

Department of Computer Science
Department of Cognitive Science

Rasmus Diederichsen

LIVE INTROSPECTION FOR NEURAL NETWORK TRAINING

Master's Thesis

First Supervisor: Prof. Dr. Oliver Vornberger
Second Supervisor: Dr. Ulf Krumnack
Project Supervisor: Anders Arpteg, PhD

ABSTRACT

Artificial neural networks have become the prevalent model class for many machine learning tasks, including image classification, segmentation, video and audio analysis, or time series prediction. With ever increasing computational resources and advances in programming infrastructure, the size of model we can train also increases. Nevertheless, it is not uncommon for training to take days or weeks, even on potent hardware. While there are many obvious causes – e.g. inherent difficulty to parallelize training with the most successful algorithms – there may still be inefficiencies in the framework we typically use for training neural networks. While the success of deep learning models has been rapid, the theoretical justification for most game-changing ideas as well as the general principles of deep learning, has not kept pace (for more information, see Sanjeev Arora’s talk at ICML 2018¹). This means that for many successful techniques, we have no clear understanding why they work so well in practice.

This explanatory gap is also a reason why developing deep learning applications is considered more of an art than a science. In contrast to traditional programming, which builds on decades of research and development in electrical engineering, logic, mathematics and theoretical computer science, there’s rarely one definitive way to solve a certain problem in deep learning. Additionally, the debugging tools available to every programmer on every level of abstraction far exceed what we currently have for differentiable programming. Simple questions like “Does my model learn what I want it to learn?” are not answerable at this point. We can thus identify a need to supply more useful tooling for deep learning practitioners. A standard approach to choosing a parametrization remains trial-and-error, or only slightly more sophisticated ways to run and test. The computational cost of training large models prohibits quick experimentation and often translates into monetary costs as well. Identifying dead ends early or points in training when to tweak certain parameters could thus provide large savings in time and money, besides enabling a more thorough understanding of what is going on.

This thesis addresses the above in two ways: A software library aiding in deep learning debugging and experimentation is designed and implemented. The same library is then used to investigate experimentally, whether heretofore unknown signals can be extracted from the training process in order to validate parameter choices while training is running.

¹ Toward Theoretical Understanding of Deep Learning

ZUSAMMENFASSUNG

Künstliche neuronale Netze sind inzwischen die populärste Art von Modell für viele Anwendungen im maschinellen Lernen, z.B. Bildklassifikation, Segmentierung oder Video- und Audioanalyse sowie Zeitreihenvorhersage. Mit immer zunehmenden Ressourcen und Fortschritten in verfügbarer Software steigt auch die Größe der Trainierbaren Modelle. Dennoch sind Trainingszeiten von Tagen und Wochen – selbst auf potenter Hardware – nicht ungewöhnlich.

Es gibt naheliegende Ursachen – z.B. die Schwierigkeit, das Training zu parallelisieren aufgrund der eingesetzten Optimierungsalgorithmen – jedoch möglicherweise auch weniger offensichtliche Verluste in der Art und Weise, wie Modelle normalerweise trainiert werden.

Während der praktische Erfolg tiefer neuronaler Netze rapide war, konnte die theoretische Erschließung nicht Schritt halten, und viele fortschrittliche Ideen sind mehr von Intuition und Experimenten gestützt als von rigorosen Formalismen. Für die meisten erfolgreichen Ideen haben wir folglich kein klares Verständnis dafür, warum sie in der Praxis gut funktionieren.

Diese Kluft zwischen Praxis und Theorie ist auch ursächlich dafür, dass Deep-Learning-Applikationen mehr als Kunst denn als Wissenschaft gelten. Anders als in der traditionellen Softwareentwicklung, die auf Jahrzehnten der Forschung in Ingenieurwissenschaften, Mathematik und theoretischer Informatik basiert, gibt es im maschinellen Lernen oft keinen definitiven Ansatz, um ein konkretes Problem zu lösen. Des Weiteren sind die verfügbaren Debugging-Werkzeuge in jeder Schicht der Softwareentwicklung denen des Deep Learning bei weitem überlegen. Einfache Fragen wie “Lernt das Modell das, was ich möchte?” sind bislang nicht zu beantworten. Wir können hier also eine Notwendigkeit für mehr nützliche Tools für Anwender konstatieren.

Der Standardweg, eine gute Parameterwahl für das Modell zu treffen, ist Trial-and-Error oder etwas intelligentere Variationen davon. Der Rechenaufwand für das Training großer Modelle verhindert schnelles Experimentieren und schlägt sich häufig auch in finanziellen Kosten für die nötige Rechenleistung nieder. Frühzeitig Sackgassen zu identifizieren oder Zeitpunkte, an denen man bestimmte Parameter ändern sollte, könnten daher große Ersparnisse bedeuten.

Diese Arbeit stellt zum Einen eine Softwarebibliothek vor, die im Debugging von Deep-Learning-Anwendungen hilft und das Experimentieren erleichtert. Zum Anderen wird dieses Werkzeug benutzt, um experimentell festzustellen, ob aus dem Trainingsprozess Signale extrahiert werden können, die Aufschluss über die Wahl von Parametern geben, während das Training läuft.

CONTENTS

I	INTRODUCTION	2
1.1	Artificial Neural Networks	2
1.2	Goals of this thesis	2
1.3	Motivation	3
1.4	Existing Applications	4
2	IKKUNA	6
2.1	Design Principles	6
2.2	Deep Learning frameworks	6
2.3	Publisher-Subscriber	8
2.4	Overview of the library	9
2.4.1	The <code>export</code> subpackage	10
2.4.2	The <code>models</code> subpackage	16
2.4.3	The <code>utils</code> subpackage	17
2.4.4	The <code>visualization</code> subpackage	17
2.4.5	Miscellaneous tools	18
2.4.6	Plugin Infrastructure	21
2.4.7	Documentation	21
3	EXPERIMENTS IN LIVE INTROSPECTION	22
3.1	Detecting Learning Rate Problems	22
3.1.1	Experiment 1	23
A	APPENDIX A — OPEN SOURCE ACKNOWLEDGMENTS	27
B	APPENDIX B — CONTRIBUTIONS TO OTHER LIBRARIES	28
	BIBLIOGRAPHY	30

LIST OF TABLES

Table 2.1	<code>ikkuna.export</code> functionalities	10
Table 2.2	Subscribable message kinds	14
Table 2.3	Pre-packaged subscriber subclasses	16
Table 2.4	Named arguments to <code>main.py</code>	20
Table B.1	Contributions to third-party projects	28

LIST OF FIGURES

Figure 1.1	Schema of a multi-layer neural network	3
Figure 2.1	Changes in popularity of different deep learning libraries in research	8
Figure 2.2	One possible implementation of the Publisher-Subscriber pattern.	9
Figure 2.3	<code>ikkuna</code> package diagram	10
Figure 2.4	<code>ikkuna.export</code> package diagram	11
Figure 2.5	Classes in the <code>ikkuna.export.messages</code> submodule	13
Figure 2.6	<code>ikkuna.export.Exporter</code> class diagram	13
Figure 2.7	Classes defined in <code>ikkuna.export.subscriber</code>	15
Figure 2.8	<code>ikkuna.models</code> package diagram	15
Figure 2.9	<code>ikkuna.utils</code> package diagram	17
Figure 2.10	Class diagram for classes in <code>ikkuna.visualization</code>	18
Figure 2.11	Exemplary view of a matplotlib figure forwarded over SSH	19
Figure 2.12	Exemplary view of a TensorBoard session	19
Figure 3.1	Simplified AlexNet architecture	23
Figure 3.2	Final accuracies after 100 epochs	24

INTRODUCTION

This thesis is concerned with *deep* neural networks, meaning architectures consisting of many layers. Typically, the number of units in such a model and thus the number of tunable parameters ($n_{\text{weights_per_neuron}} \times n_{\text{neurons}}$) exceeds the number of training examples. Training such large networks thus involves iterating over the dataset many times which incurs a high computational cost due to the massive number of matrix operations involved. Speeding up the training is thus one of the primary endeavours in deep learning research. An overview of the workings of a neural network is given in section [1.1](#).

1.1 ARTIFICIAL NEURAL NETWORKS

The artificial neural network is a class of machine learning model which can be used for regression and classification. At its core, neuronal models are simply conceptualized as an arrangement of units which receive inputs, compute a weighted sum and thus produce an output activation. The ideas date back at least to Hebbian learning (a single neuron) in the 1940s, and were developed into the Perceptron model by Rosenblatt ([Rosenblatt, 1958](#)). With the introduction of the backpropagation algorithm ([Werbos, 1975](#)) and increasing availability of computing power, this began to change, but more easily trainable models like Support Vector machines eclipsed neural nets for most applications. Architectural advances, such as the convolutional neural networks in the late 1980's ([LeCun et al., 1989](#)) and the advent of GPU-accelerated neural network implementations (pioneered in [Ciresan et al. \(2011\)](#)) – as well as subsequently, gpu-accelerated linear algebra libraries with automatic differentiation – finally made neural networks the cornerstone for many AI applications today.

