

# The Right Learning Rate at the Right Time

Anonymous submission

## Abstract

The selection of an appropriate optimizer plays a crucial role in deep learning, significantly influencing model performance. In the context of batch learning, where all data is available at a time, extensive research has been conducted to understand the implications of different optimizer choices and various optimization techniques for deep learning architectures have emerged. However, when it comes to online learning, where models need to adapt to evolving data streams, the exploration of optimizer choices is still relatively limited. This paper focuses on bridging the gap in knowledge by investigating how the choice of optimizer changes from batch learning to online learning scenarios. Our study conducts an in-depth analysis on the impact of optimizer, learning rate and batch size on the (i) predictive performance as well as the capability of adapting to changes on the underlying data pattern that might occur over time within the evolving data stream. Additionally, we explore practical choices for gradient-based online training of deep architectures, with a specific emphasis on adaptive methods. These methods dynamically adjust the learning rate based on gradient characteristics, offering potential advantages in online learning scenarios.

## Introduction

Deep learning models have demonstrated exceptional performance in various domains, with the choice of optimizer playing a crucial role in achieving outstanding results. In the context of batch learning, where all data is available simultaneously, extensive research has been conducted to explore different optimizer choices and optimization techniques for deep learning architectures. Numerous methods have emerged to effectively update the weights of these architectures. However, the investigation of optimizer choices in online learning, where models must adapt to evolving data streams, remains relatively limited.

This paper aims to bridge this knowledge gap by investigating how the choice of optimizer changes when transitioning from batch learning to online learning scenarios. Specifically, we address the following research questions:

- How does the choice for the optimizer change from batch to online learning?
- What are practical choices for gradient-based online training of deep architectures in online learning?

- Are adaptive optimization methods better suited in Online Deep Learning?

Within the first research question, we explore how the selection of an optimizer differs when moving from the traditional batch learning setting to the dynamic online learning scenario. We examine the suitability of various optimizer choices in online learning and their impact on model performance. The second research questions investigates practical choices for gradient-based online training of deep architectures. We analyze different optimization techniques and explore their effectiveness in adapting to evolving data streams while maintaining model performance. The third research question focuses on the performance of adaptive methods in online deep learning scenarios. These methods dynamically adjust the learning rate based on gradient characteristics, allowing models to adapt more effectively to changing data patterns. We compare the performance of adaptive methods against other optimization approaches to determine their suitability for online deep learning tasks. Through our in-depth analysis and experimentation, we aim to enhance our understanding of optimizer choices in online deep learning. By shedding light on the impact of optimizers, learning rates, and batch sizes, and comparing the effectiveness of adaptive methods, we aim to enable researchers and practitioners to make informed decisions when selecting optimization techniques for real-time learning tasks.

## Research Questions

- How does the choice for the optimizer change from batch to online learning? -¿ Batch learning: goal maximize batch size for best performance -¿ suitable for online learning? -¿ Learning Rate -¿ In-depth analysis on the influence of the learning rate and batch size on the effectiveness of online stochastic gradient descent -¿ Grafik SGD,
- What are practical choices for gradient-based online training of deep architectures in online learning? -¿ Description adaptive methods
- Are adaptive methods better suited in Online Deep Learning? -¿ Table with results

Motivation:

- Learning rate is one of the most important hyperparameters when it comes to training deep learning models.

- A small learning rate causes slow convergence, while an excessive learning rate causes instability (observable in the form of oscillations) or even increases in the training objective.
- Bengio (2012) goes as far as to say that out of all possible hyperparameters, if any, the learning rate is the one to be tuned.
- This applies arguably even more so to online learning. Unlike in batch learning, for instance, repeating the entire process if training diverges is infeasible, due to the restriction that previously observed data samples can, if at all, only be stored in small numbers.
- Furthermore, the trained model must be ready to predict at any given moment meaning that an inadequate learning rate has an immediate impact on the quality of generated predictions.
- In batch learning it is still common practice to determine the learning rate by running multiple training runs, beginning with a high learning rate and restarting the process with a lower learning rate every time the objective value diverges (Bengio 2012).

