

# STOCHASTIC MIRROR DESCENT FOR CONVEX OPTIMIZATION WITH CONSENSUS CONSTRAINTS

A. BOROVIKH, N. KANTAS, P. PAPPAS, G. A. PAVLIOTIS

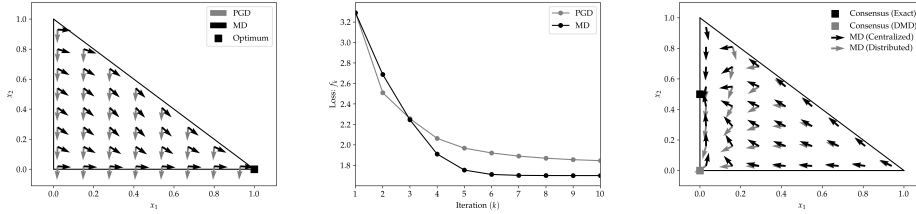
**Abstract.** The mirror descent algorithm is known to be effective in applications where it is beneficial to adapt the mirror map to the underlying geometry of the optimization model. However, the effect of mirror maps on the geometry of distributed optimization problems has not been previously addressed. In this paper we propose and study exact distributed mirror descent algorithms in continuous-time under additive noise and present the settings that enable linear convergence rates. Our analysis draws motivation from the augmented Lagrangian and its relation to gradient tracking. To further explore the benefits of mirror maps in a distributed setting we present a preconditioned variant of our algorithm with an additional mirror map over the Lagrangian dual variables. This allows our method to adapt to the geometry of the consensus manifold and leads to faster convergence. We illustrate the performance of the algorithms in convex settings both with and without constraints. We also explore their performance numerically in a non-convex application with neural networks.

**Key words.** Distributed optimization, mirror descent, pre-conditioning, interacting particles, stochastic optimization.

**AMS subject classifications.** 65K10, 68Q87, 60G07, 68W15

**1. Introduction.** The choice of mirror map has a significant impact on both the theoretical and numerical performance of the Mirror Descent (MD) algorithm [4, 9]. With an appropriate choice of the mirror map, MD captures the geometry of the optimization model more faithfully than other first-order methods. We illustrate this point in Figure 1a by plotting the vector fields generated by MD (using the negative entropy function as the mirror map) and Projected Gradient Descent (PGD) (with Euclidean projection) for a strongly convex quadratic optimization problem over the three-dimensional simplex. It is clear from Figure 1a that the PGD vector field points in the correct direction towards the unique minimum. But as soon as the PGD vector field hits the boundary, then the algorithm slows down considerably. The slowdown is due to the fact that the gradient always points towards the direction of steepest descent *for the objective function irrespective of the constraints*. When PGD hits the boundary, then the steepest descent direction is no longer appropriate for the problem's geometry. When MD hits the boundary of the feasible region, it glides across the boundary and towards the solution. This observation is reflected in the numerical performance of the two algorithms. In Figure 1b we indeed see that PGD initially makes good progress towards the solution but then stalls. MD, on the other hand, is slower in the first two iterations but converges to the optimal solution much faster. This phenomenon is not only present in problems with constraints but is also relevant in unconstrained problems, especially for ill-conditioned problems, and inverse optimization problems that have a sparsity inducing norm in the objective function. For example, in unconstrained ill-conditioned problems, the gradient descent method performs no preconditioning, whereas mirror descent uses the Hessian of the mirror map as a preconditioner (see Section 3.1).

There exists no theoretical or algorithmic framework to explain how to compute an optimal mirror map for a given problem. However, mirror maps for some particular classes of problems are well known (see Appendix A.1, [4], [9]). Mirror descent, and especially the effect of the choice of the mirror map for *distributed* optimization problems has received much less attention (see Section 1.1 for related work). Distributed optimization problems, even when otherwise unconstrained, have to satisfy a consensus



(a) Vector Fields of Centralized PGD/MD (b) PGD/MD iterations (c) Vector Fields of Distributed PGD/MD

Fig. 1: Vector fields for Projected Gradient Descent (PGD) and Mirror Descent (MD) for a quadratic function over the three dimensional simplex (we plot the two dimensional projection). MD uses the negative entropy function as the mirror map, and PGD performs the projection using the  $\ell_2$  norm.

48 constraint, and existing algorithms do not capture the geometry of the consensus  
 49 manifold. Motivated by the attractive features of the mirror descent algorithm  
 50 described above, we attempt to answer the question: *Does there exist a distributed*  
 51 *variant of Mirror Descent that can accurately capture the geometry of distributed*  
 52 *optimization models?* To answer this question, we study distributed algorithms for the  
 53 following optimization model,

$$54 \quad (1.1) \quad \min_{x_i \in \mathcal{X}} \sum_{i=1}^N f_i(x^i), \quad \text{s.t. } x^i = x^j \quad \forall (i, j) \in E.$$

55 The  $\{x^i\}_{i=1}^N$  with indices  $i = 1, \dots, N$  denote the computational nodes or particles,  
 56 as we refer to them in previous work [6]. These communicate through a strongly  
 57 connected, weighted, undirected graph  $\mathcal{G} := (V, E, A)$ ; where  $V$  represents the nodes  
 58 of the graph,  $E$  its edges and  $A$  is the adjacency matrix. Each particle has access to  
 59 its own objective function  $f_i : \mathbb{R}^d \rightarrow \mathbb{R}$ , and constraint set  $\mathcal{X} \subset \mathbb{R}^d$ .

60 For the purposes of motivating the results of our work consider the following  
 61 natural generalization of Distributed Mirror Descent (DMD),

$$62 \quad (1.2) \quad \partial_t \nabla \Phi(x_t^i) = -\nabla f_i(x_t^i) - \sum_{j=1}^N A_{ij}(x_t^i - x_t^j) \quad i = 1, \dots, N,$$

63 where  $\Phi$  is the mirror map, and  $A_{ij}$  is the weight of edge  $(i, j)$ . We call a distributed  
 64 algorithm *exact* if it converges to a solution that it is both optimal and satisfies the  
 65 *consensus constraint*. In Figure 1c we plot the vector field generated by (1.2) on a  
 66 quadratic optimization model with  $N = 2$  over the three dimensional simplex, and  
 67 the centralized MD algorithm. The centralized MD algorithm for (1.1) substitutes the  
 68 constraints in the objective function and follows the dynamics below,

$$69 \quad \partial_t \nabla \Phi(x_t) = -\sum_{j=1}^N \nabla f_j(x_t).$$

70 As expected, the two algorithms generate different vector fields. What is more  
 71 concerning is that the Distributed Mirror Descent in (1.2) does not converge to

Reference	Mirror	Rate	Exact	Noise	Step-size
Liang et al. [21]	No	Linear	Yes	No	-
[22, 16, 23, 43]	No	N/A	Yes	No	-
Sun et al. [38]	Yes	N/A	Yes	No	-
This work	Yes	Linear	Yes	Yes	-
Shi et al. [36]	No	Linear	Yes	No	Constant
Qu & Li [30]	No	Linear	Yes	No	Constant
Jakovetic et al. [19]	No	Sub-Linear	No	No	Diminishing
Pu & Nedic [29]	No	Linear	Yes	Yes	Constant
Ram et. al. [32]	No	N/A	Yes/No	Yes	Diminish./Const.
Duchi et al. [13]	Yes	Sub-Linear	Yes/No	Yes	Diminish./Const.
Nedic et al. [27]	Yes	N/A	Yes	Yes	Diminishing
Shahrampour et al. [33]	Yes	Sub-Linear	Yes	Yes	Diminishing

Table 1: Overview of convergence rates for different types of algorithms. Exact refers to whether or not the algorithm achieves exact consensus, and mirror refers to whether or not the algorithm allows for mirror maps. Continuous time methods are marked as - in the step size entry.

72 the unique solution of the problem. This observation is not surprising given that  
73 Distributed Gradient Descent (DGD) (unless suitable modifications are made to the  
74 algorithm) also fails to converge to the exact solution of distributed optimization  
75 problems [42]. The second question we seek to address in this paper is: *How should*  
76 *the dynamics of distributed mirror descent be modified, so that convergence to the*  
77 *exact solution is guaranteed?* These guarantees are meant to hold for deterministic  
78 dynamics, but in this paper we will also consider the more general case of stochastic  
79 dynamics with additive noise, where the noise is added to account for corrupted  
80 gradient information, data sub-sampling (as is the case in stochastic gradient descent)  
81 or errors due to the network, such as communication channels being corrupted.

82 **1.1. Previous work.** Distributed optimization has a variety of applications.  
83 Removing the existence of a central server and having the nodes communicate in a  
84 decentralized manner can remove both computational bottlenecks and privacy risks.  
85 A classic reference for distributed optimization is [5], and more recent applications  
86 in statistical learning are described in [8]. The authors in [10] also describe several  
87 interesting applications. A variant of distributed optimization known as federated  
88 learning [25] was proposed recently for solving optimization problems in which the  
89 data is stored across a very large number devices for privacy purposes.

90 The literature on distributed optimization algorithms is vast. Since this paper  
91 focuses on distributed first-order algorithms for convex optimization models, we will  
92 focus on this class of algorithms. Two algorithmic techniques can be used to develop  
93 exact distributed optimization algorithms. The first technique uses diminishing step-  
94 sizes, and the second one relies on gradient tracking. Gradient tracking is closely  
95 related to augmented Lagrangian methods (see Section 3 for more details). Algorithms

96 with diminishing step-sizes tend to be very slow in practice, so recent literature focuses  
 97 on using constant step-sizes. The algorithm we propose in this paper, and its variants,  
 98 are developed in continuous time. The works of (among others) [16, 22, 38, 23, 43]  
 99 also analyze decentralized optimization schemes in continuous time. The works of  
 100 [13, 27, 33, 32] focus on an analysis of the distributed mirror descent algorithm. In  
 101 Table 1, we summarize selected related works that show how this paper fits within  
 102 the existing literature. Current works on exact distributed algorithms, with a fixed  
 103 step-size, are only based on gradient or sub-gradient descent. Table 1 also lists  
 104 earlier proposed exact distributed algorithms in discrete time, which rely on the slower  
 105 mechanism of diminishing step-sizes to achieve exact convergence. Moreover, the mirror  
 106 maps in the existing literature are used only to model accurately the geometry of the  
 107 separable constraints and not the consensus constraint that is the distinguishing feature  
 108 of this paper. Finally, we note mirror descent dynamics are related to Riemannian  
 109 descent as presented in e.g. [2, 11] and preconditioning [17, 1, 34].

110 **1.2. Main results and contributions.** Our results are based on a continuous-  
 111 time analysis of stochastic mirror descent dynamics. Our contributions can be summa-  
 112 rized as follows,

- 113 • In Section 4.1 we show that without strong assumptions on the minimizers  
 114 of each  $f_i$ , the classic distributed stochastic mirror descent formulation with  
 115 constant noise achieves exponential convergence to a neighborhood around  
 116 a different point than the optimum and that the size of this neighborhood  
 117 cannot be reduced using the mirror map or reducing noise.
- 118 • To address the inexactness of the conventional DMD algorithm we propose an  
 119 exact variant called Exact Interacting Stochastic Mirror Descent (EISMD),  
 120 that is able to converge exponentially fast to a much smaller neighborhood  
 121 than the conventional distributed mirror descent (Section 4, Proposition 4.4).
- 122 • We propose a preconditioned version of EISMD, which adapts the mirror  
 123 map based on the geometry of the consensus manifold resulting in even faster  
 124 convergence (Proposition 4.6).
- 125 • In Section 5 we illustrate in detail the performance of our algorithms in  
 126 constrained and unconstrained convex optimization problems.

127 **2. Preliminaries.** In this section we fix our notation, state our main assumptions  
 128 and establish some useful technical lemmas that will be used later.

129 **2.1. Notation.** We use  $\otimes$  to denote the Kronecker product,  $I_d$  the  $d$ -dimensional  
 130 identity matrix and  $\mathbf{1}_d$  denotes the  $d$ -dimensional vector of ones.  $\text{Diag}(a)$  with  $a \in \mathbb{R}^d$   
 131 denotes a matrix with diagonal elements  $[a_1, \dots, a_d]$ . We use  $A$  to denote the  $N \times N$   
 132 weighted adjacency matrix associated with a graph  $\mathcal{G} = (V, E)$ . The graph Laplacian  
 133 is given by  $L := \text{Diag}(A\mathbf{1}_N) - A$  and we use the following notation  $\mathcal{L} := L \otimes I_d$   
 134 with  $\mathcal{L} \in \mathbb{R}^{Nd \times Nd}$  to denote the vectorized version of the graph Laplacian. We use  
 135  $\langle x, y \rangle = x^\top y$  for the standard dot product, and  $\langle x, y \rangle_Q = \langle x, Qy \rangle = x^\top Qy$  for the  
 136  $Q$ -inner product, for some positive definite matrix  $Q$ . We use  $A \succeq B$  to denote a partial  
 137 matrix ordering meaning  $A - B \succeq 0$ . We assume that  $\mathcal{X} \subseteq \mathbb{R}^d$  is a convex set. We use  $\mathcal{D}$   
 138 to denote an open set such that  $\mathcal{X} \subset \text{cl}(\mathcal{D})$ . The set  $\mathcal{D}$  will be used to denote the domain  
 139 of the mirror maps of the mirror descent algorithm. We use  $\mathcal{X}^*$  to denote the dual space  
 140 of  $\mathcal{X}$ . The normal cone of  $\mathcal{X}$  is defined as  $N_{\mathcal{X}}(x) = \{z \in \mathcal{X}^* \mid \langle z, y - x \rangle \leq 0 \ \forall y \in \mathcal{X}\}$ .

141 Given an arbitrary norm  $\|\cdot\|$  on  $\mathbb{R}^d$ , we will define  $B_{\|\cdot\|} := \{v \in \mathbb{R}^d : \|v\| \leq 1\}$ .  
 142 The dual norm  $\|\cdot\|_*$  is defined as  $\|z\|_* := \sup\{\langle z, v \rangle : v \in B_{\|\cdot\|}\}$ . If  $A$  is a matrix  
 143 then  $\|A\|_2$  denotes its spectral norm and we assume that the dual norm is compatible

144 with the spectral norm, i.e.  $\|Az\|_* \leq \|A\|_2 \|z\|_*$ . We will make use of the following  
 145 generalized Cauchy inequality,

$$146 \quad |\langle v, w \rangle| \leq \|v\|_* \|w\| \quad \forall w \in \mathcal{X}, v \in \mathcal{X}^*.$$

148 Since  $0 \leq (\|v\|_* - \|w\|)^2 = \|v\|_*^2 + \|w\|^2 - 2\|v\|_* \|w\|$ , we also have,

$$149 \quad (2.1) \quad \langle v, w \rangle \leq \frac{1}{2}\|v\|_*^2 + \frac{1}{2}\|w\|^2.$$

150 A function  $g$  is said to be  $L$ -Lipschitz continuous with respect to a norm  $\|\cdot\|$  if  
 151  $\|g(x) - g(y)\| \leq L\|x - y\|$ ,  $\forall x, y \in \mathcal{X}$ . The Bregman divergence associated with a  
 152 convex, differentiable function  $g : \mathcal{X} \rightarrow \mathbb{R}$  is defined as follows,

$$153 \quad D_g(x, y) = g(x) - g(y) - \langle \nabla g(y), x - y \rangle.$$

155 If the second-order derivative of  $g$  exists it furthermore holds,

$$156 \quad (2.2) \quad \nabla_x D_g(x, y) = \nabla g(x) - \nabla g(y), \quad \nabla_y D_g(x, y) = \nabla^2 g(y)(y - x).$$

158 The aggregate cost function will be written as  $f(\mathbf{x}) = \sum_{i=1}^N f_i(x^i)$ , where  $\mathbf{x} =$   
 159  $[x^1, \dots, x^N]^T$  denotes the stacked vector of particles and each  $x_i \in \mathcal{X}$ . We will  
 160 use  $(X^*, \Lambda^*)$  to denote the set of primal-dual variables that satisfy the first order  
 161 optimality conditions for (1.1). Unless specified otherwise gradient vectors  $\nabla f$  are  
 162 taken with respect to the joint particle vector  $\mathbf{x}$  following the usual conventions and  
 163 the same applies for Hessian matrices.

