

**Université de Montréal**

**Steepest Descent as Linear Quadratic Regulation**

par

**Simon Dufort-Labbé**

Département d'informatique et de recherche opérationnelle  
Faculté des arts et des sciences

Mémoire présenté en vue de l'obtention du grade de  
Maître ès sciences (M.Sc.)  
en informatique

31 août 2021

**Université de Montréal**

Faculté des arts et des sciences

---

Ce mémoire intitulé

**Steepest Descent as Linear Quadratic Regulation**

présenté par

**Simon Dufort-Labbé**

a été évalué par un jury composé des personnes suivantes :

*Simon Lacoste-Julien*

---

(président-rapporteur)

*Pierre-Luc Bacon*

---

(directeur de recherche)

*Ioannis Mitliagkas*

---

(membre du jury)

## Résumé

---

Concorder un modèle à certaines observations, voilà qui résume assez bien ce que l'apprentissage machine cherche à accomplir. Ce concept est maintenant omniprésent dans nos vies, entre autre grâce aux percées récentes en apprentissage profond. La stratégie d'optimisation prédominante pour ces deux domaines est la minimisation d'un objectif donné. Et pour cela, la méthode du gradient, méthode de premier-ordre qui modifie les paramètres du modèle à chaque itération, est l'approche dominante. À l'opposé, les méthodes dites de second-ordre n'ont jamais réussi à s'imposer en apprentissage profond. Pourtant, elles offrent des avantages reconnus qui soulèvent encore un grand intérêt. D'où l'importance de la méthode du col, qui unifie les méthodes de premier et second-ordre sous un même paradigme.

Dans ce mémoire, nous établissons un parallèle direct entre la méthode du col et le domaine du contrôle optimal; domaine qui cherche à optimiser mathématiquement une séquence de décisions. Et certains des problèmes les mieux compris et étudiés en contrôle optimal sont les commandes linéaires quadratiques. Problèmes pour lesquels on connaît très bien la solution optimale. Plus spécifiquement, nous démontrerons l'équivalence entre une itération de la méthode du col et la résolution d'une Commande Linéaire Quadratique (CLQ).

Cet éclairage nouveau implique une approche unifiée quand vient le temps de déployer nombre d'algorithmes issus de la méthode du col, tel que la méthode du gradient et celle des gradients naturels, sans être limitée à ceux-ci. Approche que nous étendons ensuite aux problèmes à horizon infini, tel que les modèles à équilibre profond. Ce faisant, nous démontrons pour ces problèmes que calculer les gradients via la différentiation implicite revient à employer l'équation de Riccati pour solutionner la CLQ associée à la méthode du gradient. Finalement, notons que l'incorporation d'information sur la courbure du problème revient généralement à rencontrer une inversion matricielle dans la méthode du col. Nous montrons que l'équivalence avec les CLQ permet de contourner cette inversion en utilisant une approximation issue des séries de Neumann. Surprenamment, certaines observations empiriques suggèrent que cette approximation aide aussi à stabiliser le processus d'optimisation quand des méthodes de second-ordre sont impliquées; en agissant comme un

régularisateur adaptif implicite.

**Mots-clés :** optimisation, apprentissage profond, apprentissage machine, méthode du col, réseaux de neurones, commande linéaire quadratique, algorithme du gradient, méthode des gradients naturels, équation de Riccati, contrôle optimal, modèle à équilibre profond.

# Abstract

---

Machine learning entails training a model to fit some given observations, and recent advances in the field, particularly in deep learning, have made it omnipresent in our lives. Fitting a model usually requires the minimization of a given objective. When it comes to deep learning, first-order methods like gradient descent have become a default tool for optimization in deep learning. On the other hand, second-order methods did not see widespread use in deep learning. Yet, they hold many promises and are still a very active field of research. An important perspective into both methods is steepest descent, which allows you to encompass first and second-order approaches into the same framework.

In this thesis, we establish an explicit connection between steepest descent and optimal control, a field that tries to optimize sequential decision-making processes. Core to it is the family of problems known as Linear Quadratic Regulation; problems that have been well-studied and for which we know optimal solutions. More specifically, we show that performing one iteration of steepest descent is equivalent to solving a Linear Quadratic Regulator (LQR). This perspective gives us a convenient and unified framework for deploying a wide range of steepest descent algorithms, such as gradient descent and natural gradient descent, but certainly not limited to. This framework can also be extended to problems with an infinite horizon, such as deep equilibrium models. Doing so reveals that retrieving the gradient via implicit differentiation is equivalent to recovering it via Riccati's solution to the LQR associated with gradient descent. Finally, incorporating curvature information into steepest descent usually takes the form of a matrix inversion. However, casting a steepest descent step as a LQR also hints toward a trick that allows to sidestep this inversion, by leveraging Neumann's series approximation. Empirical observations provide evidence that this approximation actually helps to stabilize the training process, by acting as an adaptive damping parameter.

**Keywords:** optimization, deep learning, machine learning, steepest descent, neural networks, linear quadratic regulator, gradient descent, natural gradient descent, Riccati's equation, optimal control, deep equilibrium models.

# Contents

---

<b>Résumé</b> .....	3
<b>Abstract</b> .....	5
<b>List of Tables</b> .....	8
<b>List of Figures</b> .....	9
<b>List of acronyms and abbreviations</b> .....	10
<b>Acknowledgments</b> .....	12
<b>Chapter 1. Introduction</b> .....	13
1.1 Deep Learning Formulation in Optimal Control . . . . .	14
1.2 Contributions. . . . .	15
<b>Chapter 2. Background</b> .....	17
2.1 Optimization Methods for Deep Learning . . . . .	17
2.1.1 First-Order Methods.....	17
2.1.2 Second-Order Methods.....	18
2.2 Steepest Descent . . . . .	19
2.2.1 Gradient Descent.....	20
2.2.2 Steepest Descent under a Quadratic Norm.....	20
2.3 Linear Quadratic Regulator. . . . .	21
2.3.1 Finite Horizon.....	22
2.3.2 Infinite Horizon.....	22
2.3.3 Cyclic Case.....	23
2.4 Deep Implicit Models . . . . .	24
2.4.1 Deep Equilibrium Models.....	25
<b>Chapter 3. Steepest Descent as Linear Quadratic Regulation</b> .....	28
3.1 Problem Statement for a Single Gradient Step. . . . .	28
3.1.1 Equivalence with Gradient Descent.....	30

