

# Data-Driven Mirror Descent with Input-Convex Neural Networks

Tan, Hong Ye; Mukherjee, Subhadip; Tang, Junqi; Schönlieb, Carola-Bibiane

DOI:

[10.1137/22M1508613](https://doi.org/10.1137/22M1508613)

License:

Creative Commons: Attribution (CC BY)

*Document Version*

Peer reviewed version

*Citation for published version (Harvard):*

Tan, HY, Mukherjee, S, Tang, J & Schönlieb, C-B 2023, 'Data-Driven Mirror Descent with Input-Convex Neural Networks', *SIAM Journal on Mathematics of Data Science*, vol. 5, no. 2, pp. 558-587.  
<https://doi.org/10.1137/22M1508613>

[Link to publication on Research at Birmingham portal](#)

## General rights

Unless a licence is specified above, all rights (including copyright and moral rights) in this document are retained by the authors and/or the copyright holders. The express permission of the copyright holder must be obtained for any use of this material other than for purposes permitted by law.

- Users may freely distribute the URL that is used to identify this publication.
- Users may download and/or print one copy of the publication from the University of Birmingham research portal for the purpose of private study or non-commercial research.
- User may use extracts from the document in line with the concept of 'fair dealing' under the Copyright, Designs and Patents Act 1988 (?)
- Users may not further distribute the material nor use it for the purposes of commercial gain.

Where a licence is displayed above, please note the terms and conditions of the licence govern your use of this document.

When citing, please reference the published version.

## Take down policy

While the University of Birmingham exercises care and attention in making items available there are rare occasions when an item has been uploaded in error or has been deemed to be commercially or otherwise sensitive.

If you believe that this is the case for this document, please contact [UBIRA@lists.bham.ac.uk](mailto:UBIRA@lists.bham.ac.uk) providing details and we will remove access to the work immediately and investigate.

# Data-Driven Mirror Descent with Input-Convex Neural Networks\*

Hong Ye Tan<sup>†</sup>, Subhadip Mukherjee<sup>†‡</sup>, Junqi Tang<sup>§</sup>, and Carola-Bibiane Schönlieb<sup>†</sup>

**Abstract.** *Learning-to-optimize* is an emerging framework that seeks to speed up the solution of certain optimization problems by leveraging training data. Learned optimization solvers have been shown to outperform classical optimization algorithms in terms of convergence speed, especially for convex problems. Many existing data-driven optimization methods are based on parameterizing the update step and learning the optimal parameters (typically scalars) from the available data. We propose a novel functional parameterization approach for learned convex optimization solvers based on the classical mirror descent (MD) algorithm. Specifically, we seek to learn the optimal Bregman distance in MD by modeling the underlying convex function using an input-convex neural network (ICNN). The parameters of the ICNN are learned by minimizing the target objective function evaluated at the MD iterate after a predetermined number of iterations. The inverse of the mirror map is modeled approximately using another neural network, as the exact inverse is intractable to compute. We derive convergence rate bounds for the proposed learned mirror descent (LMD) approach with an approximate inverse mirror map and perform extensive numerical evaluation on various convex problems such as image inpainting, denoising, learning a two-class support vector machine (SVM) classifier and a multi-class linear classifier on fixed features.

**Key words.** Mirror Descent, data-driven convex optimization solvers, input-convex neural networks, inverse problems.

**AMS subject classifications.** 46N10, 65K10, 65G50

**1. Introduction.** Convex optimization problems are pivotal in many modern data science and engineering applications. These problems can generally be formulated as

$$(1.1) \quad \min_{x \in \mathcal{X}} [f(x) + g(x)],$$

where  $\mathcal{X}$  is a Hilbert space, and  $f, g : \mathcal{X} \rightarrow \bar{\mathbb{R}}$  are proper, convex, and lower semi-continuous (l.s.c.) functions. In different scenarios,  $f$  and  $g$  have different levels of regularity such as differentiability or strong convexity. In the context of inverse problems,  $f$  can be a data fidelity loss and  $g$  a regularization function.

In the past few decades, extensive research has gone into developing efficient and provably convergent optimization algorithms for finding the minimizer of a composite objective function as in (1.1), leading to several major theoretical and algorithmic breakthroughs. For generic convex programs with first-order oracles, optimal algorithms have been proposed under different levels of regularity [24, 17, 18], which are able to match the complexity lower-bounds of the problem class. Although there exist algorithms that are optimal for generic problem classes, practitioners in different scientific areas usually only need to focus on a very narrow

---

\*Submitted to the SIAM J. on Mathematics of Data Science

<sup>†</sup>Department of Applied Mathematics and Theoretical Physics, University of Cambridge, UK ([hyt35@cam.ac.uk](mailto:hyt35@cam.ac.uk), [sm2467@cam.ac.uk](mailto:sm2467@cam.ac.uk), [jt814@cam.ac.uk](mailto:jt814@cam.ac.uk), [cbs31@cam.ac.uk](mailto:cbs31@cam.ac.uk)).

<sup>‡</sup>Department of Computer Science, University of Bath, UK ([sm3655@bath.ac.uk](mailto:sm3655@bath.ac.uk)).

<sup>§</sup>School of Mathematics, University of Birmingham, UK ([j.tang.2@bham.ac.uk](mailto:j.tang.2@bham.ac.uk)).

36 subclass, for which usually neither tight complexity lower-bounds nor optimal algorithms are  
37 known. As such, it is extremely difficult and impractical to either find tight lower-bounds or  
38 handcraft specialized optimal algorithms for every single subclass in practice.

39 The aim of this work is *learning to optimize* convex objectives of the form (1.1) in a  
40 provable manner. Learned optimization solvers have been proposed through various methods,  
41 including reinforcement learning and unsupervised learning [2, 3, 13, 19]. The goal is to  
42 minimize a fixed loss function as efficiently as possible, which can be formulated as minimizing  
43 the loss after a certain number of iterations, or minimizing the number of iterations required  
44 to attain a certain error. The common idea is to directly parameterize the update step as  
45 a neural network, taking previous iterates and gradients as arguments. These methods have  
46 been empirically shown to speed up optimization in various settings including training neural  
47 networks [19, 2]. However, many of these methods lack theoretical guarantees, and there is a  
48 lack of principled framework for integrating machine learning into existing classical algorithms.

49 Banert et al. developed a theoretically grounded method in [3] for parameterizing such  
50 update steps using combinations of proximal steps, inspired by proximal splitting methods.  
51 By learning the appropriate coefficients, the method was able to outperform the classical  
52 primal-dual hybrid gradient (PDHG) scheme [10]. However, having a fixed model limits the  
53 number of learnable parameters, and therefore the extent to which the solver can be adapted  
54 to a particular problem class. Banert et al. later drifted away from the framework of learning  
55 parameters of fixed models, and instead directly modeled an appropriate update function  
56 using a deviation-based approach, allowing for a more expressive parameterization [4].

57 Learned optimizers are sometimes modeled using classical methods, as the existing con-  
58 vergence guarantees can lead to insights on how neural networks may be incorporated with  
59 similar convergence guarantees. Even if such guarantees are not available, such as in the  
60 case of learned iterative shrinkage and thresholding algorithm (ISTA), they can still lead to  
61 better results on certain problems [13]. Conversely, Maheswaranathan et al. showed that cer-  
62 tain learned optimizers, parameterized by recurrent neural networks, can reproduce classical  
63 methods used for accelerating optimization [21]. By using a recurrent neural network taking  
64 the gradient as an input, the authors found that the learned optimizer expresses mechanisms  
65 including momentum, gradient clipping, and adaptive learning rates.

66 One related idea to our problem is meta-learning, also known as “learning to learn”. This  
67 typically concerns learning based on prior experience with similar tasks, utilizing techniques  
68 such as transfer learning, to learn how similar an optimization task is to previous tasks using  
69 statistical features [31]. Our problem setting will instead be mainly concerned with convex  
70 optimization problems, as there are concrete classical results for comparison.

71 Integrating machine learning models into classical algorithms can also be found notably  
72 in Plug-and-Play (PnP) algorithms. Instead of trying to learn a solver for a general class  
73 of optimization problems, PnP methods deal with the specific class of image restoration.  
74 By using proximal splitting algorithms and replacing certain proximal steps with generic  
75 denoisers, the PnP algorithms, first proposed by Venkatakrisnan et al. in 2013, were able to  
76 achieve fast and robust convergence for tomography problems [32]. This method was originally  
77 only motivated in an intuitive sense, with some analysis of the theoretical properties coming  
78 years later by Chan et al. [11], and more recently by Ryu et al. [29]. Most critically,  
79 many subsequent methods of showing convergence rely on classical analysis such as monotone

80 operator and fixed point theory, demonstrating the importance of having a classical model-  
81 based framework to build upon.

82 One of the main difficulties in learning to optimize is the choice of function class to learn on.  
83 Intuitively, a more constrained function class may allow for the learned method to specialize  
84 more. However, it is difficult to quantify the similarity between the geometry of different  
85 problems. Banert et al. proposed instead to use naturally or qualitatively similar function  
86 classes in [4], including regularized inverse problems such as inpainting or denoising, which  
87 will be used in this work as well.

88 **1.1. Contributions.** We propose to learn an alternative parameterization using mirror  
89 descent (MD), which is a well-known convex optimization algorithm first introduced by Ne-  
90 mirovsky and Yudin [23]. Typical applications of MD require hand-crafted mirror maps, which  
91 are limited in complexity by the requirement of a closed-form convex conjugate. We propose  
92 to replace the mirror map in MD with an input convex neural network (ICNN) [1], which has  
93 recently proved to be a powerful parameterization approach for convex functions [22]. By mod-  
94 eling the mirror map in this manner, we seek to simultaneously introduce application-specific  
95 optimization routines, as well as learn the problem geometry.

96 Using our new paradigm, we are able to obtain a learned optimization scheme with con-  
97 vergence guarantees in the form of regret bounds. We observe numerically that our learned  
98 mirror descent (LMD) algorithm is able to adapt to the structure of the class of optimization  
99 problems that it was trained on, and provide significant acceleration.

100 This paper is organized as follows. In [section 2](#), we recall the MD algorithm and the  
101 existing convergence rate bounds. [Section 3](#) presents our main results on convergence rate  
102 bounds with inexact mirror maps, and a proposed procedure of ‘learning’ a mirror map. In  
103 [section 4](#), we will show some simple examples of both MD and its proposed learned variant  
104 LMD in the setting where the inverse map is known exactly. [Section 5](#) deals with numerical  
105 experiments with inverse problems in imaging and linear classifier learning.

106 **2. Background.** In this section, we will outline the MD method as presented by Beck and  
107 Teboulle [5]. Convergence guarantees for convex optimization methods commonly involve a  
108 Lipschitz constant with respect to the Euclidean norm. However, depending on the function,  
109 this may scale poorly with dimension. Mirror descent circumvents this by allowing for this  
110 Lipschitz constant to be taken with respect to other norms such as the  $\ell^1$  norm. This has  
111 been shown to scale better with dimension compared to methods such as projected subgradient  
112 descent on problems including online learning and tomography [9, 25, 6]. Further work has  
113 been done by Gunasekar et al., showing that MD is equivalent to natural/geodesic gradient  
114 descent on certain Riemannian manifolds [14]. We will continue in the simpler setting where  
115 we have a potential given by a strictly convex  $\Psi$  to aid parameterization, but this can be  
116 replaced by a suitable Hessian metric tensor.

117 Let  $\mathcal{X} \subset \mathbb{R}^n$  be a closed convex set with nonempty interior. Let  $(\mathbb{R}^n)^*$  denote the corre-  
118 sponding dual space of  $\mathbb{R}^n$ .

119 **Definition 2.1 (Mirror Map).** We say  $\Psi : \mathcal{X} \rightarrow \mathbb{R}$  is a *mirror potential* if it is continuously  
120 differentiable and strongly convex. We call the gradient  $\nabla \Psi : \mathcal{X} \rightarrow (\mathbb{R}^n)^*$  a *mirror map*.

121 **Remark 2.2.** A mirror potential  $\Psi$  may also be referred to as a *distance generating func-*

122 tion, as a convex map induces a Bregman distance  $B_\Psi(x, y)$ , defined by  $B_\Psi(x, y) = \Psi(x) -$   
 123  $\Psi(y) - \langle \nabla \Psi(y), x - y \rangle$ . For example, taking  $\Psi(x) = \|x\|_2^2$  recovers the usual squared Euclidean  
 124 distance  $B_\Psi(x, y) = \|x - y\|_2^2$ .

125 If  $\Psi$  is a mirror potential, then the convex conjugate  $\Psi^*$  defined as

$$126 \quad \Psi^*(x^*) = \sup_{x \in \mathcal{X}} \{ \langle x^*, x \rangle - \Psi(x) \}$$

127 is differentiable everywhere, and additionally satisfies  $\nabla \Psi^* = (\nabla \Psi)^{-1}$  [5, 28]. The (forward)  
 128 mirror map  $\nabla \Psi$  *mirrors* from the primal space  $\mathcal{X}$  into a subset of the dual space  $(\mathbb{R}^n)^*$ , and  
 129 the inverse (backward) mirror map  $\nabla \Psi^*$  *mirrors* from the dual space  $\text{dom}(\nabla \Psi^*) \subseteq (\mathbb{R}^n)^*$   
 130 back into the primal space  $\mathcal{X}$ .

131 Suppose first that we are trying to minimize a convex differentiable function  $f$  over the  
 132 entire space  $\mathcal{X} = \mathbb{R}^n$ ,  $\min_{x \in \mathcal{X}} f(x)$ . Suppose further for simplicity that  $\text{dom}(\nabla \Psi^*) = (\mathbb{R}^n)^*$ .  
 133 For an initial point  $x_0 \in \mathcal{X}$  and a sequence of step-sizes  $(t_k)_{k \geq 0}$ ,  $t_k > 0$ , the mirror descent  
 134 iterations can be written as follows:

$$135 \quad (2.1) \quad y_k = \nabla \Psi(x_k) - t_k \nabla f(x_k), \quad x_{k+1} = \nabla \Psi^*(y_k).$$

136 There are two main sequences,  $(x_k)_{k=0}^\infty$  in the primal space  $\mathcal{X}$  and  $(y_k)_{k=0}^\infty$  in the dual space  
 137  $(\mathbb{R}^n)^*$ . The gradient step at each iteration is performed in the dual space, with the mirror  
 138 map  $\nabla \Psi$  mapping between them. Observe that if  $\Psi = \frac{1}{2}\|x\|_2^2$ , then  $\nabla \Psi$  is the identity  
 139 map  $\mathbb{R}^n \rightarrow (\mathbb{R}^n)^*$  and we recover the standard gradient descent algorithm. An equivalent  
 140 formulation of the MD update rule in (2.1) is the subgradient algorithm [5]:

$$141 \quad (2.2) \quad x_{k+1} = \arg \min_{x \in \mathcal{X}} \left\{ \langle x, \nabla f(x_k) \rangle + \frac{1}{t_k} B_\Psi(x, x_k) \right\}.$$

142 This can be derived by using the definitions of the Bregman distance and of the convex  
 143 conjugate  $\Psi^*$ . The convexity of  $\Psi$  implies that the induced Bregman divergence  $B_\Psi$  is non-  
 144 negative, which allows for this iteration to be defined. Observe again that if  $\Psi = \frac{1}{2}\|x\|_2^2$ , then  
 145  $B_\Psi(x, y) = \frac{1}{2}\|x - y\|_2^2$  and we recover the argmin formulation of the gradient descent update  
 146 rule.

147 MD enjoys the following convergence rate guarantees. Let  $\|\cdot\|$  be a norm on  $\mathbb{R}^n$ , and  
 148  $\|\cdot\|_* = \max\{\langle \cdot, x \rangle : x \in \mathbb{R}^n, \|x\| \leq 1\}$  be the corresponding dual norm. For a set  $\mathcal{X} \subseteq \mathbb{R}^n$ , let  
 149  $\text{int}(\mathcal{X})$  denote the interior of  $\mathcal{X}$ .

150 **Theorem 2.3.** [5, Thm 4.1] *Let  $\mathcal{X}$  be a closed convex subset of  $\mathbb{R}^n$  with nonempty interior,*  
 151 *and  $f : \mathcal{X} \rightarrow \mathbb{R}$  a convex function. Suppose that  $\Psi$  is a  $\sigma$ -strongly convex mirror potential.*  
 152 *Suppose further that the following hold:*

- 153 1.  $f$  is Lipschitz with Lipschitz constant  $L_f$  with respect to  $\|\cdot\|$ ;
- 154 2. The set of minimizers  $\min_{x \in \mathcal{X}} f(x)$  is nonempty; let  $x^*$  be a minimizer of  $f$ .

