

# Learning Rate Optimization in Online Deep Learning

Anonymous<sup>1</sup>

No Institute Given

**Abstract.** Efficient training via gradient-based optimization techniques is an essential building block to the success of deep learning. Extensive research on the impact and the effective estimation of an appropriate learning rate has partly enabled these techniques. Despite the proliferation of data streams generated by IoT devices, digital platforms, etc., previous research has been primarily focused on batch learning, which assumes that all training data is available a priori. However, characteristics such as the gradual emergence of data and the occurrence of distributional shifts also known as *concept drift* pose additional challenges. Therefore, the findings on batch learning may not be applicable to streaming environments, where the underlying model needs to adapt each time a new data instance appears. In this work, we seek to address this knowledge gap by (i) evaluating and comparing typical learning rate schedules and optimizers, (ii) exploring adaptations of these techniques, and (iii) providing insights into effective learning rate tuning in the context of stream-based deep learning.

**Keywords:** Data streams · Learning Rate · Neural Networks.

## 1 Introduction

Deep learning models have shown exceptional performance in various domains. One of the main factors leading to such outstanding results is the choice of the optimization method used to train the target model. Almost all modern deep learning applications use first-order stochastic optimization methods such as *stochastic gradient descent*, which iteratively update the parameters of the underlying model based on gradient information. One of the most important variables of such algorithms is the step size or *learning rate* (LR).

As a result, many techniques for setting and optimizing the learning rate have emerged over the years (see Figure 1). For example, based on prior knowledge, the learning rate can be set as a fixed value or as a schedule that changes the step size over time. Alternatively, one could use an adaptive learning rate technique that considers historical gradient information to modify the learning rate at each iteration. In batch learning scenarios, where all training data is available a priori, the above methods are well researched. Despite the increasing prevalence of online learning environments, where data becomes available as part of a data stream, their use in such scenarios has received little research attention.

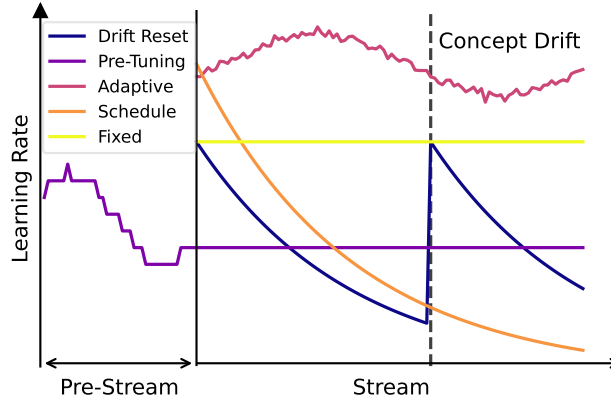


Fig. 1: Overview of different learning rate optimization approaches.

According to Bifet *et al.* [4] a machine learning model operating on a data stream 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 pose additional challenges in selecting an appropriate optimizer and learning rate. To enable more informed decisions when optimizing the learning rate, we provide insight into these challenges and empirically evaluate commonly used optimization techniques in an online learning setting (i). Furthermore, we introduce a *drift reset* mechanism to adapt the learning rate to concept drifts that commonly occur in streaming environments (ii). Finally, we propose a *pre-tuning* approach to effectively optimize the learning rate of online deep learning models *pre-stream* (see Figure 1) (iii).

## 2 Learning Rate in First-Order Optimization

In the following, we will explain the theoretical background of the learning rate hyperparameter in first-order stochastic optimization. First-order stochastic optimization algorithms, such as stochastic gradient descent, typically aim to solve the following problem

$$\min_{\theta} \mathbb{E}_{x \sim p(x)} [\mathcal{L}(x, \theta)], \quad (1)$$

where  $\mathcal{L}(x, \theta)$  is a loss function that quantifies the prediction error of the model given a mini-batch of data samples  $x$  and model parameters  $\theta$ . The blueprint process for solving this problem using first-order stochastic optimization consists of the following steps for each iteration  $t \in 0, \dots, T$ :

1. Draw a mini-batch of samples  $x_t$  from the distribution  $p(x)$ .
2. Compute the loss  $\mathcal{L}_t = \mathcal{L}(x_t, \theta_t)$  for  $x_t$  and current parameters  $\theta_t$ .
3. Compute gradient  $g_t = \nabla_{\theta_t} \mathcal{L}_t$  with respect to the parameters.
4. Update the parameters for the next iteration using  $g_t$  and possibly information from previous iterations.