3.2	Steepest Descent Objective . . . . .	31
3.2.1	General LQR Equivalence.....	32
3.2.2	Riccati’s Solution to the Steepest Descent Objective.....	35
3.2.3	Equivalence with Natural Gradients.....	36
3.2.4	Equivalence with Newton’s method . . . . .	37
3.3	Experimental validation . . . . .	39
<b>Chapter 4.</b>	<b>Extension to the Infinite Horizon Case.....</b>	<b>41</b>
4.1	Problem Statement . . . . .	41
4.1.1	Riccati’s Solution in Infinite Horizon Setting.....	44
4.2	Equivalence with Implicit Differentiation . . . . .	45
4.3	Cyclic Solution . . . . .	46
4.4	LQR-based Natural Gradient Method for DEQs . . . . .	47
<b>Chapter 5.</b>	<b>Computational Perspective.....</b>	<b>50</b>
5.1	Leveraging Neumann Series. . . . .	51
5.1.1	An Algorithm without Matrix Inversion.....	53
5.1.2	Parallel with Damping Parameter.....	54
<b>Chapter 6.</b>	<b>Conclusion.....</b>	<b>56</b>
6.1	Future Work . . . . .	57
<b>References.....</b>		<b>59</b>
<b>Appendix A.</b>	<b>Steepest descent – Minimization Without Constraints.....</b>	<b>65</b>
<b>Appendix B.</b>	<b>Riccati’s Equation – Solution with a Linear Term in the Final Cost.....</b>	<b>67</b>

# List of Tables

---

1.1	Notations used in ML vs OC.....	15
-----	---------------------------------	----



## List of Figures

---

- 3.1 **top:** A NN composed of 3 convolution layers followed by 2 dense layers trained on MNIST database with SGD (5 examples per minibatch) but implemented with the LQR approach. **left:** Accuracy over one epoch of training. **center:** Loss over one epoch. **right:** The cosine similarity, measured at every step, between the gradients calculated via the LQR approach and by traditional automatic differentiation. **bottom:** Same as above, but for NGD. .... 39
- 4.1 Zoom on the accuracy achieved (first 2500 iterations) during the training of a toy DEQ (a Resnet block with 12 inner channels) over MNIST with the NGD algorithm arising from the LQR approach. NGD starts to converge in less iteration, but it is 10x slower than GD..... 49
- 5.1 **top:** Accuracy curves of a training experiment, the same Convnet as in figure 3.1 (5 layers deep) trained on MNIST with the LQR approach equivalent to NGD and a learning rate set to  $\alpha = 0.1$ . **left:** In addition to the learning rate, a damping parameter  $\delta$  has to be fine-tuned for the training to be stable without Neumann's series trick. **right:** On the opposite, implementing Neumann's series trick also has the added benefit of adding a layer adaptive damping term that stabilizes the training. **bottom:** Same as above, but for a learning rate set to  $\alpha = 0.01$ ..... 55

## List of acronyms and abbreviations

---

CARE	Continuous Time Algebraic Riccati Equation
CLQ	Commande Linéaire Quadratique
DARE	Discrete Algebraic Riccati Equation
DL	Deep Learning
DNN	Deep Neural Nets
GD	Gradient Descent
KL	Kullback–Leibler (used for Kullback–Leibler divergence)
LQR	Linear Quadratic Regulator
MVP	Matrix-Vector Product
ML	Machine Learning

MNIST	Modified National Institute of Standards and Technology
NGD	Natural Gradient Descent
NN	Neural Network
OC	Optimal Control
OCP	Optimal Control Problem
ODEs	Ordinary Differential Equations
PMP	Pontryagin's Maximum Principle
RL	Reinforcement Learning
SGD	Stochastic Gradient Descent
SVD	Singular Values Decomposition

## Acknowledgments

---

It would be a sin not to first thank my supervisor, Dr. Pierre-Luc Bacon. It was a chance and a privilege to be given the opportunity to work with him. Without his vision and knowledge, bringing to a satisfactory completion this project in 8 months would not have been possible.

Special thanks also need to be addressed to my colleague, Evgenii Nikishin. It is hard to describe how helpful it is to be granted access to a living encyclopedia.

To Ryan D’Orazio and David Yu-Tung Hui, thank you for providing specific comments when it mattered, it was very appreciated.

It goes without saying, discussions and feedback from the members in the research group I am lucky enough to be part of were always insightful. To all of you, thank you.

Je ne peux non plus manquer de remercier ma famille et mon entourage, dont le support constant constitue le socle duquel je peux m’élancer vers les nombreux projets qui m’aspirent. À Ariane, ma conjointe, qui pallie mes défauts et équilibre ma vie. Et finalement à Léonard, mon fils, qui du haut de son an se révèle être une source inépuisable de motivation.

# Chapter 1

---

## Introduction

It took a while for Deep Learning (DL) to shine and become the workhorse of artificial intelligence, with many pieces falling into place to make it happen. One of the key components that enabled neural networks to be trained at scale is the development of efficient optimization methods tailored to the high dimensional and noisy regime under which deep neural networks typically operate. First-order optimization methods have predominantly been used for training deep neural networks due to their low computational overhead (see [Goodfellow and Vinyals 2015](#)). While computationally advantageous, first-order optimization methods fall short of leveraging the local geometry of the loss landscape; for example, first-order methods will ignore curvature while second-order methods can converge faster by leveraging this information. In order to increase the computational efficiency of second-order methods, [Pearlmutter 1994](#); [Martens and Sutskever 2011](#) proposed matrix-free variants using the "forward over reverse" approach of automatic differentiation ([Griewank and Naumann 2003](#)) combined with a matrix-free linear solver such as Conjugate Gradient method ([Hestenes and Stiefel 1952](#)).

