta_m)}{d(x, \theta_l)}, \frac{d(x, \theta_l)}{d(x, \theta_m)}\right) > s, \quad \text{where } s = \frac{1 - \omega}{1 + \omega}, \quad (3.4)$$

of relative width $0 < \omega \leq 1$. This "window rule" had to be introduced, as otherwise the prototype vectors would have diverged. [24]

3.1.2 Robust Soft Learning Vector Quantization

The LVQ algorithm may be a popular adaptive nearest prototype classifier for multiclass classification, but the algorithms from this family had only been proposed on heuristic grounds. Seo and Obermayer then derived the Robust Soft Learning Vector Quantization (RSLVQ variant) in their work [24]. It is a LVQ that is based on a Gaussian Mixture ansatz. It proposes an objective function which is based on a likelihood ratio and derives the learning rule from gradient descent.

They assume that the following objective function is maximized:

$$\begin{aligned} L_r &= \prod_{k=1}^N \frac{p(x_k, y_k | \tau)}{p(x_k, y_k | \tau) + p(x_k, \bar{y}_k | \tau)} \\ &= \prod_{k=1}^N \frac{p(x_k, y_k | \tau)}{p(x_k | \tau)} \stackrel{!}{=} \max. \end{aligned} \quad (3.5)$$

The ratio $\frac{p(x, y | \tau)}{p(x | \tau)}$ is bounded by 0 from below and bounded by 1 from above. Because of this, instead of optimizing the likelihood ratio, they then optimize the logarithm of the ratio L_r like this,

$$\log L_r = \sum_{k=1}^N \log \frac{p(x_k, y_k | \tau)}{p(x_k | \tau)} \stackrel{!}{=} \max, \quad (3.6)$$

3 Used Algorithms

were Stochastic gradient descent leads to the following learning rule:

$$\theta_l(t+1) = \theta_l(t) + \alpha(t) \frac{\partial}{\partial \theta_l} [\log \frac{p(x, y | \tau)}{p(x | \tau)}], \quad (3.7)$$

where $\alpha(t)$ is learning rate with $\sum_{t=1}^{\infty} \alpha(t) = \infty$ and $\sum_{t=1}^{\infty} \alpha^2(t) < \infty$. If the conditional density functions $p(x | j)$ are of the normalized exponential form

$$p(x | j) = K(j) \exp f(x, \theta_j), \quad (3.8)$$

the learning rule becomes

$$\begin{aligned} \frac{\partial}{\partial \theta_l} [\log \frac{p(x, y | \tau)}{p(x | \tau)}] &= \delta(c_l = y) (P_y(l | x) - P(l | x)) \frac{\partial f(x, \theta_l)}{\partial \theta_l} \\ &\quad - \delta(c_l \neq y) P(l | x) \frac{\partial f(x, \theta_l)}{\partial \theta_l}. \end{aligned} \quad (3.9)$$

where $P_y(l | x)$ and $P(l | x)$ are assignment probabilities,

$$\begin{aligned} P_y(l | x) &= \frac{p(l) \exp f(x, \theta_l)}{\sum_{j: c_j = y} p(j) \exp f(x, \theta_j)}, \\ P(l | x) &= \frac{p(l) \exp f(x, \theta_l)}{\sum_{j=1}^M p(j) \exp f(x, \theta_j)}. \end{aligned} \quad (3.10)$$

$P_y(l | x)$ describes the (posterior) probability that the data point x is assigned to the component l of the mixture, given that the data point was generated by the correct class. $P(l | x)$ describes the (posterior) probability that the data point x is assigned to the component l of the complete mixture using all classes. Using the gradient given by the equation 3.7, following learning rule is obtained:

$$\theta_l(t+1) = \theta_l(t) + \alpha(t) \begin{cases} P_y(l | x) - P(l | x) [\frac{\partial f(x, \theta_l)}{\partial \theta_l}], & c_l = y, \\ -P(l | x) [\frac{\partial f(x, \theta_l)}{\partial \theta_l}], & c_l \neq y. \end{cases} \quad (3.11)$$

3 Used Algorithms

It has to be noted that the factor $P_y(l | x) - P(l | x)$ is always positive. [24]

The Cost-function and the learning rule described until now were for training the RSLVQ for the general case, but Seo and Obermayer provided also a variant that uses a Gaussian mixture ansatz. In that version they choose a Gaussian mixture model whose components have similar width and strengths, i.e. $\sigma_l = \sigma, \forall l$, and $p(l) = \frac{1}{M}, l = 1, \dots, M$. They obtain

$$f(x, \theta_l) = \frac{-(x - \theta_l)^2}{2\sigma^2}, \quad f(x, \theta_l) = \frac{1}{\sigma^2}(x - \theta_l). \quad (3.12)$$