164 **REMARK 2.1.** *The space of the Lagrange multipliers for the consensus constraint,*  
 165  $\boldsymbol{\lambda} \in \Lambda \subset \mathbb{R}^{Nd}$ , *will play an important role in the definition of the algorithms below.*  
 166 *We note that the norm associated with  $\boldsymbol{\lambda} \in \Lambda$  will not necessarily be the same as the*  
 167 *one used for the primal variables  $\mathbf{x} \in \mathcal{X}^N$ . We will however use the same notation:*  
 168  *$\|\cdot\|$ , and its dual  $\|\cdot\|_*$  for both spaces, and it will be clear from context which norm*  
 169 *is being used. For  $\mathbf{w} = [\mathbf{x}^T, \boldsymbol{\lambda}^T]^T$  we will use the following mixed norm convention*  
 170  *$\|\mathbf{w}\| = \|\mathbf{x}\| + \|\boldsymbol{\lambda}\|$ , with the understanding that the two norms could be different. For*  
 171 *example, the norm in  $\mathcal{X}^N$  could be the  $\ell_1$ , and in  $\Lambda$  the  $Q$ -norm (for some positive*  
 172 *definite matrix  $Q$ ) so that,  $\|\mathbf{w}\| = \|\mathbf{x}\|_1 + \|\boldsymbol{\lambda}\|_Q$ .*

## 173 2.2. Assumptions and Definitions.

174 **2.2.1. Optimality Conditions and Model Assumptions.** The consensus  
 175 constraint in (1.1) is satisfied if and only if  $\mathcal{L}\mathbf{x} = 0$ , where  $\mathcal{L}$  denotes the vectorized  
 176 graph Laplacian. Therefore the optimality conditions for (1.1) are as follows,

$$177 \quad -\nabla f(\mathbf{x}^*) - \mathcal{L}\boldsymbol{\lambda}^* \in N_{\mathcal{X}}(\mathbf{x}^*).$$

179 If the solution of (1.1) is in the interior of  $\mathcal{X}$ , and if  $f$  is convex, then we must have  
 180 that  $N_{\mathcal{X}}(\mathbf{x}^*) = \{\mathbf{0}\}$  for any  $\mathbf{x}^* \in X^*$ . Because the focus of this paper is on the effect  
 181 of the consensus constraint and its impact on the dynamics of the algorithm, we  
 182 will assume that the optimal solution of (1.1) is in the interior of  $\mathcal{X}$ . Because the  
 183 consensus constraint couples all the particles together, its impact on the algorithm's  
 184 convergence is far less understood than dealing with separable constraints on  $\mathcal{X}$ . For  
 185 certain applications, especially in machine learning, the assumption that the solution  
 186 lies in the interior of the feasible set holds (e.g. [26, 37]). The extension to the general  
 187 case requires some minor technical modifications to our convergence analysis similar  
 188 to [26]. We gather our assumptions so far below.

189 ASSUMPTION 1. Each  $f_i$  in (1.1) is convex and twice differentiable. The elements  
 190 in  $X^*$  are in the interior of  $\mathcal{X}$ .

191 Derivatives of  $f$  are required so that we can apply Itô's formula. Since we assume that  
 192 the function is convex this assumption could be relaxed (see [26] Proposition C.2),  
 193 but the assumption is kept here for simplicity and brevity. We proceed with some  
 194 standard convexity and smoothness definitions.

195 DEFINITION 2.2. We say that  $f : \mathcal{X}^N \rightarrow \mathbb{R}$  is  $\mu$ -strongly convex w.r.t. some norm  
 196  $\|\cdot\|$  provided that  $\|\nabla f(\mathbf{x}) - \nabla f(\mathbf{y})\|_* \geq \mu\|\mathbf{x} - \mathbf{y}\|$ . Similarly, a function  $f$  is  $L_f$ -smooth  
 197 w.r.t. some norm  $\|\cdot\|$  when  $\|\nabla f(\mathbf{x}) - \nabla f(\mathbf{y})\|_* \leq L_f\|\mathbf{x} - \mathbf{y}\|$ .

198 Some of our results will use the notion of relative strong convexity and smoothness.  
 199 We refer the reader to [24] for more properties and [12] for the stochastic case. Below  
 200 we present some definitions and properties that will be useful later on.

DEFINITION 2.3 (Relative strong convexity). A function  $g : \mathcal{X}^N \rightarrow \mathbb{R}$  is  $\mu$ -strongly  
 convex with respect to some convex function  $h$  if for any  $\mathbf{x}, \mathbf{y} \in \mathcal{X}^N$  the following holds,

$$g(\mathbf{x}) \geq g(\mathbf{y}) + \nabla g(\mathbf{y})^T(\mathbf{x} - \mathbf{y}) + \mu D_h(\mathbf{x}, \mathbf{y}).$$

201 Or equivalently,  $\langle \mathbf{x} - \mathbf{y}, \nabla g(\mathbf{x}) - \nabla g(\mathbf{y}) \rangle \geq \mu \langle \mathbf{x} - \mathbf{y}, \nabla h(\mathbf{x}) - \nabla h(\mathbf{y}) \rangle$ .

DEFINITION 2.4 (Relative smoothness). A function  $g : \mathcal{X}^N \rightarrow \mathbb{R}$  is  $\alpha$ -smooth with  
 respect to some function  $h$  if for any  $\mathbf{x}, \mathbf{y} \in \mathcal{X}^N$  the following holds,

$$g(\mathbf{x}) \leq g(\mathbf{y}) + \nabla g(\mathbf{y})^T(\mathbf{x} - \mathbf{y}) + \alpha D_h(\mathbf{x}, \mathbf{y}).$$

202 Or equivalently,  $\langle \mathbf{x} - \mathbf{y}, \nabla g(\mathbf{x}) - \nabla g(\mathbf{y}) \rangle \leq \alpha \langle \mathbf{x} - \mathbf{y}, \nabla h(\mathbf{x}) - \nabla h(\mathbf{y}) \rangle$ .

203 If we assume that  $g$  is  $\mu$ -strongly convex and  $\alpha$ -smooth with respect to  $h$  it holds,

204 
$$\mu D_h(\mathbf{x}, \mathbf{y}) \leq D_g(\mathbf{x}, \mathbf{y}) \leq \alpha D_h(\mathbf{x}, \mathbf{y}).$$

206 We adopt the following definition for the convex conjugate of a relatively strong convex  
 207 function.

208 DEFINITION 2.5 (Convex conjugate). Let  $g : \mathcal{X}^N \rightarrow \mathbb{R}$  be a  $\mu$ -strongly convex  
 209 function with respect to some  $h$ . Then  $g^*(\mathbf{z}) := \max_{\mathbf{x} \in \mathcal{X}^N} \langle \mathbf{z}^T, \mathbf{x} \rangle - g(\mathbf{x})$  is its  
 210 Legendre-Fenchel convex conjugate. When  $g$  is differentiable, we also have  $\nabla g^*(\mathbf{z}) :=$   
 211  $\arg \max_{\mathbf{x} \in \mathcal{X}^N} \langle \mathbf{z}^T, \mathbf{x} \rangle - g(\mathbf{x})$ . and  $\nabla g \circ \nabla g^*(\mathbf{z}) = \mathbf{z}$ .

212 **2.2.2. Network Assumptions.** We first state our assumptions on the network  
 213 topology.

214 ASSUMPTION 2. The graph  $\mathcal{G}$  is connected, undirected and the adjacency matrix  
 215  $A$  is doubly stochastic.

216 These assumptions imply that the graph Laplacian  $\mathcal{L}$  is a real symmetric matrix with  
 217 nonnegative eigenvalues. We will denote the pseudo-inverse of  $\mathcal{L}$  by  $\mathcal{L}^+$  such that,

218 (2.3) 
$$\mathcal{L}\mathcal{L}^+\mathcal{L} = \mathcal{L}.$$

219 We will use the following definition of the  $\beta$ -regularized Laplacian [10],

220 (2.4) 
$$\mathcal{L}_\beta = \mathcal{L} + \frac{\beta}{N} \mathbf{1}_N \mathbf{1}_N^\top \otimes I_d.$$

6

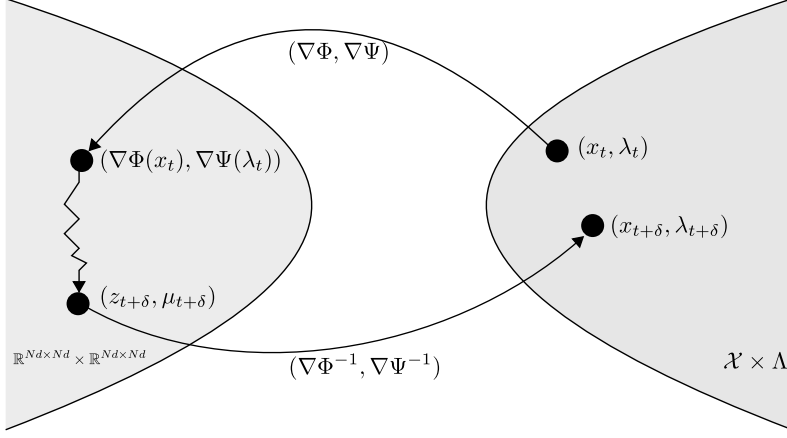


Fig. 2: Stochastic Mirror Descent with two mirror maps.  $\Phi$  maps the primal variables to the dual space, and  $\Psi$  maps the Lagrangian dual variables associated with the consensus constraint to the algebraic dual of the Lagrange multipliers.

221 Note that the  $\beta$ -regularized Laplacian is positive definite. We define the Rayleigh  
 222 quotient associated with the  $\beta$ -regularized Laplacian as follows,

$$223 \quad (2.5) \quad \kappa_{\beta, N} = \max_{\mathbf{d}_x \in \mathbb{R}^{Nd}} \frac{\|\mathcal{L}_\beta \mathbf{d}_x\|_2^2}{\|\mathbf{d}_x\|_2^2}.$$

224 It holds that [10, p.103],

$$225 \quad (2.6) \quad \mathcal{L}_\beta^{-1} = \mathcal{L}^+ + \frac{1}{\beta N} \mathbf{1}_N \mathbf{1}_N^\top \otimes I_d \succeq \mathcal{L}^+,$$

226 where the latter inequality follows from the fact that  $\mathbf{1}_d \mathbf{1}_d^\top \otimes I_N$  is positive semidefinite.

227 LEMMA 2.6. Let Assumption 2 hold and suppose that  $\kappa_{\beta, N}$  is as defined in (2.5)  
 228 then,

$$229 \quad \langle \mathbf{x}, \mathcal{L} \mathbf{x} \rangle \geq \frac{1}{\kappa_{\beta, N}} \|\mathcal{L} \mathbf{x}\|_2^2.$$

230 *Proof.* Using the definition of the pseudo-inverse in (2.3) and its relationship with  
 231 the inverse of the regularized Laplacian in (2.6) we obtain,

$$232 \quad \begin{aligned} \langle \mathbf{x}, \mathcal{L} \mathbf{x} \rangle &= \langle \mathbf{x}, \mathcal{L} \mathcal{L}^+ \mathcal{L} \mathbf{x} \rangle = \langle \mathcal{L} \mathbf{x}, (\mathcal{L}_\beta^{-1} - \frac{1}{\beta N} \mathbf{1}_d \mathbf{1}_d^\top \otimes I_N) \mathcal{L} \mathbf{x} \rangle \\ &= \langle \mathcal{L} \mathbf{x}, \mathcal{L}_\beta^{-1} \mathcal{L} \mathbf{x} \rangle, \end{aligned}$$

233 where in the last equality we used the fact that  $(\mathbf{1}_d \mathbf{1}_d^\top \otimes I_N) \mathcal{L} = 0$ . Since  $\mathcal{L}_{\beta, N} \preceq \kappa_{\beta, N} I$   
 234 then  $\mathcal{L}_{\beta, N}^{-1} \succeq \kappa_{\beta, N}^{-1} I$  and the result follows.  $\square$

235 **2.2.3. Mirror Maps.** The role of the mirror map,  $\Phi : \mathcal{D} \rightarrow \mathbb{R}$ , in the Mirror  
 236 Descent algorithm is to transform the primal  $x \in \mathcal{X}$  variables to the dual space  
 237  $\nabla \Phi(x) \subset \mathbb{R}^d$ . The dual variables will be denoted by  $z$ , i.e  $\nabla \Phi(x) = z$ . In the algorithm

238 proposed in this paper we will use two mirror maps. The first mirror map  $\Phi$ , is used  
 239 to transform the primal variables  $x$ . The second mirror map  $\Psi$ , is used to transform  
 240 the Lagrangian dual variables  $\lambda$  associated with the consensus constraint in (1.1). The  
 241 algebraic dual variables will be denoted by  $\mu$ , i.e.  $\nabla\Psi(\lambda) = \mu$ . When no confusion  
 242 arises between Lagrangian and algebraic dual variables we will refer to them simply as  
 243 dual variables. Figure 2 explains the main steps in mirror descent with the two maps.  
 244 At time-step  $t$  the primal-dual pair  $(x_t, \lambda_t)$  is mapped to  $(z_t, \mu_t) = (\nabla\Phi(x_t), \nabla\Psi(\lambda_t))$ .  
 245 The algorithm then follows the stochastic dynamics specified in Section 3. For example,  
 246 a variant of the proposed scheme performs a gradient descent on the Augmented  
 247 Lagrangian w.r.t the primal  $x$  variables, and a dual ascent w.r.t the Lagrangian dual  
 248 variables  $\lambda$  (see Section 3 for a detailed explanation). The inverse  $(\nabla\Phi^{-1}, \nabla\Psi^{-1})$   
 249 maps the algebraic duals back to the primal space  $\mathcal{X} \times \Lambda$ . The proposed algorithm,  
 250 and its variants, are described in Section 3. Below we state our related assumptions:

251 **ASSUMPTION 3 (Mirror map).**  $\Phi : \mathcal{D} \rightarrow \mathbb{R}$  is twice differentiable,  $\mu_\Phi$ -strongly  
 252 convex and  $L_\Phi$ -smooth w.r.t. some norm  $\|\cdot\|$ . The same holds for  $\Psi : \Lambda \rightarrow \mathbb{R}$   
 253 with constants  $\mu_\Psi, L_\Psi$ , respectively. We furthermore make the additional assumption  
 254  $\nabla\Phi^*(\mathbb{R}^d) = \mathcal{X}$ ,  $\nabla\Psi^*(\mathbb{R}^d) = \Lambda$ ,  $\Phi^*$  is  $L_{\Phi^*}$ -smooth and assume uniform boundedness of  
 255 the Laplacians of  $\Phi^*, \Psi^*$  such that  $\|\Delta\Phi^*\|_\infty, \|\Delta\Psi^*\|_\infty < \infty$ .