Despite these improvements, second-order optimization methods have not been adopted by practitioners. We believe that the problem is twofold: 1) the inadequacy of the local geometry implied by the chosen second-order solver, 2) the need for a layer-aware approach in improving computational efficiency.

This thesis addresses those two challenges jointly by proposing a framework in which a large class of steepest descent methods can be implemented efficiently with the same blueprints. The key insight in this work is to recognize that under certain assumptions, the problem of finding the direction of steepest descent is amenable to solving a local Linear Quadratic Regulation problem.

This provided us a direct entry point to try to address these issues more generally in the Optimal Control (OC) framework which, as in Reinforcement Learning (RL), is concerned with how to optimize sequential decision-making processes (Pontriagin *et al.* 1962; Bellman 1952). Our results offer a new perspective on the nature of the computation performed during an optimization step: the parameters update step of a whole family of steepest descent algorithms can itself be expressed as LQR problems.

The history of deep learning is intertwined with that of optimal control. Indeed, backpropagation was first described via the Optimal Control Problem (OCP) framework and studied under the Hamiltonian formalism (see Bryson and Ho 1969, Lecun 1988 and Linnainmaa 1976); this formalism being the backbone of OC. The specific case of the neural net optimization was also connected to optimal control theory explicitly by Dreyfus 1990, who applied the Kelley-Bryson procedure, a.k.a backpropagation, to optimize neural nets. More recently, Zhang *et al.* 2019 and Seidman *et al.* 2020 studied min-max optimization as an optimal control problem to improve training stability. Optimal control has also recently gained attention within the DL community for modeling continuous dynamical systems governed by Ordinary Differential Equations (ODEs). This line of work came from Chen *et al.* 2018 who developed Neural ODEs from the question: what if the discretization step in ResNets (He *et al.*, 2016) was made increasingly smaller? To efficiently optimize such models, the authors used the adjoint method from Bryson and Ho 1969 in optimal control to evaluate the derivatives. Using insights from optimal control, Haber and Ruthotto 2017 also investigated the stability of forward propagation in neural networks by studying a continuous-time counterpart that led them to propose new neural architectures. The work developed in this thesis extends those ideas by showing for the first time that each gradient step during training can be computed as the solution to a specific kind of optimal control problem.

## 1.1 Deep Learning Formulation in Optimal Control

We follow the path outlined above, and formulate some of the concepts encountered in ML in the language of OC. However, we currently restrict ourselves to the discrete setting. The training task of a NN can thus be posed as a Mayer problem (see Bliss 1968); that is, an optimal control problem (OCP) of the form:

$$\begin{aligned}
 & \min_{\mathbf{u}} J(\mathbf{x}, \mathbf{u}) \\
 & \text{subject to } \mathbf{x}_{i+1} = f_i(\mathbf{x}_i, \mathbf{u}_i) \\
 & \text{given } \mathbf{x}_0
 \end{aligned} \tag{1.1.1}$$

Indeed, when training a NN, one would receive an input ( $\mathbf{x}_0$ ) that will then be processed sequentially through all the layers  $f_i$  before evaluating the loss ( $J(\mathbf{x}, \mathbf{u})$ ) over the final representation only.

Since we will be leveraging tools from the OC setting to perform optimization in ML, we decided to align our notation with as is conventional in OC. When dealing with the specific task of optimizing a NN in ML, the correspondence between the concepts are as follows:

ML terminology		Control terminology	
Name	Notation	Name	Notation
Input	$\mathbf{x}_0$	Initial state	$\mathbf{x}_0$
Feature/Activation	$\mathbf{a}_i$	Transition state	$\mathbf{x}_i$
Output	$\mathbf{y}$	Final state	$\mathbf{x}_N$
Parameters	$\boldsymbol{\theta}$	Controls	$\mathbf{u}$
Loss function	$l(\cdot)$	Cost function	$J(\cdot)$
Layer		Transition function	$f_i$

**Table 1.1.** Notations used in ML vs OC

## 1.2 Contributions

The work in this thesis was sparked by our attempt to leverage a connection shown in Bertsekas (2016) (section 2.6.2) between a particular LQR problem and Newton’s descent step for an OCP. This particular implementation of Newton’s method was first formalized simultaneously by Dunn and Bertsekas (1989) and independently by De O. Pantoja and Mayne (1989). We realized that this was only a specific case and that an equivalence could be built more broadly for the steepest descent minimization objective. A similar connection has also been found by Jin *et al.* (2020), but only in the context of an end-to-end differentiable architecture for reinforcement learning. Our results apply more broadly to unconstrained optimization in other areas of ML.

Thus the contributions of this thesis, to the best of our knowledge, include:

- We describe how a gradient descent step, a steepest descent step under a Euclidean norm, over a discrete dynamical system with terminal cost (which encompasses neural networks) is always equivalent to the solution of a LQR (section 3.1).
- We generalize this equivalence between LQR and an individual steepest descent step for a variety of divergence measures that encompasses second-order methods such as natural gradient descent and Newton’s method (section 3.2).

- As a specific case we derive the LQR equivalent to performing a natural gradient descent update (section 3.2.3).
- We extend the above results to the infinite horizon setting, thus making the technique appropriate to train models such as deep equilibrium models. As a consequence, this gives us an algorithm to train such models with second-order methods, something that does not seem to have been attempted yet (chapter 4).
- We show that, as a special case, solving the LQR equivalent to gradient descent in the infinite horizon setting is equivalent to retrieving it via implicit differentiation (section 4.2).
- The matrix inversion involved with the general steepest procedure descent is usually computationally expensive. We show how we can sidestep this difficulty for NN that possesses ReLu activation functions by deriving a trick based on Neumann's series (section 5.1).



