

Automatic Curriculum Learning for Deep Models Using Active Learning

Ian McWilliam

Master of Science
Artificial Intelligence
School of Informatics
University of Edinburgh
2018

Abstract

Acknowledgements

Many thanks

Declaration

I declare that this thesis was composed by myself, that the work contained herein is my own except where explicitly stated otherwise in the text, and that this work has not been submitted for any other degree or professional qualification except as specified.

(Ian McWilliam)

Table of Contents

1	Introduction	1
1.1	Motivation	1
1.1.1	Supervised Learning	1
1.1.2	Curriculum Learning	1
1.1.3	Active Learning	1
1.2	Contribution	1
1.3	Document Structure	1
2	Background	5
2.1	Supervised Learning	5
2.2	Deep Learning	6
2.2.1	Feedforward Networks	7
2.2.2	Convolutional Networks	7
2.3	Stochastic Gradient Descent	7
2.4	Curriculum Learning	7
2.5	Active Learning	9
3	Related Work	11
3.1	Self Paced Learning	11
3.2	Transfer Learning	11
3.3	Reinforcement Learning	11
3.4	Active Learning	11
4	Bootstrapped Active Curricula	12
4.1	Curriculum Construction	12
4.1.1	Average Absolute Distance to Threshold (AADT)	13
4.1.2	Phase Training Epochs	14
4.1.3	Pseudocode for BAC	15

4.2	Geometric Shapes Dataset	15
4.3	Experiments	16
4.4	Results and Discussion	18
5	Dynamic Active Curricula	22
5.1	Curriculum Construction	23
5.1.1	Dynamic Task Curricula (DTC)	23
5.1.2	Biased Sampling Curricula (BSC)	24
5.1.3	Biased Task Curricula (BTC)	25
5.1.4	BALD Active Score Function	25
5.1.5	Softmax temperature	26
5.2	MNIST	27
5.3	CIFAR 10	27
5.4	Experiments	27
5.5	Results and Discussion	27
6	Conclusion and Further Work	29
	Bibliography	30

Chapter 1

Introduction

1.1 Motivation

1.1.1 Supervised Learning

1.1.2 Curriculum Learning

1.1.3 Active Learning

1.2 Contribution

1.3 Document Structure

The subsequent chapters of this thesis will be organised as follows:

- Background - Here we will introduce in more detail the topics introduced above, giving the reader the background required to understand and appreciate the rest of work. We will also develop some of the nomenclature that will be used in various other sections.
- Related Work - In this chapter we set our contributions in the context of related studies, by discussing various other works which have tested approaches for automating curriculum discovery and improving learning performance with curriculum methods.
- Bootstrapped Active Curricula - In this chapter we introduce the bootstrapped active curriculum approach, explaining the curriculum construction, the methods

used to test it, as well as the results of the experiments and subsequent discussion.

- **Dynamic Active Curricula** - Motivated by the previous chapter, we test methods for dynamically constructing and implementing learning curricula throughout training, again laying out the curriculum construction approach, the methods used to test its efficacy and a discussion of the results of the experiments.
- **Conclusion and Further Work** - Finally, we summarise the findings of the thesis and suggest directions in which future work can build on the experiments shown in this paper.

In the next section we will introduce in more detail the topics introduced above, giving the reader the background required to understand and appreciate the rest of work. We will then set our contribution in the context of related studies that have also tested methods for automating curriculum discovery and improving learning performance with curriculum methods. The

Supervised learning is the area of machine learning in which algorithms learn the relationship between a set of input features and corresponding ‘ground truth’ labels, the ultimate goal being to construct a predictive model of the relationships between the inputs and the labels in order to predict the labels of future, unseen input samples. Deep learning models perform this task by building hierarchical representations of the input features throughout a multitude of layers, often using feature maps such as convolutions or recurrent layers to construct complex representations of the inputs. When training a deep model a standard methodology is *gradient descent*, which calculates the gradient of a chosen error function with respect to the free parameters of the model so as to tune the parameters in a way that will minimise this error function on the training set of input-label pairs. A popular variant of the gradient descent algorithm is *mini-batch stochastic gradient descent*, which uniformly samples mini-batches of a preset size from the available training data, performing a gradient descent update on each batch until the all training samples have been selected, then repeating until the network converges to a solution. Sampling uniformly from the training data ensures that the mini-batch gradient is an unbiased estimation of the gradient over the whole training set, however the estimation can exhibit high variance. In this paper we analyse approaches for augmenting mini-batch stochastic gradient descent (SGD), using methods inspired by two areas of study; *active learning* and *curriculum learning*.

Active learning is generally used when there is a prohibitive cost to obtaining labels for supervised learning; in such cases it is desirable to know which samples will lead

to the greatest best improvement in algorithm performance, selected from a set of unlabeled candidate samples. As such, there is a rich literature in active learning detailing how to choose the most informative samples, in particular using *acquisition functions* to select which sample(s) to label and use for training. While active learning is usually employed to reduce labeling costs and speed up learning, curriculum learning explores the hypothesis that the overall accuracy of the network can be improved by presenting the training data to the algorithm in a meaningful order. Inspired by the way in which humans and animals learn, REF BENGIO suggest learning can be improved by emphasising easier concepts earlier on in training before introducing difficult samples, or by emphasising more difficult training samples later in training. In their paper REF BENGIO for example, the authors use a the ‘Geometric Shapes’ dataset, consisting of images of geometric shapes of different complexities, to show that by initially training on ‘easy’, regular shapes, test classification accuracy is improved.

The substantial challenge with curriculum learning however is that in many domains however it is challenging to identify a clear delineation between ‘easy’ and ‘hard’ samples through which to implement curriculum learning; in this paper we propose that the methodologies developed for the active learning approach are well suited to estimating the difficulty of training samples, allowing the automatic construction of learning curricula that will improve the training of deep networks on a wide range of tasks. Specifically, the approach set out in this paper modifies the SGD algorithm by, instead of sampling with uniform probability, sampling training examples proportionally to some measure of ‘difficulty’, as derived from an active learning style acquisition function metric. We test our methods on three image classification datasets; MNIST, CIFAR 10 and the GeoShapes dataset (a geometric shapes classification dataset with an established curriculum baseline), exploring a range of active learning metrics as well as several curriculum construction methods. Our results show consistent performance against a uniform sampling baseline, with significant reductions in test set error, robust to different network architectures, datasets and curriculum methodologies. The output of this work is a set of flexible methods for improving deep models in a wide range of tasks, as well as an investigation of how using the difficulty and uncertainty of training samples affect learning performance.