155 Let  $\{x_k\}_{k=1}^\infty$  be the sequence generated by the MD iterations (2.1) with starting point  $x_1 \in$   
 156  $\text{int}(\mathcal{X})$ . Then the iterates satisfy the following regret bound:

$$157 \quad (2.3) \quad \sum_{k=1}^s t_k (f(x_k) - f(x^*)) \leq B_\Psi(x^*, x_1) - B_\Psi(x^*, x_{s+1}) + (2\sigma)^{-1} \sum_{k=1}^s t_k^2 \|\nabla f(x_k)\|_*^2.$$

158 *In particular, we have*

$$159 \quad (2.4) \quad \min_{1 \leq k \leq s} f(x_k) - f(x^*) \leq \frac{B_{\Psi}(x^*, x_1) + (2\sigma)^{-1} \sum_{k=1}^s t_k^2 \|\nabla f(x_k)\|_*^2}{\sum_{k=1}^s t_k}.$$

160 *Remark 2.4.* The proof of [Theorem 2.3](#) depends only on the property that  $\nabla\Psi^* = (\nabla\Psi)^{-1}$ .  
 161 Therefore, the inverse mirror map  $(\nabla\Psi^*)$  as in [\(2.1\)](#) can be replaced with  $(\nabla\Psi)^{-1}$ , yielding a  
 162 formulation of MD that does not reference the convex conjugate  $\Psi^*$  of the mirror potential  $\Psi$   
 163 itself, but only the gradient  $\nabla\Psi^*$ .

164 To motivate our goal of learning mirror maps, we will demonstrate an application of MD  
 165 that drastically speeds up convergence over gradient descent. We consider optimization on  
 166 the simplex  $\Delta_d = \{x \in \mathbb{R}^d : x \geq 0, \sum_j x_j = 1\}$ , equipped with a mirror potential given by  
 167 the (negative log-) entropy map [\[5\]](#). We have the following mirror maps, where logarithms  
 168 and exponentials of vectors are to be taken component-wise:

$$169 \quad (2.5) \quad \Psi(x) = \sum_j x_j \log x_j, \quad \nabla\Psi(x) = 1 + \log(x), \quad \nabla\Psi^*(y) = \frac{\exp(y)}{\sum_j \exp(y_j)}.$$

170 This results in the *entropic mirror descent* algorithm. It can be shown to have similar conver-  
 171 gence rates as projected subgradient descent, with a  $O(1/\sqrt{k})$  convergence rate [\[5, Thm 5.1\]](#).  
 172 Given that the optimization is over a probability simplex, a natural problem class to consider  
 173 is a probabilistic distance between points, given by the KL divergence.

174 Minimizing the KL divergence is a convex problem on the simplex  $x \in \Delta_d$ . For a point  
 175  $y \in \Delta_d$ , the KL divergence is given as follows, where  $0 \log 0$  is taken to be 0 by convention:

$$176 \quad (2.6) \quad \min_{x \in \Delta_d} KL(x||y) = \sum_{i=1}^d x_i \log \left( \frac{x_i}{y_i} \right).$$

177 To demonstrate the potential of MD, we can apply the entropic MD algorithm to the problem  
 178 classes of minimizing KL divergence and of minimizing least squares loss over the simplex  $\Delta_d$ .  
 179 The function classes that we apply the entropic MD algorithm and gradient descent to are:

$$180 \quad \mathcal{F}_{KL} = \{KL(\cdot||y) : y \in \Delta_d\}, \quad \mathcal{F}_{lsq} = \{\|\cdot - y\|_2^2 : y \in \Delta_d\},$$

181 where the functions have domain  $\Delta_d$ . Note that the true minimizers of a function in either of  
 182 these function classes is given by the parameter  $y \in \Delta_d$ .

183 To compare these two optimization algorithms, we optimize 500 functions from the respec-  
 184 tive function classes, which were generated by uniformly sampling  $y$  on the simplex. [Figure 1](#)  
 185 plots the evolution of the loss for the entropic MD algorithm and gradient descent for these  
 186 two problem classes, applied with various step-sizes. The entropic MD algorithm gives lin-  
 187 ear convergence on the KL function class  $\mathcal{F}_{KL}$ , massively outperforming the gradient descent  
 188 algorithm. However, entropic MD is unable to maintain this convergence rate over the least-  
 189 squares function class  $\mathcal{F}_{lsq}$ . [The difference in convergence rate demonstrates the importance of](#)  
 190 [choosing a suitable mirror map for the target function class, as well as the potential of MD in](#)  
 191 [accelerating convergence. This relationship between the function class and mirror maps moti-](#)  
 192 [vates a learned approach to deriving mirror maps from data to replace classical hand-crafted](#)  
 193 [mirror maps.](#)



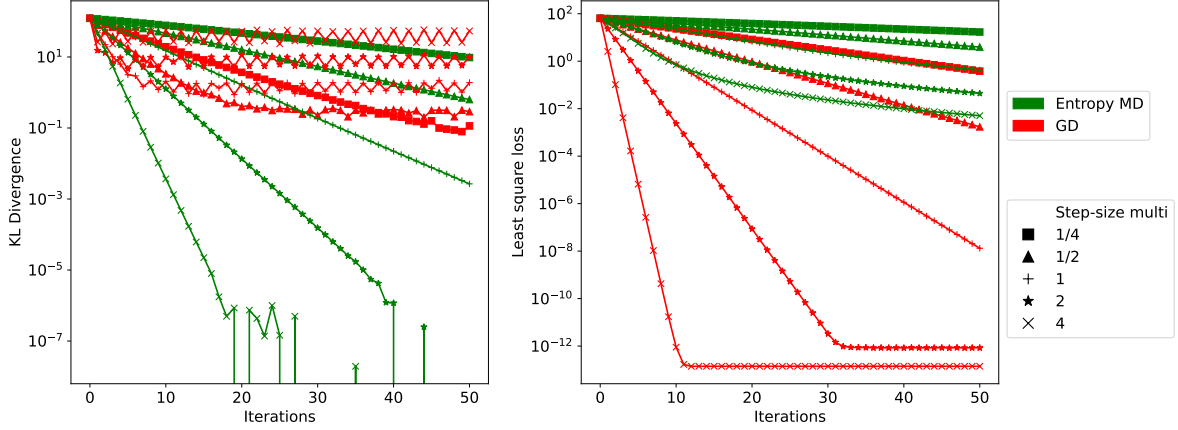


Figure 1: Effect of using the entropic MD method (2.5) to minimize KL divergence (left) and least squares loss (right). The step-sizes were taken as  $0.1 \times \text{step-size multi}$ . We can see that entropic MD (green) outperforms the gradient descent method for the KL divergence task, though loses out in the least squares task. The unstable iterations at low KL divergence are due to machine precision. The difference between the two optimization methods on each problem class demonstrates the potential of adapting to the optimization geometry using MD.

194 **3. Main Results.** We first theoretically show convergence properties of mirror descent  
 195 when the mirror map constraint  $\nabla\Psi^* = (\nabla\Psi)^{-1}$  is only approximately satisfied. Motivated  
 196 by these convergence properties, we propose our Learned Mirror Descent method, trained with  
 197 a loss function balancing empirical convergence speed and theoretical convergence guarantees.

198 We briefly explain our key objective of approximate mirror descent. Recall that MD as  
 199 given in (2.1) requires two mirror maps,  $\nabla\Psi$  and  $\nabla\Psi^*$ . We wish to parameterize both  $\Psi$  and  $\Psi^*$   
 200 using neural networks  $M_\theta$  and  $M_\theta^*$ , and weakly enforce the constraint that  $\nabla M_\theta^* = (\nabla M_\theta)^{-1}$ .  
 201 To maintain the convergence guarantees of MD, we will derive a bound on the regret depending  
 202 on the deviation between  $\nabla M_\theta^*$  and  $(\nabla M_\theta)^{-1}$  in a sense that will be made precise later. We  
 203 will call the inconsistency between the parameterized mirror maps  $\nabla M_\theta^*$  and  $(\nabla M_\theta)^{-1}$  the  
 204 *forward-backward inconsistency/loss*.

205 Recall the problem setting as in Section 2. Let  $\Psi$  be a mirror potential, i.e. a  $C^1$   $\sigma$ -  
 206 strongly-convex function with  $\sigma > 0$ . In this section, we shall work in the unconstrained case  
 207  $\mathcal{X} = \mathbb{R}^n$ . We further assume  $f$  has a minimizer  $x^* \in \mathcal{X}$ .

208 Recall the MD iteration (2.1) with step sizes  $\{t_k\}_{k=1}^\infty$  as follows. Throughout this section,  
 209  $B = B_\Psi$  is the Bregman distance with respect to  $\Psi$ , and  $\Psi^*$  is the convex conjugate of  $\Psi$ :

$$210 \quad (3.1) \quad x_{k+1} = \arg \min_{x \in \mathcal{X}} \{ \langle x, t_k \nabla f(x_k) \rangle + B(x, x_k) \} = \nabla\Psi^*(\nabla\Psi(x_k) - t_k \nabla f(x_k)).$$

211 For a general mirror map  $\Psi$ , the convex conjugate  $\Psi^*$  and the associated backward mir-  
 212 ror map  $\nabla\Psi^*$  may not have a closed form. Suppose now that we parameterize  $\Psi$  and  $\Psi^*$   
 213 with neural networks  $M_\theta$  and  $M_\theta^*$  respectively, satisfying  $\nabla M_\theta^* \approx (\nabla M_\theta)^{-1}$ . The resulting

214 *approximate mirror descent* scheme is as follows, starting from  $\tilde{x}_1 = x_1$ :

$$215 \quad (3.2) \quad \tilde{x}_{k+1} = \nabla M_{\vartheta}^*(\nabla M_{\theta}(\tilde{x}_k) - t_k \nabla f(\tilde{x}_k)), \quad k = 1, 2, \dots.$$

216 Here, we enforce that the sequence  $\{\tilde{x}_k\}$  represents an approximation of a mirror descent  
217 iteration at each step, given by

$$218 \quad (3.3) \quad x_{k+1} = \arg \min_{x \in \mathcal{X}} \{ \langle x, t_k \nabla f(\tilde{x}_k) \rangle + B(x, \tilde{x}_k) \} = (\nabla M_{\theta})^{-1}(\nabla M_{\theta}(\tilde{x}_k) - t_k \nabla f(\tilde{x}_k)).$$

219 Hereafter, we will refer to  $M_{\theta}$  and  $M_{\vartheta}^*$  as the *forward and backward (mirror) potentials*,  
220 respectively, and the corresponding gradients as the *forward and backward (mirror) maps*. For  
221 practical purposes,  $\{\tilde{x}_k\}$  should be considered as the iterations that we can compute. Typi-  
222 cally, both the argmin and  $\nabla \Psi^*$  are not easily computable, hence  $x_k$  will not be computable  
223 either. However, defining this quantity will prove useful for our analysis, as we can addition-  
224 ally use this quantity to compare how close the forward and backward maps are from being  
225 inverses of each other.

226 The following theorem puts a convergence rate bound on the approximate MD scheme  
227 (3.2) in terms of the forward-backward inconsistency. More precisely, the inconsistency is  
228 quantified by the difference of the iterates in the dual space. [This will allow us to show](#)  
229 [approximate convergence when the inverse mirror map is not known exactly](#).

230 **Theorem 3.1 (Regret Bound for Approximate MD).** *Suppose  $f$  is  $\mu$ -strongly convex with*  
231 *parameter  $\mu > 0$ , and  $\Psi$  is a mirror potential with strong convexity parameter  $\sigma$ . Let  $\{\tilde{x}_k\}_{k=0}^{\infty}$*   
232 *be some sequence in  $\mathcal{X} = \mathbb{R}^n$ , and  $\{x_k\}_{k=1}^{\infty}$  be the corresponding exact MD iterates generated*  
233 *by (3.3). We have the following regret bound:*

$$234 \quad (3.4) \quad \sum_{k=1}^K t_k (f(\tilde{x}_k) - f(x^*)) \\ \leq B(x^*, \tilde{x}_1) + \sum_{k=1}^K \left[ \frac{1}{\sigma} t_k^2 \|\nabla f(\tilde{x}_k)\|_*^2 + \left( \frac{1}{2t_k \mu} + \frac{1}{\sigma} \right) \|\nabla \Psi(\tilde{x}_{k+1}) - \nabla \Psi(x_{k+1})\|_*^2 \right].$$

235 *Proof.* We start by employing *amortization* to find an upper bound on the following ex-  
236 pression:

$$237 \quad (3.5) \quad t_k f(\tilde{x}_k) - t_k f(x^*) + (B(x^*, \tilde{x}_{k+1}) - B(x^*, \tilde{x}_k)).$$

238 From the formulation (3.2), since  $\nabla \Psi^* = (\nabla \Psi)^{-1}$ :

$$239 \quad \nabla \Psi(x_{k+1}) = \nabla \Psi(\tilde{x}_k) - t_k \nabla f(\tilde{x}_k).$$

240 We have the following bound on  $B(x^*, \tilde{x}_{k+1}) - B(x^*, \tilde{x}_k)$ :



$$\begin{aligned}
241 \quad & B(x^*, \tilde{x}_{k+1}) - B(x^*, \tilde{x}_k) = \Psi(x^*) - \Psi(\tilde{x}_{k+1}) - \langle \nabla \Psi(\tilde{x}_{k+1}), x^* - \tilde{x}_{k+1} \rangle \\
242 \quad & \quad - [\Psi(x^*) - \Psi(\tilde{x}_k) - \langle \nabla \Psi(\tilde{x}_k), x^* - \tilde{x}_k \rangle] \quad \text{[definition of } B] \\
243 \quad & = \Psi(\tilde{x}_k) - \Psi(\tilde{x}_{k+1}) - \langle \nabla \Psi(\tilde{x}_{k+1}), x^* - \tilde{x}_{k+1} \rangle + \langle \nabla \Psi(\tilde{x}_k), x^* - \tilde{x}_k \rangle \quad \text{[cancel } \Psi(x^*)] \\
244 \quad & = \Psi(\tilde{x}_k) - \Psi(\tilde{x}_{k+1}) - \langle \nabla \Psi(x_{k+1}), x^* - \tilde{x}_{k+1} \rangle \quad \text{[add/subtract} \\
245 \quad & \quad - \langle \nabla \Psi(\tilde{x}_{k+1}) - \nabla \Psi(x_{k+1}), x^* - \tilde{x}_{k+1} \rangle + \langle \nabla \Psi(\tilde{x}_k), x^* - \tilde{x}_k \rangle \quad \text{terms in blue]} \\
246 \quad & = \Psi(\tilde{x}_k) - \Psi(\tilde{x}_{k+1}) - \langle \nabla \Psi(\tilde{x}_k) - t_k \nabla f(\tilde{x}_k), x^* - \tilde{x}_{k+1} \rangle \quad \text{[MD update (3.3)} \\
247 \quad & \quad - \langle \nabla \Psi(\tilde{x}_{k+1}) - \nabla \Psi(x_{k+1}), x^* - \tilde{x}_{k+1} \rangle + \langle \nabla \Psi(\tilde{x}_k), x^* - \tilde{x}_k \rangle \quad \text{on } \nabla \Psi(x_{k+1})] \\
248 \quad & = \underbrace{\Psi(\tilde{x}_k) - \Psi(\tilde{x}_{k+1}) + \langle \nabla \Psi(\tilde{x}_k), \tilde{x}_{k+1} - \tilde{x}_k \rangle}_{-B_\Psi(\tilde{x}_{k+1}, \tilde{x}_k)} \\
249 \quad & \quad + \langle t_k \nabla f(\tilde{x}_k), x^* - \tilde{x}_{k+1} \rangle - \langle \nabla \Psi(\tilde{x}_{k+1}) - \nabla \Psi(x_{k+1}), x^* - \tilde{x}_{k+1} \rangle.
\end{aligned}$$