By substituting the derivative of $f(x, \theta_l)$ into equation 3.11 they obtained

$$\theta_l(t+1) = \theta_l(t) + \alpha(t) \begin{cases} P_y(l | x) - P(l | x)(x - \theta_l), & c_l = y, \\ -P(l | x)(x - \theta_l), & c_l \neq y, \end{cases} \quad (3.13)$$

where $\alpha(t) = \frac{\alpha(t)}{\sigma^2}$ and

$$\begin{aligned} P_y(l | x) &= \frac{\exp\frac{-(x-\theta_l)^2}{2\sigma^2}}{\sum_{j:c_j=y} \exp\frac{-(x-\theta_j)^2}{2\sigma^2}}, \\ P(l | x) &= \frac{\exp\frac{-(x-\theta_l)^2}{2\sigma^2}}{\sum_{j=1}^M \exp\frac{-(x-\theta_j)^2}{2\sigma^2}}. \end{aligned} \quad (3.14)$$

The equations 3.13 and 3.14 implement the RSLVQ algorithm. Prototypes whose label is equal to the label of the data point are attracted, while prototypes with different labels are repelled. Using deterministic annealing of σ^2 , which is a hyperparameter of the learning rule described in 3.13, the RSLVQ can be optimized. If the width σ of the Gaussian components goes to zero, both assignment probabilities (equation 3.14), become hard assignments. If the label c_q of the nearest prototype q is equal to the label y of the data point x then no update is being performed, because

$$P(l | x) = P_y(l | x) = 0, \quad \forall l \neq q, \quad P_y(q | x) - P(q | x) = 0. \quad (3.15)$$

Only if the label of the data point x differs from the label of the nearest prototype q , then

$$\begin{aligned} P(l | x) &= 0, \forall l \neq q, & P_y(l | x) &= 0, \forall l \neq q', \\ P(q | x) &= 1, & P_y(q | x) &= 1, \end{aligned} \tag{3.16}$$

where q' is the nearest prototype vector with the correct class label, and the nearest prototype vector and the nearest correct prototype vector are changed. In the "hard" version of the RSLVQ, prototypes are modified only by the data points that are not correctly classified. [24]

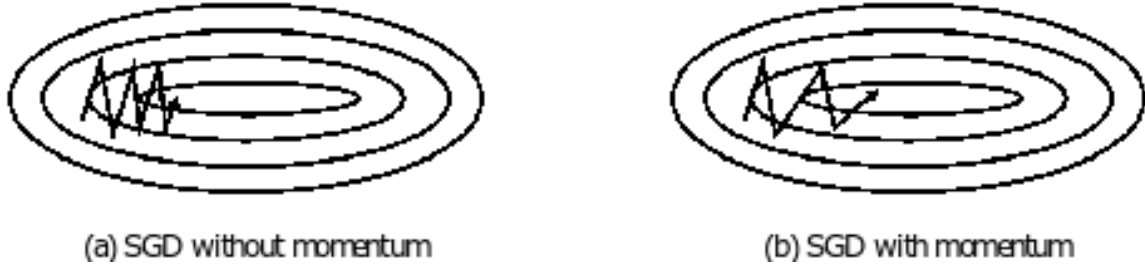
Seo and Obermayer then showed in their work [24] for LVQ 2.1 data points which are near the current class boundary and which are classified correctly have the strongest effect on the update of the prototypes. In contrast to that relies the RSLVQ on data points which are further away from the classification boundary and which are incorrectly classified, i.e. RSLVQ learns from mistakes. Additionally, the RSLVQ does not require a "window rule" like the LVQ 2.1, as the prototypes do not diverge. [24]

3.2 Momentum Based Gradient Descent

The following section will outline three algorithms that are widely used in the Deep learning Community and that will later be used to optimize the RSLVQ Algorithm in hopes to improve its performance and accuracy.

3.2.1 Momentum

SGD has trouble navigating ravines, i.e. areas where the surface curves much more steeply in one dimension than in another [26, 23], which are common around local optima. In these scenarios, SGD oscillates across the slopes of the ravine while only making hesitant progress along the bottom towards the local optimum as in figure.

Figure 3.2: Differences of a SGD with and without momentum.¹

Momentum [21, 23] is a method that helps accelerate SGD in the relevant direction and dampens oscillations as can be seen in figure . It does this by adding a fraction γ of the update vector of the past time step to the current update vector

$$\begin{aligned} v_t &= \gamma v_{t-1} + \eta \nabla_{\theta} J(\theta) \\ \theta &= \theta - v_t \end{aligned} \tag{3.17}$$

The momentum term γ is usually set to 0.9 or a similar value.

In his Article, Ruder [23] describes it like if we would push a ball down a hill. The ball accumulates momentum as long as it rolls downhill, becoming faster and faster on the way (until it reaches its terminal velocity, if there is air resistance, i.e. $\gamma < 1$). The same thing happens to our parameter updates: The momentum term increases for dimensions whose gradients point in the same directions and reduces updates for dimensions whose gradients change directions. As a result, we gain faster convergence and reduced oscillation.