256 The assumption that  $\nabla\Phi^*$  maps directly to  $\mathcal{X}$  (and similarly for  $\Psi$ ) avoids the need  
 257 for projections. Extending our results without this assumption is possible by following  
 258 a route similar to [26].

259 A useful property of the Bregman divergence induced by mirror maps that satisfy  
 260 our assumptions is the following,

$$261 \quad (2.7) \quad D_{\Phi^*}(z, z') = D_\Phi(x', x),$$

262 where  $z = \nabla\Phi(x)$  and  $z' = \nabla\Phi(x')$ . For  $x, y, z \in \mathbb{R}^d$  we have the triangle property for  
 263 Bregman divergences (see Lemma 9.11 in [4])

$$264 \quad (2.8) \quad \langle x - y, \nabla\Phi(z) - \nabla\Phi(y) \rangle = D_\Phi(x, y) + D_\Phi(y, z) - D_\Phi(x, z).$$

266 We also make use of the following property,

$$267 \quad (2.9) \quad D_f(\mathbf{x}, \mathbf{x}') \leq \alpha(\Phi) D_\Phi(\mathbf{x}', \mathbf{x})$$

268 where  $\alpha(\Phi) = \frac{L_f L_\Phi}{\mu_\Phi}$ . This property follows from the relative smoothness assumption  
 269 combined with the strong convexity and Lipschitz assumption on  $\Phi$ ,

$$270 \quad D_f(\mathbf{x}, \mathbf{x}') \leq L_f D_\Phi(\mathbf{x}, \mathbf{x}') \leq \frac{L_f L_\Phi}{2} \|\mathbf{x}' - \mathbf{x}\|^2 \leq \frac{2L_f L_\Phi}{2\mu_\Phi} D_\Phi(\mathbf{x}', \mathbf{x}),$$

272 where with slight abuse we denote  $\Phi(\mathbf{x}) = \sum_{i=1}^N \Phi(x^i)$ . We will use the following  
 273 Rayleigh quotient,

$$274 \quad (2.10) \quad \kappa_N = \max_{\mathbf{d}_x \in \mathbb{R}^{Nd}, \mathbf{d}_\lambda \in \mathbb{R}^{Nd}} \max \left\{ \frac{\|\mathcal{L}^{\frac{1}{2}} \mathbf{d}_x\|_2^2}{\|\mathbf{d}_x\|^2}, \frac{\|\mathcal{L}^{\frac{1}{2}} \mathbf{d}_\lambda\|_2^2}{\|\mathbf{d}_\lambda\|^2} \right\}.$$

275 Note that the norms for  $\mathbf{d}_x$  and  $\mathbf{d}_\lambda$  in the definition above may be different (see  
 276 Remark 2.1). We will also need the following generalized Rayleigh quotient,

$$277 \quad (2.11) \quad \kappa_g = \inf_{\mathbf{x}, \mathbf{d}_x, \mathbf{d}_\lambda \in \mathbb{R}^{Nd}} \frac{\|A(\mathbf{x})[\mathbf{d}_x^T, \mathbf{d}_\lambda^T]^\top\|_{\nabla^2 \Phi^*(\mathbf{z})}^2}{\|\mathbf{d}_x\|^2 + \|\mathbf{d}_\lambda\|^2},$$



278 where  $A(\mathbf{x}) = [\nabla^2 f(\mathbf{x}) + \mathcal{L}, \mathcal{L}] \in \mathbb{R}^{Nd \times 2Nd}$ . If strong convexity is assumed then  $\kappa_g$   
 279 is strictly positive. This fact is not obvious since  $A(\mathbf{x})$  is not a square matrix, and  
 280 the norm used in the definition of (2.11) is not standard, we therefore provide a short  
 281 proof below.

282 **LEMMA 2.7.** *Suppose Assumptions 1-3 hold and that  $f$  is relatively strongly convex*  
 283 *with respect to  $\Phi$ , then  $\kappa_g$  defined in (2.11) is positive.*

284 *Proof.* We note that  $A(\mathbf{x})$  can be obtained by removing the last  $Nd$  columns and  
 285 rows of the following matrix,

$$286 \quad B(\mathbf{x}) := \begin{bmatrix} \nabla^2 f(\mathbf{x}) + \mathcal{L} & \mathcal{L} \\ -\mathcal{L} & 0 \end{bmatrix}.$$

288 Let  $\mathbf{d} = [\mathbf{d}_x, \mathbf{d}_\lambda]^\top$  and note that  $\langle \mathbf{d}, B(\mathbf{x})\mathbf{d} \rangle = \mathbf{d}_x^\top \nabla^2 f(\mathbf{x})\mathbf{d}_x$ . It follows from  
 289 the relative strong convexity assumption that  $B(\mathbf{x}) \succeq \mu_\Phi \nabla^2 \Phi(\mathbf{x})$  and therefore  
 290  $\|B(\mathbf{x})\|_{\nabla^2 \Phi(\mathbf{z})^{-1}}^2 \succ 0$ . Since  $A(\mathbf{x})$  can be obtained by removing the last  $Nd$  columns  
 291 and rows of  $B(\mathbf{x})$  the result follows from the interlacing theorem for singular values,  
 292 see e.g. Theorem 3.1.3 in [18].  $\square$

293 Lastly, we will need the following result.

294 **LEMMA 2.8.** *Suppose that Assumptions 1-3. Then for an arbitrary optimal primal*  
 295 *dual pair  $(\mathbf{x}^*, \boldsymbol{\lambda}^*)$  we have*

$$296 \quad \begin{aligned} \|\nabla f(\mathbf{x}) + \mathcal{L}\boldsymbol{\lambda} + \mathcal{L}\mathbf{x}\|_{\nabla \Phi^*(\mathbf{z})} &\geq \frac{2\kappa_g}{\hat{\mu}} \left( \sum_{i=1}^N D_\Phi(x^*, x^i) + D_\Psi(\lambda^*, \lambda^i) \right) \\ &= \frac{2\kappa_g}{\hat{\mu}} \left( \sum_{i=1}^N D_{\Phi^*}(z^i, z^*) + D_{\Psi^*}(\mu^i, \mu^*) \right), \end{aligned}$$

298 where  $\hat{\mu} = \min\{\mu_\Phi, \mu_\Psi\}$

299 *Proof.* Since  $f$  is twice differentiable there exists an  $\mathbf{y}$  on the line segment joining  
 300  $\mathbf{x}$  and  $\mathbf{x}^*$  such that  $\nabla f(\mathbf{x}) - \nabla f(\mathbf{x}^*) = \langle \nabla^2 f(\mathbf{y}), \mathbf{x} - \mathbf{x}^* \rangle$ . We then have,

$$301 \quad \begin{aligned} \|\nabla f(\mathbf{x}) + \mathcal{L}\boldsymbol{\lambda} + \mathcal{L}\mathbf{x}\|_{\nabla^2 \Phi^*(\mathbf{z})}^2 &= \|\nabla f(\mathbf{x}) - \nabla f(\mathbf{x}^*) + \mathcal{L}(\boldsymbol{\lambda}_t - \boldsymbol{\lambda}^*) + \mathcal{L}(\mathbf{x} - \mathbf{x}^*)\|_{\nabla^2 \Phi^*(\mathbf{z})}^2 \\ 302 \quad &= \|A(\mathbf{y})[\mathbf{x}_t - \mathbf{x}^*, \boldsymbol{\lambda} - \boldsymbol{\lambda}^*]^\top\|_{\nabla^2 \Phi^*(\mathbf{z})}^2 \\ 303 \quad &\geq \kappa_g (\|\mathbf{x}^* - \mathbf{x}_t\|^2 + \|\boldsymbol{\lambda}^* - \boldsymbol{\lambda}\|^2) \\ 304 \quad &\geq \frac{2\kappa_g}{\hat{\mu}} \left( \sum_{i=1}^N D_\Phi(x^*, x^i) + D_\Psi(\lambda^*, \lambda^i) \right). \end{aligned}$$

306 We use (2.7) to obtain the bound in terms of the (algebraic) dual variables.  $\square$

307 **3. Distributed Stochastic Mirror Descent: Exact and Preconditioned**  
 308 **Dynamics.** In this section we introduce different variants of distributed MD algo-  
 309 rithms. We adopt a dynamical systems point of view for our analysis. Numerical  
 310 realizations of the proposed schemes are discussed in Section 5. For an introduction to  
 311 the original mirror descent algorithm we refer the interested reader to [4, Ch. 9].

312 *Interacting Stochastic Mirror Descent (ISMD).* The starting point of our analysis  
 313 is the Interacting Stochastic Mirror Descent (ISMD) algorithm. This algorithm was  
 314 proposed in [31] but its convergence was only established in the linear case. The case

315 where all the functions are identical (i.e.  $f_1 = \dots = f_N$ ) and strongly convex was  
 316 analyzed in [6]. The discrete time version of the algorithm for the general convex case  
 317 was discussed in [13], but exact convergence was only established under a diminishing  
 318 step-size strategy. In the continuous time setting the dynamics of ISMD are as follows,

$$319 \quad (3.1) \quad dz_t^i = -\eta \nabla f_i(x_t^i) dt + \epsilon \sum_{j=1}^N A_{ij} (z_t^j - z_t^i) dt + \sigma dB_t^i, \quad x_t^i = \nabla \Phi^*(z_t^i),$$

320 for particles  $i = 1, \dots, N$ , and where  $B_t^i$  are independent Brownian motions. The matrix  
 321  $A = \{A_{ij}\}_{i,j=1}^N$  is an  $N \times N$  doubly-stochastic matrix representing the interaction  
 322 weights and  $\eta, \epsilon$  are tuning constants representing the learning rate and interaction  
 323 strength, respectively. For simplicity in most of the subsequent analysis we set  $\eta = \epsilon =$   
 324 1, but it is straightforward to extend the results for arbitrary values of  $\eta$  and  $\epsilon$ . In the  
 325 context of modern large scale applications, we note that understanding convergence  
 326 under the presence of noise is often motivated from computational considerations such  
 327 as when sub-sampling the gradient  $f$  or the interaction graph when  $N$  is large.

328 Using the graph Laplacian, we can rewrite the evolution in vector form as

$$329 \quad (3.2) \quad d\mathbf{z}_t = (-\nabla f(\mathbf{x}_t) - \mathcal{L}\mathbf{z}_t) dt + \sigma d\mathbf{B}_t, \quad \mathbf{x}_t = \nabla \Phi^*(\mathbf{z}_t),$$

331 where  $\mathbf{B}_t := ((B_t^1)^T, \dots, (B_t^N)^T)^T$ . In the case where the mirror map above is the  $\ell_2$   
 332 norm it is known that even in the deterministic case the dynamics in (3.2) will not  
 333 converge to the exact solution (see [36, 35]). In Section 4.1 we show that in general  
 334 the dynamics in (3.2) also fails to converge to the exact solution of (1.1) and identify  
 335 that this can occur only under additional assumptions. This motivates proposing a  
 336 different dynamics below.

337 *Exact Interacting Stochastic Mirror Descent (EISMD)*. To address the limitations  
 338 of the ISMD algorithm discussed above we propose the following,

$$339 \quad (3.3) \quad \begin{aligned} d\mathbf{z}_t &= -\nabla f(\mathbf{x}_t) dt - \mathcal{L}\mathbf{x}_t dt - \mathcal{L}\boldsymbol{\lambda}_t dt + \sigma d\mathbf{B}_t, \\ d\boldsymbol{\lambda}_t &= \mathcal{L}\mathbf{x}_t dt, \end{aligned}$$

340 with  $\mathbf{x}_t = \nabla \Phi^*(\mathbf{z}_t)$  and initial conditions  $\boldsymbol{\lambda}_0 = 0$ ,  $\mathbf{z}_0 = \nabla \Phi(\mathbf{x}_0)$ . The idea behind this  
 341 method is to add historical feedback into the algorithm through the integral  $\int_0^t \mathcal{L}\mathbf{x}_s ds$ .  
 342 At optimality this will cancel out the gradient term  $\nabla f(\mathbf{x}_t)$ . In Section 3.1 we show  
 343 that the drift term in EISMD is related to the Augmented Lagrangian. We exploit  
 344 this connection in the theoretical analysis in Section 4. Compared to previous works  
 345 considered in Table 1 this algorithm integrates past information into the dynamics  
 346 through the integral term and is applicable to the mirror descent framework. For  
 347 the case where  $\sigma = 0$ , (3.3) has been considered in [38]. Here we extend the ideas  
 348 in [38] to allow  $f$  being only convex, adding Brownian noise and considering general  
 349 preconditioning.

350 *Exact Preconditioned Interacting Stochastic Mirror Descent (EPISMD)*. A poten-  
 351 tial limitation of ISMD in (3.3) is that the mirror map  $\Phi$  only captures the geometry  
 352 of the primal space  $\mathcal{X} \subset \mathbb{R}^d$ . Even if  $\mathcal{X} = \mathbb{R}^d$ , our optimization problem is still con-  
 353 strained to the consensus manifold  $\mathcal{X}^C$ . In order to incorporate information from the  
 354 consensus constraint we introduce a second mirror map  $\Psi$  that acts on the Lagrangian  
 355 dual variable (see Figure 8). The preconditioned dynamics (EPISMD) is given by,

$$356 \quad (3.4) \quad \begin{aligned} d\mathbf{z}_t &= -\nabla f(\mathbf{x}_t) dt - \mathcal{L}\mathbf{x}_t dt - \mathcal{L}\boldsymbol{\lambda}_t dt + \sigma d\mathbf{B}_t, \\ d\boldsymbol{\mu}_t &= \mathcal{L}\mathbf{x}_t dt, \quad \boldsymbol{\lambda}_t = \nabla \Psi^*(\boldsymbol{\mu}_t), \end{aligned}$$

357 where  $\boldsymbol{\mu}$  is the mirrored version of the  $\boldsymbol{\lambda}$  variable using the mirror map  $\Psi$ . As we will  
 358 show in Section 3.1, this algorithm is related to preconditioning  $\boldsymbol{\lambda}_t$  and we will later  
 359 show numerically that it can lead to faster convergence.

360 **3.1. Preconditioning and the Augmented Lagrangian.** The Augmented  
 361 Lagrangian for the standard gradient descent setting is well-known (see e.g. [28, 40, 16]).  
 362 The Alternating Direction Method of Multipliers (ADMM) is based on an Augmented  
 363 Lagrangian with a Bregman divergence [39, 41] and Riemannian primal-dual methods  
 364 over the Augmented Lagrangian were considered in [2] (see also [15] for a continuous  
 365 time analysis of ADMM). Below we discuss the relationship between the different  
 366 variants of the proposed methods.

367 Consider the Augmented Lagrangian,

$$368 \quad (3.5) \quad L(\mathbf{x}, \boldsymbol{\lambda}) = f(\mathbf{x}) + \langle \mathcal{L}\mathbf{x}, \boldsymbol{\lambda} \rangle + \frac{1}{2} \|\mathcal{L}^{\frac{1}{2}}\mathbf{x}\|_2^2.$$

370 The Augmented Lagrangian Method (ALM) proceeds by a descent step in the primal  
 371 variables and an ascent step in the Lagrangian dual variables. When Bregman  
 372 divergence is used to define the ALM, then (in discrete time) the iterates are given by  
 373 the following,

$$374 \quad (3.6) \quad \begin{aligned} \mathbf{x}_{t+\Delta t} &= \operatorname{argmin}_{\mathbf{d}} \Delta t \langle \nabla_{\mathbf{x}} L(\mathbf{x}_t, \boldsymbol{\lambda}_t), \mathbf{d} \rangle + D_{\Phi}(\mathbf{d}, \mathbf{x}_t), \\ \boldsymbol{\lambda}_{t+\Delta t} &= \operatorname{argmax}_{\mathbf{d}} \Delta t \langle \nabla_{\boldsymbol{\lambda}} L(\mathbf{x}_t, \boldsymbol{\lambda}_t), \mathbf{d} \rangle - \frac{1}{2} \|\mathbf{d} - \boldsymbol{\lambda}_t\|^2. \end{aligned}$$

375 Writing down the optimality conditions of the two subproblems above and taking the  
 376 limit  $\Delta t \rightarrow 0$ , we obtain the deterministic version of (3.3).