250 Observe that the first line in the final expression is precisely  $-B_\Psi(\tilde{x}_{k+1}, \tilde{x}_k)$ . By  $\sigma$ -strong-  
251 convexity of  $\Psi$ , we have  $-B_\Psi(\tilde{x}_{k+1}, \tilde{x}_k) \leq -\frac{\sigma}{2} \|\tilde{x}_{k+1} - \tilde{x}_k\|^2$ . Therefore, our final bound for  
252 this expression is:

$$\begin{aligned}
253 \quad (3.6) \quad & B(x^*, \tilde{x}_{k+1}) - B(x^*, \tilde{x}_k) \\
& \leq -\frac{\sigma}{2} \|\tilde{x}_{k+1} - \tilde{x}_k\|^2 + \langle t_k \nabla f(\tilde{x}_k), x^* - \tilde{x}_{k+1} \rangle - \langle \nabla \Psi(\tilde{x}_{k+1}) - \nabla \Psi(x_{k+1}), x^* - \tilde{x}_{k+1} \rangle.
\end{aligned}$$

254 Returning to bounding the initial expression (3.5), we have by substituting (3.6):

$$\begin{aligned}
255 \quad & t_k f(\tilde{x}_k) - t_k f(x^*) + (B(x^*, \tilde{x}_{k+1}) - B(x^*, \tilde{x}_k)) \\
256 \quad & \leq t_k f(\tilde{x}_k) - t_k f(x^*) + \langle t_k \nabla f(\tilde{x}_k), x^* - \tilde{x}_{k+1} \rangle \\
257 \quad & \quad - \frac{\sigma}{2} \|\tilde{x}_{k+1} - \tilde{x}_k\|^2 - \langle \nabla \Psi(\tilde{x}_{k+1}) - \nabla \Psi(x_{k+1}), x^* - \tilde{x}_{k+1} \rangle \quad \text{[by (3.6)]} \\
258 \quad & = t_k f(\tilde{x}_k) - t_k f(x^*) + \langle t_k \nabla f(\tilde{x}_k), x^* - \tilde{x}_k \rangle + \langle t_k \nabla f(\tilde{x}_k), \tilde{x}_k - \tilde{x}_{k+1} \rangle \quad \text{[add/subtract} \\
259 \quad & \quad - \frac{\sigma}{2} \|\tilde{x}_{k+1} - \tilde{x}_k\|^2 - \langle \nabla \Psi(\tilde{x}_{k+1}) - \nabla \Psi(x_{k+1}), x^* - \tilde{x}_{k+1} \rangle \quad \text{terms in blue]} \\
260 \quad & = -t_k B_f(x^*, \tilde{x}_k) + \langle t_k \nabla f(\tilde{x}_k), \tilde{x}_k - \tilde{x}_{k+1} \rangle \\
261 \quad & \quad - \frac{\sigma}{2} \|\tilde{x}_{k+1} - \tilde{x}_k\|^2 - \langle \nabla \Psi(\tilde{x}_{k+1}) - \nabla \Psi(x_{k+1}), x^* - \tilde{x}_k \rangle \\
262 \quad & \quad - \langle \nabla \Psi(\tilde{x}_{k+1}) - \nabla \Psi(x_{k+1}), \tilde{x}_k - \tilde{x}_{k+1} \rangle.
\end{aligned}$$

263 The above two equalities are obtained by writing the second term of the inner products as  
264  $x^* - \tilde{x}_{k+1} = (x^* - \tilde{x}_k) + (\tilde{x}_k - \tilde{x}_{k+1})$ , and by the definition of  $B_f$ . By  $\mu$ -strong-convexity of  
265  $f$ , we get  $-t_k B_f(x^*, \tilde{x}_k) \leq -\frac{t_k \mu}{2} \|x^* - \tilde{x}_k\|^2$ . Therefore, the bound on the quantity in (3.5)

266 reduces to

$$\begin{aligned}
& t_k f(\tilde{x}_k) - t_k f(x^*) + (B(x^*, \tilde{x}_{k+1}) - B(x^*, \tilde{x}_k)) \\
& \leq -\frac{t_k \mu}{2} \|x^* - \tilde{x}_k\|^2 + \langle t_k \nabla f(\tilde{x}_k), \tilde{x}_k - \tilde{x}_{k+1} \rangle \\
& \quad - \frac{\sigma}{2} \|\tilde{x}_{k+1} - \tilde{x}_k\|^2 - \langle \nabla \Psi(\tilde{x}_{k+1}) - \nabla \Psi(x_{k+1}), x^* - \tilde{x}_k \rangle \\
& \quad - \langle \nabla \Psi(\tilde{x}_{k+1}) - \nabla \Psi(x_{k+1}), \tilde{x}_k - \tilde{x}_{k+1} \rangle.
\end{aligned}
\tag{3.7}$$

268 Liberally applying Cauchy-Schwarz and Young's inequality to bound the inner product terms:  
269

$$\begin{aligned}
& t_k f(\tilde{x}_k) - t_k f(x^*) + (B(x^*, \tilde{x}_{k+1}) - B(x^*, \tilde{x}_k)) \\
& \leq -\frac{t_k \mu}{2} \|x^* - \tilde{x}_k\|^2 + \frac{1}{\sigma} t_k^2 \|\nabla f(\tilde{x}_k)\|_*^2 + \frac{\sigma}{4} \|\tilde{x}_k - \tilde{x}_{k+1}\|^2 \\
& \quad - \frac{\sigma}{2} \|\tilde{x}_{k+1} - \tilde{x}_k\|^2 + \frac{1}{2t_k \mu} \|\nabla \Psi(\tilde{x}_{k+1}) - \nabla \Psi(x_{k+1})\|_*^2 + \frac{t_k \mu}{2} \|x^* - \tilde{x}_k\|^2 \\
& \quad + \frac{1}{\sigma} \|\nabla \Psi(\tilde{x}_{k+1}) - \nabla \Psi(x_{k+1})\|_*^2 + \frac{\sigma}{4} \|\tilde{x}_k - \tilde{x}_{k+1}\|^2 \\
& \leq \frac{1}{\sigma} t_k^2 \|\nabla f(\tilde{x}_k)\|_*^2 + \left( \frac{1}{2t_k \mu} + \frac{1}{\sigma} \right) \|\nabla \Psi(\tilde{x}_{k+1}) - \nabla \Psi(x_{k+1})\|_*^2.
\end{aligned}
\tag{3.8}$$

271 Summing from  $k = 1$  to  $K$ , we get

$$\begin{aligned}
& \sum_{k=1}^K [t_k f(\tilde{x}_k) - t_k f(x^*) + (B(x^*, \tilde{x}_{k+1}) - B(x^*, \tilde{x}_k))] \\
& \leq \sum_{k=1}^K \left[ \frac{1}{\sigma} t_k^2 \|\nabla f(\tilde{x}_k)\|_*^2 + \left( \frac{1}{2t_k \mu} + \frac{1}{\sigma} \right) \|\nabla \Psi(\tilde{x}_{k+1}) - \nabla \Psi(x_{k+1})\|_*^2 \right].
\end{aligned}
\tag{3.9}$$

273 Observe  $\sum_{k=1}^K (B(x^*, \tilde{x}_{k+1}) - B(x^*, \tilde{x}_k)) = B(x^*, \tilde{x}_{K+1}) - B(x^*, \tilde{x}_1) \geq -B(x^*, \tilde{x}_1)$ . Apply this  
274 with (3.9) to finish the regret bound. ■

275 *Remark 3.2.* This bound may be extended to the constrained case  $\mathcal{X} \subsetneq \mathbb{R}^n$ . This can be  
276 shown by adding an extra projection step to the iterates of the form  $\pi(y) = \arg \min_{x \in \mathcal{X}} B(x, y)$ ,  
277 and having  $\tilde{x}_{k+1}$  instead approximate the projection of the exact mirror step  $\tilde{x}_{k+1} \approx \pi(x_{k+1})$   
278 in (3.2) [23]. Note that if  $y \notin \mathcal{X}$ , then  $B(x^*, \pi(y)) \leq B(x^*, y)$  for any  $x^* \in \mathcal{X}$ .

279 *Remark 3.3.* The convex function  $f$  need not be differentiable, and having a non-empty  
280 subgradient at every point is sufficient for the regret bound to hold. The proof will still work  
281 if  $\nabla f$  is replaced by a subgradient  $f' \in \partial f$ .

282 *Remark 3.4.* Observe there is a  $t_k^{-1}$  coefficient in the approximation term. This prevents  
283 us from taking  $t_k \searrow 0$  to get convergence as in the classical MD case. Intuitively, a suffi-  
284 ciently large gradient step is required to correct for the approximation. However, due to the  
285 Lipschitz condition on the objective  $f$ , the gradient step is still required to be limited above  
286 for convergence.

287 With [Theorem 3.1](#), we no longer require precise knowledge of the convex conjugate. In  
 288 particular, this allows us to parameterize the forward mirror potential with an ICNN, for  
 289 which there is no closed-form convex conjugate in general. We are thus able to approximate  
 290 the backwards mirror potential with another neural network, while maintaining approximate  
 291 convergence guarantees. While the true backward potential will be convex, these results allow  
 292 us to use a non-convex network, resulting in better numerical performance.

293 **3.1. Relative Smoothness Assumption.** We have seen that we can approximate the iter-  
 294 ations of MD and still obtain convergence guarantees. With the slightly weaker assumption of  
 295 relative smoothness and relative strong convexity, MD can be shown to converge [\[20\]](#). We can  
 296 get a similar and cleaner bound by slightly modifying the proof of convergence for classical  
 297 MD under these new assumptions.

298 **Definition 3.5 (Relative Smoothness/Convexity).** Let  $\Psi : \mathcal{X} \rightarrow \mathbb{R}$  be a differentiable con-  
 299 vex function, defined on a convex set  $\mathcal{X}$  (with non-empty interior), which will be used as a  
 300 reference. Let  $f : \mathcal{X} \rightarrow \mathbb{R}$  be another differentiable convex function.

301  $f$  is  $L$ -smooth relative to  $\Psi$  if for any  $x, y \in \text{int}(\mathcal{X})$ ,

$$302 \quad (3.10) \quad f(y) \leq f(x) + \langle \nabla f(x), y - x \rangle + LB_{\Psi}(y, x).$$

303  $f$  is  $\mu$ -strongly-convex relative to  $\Psi$  if for any  $x, y \in \text{int}(\mathcal{X})$ ,

$$304 \quad (3.11) \quad f(y) \geq f(x) + \langle \nabla f(x), y - x \rangle + \mu B_{\Psi}(y, x).$$

305 Observe that these definitions of relative smoothness and relative strong convexity extend  
 306 the usual notions of  $L$ -smoothness and strong convexity with the Euclidean norm by taking  
 307  $\Psi = \frac{1}{2}\|\cdot\|_2^2$ , recovering  $B_{\Psi}(x, y) = \frac{1}{2}\|x - y\|_2^2$ . Moreover, if  $\nabla f$  is  $L$ -Lipschitz and  $\Psi$  is  $\mu$ -  
 308 strongly convex with  $\mu > 0$ , then  $f$  is  $L/\mu$  smooth relative to  $\Psi$ . If both functions are  
 309 twice-differentiable, the above definitions are equivalent to the following [\[20, Prop 1.1\]](#):

$$310 \quad (3.12) \quad \mu \nabla^2 \Psi \preceq \nabla^2 f \preceq L \nabla^2 \Psi.$$

311 Using the relative smoothness and relative strong convexity conditions, we can show con-  
 312 vergence even when the convex objective function  $f$  is flat, as long as our mirror potential  $\Psi$   
 313 is also flat at those points. The analysis given in [\[20\]](#) readily extends to the case where our  
 314 iterations are approximate.

315 **Theorem 3.6.** Let  $f$  be relatively  $L$ -smooth and relatively  $\mu$ -strongly-convex with respect to  
 316 the mirror map  $\Psi$ , with  $L > 0, \mu \geq 0$ . Let  $\{\tilde{x}_k\}_{k \geq 0}$  be a sequence in  $\mathcal{X}$ , and consider the  
 317 iterations  $\{x_k\}_{k \geq 1}$  defined as

$$318 \quad (3.13) \quad x_{k+1} = \arg \min_{x \in \mathcal{X}} \{ \langle x, \nabla f(\tilde{x}_k) \rangle + LB(x, \tilde{x}_k) \},$$

319 i.e. the result of applying a single MD update step with fixed step size  $1/L$  to each  $\tilde{x}_k$ . We  
 320 have the following bound (for any  $x \in \mathcal{X}$ ), where the middle expression is discarded if  $\mu = 0$ :

$$321 \quad (3.14) \quad \min_{1 \leq i \leq k} f(\tilde{x}_i) - f(x) \leq \frac{\mu B(x, \tilde{x}_0)}{(1 + \frac{\mu}{L-\mu})^k - 1} + M_k \leq \frac{L - \mu}{k} B(x, \tilde{x}_0) + M_k,$$

322 where

$$323 \quad (3.15) \quad M_k = \frac{\sum_{i=1}^k \left(\frac{L}{L-\mu}\right)^i [L\langle \nabla \Psi(x_i) - \nabla \Psi(\tilde{x}_i), x - \tilde{x}_i \rangle + \langle \nabla f(x_i), \tilde{x}_i - x_i \rangle]}{\sum_{i=1}^k \left(\frac{L}{L-\mu}\right)^i}.$$

324 In particular, if  $L\langle \nabla \Psi(x_i) - \nabla \Psi(\tilde{x}_i), x - \tilde{x}_i \rangle + \langle \nabla f(x_i), \tilde{x}_i - x_i \rangle$  is uniformly bounded (from  
325 above) by  $M$ , we can replace  $M_k$  by  $M$  in (3.14).

326 *Proof.* We follow the proof of [20, Thm 3.1] very closely. We state first the three-point  
327 property ([20, Lemma 3.1], [30]).

328 **Lemma 3.7 (Three-point property).** *Let  $\phi(x)$  be a proper l.s.c. convex function. If*

$$329 \quad z_+ = \arg \min_x \{\phi(x) + B(x, z)\},$$

330 then

$$331 \quad \phi(x) + B(x, z) \geq \phi(z_+) + B(z_+, z) + B(x, z_+), \quad \text{for all } x \in \mathcal{X}.$$

332 As in [20, Eq 28], we have for any  $x \in \mathcal{X}$  and  $i \geq 1$ ,

$$333 \quad (3.16) \quad \begin{aligned} f(x_i) &\leq f(\tilde{x}_{i-1}) + \langle \nabla f(\tilde{x}_{i-1}), x_i - \tilde{x}_{i-1} \rangle + LB(x_i, \tilde{x}_{i-1}) \\ &\leq f(\tilde{x}_{i-1}) + \langle \nabla f(\tilde{x}_{i-1}), x - \tilde{x}_{i-1} \rangle + LB(x, \tilde{x}_{i-1}) - LB(x, x_i) \\ &\leq f(x) + (L - \mu)B(x, \tilde{x}_{i-1}) - LB(x, x_i). \end{aligned}$$

334 The first inequality follows from  $L$ -smoothness relative to  $\Psi$ , the second inequality from the  
335 three-point property applied to  $\phi(x) = \frac{1}{L}\langle \nabla f(\tilde{x}_{i-1}), x - \tilde{x}_{i-1} \rangle$  and  $z = \tilde{x}_{i-1}$ ,  $z_+ = x_i$ , and the  
336 last inequality from  $\mu$ -strong-convexity of  $f$  relative to  $\Psi$ . We thus have

$$337 \quad (3.17) \quad \begin{aligned} f(\tilde{x}_i) &= f(x_i) + f(\tilde{x}_i) - f(x_i) \\ &\leq f(x) + (L - \mu)B(x, \tilde{x}_{i-1}) - LB(x, x_i) + f(\tilde{x}_i) - f(x_i) \\ &= (L - \mu)B(x, \tilde{x}_{i-1}) - LB(x, \tilde{x}_i) \\ &\quad + [f(x) + LB(x, \tilde{x}_i) - LB(x, x_i) + f(\tilde{x}_i) - f(x_i)]. \end{aligned}$$

338 By induction/telescoping, we get:

$$339 \quad (3.18) \quad \begin{aligned} \sum_{i=1}^k \left(\frac{L}{L-\mu}\right)^i f(\tilde{x}_i) &\leq LB(x, \tilde{x}_0) + \sum_{i=1}^k \left(\frac{L}{L-\mu}\right)^i f(x) \\ &\quad + \sum_{i=1}^k \left(\frac{L}{L-\mu}\right)^i [L(B(x, \tilde{x}_i) - B(x, x_i)) + f(\tilde{x}_i) - f(x_i)]. \end{aligned}$$

340 The final ‘‘approximation error’’ term is

$$341 \quad (3.19) \quad \begin{aligned} &L(B(x, \tilde{x}_i) - B(x, x_i)) + f(\tilde{x}_i) - f(x_i) \\ &= L\langle \nabla \Psi(x_i) - \nabla \Psi(\tilde{x}_i), x - \tilde{x}_i \rangle - LB(\tilde{x}_i, x_i) + f(\tilde{x}_i) - f(x_i) \\ &\leq L\langle \nabla \Psi(x_i) - \nabla \Psi(\tilde{x}_i), x - \tilde{x}_i \rangle - B_f(\tilde{x}_i, x_i) + f(\tilde{x}_i) - f(x_i) \\ &= L\langle \nabla \Psi(x_i) - \nabla \Psi(\tilde{x}_i), x - \tilde{x}_i \rangle + \langle \nabla f(x_i), \tilde{x}_i - x_i \rangle, \end{aligned}$$

342 where in the inequality, we use the definition of  $L$ -relative smoothness  $B_f(x, y) \leq LB_\Psi(x, y)$ .  
 343 (Recall  $B(c, a) + B(a, b) - B(c, b) = \langle \nabla \Psi(b) - \nabla \Psi(a), c - a \rangle$  [5, Lemma 4.1].)