3.2.2 Adagrad

The first algorithm that will be described is the Adagrad [27]. It has to be noted that this algorithm will not be part of the comparison, but as this algorithm is the base from which the Adadelta and later the RMSprop were developed from, it will be shortly explained nonetheless. The Adagrad is an algorithm for gradient-based optimization that adapts to updates for each individual parameter to perform larger or smaller updates depending on their importance. The learning rate gets adapted to the parameters, where it performs larger updates if the parameters are more infrequent and smaller updates if they are frequent. Hence, it is well-suited to deal with sparse data. As other works

¹Genevieve B. Orr

3 Used Algorithms

[8] have shown, this greatly improves the robustness of SGD, which lead to its wide use to train large Neural Networks. It was also used to train GloVe word embeddings as infrequent words require much larger updates than frequent ones [20, 23].

For the reason that Adagrad uses a different learning rate for every parameter θ_i at every time step t , Adagrad's per-parameter update is first shown and then vectorized. For brevity, we set $g_{t,i}$ to be the gradient of the objective function w.r.t. the parameter θ_i at the time step t :

$$g_{t,i} = \nabla_{\theta_i} J(\theta_{t,i}) \quad (3.18)$$

The update for every SGD parameter θ_i at time t becomes to:

$$\theta_{t+1,i} = \theta_{t,i} - \eta \cdot g_{t,i} \quad (3.19)$$

In its update rule, Adagrad modifies the general learning rate η each time step t for every parameter θ_i based on the past gradients that have been computed for θ_i :

$$\theta_{t+1,i} = \theta_{t,i} - \frac{\eta}{G_{t,ii} + \epsilon} g_{t,i} \quad (3.20)$$

$G_t \in \mathbb{R}^{d \times d}$ here is a diagonal matrix where each diagonal element i , i is the sum of the squares of the gradients w.r.t. θ_i up to time step t [9], while ϵ is a smoothing term that avoids division by zero (usually on the order of $1e-8$).

As G_t contains the sum of the squares of the past gradients w.r.t. to all parameters θ along its diagonal, we can now vectorize our implementation by performing an element-wise matrix-vector multiplication \odot between G_t and g_t :

$$\Delta\theta_t = -\frac{\eta}{\sqrt{G_t + \epsilon}} \odot g_t \quad (3.21)$$

One of Adagrad's main benefits is that it eliminates the need to manually tune the learning rate. Most implementations use a default value of 0.01 and leave it at that. Adagrad's main weakness is its accumulation of the squared gradients in the denominator: Since every added term is positive, the accumulated sum keeps growing during training. This in turn causes the learning rate to shrink and eventually become infinitesimally small, at which point the algorithm is no longer able to acquire additional knowledge. [23]

3.2.3 Adadelta

One of the more often used momentum-based algorithm is the Adadelta approach. It is an extension of the Adagrad that seeks to reduce its aggressive, monotonically decreasing learning rate. Instead of accumulating all past squared gradients, Adadelta restricts the window of accumulated past gradients to some fixed size ω .

Rather than inefficiently storing ω previous squared gradients, the sum of gradients is recursively defined as a decaying average of all past squared gradients. The running average $E[g^2]_t$ at time step t then depends (as a fraction γ similarly to the Momentum term) only on the previous average and the current gradient:

$$E[g^2]_t = \gamma E[g^2]_{t-1} + (1 - \gamma) g_t^2 \quad (3.22)$$

We then set γ to a similar value as the momentum, around 0.9. The SGD update will then be rewritten in terms of the parameter update vector $\Delta\theta_t$:

$$\begin{aligned} \Delta\theta_t &= -\eta \cdot g_{t,i} \\ \theta_{t+1} &= \theta_t - \Delta\theta_t \end{aligned} \quad (3.23)$$

Because of this, the parameter update vector, that is base of the parameter update vector of the Adagrad (see equation 3.21), takes the form of:

$$\Delta\theta_t = -\frac{\eta}{\sqrt{G_t + \epsilon}} \odot g_t \quad (3.24)$$

In the next step the diagonal Matrix G_t gets replaced with the decaying average over past squared gradients $E[g^2]_t$:

$$\Delta\theta_t = -\frac{\eta}{\sqrt{E[g^2]_t + \epsilon}} g_t \quad (3.25)$$

As the denominator is just the root mean squared (RMS) error criterion of the gradient, we can replace it with the criterion short-hand:

$$\Delta\theta_t = -\frac{\eta}{RMS[g]_t} g_t \quad (3.26)$$

3 Used Algorithms

As in [9] noted, the units in this update (as well as in the SGD, Momentum and Ada-grad) do not match, i.e. the update should have the same hypothetical units as the parameter. To realize this, they first define another decaying average, this time not of squared gradients but of squared parameter updates:

$$E[\Delta\theta^2]_t = \gamma E[\Delta\theta^2]_{t-1} + (1 - \gamma)\Delta\theta_t^2 \quad (3.27)$$

This leads to that the root mean squared error of the parameter updates is:

$$RMS[\Delta\theta]_t = \sqrt{E[\Delta\theta^2]_t + \epsilon} \quad (3.28)$$

Because the $RMS[\Delta\theta]_t$ value is unknown to us, we try to approximate it with the RMS of parameter updates until the previous time step. Replacing the learning rate η in the previous update rule with $RMS[\Delta\theta]_{t-1}$ leads to the final update rule for the Adadelta:

$$\begin{aligned} \Delta\theta_t &= -\frac{RMS[\Delta\theta]_{t-1}}{RMS[g]_t} g_t \\ \theta_{t+1} &= \theta_t + \Delta\theta_t \end{aligned} \quad (3.29)$$