Contributions:

- In-depth analysis on the influence of the learning rate and batch size on the effectiveness of online stochastic gradient descent
- Experimental comparison of different SGD-based optimizers
- Investigation into possible approaches for determining an appropriate learning rate

## Learning Rate Scheduling

In the following, we will briefly outline the most important differences between the influence of the learning rate of first order gradient-based optimization methods in streaming- and in conventional batch learning environments.

First order gradient-based optimization approaches like stochastic gradient descent and its derivatives aim to iteratively minimize the error of a DL model using only first order gradient information at each step  $t$ . We denote the gradient of the prediction error for a mini-batch of training samples  $y_t, X_t \sim p_t$  with respect to model-parameters  $\theta$  as

$$g_t(\theta) = \nabla_{\theta} \mathcal{L}(y_t, f(X_t; \theta)), \quad (1)$$

where  $\mathcal{L}$  represents a loss function (e.g., cross-entropy for classification- or mean squared error for regression tasks). Using this notation, the update performed by SGD with a constant learning rate  $\eta$  at each iteration  $t$  is given by

$$\theta_{t+1} = \theta_t - \eta \cdot g_t(\theta_t). \quad (2)$$

Many previous works deal with the intricacies of SGD in the context of training deep architectures in a batch learning setting (see Bengio (2012); Bottou (2012); Goodfellow, Bengio, and Courville (2016)). The primary trade-off identified by these works when it comes to the selection of an appropriate learning rate is that between the speed of convergence and the amount of stochasticity. While increasing

the learning rate speeds up convergence, it also increases stochasticity and therefore leads to the divergence of the training criterion beyond a certain threshold. Smith and Le (2018) for instance, found that when modelling SGD as a stochastic differential equation, the “noise scale” is directly tied to  $\eta$  (Smith and Le 2018). In biological terms, increasing the learning rate increases plasticity, whereas decreasing it increases stability.

To get the best out of large and small values for  $\eta$ , a learning rate schedule  $(\eta_1, \dots, \eta_T)$ , which assigns a learning rate to each update step  $t \in 1, \dots, T$  can be used. It is for instance common to set a high learning rate initially and decrease it throughout the training process. This ensures fast convergence and a higher level of noise at the start of training, which has been claimed to help SGD skip over sharp minima with poor generalization (Hochreiter and Schmidhuber 1997; Chaudhari et al. 2017), while mitigating jumping around potential minima at later stages. Some have likened this procedure to simulated annealing, which shifts its focus from exploration at high temperatures to exploitation once temperatures have sufficiently decreased (Smith et al. 2018).

A simple way to achieve such a schedule is to exponentially decay  $\eta$  after each update by multiplying it with a factor  $\gamma < 1$ , although some practitioners have come to prefer schedules with sharper drop-offs like step schedules that decrease  $\eta$  by a larger margin after a certain number of updates (Smith et al. 2018). Other popular options include cyclic learning rate schedules which oscillate  $\eta$  between a base- and a maximum or minimum value over a predefined interval. Some studies (Smith 2017; Smith and Topin 2018) have found cyclic schedules to significantly speed up the convergence of neural networks even when compared to adaptive techniques like Adam (Kingma and Ba 2017). Wu et al. (2019) provide a detailed analysis on the effect of learning rate policies like the aforementioned ones.

In contrast to the field of conventional batch learning, the impact of the learning rate in stream-based deep learning is a lesser studied issue. According to Bifet et al. (2010) a machine learning model operating in such an environment must be able to

- R1:** process a single instance at a time,
- R2:** process each instance in a limited amount of time,
- R3:** use a limited amount of memory,
- R4:** predict at any time,
- R5:** adapt to changes in the data distribution.

These requirements lead to significant differences with respect to the problem of learning rate optimization compared to batch learning. Nevertheless, only few studies on the impact and tuning of the learning rate in stream-based deep learning exist.