A neural network consists of at least one input layer, $0 - n$ intermediate layers and at least one output layer. Data is processed in numerical form by multiplying it with the input layer's weights, mapping each product with the input layer's activation function and then propagating the resulting activations through subsequent layers in the same fashion. Figure [1.1](#) gives a graphical representation of a multi-layer feedforward network.

1.2 GOALS OF THIS THESIS

The objective of this work is twofold:

1. Create a software library that enables easy and reusable implementation of training metrics
2. Perform experiments on common datasets to investigate whether common problems in neural network training can be detected by the use of appropriate metrics. Issues to be investigated are
 - bad initializations
 - inappropriate learning rate

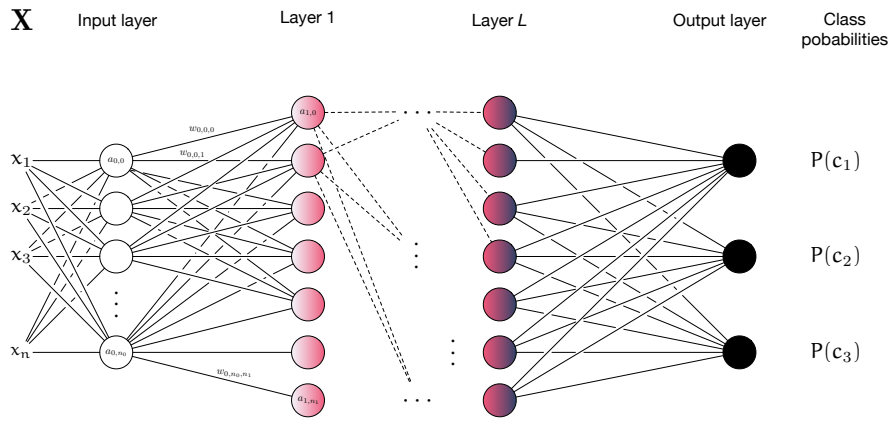


Figure 1.1: Schema of a multi-layer neural network. x_i are the input values, $a_{l_n,i}$ the i -th activation in layer l_n and $w_{l_n,i,j}$ the weight between the i -th unit in layer l_n and the j -th unit in layer l_{n+1}

- layer/model saturation
- inappropriate network architecture
- bad generalization/overfitting

1.3 MOTIVATION

In contrast to classical machine learning models, training deep neural networks requires navigating a huge parameter space. While most non-neural regression or classification algorithms only require specification of a parameter set up-front and often no more than a few, some parameters can (and should) be varied over training time for neural networks. Looking at the popular scikit-learn library, it can be seen that traditional methods such as SVMs, Gaussian Processes, Decision Trees or Gradient Boosting typically require less than 10 algorithmic parameters.¹

In neural networks the parameter space can have arbitrarily many dimensions when factoring in the fact that some parameters can change over time, such as

- learning rate (can be annealed)
- batch size²
- trainability of layers

Other parameters that need to be set initially are

- Network architecture (how many layers, how many units per layer)
- nonlinearity function
- type of loss
- optimization algorithm
- initial learning rate

¹ A look through [scikit-learn's](#) selection of regressors and classifiers shows most classes require between 5 and 10 parameters. (Pedregosa et al., 2011)

² It is not usual to change the batch size during training, but it can have an effect similar annealing the learning rate (see Smith et al. (2017))

- momentum
- weight decay

This makes finding an optimal training regimen very hard, particularly since training deep neural networks for realistic problems can take much longer than traditional methods, meaning cross-validating different models can be prohibitively expensive. It is therefore desirable to notice dead ends early during training, or be able to tweak parameters in such a way as to maximize convergence speed.

This thesis work is motivated by the scarcity of useful tools to debug and monitor deep learning training. Without years of training and a lot of mathematical intuition and expertise, it is often very hard to figure out why a network isn't learning or how to ensure timely convergence. And even with this expertise, visualizations or metrics need to be implemented over and over again because common tools do not abstract from the concrete model architecture.

There exist a some of monitoring tools (see section 1.4), but they are mostly low-level tools which provide visualization primitives (drawing and interacting with graphs). They may enable visualization of certain network metrics on top of the primitives, but There is no native support for a concept such as *Maximum singular value of the weight matrix* which can be simply applied automatically to all layers.

In contrast, the library developed in this work is geared towards modularizing introspection metrics in such a way that they are usable for any kind of model, without modifications to the model code. The secondary purpose of the library is the enablement to quickly iterate on hypothesized metrics extracted from the training in order to diagnose problems such as those outlined in section 1.2

As such, the library shall not only be useful to end users who will make use of established metrics and thus save time in their model training, but also to researchers and the author of this thesis in for evaluating hypotheses about training metrics.

1.4 EXISTING APPLICATIONS

TensorBoard

TensorBoard is a visulization toolkit originally developed for the TensorFlow (Abadi et al., 2015) deep learning framework. It is composed of a Python library for exporting data from the training process and a web server which reads the serialized data and displays it in the browser. The server can be used independently from TensorFlow, provided the data is serialized in the appropriate format. This enables, e.g., a PyTorch port, termed TensorBoardX.

For exporting data during training, the developer adds operations to the graph which write scalars, histograms, audio, or other data asynchronously to disk. This data can then be displayed in approximately real-time in the web browser. Besides scalar-valued functions, which could be e.g. the loss curve or accuracy measure, TensorBoard supports histograms, audio, and embedding data natively. However, concrete instances of these classes of training artifact must be defined by the user and can only be reused if the developer creates a separate library for the computations involved.

New kinds of visualizations can be added with plugins, which require not only writing the Python code exporting the data and for serving it from the web server, but also JavaScript for actually displaying it (the Polymer library is used for this³).

An attempt to abstract over the the programming language for talking to the server is *Crayon* which so far supports Python and Lua.

Visdom

Visdom by Facebook Research fulfills more or less the same purpose as TensorBoard, but supports Numpy and Lua Torch. In contrast to TensorBoard, Visdom includes more features for organizing the display of many visualizations at once. Still, the framework is mostly geared towards improving workflows for data scientists, and is not concerned with providing useful metrics out-of-the-box.

Others

There are other tools such as *DeepVis* for offline introspection, which offer insights into the training after the fact, but do not help guiding the training process while it's running.

³ <https://www.polymer-project.org/>

Ikkuna is the Python library developed for this thesis. It targets Python 3.6 and was designed with the following goals in mind:

1. Ease of use. Minimal configuration, maximum rewards.
2. Flexible and all-encompassing API enabling creating arbitrary metrics which act on training artifacts
3. Metrics shall be agnostic of model code.
4. Plugin architecture so metrics written once can be used for any kind of model
5. Framework agnosticism. Ideally, the library would support every deep learning framework through an extensible abstraction layer.

What it provides over the aforementioned tools is that it enables working at a higher level of abstraction, liberating the developer from having to repeat herself, exchanging visualizations and metrics and reduce the friction between development and debugging. This chapter gives a high-level overview of the library components and elaborates on the design decisions made during the creation. Throughout the chapters, UML class and package diagrams will serve as a mental map for the reader. For brevity, not all parts of the library are diagrammed down to the same level of detail.

2.1 DESIGN PRINCIPLES

Of the aforementioned goals, all except one have been accomplished. The objective of making the library agnostic to the deep learning framework being used (TensorFlow, PyTorch, PyCaffe, Chainer, etc.) has been neglected for practical reasons. Enabling this kind of support is beyond the scope of this thesis and only requires the implementation of a software layer which offers framework-agnostic access to network modules, activations, gradients and all the other necessary information. While this is certainly possible and useful, the PyTorch framework has been chosen for this work to create a proof of the concept. The choice is motivated in section 2.2.

The overarching architecture of this software must lend itself to this agnosticity goal, however. As such, a very loose coupling between model code, metric computation and visualizations is desired. Not only will this aid in extending the library to different deep learning frameworks, but it is also a prerequisite for allowing for modular, self-contained visualizations or metrics which can be installed and used separately and independently of specific model code. The Publisher-Subscriber design pattern has been chosen for these reasons (section 2.3).

2.2 DEEP LEARNING FRAMEWORKS

The currently available deep learning libraries can be located on a spectrum between define-by-run and define-and-run. The first extreme would be a framework such as

PyTorch (Paszke et al., 2017) or Chainer (Tokui et al., 2015), where there exist no two distinct execution phases – just like in an ordinary matrix library like NumPy, each statement immediately returns or operates on an actual value. By contrast, graph-based frameworks like TensorFlow¹ require specifying the model graph in a domain-specific language (TensorFlow has Python, Java and C++ APIs, Caffe uses Prototxt files), compile it to a different representation and the model is run and trained in a second phase. While this enables graph-based optimizations, the main downsides are that

- control flow cannot use the host language features, but must be done with the API used for defining models. Instead of

```
counter = torch.tensor(0)
# repeated matrix multiplication
while counter < tensor:
    counter += 1
    h = torch.matmul(W, h) + b
```

one must use a construction like this

```
counter = tf.constant(0)
while_condition = lambda counter: tf.less(counter, tensor)
# loop body
def body(counter):
    h = tf.add(tf.matmul(W, h), b)
    # increment counter
    return [tf.add(counter, 1)]

# do the actual loop
r = tf.nn.whole_loop(while_condition, body, [counter])
```

- As a corollary, the barrier of entry is higher, since a beginner cannot rely on the language feature she knows but must learn how to express many concepts without the host language.
- halting execution at arbitrary points in the training is not possible, since the actual training is not happening in the host language, but is more often handed off to lower-level implementations in its entirety.

