

IMPERIAL COLLEGE LONDON

DEPARTMENT OF COMPUTING

---

# Mirror Descent and Interacting Mirror Descent

Almost Dimension-Free Convex Optimisation  
For Non-Euclidean Spaces

---

*Author:*

George Yiasemis

*Supervisor:*

Dr Panos Parpas

Submitted in partial fulfillment of the requirements for the MSc degree in Artificial  
Intelligence of Imperial College London

September 2020

## Abstract

*Many fundamental optimisation problems in machine learning and artificial intelligence have at the core of their frameworks the search of an optimal set parameters that lead to the minimisation or maximisation of a differentiable convex or non-convex objective function over either constrained or unconstrained feasible sets. A plethora of these problems utilise iterative methods that compute the gradient of the objective function at each iteration and take a step towards the steepest ascent or descent. A widely used optimisation method is Stochastic Gradient Descent. SGD is typically a well-performing algorithm for optimisation over Euclidean vector spaces. However, SGD does not always manage to converge to the optimum due to scaling of the dimensionality of the problem on non-Euclidean spaces, or due to noisy computations of the gradient and frequently proves to be expensive if the optimisation is constrained. In this project we study an other almost dimension-free iterative optimisation algorithm called Mirror Descent and its variant Interacting particles Mirror Descent on six convex frameworks with convex objective functions. Our aim is to observe the performance and convergence properties of these methods and also compare them with the standard Stochastic Gradient Descent method. We show that the mirror descent algorithms can outperform gradient descent in cases of optimisation over constrained sets with certain geometry, such as the probability simplex or the matrix spectrahedron. Additionally, we demonstrate the benefits in terms of convergence when allowing particles to interact in the case of Mirror Descent.*

# Acknowledgements

Undertaking this MSc individual project has been an unremarkable and challenging experience and it would not have been possible without the inspiration, support and guidance that I received.

First, I would like to thank my supervisor Dr Panos Parpas of the Department of Computing at Imperial College London for taking over the supervision of my proposal and for the support, counsel and recommendations he gave me over the past few months.

Additionally, I would like to thank my family and friends for their continuous help, encouragement and support during this difficult and challenging period.

Finally, I would like to thank both, Dr Robert A Craven, the MSc in Artificial Intelligence coordinator, and Dr Oana Cocarascu, Research Associate at the Department of Computing, for their valuable help and guidance throughout the duration of my studies.

# List of Figures

1.1	Convex Optimisation (Left) and Non-convex Optimisation (Right). . .	2
1.2	Steps of Gradient Descent towards the global minimum for a simple convex function in $\mathbb{R}^2$ , Graph of the function (Left) and Contour levels (Right). . . . .	3
4.1	Examples of convex sets: An ellipsoid in $\mathbb{R}^2$ (Left) and a sphere in $\mathbb{R}^3$ (Right). . . . .	11
4.2	Illustration of Projected Gradient Descent. . . . .	20
4.3	Illustration of Mirror Descent in accordance with Equations 4.17a-4.17d. . . . .	23
5.1	Traffic Assignment Problem: Random Geographical Threshold Graph, $ \mathcal{V}  = 50$ ( $\mathcal{O}, \mathcal{T}$ ) = (46, 32). . . . .	35
5.2	Traffic Assignment Problem: <b>(a)</b> Convergence of deterministic GD, MD and IMD in the case of $n_{\mathcal{R}} = 70$ , <b>(b)</b> Convergence of SGD, SMD and SIMD in the case of $n_{\mathcal{R}} = 1169$ . . . . .	36
5.3	Traffic Assignment Problem: Convergence of stochastic optimisation with decreasing learning rate: SGD, SMD and SIMD in the case of higher dimensionality, $n_{\mathcal{R}} = 1647$ . . . . .	37
5.4	Traffic Assignment Problem: Convergence of deterministic optimisation with fixed learning rate: GD, MD and IMD in the case of higher dimensionality, $n_{\mathcal{R}} = 1647$ . . . . .	38
5.5	Linear System: SGD convergence (L) vs SMD convergence (R) for a well-conditioned $W$ with $\kappa(\mathbf{W}) = 10$ . . . . .	40
5.6	Linear System: SGD convergence (L) vs SMD convergence (R) for an ill-conditioned $W$ with $\kappa(\mathbf{W}) = 100$ . . . . .	40
5.7	Linear System: (S)MD convergence vs (S)IMD convergence for a well-conditioned $W$ with $\kappa(\mathbf{W}) = 10$ . . . . .	41
5.8	Linear System: Histogram of (S)MD and (S)IMD convergence for a well-conditioned $W$ with $\kappa(\mathbf{W}) = 10$ . . . . .	41
5.9	Linear System: Effect of number of interacting particles on variance of samples of the loss function after convergence for cases of ill and well-conditioned $\mathbf{W}$ , $\kappa(\mathbf{W}) = 10, 200$ . . . . .	42
5.10	Mini-Batch Optimisation: Loss function for SGD and SMD ( $N_{iid} = 10$ ) and SIMD ( $N_p = 10$ ) for $M_{\mathcal{B}}$ in the case of $\kappa(\mathbf{W}) = 100$ . . . . .	45

5.11 Mini-Batch Optimisation: <b>(a)</b> Loss function and <b>(b)</b> Histogram of loss samples after convergence for SIMD ( $\sigma = 1.0$ ) for increasing number of interacting particles and fixed batch size $M_{\mathcal{B}} = 5$ in the case of $\kappa(\mathbf{W}) = 10$ . . . . .	46
5.12 Mini-Batch Optimisation: <b>(a)</b> Loss function and <b>(b)</b> Histogram of loss samples after convergence for SIMD ( $\sigma = 0.01$ ) for different combinations of batch sizes ( $M_{\mathcal{B}}$ ) and number of interacting particles ( $N_p$ ) in the case of $\kappa(\mathbf{W}) = 10$ . . . . .	47
5.13 Mini-Batch Optimisation: <b>(a)</b> Loss function and <b>(b)</b> Histogram of loss samples after convergence for SIMD ( $\sigma = 0.01$ ) for different combinations of batch sizes ( $M_{\mathcal{B}}$ ) and number of interacting particles ( $N_p$ ) in the case of $\kappa(\mathbf{W}) = 200$ . . . . .	47
5.14 Laplacian $K$ -Modes for Clustering: $N = 2000$ randomly generated data points in the 3D space divided into $K = 10$ clusters (isotropic Gaussian blobs). The black dots illustrate the (true) centroid $\mathbf{c}_k$ of each cluster. . . . .	51
5.15 Laplacian $K$ -Modes for Clustering: $N = 1000$ randomly generated data points in the 2D space divided into $K = 2$ half-moon shaped clusters. . . . .	51
5.16 Laplacian $K$ -Modes for Clustering: Cluster assignment results of (a) GD and (b) IMD using a classification threshold equal to 0.1. Their achieved numerical “accuracies” are reported in Table 5.4. . . . .	52
5.17 Laplacian $K$ -Modes for Clustering: Loss convergence of GD and IMD for $N = 2000$ . . . . .	53
5.18 Laplacian $K$ -Modes for Clustering: <b>(a)</b> Loss function and <b>(b)</b> Histogram of loss samples after convergence for the three methods for $N = 1000$ . . . . .	54
5.19 Laplacian $K$ -Modes for Clustering: Cluster assignment results of <b>(a)</b> GD, <b>(b)</b> MD and <b>(c)</b> SIMD using a classification threshold equal to 0.5 in the case of two synthetically generated half-moons clusters. The black points illustrate the predicted centroids $\mathbf{c}_1$ and $\mathbf{c}_2$ . . . . .	55
5.20 Disease Outbreak Anomaly Detection: a real case scenario, Graph of NE United States. Nodes colored blue illustrate $H_0$ for the anomalous subgraph (denoted $\mathcal{V}_S$ in the previous paragraph). The anchor node is colored purple. . . . .	57
5.21 Disease Outbreak Anomaly Detection: AUC scores of MD (averaged over 5 iid runs) for disease ratios $\{1.1, 1.3, 1.5, 1.7\}$ and different values of $\gamma$ . . . . .	58
5.22 Disease Outbreak Anomaly Detection: AUC scores of IMD with 2 interacting particles for disease ratios $\{1.1, 1.3, 1.5, 1.7\}$ and different values of $\gamma$ . . . . .	58
5.23 Disease Outbreak Anomaly Detection: Best AUC scores of MD and IMD for disease ratios $\{1.1, 1.3, 1.5, 1.7\}$ . . . . .	59

5.24	Random Geometric Graphs Anomalous Clusters Classification: 3D Random generated Geometric Graph consisting of 200 vertices. Nodes with higher degrees (more neighbours) are illustrated with a larger radius and lighter colours. . . . .	60
5.25	Random Geometric Graphs Anomalous Clusters Classification: AUC performance of MD and IMD for different anomaly sizes $K$ . . . . .	62
5.26	Random Geometric Graphs Anomalous Clusters Classification: Average run-time (in seconds) per iteration for the processing of one sample for various graph sizes $n$ . . . . .	63
5.27	Mirror Langevin Diffusion: Convergence of the squared distance of the running mean $\bar{\mathbf{x}}_t$ from $\mathbf{0}$ for $d = 100$ in the case of NLA averaged over $N_{iid} = 10$ runs and INLA with $N_p = 10$ . . . . .	67
5.28	Mirror Langevin Diffusion: Convergence of the relative squared distance of the running estimate of $\hat{\Sigma}$ from $\Sigma$ for $d = 100$ in the case of NLA averaged over $N_{iid} = 10$ runs and INLA with $N_p = 10$ . . . . .	67
5.29	Mirror Langevin Diffusion: Histogram of the relative squared distance of the running estimate of $\hat{\Sigma}$ from $\Sigma$ for $d = 100$ in the case of NLA averaged over $N_{iid} = 10$ runs and INLA with $N_p = 10$ after convergence ( $T_0 = 200$ ). . . . .	68

# List of Tables

5.1	Variance after convergence (after 200 iterations of SIMD) for different number of interacting particles and condition number $\kappa(\mathbf{W})$ . Illustrated above in histograms. . . . .	43
5.2	Mini-Batch Optimisation: Variance of SIMD for fixed mini-batch size and increasing number of interacting particles after convergence. . . . .	46
5.3	Mini-Batch Optimisation: Variance of SIMD for different number of interacting particles ( $N_p$ ) and different mini-batch sizes ( $M_B$ ) after convergence. . . . .	48
5.4	Laplacian $K$ -Modes for Clustering: Classification performance of GD and IMD in the case of $K = 10$ synthetically generated Isotropic Gaussian 3D blobs. Each blob contains approximately $\frac{N}{K} = \frac{2000}{10} = 200$ data samples. . . . .	52
5.5	Laplacian $K$ -Modes for Clustering: Average run-time for one optimisation iteration for the three methods in the case of synthetic Isotropic Gaussian blobs. . . . .	53
5.6	Laplacian $K$ -Modes for Clustering: Classification performance of the three methods in the case of two synthetically generated half-moons clusters. . . . .	54
5.7	Disease Outbreak Anomaly Detection: Comparison of MD and IMD for the optimal choice of $\gamma$ as illustrated in Figures 5.21 and 5.22. . . . .	59
5.8	Disease Outbreak Anomaly Detection: Average time for the processing of one sample for MD and IMD. . . . .	59
5.9	Random Geometric Graphs Anomalous Clusters Classification: Performance of IMD and MD for different SNR values. . . . .	61
5.10	Mirror Langevin Diffusion: Variance of samples of $\log_{10} \frac{\ \hat{\Sigma} - \Sigma\ _F^2}{\ \Sigma\ _F^2}$ after convergence ( $T_0 = 200$ ) for NLA and INLA. . . . .	68
B.1	Imperial College London Ethics Checklist . . . . .	80

# List of Algorithms

1	Projected Gradient Descent Algorithm . . . . .	21
2	(Stochastic) Mirror Descent Algorithm . . . . .	25
3	(Stochastic) Interacting Mirror Descent Algorithm . . . . .	28



# List of Abbreviations

AI	Artificial Intelligence
AUC	Area Under the roc Curve
DL	Deep Learning
GD	Gradient Descent
IMD	Interacting Mirror Descent
IMLD	Interacting Mirror Langevin Diffusion
INLA	Interacting Newton Langevin Algorithm
INLD	Interacting Newton Langevin Diffusion
LD	Langevin Diffusions
MCMC	Markov Chain Monte Carlo
MD	Mirror Descent
ML	Machine Learning
MLD	Mirror Langevin Diffusion
NLA	Newton Langevin Algorithm
NLD	Newton Langevin Diffusion
PGD	Projected Gradient Descent
PSD	Positive Definite
SDE	Stochastic Differential Equation
SDP	Semidefinite Progtamming
SGD	Stochastic Gradient Descent
SIMD	Stochastic Interacting Mirror Descent
SL	Statistical Learning
SMD	Stochastic Mirror Descent
SNR	Signal-to-Noise Ratio
tr	Trace

# Contents

<b>List of Figures</b>	<b>ii</b>
<b>List of Tables</b>	<b>vi</b>
<b>List of Algorithms</b>	<b>vii</b>
<b>List of Abbreviations</b>	<b>viii</b>
<b>1 Introduction and Outline</b>	<b>1</b>
1.1 Introduction . . . . .	1
1.2 Outline . . . . .	4
<b>2 Related Literature</b>	<b>5</b>
2.1 Gradient Descent . . . . .	5
2.2 Mirror Descent . . . . .	6
2.3 Interacting Mirror Descent . . . . .	6
<b>3 Contribution</b>	<b>8</b>
<b>4 Preliminaries</b>	<b>10</b>
4.1 Notation . . . . .	10
4.2 Background and Problem Setup . . . . .	11
4.2.1 Background and Basic Definitions . . . . .	11
4.2.2 Main Problem Setup . . . . .	16
4.2.3 Examples . . . . .	16
4.3 Mirror Descent . . . . .	20
4.3.1 The (Projected) Gradient Descent Algorithm . . . . .	20
4.3.2 The Mirror Descent Algorithm . . . . .	21
4.4 Interacting Mirror Descent . . . . .	26
4.5 Applications . . . . .	28
4.5.1 Optimisation on the Probability Simplex . . . . .	28
4.5.2 Optimisation on the Spectrahedron . . . . .	32
<b>5 Experimental Results</b>	<b>33</b>
5.1 The Traffic Assignment problem . . . . .	33
5.1.1 Formulation of the traffic assignment Problem . . . . .	33
5.1.2 Numerical Results . . . . .	35

	Experiments for random graph with $ \mathcal{V}  = 50$ . . . . .	36
	Experiments for random graph with $ \mathcal{V}  = 100$ . . . . .	37
5.2	Linear System . . . . .	39
5.2.1	Formulation of the linear system problem with simplex constraints	39
5.2.2	Numerical Results . . . . .	39
	Single particle Optimisation: (S)GD vs (S)MD . . . . .	39
	Interacting particles Optimisation . . . . .	40
	Comparing S(MD) vs (S)IMD. . . . .	40
	Interacting particles for Variance Reduction . . . . .	42
5.3	Linear System: Mini-batch Optimisation . . . . .	44
5.3.1	Problem formulation . . . . .	44
5.3.2	Numerical Results . . . . .	45
	Effect of number of particles $N_p$ in Mini-Batch Optimisation . .	45
	Experiments for fixed mini-batch size $M_{\mathcal{B}}$ . . . . .	45
	Experiments for different combinations of $N_p$ and $M_{\mathcal{B}}$ . .	46
5.4	The Laplacian $K$ -Modes for Clustering . . . . .	49
5.4.1	Problem formulation . . . . .	49
5.4.2	Numerical Results . . . . .	50
	Isotropic Gaussian blobs . . . . .	52
	Two “half-moons” . . . . .	54
5.5	Optimisation on SDPs . . . . .	56
5.5.1	Connected Subgraph Detection on SDPs: Problem formulation .	56
5.5.2	Disease Outbreak Anomaly Detection: Numerical Results . . . .	57
	Performance of IMD vs MD for different values of $\gamma$ . . . . .	58
5.5.3	Random Geometric Graphs Anomalous Clusters Classification:	
	Numerical Results . . . . .	60
	Performance of MD and IMD vs SNR values . . . . .	61
	Performance of MD and IMD vs Anomaly size $K$ . . . . .	61
	Performance of MD and IMD vs graph size $n$ . . . . .	62
5.6	Langevin Diffusions . . . . .	64
5.6.1	Mirror Langevin Diffusion . . . . .	64
5.6.2	Interacting Mirror Langevin Diffusion . . . . .	65
5.6.3	Numerical Results . . . . .	66
<b>6</b>	<b>Conclusion and Future Work</b>	<b>69</b>
6.1	Conclusion . . . . .	69
6.2	Future Work . . . . .	70
	<b>Appendices</b>	<b>76</b>
<b>A</b>	<b>Code</b>	<b>77</b>
<b>B</b>	<b>Ethical Considerations</b>	<b>78</b>

# Chapter 1

## Introduction and Outline

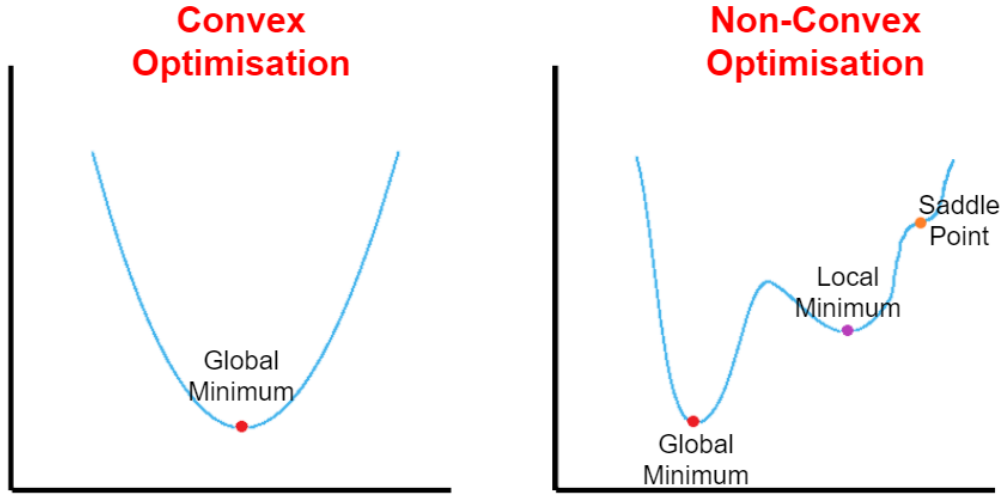
### 1.1 Introduction

The term “optimisation” in Mathematics usually refers to the process of selecting the “best” parameters out of a set of feasible parameters, for example  $\mathcal{X}$ , which optimise (maximise or minimise) an objective function  $f : \mathcal{X} \rightarrow \mathbb{R}$ . A basic optimisation problem has usually one of the two forms:

$$\max_{\mathbf{x} \in \mathcal{X}} f(\mathbf{x}), \quad \min_{\mathbf{x} \in \mathcal{X}} f(\mathbf{x}),$$

where  $\mathcal{X}$  can be bounded (e.g. a bounded subset of  $\mathbb{R}^n$ , such as a hyperball or a hypercube) or not bounded (e.g.  $\mathbb{R}^n$ ). If  $\mathcal{X}$  is bounded then the optimisation is called “constrained”, otherwise it is called “unconstrained”. Constrained optimisation techniques can be useful in a plethora of applications, not only in Mathematics, but also in Engineering, Mechanics, Economics, Finance etc. A special case of optimisation is convex optimisation which aims to optimise a convex function  $f$  over a convex set  $\mathcal{X}$ . Otherwise, we call it non-convex optimisation. See Figure 1.1 for an illustrative example.

Optimisation algorithms are key factors in almost every Machine Learning, Statistical Learning, and Deep Learning application and the optimiser is one of the most important contributors to a well-performed model. Aiming to minimise or maximise an objective function, the optimiser plays the role of the decision-making factor who decides how the values of parameters of a model should be updated, i.e. increased, decreased or remain unchanged, and thereby make the model we are training more accurate.



**Figure 1.1:** Convex Optimisation (Left) and Non-convex Optimisation (Right).

Such parameters can be coefficients  $\mathbf{m} \in \mathbb{R}^n$  in simple problems such as Linear Regression

$$f(\mathbf{m}) = \frac{1}{N} \sum_{i=1}^N (y_i - \mathbf{m}^T \mathbf{x}_i)^2,$$

or weights or biases  $(\mathbf{W}, \mathbf{b})$  of a layer of an artificial neural network

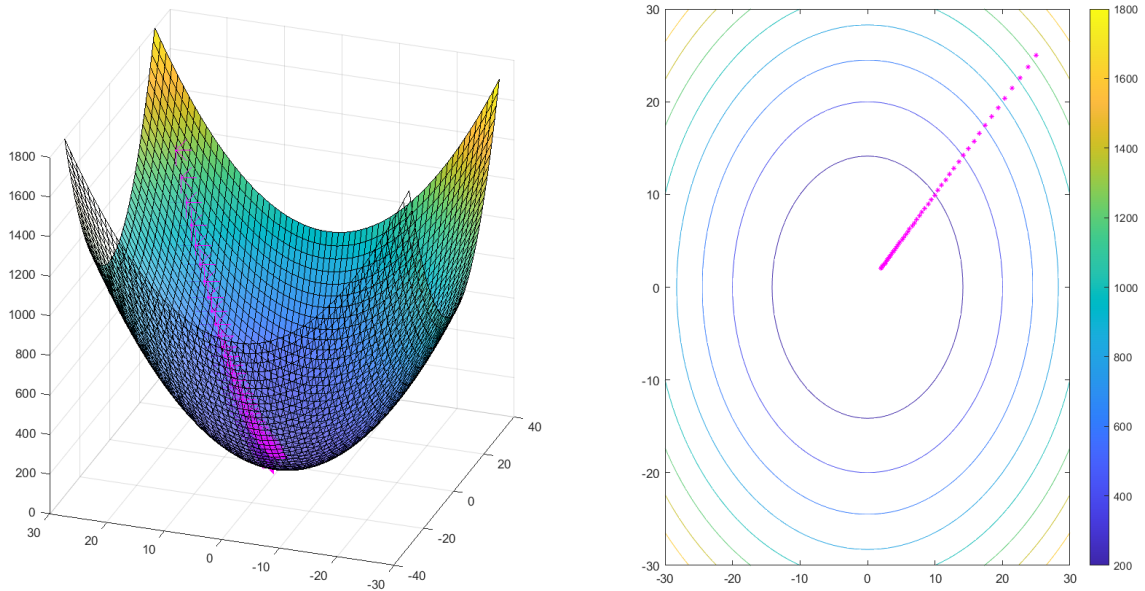
$$f(\mathbf{W}, \mathbf{B}) = \|\mathbf{Y} - g(\mathbf{W}^T \mathbf{X} - \mathbf{B})\|_2^2$$

where  $g$  is an activation function such as ReLu (rectified linear unit), the sigmoid or the hyperbolic tangent, etc.

Despite the fact that the vast majority of problems arising in Artificial Intelligence and Statistical Learning are non-convex, in this project we put our focus on convex and strictly convex cases since there is a guarantee of existence of a global (unique) minimum and therefore are simpler to solve. In contrast, an optimisation algorithm in a non-convex case can get stuck on local optima or saddle points (see Figure 1.1 for an illustration).

In AI, ML and DL contexts, Stochastic Gradient Descent is a first-order optimisation algorithm exploited to minimise an objective function by iteratively taking steps in the direction of the steepest descent utilising the negative of the function's gradient (vector of partial derivatives). Gradient Descent then updates the parameters of the model such that at each iteration the value of the objective function is decreasing. However, the size of the update is often controlled by the step size (or learning rate) which is usually treated as a hyperparameter. It should be mentioned that Stochastic Gradient Descent and all of its variants are usually the first choice for an optimiser in many applications of AI.

For an illustration of the performance of Gradient Descent with constant learning rate on a convex problem in  $\mathbb{R}^2$  see Figure 1.2.



**Figure 1.2:** Steps of Gradient Descent towards the global minimum for a simple convex function in  $\mathbb{R}^2$ , Graph of the function (Left) and Contour levels (Right).

In the case of constrained optimisation over a constraint set  $\mathcal{X}$ , Projected Gradient Descent (PGD) can be exploited. In PGD at each iteration the next point is projected onto  $\mathcal{X}$  utilising the negative gradient of the objective function at the previous point. However, projection can be computationally expensive as it requires at each iteration one computation of the gradient and one projection onto  $\mathcal{X}$  ([45]) and it is not always feasible.

One of the goals of this individual project is the investigation of another first order optimisation algorithm called Mirror Descent. Mirror Descent is a generalization of PGD that works on both, Euclidean and non-Euclidean geometries. Nemirovskii and Yudin in their work [40] introduced the idea of a “mirror map” function. As it will be explicitly described later in Chapter 4, Mirror Descent converts a constrained problem into an unconstrained problem utilising the mirror map that fits the problem and its geometry, and thus, it removes the requirements of projecting points back onto the constraints set. We also provide illustrations for both methods in Chapter 4.

In their works [5, 13] the authors provide some theoretical results concerning the convergence of both, MD and PGD. In our work we present numerical results (Chapter 5) for six problems (five are constrained and one is unconstrained) which verify their theoretical results.

Another aim of this work is the analysis of an interacting variant of MD, called Interacting Mirror Descent initially proposed in the works [9, 10]. We study the properties of the convergence of the method when it is used to optimise six problems in Chapter 5. We also make comparisons with the case of single particle optimisation (Mirror Descent and Projected Gradient Descent) which are averaged over a number of independent and identical replicas.

## 1.2 Outline

The rest of this paper is structured into five chapters, Chapter 2 (Related Literature), Chapter 3 (Contribution), Chapter 4 (Preliminaries), Chapter 5 (Experimental Results) and Chapter 6 (Conclusion and Future Work). Here we provide a short description for each of the next chapters.

- **Related Literature.** In this chapter we provide a summary of previous research performed on the topics of this individual project (Gradient Descent, Mirror Descent, Interacting Mirror Descent) as found in the work of related literature.
- **Contribution.** In this chapter we provide a synopsis of the main contributions of this project to the literature.
- **Preliminaries.** In this chapter we provide preliminaries, background notes and some theory on the optimisation methods used in this individual project. In the case of Mirror Descent we provide some examples of frameworks (optimisation on the Simplex and on the Spectrahedron) we are going to be using for our Experiments. Along with the above, we provide step-by-step algorithms for the implementation of (Projected) Gradient Descent, Mirror Descent and its variant Interacting Mirror Descent.
- **Experimental Results.** In this chapter we give insight on the formulation of the six problems we are concerned with. We define each one in detail and provide their formulations. Furthermore, we provide numerical experimental results on all six problems. We additionally, discuss and compare the performance of the three methods we use in this project.
- **Conclusion and Future Work.** In the last chapter we provide some conclusions as well as possible future directions.

# Chapter 2

## Related Literature

### 2.1 Gradient Descent

As afformentioned in 1.1 of Chapter 1, Gradient Descent is the most widely used optimisation framework in many Machine Learning and Deep Learning tasks. Zhang in his work [58] collocates some applications of Gradient Descent and its variants (Stochastic Gradient Descent, Mini-batch Gradient Descent, etc) for finding the optimal parameters of Deep Learning models.

Additionally, in the literature there have been proposed many first-order variants of Stochastic Gradient Descent used to optimise objective functions. Some of them are namely Adam, AdaGrad, RMSProp (see Chapter 8 of [22] for a complete guide) which are widely used and seem to work well in cases of unconstrained convex or non-convex optimisation.

As we stated in the previous Chapter, in the case of constrained optimisation over a constraint set  $\mathcal{X}$ , Projected Gradient Descent is widely used. In PGD at each iteration the next point is projected onto  $\mathcal{X}$  utilising the gradient of the objective function at the previous point. PGD has a variety of applications. Some interesting works found in the Machine and Deep Learning literature that utilise PGD as an optimiser are [3, 25, 31, 33, 34, 49, 56, 48, 19, 52, 23, 8, 32]. However, the operation of projecting a single point onto a constraint set can be computationally expensive and not always feasible [45].

One problem we are concerned with in this project in particular, are the works of [51] who propose a new algorithm for clustering called The Laplacian k-modes for clustering. The authors propose an efficient method based on Projected Gradient Descent in [50] for solving this problem over the probability simplex. We will later solve this problem using the Mirror Descent and Interacting Mirror Descent setups.



## 2.2 Mirror Descent

The algorithm of Mirror Descent, proposed initially by Nemirovskii and Yudin in [40], has become one of the most popular optimisation frameworks in ML and SL. MD and its variants have been used for optimising various convex and non-convex tasks by many authors in the related literature which includes the works in [5, 15, 18, 46, 39, 38, 36, 61, 55, 37, 30, 29, 53].

As mentioned in the introduction, Mirror Descent can be used to transform a constrained problem into an unconstrained one. This property of MD was utilised by the authors who used as a mirror map the “entropy” function to solve optimisation problems on simplex setups. We report here some of the works in the bibliography which we are interested in.

In [28], the authors make use of mirror descent for first-order sampling from a constrained -on the simplex- distribution and show that their method Mirrored Langevin Dynamics outperformed the previous state-of-the art method that utilised Stochastic Gradient Descent. Additionally, Mertikopoulos and Staudigl in [37] solve a version of the traffic assignment problem where the feasible space is again the probability simplex by employing MD dynamics.

Another work we are concerned with is Connected Subgraph Elevated Mean Detection on semi-definite positive programs which includes matrix constrained, convex optimisation. The authors in [1], exploiting Mirror Descent, solve this problem over the Spectrahedron (space of semi-definite positive with fixed trace, analogous to the probability simplex but in  $\mathbb{R}^{n \times n}$ ) and they show that the performance of their Mirror Descent algorithm on Elevated Mean Detection surpasses the performance of the methods previously proposed.

## 2.3 Interacting Mirror Descent

In the past years there has been a great deal of research in Statistical Theory around interacting particle systems, which they have found applications in many areas. We report some works in the related literature. The authors in [7] use  $n$  interacting particles of diffusions for continuous and discrete sampling. In [20], the authors also allow particles to interact in a system of diffusion processes and they show that an increase in interaction results in an increase of the systemic risk. Moreover, in [16], they use a mean interaction field (which we will use here as well) and they propose a maximum likelihood based inference to estimate the parameters of the interaction function induced by a potential energy in an  $n$  interacting particles system. Another interesting application of interacting particles in Finance, was presented in [14], where the authors introduce the use of interacting particle systems in the computation of probabilities of simultaneous defaults in large credit portfolios.

In our work here, we also utilise interacting particles in our optimisation setup. We adopt a method from [10], where the authors extend Stochastic Mirror Descent, by proposing a new optimisation framework that estimates the optimum of a function in which they allow  $N_p$  particles to interact. The optimisation update rule they propose is given by the following Itô Stochastic Differential Equation in continuous time:

$$d\mathbf{y}_t^i = -\nabla f(\mathbf{x}_t^i)dt + \sum_{j=1}^{N_p} A_{ij}(\mathbf{y}_t^j - \mathbf{y}_t^i) + \sigma d\mathbf{B}_t^i,$$

$$\mathbf{x}_t^i = \nabla \Psi^*(\mathbf{y}_t^i), \quad i = 1, \dots, N_p,$$

where each particle is characterised by independent Brownian motions  $\mathbf{B}_t^i$  and  $\mathbf{A}$  denotes a matrix that contains the interaction weight of each particle. Without further analysis of the equation (see Chapter 4 for more on this equation), we mention some of their findings in the next paragraph.