In the Adadelta we do not need to set a learning rate η , as it was eliminated from the update rule. [23]

3.2.4 RMSprop

RMSprop is an adaptive learning rate method proposed by Geoff Hinton in Lecture 6e of his Coursera Class. [14] As a result of the need to resolve the Adagrad's radically diminishing learning rates, both the Adadelta and the RMSprop had been developed independently from each other at the same time. RMSprop in fact is identical to the first update vector of Adadelta that we derived above:

$$E[g^2]_t = 0.9E[g^2]_{t-1} + 0.1g_t^2 \quad (3.30)$$

$$\theta_{t+1} = \theta_t - \frac{\eta}{RMS[g]_t} g_t \quad (3.31)$$

RMSprop as well divides the learning rate by an exponentially decaying average of squared gradients. Hinton suggests γ to be set to 0.9, while a good default value for the learning rate η is 0.001. [23]

3.2.5 Adam

As described in [15], the Adaptive Moment Estimation (Adam) is another way to compute the adaptive learning rate for each parameter. Additionally to storing an exponentially decaying average of past squared gradients v_t like Adadelta or RMSprop, Adam also keeps an exponentially decaying average of past gradients m_t , similar to momentum:

$$\begin{aligned} m_t &= \beta_1 m_{t-1} + (1 - \beta_1) g_t \\ v_t &= \beta_2 v_{t-1} + (1 - \beta_2) g_t^2 \end{aligned} \tag{3.32}$$

m_t and v_t are estimates of the first moment (the mean) and the second moment (the uncentered variance) of the gradients respectively, hence the name of the method. As m_t and v_t are initialized as vectors of 0's, the authors of Adam observe that they are biased towards zero, especially during the initial time steps, and especially when the decay rates are small (i.e. β_1 and β_2 are close to 1).

To counteract these biases they compute a bias-corrected estimate of the first and second moment:

$$\begin{aligned} \hat{m}_t &= \frac{m_t}{1 - \beta_1^t} \\ \hat{v}_t &= \frac{v_t}{1 - \beta_2^t} \end{aligned} \tag{3.33}$$

These are then used to update the parameters just as we have seen in Adadelta and RMSprop, which yields the Adam update rule:

$$\theta_{t+1} = \theta_t - \frac{\eta}{\sqrt{\hat{v}_t} + \epsilon} \hat{m}_t \tag{3.34}$$

The authors recommend in their work to use as default values of 0.9 for β_1 , 0.999 for β_2

and 10^{-8} for ϵ . They show empirically that Adam works well in practice and compares favorably to other adaptive learning-method algorithms. [23, 15]

Algorithm 1 Adam algorithm, as proposed in [15]. The authors recommend for default settings the value 0.9 for β_1 , 0.999 for β_2 and 10^{-8} for ϵ as well as a stepsize α of 0.001.

Require: α : Stepsize

Require: $\beta_1, \beta_2 \in [0, 1)$: Exponential decay rates for the moment estimates

Require: $f(\theta)$: Stochastic objective function with parameters θ

Require: θ_0 : Initial parameter vector

```

1: Initialize first moment vector:  $m_0 = 0$ 
2: Initialize second moment vector:  $v_0 = 0$ 
3: Initialize timestep:  $t = 0$ 
4: while  $\theta_t$  not converged do
5:  $t = t + 1$ 
6: Get the gradients at timestamp  $t$ :  $g_t \leftarrow \nabla_{\theta} J(\theta_{t-1})$  ( see equation 3.18)
7: Update first moment estimate:  $m_t \leftarrow \beta_1 m_{t-1} + (1 - \beta_1) g_t$ 
8: Update second moment estimate:  $v_t \leftarrow \beta_2 v_{t-1} + (1 - \beta_2) g_t^2$ 
9: Compute bias-corrected first moment estimate:  $\hat{m}_t \leftarrow \frac{m_t}{1 - \beta_1^t}$ 
10: Compute bias-corrected second moment estimate:  $\hat{v}_t \leftarrow \frac{v_t}{1 - \beta_2^t}$ 
11: Update the parameters:  $\theta_{t+1} \leftarrow \theta_t - \frac{\eta}{\sqrt{\hat{v}_t + \epsilon}} \hat{m}_t$ 
12: return Resulting parameters:  $\theta_t$ 

```

3.3 Implementation of the RSLVQ variants

This chapter will show how the different RSLVQ variants were implemented, so that they can work with Streaming data. The implementation is based on the version that was used in the scikit-multiflow framework [18] and was then further modified by Moritz Heusinger in [13]. It has to be said that the implementation is nearly the same as in the previously mentioned work [13].