377 Alternatively, we can rewrite the dynamics of ISMD in (3.2) in terms of  $\mathbf{x}_t$  drawing  
 378 a connection with preconditioned or Riemannian gradient descent [26]. Note that  
 379  $\nabla \Phi(\nabla \Phi^*(\mathbf{z})) = \mathbf{z}$ . Differentiating this w.r.t.  $z$  we obtain  $\nabla^2 \Phi^*(\mathbf{z}) \nabla^2 \Phi(\nabla \Phi^*(\mathbf{z})) = I_{dN}$ .  
 380 Therefore,

$$381 \quad (3.7) \quad \nabla^2 \Phi^*(\mathbf{z}) = \nabla^2 \Phi(\mathbf{x})^{-1}.$$

383 By applying Itô's lemma to  $\nabla \Phi^*(\mathbf{z}_t)$ , using the definition of the Bregman divergence  
 384 and properties (2.2) and (3.7),

$$385 \quad \begin{aligned} dx_t^i &= -\nabla^2 \Phi(x_t^i)^{-1} \nabla f_i(x_t^i) dt - \nabla^2 \Phi(x_t^i)^{-1} \sum_{j=1}^N A_{ij} (z_t^i - z_t^j) dt \\ &\quad + \frac{1}{2} \sigma^2 \nabla^2(\nabla \Phi^*(z_t^i)) dt + \sigma \nabla^2 \Phi(x_t^i)^{-1} dB_t^i. \end{aligned}$$

388 From this expression we see that mirror descent is a preconditioned algorithm where  
 389 the choice of preconditioner is determined through the function  $\Phi$ . Preconditioned  
 390 dynamics have been studied in previous work to improve communication complexity in  
 391 a distributed setting [17] or speed up mixing rates of the dynamics [20]. The additional  
 392 drift term with the third-order gradient of the mirror map arises as a correction term  
 393 due to the nonlinearity of the mirror map. A similar scheme can be derived for EISMD.

394 *Preconditioning of the Interaction.* The choice of mirror map can be used to  
 395 precondition the dynamics and accelerate convergence. Now we show, based on  
 396 the primal-dual interpretation, that the dynamics in (3.4) can also be cast into the  
 397 preconditioned setting, where preconditioning is done on the  $\boldsymbol{\lambda}_t$  variable. Motivated

398 by the augmented Lagrangian and preconditioning, as opposed to the dynamics in  
 399 (3.6), one could redefine the dual variable dynamics as,

$$400 \quad \lambda_{t+\Delta t} = \operatorname{argmax}_{\mathbf{d}} \Delta t \langle \nabla_{\lambda} L(\mathbf{x}_t, \lambda_t), \mathbf{d} \rangle - D_{\Psi}(\mathbf{d}, \lambda).$$

402 The first order optimality conditions of the problem above are,

$$403 \quad \nabla \Psi(\lambda_{t+\Delta t}) = \nabla \Psi(\lambda_t) + \Delta t \mathcal{L} \mathbf{x}_t.$$

405 Taking  $\Delta t \rightarrow 0$  we obtain,

$$406 \quad (3.8) \quad d\lambda_t = \nabla^2 \Psi(\lambda_t)^{-1} \mathcal{L} \mathbf{x}_t dt,$$

408 where we used  $\nabla^2 \Psi^*(\mu) = \nabla^2 \Psi(\lambda)^{-1}$  as in (3.7). Therefore, the method in (3.4)  
 409 allows for additional flexibility due to the preconditioning of the dual variable.

410 **4. Convergence Analysis.** In this section we present a convergence analysis  
 411 for the exact interacting mirror descent algorithm. EISMD with a strongly convex  
 412 objective is able to converge *exponentially* fast to an area of the optimum, however  
 413 the size of this area can be made arbitrarily small by decreasing  $\sigma$ .

414 **4.1. When first-order optimization fails.** Our first result shows that if there  
 415 exists an  $x^*$  such that  $\nabla f_i(x^*) = 0$  for all  $i = 1, \dots, N$ , exact consensus can be obtained  
 416 for ISMD.

417 **LEMMA 4.1.** *Let Assumptions 1-3 hold. Consider the dynamics in (3.2) with*  
 418  *$\sigma = 0$ . If*

$$419 \quad (4.1) \quad \bigcap_{i=1}^N \{\nabla f_i(x) = 0\} \neq \emptyset,$$

421 *then  $\lim_{t \rightarrow \infty} x_t^i = x^*$ .*

422 *Proof.* Let  $x_0$  be the initial point of the algorithm, and let  $z^*$  be an optimal (dual)  
 423 point closest to  $z_0 = \nabla \Phi(x_0)$  with respect to the divergence generated by  $\Phi^*$ ,

$$424 \quad z^* = \operatorname{argmin}_{z \in Z^*} D_{\Phi^*}(z, z_0)$$

425 where  $Z^* = \{z \mid z = \nabla \Phi(x), \exists x \in \mathcal{X} : \nabla f_i(x) = 0, i = 1, \dots, N\}$ . By assumption (4.1),  
 426  $Z^*$  is not empty. With a slight abuse of notation we let  $x^* = \nabla \Phi^*(z^*)$  and note that  
 427  $(x^*, z^*)$  is an equilibrium point for (3.2) (for a strongly convex function it is also the  
 428 unique equilibrium point, but here we only assume convexity of  $f$ ).

429 Define the Lyapunov candidate function  $V_t = \sum_{i=1}^N D_{\Phi^*}(z_t^i, z^*)$ . Then we obtain,

$$430 \quad dV_t = \sum_{i=1}^N (x^* - x_t^i)^T \nabla f_i(x_t^i) dt + \sum_{i=1}^N (x_t^i - x^*)^T \sum_{j=1}^N A_{ij} (z_t^j - z_t^i) dt.$$

432 Under convexity of  $f$  and optimality at  $x^*$  we have

$$433 \quad \sum_{i=1}^N (x^* - x_t^i)^T \nabla f_i(x_t^i) \leq \sum_{i=1}^N (f_i(x^*) - f_i(x_t^i)) \leq 0.$$

435 By the triangle equality of the Bregman divergence in (2.8),

$$\begin{aligned}
436 \quad (x_t^i - x^*)^T (z_t^j - z_t^i) &= -(x^* - x_t^i)^T (z_t^j - z_t^i) \\
437 \quad &= -(\nabla\Phi^*(z^*) - \nabla\Phi^*(z_t^i))^T (z_t^j - z_t^i) \\
438 \quad &= -D_{\Phi^*}(z_t^j, z_t^i) - D_{\Phi^*}(z_t^i, z^*) + D_{\Phi^*}(z_t^j, z^*).
\end{aligned}$$

440 Then,

$$\begin{aligned}
441 \quad (4.2) \quad \sum_{i=1}^N \sum_{j=1}^N A_{ij} (x_t^i - x^*)^T (z_t^j - z_t^i) \\
442 \quad \quad \quad = \sum_{i=1}^N \sum_{j=1}^N A_{ij} \left( -D_{\Phi^*}(z_t^j, z_t^i) - D_{\Phi^*}(z_t^i, z^*) + D_{\Phi^*}(z_t^j, z^*) \right) \leq 0, \\
443
\end{aligned}$$

444 where we have used  $A_{ij} \geq 0$ ,  $\sum_{i=1}^N \sum_{j=1}^N D_{\Phi^*}(z_t^i, z^*) = \sum_{i=1}^N \sum_{j=1}^N D_{\Phi^*}(z_t^j, z^*)$ , and  
445  $D_{\Phi^*}(z_t^j, z_t^i) \geq 0$ . Since  $V_t > 0$  for  $\mathbf{z} \neq \mathbf{1}_N \otimes z^*$ ,  $V_t = 0$  when  $\mathbf{z} = \mathbf{1}_N \otimes z^*$  and  $dV_t \leq 0$   
446 with equality only at  $\mathbf{z} = \mathbf{1}_N \otimes z^*$  we conclude that  $V_t$  is a Lyapunov function for  $\mathbf{z}_t$ .  
447 Since  $D_{\Phi^*}(z_t^i, z^*) = D_{\Phi}(x^*, x_t^i)$  the statement follows.  $\square$

448 The Lemma above can be extended to an if and only if statement based on the  
449 arguments of [35, Theorem 1], but precise details lie beyond the scope of this paper. If  
450 (4.1) is violated, even with the right choice of mirror map, achieving exact consensus  
451 is not possible. In general imposing  $x^*$  to satisfy (4.1) is quite restrictive as  $\nabla f(\mathbf{x}) =$   
452  $\sum_{i=1}^N \nabla f_i(x^*) = 0$  does not necessarily imply  $\nabla f_i(x^*) = 0$  for all  $i = 1, \dots, N$ . The  
453 crucial point to realise here is that if and only if (4.1) holds then  $(\mathbf{x}^*, \mathbf{z}^*)$  will also be  
454 the minimizer of  $f(\mathbf{x}) + \frac{1}{2} \mathbf{z}^T \mathcal{L} \mathbf{z}$ ; see [35, Lemma 7] for details. As a result one can  
455 establish consensus at equilibrium and  $V_t$  will approach zero at large  $t$ . If  $\mathbf{x}^*$  does  
456 not satisfy (4.1) and one has just  $\nabla f(\mathbf{x}^*) = 0$ , the arguments above can be used to  
457 establish exponential but *approximate* convergence for (3.2).

458 PROPOSITION 4.2 (Approximate convergence of (3.2)). *Let Assumptions 1-3*  
459 *hold and assume that  $f$  is  $\mu_f$ -strongly convex w.r.t.  $\Phi$ . Let  $\mathbf{x}^\dagger = \arg \min \{f(\mathbf{x}) +$   
460  $\frac{1}{2} \nabla \Phi(\mathbf{x})^T \mathcal{L} \nabla \Phi(\mathbf{x})\}$  with  $\mathbf{x}^\dagger = \mathbf{1}_N \otimes x^\dagger$  and  $V_t = \frac{1}{N} \sum_{i=1}^N D_{\Phi^*}(z_t^i, z^\dagger)$ , where  $z_t^i$  obeys  
461 the dynamics of (3.2). Then*

$$\begin{aligned}
462 \quad \mathbb{E}[V_t] &\leq e^{-\mu_f t} \frac{1}{N} \sum_{i=1}^N D_{\Phi^*}(z_0^i, z^\dagger) + \frac{\sigma^2}{2\mu_f} \|\Delta\Phi^*\|_\infty + \frac{C_f}{\mu_f}, \\
463
\end{aligned}$$

464 where  $C_f \geq 0$  is a constant depending on  $f$ .

465 The proof and details are in Appendix A.2. While the relative strong convexity of  
466 the objective function can speed up convergence, only approximate convergence can  
467 be obtained even with  $\sigma = 0$ . In this setup the additional preconditioning via the  
468 mirror map does not facilitate exact convergence nor consensus. When (4.1) holds the  
469 arguments above can be used to show that  $C_f = 0$  thus achieving exact convergence.

470 **4.2. Exact Interacting Stochastic Mirror Descent Analysis.** In this sec-  
471 tion we show that the EISMD algorithm in (3.3) allow us to converge close to the  
472 optimum and this convergence is exact when  $\sigma = 0$ .

473 We note that  $x^i = x^j$  for  $(i, j) \in E$  if and only if  $\mathcal{L}\mathbf{x} = 0$ , therefore the problem in  
 474 (1.1) can be written as,

$$475 \quad \min_{\mathbf{x} \in \mathcal{X}^N} \quad f(\mathbf{x}) + \frac{1}{2} \|\mathcal{L}^{\frac{1}{2}} \mathbf{x}\|_2^2$$

$$\quad \text{s.t.} \quad \mathcal{L}\mathbf{x} = 0.$$

476 The application of the Karush–Kuhn–Tucker (KKT) conditions to the problem above  
 477 implies that if  $(\mathbf{x}^*, \boldsymbol{\lambda}^*) = (\mathbf{1}_N \otimes x^*, \mathbf{1}_N \otimes \lambda^*)$  is an arbitrary point that satisfies the first  
 478 order optimality conditions for (1.1), then when  $\sigma = 0$ ,  $(\mathbf{x}^*, \boldsymbol{\lambda}^*)$  is also an equilibrium  
 479 point of (3.3). The connection of the dynamics of (3.3) and the augmented Lagrangian  
 480 is key to the convergence analysis below.

481 The analysis of the algorithm in (3.3) is based on the following Lyapunov function,

$$482 \quad (4.3) \quad V(\mathbf{x}, \boldsymbol{\lambda}) = c(V_1(\mathbf{x}) + V_2(\boldsymbol{\lambda})) + V_3(\mathbf{x}, \boldsymbol{\lambda}),$$

483 where,

$$484 \quad V_1(\mathbf{x}) = \sum_{i=1}^N D_{\Phi}(x^*, x^i), \quad V_2(\boldsymbol{\lambda}) = \frac{1}{2} \|\boldsymbol{\lambda} - \boldsymbol{\lambda}^*\|_2^2,$$

$$485 \quad V_3(\mathbf{x}, \boldsymbol{\lambda}) = D_f(\mathbf{x}, \mathbf{x}^*) + \langle \mathbf{x} - \mathbf{x}^*, \mathcal{L}(\boldsymbol{\lambda} - \boldsymbol{\lambda}^*) \rangle + \frac{1}{2} \|\mathcal{L}^{\frac{1}{2}} \mathbf{x}\|_2^2,$$

487 and  $c \geq \underline{c}$  with  $\underline{c} > 0$  to be specified below for different contexts. In the case that  $f$   
 488 is only convex and thus multiple minimizers might exist, then we define the optimal  
 489 primal dual pair,  $(x^*, \lambda^*)$  to be the one that is closest to the initial conditions,

$$490 \quad (x^*, \lambda^*) = \arg \min_{x, \lambda \in (X^*, \Lambda^*)} D_{\Phi}(x, x_0) + \|\lambda - \lambda_0\|_2^2.$$

491 Below we establish upper and lower bounds for (4.3) that will be useful later on.

492 LEMMA 4.3. *Let Assumptions 1-3 hold. Then (4.3) satisfies the following,*

493 **i**  $V(\mathbf{x}^*, \boldsymbol{\lambda}^*) = 0$ .

494 **ii.a** *If  $c \geq \max\{\kappa_N/\mu_{\Phi}, \kappa_N\}$  then*

$$495 \quad V(\mathbf{x}, \boldsymbol{\lambda}) \geq \frac{1}{2}(\mu_{\Phi}c - \kappa_N)\|\mathbf{x} - \mathbf{x}^*\|^2 + \frac{1}{2}(c - \kappa_N)\|\boldsymbol{\lambda} - \boldsymbol{\lambda}^*\|_2^2 \geq 0.$$

497 **iii.a** *Let  $\hat{\mu} = \min\{\mu_{\Phi}, 2\}$ . Then,*

$$498 \quad (4.4) \quad V(\mathbf{x}, \boldsymbol{\lambda}) \leq \left( c + \frac{3\kappa_N + 2\alpha(\Phi)}{\hat{\mu}} \right) \left( \sum_{i=1}^N D_{\Phi}(x^*, x_i) + \|\boldsymbol{\lambda} - \boldsymbol{\lambda}^*\|_2^2 \right),$$

499 where  $\alpha(\Phi) = L_f L_{\Phi} / \mu_{\Phi}$  was defined in (2.9), and  $\kappa_N$  in (2.10).