For basic SGD, we can define the parameter update performed at the end of each iteration as  $\theta_t = \theta_t - \eta_t \cdot g_t$ , where  $\eta_t$  denotes the step size or *learning rate* at timestep  $t$ .

The primary trade-off with respect to  $\eta$  is that increasing it speeds up convergence, but also increases stochasticity and the risk of divergence [2]. In fact, Smith and Le [23], found that the “noise scale” of SGD is tied to  $\eta$  [23].

## 2.1 Learning Rate Schedules

Often, the performance of a model can be improved by using a schedule that changes the learning rate as training progresses [28]. For example, to ensure fast convergence early in training while mitigating jumping around potential minima later, it is common to use a decaying schedule that starts with a large learning rate and decreases over time. An additional benefit of this approach is potentially better generalization, since larger learning rates can help skip sharp minima with poor generalization [13, 6].

Commonly used forms of decay are exponential decay, where  $\eta_t$  is calculated as  $\eta_t = \eta_0 \cdot \gamma^t$ , with  $\gamma < 1$ , and stepwise decay, which for a regular interval between steps of length  $s$  is given as  $\eta_0 \cdot \gamma^{\lfloor t/s \rfloor}$ . Another common approach is to decay  $\eta$  each time the training loss plateaus for a given number of iterations. Other popular schedules include cyclic learning rates that oscillate  $\eta$  between two values over a predefined interval. For a triangular cycle, the learning rate is computed as

$$\eta_t = \eta_0 + \frac{\hat{\eta} - \eta_0}{2s} \cdot \min_i \{|t - i \cdot s|\}, \quad (2)$$

where  $\hat{\eta}$  is the learning rate at the midpoint of each cycle of length  $s$ . Some studies [21, 22] have found that cyclic schedules can significantly speed up the convergence of neural networks, in some cases even compared to adaptive techniques like Adam [15]. While there are many alternatives, in this work we focus on exponential, step, and cyclic learning rates as some of the most commonly used generic schedules. For a comprehensive overview and detailed analysis of learning rate policies, see Wu *et al.* [28].

## 2.2 Adaptive Learning Rates

Several studies have proposed *adaptive optimizers* that increase the robustness of the training process with respect to the learning rate. These optimizers adjust the step size based on previous gradients at each optimization step [10].

One of the earlier techniques in this category is *AdaGrad* [10], which scales the learning rate based on the sum of squares of past gradients for each parameter,

Table 1: Overview of additional time- and space-complexity of evaluated adaptive first-order optimizers compared to basic SGD. Values are given in big O notation with respect to the number of model parameters  $D$ . For generic approaches<sup>†</sup>, we assume SGD as the base optimizer.

Optimizer	Runtime	Space	Param. specific	LR free
AdaGrad	$\mathcal{O}(5D)$	$\mathcal{O}(1D)$	✓	✗
Adam	$\mathcal{O}(12D)$	$\mathcal{O}(2D)$	✓	✗
WNGrad	$\mathcal{O}(2D)$	$\mathcal{O}(0)$	✗	✗
COCOB	$\mathcal{O}(14D)$	$\mathcal{O}(4D)$	✓	✓
HD <sup>†</sup>	$\mathcal{O}(2D)$	$\mathcal{O}(1D)$	✗	✗
Mechanic <sup>†</sup>	$\mathcal{O}(10D)$	$\mathcal{O}(1D)$	✓	✓
DoG	$\mathcal{O}(5D)$	$\mathcal{O}(1D)$	✗	✓
D-Adapt <sup>†</sup>	$\mathcal{O}(6D)$	$\mathcal{O}(2D)$	✗	✓

resulting in a parameter-specific step size. Several other approaches, such as AdaDelta [29] and RMSProp, later built on AdaGrad’s scaling approach. The same is true for the widely used Adam optimizer [15], which adds a momentum term from prior gradients to speed up convergence for parameters with consistent derivatives. Another AdaGrad based optimizer is *WNGrad* [27], which adaptively scales each parameter update based on the squared sum of past gradients.

So-called parameter-free variants of SGD aim to eliminate the learning rate altogether by optimizing it as training progresses. For example, the *COCOB* algorithm [17] models parameter optimization as a gambling problem, where the goal is to maximize the reward from betting on each gradient. The resulting strategy is equivalent to running a meta-optimization algorithm that estimates the expected optimal learning rate [17]. Several other studies [25, 1, 7] have also used the idea of learning  $\eta$  via a meta-optimization process. The *hypergradient descent* (HD) approach [1], for example, adapts the learning rate of a base optimizer like SGD using a meta-gradient descent procedure, although this does not remove the learning rate entirely, but replaces it with a less sensitive hypergradient step size. Mechanic [7] pursues the same goal by applying a meta *online convex optimization* (OCO) algorithm to an arbitrary base optimizer.