# Chapter 2

---

## Background

Steepest descent is a paradigm that encompasses a whole family of optimization algorithms, but it does not necessarily provide a unified method to solve the optimization problems formulated through it. We will attempt to cover some of this ground, by bringing a few OC results into the realm of DL. We assume most readers will be acquainted with one of the two fields, but not necessarily both. As such, we will give a very brief overview of the optimization trends in DL before going more technical and providing the details that will be leveraged from steepest descent. We will continue in this vein and explain the exploited concepts about LQR and the formulation associated with it in OC. Finally, we will succinctly introduce the deep implicit models' framework as we will be extending some of our results to this area.

## 2.1 Optimization Methods for Deep Learning

### 2.1.1 First-Order Methods

The optimization landscape in Deep Learning (DL) is largely dominated by first-order stochastic gradient descent methods. This simple approach is extremely efficient in DL, providing an unbiased estimate of the gradient and being able to escape local minima (see [Bottou 1998](#) for more information about SGD). It also scales very well with the number of examples, it is thus suited to big data (see [Robbins 2007](#)). Yet, this method still has few disadvantages: it often requires manual tuning of hyperparameters such as the learning rate and may make very slow progress in a noisy setting. These issues are well discussed in [Ruder \(2016\)](#). Multiple variants were later used for training Deep Neural Nets (DNN) in order to mitigate the challenges faced by SGD.

Momentum ([Polyak 1964](#)) is a key idea to accelerate training; it accumulates past gradients and adds them (scaled by an exponentially decaying rate) to the current update. This technique is particularly efficient in noisy high curvature settings. It is also a core idea

when it comes to accelerating learning and has many derivatives. As an example, Nesterov acceleration (Nesterov 1983) extends this idea but evaluates the gradient after applying it. This strategy brings a better rate of convergence in deterministic settings and a practical speed up in stochastic scenarios, such as the training of DNN (see Sutskever *et al.* 2013).

Even with those improvements, setting the learning rate correctly remains a difficult hyperparameter optimization problem. Momentum somewhat improved the situation but did so by introducing a new hyperparameter, namely the mass parameter. Methods were thus designed to specifically address the problem, by introducing an adaptive learning rate tailored separately to each weight. AdaGrad (Duchi *et al.* 2011) did so by proportionally scaling a learning rate decay to the size of the partial derivatives. Weights that adapt more slowly will therefore decay less and the training procedure should make better progress in the directions for which the gradient slope is smoother. RMSProp (Hinton 2012) modified AdaGrad by discarding the earlier part of the updates history when calculating the decay rate, allowing for better performance in a non-convex setting.

Finally, when it comes to first-order methods, Adam (Kingma and Ba 2015) is probably the default choice for DL practitioners nowadays, as it is considered to be more robust to the choice of hyperparameters. It incorporates momentum as an estimate of the first-order moments of the gradients and includes a bias correction term to account for the initialization at the origin. The bias correction leverages an estimate of the second-order moments along with the estimate of the first-order moments.

### 2.1.2 Second-Order Methods

In contrast, second-order methods never saw wide applications in DL, suffering drastically from the curse of dimensionality as the algorithms implementing them traditionally require computation time and memory storage scaling with  $\mathcal{O}(n^3)$ . However, there is experimental evidence (Martens 2016) that they allow training convergence in fewer iterations, even if the time to do so may not be smaller than via first-order methods (for which an individual update is computed much more quickly). Martens 2016 provides a good overview of the difficulties as well as the advantages that come with those methods when applied to DL.

Because of those challenges, most early applications of second-order methods to Machine Learning (ML) were through the methods known as Quasi-Newton that replace the Hessian inverse by an approximation. The most popular among those is the Broyden–Fletcher–Goldfarb–Shanno (BFGS) method for which a limited memory variant (Liu and Nocedal 1989) can offer better performance than plain SGD in some instances (see Le

*et al.* 2011). However, their applicability to DL remains limited as they still do not scale well with the large number of parameters that a DNN usually possesses.

Recently, another second-order method has drawn some attention: the Natural Gradient Descent (NGD). This method leverages the Fischer Information matrix to condition the gradient update and was first introduced by [Amari 1998](#). In the context of DL, empirical evidence was gathered about the robustness of the approach when facing limited training data by [Pascanu and Bengio 2014](#).

The main difficulty regarding the widespread use of NGD remains its scalability since in theory it also requires to invert a matrix scaling with the number of parameters (the Fischer information matrix). A recent approach, K-FAC ([Martens and Grosse 2015](#)), alleviated part of the problem by first performing a block-diagonal approximation of the Fischer information matrix and then using Kronecker factorization to ease the inversion task. It was shown that this strategy could be scaled and adapted to distributed systems by [Osawa et al. 2019](#), but the result requires more overheads than typical first-order methods.

## 2.2 Steepest Descent

*Steepest descent* can be seen as a generalization of gradient descent. It should be noted also that the *mirror descent* approach can be cast as the steepest direction in the Riemannian space induced by the chosen divergence (see [Raskutti and Mukherjee 2015](#); [Beck and Teboulle 2003](#)). The core idea behind those methods is that two quantities need to be determined in order to optimize a function locally. Thus, one has to decide on a direction in which to search to improve the objective as well as the magnitude of the step to take in the chosen direction.