3.3.1 RSLVQ SGD

The main differences between the different implementations of the algorithms appear in the training phase, when the prototypes get updated in the `_update_prototype` method. As the rest stays nearly the same, it will only be explained once during the explanation of how the *RSLVQ_{SGD}* variant was implemented. In the scikit-multiflow framework the `partial_fit` method, shown in listing 3.1, gets executed everytime when an algorithm initiates the learning process or when it continues training process. It will

check if the algorithm has already completed the initial learning process or if it is the first time that the method gets executed. If one of those 2 cases is true, the method will continue by calling the internal `_validate_train_parms` method, that will then carry out further checks. If none of the 2 options is true, the method will raise an error that the algorithm has to initial the learning/training process. After this check, the method will call the `_optimize` method, that is shown in listing 3.2. It has to be noted that neither the `partial_fit`, the `_validate_train_parms` and the `_optimize` methods have been altered from how they have been implemented in the scikit-multiflow framework [18].

Listing 3.1: Method `partial_fit` from the scikit-multiflow framework

```

1  def partial_fit(self, X, y, classes=None, sample_weight=None):
2      """Fit the LVQ model to the given training data and parameters
3      using
4          gradient ascent.
5
6          Parameters
7          -----
8          X : array-like, shape = [n_samples, n_features]
9              Training vector, where n_samples in the number of samples
10         and
11             n_features is the number of features.
12         y : numpy.ndarray of shape (n_samples, n_targets)
13             An array-like with the class labels of all samples in X
14         classes : numpy.ndarray, optional (default=None)
15             Contains all possible/known class labels. Usage varies
16         depending
17             on the learning method.
18         sample_weight : Not used.
19
20         Returns
21         -----
22         self
23         """
24         if set(unique_labels(y)).issubset(set(self.classes_)) or
25            self.initial_fit is True:
26             X, y = self._validate_train_parms(X, y, classes=classes)
27         else:
28             raise ValueError('Class {} was not learned - please declare
29                               all classes in first call of fit/partial_fit'.format(y))
30
31         self._optimize(X, y)
32         return self

```

The `_optimize` method will search for the nearest correct and the nearest wrong prototype, based on their euclidean distance towards the data point that has to be classified, as describe in 3.3 and visualized in 3.1. If the the correct prototype(prototype with the right label) is not the nearest to the data point, then each the model of the nearest wrong prototype as well as the model of the nearest correct prototype should be updated by calling the `_update_prototype` method

Listing 3.2: Method `_optimize` from the scikit-multiflow framework

```

1  def _optimize(self, X, y):
2      nb_prototypes = self.c_w_.size
3
4
5      n_data, n_dim = X.shape
6      prototypes = self.w_.reshape(nb_prototypes, n_dim)
7
8      for i in range(n_data):
9          xi = X[i]
10         c_xi = int(y[i])
11         best_euclid_corr = np.inf
12         best_euclid_incorr = np.inf
13
14         # find nearest correct and nearest wrong prototype
15         for j in range(prototypes.shape[0]):
16             if self.c_w_[j] == c_xi:
17                 eucl_dis = euclidean_distances(xi.reshape(1,
18                 xi.size),
19
20                 prototypes[j]
21                 .reshape(1,
22                 prototypes[j]
23                 .size))
24                 if eucl_dis < best_euclid_corr:
25                     best_euclid_corr = eucl_dis
26                     corr_index = j
27             else:
28                 eucl_dis = euclidean_distances(xi.reshape(1,
29                 xi.size),
30
31                 prototypes[j]
32                 .reshape(1,
33                 prototypes[j]
34                 .size))
35                 if eucl_dis < best_euclid_incorr:
36                     best_euclid_incorr = eucl_dis
37                     incorr_index = j

```

```

33         # Update nearest wrong prototype and nearest correct
        prototype
34         # if correct prototype isn't the nearest
35         if best_euclid_incorr < best_euclid_corr:
36             self._update_prototype(j=corr_index, c_xi=c_xi, xi=xi,
37                                   prototypes=prototypes)
38             self._update_prototype(j=incorr_index, c_xi=c_xi, xi=xi,
39                                   prototypes=prototypes)

```

The `_update_prototype`, as shown as in listing 3.3, is when the first real difference to the scikit implementation appears. The method will then compare if the prototype has the same label y_i (here named as c_{xi}) as the data point x_i . It will be then adjusted according to equation 3.13. $P_y(l | x)$ and $P(l | x)$ from 3.14 are then calculated by the `_p` method. It uses the internal method `_costf` to calculate the cost for each of the prototypes and it represents the $\frac{-(x-\theta_l)^2}{2\sigma^2}$ part of equation 3.14. Both the `_p` and the `_costf` method are unchanged from their implementation in [18]

Listing 3.3: Method `_update_prototype` as it was used in [13].

```

1
2 def _update_prototype(self, j, xi, c_xi, prototypes):
3     """SGD"""
4     d = xi - prototypes[j]
5
6     if self.c_w_[j] == c_xi:
7         # Attract prototype to data point
8         self.w_[j] += self.learning_rate * (self._p(j, xi,
9         prototypes=self.w_, y=c_xi) - self._p(j, xi, prototypes=self.w_)) *
10        d
11    else:
12        # Distance prototype from data point
13        self.w_[j] -= self.learning_rate * self._p(j, xi,
14        prototypes=self.w_) * d