In batch learning, the task of finding an optimal schedule for  $\eta$  can be defined as

$$\begin{aligned} \min_{\eta_0, \dots, \eta_T} \quad & \sum_{i=1}^V \mathcal{L}(y_i, f(X_i; \theta_T)) \\ \text{s.t.} \quad & X_i, y_i \sim p^{(v)} \quad \forall i \in 1, \dots, V, \end{aligned} \quad (3)$$

where  $p^{(v)}$  is a distribution of validation data, usually made up of a dataset split off from the training dataset. Verbally,

learning rate optimization in batch learning aims to find a schedule  $(\eta_0, \dots, \eta_T)$  leading to parameter values  $\theta_T$  at the end of a training run that in turn minimize the prediction error for validation data.

Under the requirements described above, however, we must define the task differently as

$$\begin{aligned} \min_{\eta_0, \dots, \eta_T} \quad & \sum_{t=0}^T \mathcal{L}(y_t, f(X_t; \theta_{t-1})) \\ \text{s.t.} \quad & X_t, y_t \sim p_t \quad \forall t \in 1, \dots, T. \end{aligned} \quad (4)$$

Compared to Problem (3), the most apparent difference is that there is no separate validation data. Instead, due to Requirement , the goal is to minimize  $\mathcal{L}$  with respect to the next instances  $X_{t+1}, y_{t+1}$  at each timestep  $t$ . This means that in contrast to Problem (3), not only the final parameters  $\theta_T$  but every parameter configuration  $\theta_t$  in the entire trajectory contributes equally to the objective. Therefore, speed of convergence is of much larger importance in the streaming setting, whereas the performance of the final parameters  $\theta_T$  has relatively little impact. Since memory is limited (Requirement ), it is also not possible to continue training on previously observed data as long as  $\mathcal{L}$  decreases, which puts an even greater emphasis on quick adaptation. Another difference to conventional batch learning is that the distribution  $p_t$  of the data stream might, and in practice most likely will, be subjected to change in the form of concept drift over time. Under such circumstances, the optimal parameter values  $\theta^*$  move throughout the progression of the stream increasing the distance to the model parameters.

Since the theoretically optimal learning rate  $\eta^*$  is proportional to the quadratic distance between initial and optimal parameters  $\|\theta_1 - \theta^*\|^2$  (albeit under some constraints like the absence of noise) (Carmon and Hinder 2023),  $\eta$  should be increased if concept drift occurs, increasing the model's ability to adapt to it.

Based on this concept, Kuncheva and Pluption (2008) introduced an adaptive schedule that updates the learning rate using

$$\eta_{t+1} = \eta_t^{1 + (\mathcal{L}_t - \bar{\mathcal{L}}_{t-1})}, \quad (5)$$

where  $\mathcal{L}_t$  represents the loss for the current sample, and  $\bar{\mathcal{L}}_{t-1}$  a rolling mean of past losses. By doing so, Kuncheva and Pluption (2008) increases in loss lead to increases in the learning rate and vice versa. While this approach seems intuitively sound, it bears a high risk of  $\eta$  increasing indefinitely, since increases in loss caused by an excessive learning rate would lead to a feedback loop.

We therefore propose a simple adaptation to popular decaying learning rate schedules that operates in a fixed value range and aims to increase the model's plasticity when needed by resetting  $\eta$  to its original value if a concept drift has been detected. While this does not allow for a continuous adaptation like the technique by Kuncheva and Pluption (2008) and might therefore be less suited for less severe drift, it offers more reliability since the initial learning rate cannot be exceeded.

In other works, the issue of lacking adaptability due to insufficient learning rates is often times simply avoided by using a less aggressive or cyclic learning rate.

Concept drift also complicates the tuning of  $\eta$ , since even if data is available beforehand drift would eventually cause the stream to diverge from the distribution of data used for tuning. In our experiments we investigate to which extent this effect impairs learning rate tuning on data available prior to deployment of a stream-based neural net. In this context we propose simple adaptations to the standard approach of learning rate tuning.

## Parameter-Free Optimizers

Due to the large computational cost of conventional learning rate tuning, several Parameter-free variants of SGD that eliminate the learning rate entirely have been developed. - schaul et al. - orabona et Al. - ...