In the next section we will introduce in more detail active and curriculum learning, exploring the link between the two approaches and the sometimes contradictory hypotheses they pose. We will then discuss related work where the authors implement similar methods for improving algorithm performance through biasing learning

towards certain training samples throughout training. We will then lay out the experimental methodologies and datasets used in the paper before presenting and analysing the results of the tests and concluding with a discussion and suggestions for further work.

Chapter 2

Background

2.1 Supervised Learning

Machine learning is the field of study concerned with building algorithmic systems that can automatically infer patterns and relationships in data. While machine learning has several main subfields, for example reinforcement learning REF? and unsupervised learning REF?, the most commonly studied area of the field is arguably *supervised learning*. Supervised learning is characterised by learning a function that maps inputs to outputs (or ‘labels’), using a *training set* of example input-output pairs, (supervised learning has previously been referred to as “learning from examples” REF). The training set, which we will denote by \mathcal{T} in this report, consists of a number of inputs-output pairs: $\mathcal{T} = \{(x_1, y_1), (x_2, y_2), \dots, (x_N, y_N)\}$, where N is the number of examples in the training set. We will refer to an input-output pair from the training set as a ‘sample’ or an ‘example’ interchangeably. Concretely, the goal of supervised learning is to find a function $f : X \rightarrow Y$, where X denotes the input space and Y the output space, so as to minimize the *generalization error* of the function. The generalization error is defined as the expected error of the function averaged over future input-output samples REF MURPHY, where the error is defined using a loss function $L : Y \times Y \rightarrow \mathbb{R}$, such that $L(f(x_i), y_i)$ gives the error of the function on the input-output pair $\{x_i, y_i\}$. Generalization error can therefore be defined as $\mathbb{E}_{X \times Y}(L(f(x), y))$, however, as finding a closed form solution for the generalization error is generally intractable, it is usually approximated using the empirical error on a held out *test set* of samples that were not used when fitting the function. Using a test set of M input-output examples, we therefore approximate the generalization error as $\frac{1}{M} \sum_{i=1}^M L(f(x_i), y(i))$. A supervised learning algorithm is a method for fitting the function f using the training set \mathcal{T} , usually by

minimizing the loss of the function over the input-output samples in the training set, the process of fitting the function to the training data is referred to as *training*. In order to avoid *overfitting*, where the function has a low error on the training set but a high generalization error¹, the loss function often includes regularization terms that penalise function complexity, or the range of functions that can be fitted by the algorithm is constrained. Furthermore, in addition to a training set and test set, a *validation set* of samples is often used to monitor the function error on samples outwith the training set during training, in order to measure whether or not the function is overfitting.

2.2 Deep Learning

Deep learning, as applied in the supervised learning setting, refers to algorithms which model the relationships in the training set using complex, hierarchical representations of the data, usually doing so using multiple ‘deep’ layers, as well as feature maps such as convolutions or recurrent layers REFS. Deep learning models fit functions by passing the input signals through several layers of transformations, multiplying the signal by parameterised weights and then passing the weighted signal through non-linear *activation functions*, the combined model of weights and transformations is referred to as a ‘deep network’, or ‘deep neural network’, as a result of similarities to the functionalities of neurons in brain. This approach allows the algorithm to transform the input space in order to make the closely match the output, with the weight parameters being tuned throughout training in order

In image classification in particular, the state of the art is *convolutional neural networks*, which assigns weights to areas of a predefined size of the input space, as opposed to applying weights to each input node. Doing this allows for the automatic construction of abstract features which has proven extremely effective for analysing images, for example by modeling specific shapes characteristic of certain labels.

¹More specifically, overfitting refers to the situation where the function has a low error on the training set and a higher generalization error than a similar function with a higher error on the training set.

2.2.1 Feedforward Networks

2.2.2 Convolutional Networks

2.3 Stochastic Gradient Descent

The standard method for training deep models is *gradient descent (GD)*, an optimization algorithm which varies the parameters in the model depending on the gradient of a chosen error function. To implement GD it is necessary to calculate the gradient of the error function with respect to the parameters of the model, usually this is done layer by layer, starting with the output layer, in a method referred to as *backpropagation*. Calculating this gradient however can be very computationally expensive, particularly when dealing with large training sets; to address this issue a variant of GD, *stochastic gradient descent (SGD)*, is often used. With SGD, instead of calculating the error gradient over the entire training set, only one sample is used to calculate the gradient and update the model parameters. Alternatively a selection or ‘mini-batch’ of training samples may be used to calculate the gradient, in which case the optimization algorithm is referred to as *mini-batch stochastic gradient descent*. It can be shown that that SGD and mini-batch SGD produce an unbiased estimate of the error gradient, with various convergence proofs showing that SGD will eventually converge to an optimal solution. A disadvantage to SGD however is that the gradient estimate, while unbiased, can exhibit high variance, potentially resulting in extremely slow converge times. There are many adaptations to ‘vanilla’ SGD, for example momentum based methods and other more sophisticated optimization algorithms which build on SGD to result in better and quicker fitting of the model parameters.

2.4 Curriculum Learning

While active learning uses methods to identify which samples to label and train in order to speed up training in domains with a high labeling cost, *curriculum learning* attempts to present training samples to the learner in a meaningful order that will lead to greater overall generalization performance of the model. The motivation stems from the way in which humans and other animals learn, usually beginning with easy concepts before moving onto more complex facets of the area of study. The same principle can be applied to training deep models, and the authors of REF suggest that, by initially training only on ‘easy’ samples, one can reduce overall generalization error. The authors

offer several theoretical justifications, for example comparing curriculum learning to *continuation methods* REF; it is proposed that the easier samples represent a smoother, more convex version of the error space of the overall problem, and that, by training on easier samples, the parameters of the model are effectively initialized into an area of parameter space closer to the global optimum. This argument is similar to that of unsupervised pre-training, which again has been shown to lead to better generalized models by initializing the parameters into parts of the error space closer to the global optimum. Comparisons have also been drawn between curriculum learning and *transfer learning*, with the easier samples being seen as a separate task that the model is trained on, before using the weights for a different task (i.e. the harder samples) as in transfer learning.

The example given in REF BENGIO for curriculum learning is the ‘GeoShapes’ dataset, an image classification where a network attempts to classify whether or not an image shows a rectangle, ellipse or triangle. In this case there is a natural subset of ‘easy’ samples; specifically squares (i.e. regular rectangle), circles (regular ellipses) and equilateral triangles. The authors show that, by training initially on only the regular shapes, then transitioning to training on harder shapes, the test set performance is significantly improved compared to training simply on the harder shapes for the entirety of training. One issue with this study is that it can be argued that the curriculum trained model has seen more samples overall than the baseline, as the curriculum model is trained on both an ‘easy’ training set and a ‘hard’ training set, whereas the baseline is trained only on the hard training set. A better baseline therefore is a model trained uniformly on the union of the easy and hard training sets. While the authors do comment on this issue, and claim that the curriculum method still outperform uniform sampling from the combined training set, the results we will set out in this paper did not reach the same conclusions.