This makes conditional processing and debugging much less ergonomic.

All frameworks have in common that they build a graph representation of the model, whether implicitly or explicitly. Nodes in the graph are operations while edges are data flowing between operations. This allows naturally parallelizing independent computations. To compute gradients, the graph can be traversed backwards from the output node by applying the chain rule of differentiation. Define-and-run frameworks like TensorFlow create the graph explicitly; the user uses the API to do exactly this. The graph – once compiled – is fixed for the entire training process. PyTorch on the other hand implicitly records all operations and

¹ Since version 1.4, TensorFlow gravitates toward define-by-run through the introduction of *eager execution*, which becomes the default mode in version 2.0. Graph-based execution is still available, but not the default any longer.

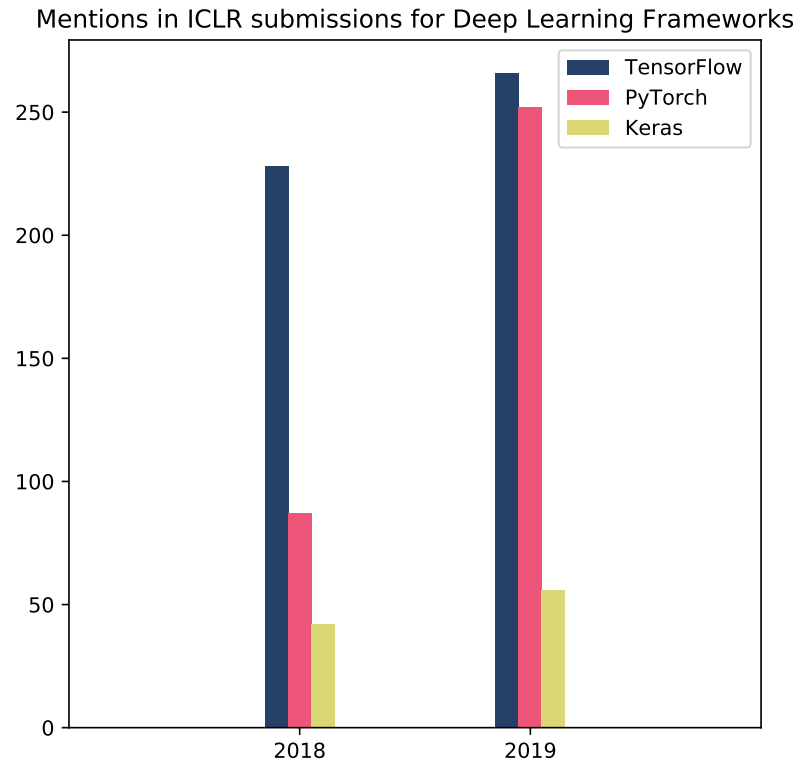


Figure 2.1: Changes in popularity of different deep learning libraries in research. Data was collected by keyword search over ICLR submissions (<http://search.iclr2019.smerity.com/search>; analogously for 2018)

also overloads operators for this purpose. The graph is thus recreated for each propagation through the network. This precludes some optimizations, but makes dynamically changing networks easily achievable.

For this work, the PyTorch framework has been chosen, due to the fact that it is growing quickly in popularity (see figure 2.1) and relatively new, so the ecosystem is not fully developed and some utilities available for e.g. TensorFlow are not available for PyTorch. Because of this, an introspection framework for training monitoring is judged to present the best value proposition for PyTorch users.

2.3 PUBLISHER-SUBSCRIBER

The Publisher-Subscriber pattern (for a detailed overview see [Eugster et al. \(2003\)](#)) is a pattern for distributed computation in which publishers publish messages either directly to any subscribers which have registered interest in them, or to a central authority orchestrating the exchange. Messages are generally associated with one or more topics and subscribers register interest in receiving messages on one or more topics.

The components are very loosely coupled; the subscribers need not even be aware of the publishers at all, and the publishers' only interaction with their subscribers is relaying messages through a uniform interface or through an optional server. A graphical schema of one possible incarnation of this pattern is shown in figure 2.2.

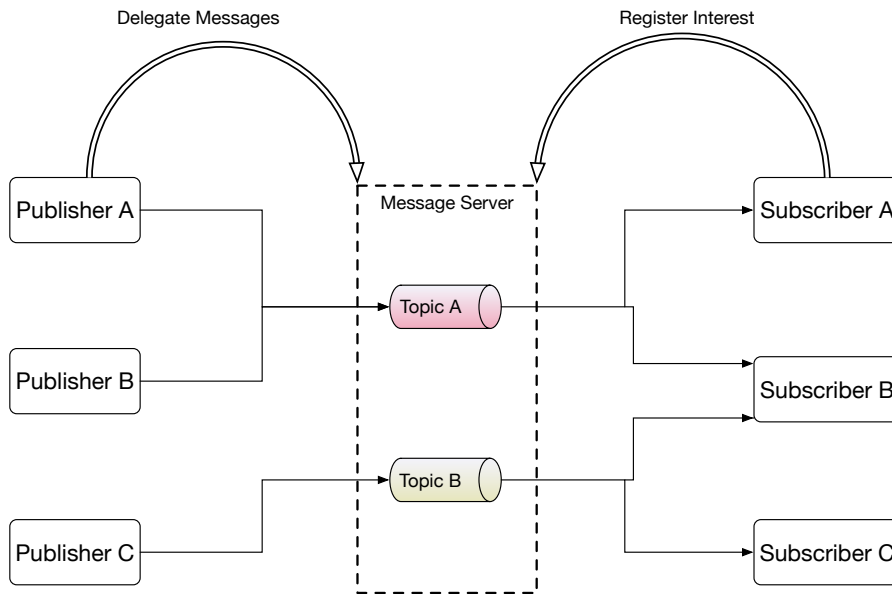


Figure 2.2: One possible implementation of the Publisher-Subscriber pattern.

This project is not distributed, but can benefit from the loose coupling in another way: Subscribers can be defined in terms of the kind of messages they need to compute their metric, without knowing anything about where the messages are coming from. Concretely, as long as the appropriate data is emitted from the training process, subscribers can work without modifications with any possible model.

Since real-world neural networks are trained on the GPU, and communication between host and GPU memory is already expensive, making this library truly distributed is not an objective. However, the design will simplify asynchronous computation of metrics in the future. The Python language does not support true multithreading², but since the expensive part of the work is running on the GPU while the host code is waiting, metric computation could happen asynchronously on the GPU as well while the expensive forward or backward passes through the network are running. This is not currently implemented but can be added later, if more computationally demanding metrics are to be explored.

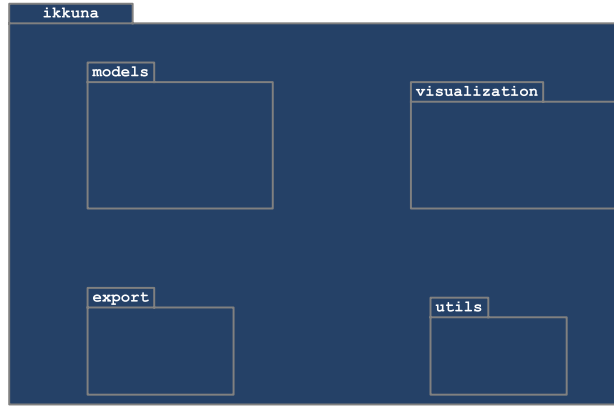
In the context of neural network training, there is only one source of information and hence only one publisher. Nevertheless, a message server is introduced to segregate responsibilities. The singular publisher extracts data from the training model, passes it on to the server which also accepts subscriber registrations and relays messages appropriately.

2.4 OVERVIEW OF THE LIBRARY

The software is structured into several packages. The root package is `ikkuna` which encapsulates all core functionality. All other packages and modules contain utilities implemented for this work specifically, but will generally not be relevant to other users. A survey of these tools will be given in section 2.4.5.

The root package diagram is shown in figure 2.3

² The `multiprocessing` module allows for truly asynchronous computation and communication, but the inter-process-communication is more expensive than memory shared between threads.