```

3.3.2 RSLVQ Adadelta

The first steps of the Adadelta implementation are the same as in the SGD implementation. In the first step the method will again calculate the adjustment for the prototype like it was described in 3.13. It will also calculate the Posterior and prior in the same way as described in 3.14 with the `_p` and `_costf` methods. In next step is now where the differences between the implementation of the SGD and Adadelta start. First, it will calculate the sum of the past squared gradients according to equation 3.22. It will then

3 Used Algorithms

calculate the gradient update at step t according to the first equation from 3.29. As the $RMS[\Delta\theta]_t$ value is unknown to us, we try to approximate it with the RMS of parameter updates until the previous time step like this: $E[\Delta\theta^2] + \epsilon$. This gives us:

$$\Delta\theta_t = \frac{(E[\Delta\theta^2] + \epsilon)}{\sqrt{(E[\Delta\theta^2] + \epsilon)}} g_t \quad (3.35)$$

Next we store the squared parameter updates according to equation 3.27. Finally we apply the update prototype according to the second equation from 3.29:

$$\theta_{t+1} = \theta_t + \Delta\theta_t. \quad (3.36)$$

This implementation of the Adadelata was also used in [13].

Listing 3.4: Adadelata implementation from [13].

```

1  """Implementation of Adadelata"""
2      d = (xi - prototypes[j])
3
4
5      if self.c_w_[j] == c_xi:
6          gradient = (self._p(j, xi, prototypes=self.w_, y=c_xi) -
7 self._p(j, xi, prototypes=self.w_)) * d
8      else:
9          gradient = - self._p(j, xi, prototypes=self.w_) * d
10
11      # Accumulate past squared gradients
12      self.squared_mean_gradient[j] = self.decay_rate
13 *self.squared_mean_gradient[j] + (1 - self.decay_rate) * gradient **
14 2
15
16      # Compute update/step
17      step = ((self.squared_mean_step[j] + self.epsilon) /
18 (self.squared_mean_gradient[j] + self.epsilon)) **0.5 * gradient
19
20      # Accumulate updates
21      self.squared_mean_step[j] = self.decay_rate *
22 self.squared_mean_step[j] + (1 - self.decay_rate) * step ** 2
23
24      # Attract/Distract prototype to/from data point
25      self.w_[j] += step

```

3.3.3 RSLVQ RMSprop

The implementation of the RMSprop is nearly the same as for the Adadelta, the difference is that we reintroduce the learning rate η , that was eliminated in the Adadelta and replaced by the root mean squared error of the parameter updates $RMS[\Delta\theta]_{t-1}$.

The steps until after the posterior and prior calculation are the exact same as in the SGD implementation. When the gradient is calculated, it will be added to the average past squared gradients like in equation 3.30 described. If we take the value 0.9 for γ like Hinton suggested in [14], we get the following equation:

$$G_t = 0.9E[g^2]_{t-1} + 0.1g_t^2 \quad (3.37)$$

The prototype update can be calculated by equation 3.31, were the $-\frac{\eta}{RMS[g_t]}g_t$ part, that calculates the gradient update, can be substituted with $\Delta\theta$. Because the RMSprop is identical to the first update vector of the Adadelta, we will then get for $\Delta\theta$ the following equation:

$$\Delta\theta = -\frac{\eta}{G_t - \epsilon}g_t. \quad (3.38)$$

This equation is then used to calculate the gradient update before it is update is performed on the prototype in the final implementation [13].

Listing 3.5: Implementation of the RMSprop from [13].

```

1  """RMSprop"""
2
3      d = xi - prototypes[j]
4
5      if self.c_w_[j] == c_xi:
6          gradient = (self._p(j, xi, prototypes=self.w_, y=c_xi) -
7 self._p(j, xi, prototypes=self.w_)) * d
8      else:
9          gradient = - self._p(j, xi, prototypes=self.w_) * d
10
11      # Accumulate gradient
12      self.squared_mean_gradient[j] = 0.9 *
13 self.squared_mean_gradient[j] + 0.1 * gradient ** 2
14
15      # Update Prototype

```

```

14     self.w_[j] += (self.learning_rate /
    ((self.squared_mean_gradient[j] + self.epsilon) ** 0.5)) * gradient

```

3.3.4 RSLVQ Adam

Listing 3.6: Implementation of Adam on basis of [13]

```

1  """Adam"""
2      d = (xi - prototypes[j])
3
4      """Calculate posterior"""
5      if self.c_w_[j] == c_xi:
6          gradient = (self._p(j, xi, prototypes=self.w_, y=c_xi) -
7 self._p(j, xi, prototypes=self.w_)) * d
8      else:
9          gradient = - self._p(j, xi, prototypes=self.w_) * d
10
11      """Compute gradients m """
12      self.gradients_m[j] = self.beta_1 * self.gradients_m[j] + (1 -
13 self.beta_1) * gradient
14
15      """Compute squared gradients v """
16      self.gradients_v_sqrt[j] = self.beta_2 *
17 self.gradients_v_sqrt[j] + (1 - self.beta_2) * gradient ** 2
18
19      """Compute m correction"""
20      m_corrected = self.gradients_m[j] / (1 - self.beta_1)
21
22      """Compute v correction"""
23      v_corrected = self.gradients_v_sqrt[j] / (1 - self.beta_2)
24
25      """Update prototype"""
26      self.w_[j] += - (self.learning_rate / (np.sqrt(v_corrected) +
27 self.epsilon)) * m_corrected

```


4 Test realisation

4.1 Used Streaming Datat sets

In this section will be a short description of the data streams, synthetic and real world data, that where used to perform the experiments.