344 Substituting  $C_k$  defined by

$$345 \quad \sum_{i=1}^k \left( \frac{L}{L-\mu} \right)^i =: \frac{1}{C_k}$$

346 and rearranging, we get

$$347 \quad (3.20) \quad \begin{aligned} & \min_{1 \leq i \leq k} f(\tilde{x}_i) - f(x) \leq C_k LB(x, \tilde{x}_0) \\ & + C_k \sum_{i=1}^k \left( \frac{L}{L-\mu} \right)^i [L \langle \nabla \Psi(x_i) - \nabla \Psi(\tilde{x}_i), x - \tilde{x}_i \rangle + \langle \nabla f(x_i), \tilde{x}_i - x_i \rangle]. \end{aligned}$$

348 In particular, if we have a uniform bound on  $[L \langle \nabla \Psi(x_i) - \nabla \Psi(\tilde{x}_i), x - \tilde{x}_i \rangle + \langle \nabla f(x_i), \tilde{x}_i - x_i \rangle]$ ,  
 349 say  $M$ , then we have

$$350 \quad (3.21) \quad \min_{1 \leq i \leq k} f(\tilde{x}_i) - f(x) \leq C_k LB(x, \tilde{x}_0) + M.$$

351 Finally, note that if  $\mu = 0$  then  $C_k = 1/k$ , and if  $\mu > 0$  then

$$352 \quad C_k = \left( \sum_{i=1}^k \left( \frac{L}{L-\mu} \right)^i \right)^{-1} = \frac{\mu}{L \left( (1 + \frac{\mu}{L-\mu})^k - 1 \right)} \leq 1/k. \quad \blacksquare$$

353 **Theorem 3.6** gives us convergence rate bounds up to an additive approximation error  
 354  $M_k$ , depending on how far the approximate iterates  $\tilde{x}_k$  are from the true MD iterates  $x_k$ .  
 355 By taking  $x$  in (3.14) to be an optimal point  $x^*$  where  $f$  attains its minimum, we can get  
 356 approximate linear convergence and approximate  $O(1/k)$  convergence if the relative strong  
 357 convexity parameters satisfy  $\mu > 0$  and  $\mu = 0$  respectively. In particular, the quantity

$$358 \quad (3.22) \quad L \langle \nabla \Psi(x_i) - \nabla \Psi(\tilde{x}_i), x - \tilde{x}_i \rangle + \langle \nabla f(x_i), \tilde{x}_i - x_i \rangle$$

359 that we would like to bound gives an interpretation in terms of how the approximate iterates  
 360  $\tilde{x}_i$  should be close to  $x_i$ . To minimize the first term,  $\nabla \Psi(x_i) - \nabla \Psi(\tilde{x}_i)$  should be small, and  
 361  $\tilde{x}_i - x_i$  should be small to minimize the second term.

362 **3.2. Training Procedure.** In this section, we will outline our general training procedure  
 363 and further detail our definitions for having faster convergence. We further propose a loss  
 364 function to train the mirror potentials  $M_\theta$  and  $M_\theta^*$  to enforce both faster convergence, as well  
 365 as forward-backward consistency in order to apply Theorem 3.1.

366 Suppose we have a fixed function class  $\mathcal{F}$  consisting of convex functions  $f : \mathcal{X} \rightarrow \mathbb{R}$ , where  
 367  $\mathcal{X} \subseteq \mathbb{R}^d$  is some convex set that we wish to optimize over. Our goal is to efficiently minimize  
 368 typical functions in  $\mathcal{F}$  by using our learned mirror descent scheme.

369 For a function  $f \in \mathcal{F}$ , suppose we have data initializations  $x \in \mathcal{X}$  drawn from a data dis-  
 370 tribution  $\mathbb{P}_{x|f}$ , possibly depending on our function. Let  $\{\tilde{x}_k\}_{k=1}^K$  be the sequence constructed

371 by applying learned mirror descent with forward potential  $M_\theta$  and backward potential  $M_\vartheta^*$ ,  
 372 with initialization  $\tilde{x}_0 = x$ :

$$373 \quad (3.23) \quad \tilde{x}_{k+1} = \nabla M_\vartheta^*(\nabla M_\theta(\tilde{x}_k) - t_k \nabla f(\tilde{x}_k)).$$

374 To parameterize our mirror potentials  $M_\theta, M_\vartheta^* : \mathbb{R}^d \rightarrow \mathbb{R}$ , we use the architecture proposed  
 375 by Amos et al. for an input convex neural network (ICNN) [1]. The input convex neural  
 376 networks are of the following form:

$$377 \quad (3.24) \quad z_{i+1} = \sigma \left( W_i^{(z)} z_i + W_i^{(x)} x + b_i \right), \quad M(x; \theta) = z_l,$$

378 where  $\sigma$  is the leaky-ReLU activation function, and  $\theta = \{W_{0:l-1}^{(x)}, W_{1:l-1}^{(z)}, b_{0:l-1}\}$  are the pa-  
 379 rameters of the network. For the forward mirror potential  $M_\theta$ , we clip the weights such that  
 380  $W_i^{(z)}$  are non-negative, so the network is convex in its input  $x$  [1, Prop 1]. This can be done  
 381 for both fully connected and convolutional layers. We note that it is not necessary for the  
 382 backwards mirror potential  $M_\vartheta^*$  to be convex, which allows for more expressivity. Using the  
 383 ICNN architecture allows for guaranteed convex mirror potentials with minimal computa-  
 384 tional overhead. By adding an additional small quadratic term  $\mu \|x\|^2$  to the ICNN, we are  
 385 able to enforce strong convexity of the mirror map as well.

386 We would like to enforce that  $f(\tilde{x}_k)$  is minimized quickly *on average*, over both the function  
 387 class and the distribution of initializations  $\tilde{x}_0 = x$  corresponding to each individual  $f$ . One  
 388 possible method is to consider the value of the loss function at or up to a particular iteration  
 389  $\tilde{x}_N$  for fixed  $N$ . We also apply a soft penalty such that  $\nabla M_\vartheta^* \approx (\nabla M_\theta)^{-1}$  in order to maintain  
 390 reasonable convergence guarantees. The loss that we would hence like to optimize over the  
 391 neural network parameter space  $(\theta, \vartheta) \in \Theta$  is thus:

$$392 \quad (3.25) \quad \arg \min_{\theta, \vartheta} \mathbb{E}_{f,x} [f(\tilde{x}_N)] + \mathbb{E}_{\mathcal{X}} [\|\nabla M_\vartheta^* \circ \nabla M_\theta - I\|].$$

393 The expectations on the first term are taken over the function class, and further on the initial-  
 394 ization distribution conditioned on our function instance. To empirically speed up training,  
 395 we find it effective to track the loss at each stage, similar to Andrychowicz et al. [2]. Moreover,  
 396 it is impractical to have a consistency loss for the entire space  $\mathcal{X}$ , so we instead limit it to  
 397 around the samples that are attained. The loss functions that we use will be variants of the  
 398 following:

$$399 \quad (3.26a) \quad \tilde{x}_{k+1} = \nabla M_\vartheta^*(\nabla M_\theta(\tilde{x}_k) - t_k \nabla f(\tilde{x}_k)),$$

$$400 \quad (3.26b) \quad L(\theta, \vartheta) = \mathbb{E}_{f,x} \left[ \sum_{k=1}^N r_k f(\tilde{x}_k) + s_k \|\nabla M_\vartheta^* \circ \nabla M_\theta - I\|(\tilde{x}_k) \right],$$

401 where  $r_k, s_k \geq 0$  are some arbitrary weights. For training purposes, we took  $r_k = r = 1$  as  
 402 constant throughout, and varied  $s_k = s_{\text{epoch}}$  to increase as training progresses. In particular,  
 403 we will take  $s_0 = 1$ , and increase the value every 50 epochs by a factor of 1.05. **To train our**  
 404 **mirror maps, we use a Monte Carlo average of (3.26b) over realizations of  $f$  and initializations**  
 405  **$x_0$  derived from the training data. This empirical average is optimized using the Adam**



406 optimizer for the network parameters  $\theta, \vartheta$ . This can be written as follows for a minibatch  
 407  $\{f^{(i)}, x_0^{(i)}\}_{i=1}^B$  of size  $B$ :

$$408 \quad (3.27) \quad \tilde{L}(\theta, \vartheta) = \frac{1}{B} \sum_{i=1}^B \left[ \sum_{k=1}^N r_k f^{(i)}(\tilde{x}_k^{(i)}) + s_k \|(\nabla M_{\vartheta}^* \circ \nabla M_{\theta} - I)(\tilde{x}_k^{(i)})\| \right].$$

409

410 The maximum training iteration was taken to be  $N = 10$ , which provided better gen-  
 411 eralization to further iterations than for smaller  $N$ . While  $N$  could be taken to be larger,  
 412 this comes at higher computational cost due to the number of MD iterates that need to be  
 413 computed. We found that endowing  $\mathcal{X} = \mathbb{R}^d$  with the  $L^1$  norm was more effective than using  
 414 the Euclidean  $L^2$  norm. The aforementioned convergence results can be then applied with  
 415 respect to the dual norm  $\|\cdot\|_* = \|\cdot\|_{\infty}$ .

416 We additionally find it useful to allow the step-sizes to vary over each iteration, rather  
 417 than being fixed. We will refer to the procedure where we additionally learn the step-sizes  
 418 as *adaptive LMD*. The learned step-sizes have to be clipped to a fixed interval to maintain  
 419 convergence and prevent instability. The LMD mirror maps are trained under this ‘adaptive’  
 420 setting, and we will have a choice between using the learned step-sizes and using fixed step-  
 421 sizes when applying LMD on test data. For testing, we will plot the methods applied with  
 422 multiple step-sizes. These step-sizes are chosen relative to a ‘base step-size’, which is then  
 423 multiplied by a ‘step-size multiplier’, denoted as ‘step-size multi’ in subsequent figures.

424 Training of LMD amounts to training the mirror potentials and applying the approximate  
 425 mirror descent algorithm. For training, target functions are sampled from a training set,  
 426 for which the loss (3.26b) is minimized over the mirror potential parameters  $\theta$  and  $\vartheta$ . After  
 427 training, testing can be done by applying the approximate mirror descent algorithm (3.2)  
 428 directly with the learned mirror maps, requiring only forward passes through the networks.  
 429 This allows for efficient forward passes with fixed memory cost, as extra iterates and back-  
 430 propagation are not required.

431 All implementations were done in PyTorch, and training was done on Quadro RTX 6000  
 432 GPUs with 24GB of memory [26]. The code for our experiments are publicly available<sup>1</sup>.

433 **4. Learned Mirror Maps With Closed-Form Inverses.** We illustrate the potential use of  
 434 LMD by learning simple mirror maps with closed-form backward maps, and how this can lead  
 435 to faster convergence rates on certain problems. We demonstrate these maps on two convex  
 436 problems: solving unconstrained least squares, and training an SVM on 50 features. We first  
 437 mention two functional mirror maps that can be parameterized using neural networks, and  
 438 describe the training setup in this scenario.

439 One possible parameterization of the mirror potential is using a quadratic form. This can  
 440 be interpreted as gradient descent, with a multiplier in front of the gradient step. The mirror  
 441 potentials and mirror maps are given as follows, where  $x \in \mathbb{R}^d$  and  $A \in \mathbb{R}^{d \times d}$ :

$$442 \quad (4.1) \quad \Psi(x) = \frac{1}{2} x^{\top} A x, \quad \nabla \Psi(x) = \left( \frac{1}{2} A + \frac{1}{2} A^{\top} \right) x, \quad \nabla \Psi^*(y) = \left[ \frac{1}{2} A + \frac{1}{2} A^{\top} \right]^{-1} y.$$

---

<sup>1</sup><https://github.com/hyt35/icnn-md>

443 The weight matrix  $A$  was initialized as  $A = I + E$ , where  $I$  is the identity matrix and  
 444  $E$  is a diagonal matrix with random  $N(0, 0.001)$  entries. For  $\Psi$  to be strictly convex, the  
 445 symmetrization  $(A + A^\top)/2$  needs to be positive definite. With this initialization of  $A$ , we  
 446 numerically found in our example that explicitly enforcing this non-negativity constraint was  
 447 not necessary, as the weight matrices  $A$  automatically satisfied this condition after training.

448 Another simple parameterization of the mirror potential is in the form of a neural network  
 449 with one hidden layer. In particular, we will consider the case where our activation function  
 450 is a smooth approximation to leaky-ReLU, given by  $g(t) := \alpha t + (1 - \alpha) \log(1 + \exp(t))$ . Here,  
 451 the binary operator  $\odot$  for two similarly shaped matrices/vectors is the Hadamard product,  
 452 defined by component-wise multiplication  $(x \odot y)_i = x_i y_i$ . Operations such as reciprocals,  
 453 logarithms, exponentials and division applied to vectors are to be taken component-wise. For  
 454  $x \in \mathbb{R}^d, A \in \mathbb{R}^{d \times d}, w \in \mathbb{R}_+^d$ , the maps are given as follows:

$$455 \quad (4.2a) \quad \Psi(x) = w^\top g(Ax) = w^\top (\alpha Ax + (1 - \alpha) \log(1 + \exp(Ax))),$$

$$456 \quad (4.2b) \quad \nabla \Psi(x) = \alpha A^\top w + (1 - \alpha) w \odot \frac{\exp(Ax)}{1 + \exp(Ax)},$$

$$457 \quad (4.2c) \quad \nabla \Psi^*(y) = A^{-1} \log \left( \frac{(1 - \alpha)^{-1} w^{-1} \odot (y - \alpha A^\top w)}{1 - (1 - \alpha)^{-1} w^{-1} \odot (y - \alpha A^\top w)} \right).$$

458 This is quite a restrictive model for mirror descent, as it requires the perturbed dual vector  
 459  $(1 - \alpha)^{-1} w^{-1} \odot (y - \alpha A^\top w) - \eta \nabla f$  to lie component-wise in  $(0, 1)$  in order for the backward  
 460 mirror map to make sense. Nevertheless, this can be achieved by clipping the resulting gradient  
 461 value to an appropriate interval inside  $(0, 1)$ .

462 The negative slope parameter was taken to be  $\alpha = 0.2$ . The weight matrix  $A$  was initialized  
 463 as the identity matrix with entry-wise additive Gaussian noise  $N(0, 0.01)$ , and the vector  $w$   
 464 was initialized entry-wise using a uniform distribution  $\text{Unif}(0, 1/d)$ .

465 **4.1. Least Squares.** The first problem class we wish to consider is that of least squares in  
 466 two dimensions. This was done with the following fixed weight matrix and randomized bias  
 467 vectors:

$$468 \quad (4.3) \quad \min_{x \in \mathbb{R}^2} \|Wx - b\|_2^2, \quad W = \begin{pmatrix} 2 & 1 \\ 1 & 2 \end{pmatrix}, \quad b \in \mathbb{R}^2.$$

469 For training LMD for least squares, the initialization vectors  $x$  and target bias vectors  $b$   
 470 were independently randomly sampled as Gaussian vectors  $b, x \sim N(0, I_2)$ . The function class  
 471 that we wish to optimize over in (3.25) is:

$$472 \quad \mathcal{F} = \{f_b(x) = \|Wx - b\|_2^2 : b \in \mathbb{R}^2\}, \quad x \sim \mathbb{P}_{x|f} = N(0, I_2),$$

473 where the expectation  $\mathbb{E}_{f,x}$  is taken over  $b, x \sim N(0, I_2)$ .