Research has shown that in an OCO problem setting with stationary data, the worst-case optimal fixed step size for SGD is

$$\eta^* = \frac{\|\theta_0 - \theta^*\|}{\sqrt{\sum_{t=0}^T \|g_t\|^2}}. \quad (3)$$

Multiple parameter-free optimizers, make use of this notion. As its name suggests, the *Distance over Gradients* (DoG) [14] algorithm estimates the unknown numerator in Equation 3 as the maximum distance  $\max_{i < t} \|\theta_0 - \theta_i\|$  between the initial parameters and the parameters of all previous iterations. DoG additionally uses polynomial decay averaging as proposed by Shamir and Zhang [20]. *D-Adaptation* by Defazio and Mishchenko [8], on the other hand, employs weighted dual averaging [9] to compute bounds on the distance between initial and optimal parameters. Although adaptive optimization techniques seem intu-

itively well suited for non-stationary data, their application to data streams has rarely been investigated. Therefore, we assess the suitability of some of the most prominent adaptive optimizers, listed in Table 1, for stream-based learning.

### 3 Learning Rate in Online Learning

In a batch learning setting, optimizing the learning rate involves minimizing the expected loss on a hold-out set of validation data. Formally, we can express this task as

$$\min_{\eta_0, \dots, \eta_T} \sum_{i=0}^V \mathcal{L}(x_i, \theta_T), \quad (4)$$

where all  $x_i$  are part of a separate validation dataset and  $\theta_T$  are the parameters at the end of training. In online learning where data is generated incrementally, this notion of learning rate optimization is not feasible. Due to requirements **R1-R5**, models operating in an online learning environment should be evaluated in a *prequential* manner [4], where each sample  $x_t$  in the data stream is first used to test and then to train the model, ensuring that testing is done on previously unobserved data.

Training in such a scenario can be more accurately modeled as an online convex optimization problem [19, 12], where the optimizer suffers a loss  $\mathcal{L}(x_t, \theta_t)$  and produces updated parameters  $\theta_{t+1}$  at each iteration of the data stream. Learning rate optimization in this setting can be formulated as

$$\min_{\eta_0, \dots, \eta_T} \sum_{t=0}^T \mathcal{L}_t(x_t, \theta_t). \quad (5)$$

Compared to Equation (4), Equation (5) features some key differences. Due to the requirement to be able to predict at any time (**R4**), the goal is to minimize the total sum of losses over all timesteps of the prequential evaluation process, instead of the validation loss for the final parameters  $\theta_T$ . Therefore, the speed of convergence is more important in the streaming setting, while the performance of the final  $\theta_T$  parameter has a much smaller impact. Since memory is limited (Requirement 1), it is also not possible to continue training on previously observed data as long as the loss decreases, which puts even more emphasis on fast adaptation. At the same time, a higher learning rate that temporarily increases the loss by skipping local minima may be suboptimal with respect to the equation (5), even if it eventually leads to a lower loss. Another difference to conventional batch learning is that the loss function  $\mathcal{L}_t$  is time-dependent, due to the fact that data streams are often subject to distributional changes in the form of so-called *concept drift* over time.

#### 3.1 Learning Rate Tuning

The differences in evaluation schemes described above may cause conventional learning rate tuning to produce poor results for stream-based learning. Therefore,

we propose a modified tuning approach that approximates the equation (5), which we call learning rate *pre-tuning*.

To emulate the target data stream, we continuously draw samples with replacement from the tuning data in a bootstrapping procedure instead of training on all data for multiple epochs. By doing so, we aim to increase the variability of the data and thus the similarity to a real data stream. We then optimize the learning rate with respect to the average prequential performance over the emulated stream using an arbitrary parameter search technique. We provide a detailed experimental evaluation of our approach in Section 4.

### 3.2 Concept Drift Adaptation

Concept drift requires repeated adaptation of the model parameters. If post-drift training is interpreted as a new online optimization problem, the worst-case optimal learning rate can be computed according to Equation 3, replacing the initial parameter values  $\theta_0$  with the values at the time of drift onset  $\theta_{t_d}$ . As a result, stronger drifts that cause  $\theta^*$  to move away from  $\theta_{t_d}$  can benefit from larger learning rates.