However, all these optimizers target batch learning scenarios, and therefore don't explicitly account for concept drift, making it unclear which of these approaches is suitable for stream-based learning. Previous work on parameter-free optimization of online deep learning models is rather limited, and mostly focused on specific applications. - Recommender systems -federated learning

## Adaptive Optimizers

In the case of adaptive approaches like SGD with momentum (Rumelhart, Hinton, and Williams 1986), AdaGrad (Duchi, Hazan, and Singer 2011) or Adam (Kingma and Ba 2017) the learning rate also varies with respect to individual model parameters, in which case

$$\eta_t \in \mathbb{R}^D \quad \forall t \in 0, \dots, T. \quad (6)$$

For SGD with momentum,

$$\eta_t = \alpha \cdot \sum_{i=1}^t \beta^{t-i} \cdot \frac{g_i}{g_t} \quad (7)$$

- Adaptive learning rate optimizers for batch learning (Kingma and Ba 2017; Zeiler 2012; Duchi, Hazan, and Singer 2011; Tieleman and Hinton 2012) - Learning-Rate free optimizers (Wu, Ward, and Bottou 2020; Schaul, Zhang, and LeCun 2013; Orabona and Tommasi 2017; Miyaguchi and Kajino 2019; van Erven and Koolen 2016; Baydin et al. 2018)

## Experiments

Type	Data Stream	Samples	Features	Classes
Synth.	RBF abrupt	20000	20	5
	RBF incremental	20000	20	5
Real	Insects abrupt	52848	33	6
	Insects incremental	57018	33	6
	Insects incr.-grad.	24150	33	6
	Covertime <sup>1</sup>	100000	54	7
	Electricity	45312	8	2

Table 1: Datasets used for experimental evaluations.

Optimizer	Schedule	RBF abrupt	RBF incr.	Covertypes	Insects abrupt	Electricity
SGD	Fixed	93.70±.90	69.33±1.36	83.08±.18	71.12±.08	73.12±.42
	Exp.	94.00±.52	69.24±1.06	82.65±.19	71.32±.19	73.06±.42
	Exp. Reset	94.28±.37	69.76±.92	82.70±.27	71.27±.14	73.05±.45
	Step	94.00±.68	69.16±.99	82.74±.10	71.39±.17	72.96±.48
	Step Reset	94.04±.70	69.75±.81	83.03±.13	71.19±.14	73.18±.50
	Cyclic	<u>94.45±.25</u>	<u>73.72±1.16</u>	<b>83.40±.21</b>	71.41±.20	67.80±1.03
	Cyclic Reset	<b>94.50±.21</b>	<b>73.78±1.15</b>	<u>83.33±.13</u>	71.39±.15	67.83±1.00
Adam	Fixed	92.77±.41	66.46±4.39	78.85±.22	<b>75.08±.13</b>	69.23±.41
	Exp.	92.17±.82	63.91±3.52	78.53±.27	74.88±.15	69.33±.40
	Exp. Reset	92.83±1.19	62.03±2.59	77.05±.08	72.04±.49	68.31±.40
AdaGrad	Fixed	91.34±.83	50.39±3.60	81.07±.22	74.31±.34	76.64±1.92
SGDH	Fixed	91.03±.45	63.38±1.55	82.33±.12	67.35±.16	73.10±.10
COCOB	Fixed	93.40±.38	63.52±2.70	82.27±.46	74.75±.11	<b>84.30±.56</b>
WNGrad	Fixed	87.23±1.24	44.79±.76	76.95±.15	66.14±.15	70.74±.59
DAdaptSGD	Fixed	74.91±4.22	45.47±2.75	76.69±.79	50.05±11.26	66.03±1.75
DoG	Fixed	92.73±.59	<u>73.17±2.72</u>	83.07±.64	70.59±.26	71.53±.70

Table 2: Average prequential accuracy [%] for the three best learning rates.