474 For this problem class, a classical MD algorithm is available. Observe that  $\nabla f_b(x) =$   
 475  $W^\top W(x - W^{-1}b)$ . By taking  $\Psi(x) = \frac{1}{2} x^\top (W^\top W)x$ , the mirror maps are  $\nabla \Psi(x) = (W^\top W)x$ ,  
 476  $\nabla \Psi^*(x) = (W^\top W)^{-1}x$ . The MD update step (2.1) applied to  $f = f_b$  becomes

$$477 \quad (4.4) \quad x_{k+1} = (W^\top W)^{-1}((W^\top W)x_k - t_k \nabla f(x_k)) = x_k - t_k(x_k - W^{-1}b).$$

478 This update step will always point directly towards the true minimizer  $W^{-1}b$ , attaining lin-  
 479 ear convergence with appropriate step-size. In [Figures 2a](#) and [3a](#), this method is added for  
 480 comparison as the “MD” method.

481 We can observe this effect graphically in [Figure 2b](#). This figure illustrates the effect of MD  
 482 on changing the optimization path from a curve for GD, to a straight line for MD. Without  
 483 loss of generality, suppose we take  $b = 0$  and work in the eigenbasis  $\{v_1, v_2\}$  of  $W$ , so the  
 484 function to minimize becomes  $f(x_1, x_2) = 9x_1^2 + x_2^2$ . From initialization  $u = (u_1, u_2)$ , the  
 485 gradient flow induces the curve  $\gamma(t) = (u_1 \exp(-9t), u_2 \exp(-t))$ . The curvature restricts the  
 486 step-size allowed for gradient descent and moreover increases the curve length compared to  
 487 the straight MD line, leading to slower convergence.

488 An alternative perspective is given using the mirror potential  $\Psi$  in [Figure 2c](#), which takes  
 489 the shape of an elliptic paraboloid. In the eigenbasis  $v_1 = (1, 1)$ ,  $v_2 = (1, -1)$  of  $W$ , the  
 490 greater curvature of  $\Psi$  in the  $v_1$  direction implies that gradients are shrunk in this direction  
 491 in the MD step. In this case, the gradient is shrunk 9 times more in the  $v_1$  direction than the  
 492  $v_2$  direction, inducing the MD curve  $\mu(t) = (u_1 \exp(-t), u_2 \exp(-t))$ , which is a straight line.

493 [Figure 2](#) and [Figure 3](#) illustrate the results of training LMD using the quadratic mirror  
 494 potential and with the one-layer NN potential respectively. These figures include the evolution  
 495 of the loss function, the iterates after 10 iterations of adaptive LMD as in the training setting,  
 496 and a visualization of the mirror map. [Figure 3b](#) shows the instabilities that occur when the  
 497 domain of the backwards map is restricted. This is an example of a problem where applying  
 498 LMD with a well-parameterized mirror map can result in significantly accelerated convergence.

499 **4.2. SVM.** The second problem class is of training an SVM on the 4 and 9 classes of  
 500 MNIST. From each image, 50 features were extracted using a small neural network  $\phi : [0, 1]^{28 \times 28} \rightarrow \mathbb{R}^{50}$ ,  
 501 created by training a neural network to classify MNIST images and re-  
 502 moving the final layer. The goal is to train an SVM on these features using the hinge loss; see  
 503 the SVM formulation in [subsection 5.1.1](#) for more details. The problem class is of the form

$$504 \quad \mathcal{F} = \left\{ f_{\mathcal{I}}(\mathbf{w}, b) = \frac{1}{2} \mathbf{w}^\top \mathbf{w} + C \sum_{i \in \mathcal{I}} \max(0, 1 - y_i(\mathbf{w}^\top \phi_i + b)) \right\}.$$

505 This is the feature class of training SVMs with certain features  $\phi_i$  and targets  $y_i$ , with  $i$  taking  
 506 values in some index set  $\mathcal{I}$ . In this case, the features and targets were taken as subsets of  
 507 features extracted from the MNIST dataset. The initializations  $(\mathbf{w}, b) \sim \mathbb{P}_{(\mathbf{w}, b)|f}$  were taken  
 508 to be element-wise standard Gaussian.

509 [Figure 4](#) and [Figure 5](#) demonstrate the evolution of the SVM hinge loss under the qua-  
 510 dratic and one-layer NN mirror potentials respectively. In [Figure 4](#), the loss evolution under  
 511 quadratic LMD is faster compared to GD and Adam. This suggests that quadratic LMD  
 512 can learn features that contribute more to the SVM hinge loss. We can see clearly the ef-  
 513 fect of learning the step-sizes for increasing convergence rate for the first 10 iterations in the  
 514 “adaptive LMD” plot, as well as the effect of a non-optimized step-size after 10 iterations  
 515 by the increase in loss. In [Figure 5](#), we can see that the one-layer NN mirror potential can  
 516 perform significantly better than both Adam and GD. However, the instability due to the  
 517 required clipping causes the hinge loss to increase for larger step-sizes. This instability further

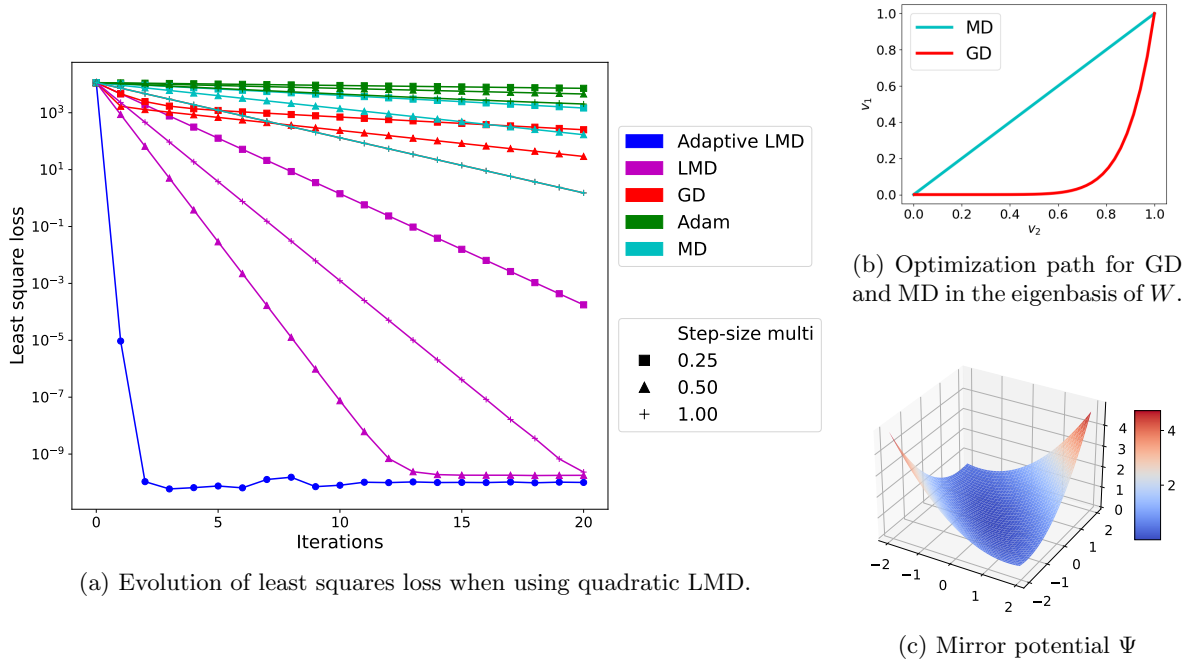


Figure 2: We observe that quadratic LMD is able to learn a map that allows for linear convergence in this case. Further learning the step-size allows for immediate convergence to machine precision. Note that  $W$  has eigenvectors  $v_1 = (1, 1), v_2 = (1, -1)$  with eigenvalues  $\lambda_1 = 3, \lambda_2 = 1$  respectively. As demonstrated in (b), the path that GD takes in the eigenbasis is curved as it minimizes the  $v_1$  direction faster than the  $v_2$  direction, whereas MD travels in a straight line to the minimizer. This is reflected in (c), where the quadratic form given by  $\Psi$  curves more in the  $v_1$  direction. Indeed, the learned weight is  $A = \begin{pmatrix} 0.69 & 0.55 \\ 0.55 & 0.69 \end{pmatrix}$ , which is almost proportional to the classical MD weight  $W^\top W = \begin{pmatrix} 5 & 4 \\ 4 & 5 \end{pmatrix}$ .

518 motivates the use of using a more expressive neural network, as well as directly modelling the  
 519 backwards mirror map.

520 Both of these methods require parameterizations of matrices, and moreover require com-  
 521 puting the inverse of these matrices, which can cause instability when performing back-  
 522 propagation. Moreover, such closed form expressions of the convex conjugate are not readily  
 523 available in general, especially for more complicated mirror potentials parameterized using  
 524 deep networks. Therefore, training LMD under this setting can not be effectively scaled up to  
 525 higher dimensions. This motivates our proposed approach and analysis of using two separate  
 526 networks instead, modeling the mirror and inverse mirror mappings separately.

527 **5. Numerical Experiments.** Motivated by the examples in the preceding section, we em-  
 528 ploy the LMD method for a number of convex problems arising in inverse problems and

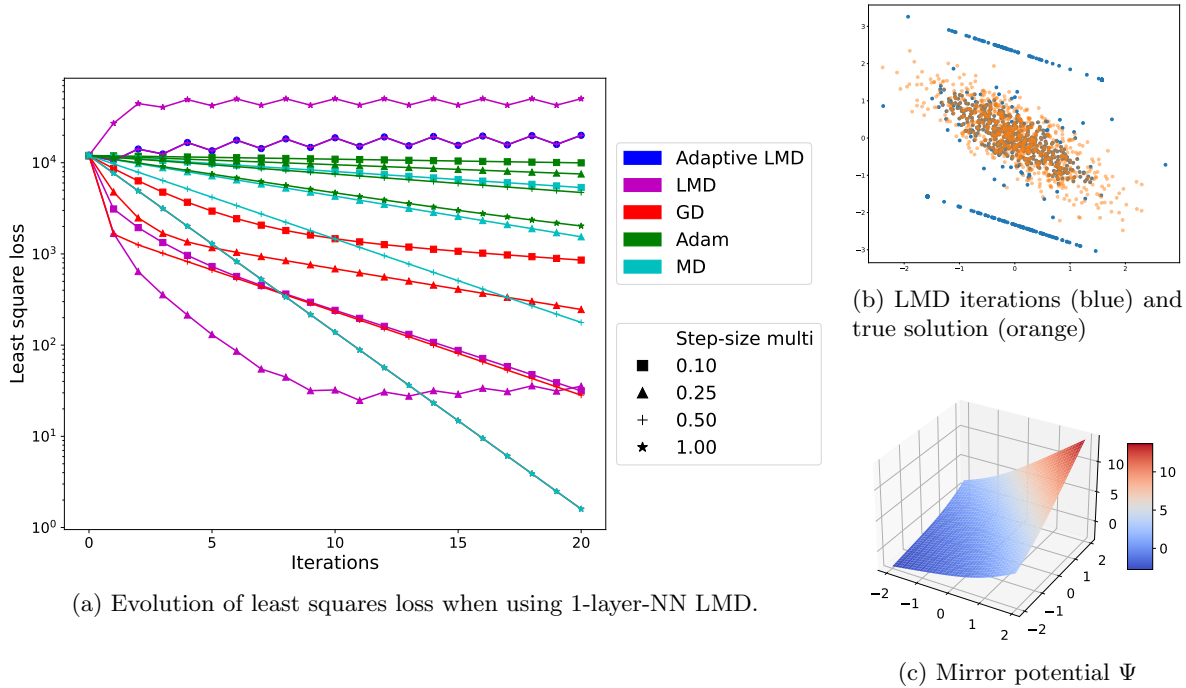


Figure 3: We can see the effect of needing to clip the dual iterates, as it creates a pair of lines (in blue). This heavily affects the performance when using certain step-sizes, and demonstrates the issues with such simple models. Note that the adaptive LMD and LMD with step-size multi 0.5 are identical. This is due to the choice of interval that the step-size is clipped to be in. The lower bound of the interval coincides with the step-size corresponding to step-size multiplier 0.5, and adaptive LMD learns the step-sizes to be this lower bound.

529 machine learning. Specifically, we use a deep ICNN for learning the optimal forward mirror  
 530 potential. However, unlike the constructions in the previous section, the convex conjugate  
 531 cannot be expressed in a closed form. We instead approximate the inverse of the mirror map  
 532 using a second neural network, which is not necessarily the gradient of an ICNN. We will  
 533 demonstrate how this can allow for learning the geometry of the underlying problems and  
 534 result in faster convergence. We will namely be applying the LMD method to the problems of  
 535 learning a two-class SVM classifier, learning a linear classifier, and model-based denoising and  
 536 inpainting on STL-10. The dimensionality of these problems, with STL-10 containing images  
 537 of size  $3 \times 96 \times 96$ , makes the matrix-based MD parameterizations proposed in the previous  
 538 section infeasible. A list of training and testing hyper-parameters can be found in [Table 2](#).

539 **5.1. SVM and Linear Classifier on MNIST.** We consider first the problem of training  
 540 an two-class SVM classifier and a multi-class linear classifier using features extracted from  
 541 MNIST. A small 5 layer neural network (2 convolutional layers, 1 dropout layer and 2 fully  
 542 connected layers) was first trained to a 97% accuracy, with the penultimate layer having 50

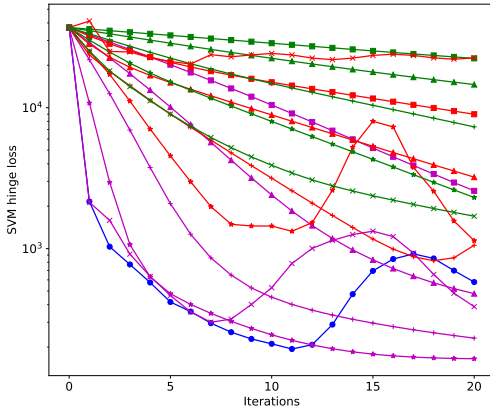


Figure 4: Evolution of SVM hinge loss under quadratic LMD. LMD outperforms GD and Adam, with nice convergence for the middle step-size multipliers. With only 3 out of 51 eigenvalues of  $A$  being greater than 1 and the rest below 0.5, this suggests that quadratic LMD is able to learn combinations of features that contribute most to the hinge loss.

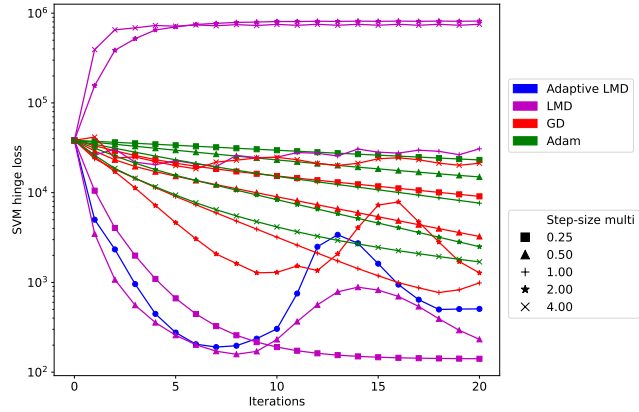


Figure 5: 1 layer NN mirror map applied to SVM training. In this case, LMD outperforms the other methods for smaller step-sizes. The two LMD lines with higher loss is due to the component-wise clipping that is required for this method.

543 features. We consider the problem of training an SVM on these features for two specific  
 544 classes. We also consider the problem of retraining the final layer of the neural network for  
 545 classification, which is equivalent to a linear classifier. Our goal is to minimize the correspond-  
 546 ing losses as quickly as possible using LMD. Let us denote the neural network that takes an  
 547 image and outputs the corresponding 50 features as  $\phi : [0, 1]^{28 \times 28} \rightarrow \mathbb{R}^{50}$ . This will work as a  
 548 feature extractor, on which we will train our SVMs and linear classifiers.