Based on this notion, we propose a simple adaptation to decaying learning rate schedules that resets  $\eta$  to its original value when a concept drift is detected. An exponential schedule modified with our approach will thus yield learning rates of

$$\eta_t = \eta_0 \cdot \gamma^{t-t_d}, \quad (6)$$

where  $t_d$  marks the timestep at which the drift was last detected. For drift detection, we apply ADWIN [3] to the prequential losses. To avoid mistakenly detecting loss decreases as concept drift, we use a one-sided ADWIN variant that tests only for increases.

Our approach is similar to some *forgetting mechanisms* commonly used in conventional non-deep online learning [11]. To improve model plasticity, such mechanisms partially or completely reset the current model parameters to their initial values. However, we hypothesize that this approach is not well suited for deep learning purposes. The reason is that, under the assumption of convexity, the newly initiated parameters must be closer to the optimal parameters  $\theta^*$  than the current parameters to be beneficial. We experimentally compare our approach with this weight-reset mechanism in Section 4.

## 4 Experiments

We empirically evaluate our hypotheses using the following setup<sup>1</sup>:

We use both synthetic and publicly available real-world classification datasets with different sizes and types of concept drift, listed in Table 2. Our evaluations include two synthetic *Random Radial Basis Function* (RBF) datasets that we manipulated to incorporate concept drift using the online learning framework *River* [16].

<sup>1</sup> Code available at [anonymous.4open.science/r/L0DL-D458/](https://anonymous.4open.science/r/L0DL-D458/).

Table 2: Datasets used for experimental evaluations. \*For Covertypes we use only the first 100,000 from a total of 581,012 instances.

	Type	Data Stream	Instances	Features	Classes
Synth.	RBF	abrupt	20,000	20	5
	RBF	incremental	20,000	20	5
Real		Insects abrupt	52,848	33	6
		Insects gradual	24,150	33	6
		Covertypes	100,000*	54	7
		Electricity	45,312	8	2

We further employ the *Electricity* and *Covertypes* [5] datasets, which are commonly used to evaluate online learning models, as well as two *Insects* datasets [24] with predefined types of concept drift. *Covertypes* is available on *OpenML* [26], while the remaining datasets are part of River. We employ a *PyTorch* [18] implementation of a single-hidden-layer MLP with units matching the number of input features. This choice is based on our experience that smaller models exhibit faster convergence and are therefore usually best suited for streaming environments.

We tune the base learning rate  $\eta_0$  of all but the parameter-free approaches using a grid search of ten geometrically spaced values and configure adaptive optimizers with their default parameter values. For HD, Mechanic and D-Adaptation we select standard SGD as the base algorithm. We select a fixed factor  $\gamma$  for decay schedules on all datasets. For the proposed learning rate resetting mechanism, we select a smaller decay factor and set the confidence level  $\delta$  for drift detection to  $10^{-4}$ . For our evaluations we process each dataset sequentially, emulating streams of mini-batches of four instances each, while recording the prequential accuracy and other metrics in intervals of 25 iterations. We report our results averaged over five random seeds.

#### 4.1 Learning Rate Schedules

To evaluate the effectiveness of our learning rate resetting mechanism for drift adaptation (see Equation (3.2)), we compare its average prequential accuracy to that of model weight resetting, commonly used in online learning.

As can be seen in Table 3, our approach clearly outperforms weight-resetting on all but *Covertypes*, but rarely yields an advantage over a static schedule. Since the oracle variant of our resetting approach, that was triggered only for timesteps with drift, shows only marginally better results, this performance gap is not caused by the drift detector. Rather, it appears that using a larger initial learning rate and slower decay is sufficient to ensure adequate adaptability to concept drift throughout a data stream, while providing better stability in later stages. Overall, a slower but static decay paired with a larger initial learning rate seems to be preferable to a more aggressive schedule with our drift resetting mechanism, unless severe concept drift as in *RBF incremental* is expected.

Table 3: Average prequential accuracy [%] for static and drift adaptive learning rate schedules with SGD. For LR-Reset Oracle we manually reset the learning rate at timesteps where concept drift occurs. Best values are shown in **bold**, values within the  $1\sigma$  interval of best values underlined.