4.1.1 Synthetic data

- Sea Concepts: !not finsihed rework required!
This dataset consists of 50000 instances with three attributes of which only two are relevant. The two class decision boundary is given by $f1 + f2 = b$, where $f1$, $f2$ are the two relevant features and b a predefined threshold. Abrupt drift is simulated with four different concepts, by changing the value of b every 12500 samples. Also included are 10% of noise. [25]
- Moving Squares:
placeholder

4.1.2 Real world data

- Weather Dataset:
This dataset contains eight different features like the minimum and maximum tempreature, visibility, wind speed, air pressure, etc. that were measured at the Offutt Air Force Base in Bellevue, Nebraska. These measuredments were measured in the period of 1949-1999. The goal of this dataset is to predict if it is going to rain on certains day or not. It contains 18159 instances with an imbalance towards no rain (69%) and was introduced by Elwell et al..
- Electicity Market Dataset:
First described by Harris et al., it is often used for performance comparison, as

it is a good benchmark for concept drift classification. It holds the information of the Australian New South Wales Electricity Market, whose prices are affected by supply and demand. Each sample is characterized by the following attributes: day of the week, time stamp, market demand, etc.. Each of the samples refers to a period of 30 minutes and the class label identifies the relative change compared to the last 24 hours. The used Dataset is the normalized version of the original dataset.

- **Poker Dataset:**

To create this dataset, one million poker hands, each represented by five cards that got encoded with its suit and rank, were randomly drawn. The class is the resulting poker hand itself such as one pair, full house and so forth. Since the poker hand can not change by definition and each instance was randomly generated, this dataset had in its original form no drift in it. As this is the version as it was introduced by [3], virtual drift is present in this dataset. This was achieved by sorting the instances of the dataset by rank and suit. Duplicate hands were also removed. Additionally, this version was also normalized.

- **Outdoor Objects:**

This Dataset was obtained by images that were recorded by a mobile in a garden environment. The goal is to classify 40 different objects, each approached ten times under varying lighting conditions affecting the color based representation. Every ten images represent one approach and each approach is in temporal order within the dataset. The objects are encoded in a normalized 21-dimensional RG-Chromaticity histogram.

4.2 Implementation of the Comparison and Results

4.2.1 Experiment begrenzung

4.2.2 Experiment aufbau

4.2.3 Experiment durchlauf

Listing 4.1: Beispiel für einen Quelltext

```
1
2 public void foo() {
3     // Kommentar
4 }
```

5 Result Evaluation

6 Summary

Abbreviations

LVQ Learning Vector Quantization

MOA Massive Online Analysis

NPC nearest prototype classification

RSLVQ Robust Soft Learning Vector Quantization

SGD Stochastic Gradient Descent

List of Figures

2.1	Overview of the Machine Learning Strategy's.	4
2.2	Process of the Supervised Learning. [22]	6
2.3	Visualisation differences of the three Gradient Descent Algorithms	8
2.4	Visualization of types of Concept drift [11]	11
2.5	Visualization of real and virtual drift [11]	11
2.6	Result Table of the example in 2.1	17
3.1	Graphic depiction of the LVQ classification task.	19
3.2	Diff. SGD moment.	24

List of Tables

Bibliography

- [1] Charu C Aggarwal. *Data streams: models and algorithms*. Vol. 31. Springer Science & Business Media, 2007.
- [2] Albert Bifet, Richard Kirkby, and Bernhardt Pfahringer. *DATA STREAM MINING A Practical Approach*. MIT Press, 2011.
- [3] Albert Bifet et al. “Efficient data stream classification via probabilistic adaptive windows”. In: *Proceedings of the 28th annual ACM symposium on applied computing*. 2013, pp. 801–806.
- [4] Albert Bifet et al. “Efficient online evaluation of big data stream classifiers”. In: *Proceedings of the 21th ACM SIGKDD international conference on knowledge discovery and data mining*. 2015, pp. 59–68.
- [5] Albert Bifet et al. *Machine Learning for Data Streams with Practical Examples in MOA*. <https://moa.cms.waikato.ac.nz/book/>. MIT Press, 2018.
- [6] Albert Bifet et al. “Pitfalls in benchmarking data stream classification and how to avoid them”. In: *Joint European Conference on Machine Learning and Knowledge Discovery in Databases*. Springer. 2013, pp. 465–479.
- [7] Jacob Cohen. “A coefficient of agreement for nominal scales”. In: *Educational and psychological measurement* 20.1 (1960), pp. 37–46.
- [8] Jeffrey Dean et al. “Large scale distributed deep networks”. In: *Advances in neural information processing systems*. 2012, pp. 1223–1231.
- [9] John Duchi, Elad Hazan, and Yoram Singer. “Adaptive Subgradient Methods for Online Learning and Stochastic Optimization”. In: *J. Mach. Learn. Res.* 12.null (July 2011), pp. 2121–2159. ISSN: 1532-4435.
- [10] Wolfgang Ertel. *Introduction to Artificial Intelligence*. Springer, 2003. ISBN: 9780857292995.
- [11] João Gama et al. “A Survey on Concept Drift Adaptation”. In: *ACM Comput. Surv.* 46.4 (Mar. 2014). ISSN: 0360-0300. DOI: 10.1145/2523813. URL: <https://doi.org/10.1145/2523813>.
- [12] Ian Goodfellow, Yoshua Bengio, and Aaron Courville. *Deep Learning*. <http://www.deeplearningbook.org>. MIT Press, 2016.