Figure 2.3: `ikkuna` package diagramTable 2.1: `ikkuna.export` functionalities

Name	Function
<code>export</code>	Publish data from an arbitrary model and send messages to registered subscribers
<code>messages</code>	Define message interface; i.e. what topics exist and which information a message must contain
<code>subscriber</code>	Define the base class for metric subscribers

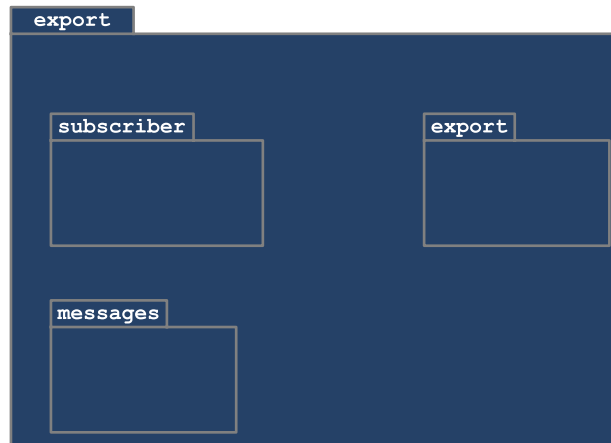
The `models` (see section 2.4.2) subpackage contains a few exemplary neural network definitions which are wired up with the library and can thus be used to showcase the library’s functionality. The `utils` (see section 2.4.3) subpackage contains miscellaneous utility classes and functions used throughout the core library. Lastly, the `visualization` subpackage (section 2.4.4) contains the plotting functionality to actually show the metrics computed during the training process.

The most important bits of the software live in the `export` subpackage (section 2.4.1). It implements the Publisher-Subscriber pattern. Extracting data from the training process, defining subscriber functionality and messages used for communication is done here.

2.4.1 The *export* subpackage

The `export` subpackage contains the core part of the library, i.e. it provides the classes that handle discovering the structure of the neural network model, attaching the appropriate callbacks and intercepting method calls on the model so the library is informed about everything entering and exiting the model and its individual layers. It also contains the definition for the subscriber API, i.e. the messages that subscribers can receive, synchronisation facilities when multiple topics are needed by a subscriber, as well as the subscriber class interface. The package diagram is displayed in figure 2.4.

The package comprises three subpackages or modules listed in table 2.1

Figure 2.4: `ikkuna.export` package diagram*The `ikkuna.export` subpackage*

In in slight deviation from the Publisher-Subscriber framework as displayed in figure 2.2, the `export.Exporter` class (figure 2.6) is the sole publisher of data. There is only one source of data during training, so it is unnecessary to accomodate multiple publishers. The `Exporter` is informed of the model with its methods `set_model()` and `set_loss()`, the latter of which is only necessary if metrics which rely on training labels should be displayed. It can accept a filter list of classes which are to be included when discovering the modules in the model. For instance, it could be desirable to only observe layers which have weights and biases associated with them, not e.g. normalisation or reshaping layers. The `Exporter` then traverses the model (which is really just a tree structure of modules) and adds to each a callback invoked when input enters the layer – in order to retrieve activations – and when gradients are computed for the layer outputs. The callbacks also use cached weights – if present – in order to publish updates to the weights. Furthermore, it replaces a few of the model’s methods with closure wrappers so it can

- be notified when the model is set to training or testing mode (this switch disables or enables layers which only make sense during one of the phases³)
- increase its own step counter automatically when a new batch is seen
- add a parameter to the model’s `forward()` method – called by the runtime when data is propagated through the model – which can be used by subscribers to temporarily turn off training mode and have it revert automatically. This is useful for subscribers which need to evaluate the model (i.e. feed data through it), but do not want to generate new messages for this occasion.
- intercept labels passed to the loss function during training and publish them as messages so the user need not concern himself with this task

³ There are two built-in layers this applies to. One is the batch normalisation layer. It normalises the output of the previous layer with the mean and variance over the entire batch of data – optionally with running means and variances over the previousprevious training steps. The variance is not defined for single data point enters the layer, as could be the case during inference/testing time. The second case is the dropout layer, which randomly zeroes out a percentage of the previous layer’s activations. This is used during training to prevent subsequent units from becoming correlated with a fixed set of units in the previous layer, instead of picking up patterns invariant of where in the input they occur. During inference time, this is turned off to make full use of the trained layers.

- intercept the final output of the network. This could be realised alternatively by identifying the last module in the network.

At every time step (training batch), the `Exporter` publishes the following information on to the message bus (see figure 2.5):

- gradients for each module
- activations for each module
- weights and biases for each module that has these properties (e.g. convolutional or fully-connected layers)
- updates to the weights and biases from the last step to the current one, provided the module has these properties
- Training labels used for the parameter updates. This requires that the `Exporter` be informed of the loss function object with `set_loss()`.
- The batch of input data passed to the network at the current training step
- The final output of the network for the current batch of training data. This is simply the tensor of activations from the last layer and is thus technically duplicated since activations are published anyway. The reason is that some subscribers may only be interested in the network predictions and it is unnecessary to determine automatically the last layer in the network as the loss function has automatic access to the activations and must be tracked anyway for the training labels

Further messages are published only at certain points in the training process

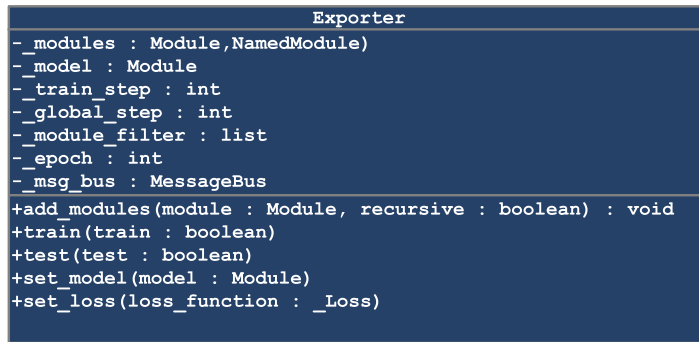
- When a batch starts or ends, a message with the current batch index is published
- When an epoch starts or ends, a message with the current batch index is published. This requires the `Exporter` be notified with `epoch_finished()` by the user, since it is impossible to determine when an epoch is over from inside the model.

The `ikkuna.export.messages` submodule

This submodule contains definitions of all permissible messages kinds, message classes and a collection class for message objects. An overview of the classes defined in this module is shown in figure 2.5.

Messages are of one of two types: They are either directly tied to a layer in the network and are thus published for each layer, or they contain information for the current training step applying to the entire network. In that case, they appear only once per training step, not once per layer. The meta-messages can carry tensor data (e.g. input data or labels), but need not to (e.g. notifications about a starting or ending epoch). All message kinds are summarised in table 2.2.

Messages can be assembled into bundles if a subscribers wants to subscribe several topics at once. The `MessageBundle` class performs all necessary error checking to ensure consistency of the contained messages.

Figure 2.5: Classes in the `ikkuna.export.messages` submoduleFigure 2.6: `ikkuna.export.Exporter` class diagram

The message server from figure 2.2 is implemented by the `MessageBus` class which publishers publish messages onto, accepts subscriber registrations, and maintains the lists of known topics. Each subscriber may in turn announce one or more new topics which can then be subscribed to by others. This is useful since it allows chaining of subscribers in order to realise arbitrary post-processing of computed metrics.

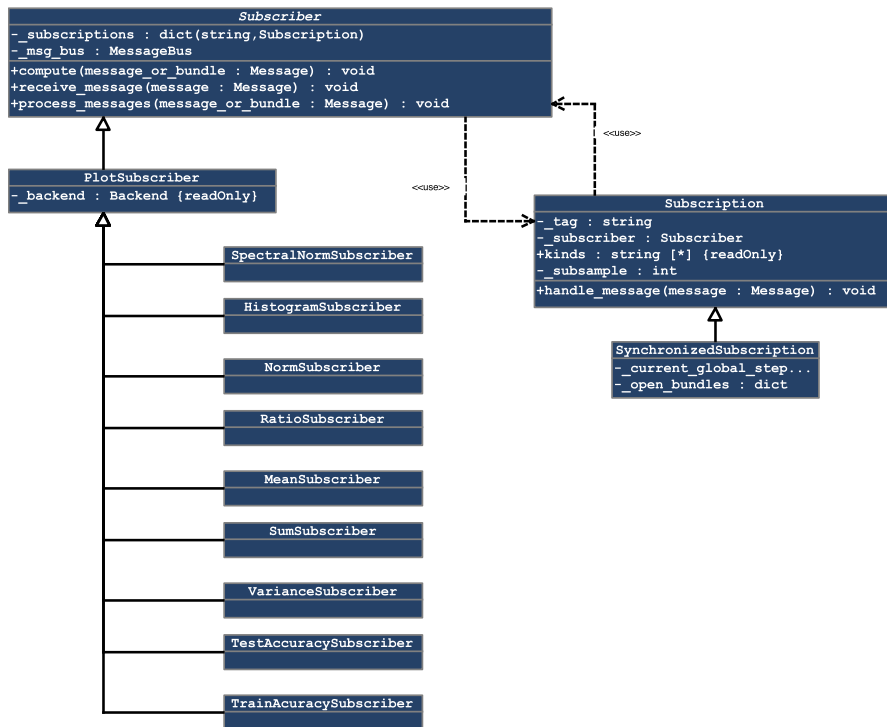
The `ikkuna.export.subscriber` subpackage

The third subpackage contained in the `ikkuna.export` package defines the subscriber part of the Publisher-Subscriber pattern. The diagram of the defined classes is shown in figure 2.7. The `Subscriber` base class is rudimentary and mandates only the implementation of the metric computation by subclasses. In the simplest case, a subscribers is interested in only one topic and therefore is coupled to a simple `Subscription` object, which handles bookkeeping tasks such as subsampling the message stream, routing only relevant messages to the subscriber and counting the received messages.