500 *If, in addition,  $f$  is  $\mu_f$ -strongly convex relative to  $\Phi$  then:*

501 **ii.b** *For any  $c \geq \max\{(\kappa_N - \mu_f L_{\Phi})/\mu_{\Phi}, \kappa_N\}$ ,*

$$502 \quad V(\mathbf{x}, \boldsymbol{\lambda}) \geq \frac{1}{2}(\mu_{\Phi}c + \mu_f L_{\Phi} - \kappa_N)\|\mathbf{x} - \mathbf{x}^*\|^2 + \frac{1}{2}(c - \kappa_N)\|\boldsymbol{\lambda} - \boldsymbol{\lambda}^*\|_2^2 \geq 0.$$

503 *Proof.* Property (i) is obvious. For (ii.a) we bound  $V_1$  using the strong convexity  
 504 of  $\Phi$ :

$$505 \quad V_1(\mathbf{x}) = \sum_{i=1}^N D_{\Phi}(x^*, x) \geq \frac{\mu_{\Phi}}{2} \|x^* - x\|^2.$$

506

507 We note that the convexity of  $f$  implies that  $D_f(\mathbf{x}, \mathbf{x}^*) \geq 0$ , and we bound  $V_3$  as  
 508 follows,

$$509 \quad V_3(\mathbf{x}, \boldsymbol{\lambda}) \geq \langle \mathbf{x} - \mathbf{x}^*, \mathcal{L}(\boldsymbol{\lambda} - \boldsymbol{\lambda}^*) \rangle \geq -\frac{1}{2}(\|\mathcal{L}^{\frac{1}{2}}\mathbf{x} - \mathbf{x}^*\|_2^2 + \|\mathcal{L}^{\frac{1}{2}}\boldsymbol{\lambda} - \boldsymbol{\lambda}^*\|_2^2)$$

$$510 \quad \geq -\frac{\kappa_N}{2}(\|\mathbf{x} - \mathbf{x}^*\|^2 + \|\boldsymbol{\lambda} - \boldsymbol{\lambda}^*\|_2^2)$$

512 where in the second inequality we used (2.1) and in the third one (2.10).

513 If, in addition,  $f$  is strongly convex relative to  $\Phi$  then,

$$514 \quad \sum_{i=1}^N D_f(x^*, x^i) \geq \mu_f \sum_{i=1}^N D_\Phi(x^*, x^i) \geq \frac{\mu_f L_\Phi}{2} \|\mathbf{x}^* - \mathbf{x}\|^2$$

516 Using the preceding inequality to bound  $V_3$  we obtain the bound (ii.b).

517 For the upper bound in (iii.a) we bound the the first term in  $V_3$  using the symmetry  
 518 bound in (2.9),

$$519 \quad \sum_{i=1}^N D_f(x^*, x_i) \leq \alpha(\Phi) \sum_{i=1}^N D_\Phi(x^*, x_i).$$

520 For the second term in  $V_3$  we use (2.1) again and for any  $\gamma > 0$ ,

$$\begin{aligned} \frac{1}{2} \langle \mathbf{x} - \mathbf{x}^*, \mathcal{L}(\boldsymbol{\lambda} - \boldsymbol{\lambda}^*) \rangle &\leq \frac{\kappa_N}{2} (\gamma \|\mathbf{x} - \mathbf{x}^*\|^2 + \frac{1}{\gamma} \|\boldsymbol{\lambda} - \boldsymbol{\lambda}^*\|_2^2) \\ &\leq \frac{\kappa_N \gamma}{\mu_\Phi} \sum_{i=1}^N D_\Phi(x^*, x^i) + \frac{\kappa_N}{2\gamma} \|\boldsymbol{\lambda} - \boldsymbol{\lambda}^*\|_2^2 \\ 521 \quad &\leq \frac{\kappa_N}{\hat{\mu}} (\gamma \sum_{i=1}^N D_\Phi(x^*, x^i) + \frac{1}{\gamma} \|\boldsymbol{\lambda} - \boldsymbol{\lambda}^*\|_2^2) \\ &\leq \frac{\kappa_N}{\hat{\mu}} \sum_{i=1}^N D_\Phi(x^*, x^i) + \frac{\kappa_N + \alpha(\Phi)}{\hat{\mu}} \|\boldsymbol{\lambda} - \boldsymbol{\lambda}^*\|_2^2, \end{aligned}$$

522 where in the second inequality we used the relative strong convexity of  $\Phi$  and for the  
 523 last inequality we set  $\gamma = \frac{\kappa_N}{\kappa_N + \alpha(\Phi)} \leq 1$ . Finally, for the last term in  $V_3$  we use the  
 524 bound from (2.10), the strong convexity of  $\Phi$  and the definition of  $\hat{\mu}$ ,

$$525 \quad \|\mathcal{L}^{\frac{1}{2}}\mathbf{x}\|_2^2 = \|\mathcal{L}^{\frac{1}{2}}(\mathbf{x} - \mathbf{x}^*)\|_2^2 \leq \kappa_N \|\mathbf{x} - \mathbf{x}^*\|^2 \leq \frac{2\kappa_N}{\hat{\mu}} \sum_{i=1}^N D_\Phi(x^*, x^i)$$

526 Using the upper bounds for the three terms in  $V_3$  we obtain the bound in (4.4).  $\square$

527 We then have the following convergence result for the dynamics in (3.3).

528 **PROPOSITION 4.4** (Convergence of the dynamics in (3.3)). *Let Assumptions 1-3*  
 529 *hold and assume  $\kappa_g > 0$ . Consider the dynamics in (3.3) and  $V_t$  as defined in (4.3).*  
 530 *Let  $\hat{\mu} = \min\{\mu_\Phi, 2\}$ . Then it holds,*

$$531 \quad (4.5) \quad \mathbb{E}[V_T] \leq e^{-rT} \mathbb{E}[V_0] + \int_0^T e^{-r(T-s)} M(\mathbf{x}_s) ds$$

532 where  $V_T := V(\mathbf{x}_T, \boldsymbol{\lambda}_T)$ ,

$$533 \quad r = \frac{2\kappa_g}{c\hat{\mu} + 2\alpha(\Phi) + 3\kappa_N},$$

$$536 \quad M(\mathbf{x}) = c \frac{\sigma^2}{2} (\text{tr}(C_1(\mathbf{x})) + \langle \mathbf{x} - \mathbf{x}^*, \Delta \cdot \nabla \Phi^*(\mathbf{z}) \rangle) + \frac{\sigma^2}{2} \text{tr}(C_2(\mathbf{x}))$$

$$537 \quad C_1(\mathbf{x}) = \sum_{i=1}^N \nabla^2 \Phi(x^i)^{-1} \nabla_{xx}^2 D_{\Phi}(x^*, x^i) \nabla^2 \Phi(x^i)^{-1}$$

$$538 \quad C_2(\mathbf{x}) = \nabla^2 \Phi(\mathbf{x})^{-1} (\nabla^2 f(x) + \mathcal{L}) \nabla^2 \Phi(\mathbf{x})^{-1},$$

540 and  $c \geq 2\kappa_{\beta, N}$ .

541 *Proof.* Let  $c$  be as in Lemma 4.3 to ensure non-negativity of the Lyapunov function.  
542 Since  $x^i = \nabla \Phi^*(z^i)$  it follows from Itô's Lemma that,

$$543 \quad dx_t^i = \nabla^2 \Phi^*(z^i) dz_t^i + \frac{1}{2} \sigma^2 \Delta \cdot \nabla \Phi^*(z^i) dt,$$

545 where the  $j^{\text{th}}$  element of the Itô correction term is  $[\Delta \cdot \nabla \Phi^*(z^i)]_j = \sum_{k=1}^d \partial_{kk}^2 \partial_j \Phi^*(z^i)$ .

546 For ease of exposition we define the following terms,

$$547 \quad dM_t^1 = \frac{\sigma^2}{2} (\text{tr}(C_1(\mathbf{x}_t)) + \langle \mathbf{x}_t - \mathbf{x}^*, \Delta \cdot \nabla \Phi^*(\mathbf{z}_t) \rangle) dt + \langle \mathbf{x}_t - \mathbf{x}^*, \sigma d\mathbf{B}_t \rangle,$$

$$550 \quad dM_t^2 = \frac{\sigma^2}{2} \text{tr}(C_2(\mathbf{x}_t)) dt + \sigma \langle \nabla_x L(\mathbf{x}_t, \boldsymbol{\lambda}_t), \nabla^2 \Phi^*(\mathbf{z}_t) d\mathbf{B}_t \rangle.$$

552 Using the equilibrium points and the fact that  $\nabla^2 \Phi^*(z) \nabla^2 \Phi(x) = I$ , we obtain,

$$553 \quad (4.6) \quad \begin{aligned} d(V_t^1 + V_t^2) &\leq -\langle \mathbf{x}_t - \mathbf{x}^*, \mathcal{L}(\mathbf{x}_t - \mathbf{x}^*) \rangle dt + dM_t^1 \\ &\leq -\frac{1}{\kappa_{\beta, N}} \|\mathcal{L}(\mathbf{x}_t - \mathbf{x}^*)\|_2^2 + dM_t^1 \end{aligned}$$

554 where for the first inequality we used the convexity of  $f$  and the symmetry of  $\mathcal{L}$  and  
555 for the second inequality we used Lemma 2.6. We also have using  $L$  from (3.5),

$$556 \quad \begin{aligned} dV_t^3 &= \langle \nabla_{\mathbf{x}} L(\mathbf{x}_t, \boldsymbol{\lambda}_t), d\mathbf{x}_t \rangle + \frac{\sigma^2}{2} \text{tr}(C_2(\mathbf{x}_t, \boldsymbol{\lambda}_t)) dt + \langle \mathcal{L}(\mathbf{x}_t - \mathbf{x}^*), d\boldsymbol{\lambda}_t \rangle \\ &= \left( -\|\nabla_x L(\mathbf{x}_t, \boldsymbol{\lambda}_t)\|_{\nabla^2 \Phi^*(\mathbf{z}_t)}^2 + \|\mathcal{L}\mathbf{x}_t\|_2^2 \right) dt + dM_t^2 \\ &\leq -\frac{2\kappa_g}{\hat{\mu}} \left( \sum_{i=1}^N D_{\Phi^*}(z_t^i, z^*) + \|\boldsymbol{\lambda}_t - \boldsymbol{\lambda}^*\|_2^2 \right) dt + \|\mathcal{L}\mathbf{x}_t\|_2^2 dt + dM_t^2, \end{aligned}$$

557 where in the first line we used the optimality conditions  $\nabla f(\mathbf{x}^*) + \mathcal{L}\boldsymbol{\lambda}^*, \mathcal{L}\mathbf{x}^* = 0$ , and  
558 to obtain the last inequality we used Lemma 2.8 with  $\Psi = \frac{1}{2} \|\cdot\|_2^2$ . If in addition  
559  $c \geq 2\kappa_{\beta, N}$  and using the bound in (4.6) we obtain,

$$560 \quad \begin{aligned} dV_t &\leq -\frac{2\kappa_g}{\hat{\mu}} \left( \sum_{i=1}^N D_{\Phi^*}(z_t^i, z^*) + \|\boldsymbol{\lambda}_t - \boldsymbol{\lambda}^*\|_2^2 \right) dt - \|\mathcal{L}\mathbf{x}_t\|_2^2 dt + cdM_t^1 + dM_t^2 \\ &\leq -\frac{2\kappa_g}{c\hat{\mu} + 2\alpha(\Phi) + 3\kappa_N} V(\mathbf{x}_t, \boldsymbol{\lambda}_t) dt + cdM_t^1 + dM_t^2, \end{aligned}$$



561 where in the last inequality we used (4.4) from Lemma 4.3. Finally, taking expectations,  
 562 integrating and using Gronwall's lemma we obtain (4.5).  $\square$

563 It is clear that  $V$  is a stochastic Lyapunov function and exact convergence can be  
 564 achieved using  $\sigma = 0$ . We note that if  $\kappa_g \geq 0$  then the result above implies that  
 565  $dV_t \leq 0$ , but we may not have an exponential convergence rate. In Lemma 2.7 we  
 566 showed that the strong convexity of the objective function implies that  $\kappa_g > 0$  and we  
 567 note that the reverse is not true.

568 **4.3. Convergence with preconditioned interaction.** The motivation behind  
 569 the EPISMD algorithm in (3.4) is that in both unconstrained and constrained settings,  
 570 additional speedup can be obtained by *preconditioning* the dual variable  $\lambda$ . The use of  
 571 the mirror map  $\Psi$  results in additional flexibility in the convergence rate; furthermore  
 572 it allows to work with the Bregman divergence as the Lyapunov function. We observe  
 573 this additional flexibility through the term  $\hat{\mu}$ , which is given by  $\hat{\mu} = \min(\mu_\Phi, \mu_\Psi)$  so  
 574 that the proper choice of mirror map  $\Psi$  can additionally improve the convergence.