549 **5.1.1. SVM.** Our objective is to train a support vector machine (SVM) on the 50 ex-  
 550 tracted features to classify two classes of digits, namely 4 and 9. Given feature vectors  $\phi_i \in \mathbb{R}^d$   
 551 and target labels  $y_i \in \{\pm 1\}$ , an SVM consists of a weight vector  $\mathbf{w} \in \mathbb{R}^d$  and bias scalar  $b \in \mathbb{R}$ .  
 552 The output of the SVM for a given feature vector is  $\mathbf{w}^\top \phi_i + b$ , and the aim is to find  $\mathbf{w}$  and  
 553  $b$  such that the prediction  $\text{sign}(\mathbf{w}^\top \phi_i + b)$  matches the target  $y_i$  for most samples. The hinge  
 554 loss formulation of the problem is as follows, where  $C > 0$  is some positive constant [7]:

$$555 \quad (5.1) \quad \min_{\mathbf{w}, b} \frac{1}{2} \mathbf{w}^\top \mathbf{w} + C \sum_i \max(0, 1 - y_i(\mathbf{w}^\top \phi_i + b)).$$

556 The function class that we wish to learn to optimize for is thus

$$557 \quad (5.2) \quad \mathcal{F} = \left\{ f_{\mathcal{I}}(\mathbf{w}, b) = \frac{1}{2} \mathbf{w}^\top \mathbf{w} + C \sum_{i \in \mathcal{I}} \max(0, 1 - y_i(\mathbf{w}^\top \phi_i + b)) \right\},$$

558 where each instance of  $f$  depends on the set of feature-target pairs, indexed by  $\mathcal{I}$ . We use  
 559  $C = 1$  in our example. For each training iteration,  $\mathcal{I}$  was sampled as a subset of 1000



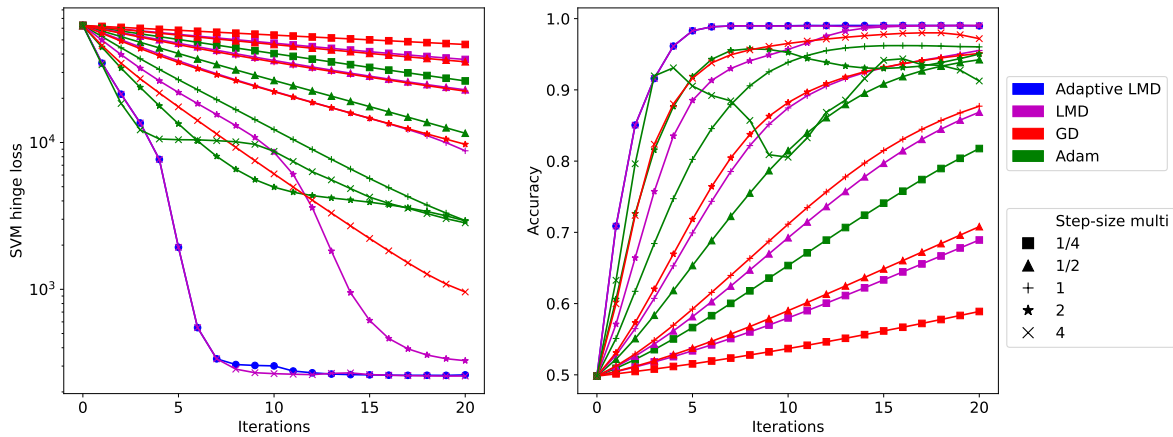


Figure 6: Plot of the SVM hinge loss (left) and SVM test accuracy (right) when optimizing from random SVM initializations. The mirror descent significantly outperforms both gradient descent and Adam, and does not exhibit as large of a decrease in accuracy for later iterations.

560 feature-target pairs from the combined 4 and 9 classes of MNIST, giving us a target function  
 561  $f_{\mathcal{I}}(\mathbf{w}, b) \in \mathcal{F}$ . A batch of 2000 initializations  $(\mathbf{w}, b)$  was then sampled according to a standard  
 562 normal distribution  $\mathbb{P}_{(\mathbf{w}, b)|f} = N(0, I_{50+1})$ . Subsets from the training fold were used for  
 563 training LMD, and subsets from the test fold to test LMD.

564 **Figure 6** shows the evolution of the hinge loss and SVM accuracy of the LMD method,  
 565 compared with GD and Adam. We can see that adaptive LMD and LMD with sufficiently  
 566 large step-size both outperform GD and Adam. In particular, considering LMD with step-  
 567 size multiplier 2, we can see accelerated convergence after around 10 iterations. One possible  
 568 interpretation is that the network is learning more about the geometry near the minima, which  
 569 is why we do not see this increased convergence for smaller step-sizes. The LMD method with  
 570 approximate backwards map is much more stable in this case, even if it performs slightly  
 571 worse than LMD with the one-layer NN-based mirror potential as in **Figure 5**.

572 **5.1.2. Linear Classifier.** We additionally consider the problem of training a multi-class  
 573 linear classifier on the MNIST features. We use the same neural network  $\phi$  to produce 50  
 574 features, and consider the task of training a linear final layer, taking the 50 features and  
 575 outputting 10 scores corresponding to each of the digits from 0-9. The task of finding the  
 576 optimal final layer with the cross entropy loss can be formulated as follows:

$$577 \quad (5.3) \quad \min_{W \in \mathbb{R}^{50 \times 10}} \mathbb{E}_{(\phi, y) \in \text{features} \times \text{target}} \left[ -\log \frac{\exp(W\phi)_y}{\sum_{i=0}^9 \exp(W\phi)_i} \right].$$

578 The corresponding feature class we wish to learn to optimize for is:

$$579 \quad (5.4) \quad \mathcal{F} = \left\{ f_{\mathcal{I}}(W) = \frac{1}{|\mathcal{I}|} \sum_{(\phi, y) \in \mathcal{I}} \left[ -\log \frac{\exp(W\phi)_y}{\sum_{i=0}^9 \exp(W\phi)_i} \right] \right\},$$

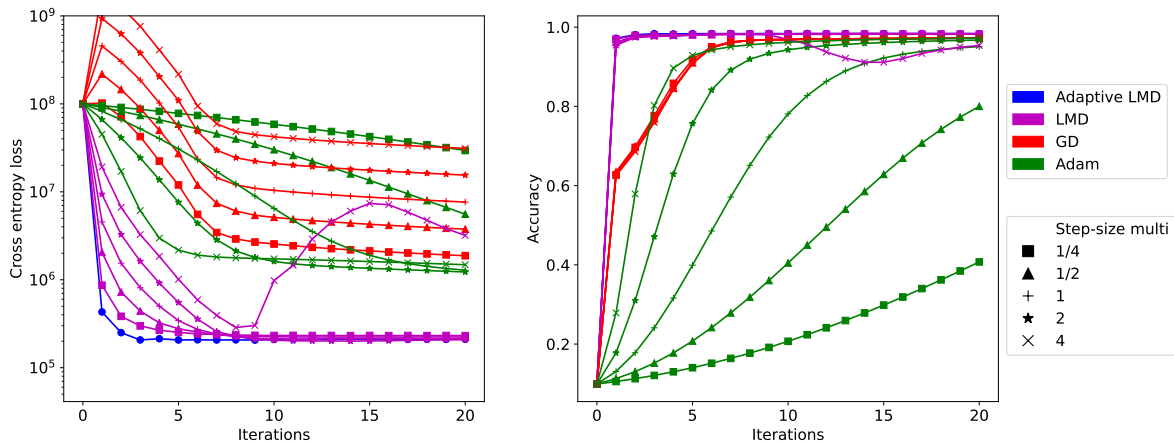


Figure 7: Plots of the linear classifier cross entropy loss (left) and classification accuracy (right). MD converges significantly faster than both GD and Adam. However, it suffers from stability issues for larger step-sizes, demonstrated by the increase in loss after 10 iterations with step-size multiplier 4. This increase in loss is also reflected in the decrease of accuracy.

580 where each instance of  $f$  depends on the set of feature-target pairs, indexed by  $\mathcal{I}$ . For each  
 581 training iteration,  $\mathcal{I}$  was sampled as a subset of 2000 feature-target pairs from MNIST, giving  
 582 a target function  $f_{\mathcal{I}}(W) \in \mathcal{F}$ . A batch of 2000 initializations  $W$  was then sampled according  
 583 to a standard normal distribution  $\mathbb{P}_{W|f} = N(0, I_{50 \times 50})$  for training. Subsets from the training  
 584 fold were used for training LMD, and subsets from the test fold to test LMD.

585 **Figure 7** shows the evolution of the cross-entropy loss and neural network classification  
 586 accuracy under our optimization schemes. All of the LMD methods converge quite quickly,  
 587 and we see that LMD with smaller step-sizes converge faster than larger step-sizes, reflecting  
 588 a similar phenomenon in gradient descent. We additionally see that for LMD with step-size  
 589 multiplier 4, the cross entropy loss has a large spike after 10 iterations. This is likely due to  
 590 the the step-size being too large for the Lipschitz constant of our problem.

591 **5.2. Image Denoising.** We further consider the problem of image denoising on the STL-10  
 592 image dataset [12]. Our goal is to have a fast solver for a single class of variational objectives  
 593 designed for denoising, rather than devise a state-of-the-art reconstruction approach. As the  
 594 reconstructions are completely model-driven and do not have a learned component, the quality  
 595 of the solution will depend completely on the chosen model.

596 The denoising problem is to minimize the distance between the reconstructed image with  
 597 an additional regularization term, which we have chosen to be total variation (TV). The  
 598 corresponding convex optimization problems can be represented as follows:

$$599 \quad (5.5) \quad \min_{x \in \mathcal{X}} \|x - y\|_{\mathcal{X}}^2 + \lambda \|\nabla x\|_{1, \mathcal{X}}.$$

600 Here,  $\mathcal{X}$  is the space of images from a pixel space  $\mathcal{S} \mapsto [0, 1]$ ,  $y$  is a noisy image,  $\lambda > 0$  is  
 601 a regularization parameter, and the gradient  $\nabla x$  is taken over the pixel space. In the case of

602 STL-10, the pixel space is  $3 \times 96 \times 96$ . The function class we wish to learn to optimize over  
 603 is thus:

$$604 \quad (5.6) \quad \mathcal{F} = \{f(x) = \|x - y\|_{\mathcal{X}}^2 + \lambda \|\nabla x\|_{1,\mathcal{X}} : \text{noisy images } y\}.$$

605 In our experiments,  $y$  was taken to have 5% random additive Gaussian noise over each  
 606 color channel, and the initializations  $x$  were taken to be the noisy images  $x = y$ . We trained  
 607 the LMD method on the training fold of STL10, and evaluated it on images in the test fold.

608 The TV regularization parameter was manually chosen to be  $\lambda = 0.3$  by visually comparing  
 609 the reconstructions after running gradient descent for 400 iterations. To parameterize the  
 610 mirror potentials, we use a convolutional neural network with an ICNN structure, as the data  
 611 is in 2D (with 3 color channels). We additionally introduce a quadratic term in each layer  
 612 for added expressiveness. The resulting models are of the following form, where the squaring  
 613 operator  $[\cdot]^2$  for a vector is to be taken element-wise, and  $\sigma$  is a leaky-ReLU activation  
 614 function:

$$615 \quad (5.7) \quad z_{i+1} = \sigma \left( W_i^{(z)} z_i + W_i^{(x,l)} x + [W_i^{(x,q)} x]^2 + b_i \right), \quad M(x; \theta) = z_l.$$

616 By clipping the kernel weights  $W_i^{(z)}$  to be non-negative, we are able to obtain an input convex  
 617 convolutional neural network.

618 [Figure 8](#) and [Figure 9](#) show the result of applying the LMD algorithm to the function class  
 619 of denoising models (5.6). In general, LMD and adaptive LMD outperform GD and Adam for  
 620 optimizing the reconstruction loss. Moreover, [Figure 9](#) shows that the reconstructed image  
 621 using LMD is very similar to the ones obtained using Adam, which is a good indicator that  
 622 LMD indeed solves the corresponding optimization problem efficiently.

623 [Figure 9a](#) shows a pixel-wise ratio between the forward map  $\nabla M_\theta(y)$  and noisy image  $y$ .  
 624 The outline of the horse demonstrates that  $\nabla M_\theta$  learns away from the identity, which should  
 625 contribute to the accelerated convergence. In particular, we observe that around the edges of  
 626 the horse, the pixel-wise ratio  $\nabla M_\theta(y)/y$  is negative. Intuitively, this corresponds to the MD  
 627 step performing gradient ascent instead of gradient descent for these pixels. As we are using  
 628 TV regularization, the gradient descent step aims to create more piecewise linear areas. If we  
 629 interpret gradient descent as a “blurring” step, then MD will instead perform a “sharpening”  
 630 step, which is more suited around the edges of the horse.

631 We additionally consider the effect of changing the noise level, and the ability of LMD to  
 632 generalize away from the training function class. We keep the LMD mirror maps trained for  
 633 5% additive Gaussian noise, and apply LMD to denoise images from STL-10 with additive  
 634 Gaussian noise levels up to 20%. We consider now the PSNR and SSIM of the denoised images  
 635 compared to a “true TV reconstruction”, which is obtained by optimizing the objective (5.5)  
 636 to a very high accuracy using gradient descent for 4000 iterations. We compare the iterates  
 637 with respect to the true TV reconstruction as opposed to the ground truth, as we want to  
 638 compare the resulting images with the minimum of the corresponding convex objective.

639 [Table 1](#) compares the PSNR and SSIM of denoised images obtained using LMD, Adam,  
 640 and GD, compared against the true TV reconstructions. We apply GD and LMD with five  
 641 fixed step-sizes ranging from  $2.5 \times 10^{-3}$  to  $4 \times 10^{-2}$  up to 20 iterations, and Adam with five

642 learning rates ranging from  $1.25 \times 10^{-2}$  to  $2 \times 10^{-1}$  for 20 iterations. We then compare the  
 643 best PSNR/SSIMs over all step-sizes and iterations for each method, and the best overall  
 644 step-sizes for the 10th and 20th iteration.

645 We see that LMD outperforms both GD and Adam when applied on the trained noise  
 646 level of 5% for the trained number of iterations  $N = 10$ , with better SSIM up to 10% noise as  
 647 well. LMD also performs well for lower noise levels, which can be attributed to good forward-  
 648 backward consistency near the true TV reconstruction. However, LMD begins to diverge for  
 649 larger noise levels. This can be attributed to the increased noise being out of the training  
 650 distribution, increasing the forward-backward loss and thereby causing instabilities.

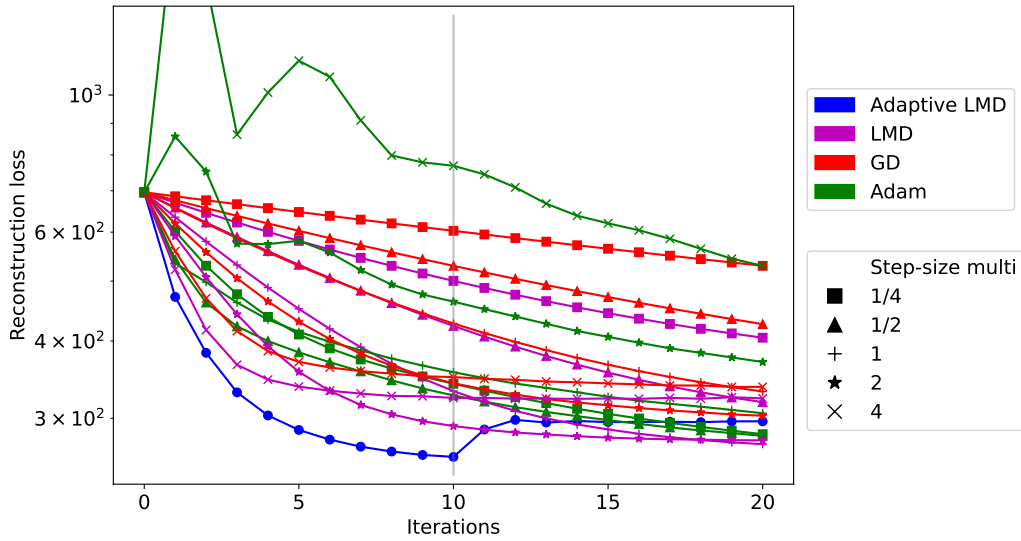


Figure 8: Denoising reconstruction loss. The vertical gray line at iteration 10 indicates the end of the training regime. After this line, the iterates are out-of-distribution for the proposed method. LMD outperforms both GD and Adam for earlier iterations, however might not reach the minimum due to forward-backward inconsistency. The sharp increase in loss for adaptive LMD after 10 iterations is due to the choice of step-size to extend the trained 10 iterations.