Schedule	RBF abrupt	RBF incr.	Covertypes	Electricity	Insects abrupt	Insects gradual
Fixed	94.79 $\pm$ .32	70.95 $\pm$ 2.89	83.42 $\pm$ .50	73.77 $\pm$ .40	71.50 $\pm$ .08	75.31 $\pm$ .21
Step	<b>94.87<math>\pm</math>.28</b>	70.19 $\pm$ 3.02	82.89 $\pm$ .37	<u>73.62<math>\pm</math>.53</u>	<b>72.23<math>\pm</math>.27</b>	<u>75.83<math>\pm</math>.21</u>
Cyclic	<u>94.79<math>\pm</math>.32</u>	<b>74.96<math>\pm</math>.86</b>	<b>83.44<math>\pm</math>.08</b>	68.38 $\pm$ .81	71.74 $\pm$ .39	<u>75.64<math>\pm</math>.06</u>
Exponential	94.85 $\pm$ .29	70.23 $\pm$ 2.40	82.95 $\pm$ .26	73.51 $\pm$ .48	72.19 $\pm$ .37	<b>75.91<math>\pm</math>.14</b>
Weight-Reset	69.96 $\pm$ .38	65.13 $\pm$ .80	83.12 $\pm$ .13	70.08 $\pm$ 1.66	51.52 $\pm$ .90	62.55 $\pm$ 2.34
LR-Reset (Ours)	<u>94.83<math>\pm</math>.26</u>	73.38 $\pm$ 2.32	82.99 $\pm$ .20	<b>73.79<math>\pm</math>.62</b>	71.73 $\pm$ .20	75.52 $\pm$ .12
LR-Reset Oracle	95.12 $\pm$ .21	—	—	—	71.88 $\pm$ .26	—

With accuracy values within  $1\sigma$  of each other on all evaluated streams, the stepwise decay shows almost identical performance to the exponential decay. The accuracy of the cyclic schedule for *RBF incremental* and *Covertypes*, on the other hand, significantly outperforms the other static and adaptive schedules, but lags behind on all other streams. We also did not find an order of magnitude improvement in convergence speed as observed by [22] for the scenario studied.

## 4.2 Adaptive Learning Rates

Our results for adaptive optimizers displayed in Table 4, show a strong data dependency as none of the evaluated algorithms significantly outperforms its competitors on average. However, since SGD yields the best accuracy on *RBF abrupt* but is clearly surpassed on *Insects abrupt*, the type of concept drift does not seem to be significant. Due to its simplicity and favorable computational efficiency, it appears that standard SGD should be selected out of the non-parameter-free approaches in most cases. The SGD variant of Hypergradient Descent (HD) [1] and WNGrad [27] on the other hand seem to rarely be optimal choices.

In the category of learning rate free optimizers COCOB [17], outperformed its competitors on all but two datasets. It comes close to or even exceeds the best tuned approaches in terms of accuracy. Although yielding lower accuracy on average, DoG also comes within reach of the tuned methods, while offering much better runtime and memory efficiency compared to COCOB (see Table 1). Mechanic [7] and D-Adaptation [8] performed significantly worse than their competitors on the evaluated streams.

The learning rate curves shown in Figure 2a, provide an indication regarding the reason for the poor performance of WNGrad and D-Adaption. Whereas the learning rate of DoG quickly approaches the tuned SGD learning rate, WNGrad and D-Adaptation diverge considerably from it.

The learning rate of the best performing Adam exhibits spikes for most change points, suggesting some form of adaptability to drift. However, since the much worse performing Mechanic shows similar spikes, this is unlikely to be contributing significantly to Adam’s high accuracy on *Insects abrupt*. Instead,