Among with some theoretical results, they exhibit results on three optimisation problems. The authors using a mean interacting field ( $A_{ij} = \frac{1}{N_p}$ ) show that interaction can speed up optimisation convergence and helps in variance reduction. Even in the case of mini-batch optimisation, they show that for large number of interacting particles the variance of the the algorithm around the optimum can be comparable to that when optimising over a full-batch. In this project we reproduce and verify some of the results of [10].

# Chapter 3

## Contribution

This MSc individual project aims at researching two optimisation algorithms, defined later on, (Stochastic) Mirror Descent (SMD/MD) and its variant (Stochastic) Interacting Mirror Descent (SIMD/IMD) proposed in [10]. We investigate their convergence properties when applied on convex optimisation frameworks and we make comparisons against the traditional optimisation technique, Projected (Stochastic) Gradient Descent (SGD/GD). We present numerical experimental results on six problems appearing in the related literature from which, for the last three, this project is the first to report results of the application of the discretised version of Interacting Mirror Descent algorithm.

We provide a synopsis of the contributions of this project:

- We provide the reader with a complete collection of the basic background (in Chapter 4) for optimisation with Mirror Descent and its variant Interacting Mirror Descent.
- We reproduce experiments on three constrained optimisation schemes on the probability simplex; for a Linear System, Mini-batch optimisation on a Linear System and the Traffic Assignment problem (see Sections 5.2, 5.3, 5.1 respectively). We compare the performance of (S)GD with that of MD and its interacting variant.
- We demonstrate how interaction in Interacting Mirror Descent can help in variance/noise reduction which is present in cases such as batch optimisation or when samples used for optimisation are drawn from a noisy distribution. Similarly to [10], we show that an increase in interacting particles can lead to a decline of noise/variance.
- Using Mirror Descent and Interacting Mirror Descent we solve numerically two clustering problems (Section 5.4) for random generated data on the probability simplex proposed in [50, 51] which they solve with Stochastic Gradient Descent. We show experimentally that the former, surpasses the latter in terms of performance, convergence and speed of convergence.

- We also solve Elevated Mean detection on two problems (of which one is a real case scenario and one is obtained from randomly generated data) on positive semi-definite matrices on the Spectrahedron previously proposed and solved in [1] with Mirror Descent (see Section 5.5 for the formulation). Here we utilise the Interacting particle Mirror Descent framework. Our numerical results show that the performance of Interacting Mirror Descent surpasses that of the original version of Mirror Descent.
- Lastly, in Section 5.6 we propose a novel algorithm called Interacting Newton Langevin Algorithm (INLA) used to sample from a target distribution. Our method allows  $N_p$  particles to interact and attains better performance than the original algorithm Newton Langevin Algorithm (NLA) proposed in the related literature (see [17]). Our algorithm outperforms NLA in convergence speed and we also show that interaction helps in reducing the variance.

# Chapter 4

## Preliminaries

### 4.1 Notation

For any two vectors  $\mathbf{x}, \mathbf{y} \in \mathbb{R}^n$ :

- $\langle \cdot, \cdot \rangle_2 : \mathbb{R}^n \times \mathbb{R}^n \rightarrow \mathbb{R}$  will denote the Euclidean inner product in  $\mathbb{R}^n$ , such that  $\langle \mathbf{x}, \mathbf{y} \rangle_2 = \mathbf{x}^T \mathbf{y}$ .

- $\| \cdot \|_2 : \mathbb{R}^n \rightarrow \mathbb{R}_0^+$  will denote the Euclidean  $\mathcal{L}2$ -norm on  $\mathbb{R}^n$ , such that

$$\| \mathbf{x} \|_2 = \sqrt{\langle \mathbf{x}, \mathbf{x} \rangle_2} = \sqrt{\mathbf{x}^T \mathbf{x}} = \sqrt{\sum_{i=1}^n x_i^2}.$$

- $\| \cdot \|_1 : \mathbb{R}^n \rightarrow \mathbb{R}_0^+$  will denote the  $\mathcal{L}1$ -norm on  $\mathbb{R}^n$ , such that

$$\| \mathbf{x} \|_1 = \sum_{i=1}^n |x_i|.$$

For any two real matrices  $\mathbf{A}, \mathbf{B} \in \mathbb{R}^{n \times n}$ :

- $\langle \cdot, \cdot \rangle_F : \mathbb{R}^{n \times n} \times \mathbb{R}^{n \times n} \rightarrow \mathbb{R}$  will denote the Frobenius inner product in  $\mathbb{R}^{n \times n}$ ,

$$\text{such that } \langle \mathbf{A}, \mathbf{B} \rangle_F = \text{tr}(\mathbf{A}^T \mathbf{B}) = \mathbf{A} \bullet \mathbf{B} = \sum_{i,j=1}^n A_{ij} B_{ij}.$$

- $\| \cdot \|_F : \mathbb{R}^{n \times n} \rightarrow \mathbb{R}_0^+$  will denote the Frobenius norm on  $\mathbb{R}^{n \times n}$ , such that

$$\| \mathbf{A} \|_F = \sqrt{\langle \mathbf{A}, \mathbf{A} \rangle_F} = \sqrt{\text{tr}(\mathbf{A}^T \mathbf{A})} = \sqrt{\sum_{i=1}^n \lambda_i(\mathbf{A})^2},$$

where  $\lambda_i(\mathbf{A})$  denotes the  $i$ -th eigenvalue of  $\mathbf{A}$ .

- $\| \cdot \|_1 : \mathbb{R}^{n \times n} \rightarrow \mathbb{R}_0^+$  will denote the nuclear norm or Shatten-1 norm on  $\mathbb{R}^{n \times n}$ , such that

$$\| \mathbf{A} \|_1 = \text{tr}(\mathbf{A}^T \mathbf{A})^{\frac{1}{2}}.$$

- $\kappa(\cdot) : \mathbb{R}^{n \times n} \rightarrow \mathbb{R}$  will denote the condition number of matrices in  $\mathbb{R}^{n \times n}$ , such that

$$\kappa(\mathbf{A}) = \frac{\lambda_{\max}(\mathbf{A})}{\lambda_{\min}(\mathbf{A})},$$

where  $\lambda_{\max}(\mathbf{A})$  and  $\lambda_{\min}(\mathbf{A})$  denote the maximum and minimum eigenvalues of matrix  $\mathbf{A}$ .

## 4.2 Background and Problem Setup

### 4.2.1 Background and Basic Definitions

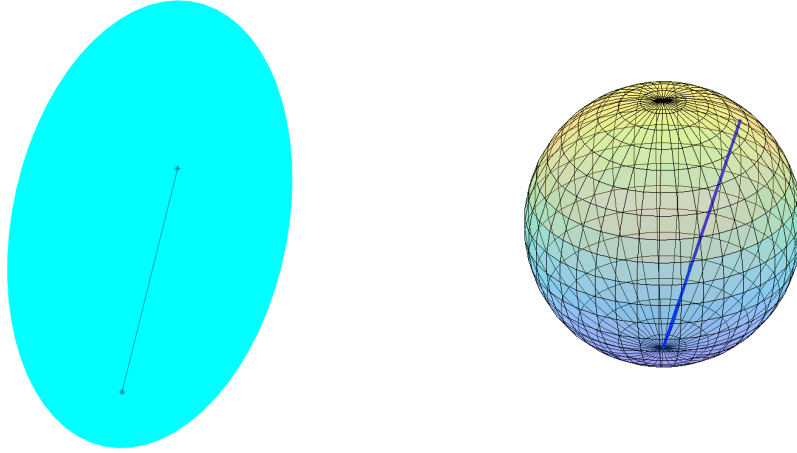
First, we provide some definitions that are needed to define our main problem in this project. We also provide all necessary definitions essential to outline the optimisation algorithms that we use.

#### Definition 1

(Convex Set)

A set  $\mathcal{X}$  is said to be convex if every point on the line connecting any two points  $\mathbf{x}$  and  $\mathbf{z}$  in  $\mathcal{X}$  is itself in  $\mathcal{X}$ . More formally

$$\alpha \mathbf{x} + (1 - \alpha) \mathbf{z} \in \mathcal{X} \quad \forall \alpha \in (0, 1). \quad (4.1)$$



**Figure 4.1:** Examples of convex sets: An ellipsoid in  $\mathbb{R}^2$  (Left) and a sphere in  $\mathbb{R}^3$  (Right).

#### Definition 2

(Convex Function)

A function  $f : \mathcal{X} \rightarrow \mathbb{R}$ , defined on a convex set  $\mathcal{X}$ , is called convex on  $\mathcal{X}$  if the line

segment connecting  $f(\mathbf{x})$  and  $f(\mathbf{z})$  at any two points  $\mathbf{x}, \mathbf{z} \in \mathcal{X}$  lies above the function between  $\mathbf{x}$  and  $\mathbf{z}$ , i.e.

$$f(\alpha\mathbf{x} + (1 - \alpha)\mathbf{z}) \leq \alpha f(\mathbf{x}) + (1 - \alpha)f(\mathbf{z}), \quad \forall \alpha \in (0, 1). \quad (4.2)$$

Now, if

$$f(\alpha\mathbf{x} + (1 - \alpha)\mathbf{z}) < \alpha f(\mathbf{x}) + (1 - \alpha)f(\mathbf{z}), \quad \forall \alpha \in (0, 1), \quad (4.3)$$

then  $f$  is called strictly convex on  $\mathcal{X}$ .

### Definition 3

(Strong Convexity)

A differentiable function  $f : \mathcal{X} \rightarrow \mathbb{R}$ , defined on a convex set  $\mathcal{X}$ , is called strongly convex with constant  $\mu > 0$  on  $\mathcal{X}$  with respect to an arbitrary norm defined on  $\mathcal{X}$ , if for any two points  $\mathbf{x}, \mathbf{z} \in \mathcal{X}$  it holds:

$$f(\mathbf{z}) \geq f(\mathbf{x}) + \nabla f(\mathbf{x})^T(\mathbf{z} - \mathbf{x}) + \frac{\mu}{2} \|\mathbf{z} - \mathbf{x}\|^2. \quad (4.4)$$

### Definition 4

(Lipschitz function)

We say that  $f$  is  $L_f$ -Lipschitz with respect to  $\|\cdot\| : \mathbb{R}^n \rightarrow \mathbb{R}_0^+$  an arbitrary norm on  $\mathbb{R}^n$  if:

$$|f(\mathbf{x}) - f(\mathbf{y})| \leq L_f \|\mathbf{x} - \mathbf{y}\|.$$

### Definition 5

(Dual Norm)

Let  $\|\cdot\| : \mathbb{R}^n \rightarrow \mathbb{R}_0^+$  be an arbitrary norm on  $\mathbb{R}^n$ . We will denote with  $\|\cdot\|_*$  the dual norm defined as:

$$\|\mathbf{g}\|_* = \max_{\mathbf{x} \in \mathbb{R}^n : \|\mathbf{x}\| \leq 1} \langle \mathbf{g}, \mathbf{x} \rangle_2 = \max_{\mathbf{x} \in \mathbb{R}^n : \mathbf{x} \neq \mathbf{0}} \frac{1}{\|\mathbf{x}\|} \langle \mathbf{g}, \mathbf{x} \rangle_2.$$

We will say that  $\nabla f$  is  $L_{\nabla f}$ -Lipschitz on  $\mathcal{X}$  with constant  $L_{\nabla f} > 0$  if  $\forall \mathbf{x}, \mathbf{y} \in \mathcal{X}$ :

$$\|\nabla f(\mathbf{x}) - \nabla f(\mathbf{y})\|_* \leq L_{\nabla f} \|\mathbf{x} - \mathbf{y}\|.$$

### Definition 6

( $n$ -Simplex)

We define as the  $n$ -Simplex the set

$$\Delta^n = \left\{ \mathbf{x} \in \mathbb{R}^n \mid \sum_{i=1}^n x_i = \|\mathbf{x}\|_1 = 1, x_i \geq 0, i = 1, 2, \dots, n \right\}.$$

$\Delta^n$  is usually called the probability simplex.

### Definition 7

( $n$ -Spectrahedron)

Let  $\mathcal{S}_+^{n \times n} = \left\{ \mathbf{X} \in \mathbb{R}^{n \times n} \mid \mathbf{X} \succeq \mathbf{0}, \mathbf{X}^T = \mathbf{X} \right\}$  be the set of all positive semi-definite (PSD) matrices in  $\mathbb{R}^{n \times n}$ .

We define as the  $n$ -Spectrahedron or the set of density matrices the set

$$\Delta^{n \times n} = \left\{ \mathbf{X} \in \mathcal{S}_+^{n \times n} \mid \text{tr}(\mathbf{X}) = 1 \right\},$$

i.e. the set of all positive semi-definite  $n \times n$  matrices with unit trace.

### Definition 8

(Convex Conjugate)

Let  $\Psi : \mathcal{D} \rightarrow \mathbb{R}$  be a convex, continuous and differentiable function. The convex conjugate of  $\Psi$  is the function  $\Psi^* : \mathbb{R}^n \rightarrow \mathbb{R}$  defined by:

$$\Psi^*(\mathbf{z}) = \max_{\mathbf{y} \in \mathcal{D}} (\mathbf{y}^T \mathbf{z} - \Psi(\mathbf{y})).$$

Some properties of  $\Psi^*$  are:

$$\nabla \Psi^*(\mathbf{z}) = \arg \max_{\mathbf{y} \in \mathcal{D}} (\mathbf{y}^T \mathbf{z} - \Psi(\mathbf{y})),$$

$$\nabla \Psi^*(\mathbb{R}^n) = \mathcal{D},$$

$$\nabla \Psi \circ \nabla \Psi^*(\mathbf{z}) = \mathbf{z},$$

$$\Psi^{**} = \Psi.$$

Additionally, if  $\Psi$  is  $\frac{1}{L}$ -strongly convex ( $L > 0$ ) with respect to a norm  $\|\cdot\|$  defined on  $\mathcal{D}$  it holds:

$$\|\nabla \Psi^*(\mathbf{z}) - \nabla \Psi^*(\mathbf{y})\| \leq L \|\mathbf{z} - \mathbf{y}\|_*, \quad \forall \mathbf{z}, \mathbf{y} \in \mathbb{R}^n$$

i.e.  $\Psi^*$  is Lipschitz with constant  $L$  [26].

We now provide the following definitions which will be important to define the Mirror Descent optimisation later on.

### Definition 9

(Mirror Map)

Assume that  $\mathcal{D} \subset \mathbb{R}^n$  is a non-empty, open and convex set such that its closure contains the constraint set  $\mathcal{X}$  and their intersection is not the empty set, i.e.  $\mathcal{X} \subseteq \overline{\mathcal{D}}$  and  $\mathcal{X} \cap \mathcal{D} \neq \emptyset$ . We say that a function  $\Psi : \mathcal{D} \rightarrow \mathbb{R}$  is a mirror map under the assumptions:

- $\Psi$  is strictly convex on  $\mathcal{D}$ , i.e.

$$\forall \mathbf{z} \neq \mathbf{y} \in \mathcal{D} \quad \forall \alpha \in (0, 1) \quad \Psi(\alpha \mathbf{z} + (1 - \alpha) \mathbf{y}) < \alpha \Psi(\mathbf{z}) + (1 - \alpha) \Psi(\mathbf{y}).$$

- $\Psi$  is continuous on and differentiable on  $\mathcal{D}$ .
- The dual space is all of  $\mathbb{R}^n$ , i.e.

$$\nabla \Psi(\mathcal{D}) = \left\{ \nabla \Psi(\mathbf{z}) \mid \mathbf{z} \in \mathcal{D} \right\} = \mathbb{R}^n.$$



- The gradient of  $\Psi$ ,  $\nabla\Psi : \mathbb{R}^n \rightarrow \mathbb{R}^n$  diverges on the boundary of  $\mathcal{D}$ , i.e.

$$\lim_{z \rightarrow \partial\mathcal{D}} \|\nabla\Psi(z)\| = +\infty.$$

Note that by equivalence of norms in finite spaces, the definition does not depend on the choice of norm.

### Definition 10

(Bregman divergence)

Let  $\Psi$  be a mirror map on  $\mathcal{D}$ . Then the Bregman divergence associated to  $\Psi$  is the function  $D_\Psi : \mathcal{D} \times \mathcal{D} \rightarrow \mathbb{R}$  given by:

$$\begin{aligned} D_\Psi(\mathbf{x}, \mathbf{y}) &= \Psi(\mathbf{x}) - \Psi(\mathbf{y}) - \langle \nabla\Psi(\mathbf{y}), \mathbf{x} - \mathbf{y} \rangle \\ &= \Psi(\mathbf{x}) - \Psi(\mathbf{y}) - \nabla\Psi(\mathbf{y})^T (\mathbf{x} - \mathbf{y}), \quad (\mathbf{x}, \mathbf{y}) \in \mathcal{D} \times \mathcal{D}. \end{aligned} \quad (4.5)$$

Some properties of the Bregman divergence are:

- 

$$\nabla_{\mathbf{x}} D_\Psi(\mathbf{x}, \mathbf{y}) = \nabla\Psi(\mathbf{x}) - \nabla\Psi(\mathbf{y})$$

- For any  $\mathbf{x}, \mathbf{y} \in \mathcal{D}$  define  $\hat{\mathbf{x}} = \nabla\Psi(\mathbf{x})$  and  $\hat{\mathbf{y}} = \nabla\Psi(\mathbf{y})$ . Then it holds:

$$D_\Psi(\mathbf{x}, \mathbf{y}) = D_{\Psi^*}(\hat{\mathbf{x}}, \hat{\mathbf{y}}),$$

where  $\Psi^*$  denotes the convex conjugate of  $\Psi$ .

Note that we define as the dual space of  $\mathcal{D}$  (under  $\Psi$ ) the space  $\mathcal{D}^* = \nabla\Psi(\mathcal{D})$ . The last equation explains how the the convex conjugate of  $\Psi$  behaves on the primal space  $\mathcal{D}$ . The gradient  $\nabla\Psi$  maps the primal space onto the dual space  $\mathcal{D}^*$  and  $\nabla\Psi^*$  maps the dual space onto the primal.

- For  $\mathbf{x}, \mathbf{y}, \mathbf{z} \in \mathcal{D}$  it holds:

$$D_\Psi(\mathbf{x}, \mathbf{y}) + D_\Psi(\mathbf{z}, \mathbf{x}) - D_\Psi(\mathbf{z}, \mathbf{y}) = (\nabla\Psi(\mathbf{x}) - \nabla\Psi(\mathbf{y}))^T (\mathbf{x} - \mathbf{z}).$$

### Definition 11

(Projection)

The projection of a point  $\mathbf{y} \in \mathcal{D}$ , onto a feasible set  $\mathcal{X} \subseteq \mathcal{D}$  is defined as

$$\Pi_{\mathcal{X}}(\mathbf{y}) = \arg \min_{\mathbf{x} \in \mathcal{X}} \frac{1}{2} \|\mathbf{x} - \mathbf{y}\|_2^2,$$

where  $\Pi_{\mathcal{X}} : \mathcal{D} \rightarrow \mathcal{X}$  is called the projection operator on  $\mathcal{X}$ .

Note that in the case that  $\mathbf{y} \in \mathcal{X}$ :

$$\Pi_{\mathcal{X}}(\mathbf{y}) = \mathbf{y}.$$

**Definition 12***(Bregman projection)**We define as the Bregman projection associated to  $\Psi$  of a point  $\mathbf{y} \in \mathcal{D}$  onto  $\mathcal{X}$  as*

$$\Pi_{\mathcal{X}}^{\Psi}(\mathbf{y}) = \arg \min_{\mathbf{x} \in \mathcal{X}} D_{\Psi}(\mathbf{x}, \mathbf{y}).$$

**Lemma 1***Let  $\Psi$  be a mirror map and  $D_{\Psi}$  its associated Bregman divergence function. Then if we restrict  $\Psi^*$  on  $\mathcal{X}$ , i.e.*

$$\Psi^*(\mathbf{z}) = \max_{\mathbf{y} \in \mathcal{X}} (\mathbf{y}^T \mathbf{z} - \Psi(\mathbf{y})),$$

*it holds:*

$$\nabla \Psi^*(\mathbf{y}) = \Pi_{\mathcal{X}}^{\Psi}(\mathbf{y}) = \arg \min_{\mathbf{x} \in \mathcal{X}} D_{\Psi}(\mathbf{x}, \nabla \Psi^*(\mathbf{y})) \quad (4.6)$$

**Proof 1**

$$\begin{aligned} \nabla \Psi^*(\mathbf{y}) &= \arg \max_{\mathbf{x} \in \mathcal{X}} (\mathbf{x}^T \mathbf{y} - \Psi(\mathbf{x})) \\ &= \arg \min_{\mathbf{x} \in \mathcal{X}} (\Psi(\mathbf{x}) - \mathbf{x}^T \mathbf{y}) \\ &= \arg \min_{\mathbf{x} \in \mathcal{X}} (\Psi(\mathbf{x}) - \Psi(\mathbf{y}) - \mathbf{y}^T \mathbf{x}) \\ &= \arg \min_{\mathbf{x} \in \mathcal{X}} (\Psi(\mathbf{x}) - \Psi(\mathbf{y}) - \nabla \Psi \circ \nabla \Psi^*(\mathbf{y})^T \mathbf{x}) \\ &= \arg \min_{\mathbf{x} \in \mathcal{X}} (\Psi(\mathbf{x}) - \Psi(\mathbf{y}) - \nabla \Psi \circ \nabla \Psi^*(\mathbf{y})^T (\mathbf{x} - \nabla \Psi^*(\mathbf{y}))) \\ &= \arg \min_{\mathbf{x} \in \mathcal{X}} D_{\Psi}(\mathbf{x}, \nabla \Psi^*(\mathbf{y})) \end{aligned}$$

*For the third and fifth equalities we added constant terms with respect to  $\mathbf{x}$ . For the fourth equality we used  $\nabla \Psi \circ \nabla \Psi^*(\mathbf{z}) = \mathbf{z}$ .*

■

### 4.2.2 Main Problem Setup

Throughout this paper we focus on convex and constrained optimisation problems as outlined below:

For  $\mathcal{X}$  a non-empty, compact and convex set (*feasible/ constraint set*), and for  $f : \mathcal{X} \rightarrow \mathbb{R}$  a convex and at least twice differentiable smooth function we are interested in solving the convex optimisation problem

$$f^* = \min_{\mathbf{x} \in \mathcal{X}} f(\mathbf{x}), \quad \mathbf{x}^* = \arg \min_{\mathbf{x} \in \mathcal{X}} f(\mathbf{x}). \quad (4.7)$$

We say that  $\mathbf{x}^*$  is the *optimum/ minimum* of  $f$  on  $\mathcal{X}$  and  $f^* = f(\mathbf{x}^*) = \min_{\mathbf{x} \in \mathcal{X}} f(\mathbf{x})$  the *optimum/ minimum* value of  $f$  on  $\mathcal{X}$ .

### 4.2.3 Examples

Below are presented some examples of mirror maps and their associated Bregman divergence which we will be using later on in this project.

#### Example 1

(Euclidean distance function)

Let  $\mathcal{X} \subseteq \mathbb{R}^n$  be any compact and convex set in  $\mathbb{R}^n$  and  $\mathcal{D} = \mathbb{R}^n$ .  $\mathcal{D}$  is clearly open and convex,  $\mathcal{X} \subseteq \overline{\mathcal{D}}$  and  $\mathcal{X} \cap \mathcal{D} \neq \emptyset$ .

The Euclidean distance function or  $\mathcal{L}2$ -norm is defined as:

$$\Psi(\mathbf{x}) = \frac{1}{2} \|\mathbf{x}\|_2^2 = \frac{1}{2} \sum_{i=1}^n x_i^2, \quad \forall \mathbf{x} \in \mathbb{R}^n. \quad (4.8)$$

The  $\mathcal{L}2$ -norm is continuous and infinitely-differentiable ( $\mathcal{C}^\infty$ ) on  $\mathbb{R}^n$ .  $\Psi$  is strictly convex on  $\mathbb{R}^n$  since for any  $\mathbf{x} \neq \mathbf{y}$ ,  $\alpha \in (0, 1)$ :

$$\begin{aligned} \Psi(\alpha \mathbf{x} + (1 - \alpha) \mathbf{y}) &= \frac{1}{2} \|\alpha \mathbf{x} + (1 - \alpha) \mathbf{y}\|_2^2 \\ &< \frac{1}{2} \alpha \|\mathbf{x}\|_2^2 + \frac{1}{2} (1 - \alpha) \|\mathbf{y}\|_2^2 \\ &= \alpha \Psi(\mathbf{x}) + (1 - \alpha) \Psi(\mathbf{y}) \end{aligned}$$

We have that

$$\nabla \Psi(\mathbf{x}) = \mathbf{x},$$

and

$$\nabla^2 \Psi(\mathbf{x}) = I_n.$$

The  $\mathcal{L}2$ -norm satisfies all assumptions of Definition 9. The Bregman Divergence induced by  $\Psi$  is given by:

$$\begin{aligned} D_{\Psi}(\mathbf{x}, \mathbf{y}) &= \frac{1}{2} \|\mathbf{x}\|_2^2 - \frac{1}{2} \|\mathbf{y}\|_2^2 - \mathbf{y}^T(\mathbf{x} - \mathbf{y}) \\ &= \frac{1}{2} \|\mathbf{x} - \mathbf{y}\|_2^2 \end{aligned} \quad (4.9)$$

i.e. the Euclidean-distance between  $\mathbf{x}$  and  $\mathbf{y}$ .

The convex conjugate of  $\Psi$  on  $\mathcal{D}$  is given by:

$$\Psi^*(\mathbf{z}) = \frac{1}{2} \|\mathbf{z}\|_2^2 = \Psi(\mathbf{z}),$$

and

$$\nabla \Psi^*(\mathbf{z}) = \mathbf{z}.$$

$\Psi$  is 1 strongly convex with respect to the L2-norm over  $\mathcal{D}$  (see 4 of [57]).

### Example 2

(Negative Entropy function)

Let  $\mathcal{X} = \Delta^n$  the  $n$ -Simplex and  $\mathcal{D} = \mathbb{R}_{0+}^n$ .  $\mathcal{X}$  is compact and convex,  $\mathcal{D}$  is open and also convex,  $\mathcal{X} \subseteq \overline{\mathcal{D}}$  and  $\mathcal{X} \cap \mathcal{D} \neq \emptyset$ . The negative entropy function on  $\mathcal{D}$  is defined as:

$$\Psi(\mathbf{x}) = \sum_{i=1}^n x_i \log(x_i) - x_i, \quad \forall \mathbf{x} \in \mathbb{R}_{0+}^n, \quad (4.10)$$

where we define  $0 \cdot \log(0) := 0$ .

It's clear that  $\Psi$  is continuous and it is doubly differentiable on  $\mathcal{D}$  with first and second derivatives:

$$\nabla \Psi(\mathbf{x}) = \log(\mathbf{x}) = \begin{bmatrix} \log(x_1) \\ \log(x_2) \\ \vdots \\ \log(x_n) \end{bmatrix}, \quad \nabla^2 \Psi(\mathbf{x}) = \text{diag}(\mathbf{x})^{-1} = \begin{bmatrix} \frac{1}{x_1} & 0 & \cdots & 0 \\ 0 & \frac{1}{x_2} & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & \frac{1}{x_n} \end{bmatrix}. \quad (4.11)$$

It is easy to see that  $\nabla^2 \Psi(\mathbf{x}) \succ \mathbf{0}$  for all  $\mathbf{x} \in \mathbb{R}_{0+}^n$  and therefore  $\Psi$  is strictly convex on  $\mathcal{D}$ . Also, note that

$$\log(\mathbb{R}_{0+}) = \mathbb{R}$$

and therefore,

$$\nabla \Psi(\mathcal{D}) = \log(\mathcal{D}) = \log(\mathbb{R}_{0+}^n) = \mathbb{R}^n.$$

The gradient of  $\Psi$  diverges on  $\partial \mathcal{D} = \{\mathbf{0} = [0 \ \cdots \ 0]^T\}$ , i.e.

$$\lim_{\mathbf{x} \rightarrow \mathbf{0}} \|\log(\mathbf{x})\|_1 = +\infty$$

since  $\lim_{z \rightarrow 0} |\log(z)| = +\infty$ .

Thus, the negative entropy function can serve as a mirror map as it satisfies all assumptions of Definition 9.

The convex conjugate of  $\Psi$  on  $\mathcal{D}$  is given by:

$$\Psi^*(\mathbf{z}) = \mathbf{1}_n^T \mathbf{exp}(\mathbf{z}) = \sum_{i=1}^n \exp(z_i), \quad \nabla \Psi^*(\mathbf{z}) = \mathbf{exp}(\mathbf{z}), \quad \forall \mathbf{z} \in \mathbb{R}^n.$$

Note that the convex conjugate of  $\Psi$  constrained on  $\mathcal{X}$  becomes:

$$\Psi^*(\mathbf{z}) = \log \left( \sum_{i=1}^n \exp(z_i) \right) + 1, \quad \nabla \Psi^*(\mathbf{z}) = \frac{1}{\sum_{i=1}^n \exp(z_i)} \mathbf{exp}(\mathbf{z}), \quad \forall \mathbf{z} \in \mathbb{R}^n$$

and  $\nabla \Psi^*(\mathbb{R}^n) = \mathcal{X}$ .

**Proof 2**

$$\nabla \Psi^*(\mathbf{z}) = \max_{\mathbf{y} \in \mathcal{X}} (\mathbf{y}^T \mathbf{z} - \Psi(\mathbf{y})) = \max_{\mathbf{y} \in \mathcal{X}} (\mathbf{y}^T \mathbf{z} - \mathbf{y}^T (\log(\mathbf{y}) - \mathbf{1}_n))$$

The Lagrangian of the problem is given by:

$$L(\mathbf{y}, \lambda, \boldsymbol{\mu}) = \mathbf{y}^T \mathbf{z} - \mathbf{y}^T (\log(\mathbf{y}) - \mathbf{1}_n) - \lambda(\mathbf{y}^T \mathbf{1}_n - 1) - \boldsymbol{\mu}^T \mathbf{y}.$$

Setting the gradient equal to  $\mathbf{0}$ :

$$\nabla_{\mathbf{y}_i} L(\mathbf{y}, \lambda, \boldsymbol{\mu}) = z_i - \log(y_i) - \lambda - \mu_i = 0 \implies y_i = \exp(z_i - \mu_i - \lambda).$$

Utilising the KKT conditions we find that  $\mu_i = 0$  and that  $e^\lambda = \sum_{i=1}^n \exp(z_i)$ .

Therefore  $y_i^* = \frac{\exp(z_i)}{\sum_{i=1}^n \exp(z_i)}$  and  $\Psi(\mathbf{y}^*) = \log(\mathbf{exp}(\mathbf{z})^T \mathbf{1}_n) + 1$  and

$$\nabla \Psi^*(\mathbf{z}) = \frac{\mathbf{exp}(\mathbf{z})}{\mathbf{exp}(\mathbf{z})^T \mathbf{1}_n}.$$

■

The Bregman Divergence induced by  $\Psi$  is given by:

$$\begin{aligned} D_\Psi(\mathbf{x}, \mathbf{y}) &= \sum_{i=1}^n (x_i \log x_i - x_i) - \sum_{i=1}^n (y_i \log y_i - y_i) - \log(\mathbf{y})^T (\mathbf{x} - \mathbf{y}) \\ &= \sum_{i=1}^n x_i \log \frac{x_i}{y_i} - \sum_{i=1}^n x_i + \sum_{i=1}^n y_i = D_{KL}(\mathbf{x} \parallel \mathbf{y}), \end{aligned} \tag{4.12}$$

i.e. the generalized Kullback–Leibler divergence between  $\mathbf{x}$  and  $\mathbf{y}$ .