If we first fix the step size, what is left to find is the optimal direction in which to take that step. For a given direction  $\mathbf{v}$ , the first-order Taylor’s approximation (see [Nocedal and Wright, 2006](#)) yields:

$$f(\mathbf{x} + \mathbf{v}) \approx f(\mathbf{x}) + \nabla f(\mathbf{x})^T \mathbf{v}$$

The *steepest direction* of norm 1 will therefore be the vector  $\mathbf{v}$  that minimizes:

$$\begin{aligned} \mathbf{v}^* &= \arg \min_{\mathbf{v}} f(\mathbf{x}) + \nabla f(\mathbf{x})^T \mathbf{v} \quad \text{s.t. } \|\mathbf{v}\| = 1 \\ &= \arg \min_{\mathbf{v}} \nabla f(\mathbf{x})^T \mathbf{v} \quad \text{s.t. } \|\mathbf{v}\| = 1 \end{aligned} \tag{2.2.1}$$

Already, we can see that the steepest descent algorithm can come in many flavors. Indeed depending on the norm choice, we can obtain very different updates. Although this problem might seem difficult to solve because of the constraint requiring  $\|\mathbf{v}\| = 1$ , we can relax the problem so that we recover the same direction, but rescaled. That is, there exists a constant  $\alpha$  such that (2.2.1) is equivalent to:

$$\mathbf{v}^* = \alpha \left[ \arg \min_{\mathbf{v}} \left\{ \nabla f(\mathbf{x})^T \mathbf{v} + \frac{1}{2} \|\mathbf{v}\|^2 \right\} \right] \quad (2.2.2)$$

which is a much more convenient form. The details for this equivalence are provided in Appendix A and a more formal introduction to steepest descent can be found in section 9.4 of [Boyd and Vandenberghe 2004](#).

### 2.2.1 Gradient Descent

Gradient descent is an instance of steepest descent, and we recover this specific algorithm by choosing the Euclidean norm in the steepest descent formulation. Indeed, in this case, we want to minimize:

$$\alpha \left[ \arg \min_{\mathbf{v}} \left\{ \nabla f(\mathbf{x})^T \mathbf{v} + \frac{1}{2} \mathbf{v}^T \mathbf{v} \right\} \right]$$

This is a simple quadratic so that we can recover the minimum by setting the derivative to zero:

$$0 = \alpha \nabla f(\mathbf{x}) + \mathbf{v}^* \implies \mathbf{v}^* = -\alpha \nabla f(\mathbf{x}) \ ,$$

where such a  $\mathbf{v}^*$  is exactly the update we want to apply when we perform gradient descent.

### 2.2.2 Steepest Descent under a Quadratic Norm

Let us now consider the quadratic norm  $\|\cdot\|_P$  defined as  $\|\mathbf{z}\|_P = (\mathbf{z}^T \mathbf{P} \mathbf{z})^{1/2}$  for a matrix  $\mathbf{P}$  that is positive semidefinite.

Then, just as above, we can seek the minimum by setting the derivative to zero, which yields:

$$\mathbf{v}^* = -\alpha \mathbf{P}^{-1} \nabla f(\mathbf{x})$$

With this formulation, we can therefore also cast NGD and Newton's descent as special instances of steepest descent, where  $\mathbf{P}$  is chosen to be the Fisher matrix,  $\mathbf{F}$ , for the former or the Hessian  $\mathbf{H}$  for the latter.

## 2.3 Linear Quadratic Regulator

Linear quadratic regulators are the canonical problems of control theory. Although far from being a simple problem, LQRs are well-studied problems that have the nice property of being analytically solvable. Moreover, LQRs have many practical applications and thus many numerical solvers have been developed for this family of problems over the years. Readers interested in a more thorough introduction can refer themselves to [Anderson and Moore 1990](#).

LQR are control problems for which the dynamics are described by a linear system. In control theory,  $\mathbf{x}_t \in \mathbb{R}^n$  is often used to denote the *state* of the system at time  $t$  while  $\mathbf{u}_t \in \mathbb{R}^m$  denotes the *control* taken at time  $t$ . The control is usually the quantity that one can vary to obtain the desired trajectory. Under this notation, the system dynamics can be described by:

$$\mathbf{x}_{t+1} = \mathbf{A}_t \mathbf{x}_t + \mathbf{B}_t \mathbf{u}_t \quad (2.3.1)$$

The goal is then to minimize a cost associated with the trajectory obtained by following the system dynamics. For LQRs problems, this cost must be at most quadratic with respect to its states and controls and of the following form:

$$J(\mathbf{x}, \mathbf{u}) = \mathbf{a}_T^T \mathbf{x}_T + \mathbf{x}_T^T \mathbf{Q}_T \mathbf{x}_T + \sum_{i=0}^{T-1} \left[ \mathbf{a}_i^T \mathbf{x}_i + \mathbf{x}_i^T \mathbf{Q}_i \mathbf{x}_i + \mathbf{b}_i^T \mathbf{u}_i + \mathbf{u}_i^T \mathbf{R}_i \mathbf{u}_i + \mathbf{u}_i^T \mathbf{M}_i \mathbf{x}_i \right] \quad (2.3.2)$$

Note that objective function is the sum of per-state costs, and the immediate cost at a given time step depends solely on the current state and control. That is, there are no cross-terms in the cost between parameters of different time steps.

For such problems, the optimal values for the controls can be retrieved via the Riccati equation. To do so, we have to solve backward for the gain matrix that will define the optimal controls. The Riccati equation is provided below for the two cases that we encounter when building the equivalence between a steepest descent step and a LQR.

Though LQR problems are well defined both for continuous and discrete settings, we will focus on the latter case. This is because DL models currently consist mostly of discrete optimization problems; as seen in section 1.1 the passage from one layer to another can be seen as moving forward by one time step. The connection made later between an optimization step and an LQR will be grounded in the discrete setting.

### 2.3.1 Finite Horizon

In OC, a finite horizon means that the dynamic system reaches a final state after  $N$  steps. For us here, the final state would be the final output given by the last layer of a NN.