Resetting Approach	RBF abrupt	RBF incr.	Covertypes	Insects abrupt	Electricity
ADWIN (Two Tailed)	<b>94.28±.37</b>	<u>69.76±.92</u>	82.70±.27	<u>71.27±.14</u>	<u>73.05±.45</u>
ADWIN (One Tailed)	94.25±.38	<u>70.22±2.58</u>	82.64±.20	<u>71.25±.15</u>	<u>73.00±.50</u>
ADWIN Weight Reset	71.79±.73	65.07±.31	82.54±.16	50.97±.36	70.19±1.00
KSWIN	<u>94.23±.55</u>	<u>70.10±1.98</u>	<b>83.01±.06</b>	<b>71.38±.16</b>	<u>73.13±.31</u>
P-KSWIN	93.86±.48	<b>70.71±1.38</b>	<b>83.01±.18</b>	<u>71.25±.15</u>	<b>73.26±.43</b>

Table 3: Average prequential accuracy [%] for the three best learning rates.

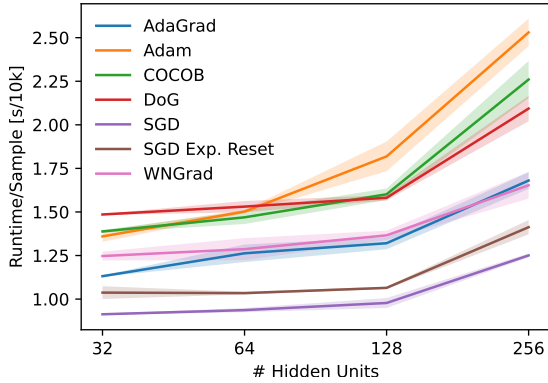


Figure 1: Runtimes in seconds per ten thousand samples for MLPs with three hidden layers of variable widths and RBF abrupt as input stream.

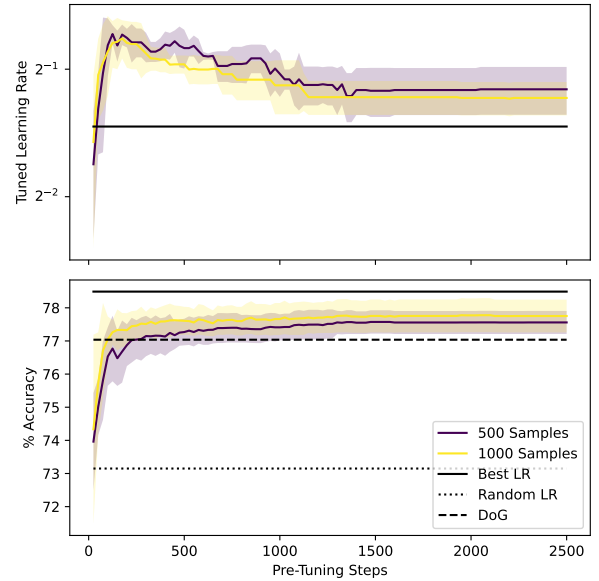


Figure 2: Pre-tuned LR (LR that maximizes accuracy on pre-tuning data) and resulting accuracy on data streams when using SGD and an exponential learning rate schedule with 500 or 1000 separate tuning samples. Results are averaged over all real-world datasets. The shaded area represents the  $1\sigma$ -interval.

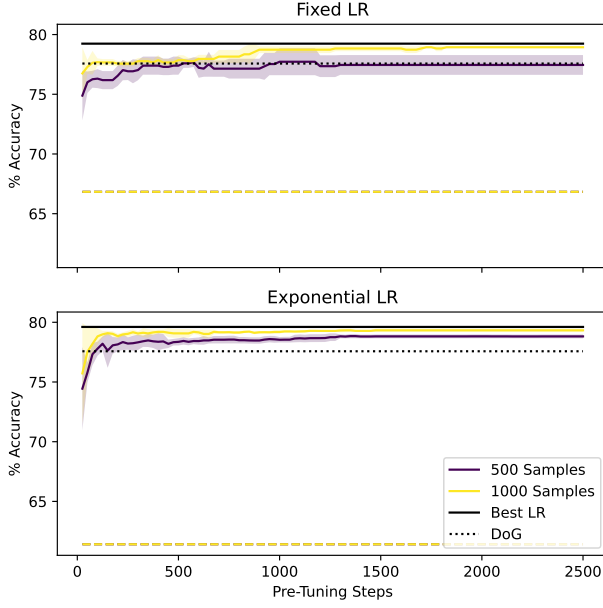


Figure 3: Accuracy achieved by pre-tuning on 500 or 1000 samples when using SGD with a fixed LR schedule (top) or an exponential schedule (bottom), averaged over all real-world datasets. The shaded area represents the  $1\sigma$ -interval.