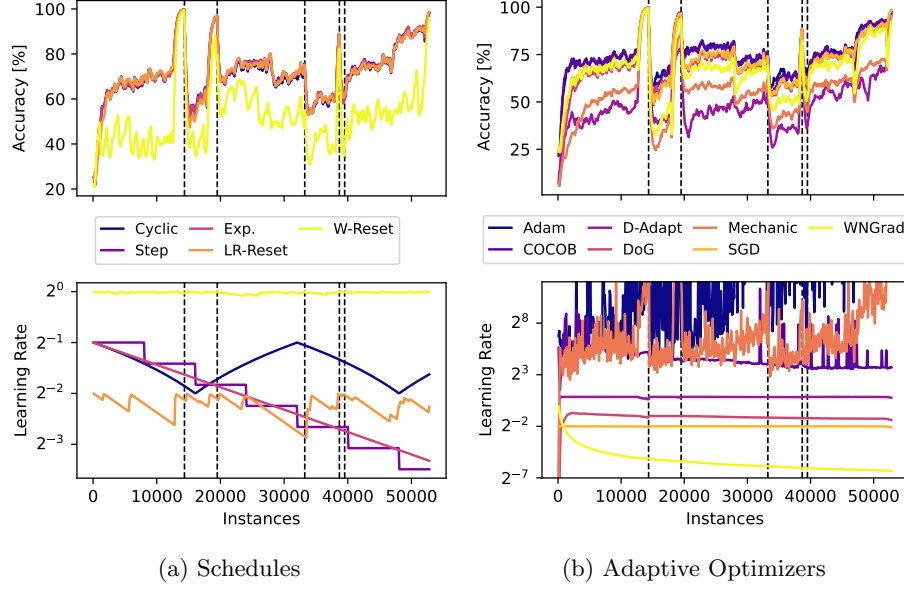


Fig. 2: Prequential accuracy and learning rate for different schedules and adaptive optimizers on *Insects abrupt* dataset. Concept drifts are marked by dashed lines. Accuracy is exponentially smoothed with a decay factor of 0.75.

it likely stems from its second moment scaling, which is also a feature of the similarly performing AdaGrad.

It may also be noted that the learning rates of the COCOB, Adam and Mechanic optimizers with parameter-specific learning rates exceed those of single value step sizes by multiple orders of magnitude. This is an effect of second moment scaling, which creates larger learning rates for parameters with small and consistent gradients [7]. Therefore, the norms of parameter updates generated by these approaches are not necessarily larger.

### 4.3 Learning Rate Tuning

We assess our pre-tuning approach using either 500 or 1000 instances held out from the beginning of each stream for tuning and evaluate MLPs with 64 or 128 hidden units per layer and either one or three hidden layers. In the following we focus on the smallest network as the one most representative of resource critical streaming applications. We select the learning rate and decay factor according to the optimal mean prequential accuracy. This is due to the potential bias towards learning rates with lower initial losses, as loss values often decrease notably during training.

Figure 3 shows the absolute difference between the learning rate resulting from the pre-tuning process and the optimal value  $|\eta_p - \eta^*|$  at each tuning step

Table 4: Average prequential accuracy [%] for adaptive optimizers and SGD. Best values are shown in **bold**, values within the  $1\sigma$  interval of best values underlined.

Optimizer	RBF abrupt	RBF incr.	Covertime	Electricity	Insects abrupt	Insects gradual
Tuned						
SGD	<b>94.79<math>\pm</math>.32</b>	70.95 $\pm$ 2.89	<b>83.42<math>\pm</math>.50</b>	73.77 $\pm$ .40	71.50 $\pm$ .08	75.31 $\pm$ .21
Adam [15]	93.45 $\pm$ .30	69.26 $\pm$ 5.14	79.01 $\pm$ .27	69.79 $\pm$ .54	<b>75.38<math>\pm</math>.24</b>	75.78 $\pm$ .74
AdaGrad [10]	92.45 $\pm$ 1.37	52.87 $\pm$ 6.62	81.68 $\pm$ .35	<b>76.99<math>\pm</math>1.20</b>	74.87 $\pm$ .40	<b>77.15<math>\pm</math>.27</b>
WNGrad [27]	87.30 $\pm$ .68	44.92 $\pm$ .73	76.98 $\pm$ .15	70.80 $\pm$ .59	66.25 $\pm$ .19	66.75 $\pm$ .40
HD [1]	93.92 $\pm$ .31	<b>72.29<math>\pm</math>2.90</b>	<u>83.36<math>\pm</math>.25</u>	73.83 $\pm$ .32	70.67 $\pm$ .06	73.37 $\pm$ .21
LR-Free						
COCOB [17]	<b>93.40<math>\pm</math>.38</b>	63.52 $\pm$ 2.70	82.27 $\pm$ .46	<b>84.30<math>\pm</math>.56</b>	<b>74.75<math>\pm</math>.11</b>	<b>77.00<math>\pm</math>.05</b>
DoG [14]	92.72 $\pm$ .59	<b>73.17<math>\pm</math>2.72</b>	<b>83.07<math>\pm</math>.64</b>	71.53 $\pm$ .70	70.59 $\pm$ .26	74.01 $\pm$ .21
D-Adapt [8]	74.91 $\pm$ 4.22	45.47 $\pm$ 2.75	76.69 $\pm$ .79	66.03 $\pm$ 1.75	50.05 $\pm$ 11.26	48.21 $\pm$ 10.62
Mechanic [7]	88.94 $\pm$ .58	49.26 $\pm$ 1.44	78.67 $\pm$ .18	50.73 $\pm$ 7.60	55.31 $\pm$ 21.47	65.80 $\pm$ .53