LQR problems that are well formulated will fit the form:

$$\begin{aligned} \min_{\mathbf{u}} \quad & J(\mathbf{x}, \mathbf{u}) = \mathbf{a}_T^T \mathbf{x}_T + \mathbf{x}_T^T \mathbf{Q}_T \mathbf{x}_T + \sum_{i=0}^{T-1} \left[ \mathbf{a}_i^T \mathbf{x}_i + \mathbf{x}_i^T \mathbf{Q}_i \mathbf{x}_i + \mathbf{b}_i^T \mathbf{u}_i + \mathbf{u}_i^T \mathbf{R}_i \mathbf{u}_i + \mathbf{u}_i^T \mathbf{M}_i \mathbf{x}_i \right] \\ \text{subject to} \quad & \mathbf{x}_{i+1} = \mathbf{A}_i \mathbf{x}_i + \mathbf{B}_i \mathbf{u}_i \\ \text{given} \quad & \mathbf{x}_0 \end{aligned} \tag{2.3.3}$$

To solve such problems, we can leverage the Riccati equation to first retrieve the gain matrices  $\mathbf{K}_i$  and the vector  $\boldsymbol{\lambda}_i$  (see [Reid, 1972](#)).

$$\mathbf{K}_i = \mathbf{A}_i^T \mathbf{K}_{i+1} \mathbf{A}_i + \mathbf{Q}_i - (\mathbf{A}_i^T \mathbf{K}_{i+1} \mathbf{B}_i + \mathbf{M}_i^T) (\mathbf{R}_i + \mathbf{B}_i^T \mathbf{K}_{i+1} \mathbf{B}_i)^{-1} (\mathbf{M}_i + \mathbf{B}_i^T \mathbf{K}_{i+1} \mathbf{A}_i) \tag{2.3.4}$$

$$\boldsymbol{\lambda}_i = \mathbf{a}_i + \mathbf{A}_i^T \boldsymbol{\lambda}_{i+1} - (\mathbf{A}_i^T \mathbf{K}_{i+1} \mathbf{B}_i + \mathbf{M}_i^T) (\mathbf{R}_i + \mathbf{B}_i^T \mathbf{K}_{i+1} \mathbf{B}_i)^{-1} (\mathbf{b}_i + \mathbf{B}_i^T \boldsymbol{\lambda}_{i+1}) \tag{2.3.5}$$

With  $\mathbf{K}_N = \mathbf{Q}_N$  and  $\boldsymbol{\lambda}_N = \mathbf{a}_N$ .

The optimal controls are then calculated with:

$$\begin{aligned} \mathbf{u}_0^* &= -(\mathbf{R}_0 + \mathbf{B}_0^T \mathbf{K}_1 \mathbf{B}_0)^{-1} (\mathbf{b}_0 + \mathbf{B}_0^T \boldsymbol{\lambda}_1) \\ \mathbf{x}_0^* &= \mathbf{x}_0 \\ \mathbf{x}_{i+1}^* &= \mathbf{A}_i \mathbf{x}_i^* + \mathbf{B}_i \mathbf{u}_i^* \\ \mathbf{u}_i^* &= -(\mathbf{R}_i + \mathbf{B}_i^T \mathbf{K}_{i+1} \mathbf{B}_i)^{-1} \left( (\mathbf{M}_i + \mathbf{B}_i^T \mathbf{K}_{i+1} \mathbf{A}_i) \mathbf{x}_i^* + \mathbf{b}_i + \mathbf{B}_i^T \boldsymbol{\lambda}_{i+1} \right) \end{aligned} \tag{2.3.6}$$

Since a full development of those equations that encompasses the case when  $\mathbf{a}_i \neq 0$  and  $\mathbf{b}_i \neq 0$  is lengthy, we provide a proof in Appendix B.

### 2.3.2 Infinite Horizon

The infinite horizon setting deals with the case when  $N \rightarrow \infty$ , with both the cost terms  $(\mathbf{a}, \mathbf{b}, \mathbf{Q}, \mathbf{M}, \mathbf{R})$  and the transition matrices  $(\mathbf{A}, \mathbf{B})$  remaining constant through time. In

such cases, it makes no sense to have a cost associated with the final state, and the problem formulation becomes:

$$\begin{aligned} \min_{\mathbf{u}} J(\mathbf{x}, \mathbf{u}) &= \sum_{i=0}^{\infty} \left[ \mathbf{a}^T \mathbf{x}_i + \mathbf{x}_i^T \mathbf{Q} \mathbf{x}_i + \mathbf{b}^T \mathbf{u}_i + \mathbf{u}_i^T \mathbf{R} \mathbf{u}_i + \mathbf{u}_i^T \mathbf{M} \mathbf{x}_i \right] \\ \text{subject to } \mathbf{x}_{i+1} &= \mathbf{A} \mathbf{x}_i + \mathbf{B} \mathbf{u}_i \\ \text{given } \mathbf{x}_0 & \end{aligned} \tag{2.3.7}$$

Again, we can use the associated Riccati equation to obtain the gain matrix  $\mathbf{K}$  and  $\boldsymbol{\lambda}$ .

$$\mathbf{K} = \mathbf{A}^T \mathbf{K} \mathbf{A} + \mathbf{Q} - (\mathbf{A}^T \mathbf{K} \mathbf{B} + \mathbf{M}^T) (\mathbf{R} + \mathbf{B}^T \mathbf{K} \mathbf{B})^{-1} (\mathbf{M} + \mathbf{B}^T \mathbf{K} \mathbf{A}) \tag{2.3.8}$$

$$\boldsymbol{\lambda} = \mathbf{a} + \mathbf{A}^T \boldsymbol{\lambda} - (\mathbf{A}^T \mathbf{K} \mathbf{B} + \mathbf{M}^T) (\mathbf{R} + \mathbf{B}^T \mathbf{K} \mathbf{B})^{-1} (\mathbf{b} + \mathbf{B}^T \boldsymbol{\lambda}) \tag{2.3.9}$$

The equation (2.3.8) is known as the Discrete Algebraic Riccati Equation (DARE), and its solution has many practical uses. Many efficient numerical DARE solvers already exist, and it is still an active area of research (see [Bini et al. 2011](#)).