More generally however, a subscriber may want to receive several pieces of information for each layer in each time step (i.e. for computing the ratio between

Table 2.2: Subscribable message kinds

Meta topics		
Identifier	Frequency	Description
batch_started	Once every batch	
batch_finished	———— " ————	
epoch_started	Once every epoch	
epoch_finished	———— " ————	
input_data	Once every batch	
input_labels	———— " ————	
network_output	———— " ————	Activations of the last layer
Data topics		
Identifier	Frequency	Description
weights	Once per layer per batch	Gradients of loss function w.r.t. layer weight matrix
weight_gradients	———— " ————	
weight_updates	———— " ————	
biases	———— " ————	Gradients of loss function w.r.t. layer bias matrix
bias_gradients	———— " ————	
bias_updates	———— " ————	
activations	———— " ————	
layer_gradients	———— " ————	Gradients of loss function w.r.t. layer output

Figure 2.7: Classes defined in `ikkuna.export.subscriber`

weight updates and weights). Since the order of messages is not guaranteed, the desired messages are unlikely to occur one after the other; instead the topics must be synchronised. A `SynchronizedSubscription` buffers messages of the relevant topics until all requested kinds have been received for the current training step, before releasing them to the subscriber.

A subscriber can thus receive a single message or a bundle of messages.

The library comes with a few subscribers already installed (they are themselves plugins, see section 2.4.6). Details are given in table 2.3

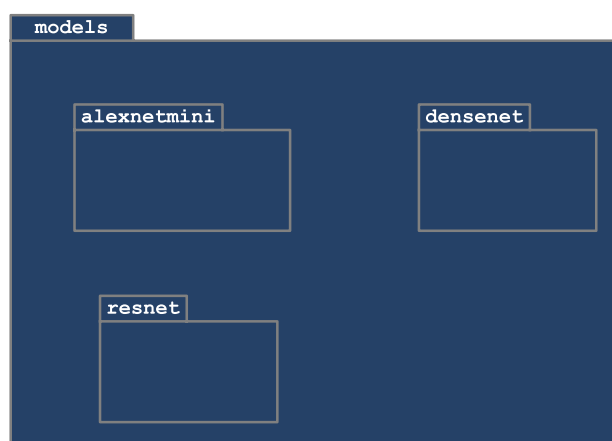
Figure 2.8: `ikkuna.models` package diagram

Table 2.3: Pre-packaged subscriber subclasses

Name	Functionality
MeanSubscriber	Computes the mean $\mu = \frac{1}{n} \sum_{i=1}^n w_i$ of a tensor
VarianceSubscriber	Computes the variance $\sum_{i=1}^n (w_i - \mu)^2$ for a tensor
SumSubscriber	Computes the sum $\sum_{i=1}^n w_i$ of a tensor
NormSubscriber	Computes the p-Norm $\sqrt[p]{\sum_{i=1}^n w_i^p}$
RatioSubscriber	Computes the average ratio $\frac{1}{n} \sum_{i=1}^n \frac{a_i}{b_i}$ of two tensors, disregarding NaN values
HistogramSubscriber	Computes the histogram of a given tensor. This is computationally heavy.
SpectralNormSubscriber	Computes the spectral norm (largest singular value) $\max_{\mathbf{h}: \ \mathbf{h}\ _2=1} \frac{\ \mathbf{A}\mathbf{h}\ _2}{\ \mathbf{h}\ _2}$ of a tensor
TestAccuracySubscriber	Computes the ratio of correctly classified examples to total examples over the test set
TrainAccuracySubscriber	Computes the ratio of correctly classified examples to total examples over current batch of training data

2.4.2 The models subpackage

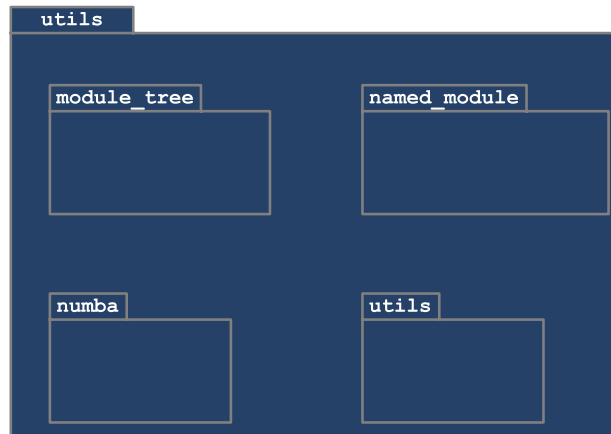
This package shown in figure 2.8 contains model definitions for demonstration purposes and for experimentation. Three architectures are currently implemented:

1. A minified version of AlexNet, since the original architecture requires larger images (Krizhevsky et al., 2012). The code is adapted from <https://github.com/sukilau/Ziff-deep-learning/blob/master/3-CIFAR10-lrate/CIFAR10-lrate.ipynb>.
2. DenseNet (Huang et al., 2017). The implementation is basically the one from (Pleiss et al., 2017)⁴ with minor modifications
3. ResNet (He et al., 2016). This implementation comes from GitHub user liukang⁵ and can handle CIFAR10-sized images of 32 pixels per side, as opposed to most implementations that are geared towards ImageNet examples which are much larger.

All models are modified such that their training can be supervised by the library.

⁴ At the time of writing, the implementation is available here: https://github.com/gpleiss/efficient_densenet_pytorch/blob/master/models/densenet.py. The licensing is unclear as the author references the original BSD-licensed implementation at <https://github.com/pytorch/vision/blob/master/torchvision/models/densenet.py> which was licensed by PyTorch core contributor Soumith Chintala. However, the code does not reproduce the BSD license text and can thus only be inspired by the original but cannot contain any of the code verbatim. It would require careful examination in order to determine whether this is the case.

⁵ The implementation is MIT-licensed. <https://github.com/kuangliu/pytorch-cifar/blob/master/models/resnet.py>

Figure 2.9: `ikkuna.utils` package diagram

2.4.3 The *utils* subpackage

As shown in figure 2.9, this package defines classes for traversing a model into a hierarchical tree of layers (called *modules* in PyTorch lingo), structures for adding information to PyTorch’s `Module` class, and a set of miscellaneous functions for

1. Seeding random number generators to make experiments reproducible (see appendix B)
2. Creating instances of weight optimizers by named
3. Initialize the weights of any model
4. Loading datasets

Additionally, it contains the `numba` module which is intended to allow interoperability with the Numba library⁶. While currently not used due to the incomplete nature of the Numba GPU array interface, it could enable leveraging Numba in the future without transferring data to the CPU. The core function was later obsoleted by an addition to the PyTorch library⁷.

2.4.4 The *visualization* subpackage

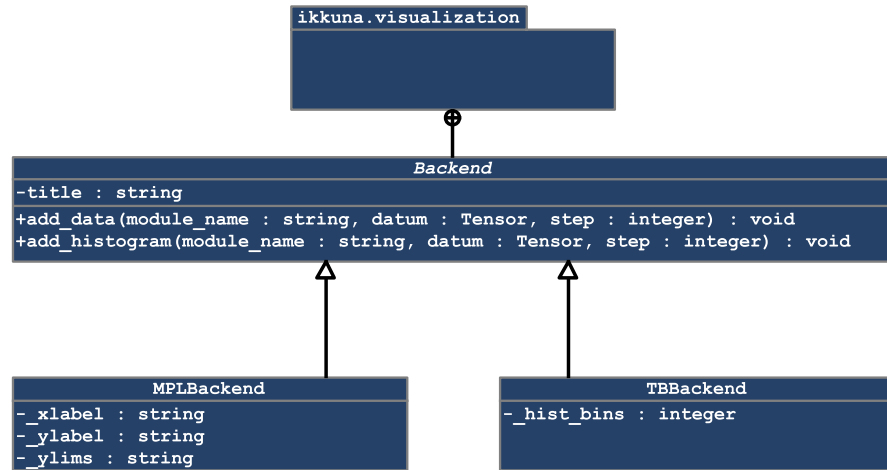
This package contains only a single module: `backend`. It defines the classes shown in figure 2.10. The module serves as an abstraction over plotting libraries so that metrics need not concern themselves with how to actually show the data.

A given metric will compute its value and dispatch it to its visualization backend, which can currently accept scalar and histogram data. The metric class itself need not care about how it is going to be displayed.