A key difficulty in implementing curriculum learning is that it is often very difficult to delineate between ‘easy’ and ‘difficult’ samples, while it is also hard to ascertain how one should transition from different difficulties. A key issue therefore is that of exploring methods for automating the construction of learning curricula, and it is towards this goal that this paper contributes; specifically investigating how active learning methods can aid such curriculum construction. Having introduced the reader to active and curriculum learning, the next section will lay out a variety of related work wherein the authors attempt to automate the process of curriculum construction or apply active learning methods with the goal of improving network performance.

2.5 Active Learning

A key component in any supervised learning effort is labeled data; in many domains it is relatively easy and cheap to obtain large volumes training samples, however in others it can be far more costly, particularly acquiring accurate labels. In medical image analysis for example one may require a domain expert to spend significant time analysing each image before assigning label, or in document tagging it can take time to read a document and assign a topic label. It can therefore be very useful for a designer to understand which samples they should go to the effort of acquiring, labeling and feeding into their chosen learning algorithm, generally measured by how much the chosen samples improve the network performance, compared to if samples were instead selected randomly. We here introduce some of the main methodologies employed for active learning, giving the reader some background to the methods that will be used in this paper.

There are a variety of approaches to the active learning problem, however most involve the use of an *acquisition function*, which selects which sample, from a set of candidate unlabeled examples, should be selected for labeling and training. As the most appropriate training examples varies depending on the learning algorithm, as well as its current state in the training process, the chosen sample is said to be ‘queried’ by the algorithm. The motivation behind different active learning approaches vary; one of the most common approaches is that of *uncertainty sampling*, wherein the samples that the learning algorithm is most uncertain about labeling are queried. This uncertainty can be captured by analysing the distance to classification threshold of the model outputs; for example one method is to select the sample about which the model is least confident in predicting: (taken from REF SETTLES:)

$$x_{LC}^* = \arg \max_x 1 - P_{\theta}(\hat{y}|x), \quad (2.1)$$

where

$$\hat{y} = \arg \max_y P_{\theta}(y|x). \quad (2.2)$$

Where x_{LC}^* is the queried training sample and $P_{\theta}(y_i|x)$ is the model’s predicted probability that sample x is of class y_i , given model parameters θ . Similarly, samples can be queried by their average distance to classification threshold or, similarly, the entropy of the algorithm prediction, again taken from REF SETTLES:

$$x_H^* = \arg \max_x - \sum_i P_{\theta}(y_i|x) \log P_{\theta}(y_i|x), \quad (2.3)$$

where the sum runs over the possible classes y_i .

An alternative approach to querying training samples is to estimate the expected change in model parameters, if trained on a given sample. As, in the active learning setting, it is assumed that the label is unavailable, this is calculated as an average across all potential labels. Model change can be estimated by the magnitude of the gradient vector produced by training on the tuple $\langle x, y \rangle$. The acquisition function then selects the sample which maximises the expected gradient size (REF SETTLES AGAIN):

$$x_{EGL}^* = \arg \max_x \sum_i P_\theta(y_i|x) \|\nabla \ell_\theta(\mathcal{L} \cup \langle x, y_i \rangle)\|, \quad (2.4)$$

where $\|\cdot\|$ is the Euclidean norm, \mathcal{L} is the current set of labeled training samples, ℓ is the objective function used to train the model and $\nabla \ell_\theta(\mathcal{L})$ is the gradient of this objective function with respect to the model parameters θ , when trained on \mathcal{L} . This approach therefore finds the sample that leads to the largest expected increase in the gradient when added to the training set \mathcal{L} .

Finally, another common approach for active learning is that of *query by committee*; here a population of different models are trained on an initial training set, then the samples about which the models exhibit the most disagreement in their predictions are queried. An example of this is *vote entropy* REF:

$$x_{VE}^* = \arg \max_x - \sum_i \frac{V(y_i)}{C} \log \frac{V(y_i)}{C}, \quad (2.5)$$

where C represent the size of the ‘committee’ (i.e. the number of models) and $V(y_i)$ is the number of models in the committee that predict label y_i . There are obvious parallels here to methods such as ensembling, boosting and bagging, indeed active learning has drawn parallels with several other learning paradigms, such as self-paced learning (REF) and curriculum learning (REF), the latter of which we shall now introduce.

Chapter 3

Related Work

3.1 Self Paced Learning

3.2 Transfer Learning

3.3 Reinforcement Learning

3.4 Active Learning

Chapter 4

Bootstrapped Active Curricula

The first curriculum construction approach we investigate is what we term *bootstrapped active curricula (BAC)*, the name is chosen as it involves ‘bootstrapping’ a curriculum from a separate model. The intuition behind this approach is that, while it may be difficult to ascertain a-priori which samples are ‘hard’ or ‘easy’ (particularly for very large datasets where it is infeasible to analyse every sample), we can automatically infer which samples are difficult by first training a model on the data and then using an active learning uncertainty metric to investigate which samples the model is uncertain about classifying. Using prediction uncertainty to approximate difficulty, we can then construct a learning curriculum which can be used to train a new model, for example by splitting the training samples into separate tasks of increasing difficulty, as detailed below.

4.1 Curriculum Construction

In the BAC approach we train two models; the first, which we term the ‘baseline model’ and denote $\theta_{baseline}$, is trained to convergence using a standard mini-batch gradient descent optimisation on the entire available training set \mathcal{T} . We then use $\theta_{baseline}$ in conjunction with the Absolute Average Distance to Threshold uncertainty function as defined below in Section 4.1.1, scoring each training sample in \mathcal{T} . This produces $\mathbf{S}^{\theta_{baseline}}$, a vector of N scores, where N is the number of samples in \mathcal{T} , such that the i^{th} element of $\mathbf{S}^{\theta_{baseline}}$ is the output of the AADT uncertainty function for the i^{th} training sample inputs:

$$S_i^{\theta_{baseline}} = AADT_{\theta_{baseline}}(x_i) \quad (4.1)$$

where x_i are the inputs of the i^{th} training sample.