575 Consider  $V_t$  as in (4.3) but with,

$$576 \quad (4.7) \quad V_t^2(\lambda_t) = \sum_{i=1}^N D_\Psi(\lambda^*, \lambda_t^i).$$

578 As before we will change the definition of the optimal point the algorithm will converge  
 579 to as follows,

$$580 \quad (x^*, \lambda^*) = \arg \min_{x, \lambda \in (X^*, \Lambda^*)} D_\Phi(x, x_0) + D_\Psi(\lambda, \lambda_0).$$

581 The convergence of (3.4) can be obtained using slight modifications of the proof of  
 582 Lemma 4.3.

583 LEMMA 4.5. *Let Assumptions 1-3 hold. Then  $V_t$  with  $V_t^2$  as in (4.7) satisfies*  
 584 *Lemma 4.3 (i) and,*

585 **ii.a** *If  $c \geq \max\{\kappa_N/\mu_\Phi, \kappa_N/\mu_\Psi\}$ ,*

$$586 \quad V(\mathbf{x}, \lambda) \geq \frac{1}{2}(\mu_\Phi c - \kappa_N)\|\mathbf{x} - \mathbf{x}^*\|^2 + \frac{1}{2}(\mu_\Psi c - \kappa_N)\|\lambda - \lambda^*\|_2^2 \geq 0.$$

588 **iii.a** *Let  $\hat{\mu} = \min\{\mu_\Phi, \mu_\Psi\}$ . Then,*

$$589 \quad (4.8) \quad V(\mathbf{x}, \lambda) \leq \left(c + \frac{3\kappa_N + 2\alpha(\Phi)}{\hat{\mu}}\right) \left(\sum_{i=1}^N D_\Phi(x^*, x^i) + \sum_{i=1}^N D_\Psi(\lambda^*, \lambda^i)\right).$$

590 *If in addition  $f$  is  $\mu_f$ -strongly convex relative to  $\Phi$  then:*

591 **ii.b** *For any  $c \geq \max\{(\kappa_N - \mu_f L_\Phi)/\mu_\Phi, \kappa_N/\mu_\Psi\}$ ,*

$$592 \quad V(\mathbf{x}, \lambda) \geq \frac{1}{2}(\mu_\Phi c + \mu_f L_\Phi - \kappa_N)\|\mathbf{x} - \mathbf{x}^*\|^2 + \frac{1}{2}(\mu_\Psi c - \kappa_N)\|\lambda - \lambda^*\|_2^2 \geq 0.$$

593 The following convergence then holds.

594 PROPOSITION 4.6 (Convergence of the preconditioned dynamics in (3.4)). *Let*  
 595 *Assumptions 1-3 hold and  $\kappa_g > 0$ . Consider the dynamics in (3.4). Let the Lyapunov*  
 596 *function  $V_t$  be defined as in (4.3) with  $V_t^2$  as in (4.7). Then the result from Proposition*  
 597 *4.4 holds with  $\hat{\mu} = \min\{\mu_\Phi, \mu_\Psi\}$  and  $c \geq 2\kappa_{\beta, N}/\mu_\Psi$ .*

598 *Proof.* Let  $c$  be as in Lemma 4.5 to ensure non-negativity of the Lyapunov function.  
 599 We follow similar steps as in the proof of Proposition 4.4. Observe that, assuming that  
 600 the Bregman divergence is differentiable in the second variable and using (3.8),

$$601 \quad dV_t^2 = \langle -\nabla^2 \Psi(\boldsymbol{\lambda}_t)(\boldsymbol{\lambda}^* - \boldsymbol{\lambda}_t), \nabla^2 \Psi(\boldsymbol{\lambda}_t)^{-1} \mathcal{L} \mathbf{x}_t \rangle$$

$$602 \quad = (\boldsymbol{\lambda}_t - \boldsymbol{\lambda}^*)^T \mathcal{L}(\mathbf{x}_t - \mathbf{x}^*) dt.$$

604 Then, a modification of Lemma 2.6 can be derived; using the positive semi-definiteness  
 605 of  $\nabla^2 \Psi$  we can derive  $\frac{\mu_\Psi}{\kappa_{\beta, N}} \nabla^2 \Psi(\boldsymbol{\lambda})^{-1} \preceq \mathcal{L}^{-1}$  and obtain,

$$606 \quad \langle \mathbf{x}_t, \mathcal{L} \mathbf{x}_t \rangle \geq \frac{\mu_\Psi}{\kappa_{\beta, N}} \langle \mathcal{L} \mathbf{x}_t, \nabla^2 \Psi(\boldsymbol{\lambda}_t)^{-1} \mathcal{L} \mathbf{x}_t \rangle.$$

608 Then,

$$609 \quad (4.9) \quad d(cV_t^1 + cV_t^2) \leq -\frac{c\mu_\Psi}{\kappa_{\beta, N}} \|\mathcal{L} \mathbf{x}_t\|_{\nabla^2 \Psi(\boldsymbol{\lambda}_t)^{-1}}^2 dt + dM_t^1.$$

611 Furthermore using (3.8),

$$612 \quad dV_t^3 = \left( -\|\nabla_{\mathbf{x}} L(\mathbf{x}_t, \boldsymbol{\lambda}_t)\|_{\nabla^2 \Phi^*(\mathbf{z}_t)}^2 + \|\mathcal{L} \mathbf{x}_t\|_{\nabla^2 \Psi(\boldsymbol{\lambda}_t)^{-1}}^2 \right) dt + dM_t^2$$

614 Consequently apply Lemma 2.8, the bound in (4.9) and use the additional assumption  
 615 that  $c \geq 2\kappa_{\beta, N}/\mu_\Psi$  to obtain

$$616 \quad dV_t \leq -\frac{2\kappa_g}{\hat{\mu}} \left( \sum_{i=1}^N D_{\Phi^*}(z_t^i, z^*) + \sum_{i=1}^N D_{\Psi}(\lambda_t^i, \lambda^*) \right) dt - \|\mathcal{L} \mathbf{x}_t\|_{\nabla^2 \Psi(\boldsymbol{\lambda})^{-1}}^2 dt$$

$$617 \quad + cdM_t^1 + dM_t^2.$$

619 Lastly use (4.8) and we complete the proof by integrating and taking expected values.  $\square$

620 The convergence of  $\boldsymbol{\lambda}_t$  is affected by how well-conditioned the interaction graph and  
 621 objective function are. In the proof of Proposition 4.6 we observe that the use of the  
 622 mirror map  $\Psi$  results in an improved rate of convergence. In particular, when  $\frac{\mu_\Phi}{\mu_\Psi} < 1$   
 623 we end up with  $r$  in Proposition 4.6 being larger than that seen in Proposition 4.4.

624 *Choosing the preconditioner.* For many interesting applications good mirror maps  
 625 are known, and the advantages of mirror descent over gradient descent are well  
 626 understood. Unfortunately, it is not clear how to select a good mirror map for the  
 627 space of Lagrange multipliers. However, if we restrict the mirror map to be quadratic,  
 628 we postulate that a positive definite approximation of the Hessian of the dual function  
 629 will work well in practice, e.g. as seen in the results of Section 5. We argue that the  
 630 extra computation associated with approximating the Hessian of the dual function  
 631 could be justified in the scenario where the Laplacian matrix is ill-conditioned. Below  
 632 we briefly outline the derivation of the Hessian for the dual function in the deterministic  
 633 setting and when  $f$  is strongly convex. In particular, the (negative) dual function  $q(\boldsymbol{\lambda})$   
 634 is defined as follows,

$$635 \quad (4.10) \quad q(\boldsymbol{\lambda}) = \max_{\mathbf{x} \in \mathcal{X}^N} \{-(f(\mathbf{x}) + \langle \boldsymbol{\lambda}, \mathcal{L} \mathbf{x} \rangle)\},$$

636 and let  $\mathbf{x}(\boldsymbol{\lambda})$  be a maximizer of (4.10). We know that for the strongly convex case,  
 637  $\mathbf{x}(\boldsymbol{\lambda})$  is unique and the gradient of the dual function is given by ([3, Theorem 6.3.3]),

$$638 \quad \nabla_{\boldsymbol{\lambda}} q(\boldsymbol{\lambda}) = -\mathcal{L} \mathbf{x}(\boldsymbol{\lambda}).$$

639 Applying the KKT conditions to (4.10) and differentiating with respect to  $\lambda$  we also  
 640 have that,

$$641 \quad \nabla^2 f(\mathbf{x}(\lambda)) \frac{d\mathbf{x}(\lambda)}{d\lambda} + \mathcal{L} = 0.$$

642 Using the preceding equation we obtain the following expression for the Hessian of the  
 643 dual function,

$$644 \quad \nabla^2 q(\lambda) = -\mathcal{L} \frac{d\mathbf{x}(\lambda)}{d\lambda} = \mathcal{L} \nabla^2 f(\mathbf{x}(\lambda))^{-1} \mathcal{L}.$$

645 Unfortunately, even if  $f$  is strongly convex, the dual function in (4.10) is only convex  
 646 (and not strongly convex). So we cannot directly use the Hessian of the dual function  
 647 above, instead we proposed to use,

$$648 \quad \nabla^2 \Psi(\lambda) = \mathcal{L}_\beta \nabla^2 f(\mathbf{x}(\lambda))^{-1} \mathcal{L}_\beta,$$

649 where  $\mathcal{L}_\beta$  is the  $\beta$ -regularized Laplacian defined in (2.4). Note that we can avoid  
 650 inverting the Hessian of  $f$  at every iteration and instead work with the convex conjugate  
 651 of  $\Psi$  which occurs a one-off cost of diagonalizing the regularized Laplacian,

$$652 \quad \nabla^2 \Psi^*(\lambda) = \mathcal{L}_\beta^{-1} \nabla^2 f(\mathbf{x}(\lambda)) \mathcal{L}_\beta^{-1}.$$

653 The argument above could be made more precise and extended to more general settings.  
 654 We report promising numerical experiments with this choice for the mirror map in  
 655 Section 5.5.

656 **5. Numerical results.** Here we study the performance of the proposed algo-  
 657 rithms in three problems:

658 **A An unconstrained ill conditioned linear system.** We set the local cost  
 659 functions as  $f_i(x) = \frac{1}{2} \|Q_i x - b_i\|_2^2$ . Unless mentioned otherwise  $\mathcal{X} = \mathbb{R}^d$  with  
 660  $d = 200$  and  $Q_i \in \mathbb{R}^{20 \times 200}$  is a random matrix with condition number 15 and  
 661  $b_i \sim \mathcal{N}(0, I_{20})$ . For the mirror map we set  $\Phi(x) = \frac{1}{2} \|x\|^2$ . We remark that in  
 662 this setting the only *constraint* in the problem comes from the consensus.

663 **B A constrained linear system.** The setup is similar to the unconstrained  
 664 linear system but we set  $\mathcal{X} = \Delta_d$ , the  $d$ -dimensional simplex. The mirror  
 665 map is set to be the negative entropy function  $\Phi(x) = \sum_{j=1}^d [x]_j \log([x]_j)$  such  
 666 that  $[z]_i = 1 + \log([x]_j)$  (and  $[z]_i = 0$  if  $[x]_j = 0$ ) and  $[x]_i = e^{[z]_i - 1}$  and the  
 667 mapping onto the simplex is done using the normalization of this negative  
 668 entropy mirror map. Here we denote with  $[x]_i$  the  $i$ -th element of vector  $x$ .

669 **C A neural network.** Motivated by federated learning applications [25] we  
 670 consider the training of a neural network with one hidden layer and 30 nodes  
 671 per layer using a ReLU activation, a softmax output and the cross-entropy  
 672 loss. The training data is the FashionMNIST data. Each particle  $i$  has access  
 673 to 10 of these samples and the cross-entropy loss over this subset defines each  
 674  $f_i$ . We will assume the solution lies in the constraint set  $\mathcal{X} = \Delta_d$  (e.g. looking  
 675 for a sparse solution) and use the negative entropy mirror map.

676 For the connectivity graph we define three options: **(i)** a cyclic graph with each node  
 677 connected to the previous and next node, **(ii)** an Erdos-Rényi graph, or **(iii)** a barbell  
 678 graph. For the dynamics we use a standard Euler discretization with  $\Delta t = 0.01$  and  
 679 50,000 epochs. We implement ISMD, EISMD and EPISMD dynamics that include

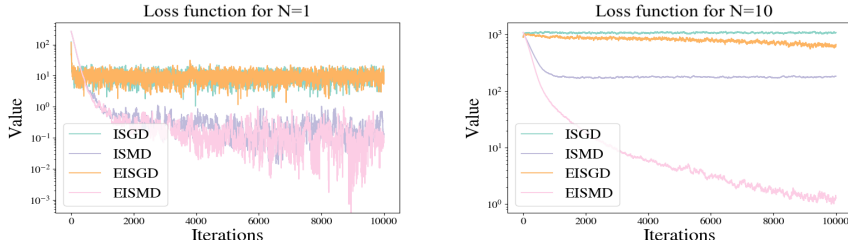


Fig. 3: Loss function for problem (B) with cyclic graph (i). Left is centralized ( $N = 1$ ), right distributed ( $N = 10$ ). ISGD, EISGD and ISMD, EISMD use (3.2), (3.3) with  $\Phi(x) = \frac{1}{2}\|x\|^2$  (i.e.  $\mathbf{x}_t = \mathbf{z}_t$ ) and the entropic mirror map, resp. projected to  $\Delta_d$ .

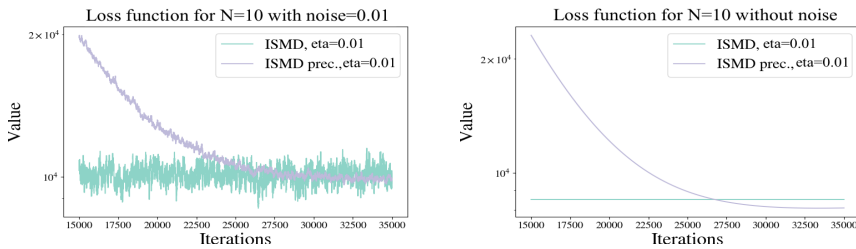


Fig. 4: Loss functions for Problem (A) and graph (i) with condition number set to 100,  $N = 10$  and preconditioned with a local objective function for different learning rates (eta in legend). Left  $\sigma = 0.01$ , right  $\sigma = 0$ .

680 hyperparameters  $\epsilon$ ,  $\eta$  for the interaction strength and learning rate; see Appendix B.1  
 681 for more details.

682 **5.1. The effect of the mirror map.** We first compare distributed mirror  
 683 descent with projected gradient descent for problem (B) using graph (i) to showcase  
 684 the benefits of the entropic mirror map in simplex constrained systems. In Figure 3  
 685 we present the results for ISMD and EISMD for  $N = 1$  (centralized implementation)  
 686 and  $N = 10$  (distributed). The benefits of the exact dynamics in (3.3) and the mirror  
 687 map are clear in both cases.

688 **5.2. The limits of the mirror map for ISMD.** As shown in Section 3.1, the  
 689 mirror map is equivalent to a preconditioner in the primal dynamics. In problem  
 690 (A) with  $\mathcal{X} = \mathbb{R}^d$ , if the mirror map is chosen to be the local objective function then  
 691 ISMD results in a local Newton-like algorithm. Here we will explore the effectiveness  
 692 of using such a local preconditioner and set  $\nabla^2\Phi(x^i) = Q_i^T Q_i$ . Figure 4 shows the  
 693 convergence results for the problem A with condition number increased to 100,  $d = 20$   
 694 and the cyclic interaction graph (i). The distributed case when preconditioned with  
 695 the local Hessian does not result in full convergence. Thus the mirror map alone cannot  
 696 facilitate convergence. The results are meant as motivation for using EISMD. We note  
 697 that other solutions are possible, e.g. extending [34] to provide an approximation of  
 698 the full preconditioner, but this will increase communication at each step.

699 **5.3. The effects of the exact algorithm.** Here we present a detailed analysis  
 700 of the benefits of EISMD compared to ISMD using both the  $\ell_2$  error between the

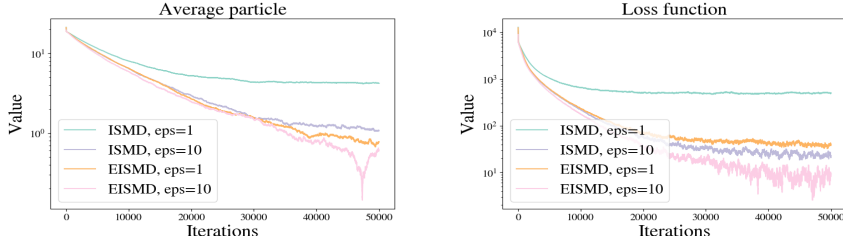


Fig. 5: ISMD and EISMD for Problem (A) and graph (i) for  $\eta = 0.01$ ,  $N = 10$ ,  $\sigma = 0.1$  and  $\epsilon = 1, 10$  (eps in legend). Left panel shows  $\|x_t^b - x_t^w\|^2$  with  $b, w$  the best and worst particle index at each time and right panel shows cost functions.

701 best and worst performing particle and the loss function computed over all particles.  
 702 Figure 5 shows results for problem (A) and graph (i). In all cases a higher interaction  
 703 strength allows to converge closer to the optimum and EISMD performs significantly  
 704 better than ISMD. Figure 6 shows similar results for problem (B). In the top panels  
 705 EISMD is clearly more effective than ISMD. The lower plots of Figure 6 show that  
 706 in the presence of noise a high interaction in ISMD is able to mitigate the process  
 707 variance due to the (appearing as oscillations) and result in convergence closer to the  
 708 optimum.