For running the library locally, a `matplotlib`-based backend has been implemented. Plotting routines from this library open a window directly on the system

⁶ <https://numba.pydata.org/>. Numba is a library for transforming high-level Python code into performant compiled code and for allowing to use the CUDA library from Python with Python arrays. This enables performance improvements for numeric calculations, but there is only a limited set of higher-level functions implemented on GPU arrays.

⁷ The main contribution of the submodule was to make PyTorch tensors accessible to Numba by monkey-patching the `__cuda_array_interface__` property. This has since been added via pull request #11984 to the PyTorch repository.

Figure 2.10: Class diagram for classes in `ikkuna.visualization`

executing the software. In practice however, deep learning code will be executed remotely on a server with adequate compute capability and the programmer connected via SSH. While it is possible to have remote windows show up locally on Linux-based systems by use of X11-Forwarding, this is generally slow and not useful for interactivity. An example is shown in figure 2.11 To remedy this issue, a plotting backend for TensorBoard (see section 1.4) is also provided. The plotting data is generated and processed on the remote system, but served over the web so it can be viewed and interacted with locally (provided the network is configured so that the server responds to HTTP requests). An example is shown in figure 2.12.

2.4.5 Miscellaneous tools

There are a few modules which simplify development with the library but are not part of the distribution obtained from PyPi or by running the setup script.

The `train` package defines a `Trainer` class which encapsulates all the logic and parameters needed to train a neural network on one of the datasets provided with PyTorch. The class's capabilities include the following

- Look up model and dataset by name
- Bundle all hyperparameters
- hook the `Exporter` into the model for publishing data
- configure the optimisation algorithm to use for training
- train the model for one batch

The `Trainer` class is used in the main script (`main.py`), which serves as a command line interface to the library while developing. When trying out the library, it can also be used as an initial starting point.

The library can be installed to the local Python environment by use of the provided `setuptools` script (`setup.py`). It can also be downloaded from the [Python Package Index](#) by use of the package manager `pip`:

```
pip install ikkuna
```

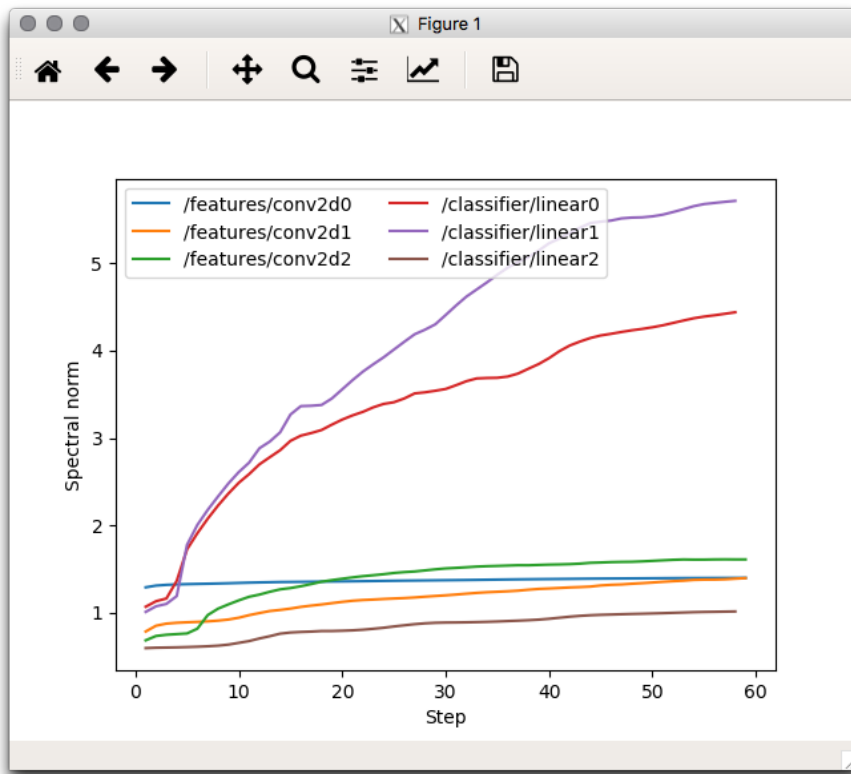



Figure 2.11: Exemplary view of a matplotlib figure forwarded over SSH

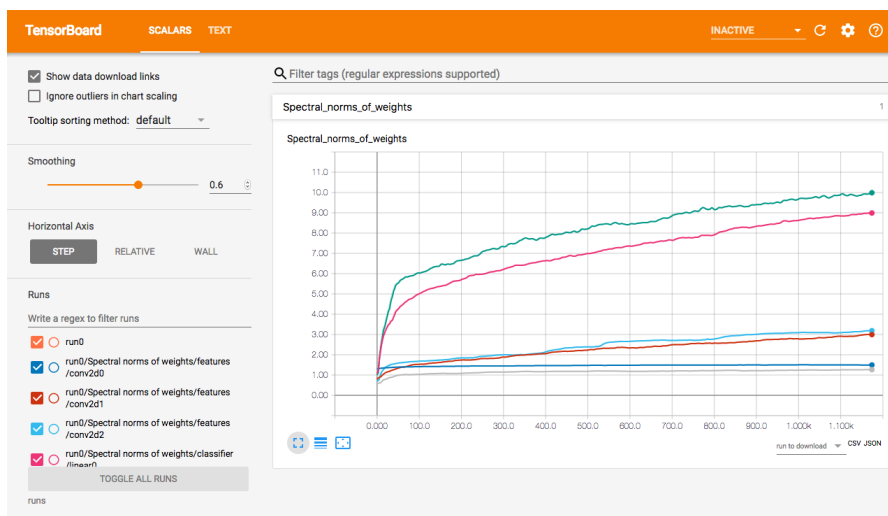


Figure 2.12: Exemplary view of a TensorBoard session

Table 2.4: Named arguments to `main.py`

Parameter	Explanation
<code>-m,--model</code>	Model class to train
<code>-d,--dataset</code>	Dataset to train on. Possible choices: MNIST, FashionMNIST, CIFAR10, CIFAR100
<code>-b,--batch-size</code>	Default: 128
<code>-e,--epochs</code>	Default: 10
<code>-o,--optimizer</code>	Optimizer to use. Default: Adam
<code>-a,--ratio-average</code>	Number of ratios to average for stability (currently unused). Default: 10
<code>-s,--subsample</code>	Number of batches to ignore between updates. Default: 1
<code>-v,--visualisation</code>	Visualisation backend to use. Possible choices: <code>tb</code> , <code>mpl</code> . Default: <code>tb</code>
<code>-V,--verbose</code>	Print training progress. Default: <code>False</code>
<code>--spectral-norm</code>	Use spectral norm subscriber on weights. Default: <code>False</code>
<code>--histogram</code>	Use histogram subscriber(s)
<code>--ratio</code>	Use ratio subscriber(s)
<code>--test-accuracy</code>	Use test set accuracy subscriber. Default: <code>False</code>
<code>--train-accuracy</code>	Use train accuracy subscriber. Default: <code>False</code>
<code>--depth</code>	Depth to which to add modules. Default: -1

2.4.6 Plugin Infrastructure

Among the main selling points of this library is the provision to add new metrics as plugins and reuse them system-wide for all architectures. Plugins in Python projects can be enabled through appropriate use of the `setuptools` library. During the setup process for installing the library, entry points are defined by the library which can be used by plugins to announce themselves. Ikkuna provides the `'ikkuna.export.subscriber'` entry point. For registering a plugin, the author must simply use that entry point to make a plugin available either through the package it is defined in, or through the `ikkuna.export.subscriber` namespace. For illustration, listing 2.1 shows how to setup a `setup.py` setuptools file. The plugin can be installed like any other Python package with

```
python setup.py install
```

which will install all required dependencies inside the current environment. The PyTorch library must be installed manually since the binary distribution is too old at the time of writing. Detailed instructions can be found in the user guide which is part of the documentation.

```
#!/usr/bin/env python

from distutils.core import setup
import setuptools

setup(name='<your package name>',
      version='<version>',
      description='<description>',
      author='<your name>',
      author_email='<your email>',
      packages=['<package name>'],
      # ... any other args
      entry_points={
          'ikkuna.export.subscriber': [
              'YourSubscriber = module.file:YourSubscriber',
          ]
      })
```

Listing 2.1: Sample setup script for subscriber plugins

2.4.7 Documentation

The entire codebase is liberally documented using the Sphinx documentation processor⁸. The documentation contains further documents with a detailed user guide, installation instructions. Sphinx allows generating documentation in many formats from the same source, most usefully HTML and PDF. At the time of writing, the HTML documentation and API reference is hosted at https://peltarion.github.io/ai_ikkuna/.