We then sort the training samples according to their score, producing an ordered training set $\mathcal{T}_{\mathbf{S}^{\theta_{baseline}}}$, such that the first training sample in $\mathcal{T}_{\mathbf{S}^{\theta_{baseline}}}$ is the sample which produces the highest value from the AADT function. We then split $\mathcal{T}_{\mathbf{S}^{\theta_{baseline}}}$ into equally sized ‘tasks’, with the number of tasks being a hyperparameter of the curriculum construction. The learning curriculum is then constructed from these tasks; the chosen approach is to split training into discrete phases, with the number of phases being equal to the chosen number of tasks. The first phase of the curriculum consists of training only on the first task, denoted $\mathcal{T}_{\mathbf{S}^{\theta_{baseline}}}^1$, in the second phase, the second task is added to the training samples, training the model on $\mathcal{T}_{\mathbf{S}^{\theta_{baseline}}}^1 \cup \mathcal{T}_{\mathbf{S}^{\theta_{baseline}}}^2$. In the third phase (if the number of tasks is greater than 2), the next task is added, and so on until all tasks have been added and the final phase consist of training on the entirety of the original training set \mathcal{T} . Having constructed the BAC learning curriculum, we train a new ‘curriculum’ model’, which we denote by $\theta_{curriculum}$, using the curriculum. Pseudocode for the BAC method is given in Section 4.1.3.

4.1.1 Average Absolute Distance to Threshold (AADT)

A popular active learning uncertainty method is to examine the proximity of the model’s outputs to the classification bounday. The assumption is that samples that the model is uncertain about classifying will produce probabilities close to the classification boundary; indeed as mentioned in Section 2.5 the authors of H-S Chang (2017) show that prediction variance is inversely proportional to the distance to the boundary. From a curriculum perspective we can estimate a sample’s difficulty by the algorithm’s uncertainty in predicting the class label, with uncertain samples being seen as hard, and vice versa. We thus define the uncertainty function $AADT$ which calculates the absolute distance to classification threshold for the model outputs, averaged across all possible classes. Note that the function is dependent on a the output of the model in question, which is denoted θ .

$$AADT_{\theta}(x_i) = \frac{\sum_{c=1}^C |P_{\theta}(y_c|x_i) - \frac{1}{C}|}{C}. \quad (4.2)$$

Where $|\cdot|$ represents the L1 norm/absolute value function. Here N is the number of training samples, C is the number of output classes and $P_{\theta}(y_c|x_i)$ is the output softmax probability for class y_c of the model θ , given input x_i . We also tested the average *square* of the distance to threshold, as opposed to the absolute distance to threshold, as well as testing the entropy of the outputs as an uncertainty measure, however results were

extremely similar in all cases.

4.1.2 Phase Training Epochs

A naive approach to the BAC approach would be to train the curriculum model for the same number of epochs as the baseline model, split equally acrossing training phases. For example if the baseline model was trained for 100 epochs and we then constructed a two task curriculum, the curriculum model would be trained on the first task for the first 50 epochs, then on both tasks for the second 50 epochs. The issue with this method however is that, as during the first 50 epochs there are only half as many training samples, the curriculum model will not have as many parameter updates as the baseline model. To emphasize this, consider that if we were using stochastic gradient descent (i.e. a mini-batch size of 1) with a full training set of 1000 samples, the baseline model in this instance would have (# epochs * # training samples) parameter updates, i.e. $100 * 1000 = 100,000$ parameter updates. The curriculum model on the other hand would have $(50*500)$ parameter updates the first training phase and $(50*1000)$ in the second, resulting in a total of $(50*500) + (50*1000) = 75,000$ parameter updates, 25% less updates than the baseline model. To address this, we increase the number of epochs in each phase by the ratio of the number of samples used in the phase to the size of the whole training set. Specifically, we set the number of training epochs in each phase as follows:

$$NumEpochs^t = \lfloor \frac{BaselineEpochs}{t} \rfloor \quad (4.3)$$

Where *BaselineEpochs* is the number of epochs used to train the baseline model and $NumEpochs^t$ denotes the number of training epochs in the t^{th} training phase. For example, in the first phase, $NumEpochs^1 = \frac{BaselineEpochs}{1} = BaselineEpochs$. $\lfloor . \rfloor$ represents the floor function; as $\frac{BaselineEpochs}{i}$ will not always be integer, we round down the number of epochs in each phase. The curriculum model will therefore be trained for a higher number of epochs (precisely, $\sum_{t=1}^{NumTasks} \frac{1}{t}$ times more epochs), however the number of parameter updates will be equalised.

