# Min-max optimisation in the continuous time limit

Luca Bracone

April 22, 2022

# 1 Introduction

In numerical optimisation, one aims to find the minimizer

$$\min_x l(x)$$

of a differentiable scalar field (also known as "loss function") $l : \mathbb{R}^n \to \mathbb{R}$. The quintessential algorithm to solve this problem is gradient descent (GD) which given a step size $s \in \mathbb{R}$ computes a sequence of iterates

$$x_{k+1} = x_k + s\nabla l(x_k)$$

whose convergence properties are suitable for a wide range of applications and have been extensively studied by thae optimization commmunity.

Nowadays, there are problems in which we would like instead to calculate equilibria or "saddle points" of our scalar field. We change the problem in the following way: the loss $l : \mathbb{R}^n \times \mathbb{R}^m \to R$ now takes two arguments $x$ and $y$ and we would like to know which are the values so that

$$\min_x \max_y l(x, y)$$

is reached.

One of the most natural approaches we could try to solve this new problem would be to adapt gradient descent in the following way

$$x_{n+1} = x_n - s\nabla_x l(x_n, y_n)$$
$$y_{n+1} = y_n + s\nabla_y l(x_n, y_n)$$

to obtain gradient descent-ascent (GDA). However we can show that even on a simple example GDA fails to find a saddle point.

*Example* 1. Take $l(x, y) = xy$ with $x, y \in \mathbb{R}$. Then $l$ has a saddle point at the origin, but GDA gives

$$\begin{bmatrix} x_{n+1} \\ y_{n+1} \end{bmatrix} = \begin{bmatrix} 1 & -s \\ s & 1 \end{bmatrix} \begin{bmatrix} x_n \\ y_n \end{bmatrix}.$$

We can see that the norm of the new iterate is $\sqrt{1 + s^2}$ times the old one, and so GDA diverges.

For the rest of this report we write $z = (x, y)$ the concatenation of $x$ and $y$ and $\nabla l(z) = (-\nabla_x l(x, y), \nabla_y l(x, y))$ so that GDA can be written simply as $z_{n+1} = z_n + s\nabla l(z_n)$ Given that, many different algorithms have been proposed to find saddle points. We give here some of the most popular ones:

- (EG) Extra gradient: $z_{n+1} = z_n + s\nabla l(z_n + s\nabla l(z_n))$

- (OGDA) Optimistic gradient descent-ascent: $z_{n+1} = z_n - 2s\nabla l(z_n) + s\nabla l(z_{n-1})$

These algorithms have been shown to converge and so we would like to understand if there is a fundamental difference that causes the difference in behaviour.

<div style="float:right; border:1px solid; background:orange; padding:4px;">Explain in what conditions, cite?</div>

## 1.1 Deriving continuous time dynamics from discrete time algorithms

Following the paper of [SBC16], there has been a lot of attention towards deriving differential equations $\dot{Z}(t) = f(Z(t), t)$ whose trajectory follows a given discrete time algorithms. In this section we show that by simply letting the timestep $s \to 0$ all of the algorithms presented so far yield the same ODE, gradient flow (GF). The main idea is to assume that the discrete algorithm is obtained from a continuous curve at fixed length times

$$z_n = Z(t_n)$$

where $t_n = t_0 + ns$. Using this, algorithms of the form

$$z_{n+1} = z_n + s\varphi(z_n)$$

can be rewritten as

$$\frac{Z(t_n + s) - Z(t_n)}{s} = \varphi(Z(t_n)).$$

Finally, taking limits $s \to 0$ gives

$$\dot{Z}(t) = \varphi(Z(t)),$$

where we dropped the index from $t_n$.

*Example* 2. Given that in GDA $\varphi(z_n)$ takes the form $\nabla l(z_n)$ we obtain a continuous time approximation

$$\dot{Z}(t) = \nabla l(Z(t)),$$

which we call gradient flow.

*Example* 3. For EG, the update step gives

$$\frac{Z(t + s) - Z(t)}{s} = \nabla l\big(Z(t) + s\nabla l(Z(t))\big),$$

and taking the limits gives gradient flow again.

<div style="float:right; border:1px solid; background:orange; padding:4px;">todo: OGDA</div>

Given that GDA diverges even in simple cases and that the other two algorithms are known to converge in a large class of situations, we are quite dissatisfied in doscovering that they present the same continuous time dynamics. In Chapter 2 we will present a new method that makes use of modified differential equations to obtain continuous time approximations that more closely follow their discrete time counterparts.