⁸ <http://www.sphinx-doc.org/>

3

EXPERIMENTS IN LIVE INTROSPECTION

In this chapter, several hypotheses for diagnosing roadblocks in the training process are introduced, motivated, and evaluated by use of the Ikkuna library. This thesis is not concerned with advancing the state of the art in classification. Instead, toy problems are developed in order to prove or disprove that the proposed metrics have the potential to be useful in guiding the training. A more thorough evaluation of the results on realistic architectures and problem sizes would require significant computational resources and is left for future work.

3.1 DETECTING LEARNING RATE PROBLEMS

The first problem to be investigated is that of choosing an appropriate learning rate. In standard stochastic gradient descent, a loss J of some the model parameters θ (here, the layer weights) is computed over the training set of m examples by forming the expectation over the sample loss L :

$$J(\theta) = E_{\mathbf{x}, \mathbf{y} \sim \hat{p}_{\text{data}}} L(\mathbf{x}, \mathbf{y}, \theta) \quad (3.1)$$

$$= \frac{1}{m} \sum_{i=1}^m L(\mathbf{x}^{(i)}, \mathbf{y}^{(i)}, \theta) \quad (3.2)$$

The cumulative loss can then be derived for θ

$$\nabla_{\theta} J = \frac{1}{m} \sum_{i=1}^m \nabla_{\theta} L(\mathbf{x}^{(i)}, \mathbf{y}^{(i)}, \theta) \quad (3.3)$$

per the sum rule of differentiation. The simplest form of parameter update rule is then

$$\theta \leftarrow \theta - \eta \nabla_{\theta} J \quad (3.4)$$

with the learning rate η . There is no hard and fast rule on what this parameter should be, and it is subject of large swathes of literature. Popular modifications to the vanilla update rule are the use of momentum (Jacobs, 1988), per-layer learning rates (ibid.), reducing the rate throughout training, or adapting the learning rate based on mean and variance of parameters across past time steps (Kingma and Ba, 2014). Nevertheless, most of the time, training begins with an arbitrarily chosen small learning rate around 0.01 which is then adapted either by the aforementioned mechanisms or by search over the parameter space when computationally feasible.

In this section, the Ikkuna library is employed to find correlates of inappropriate learning rates which could then be used to adjust its value during training.

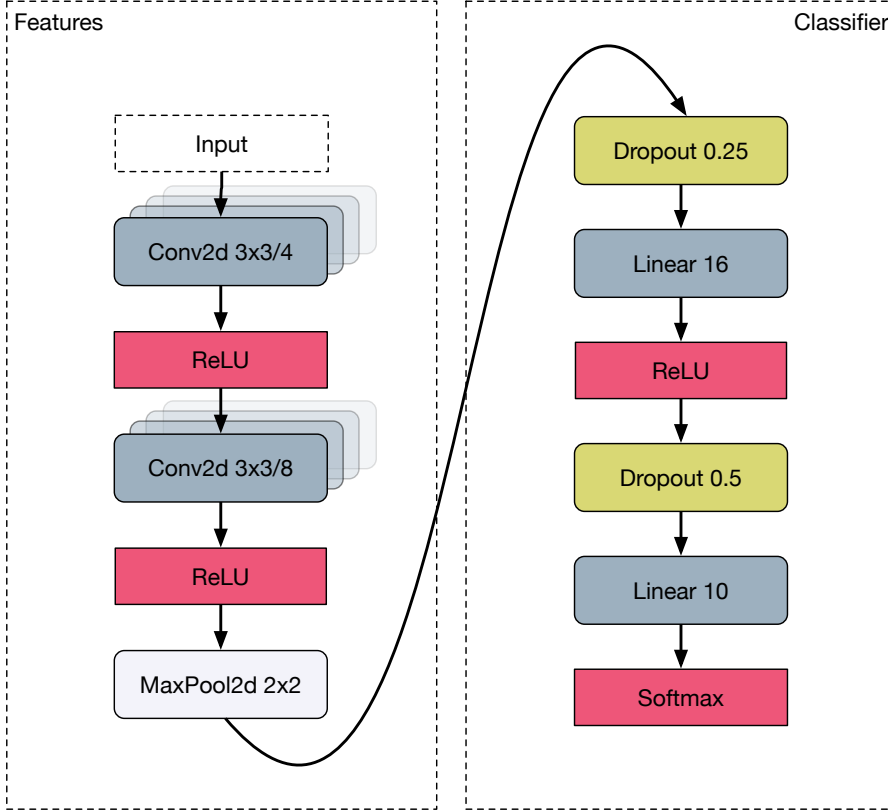


Figure 3.1: The network used in this experiment. Image features are extracted by 4 and 8 convolutional filters, respectively, with ReLU nonlinearities. Maximum pooling is applied with a filter and stride size of 2 leading to a resolution a fourth of the original size. The classifier portion employs dropout layers to reduce co-adaptation of units and a final softmax activation to map outputs to $(0, 1)$.

3.1.1 Experiment 1

The first experimental setup uses the simplified AlexNet architecture (section 2.4.2) shown in figure 3.1. The dataset used is the well-known CIFAR-10 dataset consisting of 60,000 32×32 colour images from ten object categories (Krizhevsky and Hinton, 2009). The optimization algorithm is plain stochastic gradient descent on minibatches of size 128. The dataset does not constitute a hard problem to solve; state of the art accuracies lie around 95%. For this reason, a decision must be made about how to make the problem hard enough so that improvements to the training schedule can actually be made. The learning rate has thus been fixed to a high value of 0.2 which is not the optimal value (a learning rate of 0.1 solves the problem to 45% accuracy).

In order to validate that there is room for improvement (i.e. the task is not too easy), the training has been run about twenty times for both a constant learning rate and an exponentially decaying rate according to

$$\eta_{e+1} = 0.98^{e+1} \eta_e, \quad (3.5)$$

e being the epoch index. The final accuracies after 100 epochs of training for constant learning rate, exponential decay are shown in figure 3.2. As can be seen, there is a significant improvement when decaying the learning rate over keeping it constant.

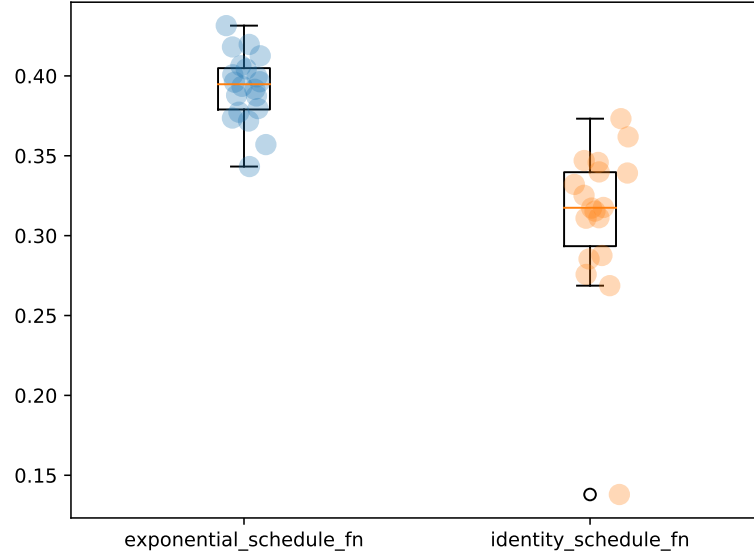


Figure 3.2: Final accuracies after 100 epochs

As a first proof of concept, an adaptive learning rate based on the ratio between weight updates and weights is implemented with the help of the library.

Let l be the number of layers with weight matrices associated with them (for instance linear or convolutional layers, but not activation functions, dropout, or the like). Let $\{W_{i,k} \mid i = 0 \dots l-1\}$ be the set of weight matrices at training step k . Let η be the base learning rate and $\frac{1}{t}$ be a target value to which we want the update-to-weight ratio to move. Furthermore, let $\gamma \in (0, 1)$ be a decay factor for exponential smoothing. Now, let

$$R_{i,k} = \frac{\|W_{i,k} - W_{i,k-1}\|_2}{\|W_{i,k}\|_2} \quad (3.6)$$

be the ratio between the L2-Norms of layer i 's weight updates before step k step and the weights at step k themselves. We then select the new learning rate for batch step $k+1$ as

$$\eta_{k+1} = \eta_k \left(t \cdot \frac{1}{l} \sum_{i=0}^{l-1} \gamma R_{i,k} + (1-\gamma) R_{i,k-1} \right)^{-1} \quad (3.7)$$

for $k \geq 2$. This is the average exponentially smoothed update-weight-ratio, divided by the target range. This learning rate is used for vanilla gradient descent without any other modifications beyond capping it to some value in case of very small ratios. The effect of adapting the learning rate according to this schedule is that the average ratio between the weight updates and the weights moves towards the target range. The idea that the ratio between updates and weights is a quantity which should be constrained has not been extensively investigated. The claim is made by [Karpathy \(2015\)](#), who states that a target of $\frac{1}{t} = 10^{-3}$ is a reasonable value. ?? displays a set of accuracy traces for each of the schedules (constant, exponential

decay, ratio-adaptive) with different base learning rates. The network was trained from scratch 5 times for each combination.