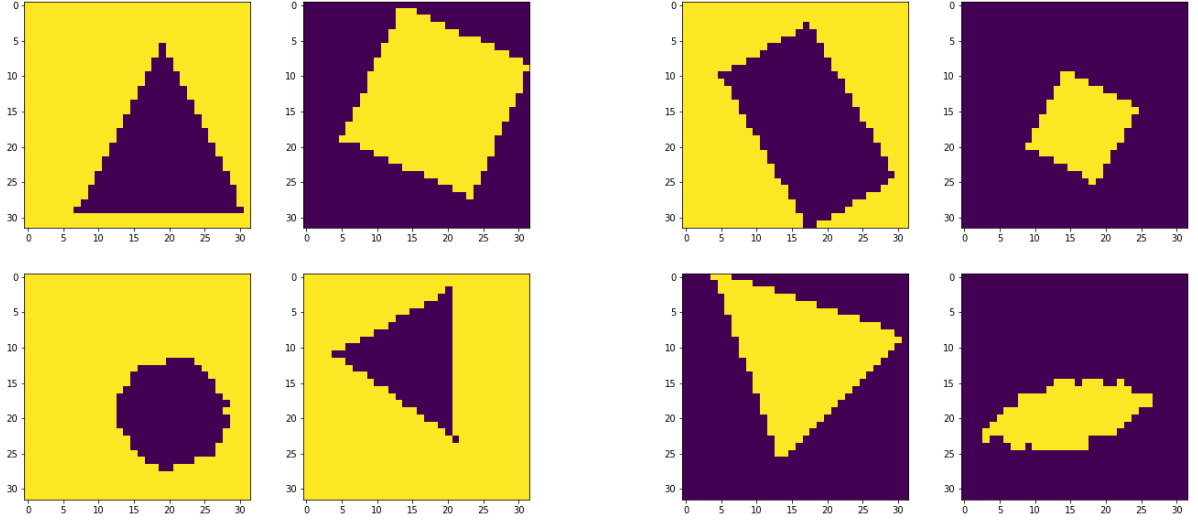
4.1.3 Pseudocode for BAC

4.2 Geometric Shapes Dataset

To test the bootstrapped active curriculum approach we use the ‘Geometric Shapes’ dataset, as used in Y.Bengio and J.Weston (2009). This dataset consists of 32x32 pixel images of geometric shapes, specifically ellipses, rectangles and triangles; the class labels are one-hot vectors which indicate to which of the three classes each sample belongs. As well as varying the shapes shown in the image, the samples also vary in colour, orientation, size and position. This dataset is used in Y.Bengio and J.Weston (2009) as it is easy to construct a predefined curriculum for geometric shapes; circles, squares and equilateral triangles represent regular, ‘easy’ versions of the broader classes of ellipses, rectangles and triangles, respectively. The authors of Y.Bengio and J.Weston (2009) train a curriculum model by first training only on an easy training set consisting only of circles, squares and equilateral triangles, before then training on a hard training set with the more general shapes. The easy and difficult training sets consist of 10,000 training samples, giving a total of 20,000 training samples, while the test set, consisting only of hard samples, contains 5,000 images. The authors demonstrate that their approach outperforms an identical model trained only on the harder shapes; the curriculum model consistently achieves greater test accuracy, with the greatest improvement coming when the first half of the training epochs train on the easier shapes, and the second half the harder ones, as shown in Figure 3 of their paper. As discussed by the authors, one potential pitfall of their experiments is that the curriculum model has seen more samples than the benchmark model, as it has been trained on both the easy and the difficult training sets, to avoid this in our experiments, the training set we use is the union of both the easy and difficult samples. Our training set therefore consists of 20,000 geometric shape images, including both the easy, regular shapes and the more difficult shapes, the test set is unchanged however, consisting of the 5,000 difficult samples.

Some example images for the Geometric Shapes dataset are given in Figure 4.1 below, with Figure 4.2a showing examples from the easy training set featuring circles, squares and equilateral triangles, while Figure 4.2b shows samples from the more difficult training set of general shapes.

In the results Section 4.4, we analyse how successful our tested curriculum method is at automatically identifying which samples come from the easy or hard training



(a) Samples from ‘easy’ training set

(b) Samples from ‘hard’ training set

Figure 4.1: Sample figures from the Geometric Shapes dataset

sets, investigating which training samples fall into the different tasks automatically constructed through the BAC curriculum method.

4.3 Experiments

In order to measure the effect of the curriculum on test performance, we set $\theta_{baseline}$ and $\theta_{curriculum}$ to have identical architectures, hyperparameters and initial weights, essentially minimizing any differences in training besides the learning curriculum. We use both the Average Absolute Distance to Threshold (AADT) function, detailed in Section 4.1.1, to score the training samples with the trained baseline model in each experiment. Furthermore, as well as testing ‘easy to hard’ curricula, where the training phases progress from $\mathcal{T}_{\theta_{baseline}}^1$ to the final task, we also test the opposite approach, with the first training phase using only the hardest task, then incorporating the other, easier tasks throughout training. We run the experiments using the Geometric Shapes dataset, as introduced in Section 4.2, allowing us to compare results with the predefined curriculum as in Y.Bengio and J.Weston (2009). To do so, as well as training the baseline and BAC curriculum models, we also train a model using the predefined curriculum method laid out in Y.Bengio and J.Weston (2009), specifically training on

only the easy samples in for the first half of the training epochs, and only the hard samples in the second half. In order to better compare the predefined curriculum method from Y.Bengio and J.Weston (2009), we also run an adjusted version of their method; unlike the BAC approach, where we correct the difference in number of parameter updates between the curriculum model and the baseline model, the predefined curriculum model will end up undergoing significantly less parameter updates than our baseline (in fact it will have half as many updates). We therefore correct for this in two ways; in the first phase of training, in which the model is trained only on the easy images, lasts for twice as many epochs as in the unadjusted method (correcting for the fact that it consists of half as many samples as the full training set). The second phase of the adjusted model is then trained on the *full* training set, consisting of the union of the hard and easy samples. This should allow for a better comparison between the BAC models, the baseline model, and the predefined curriculum approach. We therefore train 5 separate models:

- Baseline Model - Trained on the full training set of easy and hard geometric shape images for all epochs, with standard mini-batch stochastic gradient descent.
- Easy to Hard Model - Trained with a bootstrapped active curriculum construct from the above Baseline Model, beginning with the easiest/least uncertain task and including the other tasks throughout training.
- Hard to Easy Model - As above, but training begins with the hardest/most uncertain tasks before incorporating the easier tasks throughout training.
- Predefined Curriculum - As in Y.Bengio and J.Weston (2009), the first half of the training epochs use only the 10,000 easy training samples, while the second half train on only the 10,000 hard training samples.
- Adjusted Predefined Curriculum - As above, but the number of epochs in the first phase of training is doubled to account for the difference in parameter updates, and the second phase uses both easy and hard training samples for better comparison with the BAC models.

All models are tested on the same test set of 5000 images, all containing ‘hard’ geometric shapes; experiments are repeated multiple times with different weight initialisations. To analyse the effect of the curriculum on learning we report the accuracy

and cross-entropy error of the different models over the held out test set. The exact model architecture is given in tables 5.1 - 5.3 below ('FC' = fully connected Layer):

Table 4.1: Geometric Shapes Dataset Model Architecture

Geometric Shapes Dataset Model Architecture							
	Layer 1	Layer 2	Layer 3	Layer 4	Layer 5	Layer 6	Layer 7
Layer Type	FC	Dropout	FC	Dropout	FC	Dropout	FC
Units	300	NA	300	NA	300	NA	3
Activation	Tanh ¹	NA	Tanh	NA	Tanh	NA	Softmax

We also lay out the hyperparameters for the training procedures in table 4.2 below:

Table 4.2: Experiment Hyperparameters

Experiment Hyperparameters						
	Epochs	Optimiser	Learning Rate	Dropout %	Batch Size	Num Tasks
GeoShapes	350	ADAM	0.0001	0.25	32	2

Architectures and hyperparameters were chosen and tuned in order to deliver a good level of performance without prohibitive training times, robustness tests were however carried out varying the model architectures and hyperparameters resulting in little change in the relative performance of the different methods.

4.4 Results and Discussion

The results of running the bootstrapped active curriculum experiments with the Geometric Shapes dataset are summarised in Table 4.3 and Figure 4.2 which shows the test set accuracies and cross entropy errors, including standard errors, derived from 24 experiments with different initialisations.

We see from the results that the Easy to Hard curriculum model significantly outperforms the baseline model, with an average test accuracy improvement of around 2.5%, over 5 standard errors, higher than the baseline test performance. Cross-entropy error is also reduced, with the Easy to Hard curriculum test error over 6 standard errors below that of the baseline model. These results suggest that, not only is there a benefit to training the model with a curriculum, the BAC model seems to be successful at

¹The Hyperbolic Tangent activation is chosen as this was this the activation used in Y.Bengio and J.Weston (2009).