The Negative Entropy function is 1 strongly convex with respect to the L1-norm (Pinsker's Inequality) over  $\mathcal{D}$  [57].

**Example 3**

(Negative Von Neumann Entropy function)

Let  $\mathcal{X} = \Delta^{n \times n}$  the  $n$ -Spectrahedron and  $\mathcal{D} = \mathcal{S}^{n \times n} = \{\mathbf{X} \in \mathbb{R}^{n \times n} | \mathbf{X} \succ \mathbf{0}, \mathbf{X}^T = \mathbf{X}\}$ .  $\mathcal{X}$  is compact and convex,  $\mathcal{D}$  is open and also convex,  $\mathcal{X} \subseteq \overline{\mathcal{D}}$  and  $\mathcal{X} \cap \mathcal{D} \neq \emptyset$ .

**Proof 3** (Convexity of  $\mathcal{X}$  and  $\mathcal{D}$ ).

Let  $\mathbf{A}, \mathbf{B} \in \mathcal{D}$ ,  $\mathbf{Z}, \mathbf{Y} \in \mathcal{X}$ ,  $\alpha \in (0, 1)$ :

$$\begin{aligned} \alpha \mathbf{A} \succ \mathbf{0}, \quad (1 - \alpha) \mathbf{B} \succ \mathbf{0} &\implies \alpha \mathbf{A} + (1 - \alpha) \mathbf{B} \succ \mathbf{0} \\ &\implies \left( \alpha \mathbf{A} + (1 - \alpha) \mathbf{B} \right) \in \mathcal{D} \end{aligned}$$

$$\begin{aligned} \text{tr}(\alpha \mathbf{Z}) = \alpha, \quad \text{tr}((1 - \alpha) \mathbf{Y}) = 1 - \alpha &\implies \text{tr}(\alpha \mathbf{Z} + (1 - \alpha) \mathbf{Y}) = 1 \\ &\implies \left( \alpha \mathbf{Z} + (1 - \alpha) \mathbf{Y} \right) \in \mathcal{X}. \end{aligned}$$

■

The negative Von Neumann Entropy function on  $\mathcal{D}$  is defined as:

$$\Psi(\mathbf{X}) = \text{tr}(\mathbf{X} \log \mathbf{X} - \mathbf{X}) = \sum_{i=1}^n (\lambda_i(\mathbf{X}) \log(\lambda_i(\mathbf{X})) - \lambda_i(\mathbf{X})), \quad \forall \mathbf{X} \in \mathcal{D}, \quad (4.13)$$

where we define  $0 \cdot \log 0 := 0$ .

$\Psi$  is continuous and differentiable on  $\mathcal{D}$ , with gradient:

$$\nabla \Psi(\mathbf{X}) = \log \mathbf{X} = \begin{bmatrix} \log X_{11} & \log X_{12} & \dots & \log X_{1n} \\ \log X_{21} & \log X_{22} & \dots & \log X_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ \log X_{n1} & \log X_{n2} & \dots & \log X_{nn} \end{bmatrix}. \quad (4.14)$$

$\Psi$  is strictly convex on  $\mathcal{D}$ . For the proof of the convexity of  $\Psi$  refer to [2]. The Bregman Divergence induced by  $\Psi$  is given by:

$$\begin{aligned} D_\Psi(\mathbf{X}, \mathbf{Y}) &= \text{tr}(\mathbf{X} \log \mathbf{X} - \mathbf{X}) - \text{tr}(\mathbf{Y} \log \mathbf{Y} - \mathbf{Y}) - \log(\mathbf{Y}) \bullet (\mathbf{X} - \mathbf{Y}) \\ &= \text{tr}(\mathbf{X}(\log \mathbf{X} - \log \mathbf{Y}) - \mathbf{X} + \mathbf{Y}) \end{aligned} \quad (4.15)$$

For the convex conjugate of  $\Psi$  on  $\mathcal{D}$ ,  $\Psi^*(\mathbf{Z}) = \max_{\mathbf{Y} \in \mathcal{D}} (\mathbf{Y} \bullet \mathbf{Z} - \Psi(\mathbf{Y}))$ , we have:

$$\nabla \Psi^*(\mathbf{Z}) = \exp(\mathbf{Z}).$$

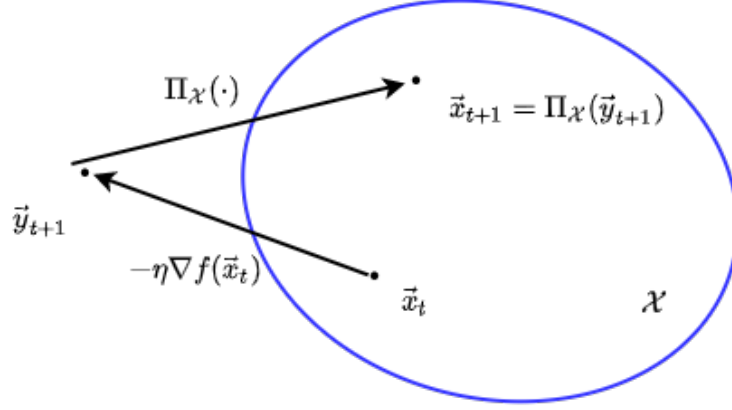
For  $\Psi^*(\mathbf{Z}) = \max_{\mathbf{Y} \in \mathcal{X}} (\mathbf{Y} \bullet \mathbf{Z} - \Psi(\mathbf{Y}))$  we have,

$$\nabla \Psi^*(\mathbf{Z}) = \frac{\exp(\mathbf{Z})}{\text{tr}(\exp(\mathbf{Z}))}.$$

The Negative Von Neumann entropy function can be proved to be  $\frac{1}{2}$  strongly convex with respect to the nuclear norm (Matrix Pinsker's Inequality) over  $\mathcal{D}$  [13, 57].

## 4.3 Mirror Descent

### 4.3.1 The (Projected) Gradient Descent Algorithm



**Figure 4.2:** Illustration of Projected Gradient Descent.

The (Stochastic) Gradient Descent is the simplest and one of the most widely used methods for unconstrained and constrained optimisation of a differentiable function. One simply takes a step in the direction of the negative gradient of  $f$  (steepest descent) and projects it in the feasible set if the optimisation is constrained using the projection operator.

The Projected (Stochastic) Gradient Descent update rule for one iteration consists of two steps and it is illustrated in Figure 4.2 and is given by:

$$\begin{aligned} \mathbf{y}_{t+1} &= \mathbf{x}_t - \eta_t \nabla f(\mathbf{x}_t) + \sigma \mathcal{N}(0, \mathbf{I}), \quad \sigma \geq 0 \\ \mathbf{x}_{t+1} &= \Pi_{\mathcal{X}}(\mathbf{y}_{t+1}), \quad t \geq 0. \end{aligned} \tag{4.16}$$

Note that if we set  $\sigma > 0$  we say that the optimisation is stochastic and deterministic otherwise.

However, for high dimensionality problems, the projection operator can slow down the convergence of projected gradient descent [4] (Chapter 9). Moreover, there are cases in which gradient descent may never converge to the true minimum. Such is the scenario when using fixed (non-adaptive) learning rate  $\eta_t$  in cases with noisy gradient calculations (e.g. Section 5.2 of [10]). For instance, in mini-batch optimisation, the gradient is calculated over a subset of the training samples set and it is usually noisy. In Chapter 5 we provide a relevant experiment. We will also see later on that Projected Gradient Descent is a special case of a “broader” optimisation algorithm called Mirror Descent. Below, we provide the reader with the PGD algorithm averaged over  $N_{iid}$  runs.

**Algorithm 1** Projected Gradient Descent Algorithm

---

**Input:**  $\mathbf{X}_0 = [\mathbf{x}_0^1, \dots, \mathbf{x}_0^{N_{iid}}] : \{\mathbf{x}_0^N\}_{N=1}^{N_{iid}} \in \mathcal{X}, \eta > 0$  learning rate  
 $\sigma \geq 0$  noise parameter,  $T > 0$  iterations

**Output:**  $\mathbf{x}^*, f^*$

**for**  $N = 1, \dots, N_{iid}$  **do**  
 $\mathbf{x}_0 = \mathbf{x}_0^N$   
 $\eta_0 = \eta$   
**for**  $t = 1, \dots, T$  **do**  
 $\mathbf{y}_{t+1} = \mathbf{x}_t - \eta_t \nabla f(\mathbf{x}_t) + \sigma \mathcal{N}(0, \mathbf{I})$   
 $\mathbf{x}_{t+1} = \Pi_{\mathcal{X}}(\mathbf{y}_{t+1}) = \arg \min_{\mathbf{x} \in \mathcal{X}} \frac{1}{2} \|\mathbf{x} - \mathbf{y}_{t+1}\|_2^2$   
**if** learning rate is adaptive **then**  
    **Update**  $\eta_t$   
**end if**  
**end for**  
 $\mathbf{x}_N^* = \frac{1}{T+1} \sum_{s=0}^T \mathbf{x}_s$  or  $\mathbf{x}_N^* = \arg \min_{\mathbf{x} \in \{\mathbf{x}_s\}_0^T} f(\mathbf{x}),$   
 $f_N^* = f(\mathbf{x}_N^*)$   
**end for**  
 $\mathbf{x}^* = \frac{1}{N_{iid}} \sum_{N=1}^{N_{iid}} \mathbf{x}_N^*$   
 $f^* = \frac{1}{N_{iid}} \sum_{N=1}^{N_{iid}} f_N^*$

---

**4.3.2 The Mirror Descent Algorithm**

The motivations of Mirror Descent are:

- Projected Gradient Descent does not discriminate between the primal and dual spaces. While at iteration  $t$   $\mathbf{x}_t$  lies in a subset of  $\mathbb{R}^n$  (the primal space), the gradient  $\nabla f(\mathbf{x}_t)$  (a linear function of  $\mathbf{x}_t$ ) lies in the dual space. Despite the fact that  $\mathbf{x}_t$  and  $\nabla f(\mathbf{x}_t)$  live in different vector spaces, GD at each iteration calculates  $\mathbf{x}_t - \eta_t \nabla f(\mathbf{x}_t)$  as if they both lie in the same vector space. In contrast, as it will be explained in detail in the next paragraph, Mirror Descent discriminates between these two spaces by utilising as a bijection between the two spaces a mirror map function  $\Psi$  and its convex conjugate  $\Psi^*$ .
- In addition, the convergence analysis on projected Gradient Descent utilises  $L$ -Lipschitz objective functions with respect to the  $L_2$ -norm. However, take for example a case where the objective function is 1-Lipschitz w.r.t. the  $L_\infty$ -norm. Then, supposing that the dimensions of the domain of  $f$  is equal to  $N$ , then  $f$  is  $\sqrt{N}$ -Lipschitz w.r.t to the  $L_2$ -norm and thus, for large  $N$  the convergence of the method can be really slow.



For a more detailed discussion of the motivations of Mirror Descent we refer the reader to [13, 24].

In continuous time dynamics, the update schema of Stochastic Mirror Descent is given by the Itô SDE (see [44]):

$$d\mathbf{y}_t = \eta_t \nabla f(\mathbf{x}_t) dt + \sigma d\mathbf{B}_t, \quad \mathbf{x}_t = \Pi_{\mathcal{X}}^{\Psi}(\nabla \Psi^*(\mathbf{y}_t)), \quad \sigma \geq 0, t \geq 0,$$

where  $\sigma$  denotes a noise parameter and  $\mathbf{B}_t$  denotes a Brownian motion. From the above equation we obtain the Euler time-discretisation:

$$\begin{aligned} \mathbf{y}_{t+1} &= \mathbf{y}_t - \epsilon \eta_t \nabla f(\mathbf{x}_t) + \sqrt{\epsilon} \sigma \mathcal{N}(0, \mathbf{I}), \quad \sigma \geq 0 \\ \mathbf{x}_{t+1} &= \Pi_{\mathcal{X}}^{\Psi}(\nabla \Psi^*(\mathbf{y}_{t+1})), \quad t \geq 0, \end{aligned}$$

where  $\epsilon$  is a discretisation parameter which we will set equal to 1 for the rest of this project.

The update schema of one iteration of (Stochastic) Mirror Descent can be described in four steps for an initial point  $\mathbf{x}_0 \in \mathcal{X} \cap \mathcal{D}$ :

$$\hat{\mathbf{x}}_t = \nabla \Psi(\mathbf{x}_t) \tag{4.17a}$$

$$\mathbf{y}_{t+1}^{\wedge} = \hat{\mathbf{x}}_t - \eta_t \nabla f(\mathbf{x}_t) + \sigma \mathcal{N}(0, \mathbf{I}), \quad \sigma \geq 0 \tag{4.17b}$$

$$\mathbf{y}_{t+1} = \nabla \Psi^*(\mathbf{y}_{t+1}^{\wedge}) \tag{4.17c}$$

$$\mathbf{x}_{t+1} = \Pi_{\mathcal{X}}^{\Psi}(\mathbf{y}_{t+1}) = \arg \min_{\mathbf{x} \in \mathcal{X}} D_{\Psi}(\mathbf{x}, \mathbf{y}_{t+1}) \tag{4.17d}$$

where  $\Psi : \mathcal{D} \rightarrow \mathbb{R}$  is a function that meets all assumptions of a mirror map (see Definition 9) and  $\Psi^*(\mathbf{z}) = \max_{\mathbf{y} \in \mathcal{D}} (\mathbf{y}^T \mathbf{z} - \Psi(\mathbf{y}))$  is its convex conjugate as described in Definition 8.

The geometric intuition behind the Mirror Descent algorithm is that in order to perform the optimization in the primal space  $\mathcal{D}$ , we first map the point  $\mathbf{x}_t \in \mathcal{X} \subset \mathcal{D}$  to its dual space  $\mathcal{D}^*$  (4.17a). Then the negative gradient step is performed in the dual space to obtain  $\mathbf{y}_{t+1}^{\wedge}$  (4.17b) which is then mapped back to the primal space via the convex conjugate function of the mirror map (4.17c).

Note that at each update step, the new point in the primal space  $\mathbf{y}_{t+1} \in \mathcal{D}$  might be outside of the feasible set  $\mathcal{X}$ . In that case it should be projected onto  $\mathcal{X}$ . The projection used to project  $\mathbf{y}_{t+1}$  onto  $\mathbf{x}_{t+1} \in \mathcal{X}$  is the Bregman Projection  $D_{\mathcal{X}}^{\Psi}$  associated with the mirror map  $\Psi$  (4.17d).

Note also that from Lemma 1, if we take

$$\Psi^*(\mathbf{z}) = \max_{\mathbf{y} \in \mathcal{X}} (\mathbf{y}^T \mathbf{z} - \Psi(\mathbf{y}))$$

then the last step can be rewritten into  $\mathbf{x}_{t+1} = \mathbf{y}_{t+1}$ .

In that case the algorithm of the Mirror Descent (without using the Bregman Divergence) is given by:

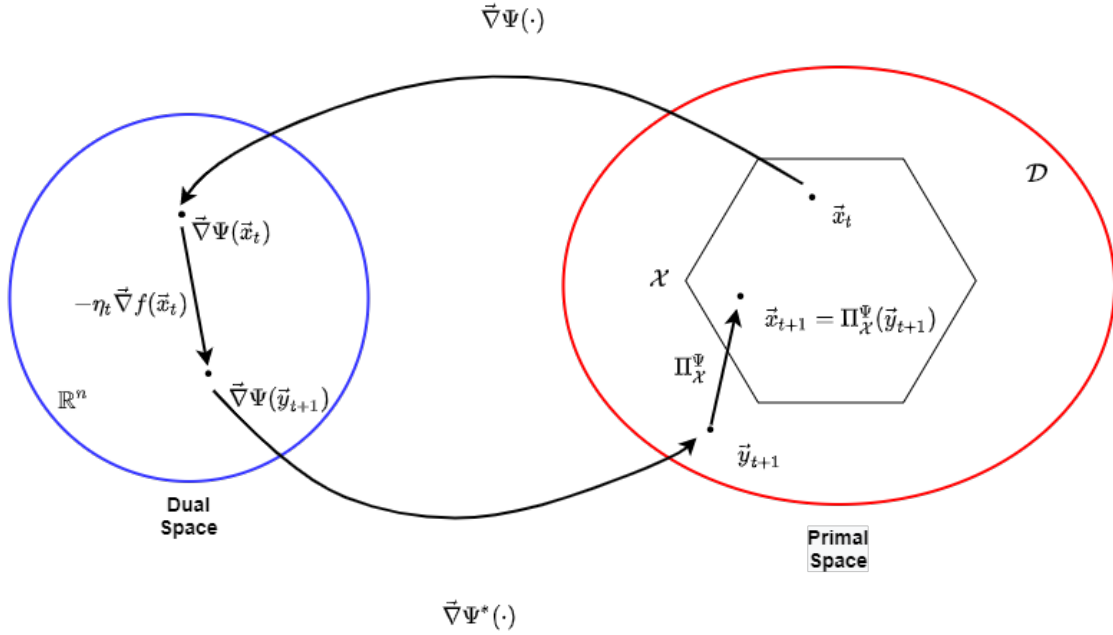
$$\begin{aligned}\mathbf{y}_{t+1} &= \mathbf{y}_t - \eta_t \nabla f(\mathbf{x}_t) + \sigma \mathcal{N}(0, \mathbf{I}), \quad \sigma \geq 0 \\ \mathbf{x}_{t+1} &= \nabla \Psi^*(\mathbf{y}_{t+1}), \quad t \geq 0,\end{aligned}$$

which is taken as the definition in [10]. However, the two definitions are equivalent. An illustration of the above schema is provided below on Figure 4.3 and the step-by-step algorithm of Mirror Descent is provided in Algorithm 2.

The update rule of (Stochastic) Mirror Descent is usually written in two steps which are often taken as the definition of the method (see [13]):

$$\begin{aligned}\nabla \Psi(\mathbf{y}_{t+1}) &= \nabla \Psi(\mathbf{x}_t) - \eta_t \nabla f(\mathbf{x}_t) + \sigma \mathcal{N}(0, \mathbf{I}), \quad \sigma \geq 0 \\ \mathbf{x}_{t+1} &= \Pi_{\mathcal{X}}^{\Psi}(\mathbf{y}_{t+1}) = \arg \min_{\mathbf{x} \in \mathcal{X}} D_{\Psi}(\mathbf{x}, \mathbf{y}_{t+1}).\end{aligned}\tag{4.18}$$

Note that, in the case of Mirror Descent we will be using the above scheme in this project.



**Figure 4.3:** Illustration of Mirror Descent in accordance with Equations 4.17a-4.17d.

We point out that the update rule of  $\mathbf{x}_{t+1}$  in Equation 4.18 can be combined in a single step and be rewritten as:

$$\mathbf{x}_{t+1} = \arg \min_{\mathbf{x} \in \mathcal{X}} \{(\eta_t \nabla f(\mathbf{x}_t) - \sigma \mathcal{N}(0, \mathbf{I}))^T \mathbf{x} + D_\Psi(\mathbf{x}, \mathbf{x}_t)\}. \quad (4.19)$$

**Proof 4** Let  $\sigma = 0$  :

$$\begin{aligned} \mathbf{x}_{t+1} &= \Pi_{\mathcal{X}}^{\Psi}(\mathbf{y}_{t+1}) = \arg \min_{\mathbf{x} \in \mathcal{X}} D_\Psi(\mathbf{x}, \mathbf{y}_{t+1}) \\ &= \arg \min_{\mathbf{x} \in \mathcal{X}} \{ \Psi(\mathbf{x}) - \Psi(\mathbf{y}_{t+1}) - \nabla \Psi(\mathbf{y}_{t+1})^T (\mathbf{x} - \mathbf{y}_{t+1}) \} \\ &= \arg \min_{\mathbf{x} \in \mathcal{X}} \{ \Psi(\mathbf{x}) - \Psi(\mathbf{y}_{t+1}) - \nabla \Psi(\mathbf{y}_{t+1})^T \mathbf{x} + \nabla \Psi(\mathbf{y}_{t+1})^T \mathbf{y}_{t+1} \} \\ &= \arg \min_{\mathbf{x} \in \mathcal{X}} \{ \Psi(\mathbf{x}) - \nabla \Psi(\mathbf{y}_{t+1})^T \mathbf{x} \} \\ &= \arg \min_{\mathbf{x} \in \mathcal{X}} \{ \Psi(\mathbf{x}) - (\nabla \Psi(\mathbf{x}_t) - \eta_t \nabla f(\mathbf{x}_t))^T \mathbf{x} \} \\ &= \arg \min_{\mathbf{x} \in \mathcal{X}} \{ \eta_t \nabla f(\mathbf{x}_t)^T \mathbf{x} + \Psi(\mathbf{x}) - \nabla \Psi(\mathbf{x}_t)^T \mathbf{x} \} \\ &= \arg \min_{\mathbf{x} \in \mathcal{X}} \{ \eta_t \nabla f(\mathbf{x}_t)^T \mathbf{x} + \Psi(\mathbf{x}) + \Psi(\mathbf{x}_t) - \nabla \Psi(\mathbf{x}_t)^T \mathbf{x} + \nabla \Psi(\mathbf{x}_t)^T \mathbf{x}_t \} \\ &= \arg \min_{\mathbf{x} \in \mathcal{X}} \{ \eta_t \nabla f(\mathbf{x}_t)^T \mathbf{x} + D_\Psi(\mathbf{x}, \mathbf{x}_t) \}. \end{aligned}$$

Note that for the first equality we expanded the update rule of Mirror Descent as portrayed by 4.18. For the fifth and eighth equalities we added/ removed constants w.r.t. to the optimisation parameter  $\mathbf{x}$ . ■

As stated above, Projected Gradient Descent is a special case of Mirror Descent. One could use as a mirror map the Euclidean  $L2$ -norm (see Example 1). In that case the update rule becomes:

$$\begin{aligned} \mathbf{x}_{t+1} &= \arg \min_{\mathbf{x} \in \mathcal{X}} \{(\eta_t \nabla f(\mathbf{x}_t) - \sigma \mathcal{N}(0, \mathbf{I}))^T \mathbf{x} + D_\Psi(\mathbf{x}, \mathbf{x}_t)\} \\ &= \arg \min_{\mathbf{x} \in \mathcal{X}} \{(\eta_t \nabla f(\mathbf{x}_t) - \sigma \mathcal{N}(0, \mathbf{I}))^T \mathbf{x} + \frac{1}{2} \|\mathbf{x} - \mathbf{x}_t\|_2^2\} \\ &= \arg \min_{\mathbf{x} \in \mathcal{X}} \frac{1}{2} \left\| \mathbf{x} - (\mathbf{x}_t - \eta_t \nabla f(\mathbf{x}_t) + \sigma \mathcal{N}(0, \mathbf{I})) \right\|_2^2 \\ &= \Pi_{\mathcal{X}}(\mathbf{y}_{t+1}) \end{aligned}$$

which is exactly the Projected Gradient Descent update schema.

---

**Algorithm 2** (Stochastic) Mirror Descent Algorithm

---

**Input:**  $\mathbf{X}_0 = [\mathbf{x}_0^1, \dots, \mathbf{x}_0^{N_{iid}}] : \{\mathbf{x}_0^N\}_{N=1}^{N_{iid}} \in \mathcal{X}, \eta > 0$  learning rate  
 $\sigma \geq 0$  noise parameter,  $T > 0$  iterations

**Output:**  $\mathbf{x}^*, f^*$

**for**  $N = 1, \dots, N_{iid}$  **do**

$\mathbf{x}_0 = \mathbf{x}_0^N$

$\eta_0 = \eta$

**for**  $t = 1, \dots, T$  **do**

$$\nabla \Psi(\mathbf{y}_{t+1}) = \nabla \Psi(\mathbf{x}_t) - \eta_t \nabla f(\mathbf{x}_t) + \sigma \mathcal{N}(0, \mathbf{I})$$

$$\mathbf{x}_{t+1} = \Pi_{\mathcal{X}}^{\Psi}(\mathbf{y}_{t+1}) = \arg \min_{\mathbf{x} \in \mathcal{X}} D_{\Psi}(\mathbf{x}, \mathbf{y}_{t+1}).$$

**if** learning rate is adaptive **then**

**Update**  $\eta_t$

**end if**

**end for**

$$\mathbf{x}_N^* = \frac{1}{T+1} \sum_{s=0}^T \mathbf{x}_s \text{ or } \mathbf{x}_N^* = \arg \min_{\mathbf{x} \in \{\mathbf{x}_s\}_0^T} f(\mathbf{x})$$

$$f_N^* = f(\mathbf{x}_N^*)$$

**end for**

$$\mathbf{x}^* = \frac{1}{N_{iid}} \sum_{N=1}^{N_{iid}} \mathbf{x}_N^*$$

$$f^* = \frac{1}{N_{iid}} \sum_{N=1}^{N_{iid}} f_N^*$$


---

As aforementioned, our study here is rather experimental rather than theoretical. For more elaborate theoretical results on the convergence of Gradient Descent and Mirror Descent please refer to Sections 3.2 and 4 of [13], respectively.

Note that if we set  $\sigma > 0$  we say that the optimisation is stochastic (Stochastic Mirror Descent) and deterministic (Mirror Descent) otherwise.

## 4.4 Interacting Mirror Descent

Optimisation algorithms such as Projected Gradient Descent may converge to non optimum solutions, such as local minima, saddle points or to neighbourhoods of the real optimum. This is due various factors. For instance, noisy data and or computing the gradient over a subset of the sample data can cause this. To combat this,  $N_{iid} \geq 1$  independent and identical versions of the same algorithm are ran in parallel (see [62, 27]) and the final result is the average of the result of all runs (see Algorithms 1 and 2), in other words:

$$f^* = \frac{1}{N_{iid}} \sum_{N=1}^{N_{iid}} f_N^*.$$

A variation of the (Stochastic) Mirror Descent Algorithm is the (Stochastic) Interacting Mirror Descent Algorithm proposed in [10] where they suggest instead of running  $N_{iid}$  replicas of the Mirror Descent algorithm, running one replica of the MD Algorithm in which they allow  $N_p \geq 1$  particles to “interact”. Assuming that  $\Psi^*(\mathbf{z}) = \max_{\mathbf{y} \in \mathcal{X}} (\mathbf{y}^T \mathbf{z} - \Psi(\mathbf{y}))$ , the optimisation update rule they propose is inspired by the following Itô Stochastic Differential Equation in continuous time for :

$$\begin{aligned} d\mathbf{y}_t^i &= -\nabla f(\mathbf{x}_t^i)dt + \sum_{j=1}^{N_p} A_{ij}(\mathbf{y}_t^j - \mathbf{y}_t^i) + \sigma d\mathbf{B}_t^i, \\ \mathbf{x}_t^i &= \nabla \Psi^*(\mathbf{y}_t^i), \quad i = 1, \dots, N_p, \quad t \geq 0, \end{aligned}$$

where each particle is characterised by independent Brownian motions  $\mathbf{B}_t^i$  and  $\mathbf{x}_t^i$  denotes the  $i$ -th particle at time-step  $t$ , where  $i = 1, \dots, N_p$ . The method is described in detail in the rest of this section.

We will be concerned with the discretised version of the Stochastic Interacting Mirror Descent Algorithm:

$$\mathbf{y}_{t+1}^i = \mathbf{y}_t^i - \eta_t \cdot \epsilon \nabla f(\mathbf{x}_t^i) + \epsilon \sum_{j=1}^{N_p} A_{ij}(\mathbf{y}_t^j - \mathbf{y}_t^i) + \sigma \sqrt{\epsilon} \mathcal{N}(0, 1) \quad (4.20a)$$

$$\mathbf{x}_{t+1}^i = \nabla \Psi^*(\mathbf{y}_{t+1}^i), \quad i = 1, \dots, N_p, \quad t \geq 0. \quad (4.20b)$$

Note that, if we use  $\Psi^*(\mathbf{z}) = \max_{\mathbf{y} \in \mathcal{D}} (\mathbf{y}^T \mathbf{z} - \Psi(\mathbf{y}))$ , then the last step becomes

$$\mathbf{x}_{t+1}^i = \Pi_{\mathcal{X}}^{\Psi} \left( \nabla \Psi^*(\mathbf{y}_{t+1}^i) \right), \quad i = 1, \dots, N_p, \quad t \geq 0.$$

We remark that  $\mathbf{A}$  represents an  $N_p \times N_p$  doubly stochastic ( $\sum_{i=1}^{N_p} A_{ij} = \sum_{j=1}^{N_p} A_{ij} = 1$ ) matrix constituting of the interaction weights  $A_{ij}$  between the particles  $i$  and  $j$  and  $\epsilon > 0$  is called the discretisation parameter. We remark that for a particle  $p$  only  $j$  particles  $1 \leq j \leq N_p$  have impact on such that  $A_{pj} \neq 0$ . In case that  $\mathbf{A} = \mathbf{0}$  (i.e.

there is no interaction), then the problem is equivalent to the MD framework averaged over  $N_p$  iid copies.

We also remark that if  $\sigma > 0$  the optimisation is called Stochastic Interacting Mirror Descent, and Interacting Mirror Descent otherwise.

Let  $\epsilon = 1$ . The update rule of  $\mathbf{x}_{t+1}^i$  in the previous update schema is equivalent to

$$\mathbf{x}_{t+1}^i = \arg \min_{\mathbf{x} \in \mathcal{X}} \left\{ (\eta_t \nabla f(\mathbf{x}_t^i) - \sigma \mathcal{N}(0, \mathbf{I}))^T \mathbf{x} + \sum_{j=1}^{N_p} A_{ij} D_\Psi(\mathbf{x}, \mathbf{x}_t^j) \right\} \quad (4.21)$$

**Proof 5** Let  $i \in \{1, \dots, N_p\}$ . For  $\sigma = 0$  we have that

$$\mathbf{y}_{t+1}^i = -\eta_t \nabla f(\mathbf{x}_t^i) + \sum_{j=1}^{N_p} A_{ij} \mathbf{y}_t^j.$$

Then,

$$\begin{aligned} \mathbf{x}_{t+1}^i &= \nabla \Psi^*(\mathbf{y}_{t+1}^i) \\ &= \arg \min_{\mathbf{x} \in \mathcal{X}} \left\{ D_\Psi(\mathbf{x}, \nabla \Psi^*(\mathbf{y}_{t+1}^i)) \right\} \\ &= \arg \min_{\mathbf{x} \in \mathcal{X}} \left\{ \Psi(\mathbf{x}) - \Psi \circ \nabla \Psi^*(\mathbf{y}_{t+1}^i) - (\mathbf{y}_{t+1}^i)^T (\mathbf{x} - \nabla \Psi^*(\mathbf{y}_{t+1}^i)) \right\} \\ &= \arg \min_{\mathbf{x} \in \mathcal{X}} \left\{ \Psi(\mathbf{x}) - (\mathbf{y}_{t+1}^i)^T \mathbf{x} \right\} \\ &= \arg \min_{\mathbf{x} \in \mathcal{X}} \left\{ \Psi(\mathbf{x}) - \Psi(\mathbf{x}_{t+1}^i) + \eta_t \nabla f(\mathbf{x}_t^i)^T \mathbf{x} - \sum_{j=1}^{N_p} A_{ij} (\mathbf{y}_t^j)^T \mathbf{x} \right\} \\ &= \arg \min_{\mathbf{x} \in \mathcal{X}} \left\{ \eta_t \nabla f(\mathbf{x}_t^i)^T \mathbf{x} + \sum_{j=1}^{N_p} A_{ij} (\Psi(\mathbf{x}) - \Psi(\mathbf{x}_{t+1}^i) - (\mathbf{y}_t^j)^T \mathbf{x}) \right\} \\ &= \arg \min_{\mathbf{x} \in \mathcal{X}} \left\{ \eta_t \nabla f(\mathbf{x}_t^j)^T \mathbf{x} + \sum_{j=1}^{N_p} A_{ij} (\Psi(\mathbf{x}) - \Psi(\mathbf{x}_t^j) - \nabla \Psi(\mathbf{x}_t^j)^T (\mathbf{x} - \mathbf{x}_t^j)) \right\} \\ &= \arg \min_{\mathbf{x} \in \mathcal{X}} \left\{ \eta_t \nabla f(\mathbf{x}_t^i)^T \mathbf{x} + \sum_{j=1}^{N_p} A_{ij} D_\Psi(\mathbf{x}, \mathbf{x}_t^j) \right\} \end{aligned}$$

For the second equality we utilise Lemma 1. For the third equality we use  $\nabla \Psi \circ \nabla \Psi^*(\mathbf{z}) = \mathbf{z}$  and for the fourth and fifth equality we add or remove constants since the optimisation is w.r.t. to  $\mathbf{x}$ . Additionally, in the fifth equality we expanded  $\mathbf{y}_{t+1}^i$ .