The experiment shows surprising results, namely that a this simple scheme of dynamically adapting the learning rate to the magnitude of weight change at each training step easily outperforms any constant value for η and achieves almost state-of-the-art accuracy in some cases (the best reported performance by [Yamada et al. \(2018\)](#) is 97.69%).

As an impression of how the library presented in chapter 2 simplifies a general implementation of such a learning rate schedule, code is provided here.

When an `Exporter\` is configured for a given model, a `RatioSubscriber` (see table 2.3) must be added to the message bus in order for the update-weight-ratio ($R_{i,k}$ in the above equations) to be published. One can then subscribe them and process the information with this subscriber:

```
class RatioLRSubscriber(PlotSubscriber):
    def __init__(self, base_lr, smoothing=0.9, target=1e-3,
                 max_factor=500):
        subscription = Subscription(self, ['
            weight_updates_weights_ratio', 'batch_started'],
                                   tag=None, subsample=1)
        super().__init__([subscription], get_default_bus(),
                         {'title': 'learning_rate',
                          'ylabel': 'Adjusted learning rate',
                          'ylims': None,
                          'xlabel': 'Train step'})

        # exponential moving avg of  $R_{i,k}$ 
        self._ratios = defaultdict(float)
        # maximum multiplier for base learning rate (in
        # pathological cases)
        self._max_factor = max_factor
        # exp smoothing factor
        self._smoothing = smoothing
        # target ratio
        self._target = target
        # this factor is always returned to the learning rate
        # scheduler
        self._factor = 1
        self._base_lr = base_lr

    def _compute_lr_multiplier(self):
        '''Compute learning rate multiplicative. Will output 1
        for the first batch since no layer
        ratios have been recorded yet. Will also output 1 if the
        average ratio is close to 0.
        Will clip the factor to some max limit'''
        n_layers = len(self._ratios)
        if n_layers == 0: # before first batch
            return 1
        else:
            mean_ratio = sum(ratio for ratio in self._ratios.
                             values()) / n_layers
            if mean_ratio <= 1e-9:
                return 1
            else:
```

```

        factor = self._target / mean_ratio
        if factor > self._max_factor:
            return self._max_factor
        else:
            return factor

# invoked by the runtime for each incoming message
def compute(self, message):
    if message.kind == 'weight_updates_weights_ratio':
        # the 'key' property for these messages will be the
        # module/layer
        # here we compute the exponential moving average of
        # ratios
        i = message.key
        R_ik = message.data
        R_ik_1 = self._ratios[i]
        gamma = self._smoothing
        self._ratios[i] = (gamma * R_ik + (1 - gamma) *
                           R_ik_1)
    elif message.kind == 'batch_started':
        # before a batch starts, update the lr multiplier
        self._factor = self._compute_lr_multiplier()

def __call__(self, epoch):
    return self._factor

```

The subscriber implements the `__call__()` method so it can be dropped into PyTorch's learning rate scheduler (`torch.optim.lr_scheduler.LambdaLR`). This learning rate schedule can thus be used in every model, without modification.



APPENDIX A — OPEN SOURCE ACKNOWLEDGMENTS

The following presents a non-exhaustive list of open source tools used in creating the software, experiments and this document:

1. \LaTeX

B

APPENDIX B — CONTRIBUTIONS TO OTHER LIBRARIES

This appendix lists all code and documentation contributions made to other projects during development.

Table B.1: Contributions to third-party projects

Project	Contribution
PyTorch	Added a document to the official documentation detailing how to make experiments reproducible (see https://pytorch.org/docs/master/notes/randomness.html).
CuPy	Added NumPy-style ordering and comparison functionality for complex number types, thereby enabling all sorts of operations for complex matrices which did not work previously.
TensorBoardX	Documentation improvements

BIBLIOGRAPHY

- Abadi, M., Agarwal, A., Barham, P., Brevdo, E., Chen, Z., Citro, C., Corrado, G. S., Davis, A., Dean, J., Devin, M., Ghemawat, S., Goodfellow, I., Harp, A., Irving, G., Isard, M., Jia, Y., Jozefowicz, R., Kaiser, L., Kudlur, M., Levenberg, J., Mané, D., Monga, R., Moore, S., Murray, D., Olah, C., Schuster, M., Shlens, J., Steiner, B., Sutskever, I., Talwar, K., Tucker, P., Vanhoucke, V., Vasudevan, V., Viégas, F., Vinyals, O., Warden, P., Wattenberg, M., Wicke, M., Yu, Y., and Zheng, X. (2015). TensorFlow: Large-scale machine learning on heterogeneous systems. Software available from tensorflow.org.
- Ciresan, D. C., Meier, U., Masci, J., Gambardella, L. M., and Schmidhuber, J. (2011). Flexible, high performance convolutional neural networks for image classification.
- Eugster, P. T., Felber, P. A., Guerraoui, R., and Kermarrec, A.-M. (2003). The many faces of publish/subscribe. *ACM Comput. Surv.*, 35(2):114–131.
- He, K., Zhang, X., Ren, S., and Sun, J. (2016). Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 770–778.
- Huang, G., Liu, Z., Van Der Maaten, L., and Weinberger, K. Q. (2017). Densely connected convolutional networks. In *CVPR*, volume 1, page 3.
- Jacobs, R. A. (1988). Increased rates of convergence through learning rate adaptation. *Neural networks*, 1(4):295–307.
- Karpathy, A. (2015). Convolutional neural networks for visual recognition lecture notes. <http://cs231n.github.io/neural-networks-3/#ratio>.
- Kingma, D. P. and Ba, J. (2014). Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*.
- Krizhevsky, A. and Hinton, G. (2009). Learning multiple layers of features from tiny images. Technical report, Citeseer.
- Krizhevsky, A., Sutskever, I., and Hinton, G. E. (2012). Imagenet classification with deep convolutional neural networks. In *Advances in neural information processing systems*, pages 1097–1105.
- LeCun, Y., Boser, B., Denker, J. S., Henderson, D., Howard, R. E., Hubbard, W., and Jackel, L. D. (1989). Backpropagation applied to handwritten zip code recognition. *Neural Comput.*, 1(4):541–551.
- Paszke, A., Gross, S., Chintala, S., Chanan, G., Yang, E., DeVito, Z., Lin, Z., Desmaison, A., Antiga, L., and Lerer, A. (2017). Automatic differentiation in pytorch. In *NIPS-W*.
- Pedregosa, F., Varoquaux, G., Gramfort, A., Michel, V., Thirion, B., Grisel, O., Blondel, M., Prettenhofer, P., Weiss, R., Dubourg, V., Vanderplas, J., Passos,

- A., Cournapeau, D., Brucher, M., Perrot, M., and Duchesnay, E. (2011). Scikit-learn: Machine learning in Python. *Journal of Machine Learning Research*, 12:2825–2830.
- Pleiss, G., Chen, D., Huang, G., Li, T., van der Maaten, L., and Weinberger, K. Q. (2017). Memory-efficient implementation of densenets. *arXiv preprint arXiv:1707.06990*.
- Rosenblatt, F. (1958). The perceptron: A probabilistic model for information storage and organization in the brain. *Psychological Review*, pages 65–386.
- Smith, S. L., Kindermans, P., and Le, Q. V. (2017). Don’t decay the learning rate, increase the batch size. *CoRR*, abs/1711.00489.
- Tokui, S., Oono, K., Hido, S., and Clayton, J. (2015). Chainer: a next-generation open source framework for deep learning. In *Proceedings of workshop on machine learning systems (LearningSys) in the twenty-ninth annual conference on neural information processing systems (NIPS)*, volume 5, pages 1–6.
- Werbos, P. (1975). *Beyond Regression: New Tools for Prediction and Analysis in the Behavioral Sciences*. Harvard University.
- Yamada, Y., Iwamura, M., and Kise, K. (2018). Shakedrop regularization. *CoRR*, abs/1802.02375.

ACKNOWLEDGMENTS

I wish to acknowledge the contributions of many people who—directly or indirectly—supported this work.

I thank Justin Shenk for providing the introduction to the Peltarion team and his generosity while traveling with me and hosting me for my visits to Stockholm.

I am also grateful to Anders Arpteg who has been a constant source of support and advice, for his willingness to let me work on my own terms on what I saw fit.

I owe further thanks to Professor Oliver Vornberger who agreed to act as first examiner for this thesis, instead of enjoying retirement, and Dr. Ulf Krumnack for acting as co-examiner.

I thank the entire team at Peltarion AB for creating a welcoming and very entertaining environment for working on this thesis.

Lastly, I wish to acknowledge the countless unnamed developers of the myriads of open source tools that form the bedrock of any productive scientific endeavour. A list of software used during the creation of this thesis can be found in the appendix.

DECLARATION OF AUTHORSHIP

I hereby certify that the work presented here is—to the best of my knowledge and belief—original and the result of my own investigations, except as acknowledged, and has not been submitted, either in part or whole, for a degree at this or any other university.

Rasmus Diederichsen

Osnabrück, November 7, 2018