averaged over all real-world datasets. Batch tuning with 800 training and 200 validation samples initially yields a better approximation of the optimal learning rate. However, our streaming-specific approach undercuts the baseline after 1000 tuning steps, consistently decreasing and ultimately reaching approximately half of the approximation error observed in batch tuning. The performance also remains nearly identical, even when using only 500 samples for tuning, which demonstrates the data efficiency of pre-tuning. The superior performance of our approach is also reflected in the accuracy scores depicted in the right subfigure of Figure 3. After less than 1000 steps, our approach achieves notably higher accuracy than both conventional tuning and DoG, which we selected as a baseline due to being best performing parameter-free optimizer with a global learning rate.

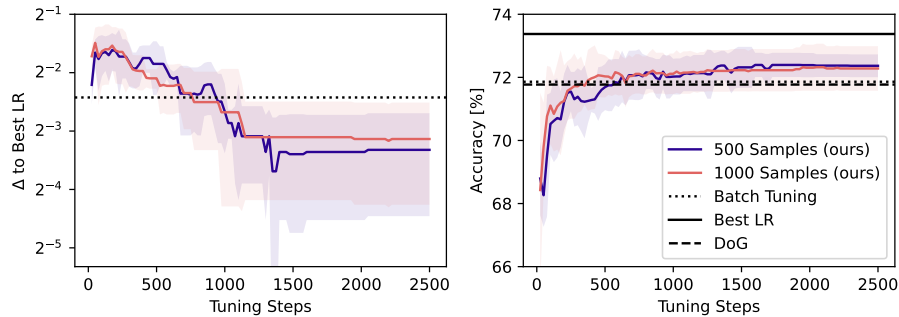


Fig. 3: Absolute difference between pre-tuned and optimal learning rate and resulting accuracy on data streams for SGD and an exponential learning rate schedule with 500 or 1000 tuning samples. DoG is not included in first subplot since it is not intended for use with LR decay. Results are averaged over all real-world datasets. The shaded area represents the  $1\sigma$ -interval.

In conclusion, our proposed tuning approach enables significantly better learning rate selection for prequential evaluation on data streams compared to both conventional tuning and DoG. Additionally, pre-tuning has the benefit that once completed, no additional memory or runtime costs are incurred. In streaming applications, where computing resources are often times a limiting factor, this could be a critical advantage. Although, if computational efficiency is insignificant, the highly performant but expensive COCOB [17] or the slightly less performant and much less expensive DoG [14] may be more appropriate.

## 5 Conclusion

In this work, we investigate the influence and selection of the learning rate and optimization procedure with respect to deep learning in streaming environments. We first provide theoretical background on discrepancies between learning rate optimization in conventional batch learning and online learning. Based on these differences, we derive a simple mechanism resetting the learning rate on concept drift occurrences. We then give an overview learning rate free algorithms popular in batch learning, which we compare experimentally on multiple synthetic and real-world datasets, finding both COCOB and DoG to come close to the performance of optimizers with tuned learning rates. Lastly, we introduce a streaming specific learning rate tuning approach that grants significant performance increases over conventional tuning via a train-validation split.

## References

1. Baydin, A.G., Cornish, R., Rubio, D.M., Schmidt, M., Wood, F.: Online Learning Rate Adaptation with Hypergradient Descent. In: ICLR Proceedings (2018)
2. Bengio, Y.: Practical Recommendations for Gradient-Based Training of Deep Architectures, (2012).
3. Bifet, A., Gavaldà, R.: Learning from Time-Changing Data with Adaptive Windowing. In: Proceedings of the 2007 SIAM International Conference on Data Mining, pp. 443–448. Society for Industrial and Applied Mathematics (2007). <https://doi.org/10.1137/1.9781611972771.42>
4. Bifet, A., Holmes, G., Kirkby, R., Pfahringer, B.: MOA: Massive Online Analysis. *Journal of Machine Learning Research* **11** (2010)
5. Blackard, J.: Coverttype, UCI Machine Learning Repository (1998).
6. Chaudhari, P. *et al.*: Entropy-SGD: Biasing Gradient Descent Into Wide Valleys, (2017).
7. Cutkosky, A., Defazio, A., Mehta, H.: Mechanic: A Learning Rate Tuner, (2023). <http://arxiv.org/abs/2306.00144>.
8. Defazio, A., Mishchenko, K.: Learning-Rate-Free Learning by D-Adaptation, (2023). <http://arxiv.org/abs/2301.07733>.
9. Duchi, J.C., Agarwal, A., Wainwright, M.J.: Dual Averaging for Distributed Optimization: Convergence Analysis and Network Scaling. *IEEE Transactions on Automatic Control* **57**(3), 592–606 (2012). <https://doi.org/10.1109/TAC.2011.2161027>