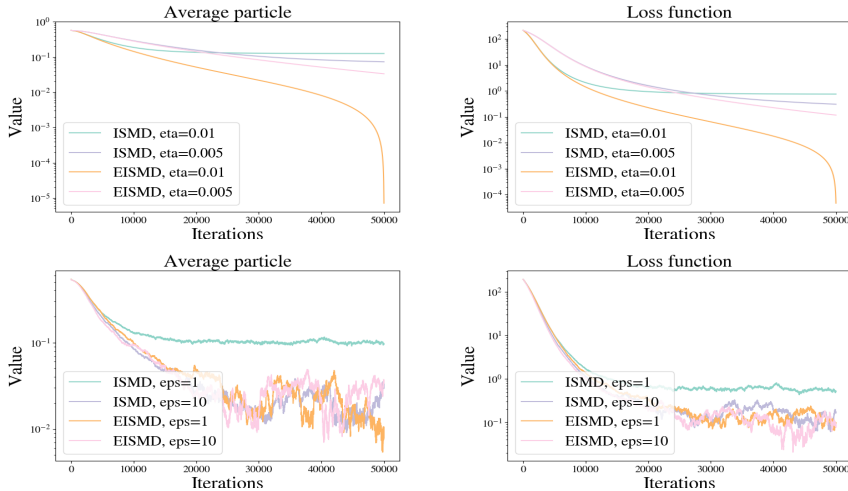


Fig. 6: ISMD and EISMD for Problem (B) and graph (i). Bottom row details as in Figure 5. Top row panels same but with  $\sigma = 0$ ,  $\epsilon = 1$  and  $\eta = 0.01, 0.005$ .

709 **5.4. The choice of interaction graph.** Here we set the communication struc-  
 710 ture to be an Erdos-Rényi communication graph (ii). This is a random graph where  
 711 each edge is chosen with a certain probability  $p$ . The communication between the  
 712 nodes is thus determined by this connectivity probability with on average each node  
 713 being connected to  $p \times N$  other nodes. The weight matrix  $A$  is set to be a doubly  
 714 stochastic version of this communication graph and additionally we use [7] to optimize

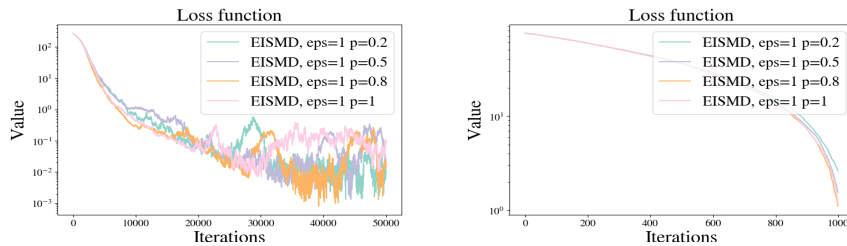


Fig. 7: Loss functions of EISMD for problem (B) and graph (ii). Left panel contains full results and right magnifies the initial epochs. We set  $\epsilon = 1$ ,  $\sigma = 0.1$ ,  $\eta = 0.01$  and vary the number of nodes each particle communicates with to 2, 5, 8 or 10.

715 the coefficients to get the fastest mixing in (3.3) whilst maintaining minimum commu-  
 716 nications. Figure 7 shows results for problem (B) and different graph connectivities.  
 717 The benefits of EISMD become even more clear from these results; using a graph  
 718 with an optimal interaction structure while maintaining minimal interaction costs  
 719 in combination with EISMD can result in very fast and computationally efficient  
 720 convergence. Even with a very low communication cost per round, specifically each  
 721 node communicating with approximately two other nodes per round, it results in a  
 722 convergence speed comparable to full communication seen earlier.

723 **5.5. Accelerated dual convergence.** The right choice of mirror map (pre-  
 724 conditioner) in the unconstrained setting can also result in accelerated convergence  
 725 of the  $\lambda$  variable. We know from the discussion in Section 4.3 the form of a good  
 726 preconditioner. For the linear case, it is given by  $\nabla^2 \Psi(\lambda) = -\mathcal{L}_\beta (QQ^T)^{-1} \mathcal{L}_\beta$  with  
 727  $Q$  being block diagonal composition of  $Q_i$ -s. We use the regularized Laplacian here,  
 728 specifically  $\mathcal{L}_\beta = \mathcal{L} + 0.01 \cdot \mathbf{1}_N \mathbf{1}_N^T \otimes I_d$ . We show the benefits of this preconditioner  
 729 in EPISMD (shown in (3.4)) numerically in Figure 8 for a barbell graph with two  
 730 clusters (iii). Clearly, preconditioning allows to converge much *faster and closer*  
 731 to the optimum. We observe that preconditioning enables the algorithm to converge  
 732 much faster and closer to the optimum. In the case of the barbell graph (iii) where  
 733 the graph Laplacian is not well-conditioned, the ability to speed up the convergence  
 734 using the preconditioning of the dual variable can lead to very substantial benefits in  
 735 many real-world clustered systems.

736 **5.6. Federated learning in a nonconvex example.** In this example we study  
 737 the performance of ISMD and EISMD for problem (C) and graph (i). In Figure 9 we  
 738 compute the gradient over a random batch of data points. Each particle computes  
 739 the gradient in each iteration over a batch of size 10, which results in implicit noise in  
 740 the algorithm. We observe that the exact method results converge significantly faster  
 741 even in this nonconvex setting.

742 **6. Discussion.** Our work successfully presented the benefits of mirror maps and  
 743 how to use them while still achieving consensus in a distributed setting. For future  
 744 work several interesting extensions can be considered. If one discretizes the continuous  
 745 time dynamics, e.g. as in constant step size Euler-Maruyama schemes, then a bias will  
 746 be present that depends on the discretization interval. In this sense exactness is lost,  
 747 but could be recovered again using decreasing step sizes or constant ones combined  
 748 with ideas from sampling methods [14]. Addressing this in detail and comparing with

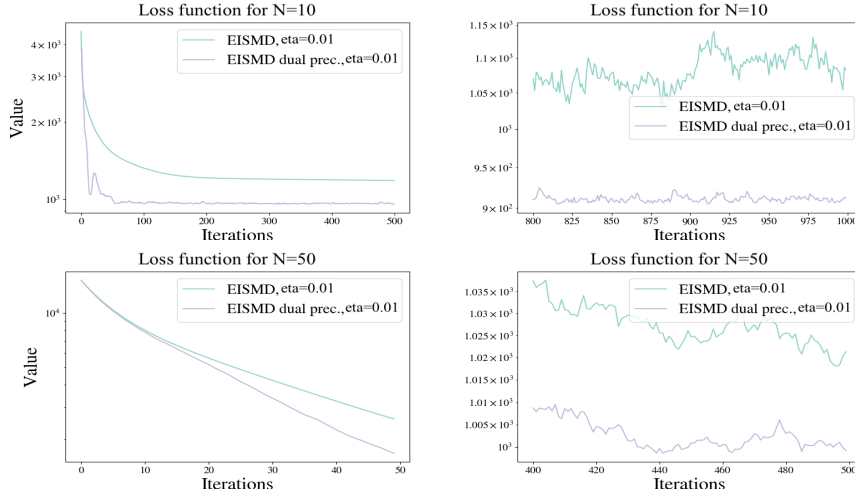


Fig. 8: Loss functions comparing EPISMD and EISMD for Problem A and graph (iii). We use  $\sigma = 0.01$  and top row panels use a barbell graph with two clusters and 10 nodes and bottom panels two clusters and 50 nodes.

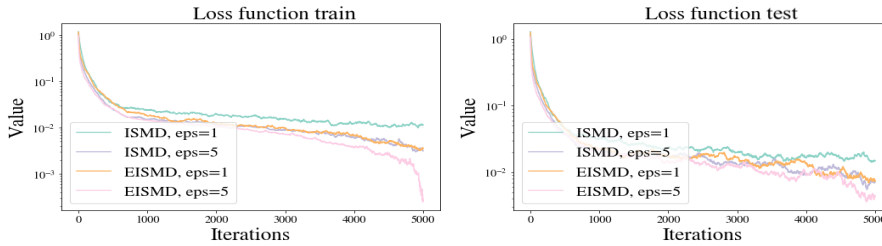


Fig. 9: ISMD and EISMD for problem (C) and graph (i). Left shows training loss against epoch and right shows test loss.

749 the discrete time methods in Table 1 is left for further work. The analysis shown  
750 here for EISMD and EPISMD can be extended to a setting where the solution of  
751 (1.1) does not lie in the constraint space  $\mathcal{X}$  using steps similar to [26]. Furthermore,  
752 we only numerically explored the nonconvex setting in EISMD using a small neural  
753 network with a negative entropy mirror map; future work can include a theoretical  
754 analysis of the nonconvex case as well as exploring further the benefits of different  
755 mirror maps using EPISMD in the distributed training of neural networks. Last, the  
756 algorithm could be modified to one applying the Laplacian onto the mirrored variables;  
757 preliminary numerical results showed that in high-dimensional and/or nonconvex  
758 settings this could be of benefit; a more detailed converge analysis of this phenomenon  
759 may be of interest which would require the definition of a different Lyapunov function.

760 *Acknowledgements.* This project was funded by JPMorgan Chase & Co under J.P.  
761 Morgan A.I. Research Awards in 2019 and 2021. G.A.P. was partially supported by  
762 the EPSRC through grant number EP/P031587/1.

- 764 [1] Y. ARJEVANI AND O. SHAMIR, *Communication complexity of distributed convex learning and*  
765 *optimization*, arXiv preprint arXiv:1506.01900, (2015).
- 766 [2] M. BADIOI KHUZANI AND N. LI, *Stochastic Primal-Dual Method on Riemannian Manifolds*  
767 *with Bounded Sectional Curvature*, arXiv e-prints, (2017), pp. arXiv-1703.
- 768 [3] M. S. BAZARAA, H. D. SHERALI, AND C. M. SHETTY, *Nonlinear programming: theory and*  
769 *algorithms*, John Wiley & Sons, 2013.
- 770 [4] A. BECK, *First-order methods in optimization*, vol. 25, SIAM, 2017.
- 771 [5] D. P. BERTSEKAS AND J. N. TSITSIKLIS, *Parallel and distributed computation: numerical*  
772 *methods*, 2015.
- 773 [6] A. BOROVYKH, P. PAPPAS, N. KANTAS, AND G. PAVLIOTIS, *On stochastic mirror descent with*  
774 *interacting particles: convergence properties and variance reduction*, Physica D.Nonlinear  
775 Phenomena 418, (2021).
- 776 [7] S. BOYD, P. DIACONIS, AND L. XIAO, *Fastest mixing Markov Chain on a graph*, SIAM review,  
777 46 (2004), pp. 667–689.
- 778 [8] S. BOYD, N. PARIKH, E. CHU, B. PELEATO, AND J. ECKSTEIN, *Distributed optimization and*  
779 *statistical learning via the alternating direction method of multipliers*, Found. Trends Mach.  
780 Learn., 3 (2011), p. 1–122.
- 781 [9] S. BUBECK, *Convex optimization: Algorithms and complexity*, arXiv preprint arXiv:1405.4980,  
782 (2014).
- 783 [10] F. BULLO, *Lectures on Network Systems*, Kindle Direct Publishing, 1.5 ed., 2021.
- 784 [11] S. CHEN, A. GARCIA, M. HONG, AND S. SHAHRAMPOUR, *Decentralized Riemannian Gradient*  
785 *Descent on the Stiefel Manifold*, arXiv preprint arXiv:2102.07091, (2021).
- 786 [12] R. D’ORAZIO, N. LOIZOU, I. LARADJI, AND I. MITLIAGKAS, *Stochastic mirror descent:*  
787 *Convergence analysis and adaptive variants via the mirror stochastic polyak stepsize*, arXiv  
788 preprint arXiv:2110.15412, (2021).
- 789 [13] J. C. DUCHI, A. AGARWAL, AND M. J. WAINWRIGHT, *Dual averaging for distributed optimiza-*  
790 *tion: Convergence analysis and network scaling*, IEEE Transactions on Automatic control,  
791 57 (2011), pp. 592–606.
- 792 [14] A. DURMUS AND E. MOULINES, *Sampling from strongly log-concave distributions with the*  
793 *unadjusted Langevin algorithm*, arXiv preprint arXiv:1605.01559, 5 (2016).
- 794 [15] G. FRANCA, D. ROBINSON, AND R. VIDAL, *ADMM and accelerated ADMM as continuous*  
795 *dynamical systems*, in International Conference on Machine Learning, PMLR, 2018, pp. 1559–  
796 1567.
- 797 [16] B. GHARESIFARD AND J. CORTÉS, *Distributed continuous-time convex optimization on weight-*  
798 *balanced digraphs*, IEEE Transactions on Automatic Control, 59 (2013), pp. 781–786.
- 799 [17] H. HENDRIKX, L. XIAO, S. BUBECK, F. BACH, AND L. MASSOULIE, *Statistically preconditioned*  
800 *accelerated gradient method for distributed optimization*, in International Conference on  
801 Machine Learning, PMLR, 2020, pp. 4203–4227.
- 802 [18] R. A. HORN AND C. R. JOHNSON, *Topics in matrix analysis*, Cambridge university press, 1994.
- 803 [19] D. JAKOVETIĆ, J. XAVIER, AND J. M. MOURA, *Fast distributed gradient methods*, IEEE  
804 Transactions on Automatic Control, 59 (2014), pp. 1131–1146.
- 805 [20] C. LI, C. CHEN, D. CARLSON, AND L. CARIN, *Preconditioned stochastic gradient Langevin*  
806 *dynamics for deep neural networks*, in Proceedings of the AAAI Conference on Artificial  
807 Intelligence, vol. 30, 2016.
- 808 [21] S. LIANG, L. WANG, AND G. YIN, *Exponential convergence of distributed primal–dual convex*  
809 *optimization algorithm without strong convexity*, Automatica, 105 (2019), pp. 298–306.
- 810 [22] P. LIN, W. REN, AND J. A. FARRELL, *Distributed continuous-time optimization: nonuniform*  
811 *gradient gains, finite-time convergence, and convex constraint set*, IEEE Transactions on  
812 Automatic Control, 62 (2016), pp. 2239–2253.
- 813 [23] S. LIU, Z. QIU, AND L. XIE, *Continuous-time distributed convex optimization with set*  
814 *constraints*, IFAC Proceedings Volumes, 47 (2014), pp. 9762–9767.
- 815 [24] H. LU, R. M. FREUND, AND Y. NESTEROV, *Relatively smooth convex optimization by first-order*  
816 *methods, and applications*, SIAM Journal on Optimization, 28 (2018), pp. 333–354.
- 817 [25] B. McMAHAN, E. MOORE, D. RAMAGE, S. HAMPSON, AND B. A. Y ARCAS, *Communication-*  
818 *efficient learning of deep networks from decentralized data*, in Artificial Intelligence and  
819 Statistics, 2017, pp. 1273–1282.
- 820 [26] P. MERTIKOPOULOS AND M. STAUDIGL, *On the convergence of gradient-like flows with noisy*  
821 *gradient input*, SIAM Journal on Optimization, 28 (2018), pp. 163–197.
- 822 [27] A. NEDIĆ, S. LEE, AND M. RAGINSKY, *Decentralized online optimization with global objectives*  
823 *and local communication*, in 2015 American Control Conference (ACC), IEEE, 2015,