■

**Algorithm 3** (Stochastic) Interacting Mirror Descent Algorithm

---

**Input:**  $\mathbf{X}_0 = [\mathbf{x}_0^1, \dots, \mathbf{x}_0^{N_p}]$ ,  $\eta > 0$  learning rate,  $\sigma \geq 0$  noise parameter,  
 $\mathbf{A}$  an  $N_p \times N_p$  doubly stochastic matrix,  $T > 0$  iterations  
**Output:**  $\mathbf{X}^*$ ,  $f^*$

$\eta_0 = \eta$   
**for**  $t = 1, \dots, T$  **do**  
  **for**  $i = 1, \dots, N_p$  **do**  
     $\mathbf{x}_{t+1}^i = \arg \min_{\mathbf{x} \in \mathcal{X}} \{(\eta_t \cdot \nabla f(\mathbf{x}_t^i) - \sigma \mathcal{N}(0, \mathbf{I}))^T \mathbf{x} + \sum_{j=1}^{N_p} \mathbf{A}_{ij} D_\Psi(\mathbf{x}, \mathbf{x}_{t+1}^j)\}$   
    **if** learning rate is adaptive **then**  
      **Update**  $\eta_t$   
    **else**  
       $\eta_t = \eta$   
    **end if**  
  **end for**  
**end for**  
 $\mathbf{X}^* = \mathbf{X}_T = [\mathbf{x}_T^1, \dots, \mathbf{x}_T^{N_p}]$   
 $f^* = \frac{1}{N_p} \sum_{i=1}^{N_p} f(\mathbf{x}_T^i)$

---

## 4.5 Applications

In Chapter 5 we present results for the optimisation of six convex problems, out of which the five are constrained. We are concerned with two specific constrained optimisation frameworks: one for optimising a convex function

$$f : \mathcal{X} \subseteq \mathbb{R}^n \longrightarrow \mathbb{R},$$

where  $\mathcal{X}$  is the  $n$ -Simplex and one for optimising a convex function

$$f : \mathcal{X} \subseteq \mathbb{R}^{n \times n} \longrightarrow \mathbb{R},$$

where  $\mathcal{X}$  is the  $n$ -Spectrahedron. In the rest of this Chapter we provide the update rules for the Mirror Descent and Interacting Mirror Descent algorithms for these two cases.

### 4.5.1 Optimisation on the Probability Simplex

First, we consider:

$$\min_{\mathbf{x} \in \Delta^n} f(\mathbf{x})$$

and let  $\Psi$  be such that  $\Psi(\mathbf{x}) = \sum_{i=1}^n x_i \log(x_i) - x_i$  (i.e. the negative entropy function).

As described in Example 2,  $\Psi$  attains all assumptions to be a mirror map.

**Lemma 2**

The Mirror Descent update rule for  $\sigma = 0$  at time-step  $t + 1$  takes the form:

$$\begin{aligned} y_i^{t+1} &= x_i^t \exp \left( -\eta_t (\nabla f(\mathbf{x}_t))_i \right) \\ x_i^{t+1} &= \frac{y_i^{t+1}}{\sum_{j=1}^n y_j^{t+1}}, \quad i = 1, \dots, n, \end{aligned} \quad (4.22)$$

where  $i$  represents the  $i$ -th component of each vector.

**Proof 6** We remind that the associated with  $\Psi$  Bregman divergence is given by  $D_\Psi(\mathbf{x}, \mathbf{y}) = \sum_{i=1}^n x_i \log \frac{x_i}{y_i} - \sum_{i=1}^n x_i + \sum_{i=1}^n y_i$ . From Equation 4 for the update rule at time-step  $t + 1$  we have:

$$\begin{aligned} \mathbf{x}_{t+1} &= \arg \min_{\mathbf{x} \in \mathcal{X}} \{ \eta_t \cdot \nabla f(\mathbf{x}_t)^T \mathbf{x} + D_\Psi(\mathbf{x}, \mathbf{x}_t) \} \\ &= \arg \min_{\mathbf{x} \in \mathcal{X}} \left\{ \eta_t \cdot \nabla f(\mathbf{x}_t)^T \mathbf{x} + \sum_{i=1}^n x_i \log \frac{x_i}{x_i^t} - \sum_{i=1}^n x_i + \sum_{i=1}^n x_i^t \right\} \\ &= \arg \min_{\mathbf{x} \in \mathcal{X}} \left\{ \eta_t \cdot \nabla f(\mathbf{x}_t)^T \mathbf{x} + \sum_{i=1}^n x_i \log \frac{x_i}{x_i^t} - \sum_{i=1}^n x_i \right\} \end{aligned}$$

. The Lagrangian of this optimisation problem is given by:

$$L(\mathbf{x}, \lambda, \boldsymbol{\mu}) = \eta_t \cdot \nabla f(\mathbf{x}_t)^T \mathbf{x} + \sum_{i=1}^n x_i \log \frac{x_i}{x_i^t} - \sum_{i=1}^n x_i - \lambda \cdot \left( \sum_{i=1}^n x_i - 1 \right) + \boldsymbol{\mu}^T \mathbf{x},$$

where due to KKT conditions it holds:

$$\boldsymbol{\mu} \geq 0 \quad (4.23a)$$

$$\mu_i \cdot x_i = 0, \quad i = 1, \dots, n \quad (4.23b)$$

$$\sum_{i=1}^n x_i = 1 \quad (4.23c)$$

$$\nabla_{\mathbf{x}} L(\mathbf{x}, \lambda, \boldsymbol{\mu}) = \mathbf{0}. \quad (4.23d)$$

We take the derivative of the Lagrangian with respect to the  $i$ -th component of  $\mathbf{x}$  and set it equal to 0:

$$\begin{aligned} \nabla_{x_i} L(\mathbf{x}, \lambda, \boldsymbol{\mu}) &= \eta_t \cdot (\nabla f(\mathbf{x}_t))_i + \log \frac{x_i}{x_i^t} - \lambda + \mu_i = 0 \\ \implies \log(x_i) &= \log(x_i^t) - \eta_t \cdot (\nabla f(\mathbf{x}_t))_i + \lambda - \mu_i \\ \implies x_i &= e^\lambda \exp \left( \log(x_i^t) - \eta_t \cdot (\nabla f(\mathbf{x}_t))_i - \mu_i \right) \\ \implies x_i > 0 &\implies \mu_i = 0 \end{aligned}$$



Utilising 4.23c we get that:

$$e^\lambda = \frac{1}{\sum_{i=1}^n x_i^t \exp(-\eta_t \cdot (\nabla f(\mathbf{x}_t))_i)}$$

and therefore:

$$x_i^{t+1} = \frac{x_i^t \exp(-\eta_t \cdot (\nabla f(\mathbf{x}_t))_i)}{\sum_{j=1}^n x_j^t \exp(-\eta_t \cdot (\nabla f(\mathbf{x}_t))_j)}. \quad (4.24)$$

■

### Lemma 3

The Interacting Mirror Descent update rule for  $\sigma = 0$  for the  $i$ -th particle at time-step  $t + 1$  takes the form:

$$(x_{t+1}^i)_m = \frac{B_m^i \exp\left(-\eta_t \cdot (\nabla f(\mathbf{x}_t^i))_m\right)}{\sum_{k=1}^n B_k^i \exp\left(-\eta_t \cdot (\nabla f(\mathbf{x}_t^i))_k\right)}, \quad m = 1, \dots, n, \quad (4.25)$$

where  $B_m^i = \prod_{j=1}^{N_p} ((x_t^j)_m)^{A_{ij}}$  and  $m$  represents the  $m$ -th component of each vector.

**Proof 7** Let  $\eta_t = 1$ :

$$\begin{aligned} \mathbf{x}_{t+1}^i &= \arg \min_{\mathbf{x} \in \mathcal{X}} \left\{ \nabla f(\mathbf{x}_t^i)^T \mathbf{x} + \sum_{j=1}^{N_p} A_{ij} D_\Psi(\mathbf{x}, \mathbf{x}_t^j) \right\} \\ &= \arg \min_{\mathbf{x} \in \mathcal{X}} \left\{ \nabla f(\mathbf{x}_t^i)^T \mathbf{x} + \sum_{j=1}^{N_p} A_{ij} \left( \sum_{m=1}^n x_m \log \left( \frac{x_m}{(x_t^j)_m} \right) - \sum_{m=1}^n x_m \right) \right\} \end{aligned} \quad (4.26)$$

The Lagrangian of this optimisation problem is given by:

$$L(\mathbf{x}, \lambda, \boldsymbol{\mu}) = \nabla f(\mathbf{x}_t^i)^T \mathbf{x} + \sum_{j=1}^{N_p} A_{ij} \left( \sum_{m=1}^n x_m \log \left( \frac{x_m}{(x_t^j)_m} \right) - \sum_{m=1}^n x_m \right) - \lambda \cdot \left( \sum_{m=1}^n x_m - 1 \right) + \boldsymbol{\mu}^T \mathbf{x},$$

where due to KKT conditions it holds:

$$\boldsymbol{\mu} \geq 0 \quad (4.27a)$$

$$\mu_m \cdot x_m = 0, \quad m = 1, \dots, n \quad (4.27b)$$

$$\sum_{m=1}^n x_m = 1 \quad (4.27c)$$

$$\nabla_{\mathbf{x}} L(\mathbf{x}, \lambda, \boldsymbol{\mu}) = \mathbf{0}. \quad (4.27d)$$

We take the derivative of the Lagrangian with respect to the  $m$ -th,  $m = 1, \dots, n$  component of  $\mathbf{x}$  and set it equal to 0:

$$\begin{aligned} \nabla_{x_m} L(\mathbf{x}, \lambda, \boldsymbol{\mu}) &= (\nabla f(\mathbf{x}_t^i))_m + \sum_{j=1}^{N_p} \left( \log \frac{x_m}{(x_t^j)_m} \right) - \lambda - \mu_m \\ &= (\nabla f(\mathbf{x}_t^i))_m + \log(x_m) + \sum_{j=1}^{N_p} A_{ij} \left( -\log(x_t^j)_m \right) - \lambda + \mu_m \\ &= (\nabla f(\mathbf{x}_t^i))_m + \log(x_m) - \sum_{j=1}^{N_p} \left( \log((x_t^j)_m)^{A_{ij}} \right) - \lambda + \mu_m \\ &= (\nabla f(\mathbf{x}_t^i))_m + \log(x_m) - \log \left( \prod_{j=1}^{N_p} ((x_t^j)_m)^{A_{ij}} \right) - \lambda + \mu_m = 0 \end{aligned}$$

Similarly to the MD case, we derive

$$x_m = \prod_{j=1}^{N_p} ((x_t^j)_m)^{A_{ij}} \exp \left\{ -(\nabla f(\mathbf{x}_t^i))_m - \mu_m + \lambda \right\}$$

and therefore  $\mu_m = 0$ . Let  $B_m^i = \prod_{j=1}^{N_p} ((x_t^j)_m)^{A_{ij}}$ , utilising the third KKT condition:

$$e^\lambda = \frac{1}{\sum_{m=1}^n B_m^i \exp \left( -(\nabla f(\mathbf{x}_t^i))_m \right)}.$$

Hence, for  $\eta_t = 1$ :

$$(x_{t+1}^i)_m = x_m = \frac{B_m^i \exp \left( -(\nabla f(\mathbf{x}_t^i))_m \right)}{\sum_{k=1}^n B_k^i \exp \left( -(\nabla f(\mathbf{x}_t^i))_k \right)}, \quad m = 1, \dots, n.$$

In the general case:

$$(x_{t+1}^i)_m = \frac{\prod_{j=1}^{N_p} ((x_t^j)_m)^{A_{ij}} \exp \left( -\eta_t \cdot (\nabla f(\mathbf{x}_t^i))_m \right)}{\sum_{k=1}^n \left( \prod_{j=1}^{N_p} ((x_t^j)_k)^{A_{ij}} \exp \left( -\eta_t \cdot (\nabla f(\mathbf{x}_t^i))_k \right) \right)}, \quad m = 1, \dots, n. \quad (4.28)$$

■

### 4.5.2 Optimisation on the Spectrahedron

Now consider:

$$\min_{\mathbf{X} \in \Delta^{n \times n}} f(\mathbf{X})$$

and let  $\Psi$  be such that  $\Psi(\mathbf{X}) = \text{tr}(\mathbf{X} \log(\mathbf{X}))$  (i.e. the negative von Neumann entropy function). As described in Example 3  $\Psi$  attains all assumptions to be a mirror map.

**Lemma 4**

The mirror descent update rule for  $\mathbf{X}_0 \in S^{n \times n}$  for  $\sigma = 0$  at time-step  $t+1$  takes the form:

$$\begin{aligned} \mathbf{Y}_{t+1} &= \mathbf{X}_t \circ \exp(-\eta_t \nabla f(\mathbf{X}_t)) \\ \mathbf{X}_{t+1} &= \frac{1}{\text{tr}(\mathbf{Y}_{t+1})} \mathbf{Y}_{t+1}, \end{aligned} \tag{4.29}$$

where the operator  $\circ$  represents the element-wise product, i.e.

$$(\mathbf{A} \circ \mathbf{B})_{ij} = A_{ij} \cdot B_{ij}.$$

**Lemma 5**

The Interacting Mirror Descent update rule of  $\mathbf{X}_i^{t+1}$  for  $\mathbf{X}_i^0 \in S^{n \times n}$ ,  $i = 1, \dots, N_p$  for  $\sigma = 0$  at time-step  $t+1$  takes the form:

$$\begin{aligned} \mathbf{Y}_{t+1}^i &= \prod_{j=1}^{N_p} \circ \left( (\mathbf{X}_t^j)^{\circ A_{ij}} \right) \circ \exp(-\eta_t \nabla f(\mathbf{X}_t^i)) \\ \mathbf{X}_{t+1}^i &= \frac{1}{\text{tr}(\mathbf{Y}_{t+1}^i)} \mathbf{Y}_{t+1}^i, \end{aligned} \tag{4.30}$$

where the operator  $\prod_{j=1}^{N_p} \circ$  denotes the elementwise product of matrices and the power  $\circ A_{ij}$  is an elementwise exponentiation.

# Chapter 5

## Experimental Results

In this section we present experimental results of six problems which fit our framework:

- five (convex) constrained optimisation problems (four on the probability simplex and one on the Spectrahedron) (Sections 5.1 - 5.5), and
- one convex unconstrained problem (Section 5.6).

Each problem is explained and formulated in detail below. We solve these problems using three optimisation methods; the standard method, (Stochastic) Gradient Descent, (Stochastic) Mirror Descent and (Stochastic) Interacting Mirror Descent. We discuss and compare our numerical results for the three methods. In case of the latter the interaction matrix  $A$  is assumed to be a mean-field interaction matrix, i.e.  $A_{ij} = \frac{1}{N_p}$  where  $N_p$  denotes the number of particles used in each case. Note that it holds that  $A$  is doubly stochastic. In the case of GD/SGD or MD/SMD the number of independent and identical runs  $N_{iid}$  is stated for each problem.

### 5.1 The Traffic Assignment problem

The first application we consider is the Traffic Assignment problem, which is concerned with selecting the best path in a traffic system between a given starting point and a destination. The best path is defined as the path which incurs the minimum cost of transporting loads from the origin to the destination. Below, we define the problem more formally. For a more detailed discussion of the formulation of the problem refer to [47].

#### 5.1.1 Formulation of the traffic assignment Problem

Let  $\mathcal{G} = (\mathcal{V}, \mathcal{E})$  denote a directed, weighted graph, where  $\mathcal{V}$  denotes a set of  $n$  vertices (nodes) and  $\mathcal{E}$  the set of edges connecting the vertices. Assuming that  $\mathcal{V}$  contains  $n$  vertices ( $|\mathcal{V}| = n$ ), then  $|\mathcal{E}| \leq n \times n$ . We denote with  $\mathcal{W} \in \mathbb{R}^{n \times n}$  the weights-matrix, where  $w_{ij} \neq 0$  if and only if nodes  $n_i, n_j \in \mathcal{V}$  are connected with an edge, i.e.  $n_i, n_j \in \mathcal{V}$  are not connected if and only if  $w_{ij} = 0$ . We further denote with  $\mathcal{O} \in \mathcal{V}$  the starting point and with  $\mathcal{T} \in \mathcal{V}$  the destination. Let also  $\mathcal{R}$  be the set of simple

routes (paths) connecting  $\mathcal{O}$  and  $\mathcal{T}$  within the graph and let  $|\mathcal{R}| = n_{\mathcal{R}}$ . We are concerned with sending  $\mu$  units of traffic from  $\mathcal{O}$  to  $\mathcal{T}$  via  $r \in \mathcal{R}$ . Note that for the experiments performed in this project we set  $\mu = 1$  for simplicity. Let  $\mathcal{X}$  be the feasible set defined as:

$$\mathcal{X} = \left\{ \mathbf{x} \in \mathbb{R}^{n_{\mathcal{R}}} \mid x_r \geq 0, \forall r = 1, \dots, n_{\mathcal{R}}, \sum_{r=1}^{n_{\mathcal{R}}} x_r = \mu = 1 \right\} = \Delta^{n_{\mathcal{R}}} \subset \mathbb{R}^{n_{\mathcal{R}}}.$$

A vector  $[x_1, \dots, x_{n_{\mathcal{R}}}]^T = \mathbf{x} \in \mathcal{X}$ ,  $\mathbf{x}$  is called a routing flow. For a routing flow  $\mathbf{x} \in \mathcal{X}$ , the load on an edge  $\epsilon \in \mathcal{E}$  is given by  $l_{\epsilon}(\mathbf{x}) = \sum_{r \subset \mathcal{R} : \epsilon \in r} x_r$ . The cost (delay) occurring on an edge  $\epsilon$  is denoted as  $c_{\epsilon}(\mathbf{x})$  where  $c_{\epsilon} : \mathbb{R} \rightarrow \mathbb{R}_+$  is a non decreasing convex function in  $\mathbb{R}$ .

For the purposes of this project we assume that

$$c_{\epsilon}(\mathbf{x}) = \frac{1}{2} w_{\epsilon} (1 + l_{\epsilon}(\mathbf{x}))^2, \quad \forall \epsilon \in r \subset \mathcal{R},$$

where  $w_{\epsilon}$  denotes the weight of the two vertices of edge  $\epsilon$ . Note that  $c_{\epsilon}$  is convex on  $\mathbb{R}^{n_{\mathcal{R}}}$ :

**Proof 8** Let  $\mathbf{x}, \mathbf{y} \in \mathbb{R}^{n_{\mathcal{R}}}$ ,  $\alpha \in (0, 1)$ :

$$\begin{aligned} c_{\epsilon}(\alpha \mathbf{x} + (1 - \alpha) \mathbf{y}) &= \frac{1}{2} w_{\epsilon} (1 + l_{\epsilon}(\alpha \mathbf{x} + (1 - \alpha) \mathbf{y}))^2 \\ &= \frac{1}{2} w_{\epsilon} (1 + \alpha l_{\epsilon}(\mathbf{x}) + (1 - \alpha) l_{\epsilon}(\mathbf{y}))^2 \\ &= \frac{1}{2} w_{\epsilon} (1 + \alpha l_{\epsilon}(\mathbf{x}) + (1 - \alpha) l_{\epsilon}(\mathbf{y}) + \alpha - \alpha)^2 \\ &= \frac{1}{2} w_{\epsilon} (\alpha(1 + l_{\epsilon}(\mathbf{x})) + (1 - \alpha)(1 + l_{\epsilon}(\mathbf{y})))^2 \\ &\leq \frac{1}{2} w_{\epsilon} \alpha^2 (1 + l_{\epsilon}(\mathbf{x}))^2 + \frac{1}{2} w_{\epsilon} (1 - \alpha)^2 (1 + l_{\epsilon}(\mathbf{y}))^2 \\ &\leq \alpha c_{\epsilon}(\mathbf{x}) + (1 - \alpha) c_{\epsilon}(\mathbf{y}). \end{aligned}$$

■

The cost through route  $r \in \mathcal{R}$  is defined as the sum of costs over all edges on the path, i.e.:

$$c_r(\mathbf{x}) = \sum_{\epsilon \in r} c_{\epsilon}(\mathbf{x})$$

and the mean cost of the network is obtained by

$$f(\mathbf{x}) = \sum_{r \in \mathcal{R}} x_r c_r(\mathbf{x}) = \sum_{r \in \mathcal{R}} \sum_{\epsilon \in \mathcal{E}} x_r c_{\epsilon}(\mathbf{x}) = \sum_{\epsilon \in \mathcal{E}} l_{\epsilon}(\mathbf{x}) c_{\epsilon}(\mathbf{x}).$$

For our setting the loss function is given by:

$$f(\mathbf{x}) = \frac{1}{2} \sum_{\epsilon \in \mathcal{E}} l_{\epsilon}(\mathbf{x}) w_{\epsilon} (1 + l_{\epsilon}(\mathbf{x}))^2, \quad \mathbf{x} \in \Delta^{n_{\mathcal{R}}}. \quad (5.1)$$

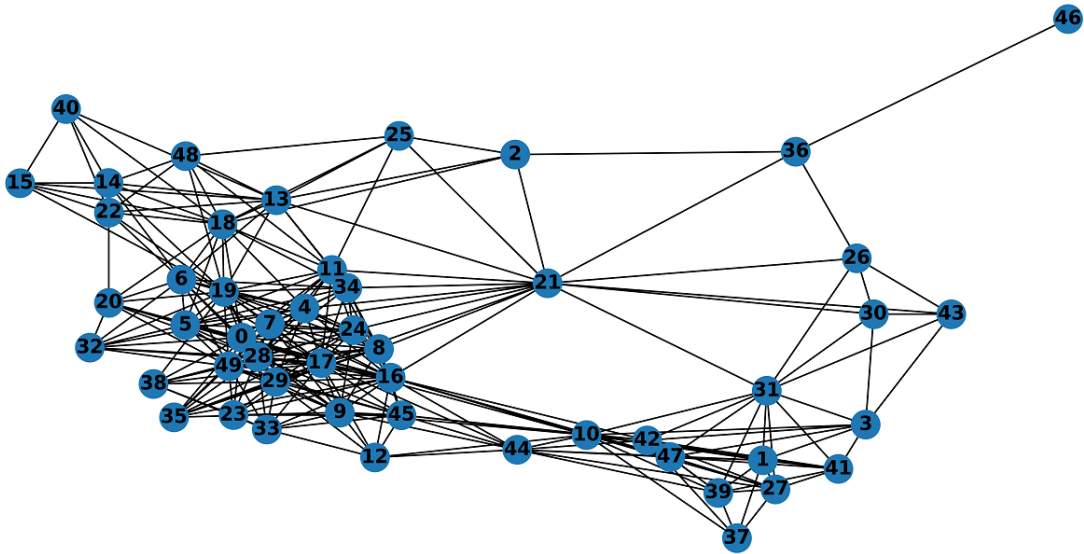
The gradient of  $f$  as obtained by the chain rule is given by:

$$\begin{aligned} \frac{\partial f}{\partial x_r}(\mathbf{x}) &= \sum_{\epsilon \in r} \frac{d(l_\epsilon(\mathbf{x}) c_\epsilon(\mathbf{x}))}{dl_\epsilon(\mathbf{x})} = \sum_{\epsilon \in r} c_\epsilon(\mathbf{x}) + l_\epsilon(\mathbf{x}) c'_\epsilon(l_\epsilon(\mathbf{x})) \\ &= \sum_{\epsilon \in r} \left( \frac{1}{2} w_\epsilon (1 + l_\epsilon(\mathbf{x}))^2 + l_\epsilon(\mathbf{x}) w_\epsilon (1 + l_\epsilon(\mathbf{x})) \right), \quad r = 1, \dots, n_{\mathcal{R}}. \end{aligned} \quad (5.2)$$

### 5.1.2 Numerical Results

For this problem we generate two geographical threshold graphs consisting of  $|\mathcal{V}| = 50, 100$  vertices placed uniformly at random in a rectangular domain. The resulting generated graph consisting of 50 vertices is illustrated in Figure 5.1. Each vertex is assigned a random weight generated from the exponential distribution. Given a threshold  $\theta$ , the vertices are joined by an edge as proposed in [35, 12]. Origin and destination nodes  $\mathcal{O} (= 46, 8)$  and  $\mathcal{T} (= 32, 74)$  respectively, are chosen randomly.

We then calculate all simple paths between the origin and destination of each graph with maximum length (henceforth called cutoff and denoted  $r_{\max}$ ). We choose  $r_{\max} = 5$  or 6 in the case of  $|\mathcal{V}| = 50$  and  $r_{\max} = 7$  in the case of  $|\mathcal{V}| = 100$ , which amount to a total of 70 and 1169, and 1647 distinct paths respectively. Therefore, the dimension of each problem is equal to  $n_{\mathcal{R}} = 70, 1169$  and  $n_{\mathcal{R}} = 1647$ . We perform experiments on both problems and plot the convergence of the loss function for the three methods.



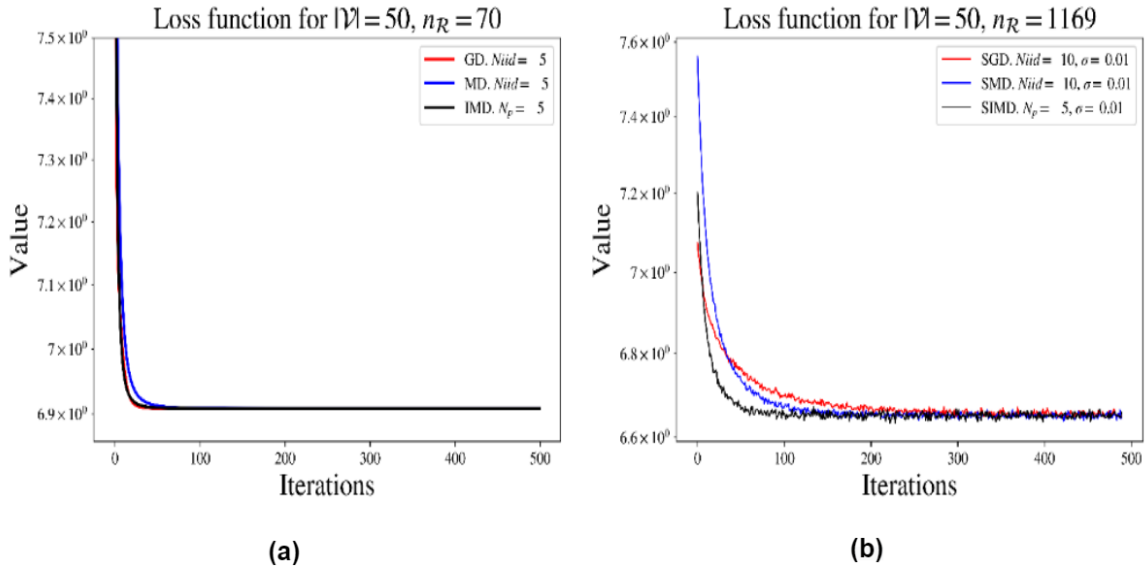
**Figure 5.1:** Traffic Assignment Problem: Random Geographical Threshold Graph,  $|\mathcal{V}| = 50$  ( $\mathcal{O}, \mathcal{T}$ ) = (46, 32).

### Experiments for random graph with $|\mathcal{V}| = 50$

We begin by demonstrating results on the random geographical threshold graph illustrated in Figure 5.1 above which contains 50 vertices. In Figure 5.2 (a) and (b) we plot the convergence of the loss function of the problem resulting when  $r_{\max} = 5$  ( $n_{\mathcal{R}} = 70$ ) and  $r_{\max} = 6$  ( $n_{\mathcal{R}} = 1169$ ) respectively, after running the three algorithms for  $T = 500$  iterations.

For the experiment in (a) we use fixed learning rate  $\eta = 0.01$  for the deterministic Gradient Descent method and fixed learning rate  $\eta = 0.2$  for deterministic Mirror Descent and its variant IMD. For GD and MD we run 5 identical and independent copies, whilst for IMD we use 5 interacting particles. We can observe that all three methods perform comparably. As expected, since the dimensionality of the problem is low ( $n_{\mathcal{R}} = 70$ ), all methods converge to the minimum of the objective function in less than 50 iterations.

For the experiment in (b) we set  $\sigma = 0.01$  and run 10 iid copies of SGD and SMD and use 5 interacting particles in the SIMD case. All other choices of hyperparameters are the same as (a). All three algorithms manage to converge to the optimum even with the added noise. However, we can notice that SIMD with 5 interacting particles and SMD outperform the SGD as they converge faster to the optimum.



**Figure 5.2:** Traffic Assignment Problem: (a) Convergence of deterministic GD, MD and IMD in the case of  $n_{\mathcal{R}} = 70$ , (b) Convergence of SGD, SMD and SIMD in the case of  $n_{\mathcal{R}} = 1169$ .