We ran prequential evaluations using basic SGD with variable batch sizes and learning rates for synthetic data streams with and without incremental concept drift, the results of which are displayed in Figure 5. For static data, the average prequential accuracy over the entire stream gradually improves when moving up from an inadequately low learning rate until a certain point where training begins to diverge and performance consequently crashes. Based on our results, there seems to be an inverse relationship between batch size and both the optimal learning rate and the optimal accuracy, with larger batch sizes seemingly increasing the risk of divergence.

The primary cause of this relationship can be seen in Figure 7a, which depicts the trajectories of parameters of two identical MLPs trained on a bivariate data stream using batch sizes of either 4 or 16 samples. With a batch size of 16, the model has to “wait” four times longer until it receives new gradient information than the one trained with a batch size of 4. As a result, the gradient information gained from each sample on average becomes less recent and therefore less accurate with a larger batch size. Example: in Figure 7a

(8)

This effect is much stronger in the presence of concept drift as the results for RBF Incremental show.

It could be explained by the fact that the presence of concept drift exacerbates the gradient stochasticity caused by the delay between observation and learning of samples.

<sup>1</sup>We used the first 100k from a total of 581k examples only.

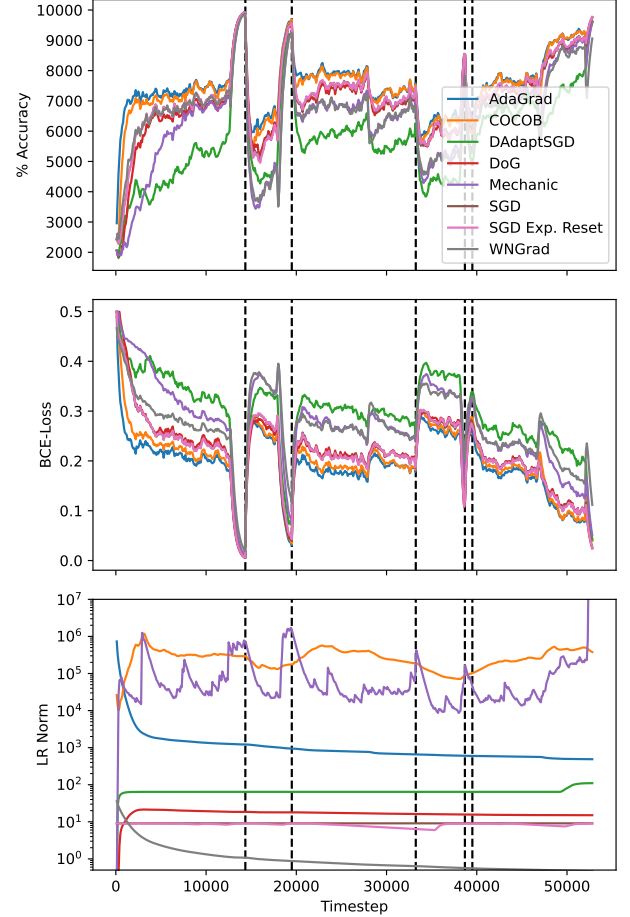


Figure 4: Prequential accuracy, binary cross-entropy loss and LR norms  $\|\eta_t\|$  over time for various optimization algorithms on Insects abrupt. Each dashed vertical line represents a concept drift. Lines are exponentially smoothed with a factor of 0.8.