<div style="float:right; border:1px solid; background:orange; padding:4px;">todo: natural descent from mirror descent</div>

# 2 Deriving ODEs with the $O(s^r)$-resolution framework

## 2.1 The process

In light of [Lu21] we present here a method to obtain differential equations from discrete algorithms. We assume that the velocity depends on the current position in the following way

$$\dot{Z}(t) = f_0(Z(t)) + sf_1(Z(t)) + s^2 f_2(Z(t)) + \dots$$

We want to find what are the $f_i$ functions. To alleviate the notation we will write $Z$ instead of $Z(t)$. We obtain a Taylor expansion

$$Z(t+s) = Z + s(f_0(Z) + sf_1(Z) + \dots)$$
$$+ \frac{s^2}{2}(f_0'(Z) + sf_1'(Z) + \dots)(f_0(Z) + sf_1(Z) + \dots) + \dots + o(s^r)$$

On the other hand, given a numerical algorithm $z_{n+1} = \varphi(z_n, s)$ we can also express it as a Taylor expansion on $s$,

$$\varphi(z, s) = \varphi(z, 0) + s\partial_s\varphi(z, s) + \frac{s^2}{2}\partial_{ss}\varphi(z, s) + \dots + o(s^r).$$

In order to obtain $Z(t_n) = z_n$ like the intuition in the introduction, we need to have $\varphi(z, s) = Z(t + s)$. By equating the coefficients in $s$ for both polynomials we find that

$$f_0 = \partial_s\varphi$$
$$2f_1 = \partial_{ss}\varphi - f_0'f_0,$$

and the other coefficients are found in the same manner.

*Example* 4. With GDA we find that $\partial_s\varphi(z, s) = \nabla l(z)$ and that $\partial_{ss}\varphi(z, s) = 0$. So, the first two coefficients are

$$f_0 = \nabla l(z)$$
$$f_1 = -\frac{1}{2}\nabla^2 l(z)\nabla l(z)$$

*Example* 5. With EG the update rule is $\varphi(z, s) = z + s\nabla l(z + s\nabla l(z))$. By Taylor expansion of the second term, we find that this expression equals

$$z + s\nabla l(z) + s^2\nabla^2 l(z)\nabla l(z) + \frac{s^3}{2}\nabla^2 l(z)\nabla l(z)\nabla l(z).$$

By equating the two polynomials as usual we find that

$$f_0 = \nabla l(z)$$

$$f_1 = \frac{1}{2}\nabla^2 l(z)\nabla l(z)$$

## 2.2 Example: Mirror descent $O(s^r)$-resolution ODE

Let $\psi : \mathbb{R}^n \to \mathbb{R}$ be what we call a mirror map. Then, mirror descent uses as update rule

$$g(z, s) = \nabla\psi^{-1}(\nabla\psi(z) + s\nabla l(z)).$$

We denote $\alpha_n(z) = \nabla^n\psi^{-1}(\nabla\psi(x))$ and $L^k(z) = \prod_{i=1}^{k} \nabla l(z)$ which is an expression that makes sense only when multiplied on the left by a tensor of sufficient order. So we see that $\nabla\alpha_n(z) = \alpha_{n+1}(z)\nabla^2\psi(z)$. Given that, the first derivative with respect to $s$ is

$$\nabla^2\psi^{-1}(\nabla\psi(z) + s\nabla l(z))\nabla l(z)$$

$$= \alpha_2(z)L^1(z) + s\alpha_3(z)\nabla^2\psi(z)L^2(z) + \frac{s^2}{2}\left(\alpha_4(z)(\nabla^2\psi(z))^2 + \alpha_3(z)\nabla^3\psi(z)\right)L^3(z)$$

where we performed a Taylor expansion of the first term.

# Bibliography

[Lu21]   Haihao Lu.   An $O(s\hat{r})$-Resolution ODE Framework for Understanding Discrete-Time Algorithms and Applications to the Linear Convergence of Minimax Problems. *arXiv:2001.08826 [cs, math]*, July 2021.

[SBC16]  Weijie Su, Stephen Boyd, and Emmanuel J Candes. A Differential Equation for Modeling Nesterov's Accelerated Gradient Method: Theory and Insights. *Journal of Machine Learning Research*, 2016.