Table 4.3: Geometric Shapes BAC Results

GeoShapes BAC Results		
	Test Accuracy (%)	Test Cross-Entropy Error
Baseline	0.825 ± 0.00269	0.440 ± 0.00614
Easy to Hard Curriculum	0.840 ± 0.00267	0.398 ± 0.00670
Hard to Easy Curriculum	0.815 ± 0.00278	0.467 ± 0.00699
Adjusted Predefined Curriculum	0.800 ± 0.00143	0.486 ± 0.00409
Predefined Curriculum	0.757 ± 0.00279	0.576 ± 0.00489

identifying which samples should be used in the different curriculum learning phases. Conversely, the other curriculum methods all underperform the baseline model; the Hard to Easy curriculum in some ways acts as a control, showing that the improvement in the Easy to Hard model is driven by the order in which the tasks are included in training, as opposed to being a consequence of another part of the BAC method. It is interesting to note that both of the predefined curriculum models underperform the baseline model; the lower accuracy of the unadjusted method would seem to support the authors' suggestion in Y.Bengio and J.Weston (2009) that the performance differences in their experiments may have been a result of their curriculum model having seen overall more training samples than their baseline. However, even our adjusted version of the predefined curriculum results in a significantly lower test accuracy than the baseline. This is interesting as this method is effectively trained in the same way as the BAC approach, however the tasks are predetermined by an intuitive sense of difficulty, as opposed to being derived from the baseline. This would suggest that, even for datasets where there is a somewhat clear distinction between hard and easy samples, a bootstrapped curriculum may be more useful. Potentially, this may be because the BAC method specifically calculates which samples the model itself is uncertain in classifying, as opposed to assuming that our intuitive sense of difficulty corresponds to the best curriculum for the model. We would note however, that given the relatively low resolution of the images (as illustrated in Figure 4.1), there is perhaps less observable difference between the easy and hard training sets than might be assumed, and that this may contribute to the poor performance of the predefined curriculum.

With that in mind, we can analyse one of the the Easy to Hard curriculum experiments in more depth, in order to ascertain which samples are included in the two tasks. Recall that in the BAC method we train a baseline model, then use the uncertainty

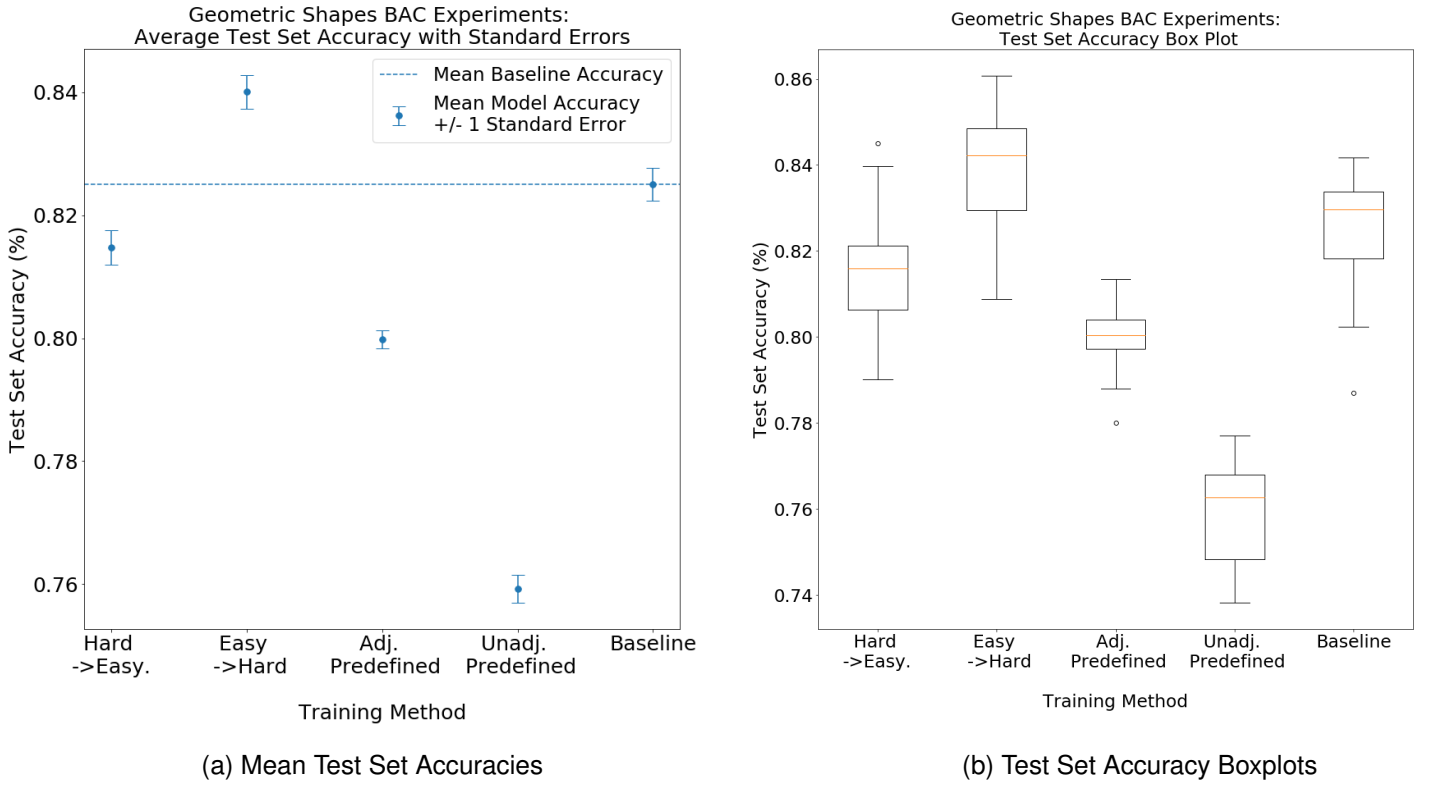


Figure 4.2: Geometric Shapes Bootstrapped Active Curriculum (BAC) Experiment Results

in the model’s output class probabilities for the training samples, calculated using the AADT function defined in Section 4.1.1, to score and rank the training samples. The training samples are then split into two tasks, the first consisting of the training samples that the model is least uncertain in classifying, and the second task containing the samples about which the model’s outputs are most uncertain. We can analyse the composition of the two tasks to infer which samples the model is least/most uncertain about. Doing so for one of the experiments presents some interesting conclusions; first of all, there is a significant lack of balance in the target classes in the two tasks. In task 1, the easier task, 43% of the images are of triangles, while 31% are of ellipses and 26% are of rectangles, implying that the baseline model is significantly more confident in classifying whether or not an image shows a triangle than it is in classifying the other shapes. The bottom right image in Figure 4.2b illustrates why this may be, with an example of an ellipse that looks quite similar to a mis-shapen rectangle, again this is probably the result of the low resolution of the images. We also observe that 64% of the samples in the first task are from the easy training set of squares, circles and equi-

lateral triangles, showing that the AADT function is somewhat succesful in separating the two training sets, however the first task still contains a substantial number of the harder samples. In appendix REF APPENDIX, we show a number of samples from either task, illustrating how the first, easier task is predominantly made up of regular shapes, as well as triangle, which appear to be easier to classify than rectangles and ellipses.