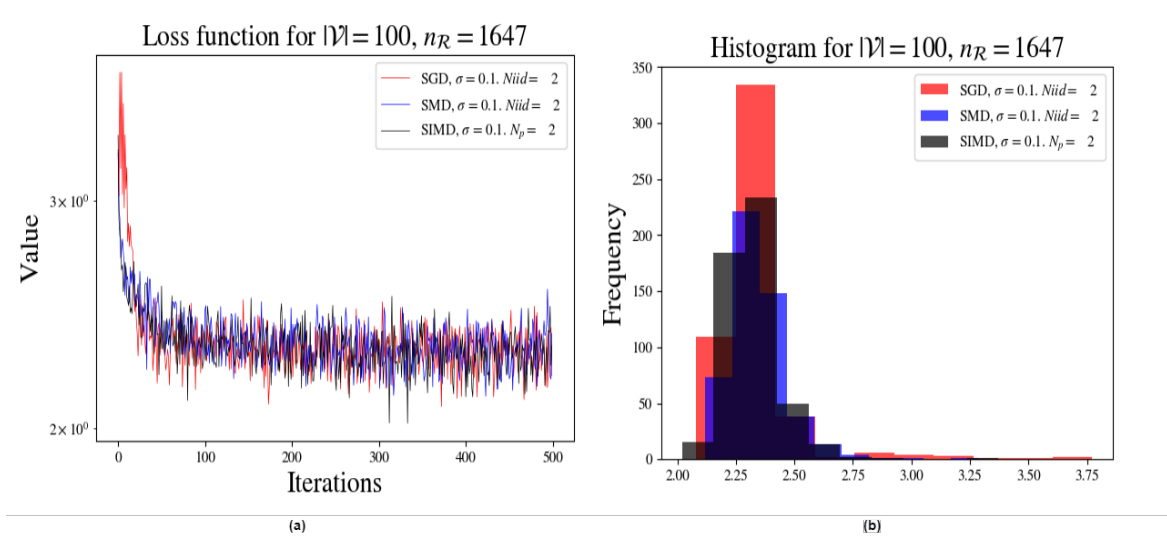
### Experiments for random graph with $|\mathcal{V}| = 100$

In this paragraph we provide the results of our experiments for a graph consisting of 100 vertices and cutoff parameter  $r_{\max} = 7$ . Note that the dimension of the problem is now  $n_{\mathcal{R}} = 1647$ .

We run all three methods in the deterministic case with noise parameter  $\sigma = 0.1$  for  $T = 500$  iterations and compare their performance. In the deterministic case we set a fixed learning rate  $\eta = 0.005$  for GD and,  $\eta = 0.05$  for MD and IMD. In the stochastic case we use a decreasing learning rate  $\eta_t = \frac{0.01}{\sqrt{t}}$  for SGD and,  $\eta_t = \frac{0.5}{\sqrt{t}}$  for SMD and SIMD.

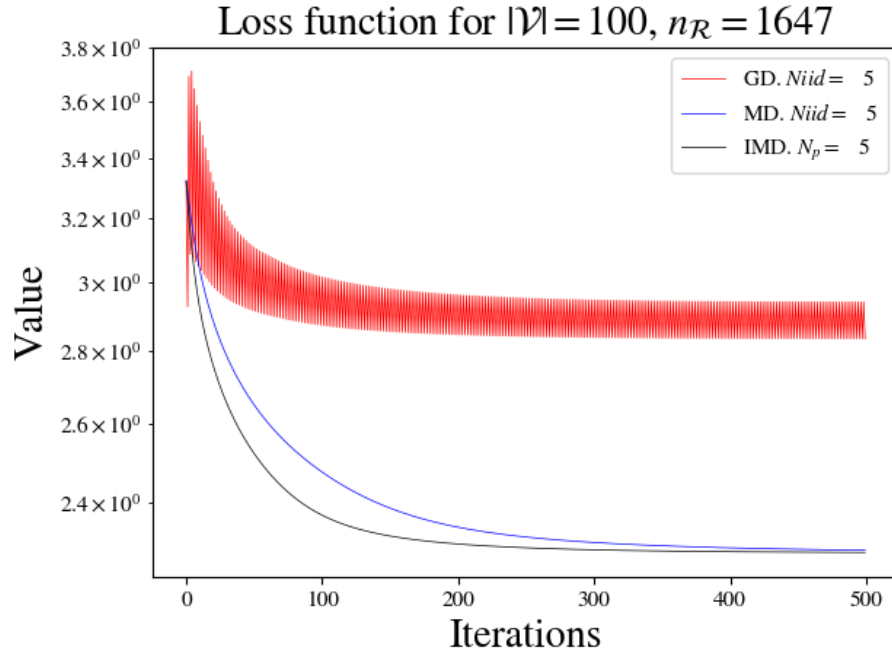
In Figure 5.3 (a), we plot the convergence of the loss function w.r.t. the iterations in the stochastic case ( $\sigma = 0.1$ ) with a decreasing learning rate. Despite the fact that SIMD and SMD converge faster than SGD, all three methods, seem to converge to a neighbourhood of the solution (the fluctuation term is governed by the noise parameter  $\sigma$ ).

However, as shown in Figure 5.4, when using fixed learning rates in the deterministic case, IMD outruns the other two methods again as it seems to converge faster to a minimum solution. In contrast, as anticipated due to the high dimensionality of the problem, GD performs much worse than MD and IMD, as it never converges to the optimum. Moreover, the samples of the loss function in the case of GD are noisier and have higher variance.



**Figure 5.3:** Traffic Assignment Problem: Convergence of stochastic optimisation with decreasing learning rate: SGD, SMD and SIMD in the case of higher dimensionality,  $n_{\mathcal{R}} = 1647$ .





**Figure 5.4:** Traffic Assignment Problem: Convergence of deterministic optimisation with fixed learning rate: GD, MD and IMD in the case of higher dimensionality,  $n_{\mathcal{R}} = 1647$ .

In Figure 5.3 (b), we provide a histogram of the samples of the function value in the stochastic case. In the next sections we are going to be studying more about the variance of the samples for each method and we will study the effect of interaction in noise reduction.

## 5.2 Linear System

In this section we are concerned with a standard optimisation problem of a linear system, often referred to as (constrained) Linear Regression, with similar a setup as Section 5.1 of [10] and Example 9.19 of [4] (where they solve with PGD).

### 5.2.1 Formulation of the linear system problem with simplex constraints

Let  $f(\mathbf{x}) = \frac{1}{2} \|\mathbf{W}\mathbf{x} - \mathbf{b}\|_2^2$ ,  $\mathbf{x} \in \mathbb{R}^n$ , where  $\mathbf{W} \in \mathbb{R}^{m \times n}$ ,  $\mathbf{b} \in \mathbb{R}^m$ ,  $m \geq n$ . We assume that  $\text{rank}(\mathbf{W}) = n$ . We are concerned with optimising  $f$  over the  $n$ -Simplex:

$$\min_{\mathbf{x} \in \mathcal{X} = \Delta^n} f(\mathbf{x}). \quad (5.3)$$

Note that  $f$  is a continuous and twice differentiable function:

$$\nabla f(\mathbf{x}) = \mathbf{W}^T (\mathbf{W}\mathbf{x} - \mathbf{b}), \quad \nabla^2 f(\mathbf{x}) = \mathbf{W}^T \mathbf{W}.$$

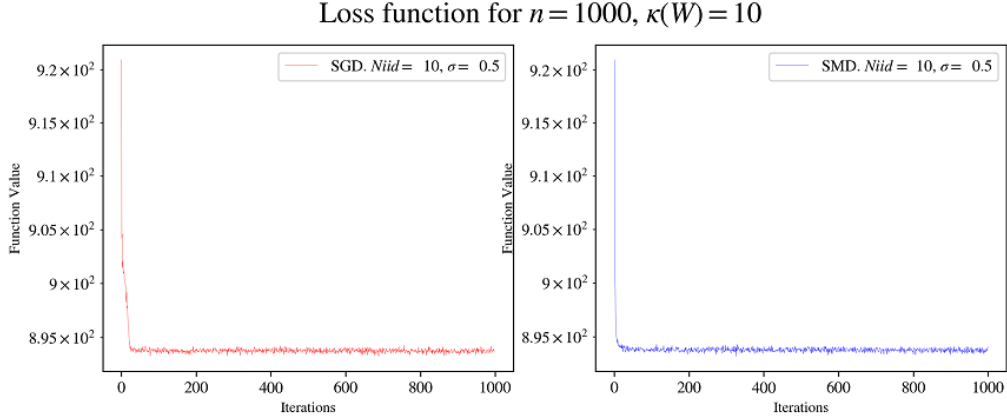
Since  $\text{rank}(\mathbf{W}) = n$ , the Hessian of  $f$  is positive definite  $\forall \mathbf{x} \in \mathbb{R}^n$  and thus,  $f$  is strictly convex on  $\mathcal{X} \subset \mathbb{R}^n$ .

### 5.2.2 Numerical Results

Here we present convergence results of the three methods (PGD, MD, IMD) for  $n = m = 1000$ . For each case, we generate randomly  $\mathbf{b} \sim \mathcal{N}(0, \mathbf{I}_n)$  and we also generate randomly one well and two ill-conditioned  $\mathbf{W}$ , with  $\kappa(\mathbf{W}) = 10$  and  $\kappa(\mathbf{W}) = 100, 200$  using the methods proposed by the authors in [6].

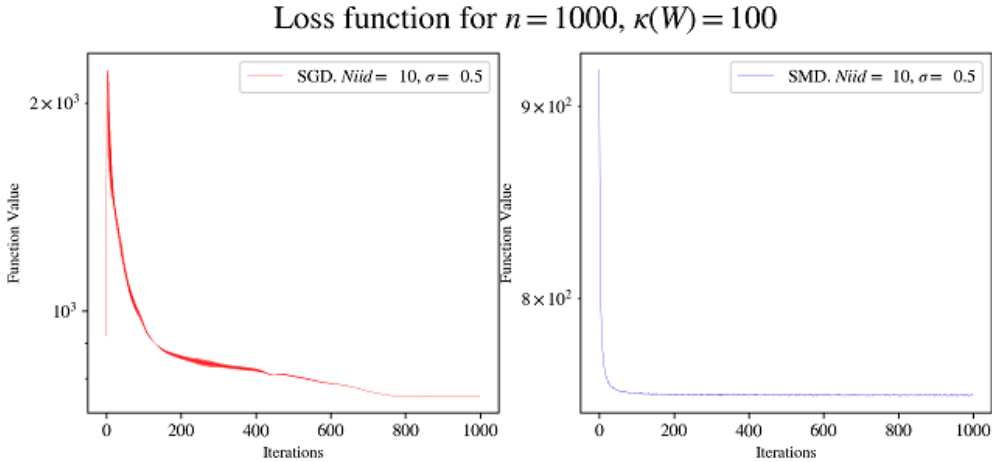
#### Single particle Optimisation: (S)GD vs (S)MD

First, we study the convergence of Stochastic Gradient Descent and Stochastic Mirror Descent (we set as the noise parameter  $\sigma = 0.5$ ) executed in 10 identical and independent runs ( $N_{iid} = 10$ ). We run both optimisation algorithms on two problems with condition of  $\mathbf{W}$ ,  $\kappa(\mathbf{W}) = 10$  and  $\kappa(\mathbf{W}) = 100$ . We then draw in parallel the loss function of each method for both cases and compare their convergence. Here, we run each algorithm for  $T = 1000$  iterations and use decreasing learning rate  $\eta_t = \frac{0.1}{\sqrt{t}}$ .



**Figure 5.5:** Linear System: SGD convergence (L) vs SMD convergence (R) for a well-conditioned  $W$  with  $\kappa(W) = 10$ .

In Figures 5.5 and 5.6 we plot the convergence of the loss function for both cases. We observe that when  $\kappa(W) = 10$ , the two methods attain similar performance as they both converge to the optimum in just a few iterations. When  $W$  is ill-conditioned, Mirror Descent seems to outperform Gradient Descent as shown in Figure 5.6, as the former converges in a few iterations, while the latter converges much slower and converges after approximately 800 iterations.



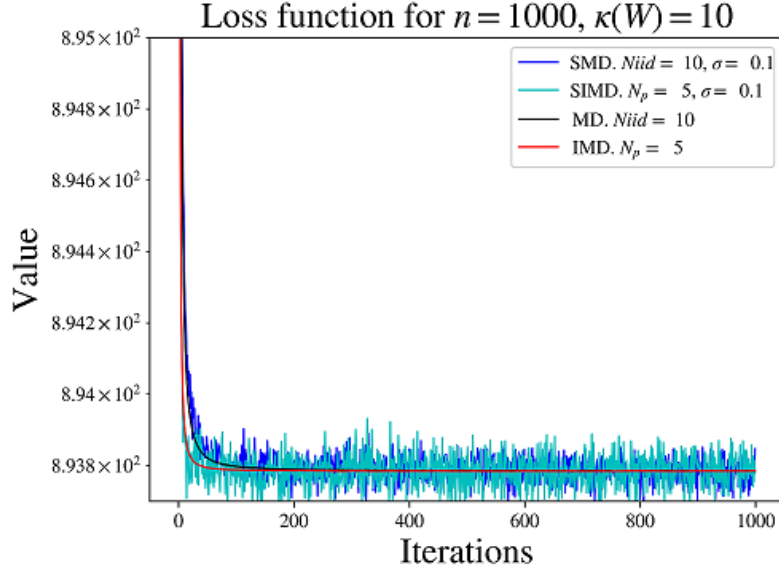
**Figure 5.6:** Linear System: SGD convergence (L) vs SMD convergence (R) for an ill-conditioned  $W$  with  $\kappa(W) = 100$ .

### Interacting particles Optimisation

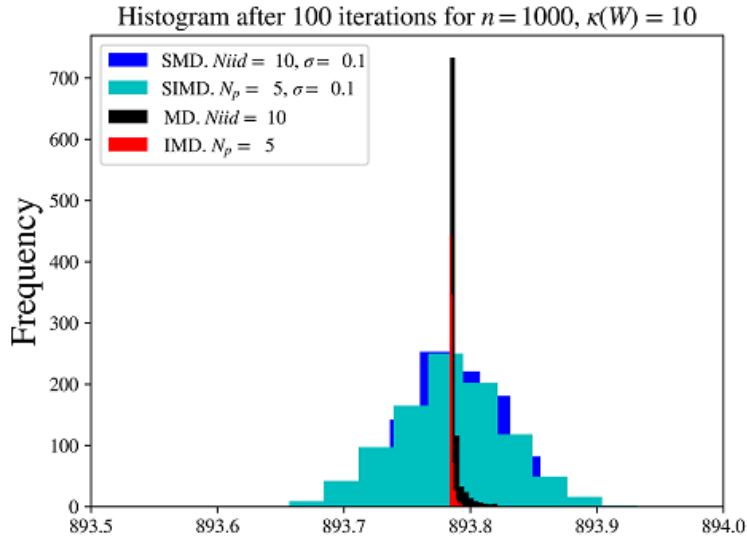
Next, we conduct experiments to investigate the convergence properties of Interacting Mirror Descent. We run all algorithms in this subsection for  $T = 1000$  iterations.

**Comparing S(MD) vs (S)IMD.** We compare the iterations of  $N_{iid} = 10$  runs of MD and SMD with single particle against one run of IMD and SIMD with  $N_p = 5$  interacting particles. Figure 5.7 shows that both methods converge to the optimum

in less than 100 iterations in the case of a well-conditioned  $W$ ,  $\kappa(W) = 10$ . Both algorithms are also able to converge even when adding some noise ( $\sigma = 0.1$ ). In Figure 5.8, we illustrate a histogram of samples after convergence. Both methods perform similarly in both, stochastic and deterministic settings.



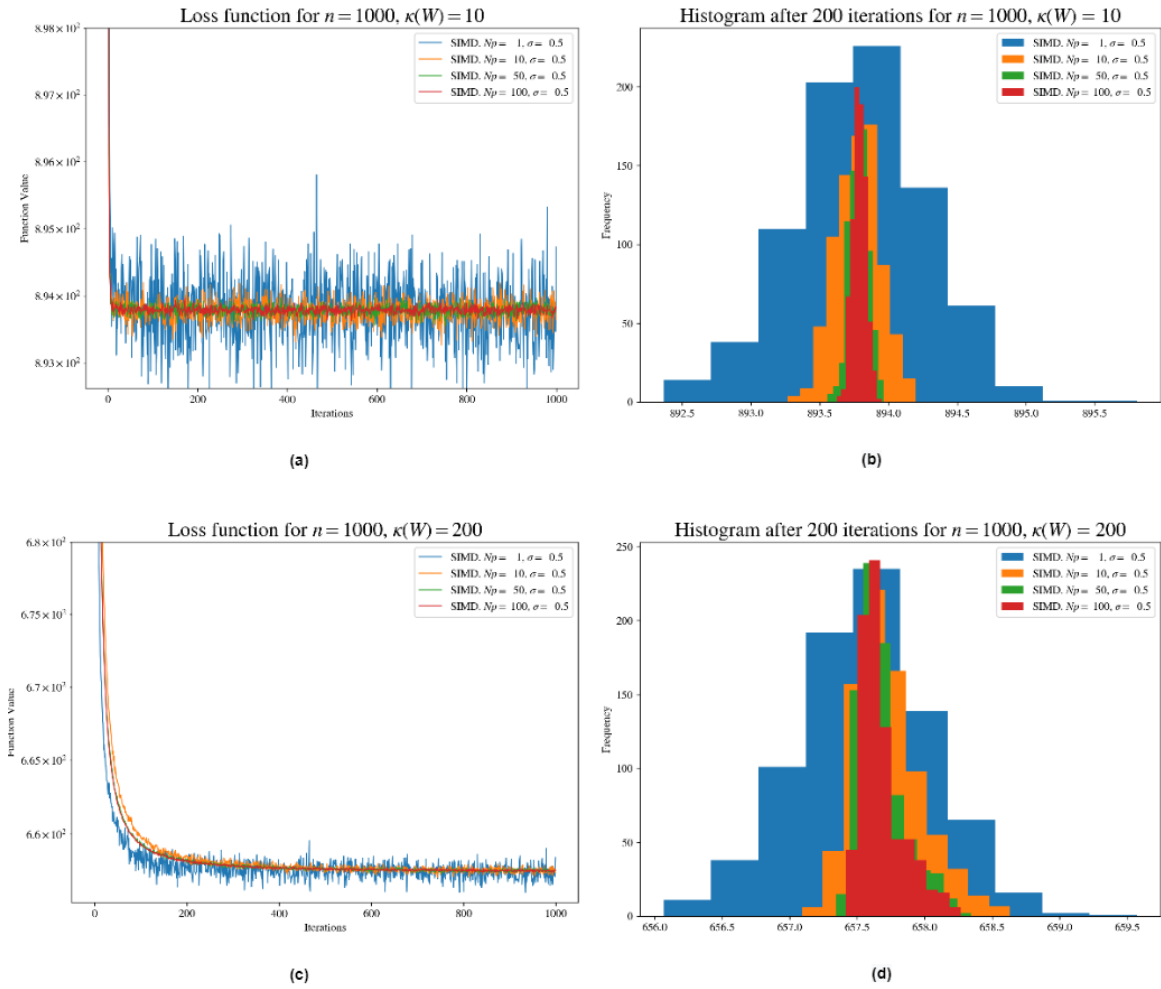
**Figure 5.7:** Linear System: (S)MD convergence vs (S)IMD convergence for a well-conditioned  $W$  with  $\kappa(W) = 10$ .



**Figure 5.8:** Linear System: Histogram of (S)MD and (S)IMD convergence for a well-conditioned  $W$  with  $\kappa(W) = 10$ .

**Interacting particles for Variance Reduction** Theoretical results in Section 3 of [10] imply that better convergence is expected for larger number of interacting particles ( $N_p$ ) used in Equation 4.21. In this paragraph we also investigate the effects of interaction in terms of the number of particles  $N_p$  against the variance of the objective function after some  $t > T_0$  iterations of the algorithm.

Using decreasing learning rate  $\eta_t = \frac{0.1}{\sqrt{t}}$  we run the stochastic case of IMD for  $T = 1000$  iterations for various numbers of interacting particles  $N_p \in \{1, 10, 50, 100\}$  in the cases when  $\kappa(\mathbf{W}) = 10$  and when  $\kappa(\mathbf{W}) = 200$  and plot the results in Figure 5.9.



**Figure 5.9:** Linear System: Effect of number of interacting particles on variance of samples of the loss function after convergence for cases of ill and well-conditioned  $\mathbf{W}$ ,  $\kappa(\mathbf{W}) = 10, 200$ .

In Table 5.1 we present the sample variance of the loss function for the two problems illustrated in Figure 5.9. Our calculations of the variance were performed on the samples obtained after  $T_0 = 200$  iterations when all methods have converged to the optimum.

First, we observe that the variance of the samples is also dependent on the condition number of the matrix  $\mathbf{W}$ . Table 5.1 suggests that the variance is higher in the ill-conditioned  $\mathbf{W}$  than the well-conditioned  $\mathbf{W}$ . Additionally, we remark (as preserved also from the illustration and as expected), the variance of the samples decreases as the number of interacting particles increases in both cases.

Variance of Interacting Mirror Descent		
$N_p$	$\kappa(\mathbf{W}) = 10$	$\kappa(\mathbf{W}) = 200$
<b>1</b>	0.2324	0.2417
<b>10</b>	0.02553	0.0652
<b>50</b>	0.00480	0.0290
<b>100</b>	0.00242	0.0251

**Table 5.1:** Variance after convergence (after 200 iterations of SIMD) for different number of interacting particles and condition number  $\kappa(\mathbf{W})$ . Illustrated above in histograms.

## 5.3 Linear System: Mini-batch Optimisation

### 5.3.1 Problem formulation

In this section we investigate a problem with similar settings as the previous section with the difference that here we are interested in minimising the average over  $M > 1$  data samples instead of a single one with Simplex constraints, where  $M$  is the sample size. Such setups are widely used in Machine and Deep Learning problems where it is standard to minimise an objective function of the form

$$f(\mathbf{x}) = \frac{1}{M} \sum_{m=1}^M f_m(\mathbf{x}).$$

Specifically, our optimisation framework here is:

$$\min_{\mathbf{x} \in \mathcal{X} = \Delta^n} f(\mathbf{x}) = \min_{\mathbf{x} \in \mathcal{X}} \frac{1}{M} \sum_{m=1}^M f_m(\mathbf{x}) \quad (5.4)$$

where

$$f_m(\mathbf{x}) = \frac{1}{2} \|\mathbf{W}_m \mathbf{x} - \mathbf{b}_m\|_2^2$$

and  $\mathbf{W}_m \in \mathbb{R}^{n \times n}$ ,  $\mathbf{b}_m \in \mathbb{R}^n$  for  $m = 1, \dots, M$ . It's important to note that since  $f_m$  is convex and continuously differentiable (see Section 5.2 for more details),  $f$  is also convex and continuously differentiable as the sum of convex and continuously differentiable functions.

Assuming that  $\mathcal{O}$  is the set of all observations ( $|\mathcal{O}| = M$ ), the gradient of the objective function in this setting is usually computed over a subset  $\mathcal{B} \subseteq \mathcal{O}$  of size  $M_{\mathcal{B}} \leq M$  i.e.:

$$\begin{aligned} \nabla f(\mathbf{x}) &\approx \nabla f_{\mathcal{B}}(\mathbf{x}) = \frac{1}{M_{\mathcal{B}}} \sum_{m: (\mathbf{W}_m, \mathbf{b}_m) \in \mathcal{B}} \nabla f_m(\mathbf{x}) \\ &= \frac{1}{M_{\mathcal{B}}} \sum_{m: (\mathbf{W}_m, \mathbf{b}_m) \in \mathcal{B}} \mathbf{W}_m^T (\mathbf{W}_m \mathbf{x} - \mathbf{b}_m) \end{aligned} \quad (5.5)$$

where

$$f_{\mathcal{B}}(\mathbf{x}) = \frac{1}{M_{\mathcal{B}}} \sum_{m: (\mathbf{W}_m, \mathbf{b}_m) \in \mathcal{B}} f_m(\mathbf{x}).$$

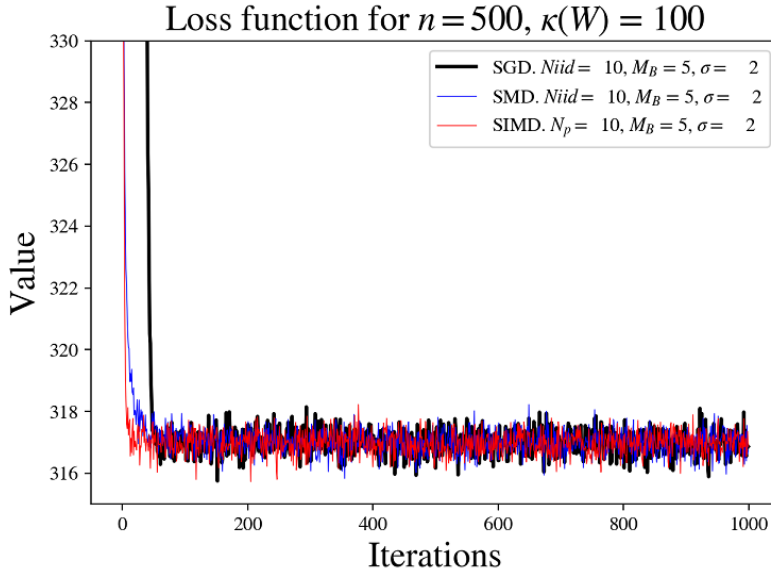
This is usually called Mini-Batch Optimisation. The need for Mini-Batch Optimisation is due to the fact that for large  $M$ , processing of the entire observation set  $\mathcal{O}$  before taking a step in gradient descent is expensive in terms of both time and computational power [11, 21].

### 5.3.2 Numerical Results

For the purposes of our experiments in this section all  $\mathbf{W}_m$  and  $\mathbf{b}_m$ ,  $m = 1, \dots, M$  are generated by adding some noise to a random non-singular full rank matrix  $\mathbf{W} \in \mathbb{R}^{n \times n}$  with a specific condition number  $\kappa(\mathbf{W})$  and a random bias vector  $\mathbf{b} \in \mathbb{R}^n$ , similarly to the previous section.

For our experiments we choose as the size of our observation set  $|\mathcal{O}| = M = 100$ . We plot and discuss results for  $n = 500$  and  $\kappa(\mathbf{W}) = 10, 100, 200$ . For our calculations here we use a fixed learning rate  $\eta = 0.05$ .

First, we execute all three algorithms on a problem with an ill-conditioned  $\mathbf{W}$ ,  $\kappa(\mathbf{W}) = 100$  and  $M_B = 5$  with stochastic settings ( $\sigma = 2$ ). We provide a plot in Figure 5.10 and compare their convergence. We observe that SIMD with 10 interacting particles governs the other two methods (averaged over 10 iid copies) as it manages to converge faster and with lower variance.

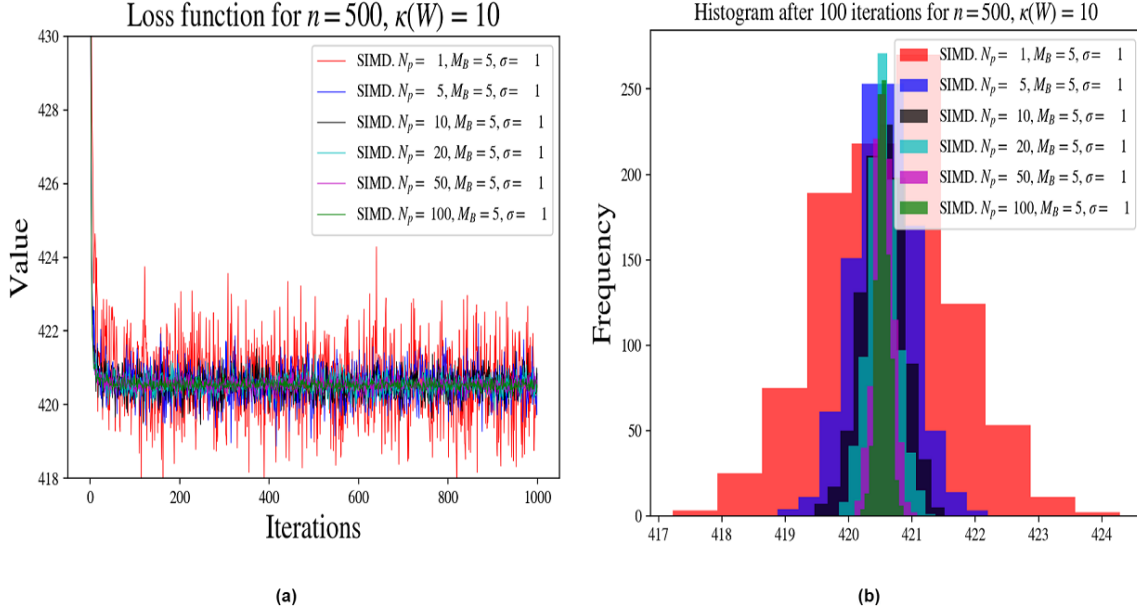


**Figure 5.10:** Mini-Batch Optimisation: Loss function for SGD and SMD ( $N_{id} = 10$ ) and SIMD ( $N_p = 10$ ) for  $M_B$  in the case of  $\kappa(\mathbf{W}) = 100$ .

#### Effect of number of particles $N_p$ in Mini-Batch Optimisation

**Experiments for fixed mini-batch size  $M_B$**  In this paragraph we present results for a well-conditioned  $\mathbf{W}$  with condition number  $\kappa(\mathbf{W}) = 10$ . We study the effect of the number of interacting particles used in (S)IMD on the convergence of the loss function when the mini-batch size is rather small. We illustrate the results in Figures 5.11, and in Table 5.2 we calculate the variance of the samples for each method.





**Figure 5.11:** Mini-Batch Optimisation: (a) Loss function and (b) Histogram of loss samples after convergence for SIMD ( $\sigma = 1.0$ ) for increasing number of interacting particles and fixed batch size  $M_B = 5$  in the case of  $\kappa(\mathbf{W}) = 10$ .

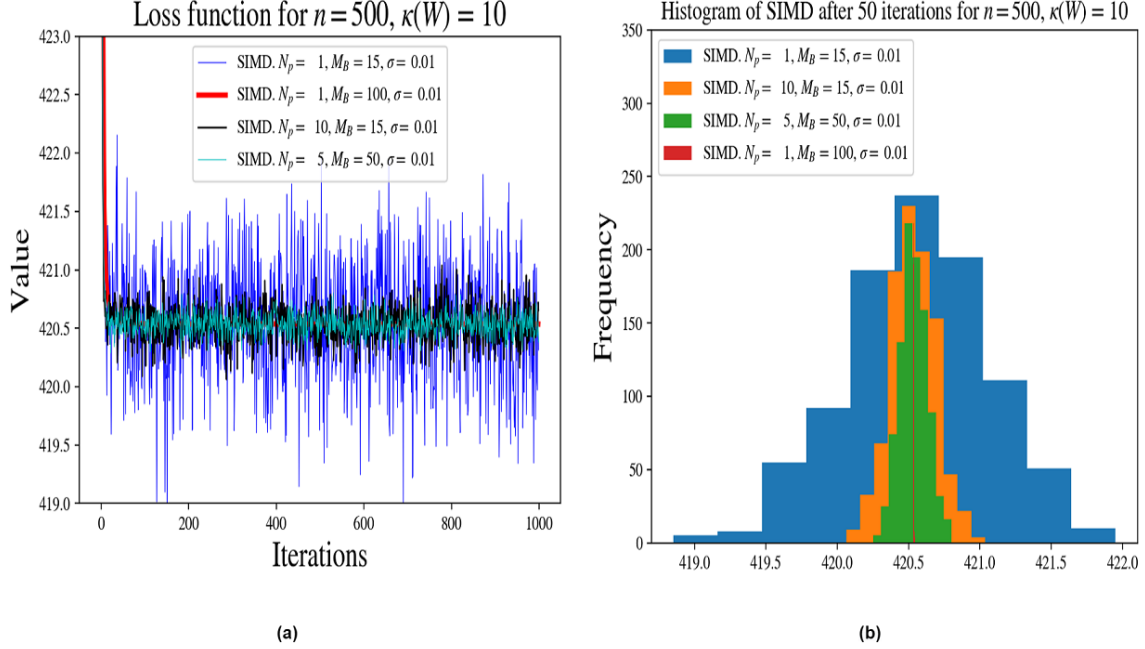
We execute the Stochastic Interacting Mirror Descent algorithm six times for different number of particles, one for each  $N_p \in \{1, 5, 10, 20, 50, 100\}$  and fixed batch size  $M_B = 5$  and we plot the results. Figure 5.11 and Table 5.2 both show that a higher number of interacting particles ( $N_p$ ) can decrease the noise of the calculations of the loss function even for a small sized batch.