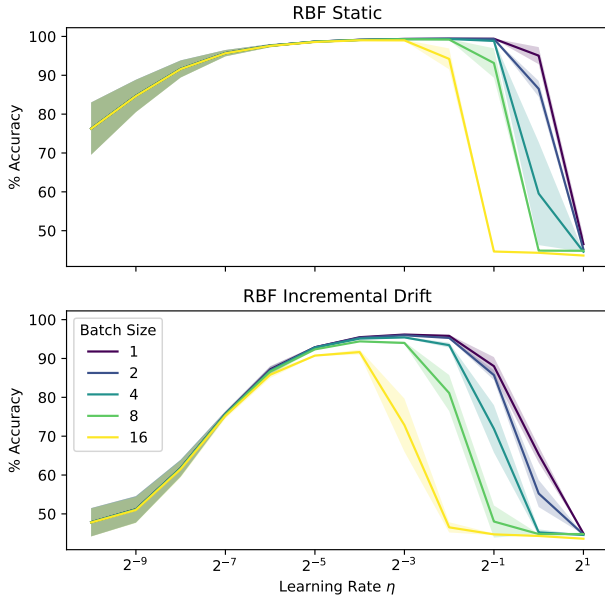


Figure 5: Average accuracy over static and incrementally drifting synthetic data streams in relation to SGD mini-batch size and learning rate  $\eta$ . Shaded areas mark the  $1\sigma$  interval.

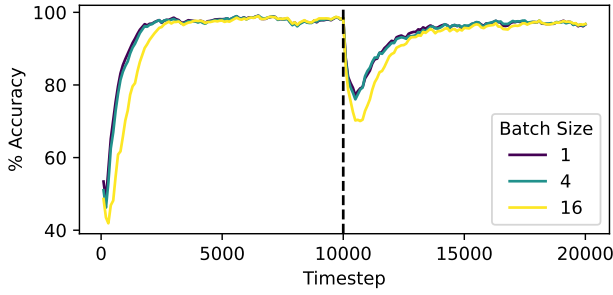


Figure 6: Prequential accuracy for different batch sizes on RBF abrupt data stream.

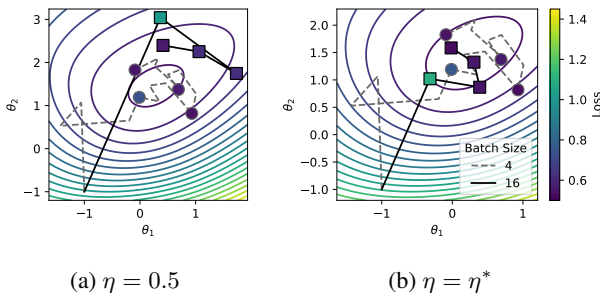


Figure 7: Trajectories of parameters  $\theta_1, \theta_2$  of a single layer MLP resulting from performing mini-batch SGD on an artificial stream of two-dimensional classification data. Markers are plotted after every 16th observed sample, with colors representing the expected prequential loss since the previous marker.

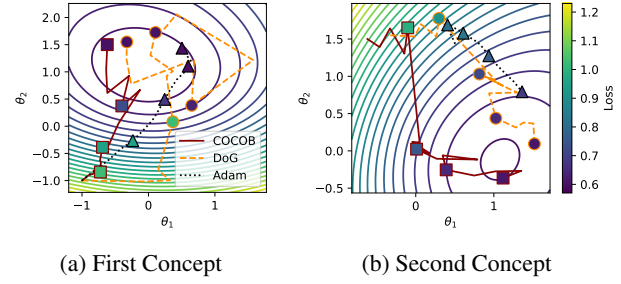


Figure 8: Weight trajectories of different adaptive optimizers for a synthetic data stream with an abrupt drift to a different concept after 64 samples. Markers highlight the parameter values after every 16th sample, with colors showing the expected prequential loss over the previous 16 samples.