- 824 pp. 4497–4503.  
825 [28] A. NEDIC, A. OLSHEVSKY, AND W. SHI, *Achieving geometric convergence for distributed*  
826 *optimization over time-varying graphs*, SIAM Journal on Optimization, 27 (2017), pp. 2597–  
827 2633.  
828 [29] S. PU AND A. NEDIĆ, *Distributed stochastic gradient tracking methods*, Mathematical Program-  
829 ming, 187 (2021), pp. 409–457.  
830 [30] G. QU AND N. LI, *Harnessing smoothness to accelerate distributed optimization*, IEEE Trans-  
831 actions on Control of Network Systems, 5 (2017), pp. 1245–1260.  
832 [31] M. RAGINSKY AND J. BOUVRIE, *Continuous-time stochastic mirror descent on a network:*  
833 *Variance reduction, consensus, convergence*, in 2012 IEEE 51st IEEE Conference on Decision  
834 and Control (CDC), IEEE, 2012, pp. 6793–6800.  
835 [32] S. S. RAM, A. NEDIĆ, AND V. V. VEERAVALLI, *Distributed stochastic subgradient projection*  
836 *algorithms for convex optimization*, Journal of optimization theory and applications, 147  
837 (2010), pp. 516–545.  
838 [33] S. SHAHRAMPOUR AND A. JADBABAIE, *Distributed online optimization in dynamic environments*  
839 *using mirror descent*, IEEE Transactions on Automatic Control, 63 (2017), pp. 714–725.  
840 [34] O. SHAMIR, N. SREBRO, AND T. ZHANG, *Communication-efficient distributed optimization*  
841 *using an approximate newton-type method*, in International conference on machine learning,  
842 PMLR, 2014, pp. 1000–1008.  
843 [35] G. SHI, A. PROUTIERE, AND K. H. JOHANSSON, *Network synchronization with convexity*,  
844 SIAM Journal on Control and Optimization, 53 (2015), pp. 3562–3583.  
845 [36] W. SHI, Q. LING, G. WU, AND W. YIN, *EXTRA: An exact first-order algorithm for decen-*  
846 *tralized consensus optimization*, SIAM Journal on Optimization, 25 (2015), pp. 944–966.  
847 [37] S. SRA, S. NOWOZIN, AND S. J. WRIGHT, *Optimization for machine learning*, Mit Press, 2012.  
848 [38] Y. SUN AND S. SHAHRAMPOUR, *Distributed mirror descent with integral feedback: Asymptotic*  
849 *convergence analysis of continuous-time dynamics*, arXiv preprint arXiv:2009.06747, (2020).  
850 [39] H. WANG AND A. BANERJEE, *Bregman alternating direction method of multipliers*, Advances  
851 in Neural Information Processing Systems, (2014).  
852 [40] J. WANG AND N. ELIA, *A control perspective for centralized and distributed convex optimization*,  
853 in 2011 50th IEEE conference on decision and control and European control conference,  
854 IEEE, 2011, pp. 3800–3805.  
855 [41] A. WIBISONO, A. C. WILSON, AND M. I. JORDAN, *A variational perspective on accelerated*  
856 *methods in optimization*, proceedings of the National Academy of Sciences, 113 (2016),  
857 pp. E7351–E7358.  
858 [42] K. YUAN, Q. LING, AND W. YIN, *On the convergence of decentralized gradient descent*, SIAM  
859 Journal on Optimization, 26 (2016), pp. 1835–1854.  
860 [43] X. ZENG, P. YI, AND Y. HONG, *Distributed continuous-time algorithm for constrained convex*  
861 *optimizations via nonsmooth analysis approach*, IEEE Transactions on Automatic Control,  
862 62 (2016), pp. 5227–5233.

## 863 Appendix A. Auxiliary results.

**A.1. The benefits of the mirror map.** To show the ability of mirror descent to adapt to a particular geometry we present the classical mirror descent proof which can be interpreted as having a central server; this server broadcasts  $x_t \in \mathcal{X}$  to all nodes which compute  $\nabla f_i(x_t)$  and send this back to the server. In continuous-time the algorithm at the server is given by

$$dz_t = - \sum_{i=1}^N \nabla f_i(x_t) dt + \sigma dB_t.$$

864 The aim is to converge to the optimal point  $(x^*, z^*)$  given as

$$865 \quad x^* = \arg \min_{x \in \mathcal{X}} f(x), \quad z^* = \nabla \Phi(x^*).$$

866 We will make use of the inequality,

$$867 \quad (\text{A.1}) \quad \int_0^t e^{-\alpha(t-s)} K ds \leq \frac{K}{\alpha},$$

868

869 and  $\|\Delta\Phi^*(z_t)\|_\infty \leq \infty$  (by Assumption 3).

870 LEMMA A.1 (Convergence of centralized mirror descent). *Let  $f(x) = \sum_{i=1}^N f_i(x)$*   
 871 *be  $\mu_f$ -strongly convex with respect to  $\Phi$ . Then,*

$$872 \quad \mathbb{E}[D_\Phi(x^*, x_t)] \leq e^{-\mu_f t} D_\Phi(x^*, x_0) + \frac{\sigma^2}{2\mu_f} \|\Delta\Phi^*\|_\infty.$$

874 *Proof.* Let  $V_t = D_{\Phi^*}(z_t, z^*)$ . We have through Itô's lemma,

$$875 \quad dV_t = - \sum_{i=1}^N (x^* - x_t)^T \nabla f_i(x_t) dt + \frac{1}{2} \sigma^2 \text{tr}(\Delta\Phi^*(z_t)) dt + \sigma(x_t - x^*)^T dB_t$$

$$876 \quad \leq (-\mu_f D_\Phi(x^*, x_t) - f(x_t) + f(x^*) + \frac{\sigma^2}{2} \|\Delta\Phi^*\|_\infty) dt + \sigma(x_t - x^*)^T dB_t$$

$$877 \quad \leq (-\mu_f D_{\Phi^*}(z_t, z^*) + \frac{\sigma^2}{2} \|\Delta\Phi^*\|_\infty) dt + \sigma(x_t - x^*)^T dB_t,$$

879 where in the first inequality we have used the  $\mu_f$ -strong convexity w.r.t.  $\Phi$  and in the  
 880 second inequality that by the properties of the mirror map  $D_\Phi(x^*, x_t) = D_{\Phi^*}(z_t, z^*)$   
 881 and that  $f(x^*) - f(x_t) \leq 0$ . Using (A.1), taking expectations and applying Grönwall  
 882 inequality the result follows.  $\square$

883 Comparing the results of mirror (preconditioned) gradient descent to the Euclidean  
 884 setting where  $\Phi_E(x) = \frac{1}{2}\|x\|_2^2$ , the effectiveness of the algorithm hinges on how large  
 885  $\mu_f$  is compared to  $\mu_E$  where  $\mu_E$  is such that  $f$  is  $\mu_E$ -strongly convex with respect  
 886 to  $\frac{1}{2}\|x\|_2^2$ . An optimal mirror map is then the one that is given by  $\Phi(x) = f(x)$ , in  
 887 which case it holds that  $\mu_f = 1$ . Discretizing the continuous-time algorithm would  
 888 then result in achieving convergence with a larger discretization step size. The works  
 889 of [17, 34, 1] study the discrete-time convergence if the mirror map is chosen to be  
 890 one of the objective functions, i.e.  $\Phi(x) = f_i(x)$  for some  $i$  and show that convergence  
 891 can be obtained in one time step with the right (estimate of the) preconditioner. One  
 892 can seek to choose the mirror map  $\Phi$  in such a way that  $D_\Phi(x_0, x^*)$  is smaller than  
 893  $D_{\Phi_E}(x_0, x^*)$  and get faster convergence, and in such a way that  $\|\Delta\Phi^*\|_\infty$  is smaller  
 894 than  $\|\Delta\Phi_E^*\|_\infty$  for closer convergence.

895 **A.2. Approximate convergence of distributed mirror descent.** Recall the  
 896 vectorized ISMD dynamics of (3.1) are

$$897 \quad (\text{A.2}) \quad d\mathbf{z}_t = (-\eta \nabla f(\mathbf{x}_t) - \epsilon \mathcal{L} \mathbf{z}_t) dt + \sigma d\mathbf{B}_t, \quad \mathbf{x}_t = \nabla \Phi^*(\mathbf{z}_t).$$

899 If  $x^*$  as per Lemma 4.1 does not exist, even in the deterministic case exact consensus  
 900 and optimality at convergence can no longer be achieved. We present a result which  
 901 shows that exponential convergence holds only up to a certain neighborhood of  $(\mathbf{x}^\dagger, \mathbf{z}^\dagger)$   
 902 both minimizing  $f(\mathbf{x}) + \frac{1}{2} \mathbf{z}^T \mathcal{L} \mathbf{z}$ , the size of this neighborhood depending on the noise  
 903 and the distance between the  $f(x^\dagger)$  and  $f(x^*)$ .

904 PROPOSITION A.2 (Approximate convergence of (A.2)). *Let Assumptions 1-3*  
 905 *hold and assume that  $f$  is  $\mu_f$ -strongly convex w.r.t.  $\Phi$ . Let  $V_t = \frac{1}{N} \sum_{i=1}^N D_{\Phi^*}(z_t^i, z^\dagger)$ ,*  
 906 *where  $z_t^i$  obeys the dynamics of (3.1) (or (A.2)). Then we have:*

$$907 \quad \mathbb{E}[V_t] \leq e^{-\eta \mu_f t} \frac{1}{N} \sum_{i=1}^N D_{\Phi^*}(z_0^i, z^\dagger) + \frac{\sigma^2}{2\eta \mu_f} \|\Delta\Phi^*\|_\infty + \frac{1}{\mu_f} (f(\mathbf{x}^\dagger) - f(\mathbf{x}^\circ)),$$

909 where  $x^\circ = \arg \min_{\mathbf{x} \in \mathcal{X}} f(\mathbf{x})$ .

910 *Proof.* Denote  $\mathbf{x}^\dagger = \arg \inf \{f(\mathbf{x}) + \frac{1}{2} \nabla \Phi(\mathbf{x})^T \mathcal{L} \nabla \Phi(\mathbf{x})\}$ ,  $\mathbf{z}^\dagger = \nabla \Phi(\mathbf{x}^\dagger)$ . Then

$$\begin{aligned}
911 \quad dV_t &= -\frac{1}{N} \sum_{i=1}^N (x_t^i - x^\dagger)^T \eta \nabla f_i(x_t^i) dt + \epsilon \frac{1}{N} \sum_{i=1}^N (x_t^i - x^\dagger)^T \sum_{j=1}^N A_{ij} (z_t^j - z_t^i) dt \\
912 \quad (A.3) \quad &+ \frac{1}{2} \sigma^2 \frac{1}{N} \sum_{i=1}^N \text{tr}(\Delta \Phi^*(z_t^i)) dt + \frac{1}{N} \sum_{i=1}^N (x_t^i - x^\dagger)^T dB_t^i. \\
913
\end{aligned}$$

914 Using the  $\mu_f$ -strong convexity of  $f$  as in the proof of Lemma A.1 we obtain,

$$\begin{aligned}
915 \quad \nabla f(\mathbf{x}_t)^T (\mathbf{x}^\dagger - \mathbf{x}_t) &\leq f(\mathbf{x}^\dagger) - f(\mathbf{x}_t) - \mu_f D_\phi(\mathbf{x}^\dagger, \mathbf{x}_t) \\
916 &\leq f(\mathbf{x}^\dagger) - f(\mathbf{x}^\circ) + f(\mathbf{x}^\circ) - f(\mathbf{x}_t) - \mu_f D_\phi(\mathbf{x}^\dagger, \mathbf{x}_t) \\
917 &\leq f(\mathbf{x}^\dagger) - f(\mathbf{x}^\circ) - \mu_f D_\phi(\mathbf{x}^\dagger, \mathbf{x}_t)
\end{aligned}$$

918 Substituting in (A.3) and using (2.7), (4.2) and taking expectations gives

$$\begin{aligned}
920 \quad \frac{d\mathbb{E}[V_t]}{dt} &\leq -\eta \mu_f \mathbb{E}[V_t] + \eta (f(\mathbf{x}^\dagger) - f(\mathbf{x}^\circ)) + \frac{1}{2} \sigma^2 \|\Delta \Phi^*\|_\infty. \\
921
\end{aligned}$$

922 Standard Grönwall arguments gives the result.  $\square$

923 Note that when (4.1)  $\mathbf{x}^\circ$ ,  $\mathbf{x}^*$  and  $\mathbf{x}^\dagger$  coincide [35, Lemma 7], so  $C_f = 0$ .

## 924 Appendix B. Additional information for the numerical results.

925 **B.1. Dynamics with hyperparameters and their discretization.** In (A.2)  
926 and (3.1) we included  $\eta, \epsilon$  to allow tuning of the relative effect of the gradient and the  
927 interaction. The expression of Proposition A.2 can be sharpened to include more precise  
928 contributions from  $\epsilon \mathcal{L} \mathbf{z}_t$  using similar Grönwall arguments to bound  $\mathbb{E}[\|z_t - z^\dagger\|^2]$   
929 instead of simply using (4.2). This is omitted here for brevity and we note  $\epsilon$  does  
930 have an effect in the numerical results. For the discretization, let  $\Delta t$  be a constant  
931 discretization interval and  $B_k^{i,j} \sim \mathcal{N}(0, \sigma^2 \Delta t)$  form i.i.d. sequences for  $j = 1, 2$  and  
932  $i = 1, \dots, N$ . Then the Euler discretization of the dynamics in (A.2) are given by

$$\begin{aligned}
933 \quad z_{(k+1)\Delta t}^i &= z_{k\Delta t}^i - \eta \nabla f_i(x_{k\Delta t}^i) \Delta t + \epsilon \sum_{j=1}^N A_{ij} (z_{k\Delta t}^j - z_{k\Delta t}^i) \Delta t + B_k^{i,1}, \\
934 \quad x_{(k+1)\Delta t}^i &= \nabla \Phi^*(z_{(k+1)\Delta t}^i). \\
935
\end{aligned}$$

936 Similarly for the EISMD dynamics of (3.3) we can include a learning rate  $\eta$  and  
937 interaction strength  $\epsilon$

$$\begin{aligned}
938 \quad (B.1) \quad d\mathbf{z}_t &= -\eta \nabla f(\mathbf{x}_t) dt - \epsilon \mathcal{L} \mathbf{z}_t dt - \mathcal{L} \boldsymbol{\lambda}_t dt + \sigma d\mathbf{B}_t, \quad d\boldsymbol{\lambda}_t = \mathcal{L} \mathbf{x}_t dt.
\end{aligned}$$

939 Similarly, the discretized dynamics of (B.1) are given by

$$\begin{aligned}
940 \quad z_{(k+1)\Delta t}^i &= z_{k\Delta t}^i - \Delta t \eta \nabla f_i(x_{k\Delta t}^i) + \Delta t \epsilon \sum_{j=1}^N A_{ij} (z_{k\Delta t}^j - z_{k\Delta t}^i) \\
&\quad + \Delta t \sum_{j=1}^N A_{ij} (\lambda_{k\Delta t}^j - \lambda_{k\Delta t}^i) + B_k^{i,2}, \\
941 \quad \lambda_{(k+1)\Delta t}^i &= \lambda_{k\Delta t}^i - \sum_{j=1}^N A_{ij} (x_{k\Delta t}^j - x_{k\Delta t}^i) \Delta t
\end{aligned}$$

941 The rest of the cases in Section 3 follow similarly.