$N_p$	$M_B$	Variance of Interacting Mirror Descent for $\kappa(\mathbf{W}) = 10$
1	5	1.05135
5	5	0.21831
10	5	0.10661
20	5	0.05277
50	5	0.02210
100	5	0.01010

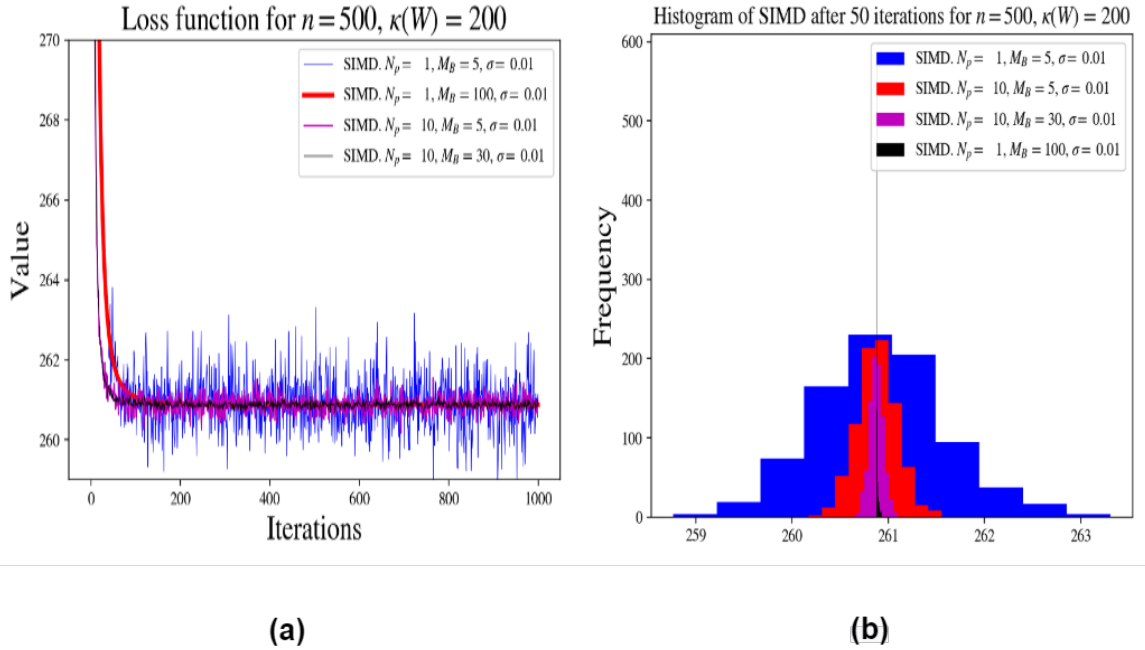
**Table 5.2:** Mini-Batch Optimisation: Variance of SIMD for fixed mini-batch size and increasing number of interacting particles after convergence.

**Experiments for different combinations of  $N_p$  and  $M_B$**  Ideally, the gradient of the loss function should be computed over the whole training batch. However, this can be expensive both, in terms of computation time and resources. Hence, we are interested in the performance of (S)IMD with respect to the size of the mini-batch since the gradient of the loss function is calculated over mini-batches and is expected to be noisy. We carry out experiments on two problems with  $\kappa(\mathbf{W}) = 10, 200$  for

different combinations of the number of interacting particles and the mini-batch size and illustrate the results in Figures 5.12 and 5.13.



**Figure 5.12:** Mini-Batch Optimisation: (a) Loss function and (b) Histogram of loss samples after convergence for SIMD ( $\sigma = 0.01$ ) for different combinations of batch sizes ( $M_B$ ) and number of interacting particles ( $N_p$ ) in the case of  $\kappa(\mathbf{W}) = 10$ .



**Figure 5.13:** Mini-Batch Optimisation: (a) Loss function and (b) Histogram of loss samples after convergence for SIMD ( $\sigma = 0.01$ ) for different combinations of batch sizes ( $M_B$ ) and number of interacting particles ( $N_p$ ) in the case of  $\kappa(\mathbf{W}) = 200$ .

$N_p$	$M_{\mathcal{B}}$	$\kappa(W)$	Variance of Interacting Mirror Descent
1	15	10	0.25049
10	15	10	0.02466
5	50	10	0.00895
1	100 (Full batch)	10	0.00001
1	5	200	0.45695
10	5	200	0.04244
10	30	200	0.00515
1	100 (Full batch)	200	0.00066

**Table 5.3:** Mini-Batch Optimisation: Variance of SIMD for different number of interacting particles ( $N_p$ ) and different mini-batch sizes ( $M_{\mathcal{B}}$ ) after convergence.

As expected, computing the gradient over a full batch ( $M_{\mathcal{B}} = 100$  here) exhibits less noise around the optimum in both cases. However, as illustrated in 5.12 and 5.13 when computing the gradient over mini-batches, using a higher number of interacting particles can result in reduction of the variance of the samples of the loss function. We also provide in Table 5.3 the variance of each method plotted in these figures.

Our calculations in the two previous paragraphs suggest that interaction is a way of combating noisy gradients and helps in reducing the variance. This result can prove helpful in cases in which it is more expensive to use large batch sizes rather than large number of interacting particles.

## 5.4 The Laplacian $K$ -Modes for Clustering

### 5.4.1 Problem formulation

In this section we consider the Laplacian  $K$ -modes algorithm for clustering described in [51]. For  $N$  samples  $\mathbf{x}_m \in \mathbb{R}^D$ ,  $m = 1, \dots, N$ , we define  $\mathbf{X} = [\mathbf{x}_1 \dots \mathbf{x}_N]^T$ . Let  $\mathbf{W}$  be a symmetric matrix that denotes the affinity matrix of the undirected graph defined by  $\mathbf{X}$ ,  $\mathcal{G} = \mathcal{G}(\mathbf{X})$ . Here, we assume that  $\mathbf{W}$  is either binary (the Adjacency matrix of  $\mathcal{G}$ ):

$$W_{mn} = \begin{cases} 1 & \text{if } \mathbf{x}_m \in kNN(\mathbf{x}_n) \\ 0 & \text{otherwise,} \end{cases}$$

or obtained by the Radial Basis Function (RBF):

$$W_{mn} = \exp(-\gamma \|\mathbf{x}_m - \mathbf{x}_n\|_2^2)$$

for  $n, m = 1, \dots, N$ . Given the above, we want to optimise the objective function w.r.t.  $\mathbf{Z} = [\mathbf{z}_1^T, \dots, \mathbf{z}_N^T]^T \in \mathbb{R}^{N \times K}$  and  $\mathbf{C} = [\mathbf{c}_1^T, \dots, \mathbf{c}_K^T]^T \in \mathbb{R}^{K \times D}$ :

$$\begin{aligned} \min_{\mathbf{Z}, \mathbf{C}} \quad & \frac{\lambda}{2} \sum_{m=1}^N \sum_{n=1}^N W_{mn} \|\mathbf{z}_m - \mathbf{z}_n\|_2^2 - \sum_{m=1}^N \sum_{k=1}^K z_{mk} \mathcal{N}\left(\frac{\|\mathbf{x}_m - \mathbf{c}_k\|_2^2}{\sigma^2}\right) \\ \text{s.t.} \quad & \sum_{k=1}^K z_{mk} = 1, \quad m = 1, \dots, N, \\ & z_{mk} \geq 0, \quad m = 1, \dots, N \quad k = 1, \dots, K, \end{aligned} \tag{5.6}$$

where  $\lambda \geq 0$  represents a trade-off parameter, each  $[z_{m1} \dots, z_{mK}]^T = \mathbf{z}_m \in \mathbb{R}^K$  represents soft assignments of sample  $\mathbf{x}_m$  to  $K$  clusters and each  $\mathbf{c}_k \in \mathbb{R}^D$ ,  $k = 1, \dots, K$  represents modes (centers) of each cluster. Note that  $\mathcal{N}(\cdot)$  denotes the Normal probability distribution function.

Let further  $\mathbf{L} = \mathbf{D} - \mathbf{W}$  denote the Laplacian of  $\mathcal{G}$ , where  $\mathbf{D} = \text{diag}([\sum_{n=1}^N W_{mn}]_{m=1}^N)$  is the degree matrix of  $\mathcal{G}$ . The matrix form of 5.6 is given by:

$$\begin{aligned} \min_{\mathbf{Z}, \mathbf{C}} f(\mathbf{Z}, \mathbf{C}) &= \min_{\mathbf{Z}, \mathbf{C}} \lambda \text{tr}(\mathbf{Z}^T \mathbf{L} \mathbf{Z}) - \text{tr}(\mathbf{B}^T \mathbf{Z}) \\ \text{s.t.} \quad & \mathbf{Z} \mathbf{1}_K = \mathbf{1}_N, \\ & z_{mk} \geq 0, \quad m = 1, \dots, N \quad k = 1, \dots, K, \end{aligned} \tag{5.7}$$

where the elements of  $\mathbf{B}$  are

$$B_{mk} = \mathcal{N}\left(\frac{\|\mathbf{x}_m - \mathbf{c}_k\|_2^2}{\sigma^2}\right), \quad m = 1, \dots, N, \quad k = 1, \dots, K.$$

Note that since  $\mathbf{L}$  is positive semi-definite, it can be decomposed as a product  $\mathbf{L} = \mathbf{A} \mathbf{A}^T$  such that  $\mathbf{A}$  has full row rank  $r = \text{rank}(\mathbf{L})$  and thus,  $\text{tr}(\mathbf{Z}^T \mathbf{L} \mathbf{Z}) = \|\mathbf{Z}^T \mathbf{A}\|_F^2$ .

The optimisation problem 5.7 is convex w.r.t  $\mathbf{Z}$ :

**Proof 9** Assume a fixed  $\mathbf{B}$ . Let  $g : \mathcal{X} \rightarrow \mathbb{R}$  be defined by:

$$g(\mathbf{Z}) = \lambda \operatorname{tr}(\mathbf{Z}^T \mathbf{L} \mathbf{Z}) - \operatorname{tr}(\mathbf{B}^T \mathbf{Z})$$

and let  $\mathbf{X}, \mathbf{Y} \in \mathbb{R}^{N \times K}$ , such that  $\mathbf{x}_m, \mathbf{y}_m \in \mathcal{X} = \Delta^N$ ,  $m = 1, \dots, N$ ,  $\alpha \in (0, 1)$ .

$$\begin{aligned} g(\alpha \mathbf{X} + (1 - \alpha) \mathbf{Y}) &= \frac{\lambda}{2} \|(\alpha \mathbf{X} + (1 - \alpha) \mathbf{Y})^T \mathbf{A}\|_F^2 - \operatorname{tr}(\mathbf{B}^T (\alpha \mathbf{X} + (1 - \alpha) \mathbf{Y})) \\ &\leq \frac{\lambda}{2} \|\alpha \mathbf{X}^T \mathbf{A}\|_F^2 + \frac{\lambda}{2} \|(1 - \alpha) \mathbf{Y}^T \mathbf{A}\|_F^2 - \operatorname{tr}(\mathbf{B}^T (\alpha \mathbf{X} + (1 - \alpha) \mathbf{Y})) \\ &= \alpha g(\mathbf{X}) + (1 - \alpha) g(\mathbf{Y}). \end{aligned}$$

■

The gradient of  $f$  w.r.t  $\mathbf{Z}$  is given by:

$$\nabla_{\mathbf{Z}} f(\mathbf{Z}, \mathbf{C}) = 2 \cdot \lambda \mathbf{L} \mathbf{Z} - \mathbf{B} \quad (5.8)$$

and w.r.t.  $\mathbf{c}_k$ ,  $k = 1, \dots, K$ :

$$\begin{aligned} \nabla_{\mathbf{c}_k} f(\mathbf{Z}, \mathbf{C}) &= \sum_{m=1}^N \frac{2}{\sigma^2} \cdot z_{mk} \cdot \mathcal{N}'\left(\frac{\|\mathbf{x}_m - \mathbf{c}_k\|_2^2}{\sigma^2}\right) (\mathbf{x}_m - \mathbf{c}_k) \\ &= - \sum_{m=1}^N \frac{2}{\sigma^2} \cdot z_{mk} \cdot \frac{\|\mathbf{x}_m - \mathbf{c}_k\|_2^2}{\sigma^2} \cdot \mathcal{N}\left(\frac{\|\mathbf{x}_m - \mathbf{c}_k\|_2^2}{\sigma^2}\right) (\mathbf{x}_m - \mathbf{c}_k). \end{aligned} \quad (5.9)$$

The optimisation with respect to each row  $\mathbf{z}_m$ ,  $m = 1, \dots, N$ , of  $\mathbf{Z}$  (for a fixed  $\mathbf{C}$ ) is a convex constrained optimisation problem on the  $K$ -Simplex and therefore fits our framework. We utilise the matrix form of the problem and update  $\mathbf{Z}$  “altogether” in the cases of MD and IMD. The optimisation with respect to  $\mathbf{C}$  for a fixed  $\mathbf{Z}$  is unconstrained and solved with the traditional method, GD.

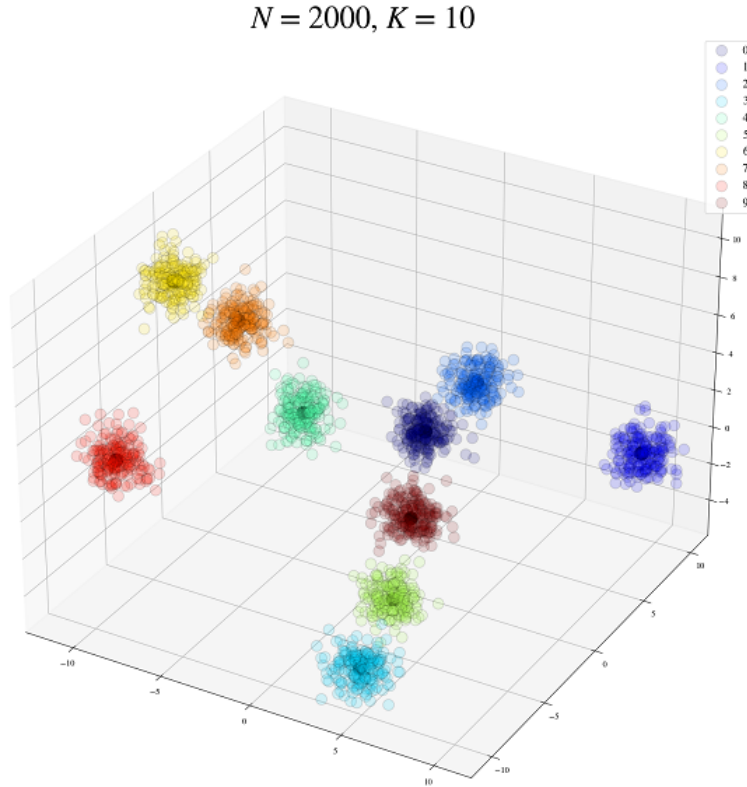
### 5.4.2 Numerical Results

To solve (5.7) we alternate between optimising w.r.t.  $\mathbf{Z}$  for a fixed  $\mathbf{C}$  and optimising w.r.t.  $\mathbf{C}$  for a fixed  $\mathbf{Z}$ . It should be mentioned that in [51] they suggest that there is a trade-off between the two hyperparameters  $\lambda$  and  $\sigma$  which we do not investigate here as it is out of the scope of this project. We specify  $\lambda$  and  $\sigma$  for each experiment.

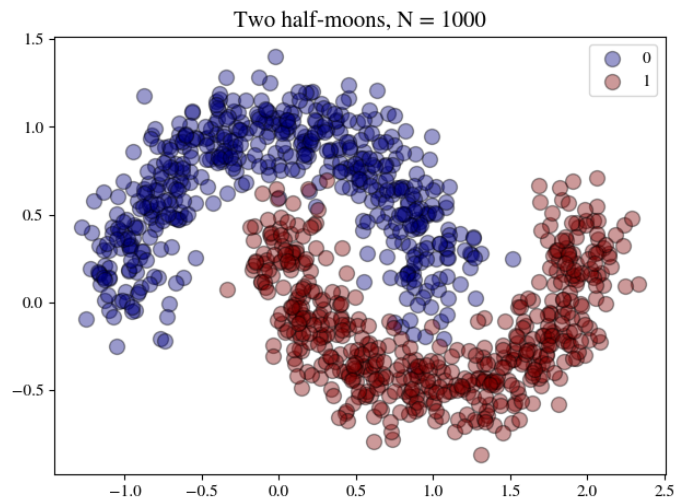
We present results for two problems by generating synthetic data. For the first, we generate  $N$  random data-points in  $K$  isotropic Gaussian blobs using the method `make_blobs` from [42]. For the second problem, we generate  $N$  two-dimensional data-points split into two “half-moon” shaped clusters (i.e.  $D = K = 2$ ) using the method `make_moons` from [42].

In Figures 5.14 and 5.15 are illustrated examples of generated datasets for each of

the problems. Note that for both problems we choose randomly the initial centroids  $\mathbf{c}_k^0$ ,  $k = 1, \dots, K$  which we optimise.



**Figure 5.14:** Laplacian  $K$ -Modes for Clustering:  $N = 2000$  randomly generated data points in the 3D space divided into  $K = 10$  clusters (isotropic Gaussian blobs). The black dots illustrate the (true) centroid  $\mathbf{c}_k$  of each cluster.



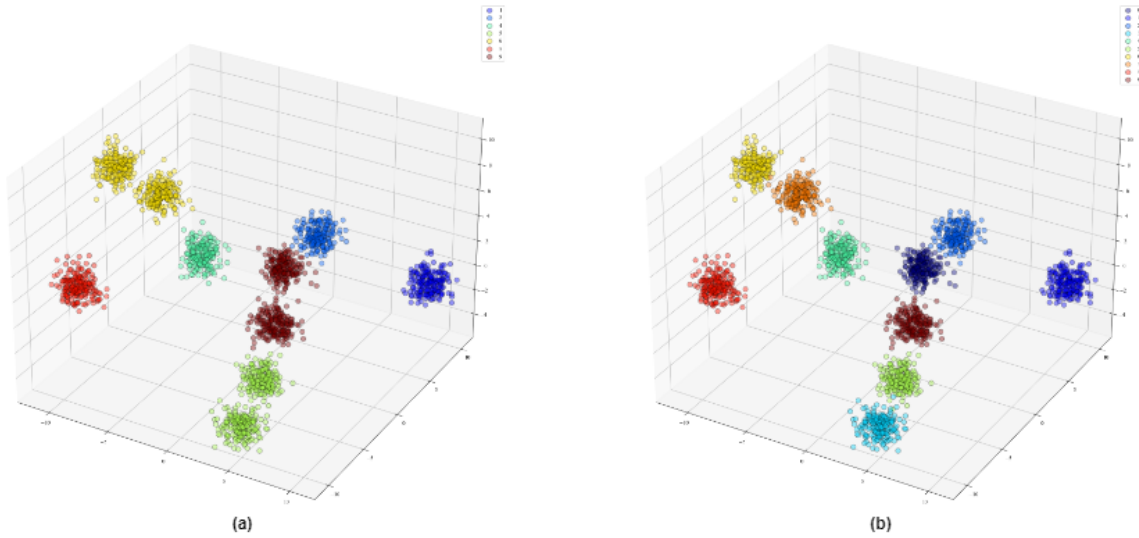
**Figure 5.15:** Laplacian  $K$ -Modes for Clustering:  $N = 1000$  randomly generated data points in the 2D space divided into  $K = 2$  half-moon shaped clusters.

### Isotropic Gaussian blobs

In this paragraph, we report results for the optimisation of the problem proposed in Equation 5.7 for  $K = 10$  3D isotropic Gaussian blobs as shown in Figure 5.14. First, we study the performance of the deterministic Gradient Descent ran in two identical and independent runs and the deterministic Interacting Mirror Descent with two particles. For our experiments here as hyperparameters we set  $\lambda = 100$  and  $\sigma = 10$ . Note that as an initialisation of  $\mathbf{c}_k$  we randomly generate  $K$  vectors in the box bounded by all  $N$  samples.

We demonstrate the assignment results after running each algorithm for  $T = 1000$  iterations in Figure 5.16 and Table 5.4. It should be mentioned that we assign each sample to a cluster by using a classification threshold  $\theta = 0.1$ , i.e. we assign a sample  $\mathbf{x}_m$ ,  $m = 1, \dots, N$  in cluster  $k$ ,  $k = 1, \dots, K = 10$  if:

$$Z_{mk} \in \left[ \frac{k-1}{K}, \frac{k}{K} \right).$$

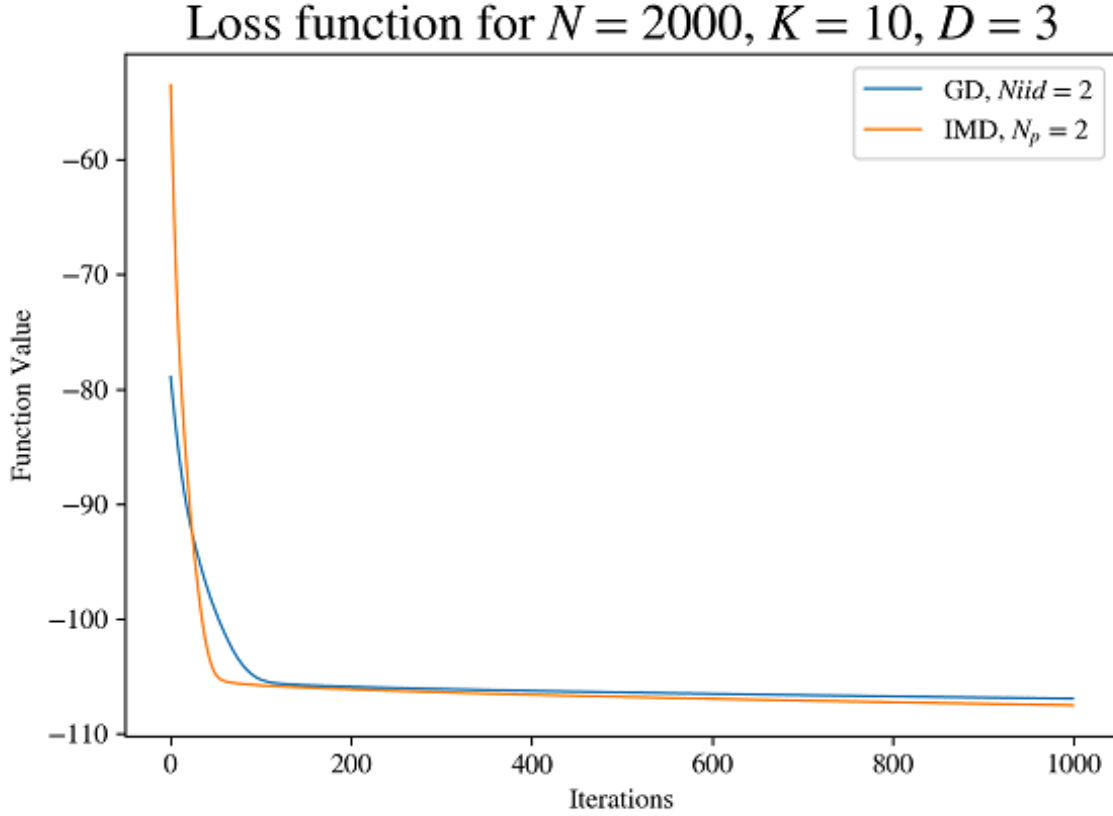


**Figure 5.16:** Laplacian  $K$ -Modes for Clustering: Cluster assignment results of (a) GD and (b) IMD using a classification threshold equal to 0.1. Their achieved numerical “accuracies” are reported in Table 5.4.

	AUC	“Accuracy”
GD	0.9833	0.8493
IMD	0.9996	0.9980

**Table 5.4:** Laplacian  $K$ -Modes for Clustering: Classification performance of GD and IMD in the case of  $K = 10$  synthetically generated Isotropic Gaussian 3D blobs. Each blob contains approximately  $\frac{N}{K} = \frac{2000}{10} = 200$  data samples.

We also plot in Figure 5.17 the performance of both methods on the optimisation problem for  $N = 2000$ . Once again, similarly to the previous problems, we remark that IMD reaches the global optimum faster than GD.



**Figure 5.17:** Laplacian  $K$ -Modes for Clustering: Loss convergence of GD and IMD for  $N = 2000$ .

Below, in Table 5.5, are indicated the average run-times for all three methods studied in this project. As anticipated, Gradient Descent is much more computationally expensive, as at each iteration we project  $N = 2000$  vectors on the  $K$ -Simplex which is costly in terms of time and computational power. In contrast, Interacting Mirror Descent seems to be robust to large dimensionality problems, as expected.

	Average Runtime for a single iteration
<b>GD</b>	1.215 s for 1 iid run
<b>IMD</b>	0.840 s per particle

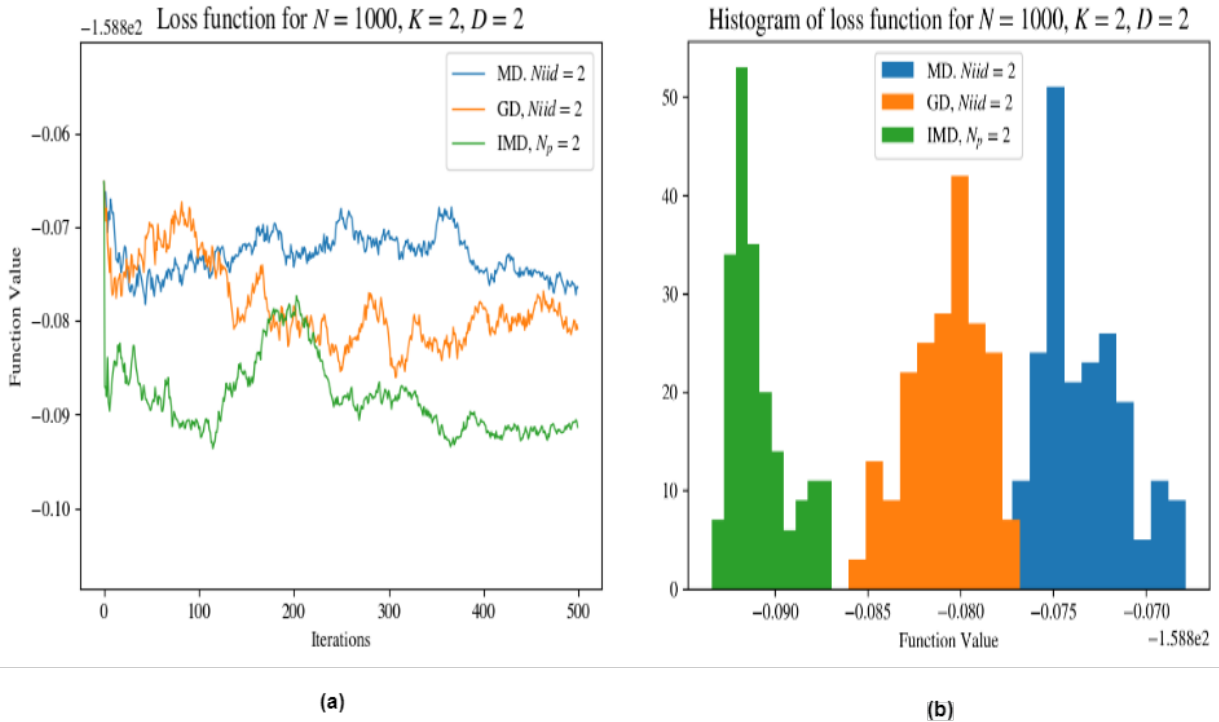
**Table 5.5:** Laplacian  $K$ -Modes for Clustering: Average run-time for one optimisation iteration for the three methods in the case of synthetic Isotropic Gaussian blobs.



### Two “half-moons”

We study the performance of the deterministic versions of Gradient and Mirror Descent ran in two identical and independent runs and Interacting Mirror Descent with two particles. For our experiments here, as hyperparameters we set  $\lambda = 2$  and  $\sigma = 0.65$ . We assign each sample to each cluster similarly to the previous problem with threshold  $\theta = 0.5$ , i.e. we assign  $\mathbf{x}_n$  to cluster  $k = 0$  if  $z_{n0} < 0.5$  and to cluster  $k = 1$  otherwise.

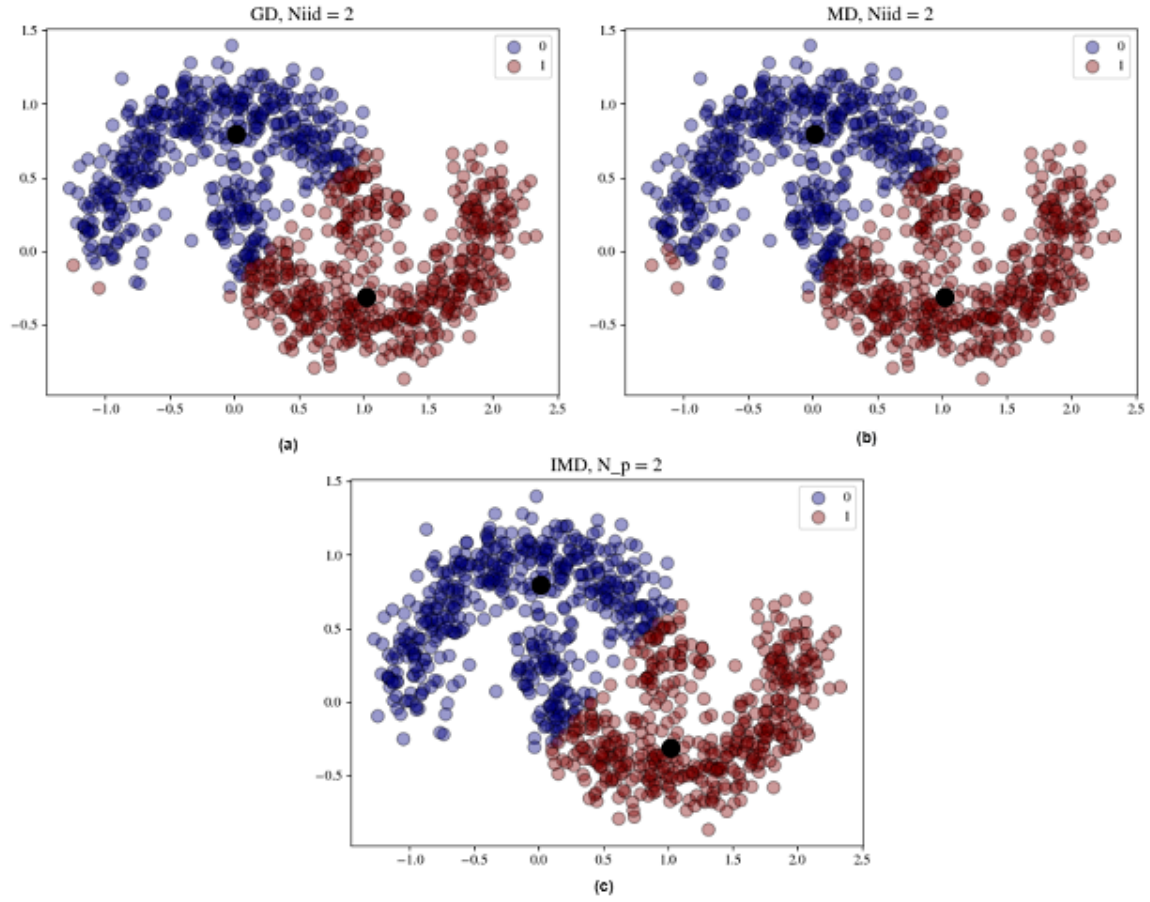
Below we include a plot in Figure 5.19 of the assignments of each method and we provide a table with their respective scores (Table 5.6). Also, in Figure 5.18, we plot the loss functions achieved by the three methods. We observe that in all three cases the performance is comparable. However, IMD achieves a higher “accuracy” and a better convergence to the minimum than the other methods.