A potential enhancement to the BAC method would be to control the balance of the classes in either task, however if one class is indeed easier than the others to classify then it may be better for it to be over-represented in the easier task(s). While we used two tasks for the Geometric Shapes dataset as this corresponded to the number of tasks of the predefined curriculum, we also ran some initial experiments with a larger number of tasks. We found some variation in results with a larger number of tasks, with some choices of task number underperforming the baseline model, however further tests could be carried out in future work to ascertain how robust these results are and investigate what is driving the different performances.

4.5 Summary

To summarize, these experiments seem to support the hypothesis that training a deep model with a curriculum can reduce generalization error, and furthermore that using an active learning approach to inferring sample difficulty through model prediction uncertainty can be an effective way of automatically constructing such curricula. In particular such curriculum construction methods may even outperform a preconstructed curriculum using an intuitive sense of sample difficulty, as model uncertainty incorporates information about which samples the model itself finds difficult. One of the arguable drawbacks to the BAC method is that it requires that a baseline model is trained to convergence, in order to use its outputs to construct the curriculum to train the curriculum model. This effectively doubles the training time and, while this may not be too significant a computational burden in some problems, motivates the approach of the next chapter, which investigates whether or not a curriculum can be constructed dynamically throughout training, without first training a baseline model.

Chapter 5

Dynamic Active Curricula

The results laid out in Chapter 4 suggest that generalisation performance can indeed be improved through the use of a learning curriculum, and that using active learning approaches to score which training samples are hard or difficult can be an effective way of automatically constructing such curricula without manually analysing the training samples. The disadvantage with the BAC method from Chapter 4 however is that it is necessary to first train a ‘baseline’ model that can be used to derive a curriculum based on which samples it is uncertain about classifying. While in some cases this may not too significant a computational burden, it does effectively double the overall training time, motivating the approach laid out in this section. Specifically, we wish to investigate methods that can dynamically construct a curriculum throughout training, without needed to reference another model. We do this by calculating the model uncertainty on the training samples throughout training, using the uncertainty scores to construct dynamic curricula which evolve throughout training to focus on the samples that the model is more, or less, uncertain in its predictions. If succesful, these methods should lead to model performance which beats a the benchmark test set performance of a baseline model trained with normal mini-batch stochastic gradient descent on the entire training set, and hopefully will achieve results similar to those achieved by the bootstrapped active curricula from Chapter 4.

5.1 Curriculum Construction

5.1.1 Dynamic Task Curricula (DTC)

The first dynamic curriculum method we test we term *dynamic task curricula*; this approach is very similar to that of the BAC method laid out in Chapter 4, however instead of constructing tasks based on the uncertainty scores of a fully trained baseline model, tasks are constructed dynamically using the uncertainty of the curriculum model itself throughout training. To do this, we calculate the model's uncertainty in predicting the classes of the training samples in the full training set \mathcal{T} using the AADT function 4.1.1 (or another active learning uncertainty measure) at the end of every epoch, then ranking the training sample by the the model's uncertainty. We then split the training data into separate tasks, with the first task consisting of the first $\frac{1}{T}$ samples where T is the chosen number of tasks. At the start of each epoch, we will obtain a new set of ordered tasks, $\mathcal{T}_i^1, \mathcal{T}_i^2, \dots, \mathcal{T}_i^T$, where i is the current epoch, and T is the chosen number of tasks. Training is divided into T phases, such that, during epoch j of phase t , the model is trained on a training set consisting $\mathcal{T}_j^1 \cup \mathcal{T}_j^2 \cup \dots \cup \mathcal{T}_j^t$, where t denotes the current training phase. The final training phase will therefore consist of training on the full training set, as $\mathcal{T} = \mathcal{T}_i^1 \cup \mathcal{T}_i^2 \cup \dots \cup \mathcal{T}_i^T, \forall i$. As in the BAC curriculum method 4, we normalise the number of parameter updates in each phase by scaling the number of epochs relative to the number of samples in the training set for that phase. Similarly in Equation 4.3, the number of epochs in training phase t is set as follows:

$$NumEpochs^t = \lfloor \frac{TotalEpochs}{t} \rfloor \quad (5.1)$$

where *TotalEpochs* is the number of epochs the model would need to be trained for on the whole training set in order to achieve the same number of parameter updates as the DTC curriculum method.

Note that we can rank the samples from low to high uncertainty or from high to low uncertainty; in the low to high case the first tasks will contain samples about which the model is confident in predicting, corresponding to dynamically constructing an easy to hard curriculum. Alternatively, ranking the samples from high to low uncertainty will produce a curriculum more similar to an active learning approach, focussing on training the model on areas of the input space that it is more uncertain about classifying. While the latter method may contradict some of the underlying motivations for curriculum learning, for example that learning with an easy to hard curriculum allows the model to explore a smoother version of the error space and better initialise the param-

eters, we would still argue that the hard to easy approach still qualifies as a curriculum method, and may be appropriate in instances where the model already has achieved a high level of expertise in the problem, and improvements are therefore more likely to be made by studying harder samples than easy ones. We therefore have two variations of the Dynamic Task Curriculum method; Easy to Hard DTC and Hard to Easy DTC. This approach is similar to that of *Self-Paced Learning* REF KOLLER ET AL, where they apply a similar approach with a latent SSVM model REF, however they implement their method by introducing a regularization term reflecting the difficulty of each sample into the objective function of the model, as opposed to employing an actual curriculum. Pseudocode for the DTC method is given in section 5.1.1.1.

5.1.1.1 Pseudocode for the DTC curriculum

5.1.2 Biased Sampling Curricula (BSC)

The second dynamic curriculum construction method we test is *Biased Sampling Curricula (BSC)*; one of the potential problems with some curriculum and active learning methods is that of *diversity* REF PAPER ON DIVERSITY. Diversity refers to how well the whole training set is represented in the samples used to train the model; if a sample selection method is extremely undiverse then it may end up selecting only a very small selection of the training samples, leading to the model not training on a suitably broad set of training samples, leading to high generalization error. One way to approach this problem is to avoid deterministic methods of selecting training examples, and instead sample from the training samples using some probability distribution. Mini-batch stochastic gradient descent REF? is one such method, sampling batches of examples from the training data using a uniform probability distribution, however by varying the sampling distribution away from a uniform distribution we can bias training towards different samples, resulting in them being selected more or less often than others. We use this approach in the BSC curriculum method by biasing the sampling probability proportionally to the model uncertainty in predicting the labels of the sample. In particular, at the start of every epoch we use the AADT uncertainty function 4.1.1 to infer the model's current uncertainty in predicting the labels of the training data, giving $\mathbf{S}^{\theta_{baseline}}$, a vector of N scores as in Equation ??, where the N is the number of samples in the full training set \mathcal{T} and the i^{th} element of \mathbf{S} corresponds to the output of the AADT function for the i^{th} training input. We then pass the vector of scores through a softmax function, effectively transforming the scores into probabilities that can be

used as the sampling probability function:

$$\tilde{p}(x_i|\theta) = \frac{\exp(\frac{AADT_{\theta}(x_i)}{\tau})}{\sum_j^N \exp(\frac{AADT_{\theta}(x_j)}{\tau})} \quad (5.2)$$

where θ is the model being trained with the curriculum, $\tilde{p}(x_i|\theta)$ is the sampling probability of training sample i and τ is the softmax temperature parameter that can be set as a hyperparameter or tuned each epoch to achieve a target diversity level in the sampling probability.