Bibliography

- [13] Moritz Heusinger, Christoph Raab, and Frank-Michael Schleif. “Passive Concept Drift Handling via Momentum Based Robust Soft Learning Vector Quantization”. In: *Advances in Self-Organizing Maps, Learning Vector Quantization, Clustering and Data Visualization*. Ed. by Alfredo Vellido et al. Cham: Springer International Publishing, 2020, pp. 200–209. ISBN: 978-3-030-19642-4.
- [14] Geoff Hinton. “Overview of mini-batch gradient descent”. In: (). URL: http://www.cs.toronto.edu/~tijmen/csc321/slides/lecture_slides_lec6.pdf.
- [15] Diederik P. Kingma and Jimmy Ba. “Adam: A Method for Stochastic Optimization”. In: *CoRR* abs/1412.6980 (2014).
- [16] Takashi Komori and Shigeru Katagiri. “Application of a generalized probabilistic descent method to dynamic time warping-based speech recognition”. In: *[Proceedings] ICASSP-92: 1992 IEEE International Conference on Acoustics, Speech, and Signal Processing*. Vol. 1. IEEE. 1992, pp. 497–500.
- [17] Erik McDermott and Shigeru Katagiri. “Prototype-based minimum classification error/generalized probabilistic descent training for various speech units”. In: *Computer Speech & Language* 8.4 (1994), pp. 351–368.
- [18] Jacob Montiel et al. “Scikit-Multiflow: A Multi-output Streaming Framework”. In: *Journal of Machine Learning Research* 19.72 (2018), pp. 1–5. URL: <http://jmlr.org/papers/v19/18-251.html>.
- [19] Nils J. Nilsson. *Introduction to Machine Learning*. 1998.
- [20] Jeffrey Pennington, Richard Socher, and Christopher D Manning. “Glove: Global vectors for word representation”. In: *Proceedings of the 2014 conference on empirical methods in natural language processing (EMNLP)*. 2014, pp. 1532–1543.
- [21] Ning Qian. “On the momentum term in gradient descent learning algorithms”. In: *Neural Networks* 12.1 (1999), pp. 145–151. ISSN: 0893-6080. DOI: [https://doi.org/10.1016/S0893-6080\(98\)00116-6](https://doi.org/10.1016/S0893-6080(98)00116-6). URL: <http://www.sciencedirect.com/science/article/pii/S0893608098001166>.
- [22] Sebastian Raschka. *Python MACHine Learning*. Packt Publishing, 2015. ISBN: 1783555130.
- [23] Sebastian Ruder. *An overview of gradient descent optimization algorithms*. 2016.
- [24] Sambu Seo and Klaus Obermayer. “Soft learning vector quantization”. In: *Neural computation* 15.7 (2003), pp. 1589–1604.
- [25] W. Nick Street and YongSeog Kim. “A Streaming Ensemble Algorithm (SEA) for Large-Scale Classification”. In: *Proceedings of the Seventh ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*. KDD ’01. San Francisco, California: Association for Computing Machinery, 2001, pp. 377–382. ISBN: 158113391X. DOI: 10.1145/502512.502568. URL: <https://doi.org/10.1145/502512.502568>.
- [26] Richard S. Sutton. *Two problems with backpropagation and other steepest-descent learning procedures for networks*. 1986.

Bibliography

- [27] Matthew D. Zeiler. “ADADELTA: An Adaptive Learning Rate Method”. In: *ArXiv* abs/1212.5701 (2012).
- [28] Indrė Žliobaitė et al. “Evaluation methods and decision theory for classification of streaming data with temporal dependence”. In: *Machine Learning* 98.3 (2015), pp. 455–482.

Statutory Declaration

Hiermit versichere ich, dass ich die vorgelegte Bachelorarbeit selbstständig verfasst und noch nicht anderweitig zu Prüfungszwecken vorgelegt habe. Alle benutzten Quellen und Hilfsmittel sind angegeben, wörtliche und sinngemäße Zitate wurden als solche gekennzeichnet.

Florian Hohn, am February 11, 2020