10. Duchi, J.C., Hazan, E., Singer, Y.: Adaptive Subgradient Methods for Online Learning and Stochastic Optimization. *Journal of Machine Learning Research* **12**(61), 2121–2159 (2011)
11. Gama, J., Žliobaitė, I., Bifet, A., Pechenizkiy, M., Bouchachia, A.: A Survey on Concept Drift Adaptation. *ACM Computing Surveys* **46**(4), 1–37 (2014). <https://doi.org/10.1145/2523813>
12. Hazan, E.: Introduction to Online Convex Optimization. *Foundations and Trends® in Optimization* **2**(3-4), 157–325 (2016). <https://doi.org/10.1561/24000000013>
13. Hochreiter, S., Schmidhuber, J.: Flat Minima. *Neural Computation* **9**(1), 1–42 (1997). <https://doi.org/10.1162/neco.1997.9.1.1>
14. Ivgi, M., Hinder, O., Carmon, Y.: DoG Is SGD’s Best Friend: A Parameter-Free Dynamic Step Size Schedule, (2023). <http://arxiv.org/abs/2302.12022>.
15. Kingma, D.P., Ba, J.: Adam: A Method for Stochastic Optimization, (2017).
16. Montiel, J. *et al.*: River: Machine Learning for Streaming Data in Python. (2021)
17. Orabona, F., Tommasi, T.: Training Deep Networks without Learning Rates Through Coin Betting. In: NIPS (2017)
18. Paszke, A. *et al.*: PyTorch: An Imperative Style, High-Performance Deep Learning Library. In: Wallach, H. (ed.) *NeurIPS Proceedings*, pp. 8024–8035. Curran Associates, Inc. (2019). <http://papers.neurips.cc/paper/9015-pytorch-an-imperative-style-high-performance-deep-learning-library.pdf>
19. Shalev-Shwartz, S.: Online Learning and Online Convex Optimization. *Foundations and Trends® in Machine Learning* **4**(2), 107–194 (2011). <https://doi.org/10.1561/22000000018>
20. Shamir, O., Zhang, T.: Stochastic Gradient Descent for Non-smooth Optimization: Convergence Results and Optimal Averaging Schemes, (2012). <http://arxiv.org/abs/1212.1824>.
21. Smith, L.N.: Cyclical Learning Rates for Training Neural Networks, (2017). <https://doi.org/10.48550/arXiv.1506.01186>.
22. Smith, L.N., Topin, N.: Super-Convergence: Very Fast Training of Neural Networks Using Large Learning Rates, (2018). <http://arxiv.org/abs/1708.07120>.
23. Smith, S.L., Le, Q.V.: A Bayesian Perspective on Generalization and Stochastic Gradient Descent, (2018). <https://doi.org/10.48550/arXiv.1710.06451>.
24. Souza, V.M.A., dos Reis, D.M., Maletzke, A.G., Batista, G.E.A.P.A.: Challenges in Benchmarking Stream Learning Algorithms with Real-world Data. *Data Mining and Knowledge Discovery* **34**(6), 1805–1858 (2020). <https://doi.org/10.1007/s10618-020-00698-5>
25. van Erven, T., Koolen, W.M.: MetaGrad: Multiple Learning Rates in Online Learning, (2016). <https://doi.org/10.48550/arXiv.1604.08740>.
26. Vanschoren, J., van Rijn, J.N., Bischl, B., Torgo, L.: OpenML: Networked Science in Machine Learning. *ACM SIGKDD Explorations Newsletter* **15**(2), 49–60 (2014). <https://doi.org/10.1145/2641190.2641198>
27. Wu, X., Ward, R., Bottou, L.: WNGrad: Learn the Learning Rate in Gradient Descent, (2020).
28. Wu, Y. *et al.*: Demystifying Learning Rate Policies for High Accuracy Training of Deep Neural Networks, (2019). <http://arxiv.org/abs/1908.06477>.
29. Zeiler, M.D.: ADADELTA: An Adaptive Learning Rate Method, (2012).