**Figure 5.18:** Laplacian  $K$ -Modes for Clustering: (a) Loss function and (b) Histogram of loss samples after convergence for the three methods for  $N = 1000$ .

	AUC	“Accuracy”
<b>GD</b>	0.815	0.8664
<b>MD</b>	0.814	0.8539
<b>IMD</b>	0.805	0.9086

**Table 5.6:** Laplacian  $K$ -Modes for Clustering: Classification performance of the three methods in the case of two synthetically generated half-moons clusters.



**Figure 5.19:** Laplacian  $K$ -Modes for Clustering: Cluster assignment results of (a) GD, (b) MD and (c) SIMD using a classification threshold equal to 0.5 in the case of two synthetically generated half-moons clusters. The black points illustrate the predicted centroids  $c_1$  and  $c_2$ .

## 5.5 Optimisation on SDPs

In this section we investigate a convex optimisation problem on SDPs for PSD matrices over the  $n$ -Spectrahedron. Specifically, we are concerned with connected subgraph anomaly detection in graph-structured data. Such problems can be applied in numerous instances, such as detection of anomalies in disease outbreaks, detection of anomalous incidents in social networks, identification of anomalous groups, clusters or paths in graphs, etc. We investigate a convex relaxation of a problem that aims in detecting whether there exist any connected sets of vertices in a graph that transmit abnormal signal values [1, 60].

### 5.5.1 Connected Subgraph Detection on SDPs: Problem formulation

Let  $\mathcal{G} = (\mathcal{V}, \mathcal{E})$  denote a directed, weighted graph, where  $\mathcal{V}$  denotes a set of  $n$  vertices and  $\mathcal{E}$  the set of edges connecting the vertices. For root vertex  $r \in S$ , we consider an Elevated Mean Detection, i.e. we want to detect the existence of a connected in  $G$  subgraph  $S \in \Lambda_r = \{S \subset V : r \in S, G_S \text{ is connected in } G\}$  with vertices transmitting an elevated mean compared to the other nodes. See [1] for a detailed formulation of the problem. In particular, we aim to solve a relaxed convex SDP optimisation problem as shown in Section 4 of [1]:

$$\begin{aligned} \max_{s \geq 0, \mathbf{M} \in \Delta^{n \times n}} \quad & \mathbf{C} \bullet \mathbf{M} - s \\ \text{s.t.} \quad & \mathbf{Q}_\gamma(\mathbf{M}) + s \cdot \delta \cdot \mathbf{D} \succeq \mathbf{0} \end{aligned} \quad (5.10)$$

where  $\mathbf{C} = \mathbf{x} \mathbf{x}^T$  for  $\mathbf{x} \in \mathbb{R}^n$ ,  $\mathbf{x} \geq \mathbf{0}$ ,  $s$  a slack variable,  $\delta \geq 0$  a fixed margin and  $\mathbf{D}$  is the degree matrix of the given graph and  $\mathbf{Q}_\gamma(\mathbf{M}) = \mathbf{L}_{G[\mathbf{M}]} - \frac{\gamma^2}{2} \mathbf{L}_{Star(r)[\mathbf{M}]}$ .

For more details on the operators  $\mathbf{L}_{G[\cdot]}$ ,  $\frac{\gamma^2}{2} \mathbf{L}_{Star(r)[\cdot]}$ ,  $\mathbf{Q}_\gamma(\cdot)$  refer to [1]. Note that the value of  $\gamma$  is used to quantify the size and conductance of the anomalous graph. In our experiments in the next section we use different values of  $\gamma$  to get our results.

Now, let  $\mathbf{Y} \succeq \mathbf{0}$  be the Lagrange multiplier in the Lagrangian of the optimisation problem associated to the constraint  $\mathbf{Q}_\gamma(\mathbf{M}) + s \cdot \delta \cdot \mathbf{D}$ . Then the dual of Equation 5.10 is given by the following problem :

$$\min_{\mathbf{Y} \in \Delta_D^{n \times n}} f(\mathbf{Y}) \quad (5.11)$$

where

$$f(\mathbf{Y}) = \max_{\mathbf{M} \in \Delta^{n \times n}} \left( \mathbf{C} + \frac{1}{\delta} \cdot \mathbf{P}_\gamma(\mathbf{Y}) \right) \bullet \mathbf{M}$$

and

$$\Delta_D^{n \times n} = \left\{ \mathbf{X} \in \mathcal{S}_+^{n \times n} \mid \text{tr}(\mathbf{D} \cdot \mathbf{X}) = 1 \right\}.$$

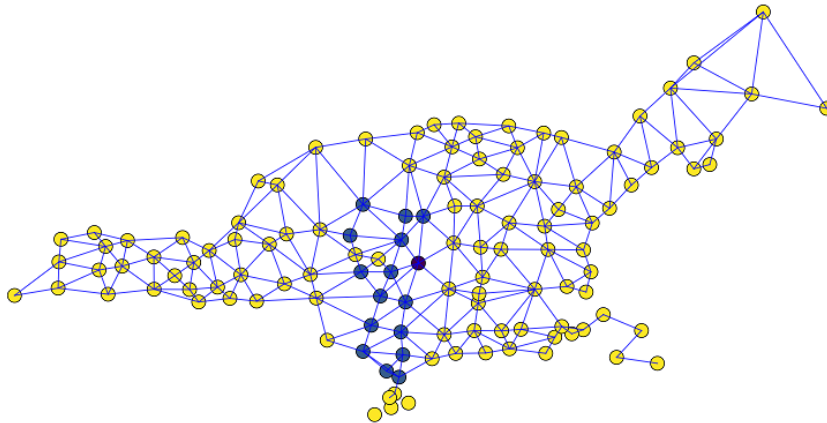
Note that  $\Delta_D^{n \times n}$  is a convex set in  $\mathbb{R}^{n \times n}$  and therefore the optimisation over the Lagrange multiplier  $\mathbf{Y}$  fits our framework. We utilise the MD and IMD algorithms for PSD matrices presented in Chapter 3 over the  $n$ -Spectrahedron. It should also be mentioned that for the optimisation of the primal variable  $M$  we utilise the update rule presented in Section 4 of [1] which is out of the scope of this project.

We present performance results of our methods on two problems: Disease Outbreak Anomaly Detection for a real case scenario and Elevated Mean anomaly classification on randomly generated geometric graphs. We will be referring to  $\text{SNR} = \frac{\lambda_1}{\lambda_0}$  as signal-to-noise ratio where  $\lambda_0$  denotes the base disease ratio and  $\lambda_1$  the anomalous disease ratio.

### 5.5.2 Disease Outbreak Anomaly Detection: Numerical Results

Here we study a real-world graph  $G = (\mathcal{V}, \mathcal{E})$ , also studied in [1].  $\mathcal{V}$  constitutes of 129 vertices corresponding to 129 counties of Northeastern United States. For our null hypothesis ( $H_0$ ) we consider a subgraph  $S \subset G$  of size  $|\mathcal{V}_S| = 16$  as the abnormal case, where  $\mathcal{V}_S = \{69, 66, 105, 93, 106, 95, 49, 90, 68, 59, 104, 62, 88, 108, 84, 92\}$ .

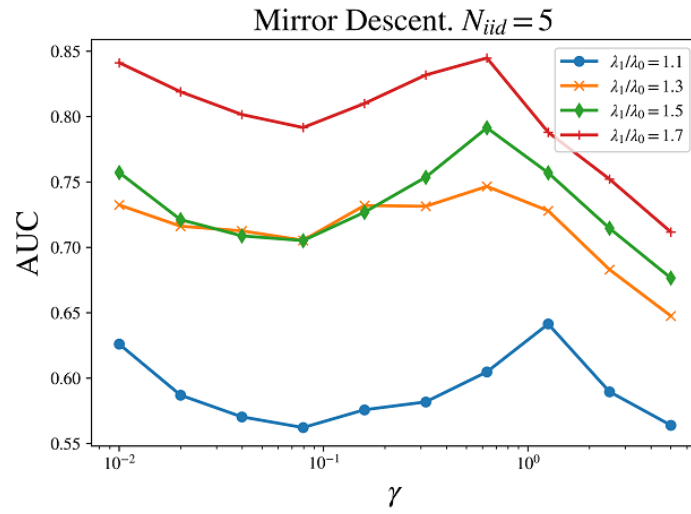
Letting  $p_j$  denote the population of each county, we assume the number of disease cases for each county as  $y_j \sim \text{Poisson}(p_j \lambda_0)$  and  $y_j \sim \text{Poisson}(p_j \lambda_1)$  for the non-anomalous and anomalous counties respectively. We then let  $\mathbf{C} = \mathbf{x}\mathbf{x}^T$  where  $x_j = \frac{y_j}{p_j}$  for  $j = 1, \dots, 129$  denotes the rate of infection in each county. See Figure 5.20 for an illustration of the graph and the ground truth subgraph.



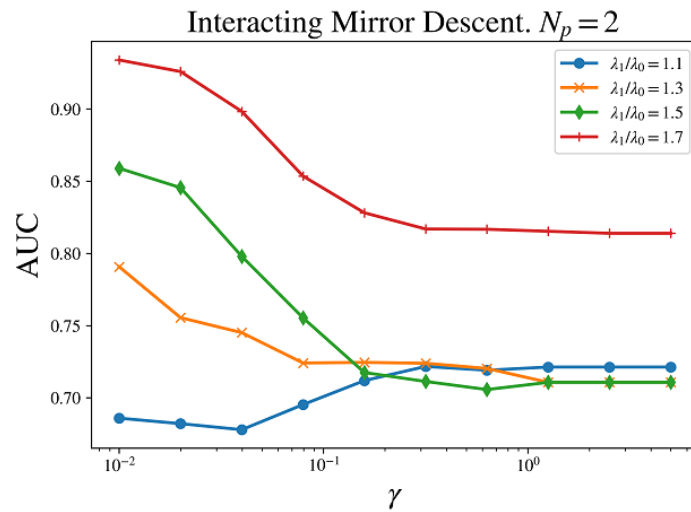
**Figure 5.20:** Disease Outbreak Anomaly Detection: a real case scenario, Graph of NE United States. Nodes colored blue illustrate  $H_0$  for the anomalous subgraph (denoted  $\mathcal{V}_S$  in the previous paragraph). The anchor node is colored purple.

### Performance of IMD vs MD for different values of $\gamma$

In Figures 5.21 and 5.22 we illustrate the detection performance results of both, Mirror Descent and Interacting Mirror Descent respectively, for 10 different values of  $\gamma$  evenly spaced on log scale between 0.01 and 5.00. To measure the performance we use different values of threshold and calculate the AUC score by calculating for each threshold the true positive and false positive rates of the classification of anomalous clusters over a sample set generated randomly from the null hypothesis  $H_0$  and alternative hypothesis  $H_1$ . We execute all optimisation algorithms over a random sample set of size  $x_{samp} = 50$ . All algorithms were ran with fixed learning rate  $\eta = 5$  for  $T = 50$  iterations.



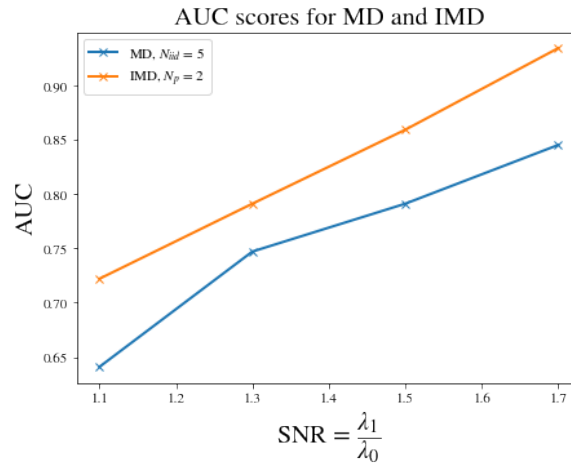
**Figure 5.21:** Disease Outbreak Anomaly Detection: AUC scores of MD (averaged over 5 iid runs) for disease ratios  $\{1.1, 1.3, 1.5, 1.7\}$  and different values of  $\gamma$ .



**Figure 5.22:** Disease Outbreak Anomaly Detection: AUC scores of IMD with 2 interacting particles for disease ratios  $\{1.1, 1.3, 1.5, 1.7\}$  and different values of  $\gamma$ .

In Figures 5.21 and 5.22 we illustrate the AUC scores attained by Interacting Mirror Descent with 2 interacting particles and Mirror Descent, respectively, considering for each method four different SN-ratios  $\frac{\lambda_1}{\lambda_0} = 1.1, 1.3, 1.5, 1.7$ . We mention that for the base disease rate for all cases we set  $\lambda_0 = 5 \times 10^{-5}$ .

In Table 5.7 we report the best AUC scores for each disease ratio and plot them in Figure 5.23. Observing both, the Figures and the Table we can see that IMD with only two interacting particles outperforms an averaged over 5 runs MD for all values of ratios  $\frac{\lambda_1}{\lambda_0}$ . We also report the average run-time for the processing of a single particle for each method in Table 5.8. It is clear that IMD not only exhibits better performance than MD, it can also be computationally cheaper for a small number of interacting particles.



**Figure 5.23:** Disease Outbreak Anomaly Detection: Best AUC scores of MD and IMD for disease ratios  $\{1.1, 1.3, 1.5, 1.7\}$ .

AUC	SNR = $\lambda_1/\lambda_0$			
	1.1	1.3	1.5	1.7
MD	0.641	0.747	0.791	0.845
IMD	0.722	0.791	0.859	0.934

**Table 5.7:** Disease Outbreak Anomaly Detection: Comparison of MD and IMD for the optimal choice of  $\gamma$  as illustrated in Figures 5.21 and 5.22.

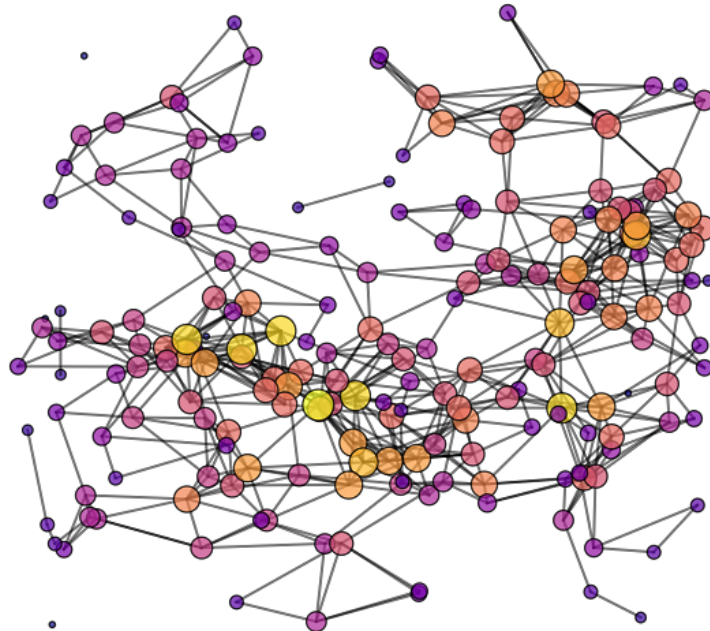
	$N_{iid}$	$N_p$	Iters	Run-time for processing a single sample
MD	5	-	50	2.02 secs
IMD	-	2	50	1.92 secs

**Table 5.8:** Disease Outbreak Anomaly Detection: Average time for the processing of one sample for MD and IMD.

### 5.5.3 Random Geometric Graphs Anomalous Clusters Classification: Numerical Results

In this section we present our results on the performance of Mirror Descent and Interacting Mirror Descent in the case of geometric graphs. We generate uniformly 3-dimensional geometric graphs with  $n$  nodes in the 3D cube  $\mathcal{C}^3 = [-1, 1] \times [-1, 1] \times [-1, 1]$  with the methods presented in [43]. We then produce anomalous clusters by selecting groups of vertices that fall in origin-centered spheres or hyper-ellipsoids with specified radii. We denote with  $(r_x, r_y, r_z)$  the radius of each axis of the sphere or hyper-ellipsoid and with  $K$  the mean number of abnormal vertices in those clusters. In Figure 5.24 is illustrated a random generated graph consisting of  $n = 200$  vertices placed uniformly in  $\mathcal{C}^3$ .

Similarly to the previous section, our aim here is to classify each cluster as abnormal or non-abnormal with the elevated-mean detection method. For  $i = 1, \dots, n$  we assign randomly to anomalous vertices  $x_i \sim \text{Poisson}(\lambda_1)$  and normal vertices  $x_i \sim \text{Poisson}(\lambda_0)$  (we set  $\mathbf{C} = \mathbf{x}\mathbf{x}^T$ ). Below, we investigate the AUC performance of our methods with respect to different values of SNR ratios, increasing number of graph vertices and increasing  $K$  and different number of iterations. All graphs presented here, are generated with radii  $r_x = r_y = r_z = r(K)$  and for all our computations we use a fixed  $\gamma = \sqrt{\frac{1}{500}}$  and run our optimisation for  $x_{\text{samp}} = 40$  samples. Any other parameter is specified in each problem.



**Figure 5.24:** Random Geometric Graphs Anomalous Clusters Classification: 3D Random generated Geometric Graph consisting of 200 vertices. Nodes with higher degrees (more neighbours) are illustrated with a larger radius and lighter colours.

### Performance of MD and IMD vs SNR values

We begin our experiments for the case of random geographic graphs by investigating the performance of Interacting Mirror Descent and Mirror Descent for various values of SNR. To show the scalability of Interacting Mirror Descent we generate a big-sized random geographical graph of  $n = 1000$  vertices with abnormal cluster size  $K = 40$ . We run for  $T = 50$  iterations each framework with base rate  $\lambda_0 = 100$  and  $\lambda_1 = 110, 130, 150$  and  $170$ . In the case of IMD we use 2 interacting particles, while for MD we average the results over 2 iid copies.

We report the AUC performance results for both algorithms in Table 5.9. As observed in the previous section, we observe that both algorithms perform better for higher SNR values. Also, we can discern that for lower SNRs ( $\lambda_1/\lambda_0 = 1.1$ ), IMD outperforms MD, while for higher SNRs ( $\lambda_1/\lambda_0 = 1.3, 1.5, 1.7$ ) the performance of both methods is similar.

AUC	SNR = $\frac{\lambda_1}{\lambda_0}$			
	1.1	1.3	1.5	1.7
IMD	0.930	1.000	1.000	1.000
MD	0.831	0.999	1.000	1.000

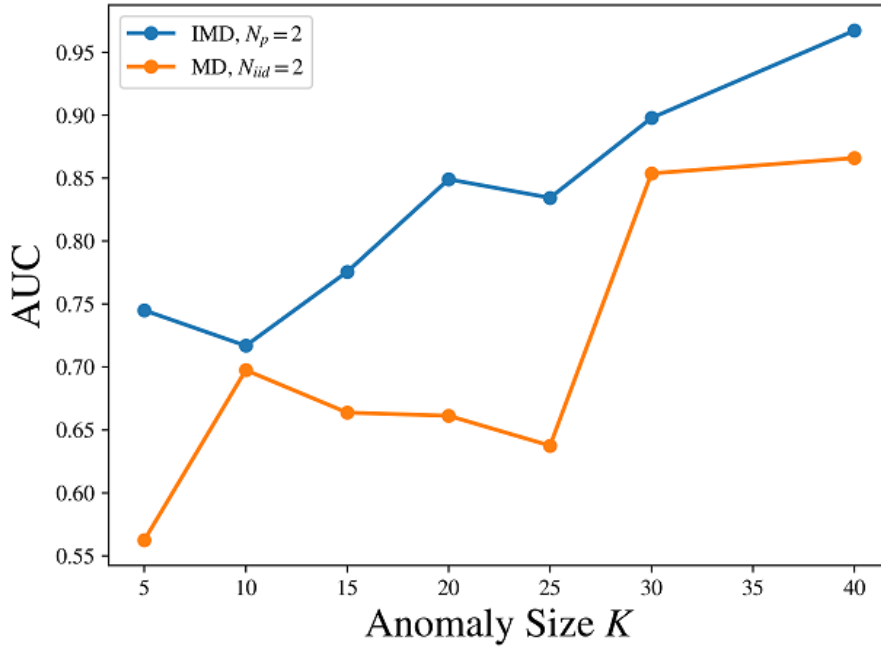
**Table 5.9:** Random Geometric Graphs Anomalous Clusters Classification: Performance of IMD and MD for different SNR values.

### Performance of MD and IMD vs Anomaly size $K$

In this paragraph we investigate the performance of Interacting Mirror Descent and Mirror Descent for various sizes of abnormal clusters  $K$ . We experiment with values of  $K$  between 5 and 40 and illustrate the results in Figure 5.25. We set  $\lambda_0 = 100$  and  $\lambda_1 = 110$  and we generate a random geographic graph consisting of  $n = 1000$  nodes and run each optimiser for  $T = 50$  iterations with fixed learning rate  $\eta = 5$ . In the case of IMD we use 2 interacting particles, while for MD we average the results over 2 iid copies.

We observe that the performance of both algorithms improves as the size of anomaly  $K$  increases. This result is consistent with the results presented in Section 5.2 of [1]. However, here we show that for all anomaly sizes IMD outperforms MD.



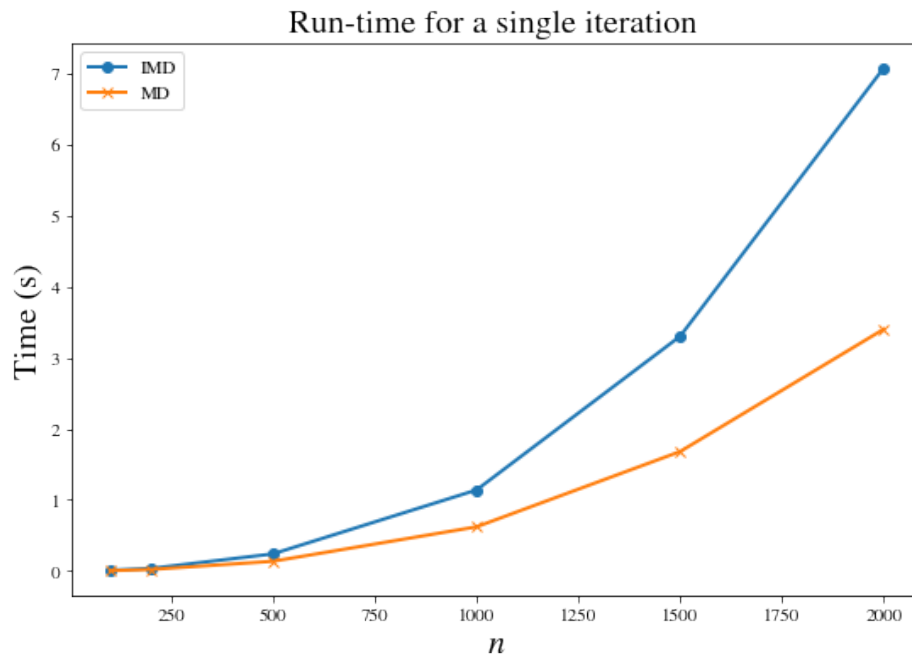


**Figure 5.25:** Random Geometric Graphs Anomalous Clusters Classification: AUC performance of MD and IMD for different anomaly sizes  $K$ .

### Performance of MD and IMD vs graph size $n$

We further examine the time complexity of Mirror Descent and Interacting Mirror Descent algorithms with respect to the size  $n$  of the generated graph  $G$ . We use  $K = 40$  and  $T = 50$ . We plot the average run-time per iteration for various graph sizes in semi-log axes in Figure 5.26. Note that in the case of the latter we use two interacting particles and we average the times with respect to the particles as well.

From the Figure we can observe that the average run-time per iteration for both methods grows exponentially with respect to the size of the graph. Additionally, we observe that for small-sized graphs, both methods perform similarly. However, for larger graph sizes, IMD takes more time to complete a single iteration than MD. Nevertheless, that was expected due to the fact that interaction requires the element-wise multiplication and exponentiation to the  $A_{ij}$ -power as per the update rule shown in Equation 4.30.



**Figure 5.26:** Random Geometric Graphs Anomalous Clusters Classification: Average run-time (in seconds) per iteration for the processing of one sample for various graph sizes  $n$ .

## 5.6 Langevin Diffusions

The last problem we are concerned with in this project is a popular optimisation problem arising in statistical and machine learning, called Markov Chain Monte Carlo (MCMC) which is used for sampling. Specifically, the sampling problem is as follows:

Given  $V : \mathbb{R}^d \rightarrow \mathbb{R}$  a convex, twice differentiable potential function, we want to sample from a negative log-convex target probability distribution  $\pi$  defined on  $\mathbb{R}^d$  such that

$$d\pi = e^{-V(\mathbf{x})} d\mathbf{x}, \quad \mathbf{x} \in \mathbb{R}^d.$$

A widely used algorithm for MCMC sampling is Langevin Diffusions. LD is the solution  $\{\mathbf{x}_t\}_{t=0}^T$  of the following SDE:

$$d\mathbf{x}_t = \nabla V(\mathbf{x}_t) dt + \sqrt{2} d\mathbf{B}_t,$$

where  $\{\mathbf{B}_t\}_{t=0}^T$  are standard Brownian motions in  $\mathbb{R}^d$ . For more information on Langevin Diffusions we refer the reader to [54].

### 5.6.1 Mirror Langevin Diffusion

Now we describe a variant of LD, called Mirror Langevin Diffusion (MLD), which was proposed initially in [59] and utilises Mirror Descent as the optimisation framework. The problem is as follows:

Let  $\Psi : \mathbb{R}^d \rightarrow \mathbb{R}$  be a mirror map as defined earlier in Chapter 3. The MLD then, is the SDE:

$$\begin{aligned} \mathbf{x}_t &= \nabla \Psi^*(\mathbf{y}_t) \\ d\mathbf{y}_t &= -\nabla V(\mathbf{x}_t) dt + \sqrt{2} \cdot \nabla^2 \Psi(\mathbf{x}_t)^{\frac{1}{2}} d\mathbf{B}_t, \end{aligned} \tag{5.12}$$

where  $\Psi^*$  denotes the convex conjugate of  $\Psi$  and  $\nabla^2 \Psi(\mathbf{x}_t)^{\frac{1}{2}}$  is the solution  $\mathbf{A}_t$  of

$$\mathbf{A}_t^2 = \nabla^2 \Psi(\mathbf{x}_t).$$

In [17] the authors suggest picking as a mirror map  $\Psi = V$ , a potential function that meets all assumptions of a mirror map and propose a sampling method called Newton Langevin Diffusion. NLD has the form:

$$\begin{aligned} \mathbf{x}_t &= \nabla V^*(\mathbf{y}_t) \\ d\mathbf{y}_t &= -\nabla V(\mathbf{x}_t) dt + \sqrt{2} \cdot \nabla^2 V(\mathbf{x}_t)^{\frac{1}{2}} d\mathbf{B}_t. \end{aligned} \tag{5.13}$$

We get the Newton Langevin Algorithm (NLA) by utilising the Euler discretization of Equation 5.13:

$$\begin{aligned} \nabla V(\mathbf{x}_{t+1}) &= (1 - \epsilon) \nabla V(\mathbf{x}_t) + \sqrt{2\epsilon} \nabla^2 V(\mathbf{x}_t)^{\frac{1}{2}} \boldsymbol{\zeta}_t \\ \mathbf{x}_{t+1} &= \nabla V^* \circ \nabla V(\mathbf{x}_{t+1}), \quad t \geq 0, \end{aligned} \tag{5.14}$$

where  $\epsilon$  denotes the discretisation parameter (here used as a learning rate) and  $\boldsymbol{\zeta}_t \sim \mathcal{N}(0, \mathbf{I}_d)$ .

### 5.6.2 Interacting Mirror Langevin Diffusion

As an alternative to 5.12 we now propose an interacting particles SDE similar to the previous frameworks of this work, inspired by [9, 10], called **Interacting Mirror Langevin Diffusion** (IMLD) outlined by the SDE:

$$\begin{aligned} \mathbf{x}_t^i &= \nabla \Psi^*(\mathbf{y}_t^i) \\ d\mathbf{y}_t^i &= -\nabla V(\mathbf{x}_t^i) dt + \sum_{j=1}^{N_p} A_{ij}(\mathbf{y}_t^j - \mathbf{y}_t^i) + \sqrt{2} \cdot \nabla^2 \Psi(\mathbf{x}_t^i)^{\frac{1}{2}} d\mathbf{B}_t^i \\ i &= 1, 2, \dots, N_p, \end{aligned} \quad (5.15)$$

where  $N_p$  denotes the number of interacting particles and  $\mathbf{A}$  is a doubly stochastic matrix in  $\mathbb{R}^{N_p \times N_p}$ .

Similarly to [17], we use as a mirror map the potential function  $V$  and propose the Interacting Newton Langevin Diffusion (INLD), i.e.:

$$\begin{aligned} \mathbf{x}_t^i &= \nabla V^*(\mathbf{y}_t^i) \\ d\mathbf{y}_t^i &= -\nabla V(\mathbf{x}_t^i) dt + \sum_{j=1}^{N_p} A_{ij}(\mathbf{y}_t^j - \mathbf{y}_t^i) + \sqrt{2} \cdot \nabla^2 V(\mathbf{x}_t^i)^{\frac{1}{2}} d\mathbf{B}_t^i \\ i &= 1, 2, \dots, N_p. \end{aligned} \quad (5.16)$$

The Euler ization of Equation 5.16 is gives our proposal algorithm, the Interacting Newton Langevin Algorithm (INLA):

$$\begin{aligned} \nabla V(\mathbf{x}_{t+1}^i) &= (1 - 2\epsilon) \nabla V(\mathbf{x}_t^i) + \epsilon \sum_{j=1}^{N_p} A_{ij} \mathbf{x}_t^j + \sqrt{2\epsilon} \nabla^2 V(\mathbf{x}_t^i)^{\frac{1}{2}} \boldsymbol{\zeta}_t^i \\ \mathbf{x}_{t+1}^i &= \nabla V^* \circ \nabla V(\mathbf{x}_{t+1}^i), \quad i = 1, \dots, N_p, \quad t \geq 0, \end{aligned} \quad (5.17)$$

where again,  $\epsilon$  denotes the discretization parameter and  $\boldsymbol{\zeta}_t^i \sim \mathcal{N}(0, \mathbf{I}_d)$ .

Our results from the previous sections, suggest that interaction improves the performance of the optimisation problem. Hence, we expect that our numerical experiments in this section will show that our method INLD outperforms NLD. Note that in contrast with the previous problems, the optimisation problem here is unconstrained as we let  $\mathbf{x}$  to lie in the whole space, i.e.  $\mathbf{x} \in \mathbb{R}^d$ .