## Conclusion

## References

- Baydin, A. G.; Cornish, R.; Rubio, D. M.; Schmidt, M.; and Wood, F. 2018. Online Learning Rate Adaptation with Hypergradient Descent. In *ICLR Proceedings*.
- Bengio, Y. 2012. Practical Recommendations for Gradient-Based Training of Deep Architectures. arxiv:1206.5533.
- Bifet, A.; Holmes, G.; Kirkby, R.; and Pfahringer, B. 2010. MOA: Massive Online Analysis. *Journal of Machine Learning Research*, 11.
- Bottou, L. 2012. Stochastic Gradient Descent Tricks. In Montavon, G.; Orr, G. B.; and Müller, K.-R., eds., *Neural Networks: Tricks of the Trade: Second Edition*, Lecture Notes in Computer Science, 421–436. Berlin, Heidelberg: Springer. ISBN 978-3-642-35289-8.
- Carmon, Y.; and Hinder, O. 2023. Making SGD Parameter-Free. arxiv:2205.02160.
- Chaudhari, P.; Choromanska, A.; Soatto, S.; LeCun, Y.; Baldassi, C.; Borgs, C.; Chayes, J.; Sagun, L.; and Zecchina, R. 2017. Entropy-SGD: Biasing Gradient Descent Into Wide Valleys. arxiv:1611.01838.
- Duchi, J. C.; Hazan, E.; and Singer, Y. 2011. Adaptive Sub-gradient Methods for Online Learning and Stochastic Optimization. *Journal of Machine Learning Research*, 12(61): 2121–2159.
- Goodfellow, I.; Bengio, Y.; and Courville, A. 2016. *Deep Learning*. MIT Press.
- Hochreiter, S.; and Schmidhuber, J. 1997. Flat Minima. *Neural Computation*, 9(1): 1–42.
- Kingma, D. P.; and Ba, J. 2017. Adam: A Method for Stochastic Optimization. arxiv:1412.6980.
- Kuncheva, L. I.; and Plampton, C. O. 2008. Adaptive Learning Rate for Online Linear Discriminant Classifiers. In Da Vitoria Lobo, N.; Kasparis, T.; Roli, F.; Kwok, J. T.; Georgiopoulos, M.; Anagnostopoulos, G. C.; and Loog, M., eds., *Structural, Syntactic, and Statistical Pattern Recognition*, volume 5342, 510–519. Berlin, Heidelberg: Springer Berlin Heidelberg. ISBN 978-3-540-89688-3 978-3-540-89689-0.



Miyaguchi, K.; and Kajino, H. 2019. Cogra: Concept-Drift-Aware Stochastic Gradient Descent for Time-Series Forecasting. *Proceedings of the AAAI Conference on Artificial Intelligence*, 33(01): 4594–4601.

Orabona, F.; and Tommasi, T. 2017. Training Deep Networks without Learning Rates Through Coin Betting. In *NIPS*.

Rumelhart, D. E.; Hinton, G. E.; and Williams, R. J. 1986. Learning Representations by Back-Propagating Errors. *Nature*, 323(6088): 533–536.

Schaul, T.; Zhang, S.; and LeCun, Y. 2013. No More Pesky Learning Rates. arxiv:1206.1106.

Smith, L. N. 2017. Cyclical Learning Rates for Training Neural Networks. arxiv:1506.01186.

Smith, L. N.; and Topin, N. 2018. Super-Convergence: Very Fast Training of Neural Networks Using Large Learning Rates. arxiv:1708.07120.

Smith, S. L.; Kindermans, P.-J.; Ying, C.; and Le, Q. V. 2018. Don't Decay the Learning Rate, Increase the Batch Size. arxiv:1711.00489.

Smith, S. L.; and Le, Q. V. 2018. A Bayesian Perspective on Generalization and Stochastic Gradient Descent. arxiv:1710.06451.

Tieleman, T.; and Hinton, G. 2012. Lecture 6.5-Rmsprop: Divide the Gradient by a Running Average of Its Recent Magnitude. In *COURSERA: Neural Networks for Machine Learning*. Coursera.

van Erven, T.; and Koolen, W. M. 2016. MetaGrad: Multiple Learning Rates in Online Learning. arxiv:1604.08740.

Wu, X.; Ward, R.; and Bottou, L. 2020. WNGrad: Learn the Learning Rate in Gradient Descent. arxiv:1803.02865.

Wu, Y.; Liu, L.; Bae, J.; Chow, K.-H.; Iyengar, A.; Pu, C.; Wei, W.; Yu, L.; and Zhang, Q. 2019. Demystifying Learning Rate Policies for High Accuracy Training of Deep Neural Networks. arxiv:1908.06477.

Zeiler, M. D. 2012. ADADELTA: An Adaptive Learning Rate Method. arxiv:1212.5701.