651 **5.3. Image Inpainting.** We additionally consider the problem of image inpainting with  
 652 added noise on STL10, in a similar setting to image denoising. 20% of the pixels in the image  
 653 were randomly chosen to be zero to create a fixed mask  $Z$ , and 5% Gaussian noise was added  
 654 to the masked images to create noisy masked images  $y$ . The inpainting problem is to minimize  
 655 the distance between the masked reconstructed image and the noisy masked image, including  
 656 TV regularization. The corresponding convex optimization problem is

$$657 \quad (5.8) \quad \min_{x \in \mathcal{X}} \|Z \circ (x - y)\|_{\mathcal{X}}^2 + \lambda \|\nabla x\|_{1, \mathcal{X}},$$

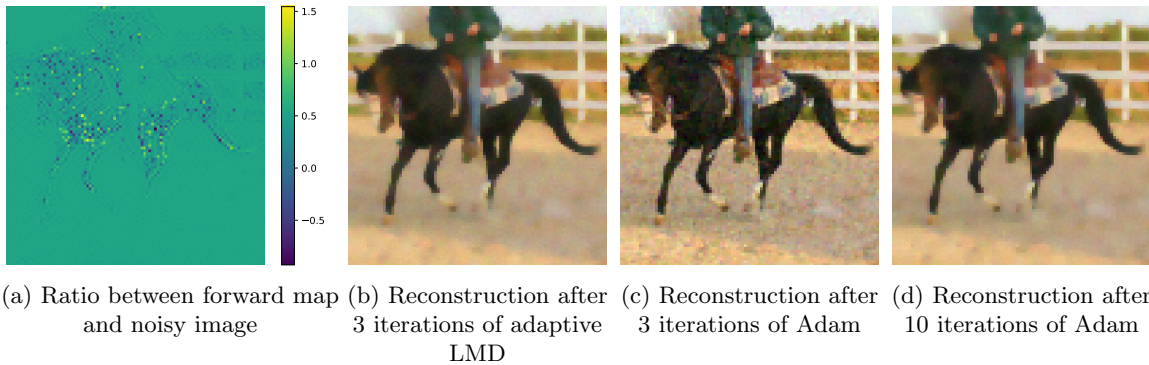


Figure 9: Visualization of outputs when applying LMD for TV model-based denoising. We can see a faint outline of the horse when taking a pixel-wise ratio between the forward and noisy image [indicating a region of interest](#). LMD allows for much faster convergence compared to Adam here, reaching a comparable reconstruction in only 3 iterations compared to 10 for Adam.

658 where  $Z$  denotes the masking map  $\mathcal{S} \mapsto \{0, 1\}^d$ , and the image difference  $x - y$  is taken  
 659 pixel-wise. The corresponding function class that we wish to learn to optimize over is:

$$660 \quad (5.9) \quad \mathcal{F} = \{f(x) = \|Z \circ (x - y)\|_{\mathcal{X}}^2 + \lambda \|\nabla x\|_{1, \mathcal{X}} : \text{noisy masked images } y\}.$$

661 The initializations  $x$  were taken to be the noisy masked images  $x = y$ . We trained the  
 662 LMD method on the training fold of STL10, and evaluated it on images in the test fold.  
 663 The TV regularization parameter was chosen to be  $\lambda = 0.3$  as in the denoising case, and the  
 664 mirror potentials are parameterized with a convolutional neural network similar to that used  
 665 in the denoising experiment. We trained the LMD method on the training fold of STL10, and  
 666 evaluated it on images in the test fold.

667 [Figure 10](#) shows the loss evolution of applying the LMD algorithm to the function class  
 668 of inpainting models (5.9). LMD with sufficiently large step-size outperforms GD and Adam,  
 669 however having too small of a step-size can lead to instability. We can also clearly see the  
 670 effect of approximating our backward maps, as some of the LMD methods result in asymptotic  
 671 reconstruction loss that is higher than a minimum. Nonetheless, adaptive LMD results in the  
 672 best convergence out of the tested methods.

673 [Figure 11](#) provides a visualization of the resulting iterations. [Figure 11a](#) plots the ratio  
 674 between the forward mapped masked image  $\nabla M_{\theta}(y)$  and masked image  $y$ , with clipped values  
 675 to prevent blowup in the plot. We can again see a faint outline of the horse [indicating a](#)  
 676 [region of interest, with some speckling due to the image mask](#). [Figure 11b](#) is a plot of the  
 677 result after 20 iterations of adaptive LMD, and it is qualitatively quite similar to the result  
 678 after 20 iterations of Adam, demonstrating the feasibility of LMD as a solver for model-based  
 679 reconstruction.

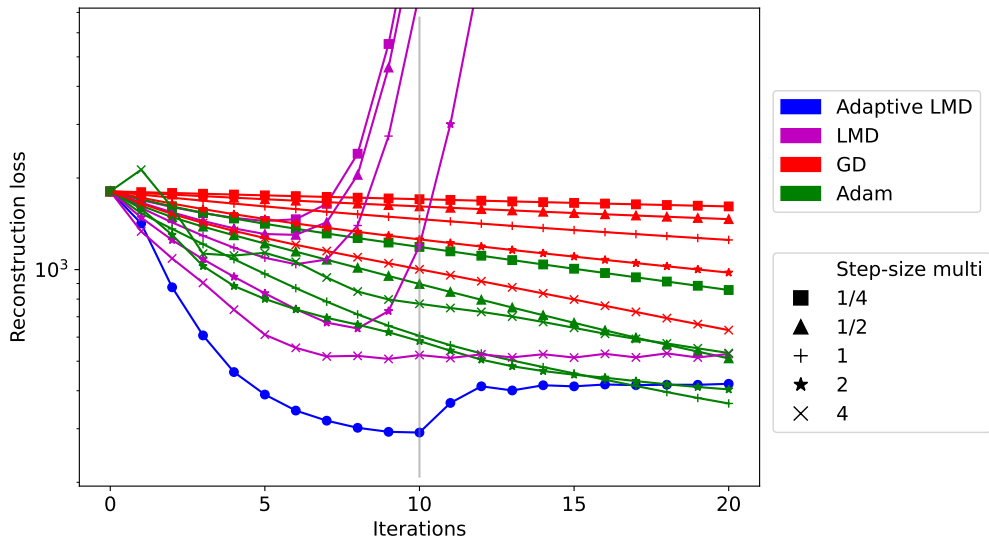


Figure 10: Inpainting reconstruction loss. The vertical gray line at iteration 10 indicates the end of the training regime. LMD outperforms both GD and Adam, however suffers from instability when the step-size is small, as remarked in Remark 3.4. The increase in loss after 10 iterations for adaptive LMD is due to the choice of step-size to extend the trained 10 iterations.

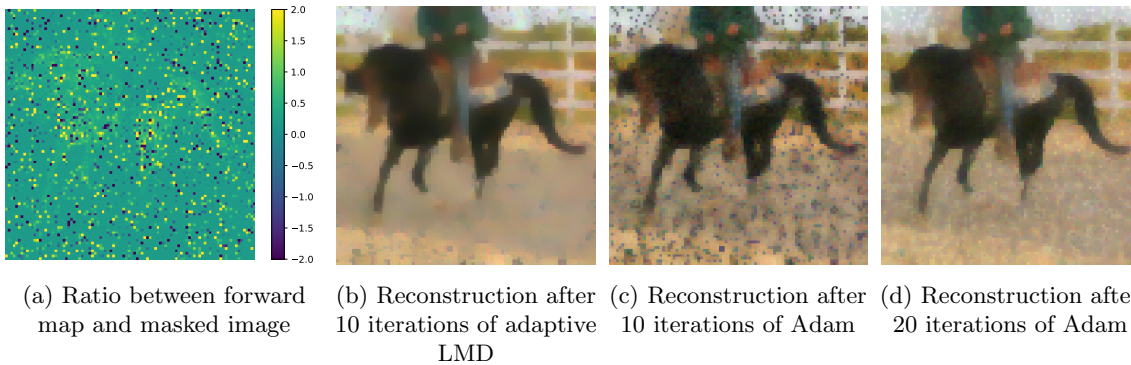


Figure 11: Visualization of some LMD on TV model-based inpainting. While a faint outline of the horse is visible, it is not as clear as in Figure 9 with speckling due to the zeroing mask. LMD is able to reach a reasonable reconstruction in fewer iterations compared to Adam. While the LMD reconstruction has artifacts around the edges, the Adam reconstruction is generally noisy.



Table 1: Table of PSNR and SSIM, compared to the true TV reconstruction. As our goal is to minimize the TV-regularized loss function, we compare with the loss-minimizing image as opposed to the ground truth image. LMD outperforms both GD and Adam when applied for noise levels up to 5% for the trained  $N = 10$  iterations, but is unstable for noise levels above 10%, which are out-of-distribution. Values are taken as the best over five step-sizes.

Gaussian Noise %	Best			Iteration 10			Iteration 20			
	GD	Adam	LMD	GD	Adam	LMD	GD	Adam	LMD	
PSNR	1	30.91	34.13	<b>34.25</b>	27.92	31.04	<b>33.27</b>	30.91	<b>34.03</b>	32.88
	2	30.93	34.06	<b>34.21</b>	27.90	30.86	<b>33.21</b>	30.93	<b>34.00</b>	32.87
	5	31.09	33.36	<b>34.22</b>	27.73	29.91	<b>32.92</b>	31.09	<b>33.44</b>	33.11
	10	31.08	<b>32.59</b>	28.21	26.45	<b>29.00</b>	27.88	31.08	<b>32.40</b>	25.92
	15	29.68	<b>32.56</b>	21.39	23.65	<b>28.89</b>	19.84	29.68	<b>32.30</b>	13.25
	20	28.96	<b>33.68</b>	20.12	22.97	<b>30.32</b>	10.94	28.96	<b>33.37</b>	-21.09
SSIM	1	0.905	0.960	<b>0.963</b>	0.862	0.914	<b>0.956</b>	0.905	0.956	<b>0.961</b>
	2	0.905	0.955	<b>0.963</b>	0.858	0.908	<b>0.955</b>	0.905	0.951	<b>0.961</b>
	5	0.898	0.935	<b>0.962</b>	0.857	0.880	<b>0.950</b>	0.898	0.932	<b>0.961</b>
	10	0.893	0.907	<b>0.950</b>	0.817	0.831	<b>0.908</b>	0.893	0.902	<b>0.950</b>
	15	0.876	<b>0.893</b>	0.849	0.698	<b>0.799</b>	0.689	0.876	<b>0.889</b>	0.849
	20	0.850	<b>0.917</b>	0.887	0.662	<b>0.841</b>	0.772	0.850	<b>0.915</b>	0.878

Table 2: Hyper-parameters for the problem classes considered.

	SVM	Linear Classifier	Denoising	Inpainting
Batch size	2000	2000	10	10
Epochs	10,000	10,000	1300	1100
All				
ICNN training parameters (Adam)	$\alpha = 10^{-5}, \beta = (0.9, 0.99)$			
Learned iterations $N$	10			
Learned step-size initialization	$10^{-2}$			
Learned step-size range	$(10^{-3}, 10^{-1})$			
Testing base step-size (LMD,GD)	$10^{-2}$			
Testing base step-size (Adam)	$5 \times 10^{-2}$			

680 **5.4. Effect of Regularization Parameter.** We now turn to studying the effect of the reg-  
681 ularization parameters used to enforce consistency of the forward and backward mirror maps.  
682 The regularization parameter  $s_k = s_{\text{epoch}}$  as in (3.26b) was initialized as 1, and subsequently  
683 multiplied by 1.05 every 50 epochs.

684 Under the assumption that the model is trained well for each regularization parameter,  
685 the training loss gives a perspective into the trade-off between the loss and the forward-

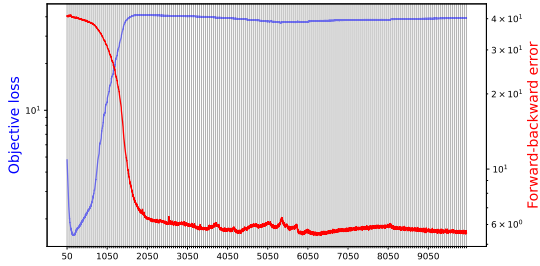


Figure 12: Training loss and forward-backward consistency loss when training an SVM, plotted against training epochs. We can see clearly the tradeoff between the loss and forward-backward loss at the earlier iterations. Each vertical grey line corresponds to an epoch where the forward-backward loss regularization is increased.

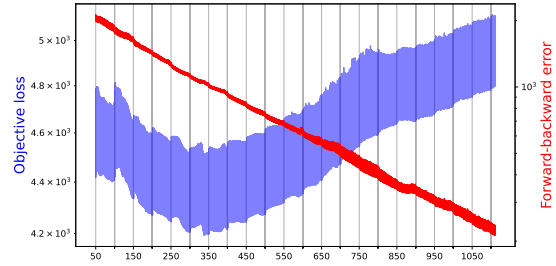


Figure 13: Training loss and forward-backward consistency loss when training inpainting on STL10, plotted against training epochs. We can see the effect of increasing the forward-backward regularization parameter as the forward-backward loss continues to decrease along the iterations, while the loss begins to increase.

686 backward consistency of the learned mirror maps. Informally, the model will try to learn  
 687 a one-shot method similar to an end-to-end encoder-decoder model. Increasing the forward-  
 688 backward regularization parameter  $s_{\text{epoch}}$  reduces this one-shot effect, and encourages a proper  
 689 optimization scheme to emerge. Therefore, it is natural that the objective loss will increase  
 690 as the forward-backward loss decreases. This effect can be seen in [Figure 12](#), where the  
 691 objective loss starts very low but then increases as the forward-backward error decreases.  
 692 This could be interpreted as the LMD learning a single good point, then switching to learning  
 693 how to optimize to a good point. In addition to encouraging a proper optimization scheme,  
 694 increasing the forward-backward regularization parameter has the added effect of encouraging  
 695 the forward-backward loss to continue decreasing. This can be seen in [Figure 13](#), where the  
 696 objective loss also decreases before increasing again.

697 **5.5. Ablation Study.** In this section, we will compare the effect of various design choices  
 698 on LMD. In particular, we will consider (i) the effect of the number of training iterations  
 699  $N$ , (ii) the effect of not enforcing the forward-backward consistency by setting  $s_k = 0$ , and  
 700 (iii) a further comparison against GD with learned step-sizes (LGD). In particular, the first  
 701 experiment will be LMD trained with  $N = 2$ . The latter experiment is equivalent to our LMD  
 702 with both mirror maps fixed to be the identity. We will compare these three experiments on  
 703 the inpainting setting as in [Subsection 5.3](#). [Figure 14](#) compares the forward-backward incon-  
 704 sistency and the loss for these three experiments with LMD trained for inpainting, detailed  
 705 in [Subsection 5.3](#). For each of these methods, we choose to extend the learned step-sizes by a  
 706 constant, up to 20 iterations.

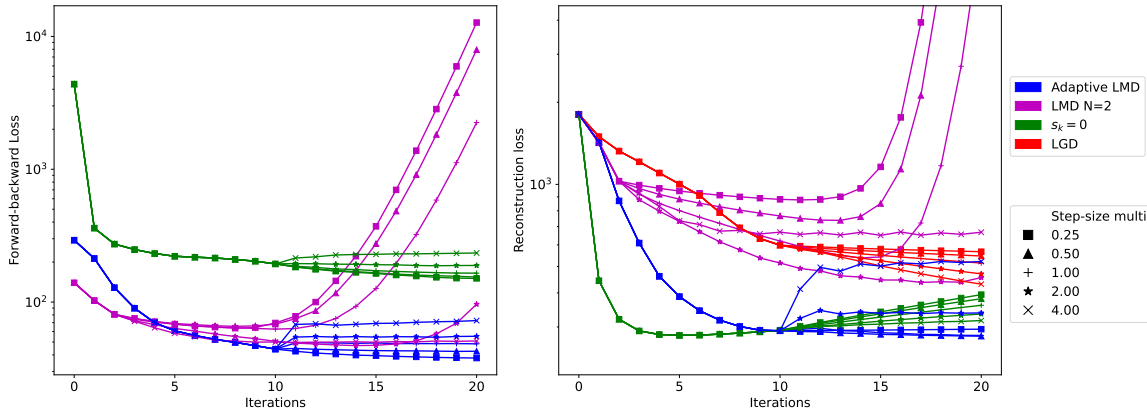
707 For experiment (i), decreasing the number of training iterations  $N$  severely impacts the  
 708 forward-backward inconsistency. Moreover, the number of training iterations is insufficient to

709 be close to the minimum of the problem. These problems coupled together lead to the loss  
 710 converging to a poor value, or diverging depending on the step-size extension.