### 5.6.3 Numerical Results

In this paragraph we provide results for our proposed method Interacting Newton Langevin Diffusion and we compare it with Newton Langevin Diffusion from [17]. Similarly to Section E.1 of [17], we use as a potential function

$$V(\mathbf{x}) = \frac{1}{2} \mathbf{x}^T \Sigma^{-1} \mathbf{x}, \quad \mathbf{x} \in \mathbb{R}^d,$$

which is utilised to sample from a Gaussian Distribution. The matrix  $\Sigma$  denotes the scatter matrix of the desired distribution and it holds that  $\Sigma = \Sigma^T$  and  $\Sigma \succ \mathbf{0}$ .

Note that  $V$  is twice differentiable with gradient and Hessian:

$$\nabla V(\mathbf{x}) = \Sigma^{-1} \mathbf{x}, \quad \nabla^2 V(\mathbf{x}) = \Sigma^{-1}.$$

Since  $\Sigma \succ \mathbf{0}$ , it holds that  $\Sigma^{-1} \succ \mathbf{0}$ , and therefore  $V$  is a convex function defined on a convex domain ( $\mathbb{R}^d$ ). It's also easy to proof that

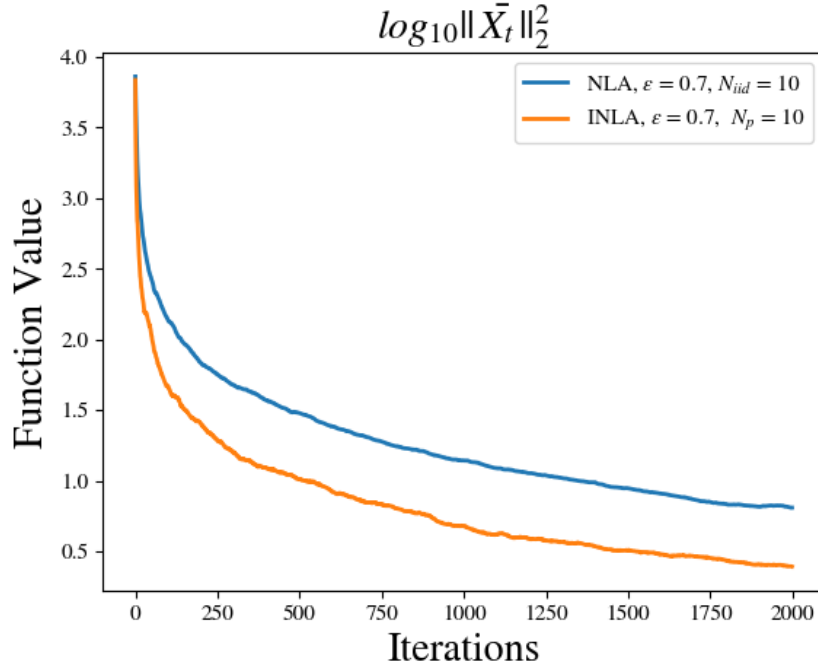
$$V^*(\mathbf{z}) = \frac{1}{2} \mathbf{z}^T \Sigma^{-1} \mathbf{z}, \quad \nabla V^*(\mathbf{z}) = \Sigma^{-1} \mathbf{z}.$$

We quantify the performance of each algorithm similarly to [17], by computing the  $\log_{10}$  of running estimates at each iteration  $t = 1, \dots, T$  of the mean  $\bar{\mathbf{x}}_t = \frac{1}{t+1} \sum_{s=0}^t \mathbf{x}_s$  and by computing the  $\log_{10}$  of the relative squared error of the covariance (scatter) matrix  $\frac{\|\hat{\Sigma} - \Sigma\|_F^2}{\|\Sigma\|_F^2}$ , where  $\hat{\Sigma} = \frac{1}{t+1} \sum_{s=0}^t \text{cov}(\bar{\mathbf{x}}_s)$ .

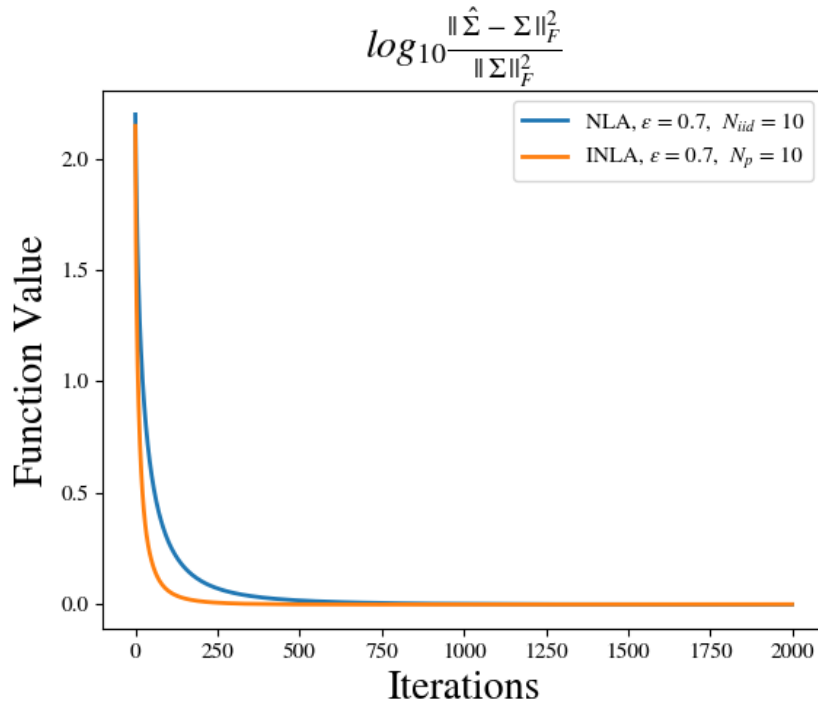
We plot the results for each metric in Figures 5.27 and 5.28, respectively, in the case

$$\text{of } \Sigma = \begin{bmatrix} 1 & 0 & \cdots & 0 \\ 0 & 2 & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & d \end{bmatrix} \text{ and } d = 100.$$

In the case of NLA we average the results over 10 independent and identical runs and in the case of INLA we use 10 interacting particles. For both algorithms we set a constant discretization parameter  $\epsilon = 0.7$ .



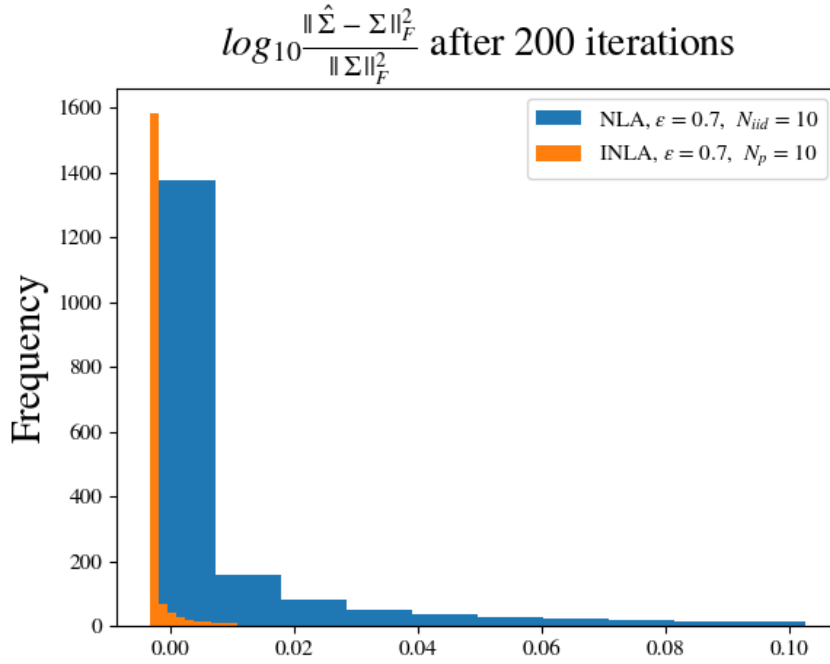
**Figure 5.27:** Mirror Langevin Diffusion: Convergence of the squared distance of the running mean  $\bar{x}_t$  from  $\mathbf{0}$  for  $d = 100$  in the case of NLA averaged over  $N_{iid} = 10$  runs and INLA with  $N_p = 10$ .



**Figure 5.28:** Mirror Langevin Diffusion: Convergence of the relative squared distance of the running estimate of  $\hat{\Sigma}$  from  $\Sigma$  for  $d = 100$  in the case of NLA averaged over  $N_{iid} = 10$  runs and INLA with  $N_p = 10$ .

Observing Figures 5.27 and 5.28 we can infer that INLA achieves a much greater performance than NLA, as the former manages to generate samples from the desired distribution significantly faster. Also, the speed of convergence of our proposed algorithm is greater than the speed of NLA.

Additionally, we plot in Figure 5.29 the histogram of the relative squared error and we present in Table 5.10 the variance for the two methods after 200 iterations. We then, compare the two methods. As expected from our previous numerical results, interaction aids in variance reduction. INLA attains a much lower variance than NLA as observed from Figure 5.29 and Table 5.10.



**Figure 5.29:** Mirror Langevin Diffusion: Histogram of the relative squared distance of the running estimate of  $\hat{\Sigma}$  from  $\Sigma$  for  $d = 100$  in the case of NLA averaged over  $N_{iid} = 10$  runs and INLA with  $N_p = 10$  after convergence ( $T_0 = 200$ ).

Variance after convergence	
NLA	0.000359875034
INLA	4.10370467e-06

**Table 5.10:** Mirror Langevin Diffusion: Variance of samples of  $\log_{10} \frac{\|\hat{\Sigma} - \Sigma\|_F^2}{\|\Sigma\|_F^2}$  after convergence ( $T_0 = 200$ ) for NLA and INLA.

# Chapter 6

## Conclusion and Future Work

### 6.1 Conclusion

In this individual project, we have reviewed methods for optimising a convex function as an objective. We have applied Gradient Descent, Mirror Descent and Interacting Mirror Descent on six problems and we have compared their performance. Utilising the latter method, we have solved numerically for the first time two problems from the existing related literature (see subsections 5.4 and 5.5).

Additionally, in Section 5.6 we have proposed a novel Langevin diffusion-based method for sampling from a target distribution called Interacting Newton Langevin Algorithm. Our proposed method utilises interaction between particles as a variant of Newton Langevin Algorithm proposed in the related literature. Our numerical results suggest that our algorithm outperforms the original method as it establishes better convergence performance.

Moreover, throughout Chapter 5, we have analysed the benefits of using Mirror Descent and its variant Interacting Mirror Descent, over the Projected Gradient Descent. Our numerical experiments indicate that the first two surpass the latter algorithm.

We have also shown that interaction in IMD can speed up convergence of the optimiser to the optimum solution. Furthermore, we have seen that interaction can result in a reduction of noise even in cases such as Mini-batch optimisation. We have shown that using more interacting particles can decrease the variance of the samples of the loss function significantly, performing almost as good as full-batch optimisation. We also remark that we have ran experiments with fixed and decreasing learning rate. In both cases, MD and IMD exhibit better convergence properties than GD.



## 6.2 Future Work

We conclude this work by discussing several interesting directions for further future implementation and research. The results shown in this project spark the idea for the implementation of some possible future directions:

- For our experiments performed in Chapter 5, in the case of Interacting Mirror Descent we used as an interaction matrix a mean-field interaction, i.e.  $A_{ij} = \frac{1}{N_p}$  where  $N_p$  denotes the number of interacting particles. We therefore allowed all particles to interact in the same way. An extension would be to study the convergence properties of the IMD method with respect to  $\mathbf{A}$  by experimenting different interaction matrices by allowing some particles to interact more with some particles and less with others.
- A future direction could be the implementation of the Mirror Descent algorithm with interacting particles on a famous problem called Langevin Dynamics for Latent Dirichlet Allocation (LDA). This is a similar setup to the last problem of this project (Langevin Diffusions) and it is usually applied for topic modelling problems. Such examples are [28, 41] where the authors utilise Stochastic projected Gradient Descent and Stochastic Mirror Descent on the probability simplex.
- In Section 5.6 of Chapter 5 we proposed a new algorithm called Interacting Mirror Langevin Diffusion. For our numerical results we used as a potential function  $V(\mathbf{x}) = \frac{1}{2} \mathbf{x}^T \Sigma^{-1} \mathbf{x}$ . Another future direction would be to apply our method for different choices of  $V$  that aim in sampling from ill-conditioned Gaussian distributions, such as  $V(\mathbf{x}) = \frac{1}{2} (\mathbf{x}^T \Sigma^{-1} \mathbf{x})^\alpha$ ,  $\alpha \in (0, 1)$ , similarly to [17] where they use  $\alpha = 3/4$ .
- Most of problems in ML and DL are non-convex and quite often optimisers get stuck in local optima. Stochasticity can usually help escaping such situations. A possible direction would be the theoretical and experimental investigation of (Stochastic) Mirror Descent and (Stochastic) Interacting Mirror Descent for the optimisation of non-convex constrained problems.

# Bibliography

- [1] Cem Aksoylar, Lorenzo Orecchia, and Venkatesh Saligrama. Connected subgraph detection with mirror descent on SDPs. In Doina Precup and Yee Whye Teh, editors, *Proceedings of the 34th International Conference on Machine Learning*, volume 70 of *Proceedings of Machine Learning Research*, pages 51–59, International Convention Centre, Sydney, Australia, 06–11 Aug 2017. PMLR. pages 6, 9, 56, 57, 61
- [2] Jürgen Audretsch. *The von Neumann Entropy and Quantum Information*, chapter 6, pages 103–113. John Wiley & Sons, Ltd, 2008. pages 19
- [3] S. Bahmani and B. Raj. A unifying analysis of projected gradient descent forp-constrained least squares. *Applied and Computational Harmonic Analysis*, 34(3):366–378, May 2013. pages 5
- [4] Amir Beck. *First-Order Methods in Optimization*. Society for Industrial and Applied Mathematics, Philadelphia, PA, 2017. pages 20, 39
- [5] Amir Beck and Marc Teboulle. Mirror descent and nonlinear projected subgradient methods for convex optimization. *Operations Research Letters*, 31:167–175, 05 2003. pages 3, 6
- [6] M. Bierlaire, Ph.L. Toint, and D. Tuytens. On iterative algorithms for linear least squares problems with bound constraints, 1995. pages 39
- [7] Jaya Bishwal. Estimation in interacting diffusions: Continuous and discrete sampling. *Applied Mathematics*, 02:1154–1158, 01 2011. pages 6
- [8] Eliot Bolduc, George Knee, Erik Gauger, and Jonathan Leach. Projected gradient descent algorithms for quantum state tomography, 2016. pages 5
- [9] Anastasia Borovykh, Nikolas Kantas, Panos Parpas, and Grigorios Pavliotis. To interact or not? the convergence properties of interacting stochastic mirror descent, 07 2020. pages 4, 65
- [10] Anastasia Borovykh, Nikolas Kantas, Panos Parpas, and Grigorios A. Pavliotis. On stochastic mirror descent with interacting particles: convergence properties and variance reduction, 2020. pages 4, 7, 8, 20, 23, 26, 39, 42, 65

- 
- [11] Léon Bottou. Large-scale machine learning with stochastic gradient descent. In Yves Lechevallier and Gilbert Saporta, editors, *Proceedings of the 19th International Conference on Computational Statistics (COMPSTAT'2010)*, pages 177–187, Paris, France, August 2010. Springer. pages 44
  - [12] Milan Bradonjić, Aric Hagberg, and Allon G. Percus. Giant component and connectivity in geographical threshold graphs. In *Proceedings of the 5th International Conference on Algorithms and Models for the Web-Graph, WAW'07*, page 209–216, Berlin, Heidelberg, 2007. Springer-Verlag. pages 35
  - [13] Sébastien Bubeck. Convex optimization: Algorithms and complexity, 2014. pages 3, 19, 22, 23, 25
  - [14] Rene Carmona, Jean-Pierre Fouque, and Douglas Vestal. Interacting particle systems for the computation of rare credit portfolio losses. *Finance and Stochastics*, 13:613–633, 09 2009. pages 6
  - [15] Xi Chen, Qihang Lin, and Javier Pena. Optimal regularized dual averaging methods for stochastic optimization. In F. Pereira, C. J. C. Burges, L. Bottou, and K. Q. Weinberger, editors, *Advances in Neural Information Processing Systems 25*, pages 395–403. Curran Associates, Inc., 2012. pages 6
  - [16] Xiaohui Chen. Maximum likelihood estimation of potential energy in interacting particle systems from single-trajectory data, 2020. pages 6
  - [17] Sinho Chewi, Thibaut Le Gouic, Chen Lu, Tyler Maunu, Philippe Rigollet, and Austin J. Stromme. Exponential ergodicity of mirror-langevin diffusions, 2020. pages 9, 64, 65, 66, 70
  - [18] John C. Duchi, Alekh Agarwal, Mikael Johansson, and Michael I. Jordan. Ergodic mirror descent, 2011. pages 6
  - [19] Dan Garber. On the convergence of stochastic gradient descent with low-rank projections for convex low-rank matrix problems, 2020. pages 5
  - [20] Josselin Garnier, George Papanicolaou, and Tzu-Wei Yang. Large deviations for a mean field model of systemic risk, 2012. pages 6
  - [21] Rong Ge, Furong Huang, Chi Jin, and Yang Yuan. Escaping from saddle points — online stochastic gradient for tensor decomposition. In Peter Grünwald, Elad Hazan, and Satyen Kale, editors, *Proceedings of The 28th Conference on Learning Theory*, volume 40 of *Proceedings of Machine Learning Research*, pages 797–842, Paris, France, 03–06 Jul 2015. PMLR. pages 44
  - [22] Ian Goodfellow, Yoshua Bengio, and Aaron Courville. *Deep Learning*. MIT Press, 2016. <http://www.deeplearningbook.org>. pages 5
  - [23] Harshit Gupta, Kyong Hwan Jin, Ha Q. Nguyen, Michael T. McCann, and Michael Unser. Cnn-based projected gradient descent for consistent ct image reconstruction. *IEEE Transactions on Medical Imaging*, 37(6):1440–1453, Jun 2018. pages 5

- [24] Nicholas Harvey. Machine learning theory: Mirror descent, 2018. pages 22
- [25] Adrian Hauswirth, Saverio Bolognani, Gabriela Hug, and Florian Dörfler. Projected gradient descent on riemannian manifolds with applications to online power system optimization, 09 2016. pages 5
- [26] Jean-Baptiste Hiriart-Urruty and Claude Lemaréchal. Fundamentals of convex analysis / j.b. hiriart-urruty, c. lemaréchal., 01 2001. pages 13
- [27] Tito Homem-de Mello. On rates of convergence for stochastic optimization problems under non-independent and identically distributed sampling. *SIAM Journal on Optimization*, 19(2):524–551, 2008. pages 26
- [28] Ya-Ping Hsieh, Ali Kavis, Paul Rolland, and Volkan Cevher. Mirrored langevin dynamics, 2018. pages 6, 70
- [29] Gautam Kunapuli and Jude Shavlik. Mirror descent for metric learning: A unified approach. In Peter A. Flach, Tijl De Bie, and Nello Cristianini, editors, *Machine Learning and Knowledge Discovery in Databases*, pages 859–874, Berlin, Heidelberg, 2012. Springer Berlin Heidelberg. pages 6
- [30] Guanghai Lan, Arkadi Nemirovski, and Alexander Shapiro. Validation analysis of mirror descent stochastic approximation method. *Mathematical Programming*, 134:1–34, 09 2011. pages 6
- [31] Måns Larsson, Anurag Arnab, Fredrik Kahl, Shuai Zheng, and Philip Torr. A projected gradient descent method for crf inference allowing end-to-end training of arbitrary pairwise potentials, 2017. pages 5
- [32] Jie Liu and Martin Takac. Projected semi-stochastic gradient descent method with mini-batch scheme under weak strong convexity assumption, 2016. pages 5
- [33] S. Liu, X. Li, P. Chen, J. Haupt, and L. Amini. Zeroth-order stochastic projected gradient descent for nonconvex optimization. In *2018 IEEE Global Conference on Signal and Information Processing (GlobalSIP)*, pages 1179–1183, 2018. pages 5
- [34] Mehrdad Mahdavi, Tianbao Yang, Rong Jin, Shenghuo Zhu, and Jinfeng Yi. Stochastic gradient descent with only one projection. In F. Pereira, C. J. C. Burges, L. Bottou, and K. Q. Weinberger, editors, *Advances in Neural Information Processing Systems 25*, pages 494–502. Curran Associates, Inc., 2012. pages 5
- [35] Naoki Masuda, Hiroyoshi Miwa, and Norio Konno. Geographical threshold graphs with small-world and scale-free properties. *Phys. Rev. E*, 71:036108, Mar 2005. pages 35
- [36] P. Mertikopoulos, E. V. Belmega, R. Negrel, and L. Sanguinetti. Distributed stochastic optimization via matrix exponential learning. *IEEE Transactions on Signal Processing*, 65(9):2277–2290, 2017. pages 6

- 
- [37] Panayotis Mertikopoulos and Mathias Staudigl. On the convergence of gradient-like flows with noisy gradient input. *SIAM Journal on Optimization*, 28:163–197, 01 2018. pages 6
- [38] Angelia Nedich and Soomin Lee. On stochastic subgradient mirror-descent algorithm with weighted averaging, 2013. pages 6
- [39] Arkadi Nemirovski. Prox-method with rate of convergence  $o(1/t)$  for variational inequalities with lipschitz continuous monotone operators and smooth convex-concave saddle point problems. *SIAM Journal on Optimization*, 15:229–251, 01 2004. pages 6
- [40] A.S. Nemirovskii and D.B. Yudin. *Problem Complexity and Method Efficiency in Optimization*. A Wiley-Interscience publication. Wiley, 1983. pages 3, 6
- [41] Sam Patterson and Yee Whye Teh. Stochastic gradient riemannian langevin dynamics on the probability simplex. In C. J. C. Burges, L. Bottou, M. Welling, Z. Ghahramani, and K. Q. Weinberger, editors, *Advances in Neural Information Processing Systems 26*, pages 3102–3110. Curran Associates, Inc., 2013. pages 70
- [42] F. Pedregosa, G. Varoquaux, A. Gramfort, V. Michel, B. Thirion, O. Grisel, M. Blondel, P. Prettenhofer, R. Weiss, V. Dubourg, J. Vanderplas, A. Passos, D. Cournapeau, M. Brucher, M. Perrot, and E. Duchesnay. Scikit-learn: Machine learning in Python. *Journal of Machine Learning Research*, 12:2825–2830, 2011. pages 50
- [43] D.M.S.M. Penrose, M. Penrose, and Oxford University Press. *Random Geometric Graphs*. Oxford studies in probability. Oxford University Press, 2003. pages 60
- [44] M. Raginsky and J. Bouvrie. Continuous-time stochastic mirror descent on a network: Variance reduction, consensus, convergence. In *2012 IEEE 51st IEEE Conference on Decision and Control (CDC)*, pages 6793–6800, 2012. pages 22
- [45] Clement Royer. Optimization for machine learning: Basics of constrained optimization. Dauphine University Paris -M2 IASD, dec 2019. pages 3, 5
- [46] Shai Shalev-Shwartz. *Online Learning and Online Convex Optimization*. Now Publishers Inc, 01 2011. pages 6
- [47] R. J. Smeed. Studies in the Economics of Transportation. *The Economic Journal*, 67(265):116–118, 03 1957. pages 33
- [48] Tom Tirer and Raja Giryes. On the convergence rate of projected gradient descent for a back-projection based objective, 2020. pages 5
- [49] Yann Traonmilin, Jean-François Aujol, and Arthur Leclaire. Projected gradient descent for non-convex sparse spike estimation, 2020. pages 5
- [50] Weiran Wang and Miguel Á. Carreira-Perpiñán. Projection onto the probability simplex: An efficient algorithm with a simple proof, and an application, 2013. pages 5, 8
-

- [51] Weiran Wang and Miguel Á. Carreira-Perpiñán. The laplacian k-modes algorithm for clustering, 2014. pages 5, 8, 49, 50
- [52] Fanyou Wu, Rado Gazo, Eva Haviarova, and Bedrich Benes. Efficient project gradient descent for ensemble adversarial attack, 2019. pages 5
- [53] Lin Xiao. Dual averaging methods for regularized stochastic learning and online optimization. *Journal of Machine Learning Research*, 11(88):2543–2596, 2010. pages 6
- [54] T. Xifara, C. Sherlock, S. Livingstone, S. Byrne, and M. Girolami. Langevin diffusions and the metropolis-adjusted langevin algorithm. *Statistics & Probability Letters*, 91:14–19, Aug 2014. pages 64
- [55] Pan Xu, Tianhao Wang, and Quanquan Gu. Accelerated stochastic mirror descent: From continuous-time dynamics to discrete-time algorithms. In Amos Storkey and Fernando Perez-Cruz, editors, *Proceedings of the Twenty-First International Conference on Artificial Intelligence and Statistics*, volume 84 of *Proceedings of Machine Learning Research*, pages 1087–1096, Playa Blanca, Lanzarote, Canary Islands, 09–11 Apr 2018. PMLR. pages 6
- [56] Sheng Xu, Zhou Fan, and Sahand Negahban. Tree-projected gradient descent for estimating gradient-sparse parameters on graphs, 2020. pages 5
- [57] Yao-Liang Yu. The strong convexity of von neumann’s entropy, 2013. pages 17, 18, 19
- [58] Jiawei Zhang. Gradient descent based optimization algorithms for deep learning models training, 2019. pages 5
- [59] Kelvin Shuangjian Zhang, Gabriel Peyré, Jalal Fadili, and Marcelo Pereyra. Wasserstein control of mirror langevin monte carlo, 2020. pages 64
- [60] B. Zhou and F. Chen. Graph-structured sparse optimization for connected sub-graph detection. In *2016 IEEE 16th International Conference on Data Mining (ICDM)*, pages 709–718, 2016. pages 56
- [61] Zhengyuan Zhou, Panayotis Mertikopoulos, Nicholas Bambos, Stephen Boyd, and Peter Glynn. On the convergence of mirror descent beyond stochastic convex programming, 2017. pages 6
- [62] Martin Zinkevich, Markus Weimer, Alexander Smola, and Lihong Li. Parallelized stochastic gradient descent. *Advances in Neural Information Processing Systems*, 23:2595–2603, 01 2010. pages 26

# Appendices

# Appendix A

## Code

The interested reader can find the Gitlab repository of this individual project in the link below:

<https://gitlab.doc.ic.ac.uk/gy615/mirror-descent-optimisation> .

In our Gitlab repository we provide all the code written for the implementation of all six optimisation problems solved in this project along with the results presented throughout Chapter 5. All code was written in Python3. A list of package-dependencies is also provided needed to run the code on a local machine.

The repository contains six project folders, one for each optimisation problem we solve in Chapter 5. The main functions that solve each optimisation problem with the methods afformentioned in the project are listed below:

- The Traffic Assignment Problem: `\run_opt.py`
- Linear System: `\run_opt.py`
- Linear System: Mini-batch Optimisation: `\run_opt.py`
- The Laplacian K-Modes for Clustering: `\run_opt.py`
- Optimisation on SDPs: `\startOpt\startOpt.py`
- Langevin Diffusions: `\run_opt.py`

The repository also includes all numerical results for the problems presented in Chapter 5 in the `\saved_items` subfolder of each problem, as well as all illustrations `\figures` subfolders.



# Appendix B

## Ethical Considerations

We declare that the work performed here poses no ethical, social or legal considerations, since we are concerned with the algorithmic performance and numerical results on mathematical problems. Additionally, all data used for this project were randomly generated and/or did not include any kind of personal data.

However, optimisation is a key component for every problem in artificial intelligence and machine learning and finding the best optimiser is always an open problem in every research field. Thus, the case of such algorithm being used and applied in scenarios that would cause harm to any extent to the environment or to any natural person, etc or for dual or medicinal use would not be improbable. Nevertheless, such cases are far from probable to happen and therefore we assume that they are no considerations of such.

In B.1 we review some ethical, social or legal considerations. I declare that everything stated on this table is true and correct and was completed solemnly with my own free will.

	Yes	No
<b>Section 1: HUMAN EMBRYOS/FOETUSES</b>		
Does your project involve Human Embryonic Stem Cells?		✓
Does your project involve the use of human embryos?		✓
Does your project involve the use of human foetal tissues / cells?		✓
<b>Section 2: HUMANS</b>		
Does your project involve human participants?		✓
<b>Section 3: HUMAN CELLS / TISSUES</b>		
Does your project involve human cells or tissues? (Other than from “Human Embryos/Foetuses” i.e. Section 1)?		✓
<b>Section 4: PROTECTION OF PERSONAL DATA</b>		
Does your project involve personal data collection and/or processing?		✓
Does it involve the collection and/or processing of sensitive personal data (e.g. health, sexual lifestyle, ethnicity, political opinion, religious or philosophical conviction)?		✓
Does it involve processing of genetic information?		✓

Does it involve tracking or observation of participants? It should be noted that this issue is not limited to surveillance or localization data. It also applies to Wan data such as IP address, MACs, cookies etc.		✓
Does your project involve further processing of previously collected personal data (secondary use)? For example Does your project involve merging existing data sets?		✓
<b>Section 5: ANIMALS</b>		
Does your project involve animals?		✓
<b>Section 6: DEVELOPING COUNTRIES</b>		
Does your project involve developing countries?		✓
If your project involves low and/or lower-middle income countries, are any benefit-sharing actions planned?		✓
Could the situation in the country put the individuals taking part in the project at risk?		✓
<b>Section 7: ENVIRONMENTAL PROTECTION AND SAFETY</b>		
Does your project involve the use of elements that may cause harm to the environment, animals or plants?		✓
Does your project deal with endangered fauna and/or flora /protected areas?		✓
Does your project involve the use of elements that may cause harm to humans, including project staff?		✓
Does your project involve other harmful materials or equipment, e.g. high-powered laser systems?		✓
<b>Section 8: DUAL USE</b>		
Does your project have the potential for military applications?		✓
Does your project have an exclusive civilian application focus?		✓
Will your project use or produce goods or information that will require export licenses in accordance with legislation on dual use items?		✓
Does your project affect current standards in military ethics – e.g., global ban on weapons of mass destruction, issues of proportionality, discrimination of combatants and accountability in drone and autonomous robotics developments, incendiary or laser weapons?		✓
<b>Section 9: MISUSE</b>		
Does your project have the potential for malevolent/criminal/terrorist abuse?		✓
Does your project involve information on/or the use of biological-, chemical-, nuclear/radiological-security sensitive materials and explosives, and means of their delivery?		✓
Does your project involve the development of technologies or the creation of information that could have severe negative impacts on human rights standards (e.g. privacy, stigmatization, discrimination), if misapplied?		✓

Does your project have the potential for terrorist or criminal abuse e.g. infrastructural vulnerability studies, cybersecurity related project?		✓
<b>SECTION 10: LEGAL ISSUES</b>		
Will your project use or produce software for which there are copyright licensing implications?		✓
Will your project use or produce goods or information for which there are data protection, or other legal implications?		✓
<b>SECTION 11: OTHER ETHICS ISSUES</b>		
Are there any other ethics issues that should be taken into consideration?		✓

**Table B.1:** Imperial College London Ethics Checklist