A solution exists to (2.3.8) if the *closed loop state transfer matrix*,  $\mathbf{A} - \mathbf{B}(\mathbf{R} + \mathbf{B}^T \mathbf{K} \mathbf{B})^{-1} (\mathbf{M} + \mathbf{B}^T \mathbf{K} \mathbf{A})$ , has a spectral radius smaller than one (it is stable). The optimal trajectory can then be retrieved in the same fashion:

$$\begin{aligned} \mathbf{u}_0^* &= -(\mathbf{R} + \mathbf{B}^T \mathbf{K} \mathbf{B})^{-1} (\mathbf{b} + \mathbf{B}^T \boldsymbol{\lambda}) \\ \mathbf{x}_0^* &= \mathbf{x}_0 \\ \mathbf{x}_{i+1}^* &= \mathbf{A} \mathbf{x}_i^* + \mathbf{B} \mathbf{u}_i^* \\ \mathbf{u}_i^* &= -(\mathbf{R} + \mathbf{B}^T \mathbf{K} \mathbf{B})^{-1} \left( (\mathbf{M} + \mathbf{B}^T \mathbf{K} \mathbf{A}) \mathbf{x}_i^* + \mathbf{b} + \mathbf{B}^T \boldsymbol{\lambda} \right) \end{aligned} \tag{2.3.10}$$

Note that the optimal controls are guaranteed to stabilize when the closed-loop state transfer matrix is stable.

### 2.3.3 Cyclic Case

Finally, as it will be relevant later, we mention that a *cyclic* case exists for LQR in the infinite horizon setting. If a given LQR cycles through  $N$  different linear system, it takes the form:

$$\begin{aligned}
\min_{\mathbf{u}} J(\mathbf{x}, \mathbf{u}) &= \sum_{i=0}^{\infty} \sum_{n=1}^N \left[ \mathbf{a}_n^T \mathbf{x}_{in} + \mathbf{x}_{in}^T \mathbf{Q}_n \mathbf{x}_{in} + \mathbf{b}_n^T \mathbf{u}_{in} + \mathbf{u}_{in}^T \mathbf{R}_n \mathbf{u}_{in} + \mathbf{u}_{in}^T \mathbf{M}_n \mathbf{x}_{in} \right] \\
\text{subject to } \mathbf{x}_{i(n+1)} &= \mathbf{A}_n \mathbf{x}_{in} + \mathbf{B}_n \mathbf{u}_{in} \\
\text{with } \mathbf{x}_{iN} &= \mathbf{x}_{(i+1)0} \\
\text{given } \mathbf{x}_{00} &
\end{aligned} \tag{2.3.11}$$

The Riccati solution is similar, the difference being that we now need to solve a system of gain matrices.

$$\begin{aligned}
\mathbf{K}_{N-1} &= \mathbf{A}_{N-1}^T \mathbf{K}_N \mathbf{A}_{N-1} + \mathbf{Q}_{N-1} - \\
&\quad (\mathbf{A}_{N-1}^T \mathbf{K}_N \mathbf{B}_{N-1} + \mathbf{M}_{N-1}^T) (\mathbf{R}_{N-1} + \mathbf{B}_{N-1}^T \mathbf{K}_N \mathbf{B}_{N-1})^{-1} (\mathbf{M}_{N-1} + \mathbf{B}_{N-1}^T \mathbf{K}_N \mathbf{A}_{N-1}) \\
&\quad \vdots \\
\mathbf{K}_1 &= \mathbf{A}_1^T \mathbf{K}_2 \mathbf{A}_1 + \mathbf{Q}_1 - (\mathbf{A}_1^T \mathbf{K}_2 \mathbf{B}_1 + \mathbf{M}_1^T) (\mathbf{R}_1 + \mathbf{B}_1^T \mathbf{K}_2 \mathbf{B}_1)^{-1} (\mathbf{M}_1 + \mathbf{B}_1^T \mathbf{K}_2 \mathbf{A}_1) \\
\mathbf{K}_N &= \mathbf{A}_N^T \mathbf{K}_1 \mathbf{A}_N + \mathbf{Q}_N - (\mathbf{A}_N^T \mathbf{K}_1 \mathbf{B}_N + \mathbf{M}_N^T) (\mathbf{R}_N + \mathbf{B}_N^T \mathbf{K}_1 \mathbf{B}_N)^{-1} (\mathbf{M}_N + \mathbf{B}_N^T \mathbf{K}_1 \mathbf{A}_N)
\end{aligned} \tag{2.3.12}$$

Likewise, we need to solve the cyclic system for  $\lambda_i$ .

The good news for us is that, again, there are multiple cyclic DARE solvers already available to solve such problems.

## 2.4 Deep Implicit Models

Implicit models were recently brought into light within the deep learning community first by its use in modeling Ordinary Differential Equations (ODEs) system [Chen \*et al.\* 2018](#) and then the introduction of the Deep Equilibrium Models (DEQ) [Bai \*et al.\* 2019](#) that leveraged the same ideas over modern architectures. Let us remark that those papers made popular the uses of implicit layers within the ML community, but the core ideas behind date back to the '80s (e.g. [Almeida 1990](#) and [Pineda 1987](#)).

The key point to the implicit layers is that instead of representing a typical transformation of the form :

$$\mathbf{z} = f(\mathbf{x}) \tag{2.4.1}$$



it instead focus on its *implicit* representation of the form :

$$g(\mathbf{z}, \mathbf{x}) = f(\mathbf{x}) - \mathbf{z} = 0 \tag{2.4.2}$$

and the goal is now to find the root of the above equation.

This formulation becomes particularly useful when the transformation  $f$  has a dependency on  $\mathbf{z}$  ( $f : \mathcal{X} \times \mathcal{Z} \rightarrow \mathcal{Z}$ ) and is applied repeatedly until convergence to a fixed point, that we can denote  $\mathbf{z}^*$ .

The root-finding problem is then:

$$g(\mathbf{z}, \mathbf{x}) = f(\mathbf{x}, \mathbf{z}) - \mathbf{z} = 0 \tag{2.4.3}$$

With a solution at the fixed point  $\mathbf{z}^* = f(\mathbf{x}, \mathbf{z}^*)$ . The fixed point has an implicit dependency over  $\mathbf{x}$ , and we can design the function  $\varphi : \mathcal{X} \rightarrow \mathcal{Z}$  to be the function returning the fixed point of  $f(\mathbf{x}, \mathbf{z})$ , i.e.  $\varphi(\mathbf{x}) = \mathbf{z}^*$ .