711 For experiment (ii), setting  $s_k = 0$  in (3.26b) and not enforcing forward-backward con-  
 712 sistency results in high forward-backward loss. Nonetheless, the loss rapidly decreases in the  
 713 first couple iterations, faster than LMD. This is consistent with the view that the pair of mir-  
 714 ror potentials acts as an encoder-decoder network, rapidly attaining close to the minimum.  
 715 Due to the higher forward-backward loss, this method has looser bounds on the convergence,  
 716 resulting in the increase in reconstruction loss in the later iterations compared to LMD.

717 For experiment (iii), learning the step-sizes for GD directly results in significantly worse  
 718 performance compared to LMD. This can be attributed to LMD learning the direction of  
 719 descent via the mirror maps, in addition to the speed of descent given by the learned step-  
 720 sizes. This demonstrates that a better direction than steepest descent exists, can be learned  
 721 by LMD and results in faster convergence rates.

722 These experiments demonstrate the effect of the variables of LMD. In particular, we show  
 723 that a sufficient number of training iterates is required to maintain longer term convergence,  
 724 and that enforcing forward-backward consistency sacrifices some early-iterate convergence rate  
 725 for better stability for longer iterations. Moreover, the LGD experiment shows that learning  
 726 a direction via the mirror maps in addition to the speed of convergence allows for faster  
 727 convergence.



(a) Evolution of forward-backward loss.

(b) Evolution of inpainting reconstruction loss.

Figure 14: Ablation study considering the forward-backward loss and reconstruction loss for image inpainting. We consider (i) training a small number of iterations  $N = 2$ , (ii) training without enforcing the forward-backward inconsistency  $s_k = 0$ , and (iii) training where the mirror maps are fixed to be the identity, corresponding to GD with learned step-sizes (LGD). Adaptive LMD is trained for  $N = 10$  iterations, as in [Subsection 5.3](#). Note that LGD does not have a forward-backward loss, as the iterates are exact.

728 **5.6. Computational Complexity.** In this subsection, we discuss the computational com-  
 729 plexity of the LMD method in terms time and memory, for both training and testing time.

730 For a single backward and forward pass, the proposed method scales linearly with the  
 731 dimension and number of iterations. In particular, suppose the space we wish to optimize is  
 732 over  $\mathcal{X} \subseteq \mathbb{R}^d$  with dimension  $d$  and batch size  $n$ , and we run  $N$  iterations of LMD. Assuming  
 733 that backpropagating and taking gradients scales linearly with the number of parameters  $P$ ,  
 734 the backwards pass takes  $\mathcal{O}(n \times d \times N \times P)$  time and  $\mathcal{O}(n \times d \times N \times P)$  memory. The forwards  
 735 pass takes  $\mathcal{O}(n \times d \times N \times P)$  time and  $\mathcal{O}(n \times d \times P)$  memory, where we drop a factor of  $N$   
 736 as holding intermediate iterates is not required.

737 **Table 3** compares the GPU wall-times and memory consumption for various numbers  
 738 of training iterations  $N$ , tested for the STL-10 inpainting experiment for both training and  
 739 testing. We find that the times and memory consumption are as expected, with near-linear  
 740 increase in time and train memory, and near-constant test memory.

Table 3: Table of GPU wall time and memory consumption for training and testing LMD, with various iteration counts. Times are per batch, with a batch-size of 25 on STL-10 images with dimension  $3 \times 96 \times 96$ . Training and testing was done on Quadro RTX 6000 GPUs with 24GB of memory.

Iterations	Train time (s)	Test time (s)	Train memory (GB)	Test memory (GB)
$N = 2$	5.13	0.422	8.24	1.52
$N = 5$	8.26	1.05	12.68	1.53
$N = 10$	13.41	2.11	20.09	1.54
$N = 20$	-	4.22	-	1.57
$N = 50$	-	10.55	-	1.64

741 **6. Discussion and Conclusions.** In this work, we proposed a new paradigm for learning-  
 742 to-optimize with theoretical convergence guarantees, interpretability, and improved numerical  
 743 efficiency for convex optimization tasks in data science, based on learning the optimal Breg-  
 744 man distance of mirror descent modeled by input-convex neural networks. Due to this novel  
 745 functional parameterization of the mirror map, and by taking a structured and theoretically-  
 746 principled approach, we are able to provide convergence guarantees akin to the standard  
 747 theoretical results of classical mirror descent. We then demonstrate the effectiveness of our  
 748 LMD approach via extensive experiments on various convex optimization tasks in data sci-  
 749 ence, comparing to classical gradient-based optimizers. The provable LMD approach achieves  
 750 competitive performance with Adam, a heuristically successful method. However, Adam lacks  
 751 convergence guarantees for the convex case, achieving only local convergence [27, 8, 34]. LMD  
 752 is able to achieve the fast convergence rates from Adam, while retaining convergence guaran-  
 753 tees from slower classical methods such as GD.

754 In this paper, we have only considered the most basic form of mirror descent as our starting  
 755 point. There is still much potential for further improvements on both theoretical results and  
 756 numerical performance of the algorithm. If a deep parameterization of convex functions with

757 closed form convex conjugate exists, then this would allow for exact convergence. One open  
 758 question is what an optimal mirror map should look like for a particular problem class such as  
 759 image denoising, and how well a deep network is able to approximate it. Our ongoing works  
 760 include accelerating the convergence rates of LMD with momentum acceleration techniques  
 761 which have been developed for accelerating classical mirror descent [15, 16], and stochastic  
 762 approximation schemes [33].

763

## REFERENCES

- 764 [1] B. AMOS, L. XU, AND J. Z. KOLTER, *Input convex neural networks*, in Proceedings of the 34th Interna-  
 765 tional Conference on Machine Learning, vol. 70 of Proceedings of Machine Learning Research, PMLR,  
 766 2017, pp. 146–155.
- 767 [2] M. ANDRYCHOWICZ, M. DENIL, S. GÓMEZ, M. W. HOFFMAN, D. PFAU, T. SCHAUL, B. SHILLINGFORD,  
 768 AND N. DE FREITAS, *Learning to learn by gradient descent by gradient descent*, in Advances in Neural  
 769 Information Processing Systems, D. Lee, M. Sugiyama, U. Luxburg, I. Guyon, and R. Garnett, eds.,  
 770 vol. 29, Curran Associates, Inc., 2016.
- 771 [3] S. BANERT, A. RINGH, J. ADLER, J. KARLSSON, AND O. ÖKTEM, *Data-driven nonsmooth optimization*,  
 772 SIAM Journal on Optimization, 30 (2020), pp. 102–131.
- 773 [4] S. BANERT, J. RUDZUSIKA, O. ÖKTEM, AND J. ADLER, *Accelerated forward-backward optimization using*  
 774 *deep learning*, 2021, <https://doi.org/10.48550/ARXIV.2105.05210>.
- 775 [5] A. BECK AND M. TEBoulLE, *Mirror descent and nonlinear projected subgradient methods for convex*  
 776 *optimization*, Operations Research Letters, 31 (2003), pp. 167–175.
- 777 [6] A. BEN-TAL, T. MARGALIT, AND A. NEMIROVSKI, *The ordered subsets mirror descent optimization*  
 778 *method with applications to tomography*, SIAM Journal on Optimization, 12 (2001), pp. 79–108,  
 779 <https://doi.org/10.1137/S1052623499354564>, <https://doi.org/10.1137/S1052623499354564>, [https://](https://arxiv.org/abs/https://doi.org/10.1137/S1052623499354564)  
 780 [arxiv.org/abs/https://doi.org/10.1137/S1052623499354564](https://arxiv.org/abs/https://doi.org/10.1137/S1052623499354564).
- 781 [7] C. M. BISHOP, *Pattern Recognition and Machine Learning (Information Science and Statistics)*, Springer-  
 782 Verlag, Berlin, Heidelberg, 2006.
- 783 [8] S. BOCK AND M. WEISS, *A proof of local convergence for the adam optimizer*, in 2019 International  
 784 Joint Conference on Neural Networks (IJCNN), 2019, pp. 1–8, [https://doi.org/10.1109/IJCNN.2019.](https://doi.org/10.1109/IJCNN.2019.8852239)  
 785 [8852239](https://doi.org/10.1109/IJCNN.2019.8852239).
- 786 [9] S. BUBECK ET AL., *Convex optimization: Algorithms and complexity*, Foundations and Trends® in  
 787 Machine Learning, 8 (2015), pp. 231–357.
- 788 [10] A. CHAMBOLLE AND T. POCK, *A first-order primal-dual algorithm for convex problems with applications*  
 789 *to imaging*, J. Math. Imaging and Vision, 40 (2010), pp. 120–145.
- 790 [11] S. H. CHAN, X. WANG, AND O. A. ELGENDY, *Plug-and-play ADMM for image restoration: Fixed*  
 791 *point convergence and applications*, CoRR, abs/1605.01710 (2016), <http://arxiv.org/abs/1605.01710>,  
 792 <https://arxiv.org/abs/1605.01710>.
- 793 [12] A. COATES, A. NG, AND H. LEE, *An analysis of single-layer networks in unsupervised feature learning*,  
 794 in Proceedings of the Fourteenth International Conference on Artificial Intelligence and Statistics,  
 795 G. Gordon, D. Dunson, and M. Dudík, eds., vol. 15 of Proceedings of Machine Learning Research,  
 796 Fort Lauderdale, FL, USA, 11–13 Apr 2011, PMLR, pp. 215–223.
- 797 [13] K. GREGOR AND Y. LECUN, *Learning fast approximations of sparse coding*, in Proceedings of the 27th  
 798 International Conference on International Conference on Machine Learning, ICML’10, Omnipress,  
 799 2010, p. 399–406.
- 800 [14] S. GUNASEKAR, B. WOODWORTH, AND N. SREBRO, *Mirrorless mirror descent: A natural derivation of*  
 801 *mirror descent*, in Proceedings of The 24th International Conference on Artificial Intelligence and  
 802 Statistics, A. Banerjee and K. Fukumizu, eds., vol. 130 of Proceedings of Machine Learning Research,  
 803 PMLR, 13–15 Apr 2021, pp. 2305–2313.
- 804 [15] F. HANZELY, P. RICHTARIK, AND L. XIAO, *Accelerated bregman proximal gradient methods for relatively*  
 805 *smooth convex optimization*, Computational Optimization and Applications, 79 (2021), pp. 405–440.
- 806 [16] W. KRICHENE, A. BAYEN, AND P. L. BARTLETT, *Accelerated mirror descent in continuous and discrete*

- 807 *time*, Advances in neural information processing systems, 28 (2015).
- 808 [17] G. LAN, *An optimal method for stochastic composite optimization*, Mathematical Programming, 133  
809 (2012), pp. 365–397.
- 810 [18] G. LAN AND Y. ZHOU, *An optimal randomized incremental gradient method*, Mathematical programming,  
811 171 (2018), pp. 167–215.
- 812 [19] K. LI AND J. MALIK, *Learning to optimize*, CoRR, abs/1606.01885 (2016), <https://arxiv.org/abs/1606.01885>.
- 813
- 814 [20] H. LU, R. M. FREUND, AND Y. NESTEROV, *Relatively smooth convex optimization by first-order methods,  
815 and applications*, SIAM Journal on Optimization, 28 (2018), pp. 333–354.
- 816 [21] N. MAHESWARANATHAN, D. SUSSILLO, L. METZ, R. SUN, AND J. SOHL-DICKSTEIN, *Reverse engineering  
817 learned optimizers reveals known and novel mechanisms*, CoRR, abs/2011.02159 (2020), <https://arxiv.org/abs/2011.02159>.
- 818
- 819 [22] S. MUKHERJEE, S. DITTMER, Z. SHUMAYLOV, S. LUNZ, O. ÖKTEM, AND C.-B. SCHÖNLIEB, *Learned con-  
820 vex regularizers for inverse problems*, arXiv:2008.02839v2, (2020), <https://doi.org/10.48550/ARXIV.2008.02839>, <https://arxiv.org/abs/2008.02839>.
- 821
- 822 [23] A. S. A. S. NEMIROVSKY, *Problem complexity and method efficiency in optimization / A.S. Nemirovsky,  
823 D.B. Yudin ; translated by E.R. Dawson.*, Wiley-Interscience series in discrete mathematics, Wiley,  
824 Chichester, 1983.
- 825 [24] Y. NESTEROV, *Gradient methods for minimizing composite functions*, Mathematical programming, 140  
826 (2013), pp. 125–161.
- 827 [25] F. ORABONA, K. CRAMMER, AND N. CESA-BIANCHI, *A generalized online mirror descent with applica-  
828 tions to classification and regression*, Machine Learning, 99 (2015), pp. 411–435.
- 829 [26] A. PASZKE, S. GROSS, F. MASSA, A. LERER, J. BRADBURY, G. CHANAN, T. KILLEEN, Z. LIN,  
830 N. GIMELSHEIN, L. ANTIGA, A. DESMAISON, A. KOPF, E. YANG, Z. DEVITO, M. RAISON,  
831 A. TEJANI, S. CHILAMKURTHY, B. STEINER, L. FANG, J. BAI, AND S. CHINTALA, *Pytorch: An  
832 imperative style, high-performance deep learning library*, in Advances in Neural Information Pro-  
833 cessing Systems 32, H. Wallach, H. Larochelle, A. Beygelzimer, F. d'Alché-Buc, E. Fox, and  
834 R. Garnett, eds., Curran Associates, Inc., 2019, pp. 8024–8035, [http://papers.neurips.cc/paper/  
835 9015-pytorch-an-imperative-style-high-performance-deep-learning-library.pdf](http://papers.neurips.cc/paper/9015-pytorch-an-imperative-style-high-performance-deep-learning-library.pdf).
- 836 [27] S. J. REDDI, S. KALE, AND S. KUMAR, *On the convergence of adam and beyond*, CoRR, abs/1904.09237  
837 (2019), <http://arxiv.org/abs/1904.09237>, <https://arxiv.org/abs/1904.09237>.
- 838 [28] R. ROCKAFELLAR AND R. J.-B. WETS, *Variational Analysis*, Springer Verlag, Heidelberg, Berlin, New  
839 York, 1998.
- 840 [29] E. K. RYU, J. LIU, S. WANG, X. CHEN, Z. WANG, AND W. YIN, *Plug-and-play methods provably  
841 converge with properly trained denoisers*, CoRR, abs/1905.05406 (2019), [http://arxiv.org/abs/1905.  
842 05406](http://arxiv.org/abs/1905.05406), <https://arxiv.org/abs/1905.05406>.
- 843 [30] P. TSENG, *On accelerated proximal gradient methods for convex-concave optimization*, tech. report, Uni-  
844 versity of Washington, Seattle, 2008.
- 845 [31] J. VANSCHOREN, *Meta-learning: A survey*, CoRR, abs/1810.03548 (2018), [https://arxiv.org/abs/1810.  
846 03548](https://arxiv.org/abs/1810.03548).
- 847 [32] S. V. VENKATAKRISHNAN, C. A. BOUMAN, AND B. WOHLBERG, *Plug-and-play priors for model based  
848 reconstruction*, in 2013 IEEE Global Conference on Signal and Information Processing, 2013, pp. 945–  
849 948, <https://doi.org/10.1109/GlobalSIP.2013.6737048>.
- 850 [33] P. XU, T. WANG, AND Q. GU, *Accelerated stochastic mirror descent: From continuous-time dynamics to  
851 discrete-time algorithms*, in International Conference on Artificial Intelligence and Statistics, PMLR,  
852 2018, pp. 1087–1096.
- 853 [34] F. ZOU, L. SHEN, Z. JIE, W. ZHANG, AND W. LIU, *A sufficient condition for convergences of adam and  
854 rmsprop*, CoRR, abs/1811.09358 (2018), <https://arxiv.org/abs/1811.09358>.