5.1.3 Biased Task Curricula (BTC)

5.1.4 BALD Active Score Function

Recent advances in Bayesian neural networks and variational inference motivate an alternative approach to measuring uncertainty; while the distance to classification threshold may encapsulate classification uncertainty for samples close to the boundary, it does not consider the uncertainty associated with analysing samples from parts of the feature space that is not represented in the training data. Consider the toy example shown in FIGURE, where a sample may be far from the classification boundary, but so dissimilar from the training samples that we would want to measure the sample as having high classification uncertainty. One approach to this problem is given by REF GAL, where they motivate using the *Monte Carlo dropout* method as a way of approximating variational inference in neural networks. As with the usual dropout procedure, weights are randomly set to zero throughout the training phase, however, unlike the usual approach, dropout is maintained at the test stage, and a number of forward passes are carried out, resulting in a distribution of outputs. The resultant distribution can subsequently be analysed to infer which test samples the model is more or less confident in predicting, for example by comparing the variance of the output distributions. In REF GAL ACTIVE LEARNING, the authors use the MC dropout method to construct an active learning acquisition function *Bayesian Active Learning by Disagreement (BALD)*, which queries points which “maximise the mutual information between predictions and model posterior”, identifying samples that have a high probability of being placed into different classes in the different stochastic forward passes. One interpretation of the BALD method is that it is similar to the ‘Query by Committee’ active learning methods, with the different forward passes representing different models’ votes.

We calculate the BALD active score function as follows:

$$P_{\theta}^{BALD}(x_i) = \frac{\exp(\frac{S_{\theta}^{BALD}(x_i)}{\tau})}{\sum_j^N \exp(\frac{S_{\theta}^{BALD}(x_j)}{\tau})}, \quad (5.3)$$

where

$$S_{\theta}^{BALD}(x_i) = -\sum_c^C \bar{P}_{\theta}(y_c|x_i) \log(\bar{P}_{\theta}(y_c|x_i)) + \frac{1}{M} \sum_m^M (\sum_c^C P_{\theta}^m(y_c|x) \log(P_{\theta}^m(y_c|x))), \quad (5.4)$$

and

$$\bar{P}_{\theta}(y_c|x_i) = \frac{\sum_m^M P_{\theta}^m(y_c|x)}{M}. \quad (5.5)$$

Here M is the number of stochastic forward passes carried out and $P_{\theta}^m(y_c|x)$ is the softmax probability of class c from the m^{th} forward pass. The score can therefore be interpreted as the difference between the entropy of the average softmax output and the average entropy of the output of each forward pass.

5.1.5 Softmax temperature

In order to homogenize the outputs of the different active score functions, we pass the the scores through a softmax functions, resulting in an output of softmax probabilities summing to 1. Using the softmax function also allows us to use the softmax temperature in order to control the diversity of the sampling probabilities. A common issue with active learning is that the acquisition functions can end up sampling from an unrepresentative subset of the input space, resulting in significant bias in the training of the model (REF!). Indeed, GIVE EXAMPLE OF MAX RATIO FOR DIST2THRESH

We control this effect by using the softmax temperature to target a preset *maximum probability ratio*, defined as follows:

$$MaxRatio = \frac{\max_i P_{\theta}(x_i)}{\min_i P_{\theta}(x_i)}. \quad (5.6)$$

To do this we begin with a softmax temperature of 1, then calculate the current maximum probability ratio and increment the temperature until the target ratio is achieved.

Pseudocode is given below:

PSUEDOCODE FOR SOFTMAX CONTROL

The downside to this method is it could be skewed by outlier probabilities; if there were one sample with an extremely large sampling probability this method would still achieve a target maximum probability ratio, without achieving sufficient diversity in

the sampling probabilities. We investigated using winsorization (REF) as an outlier removal technique prior to tuning the temperature parameter, however, as the active score functions we use generally did not result in outliers, this did not affect results. We show in section REF how controlling the maximum probability ratio affects training.

5.2 MNIST

5.3 CIFAR 10

5.4 Experiments

Table 5.1: Geometric Shapes Dataset Model Architecture

Geometric Shapes Dataset Model Architecture							
	Layer 1	Layer 2	Layer 3	Layer 4	Layer 5	Layer 6	Layer 7
Layer Type	FC	Dropout	FC	Dropout	FC	Dropout	FC
Units	300	NA	300	NA	300	NA	3
Activation	Tanh	NA	Tanh	NA	Tanh	NA	Softmax

Table 5.2: MNIST Dataset Model Architecture

MNIST Dataset Model Architecture							
	Layer 1	Layer 2	Layer 3	Layer 4	Layer 5	Layer 6	Layer 7
Layer Type	FC	Dropout	FC	Dropout	FC	Dropout	FC
Units	300	NA	300	NA	300	NA	3
Activation	ReLU	NA	ReLU	NA	ReLU	NA	Softmax

5.5 Results and Discussion

Table 5.3: CIFAR Dataset Model Architecture

CIFAR Dataset Model Architecture							
	Layer 1	Layer 2	Layer 3	Layer 4	Layer 5	Layer 6	Layer 7
Layer Type	Conv	Conv	Flatten	Dropout	FC	Dropout	FC
Units	50	50	NA	NA	100	NA	10
Activation	ReLU	ReLU	NA	NA	ReLU	NA	Softmax
Kernel Size	3x3	3x3	NA	NA	NA	NA	NA

Chapter 6

Conclusion and Further Work

Bibliography

H-S Chang, E.Learned-Miller, A. (2017). Active bias: Training more accuracy neural networks by emphasizing high variance samples. *Advances in Neural Information Processing Systems*, 31:1–122.

Y.Bengio, J.Louradour, R. and J.Weston (2009). Curriculum learning. *Procedures of the International Conference on Machine Learning*, 26:41–18.