For such problems, when a fixed point exists, we can then leverage the Implicit Function Theorem to find the derivative to (2.4.3) around the fixed point  $\mathbf{z}^*$ , given by:

$$\nabla_{\mathbf{x}}^T \varphi(\mathbf{x}) = \left( \nabla_{\mathbf{z}}^T g(\mathbf{x}, \mathbf{z}) \Big|_{\mathbf{z}=\mathbf{z}^*} \right)^{-1} \nabla_{\mathbf{x}}^T g(\mathbf{x}, \mathbf{z}^*) \tag{2.4.4}$$

Under the condition that the matrix  $\left( \nabla_{\mathbf{z}}^T g(\mathbf{x}, \mathbf{z}) \Big|_{\mathbf{z}=\mathbf{z}^*} \right)^{-1}$  is invertible.

This relationship has practical consequences: the Jacobian of the implicit layer  $f$  can be recovered *without* having to backpropagate through the fixed-point solver/iterator.

This means that we are free to use any fixed point solver to find  $\mathbf{z}^*$  and then recover analytically the Jacobian of interest via (2.4.4): without representing explicitly the computational graph (as build by typical automatic differentiation tools) of the fixed-point solver.

### 2.4.1 Deep Equilibrium Models

DEQ models deploy the above ideas in the realm of deep learning. Indeed, conceptually, one could think of a DEQ as a recurrent neural network of infinite depth. If the representation induced by the infinitely repeated layer eventually reaches a fixed point, such a NN could be written as:

$$\begin{aligned}
z_0 &= 0 \\
z_{i+1} &= f(\mathbf{x}, \boldsymbol{\theta}, z_i) \quad \text{until } z_{i=1} = z_i \\
J(\mathbf{x}) &= l(z^*(\mathbf{x}, \boldsymbol{\theta}))
\end{aligned} \tag{2.4.5}$$

Where  $z^*$  is the fixed point of the layer  $f(\mathbf{x}, \boldsymbol{\theta}, z)$ ,  $\mathbf{x}$  is the input,  $\boldsymbol{\theta}$  are the layer parameters and  $l(\cdot)$  is the loss calculated with respect to the NN output.

In practice, we often want to optimize such a network with a variant of the SGD algorithm. In all cases, the optimization step will most likely require the computation of the quantity:

$$\nabla_{\boldsymbol{\theta}} \varphi(\mathbf{x}, \boldsymbol{\theta}) \nabla_{z^*} l(z^*) = \nabla_{\boldsymbol{\theta}} l(z^*)$$

That is, the vector-Jacobian product between the Jacobian of the fixed point with respect to the parameters and the gradient of the loss with respect to the final representation. To recover such a quantity, we can leverage the Implicit Function Theorem showcased above that tells us:

$$\nabla_{\boldsymbol{\theta}}^T \varphi(\mathbf{x}, \boldsymbol{\theta}) = \left( \nabla_z^T f(\mathbf{x}, \boldsymbol{\theta}, z) \Big|_{z=z^*} - \mathbf{I} \right)^{-1} \nabla_{\boldsymbol{\theta}}^T f(\mathbf{x}, \boldsymbol{\theta}, z^*) \tag{2.4.6}$$

For any vector  $\mathbf{b}$  (such as  $\nabla_{z^*} l(z^*)$ ), we have that:

$$\nabla_{\boldsymbol{\theta}} \varphi(\mathbf{x}, \boldsymbol{\theta}) \mathbf{b} = \nabla_{\boldsymbol{\theta}} f(\mathbf{x}, \boldsymbol{\theta}, z^*) \left( \nabla_z^T f(\mathbf{x}, \boldsymbol{\theta}, z) \Big|_{z=z^*} - \mathbf{I} \right)^{-T} \mathbf{b} \tag{2.4.7}$$

Finally, to be compatible with the automatic differentiation framework, one will want to retrieve the quantity  $\left( \nabla_z^T f(\mathbf{x}, \boldsymbol{\theta}, z) \Big|_{z=z^*} - \mathbf{I} \right)^{-T} \mathbf{b}$  by solving instead the system:

$$\bar{\mathbf{x}} = \left( \nabla_z^T f(\mathbf{x}, \boldsymbol{\theta}, z) \Big|_{z=z^*} \right) \bar{\mathbf{x}} - \mathbf{b} \tag{2.4.8}$$

The system above also describes a fixed-point equation, and one key concept of the DEQ networks is that we have access to a fixed-point solver. We thus have all the tools needed to derive the gradient update with automatic differentiation. Moreover, if a fixed point exists for the implicit layer of the DEQ (2.4.5), then we are guaranteed that a fixed point also exists for (2.4.8). Indeed, the local convergence condition for both system is the same : that  $\left( \nabla_z^T f(\mathbf{x}, \boldsymbol{\theta}, z) \Big|_{z=z^*} \right)$  has a spectral radius smaller than 1 (and is thus a contraction,

see [Nayfeh and Balachandran 1995](#) section 2.2).

We now have in hand everything needed to derive the equivalences proposed in the introduction.

# Chapter 3

---

## Steepest Descent as Linear Quadratic Regulation

We begin by showing that performing a single GD step with steepest descent is equivalent to solving a corresponding unique LQR. From there, we will be able to move toward our main result, which extends this equivalence to steepest descent under various divergence measures. Please note that the results presented in this chapter are developed for finite horizon, discrete OCP – a setting that holds when trying to optimize a traditional NN made from a finite number of layers.

### 3.1 Problem Statement for a Single Gradient Step

For consistency, we will now denote the various quantities via their typical OC denomination, following the equivalences put forward in table 1.1. Let  $F : \mathbb{R}^n \rightarrow \mathbb{R}^d$  be a NN mapping the input  $\mathbf{x}_0^k$  to the output  $\mathbf{x}_N^k$  through the applications of  $N - 1$  layers labelled  $f_0, f_1, \dots, f_{N-1}$ . Each layers has its own parameters  $\mathbf{u}_i$  so that  $\mathbf{x}_{i+1}^k = f(\mathbf{x}_i^k, \mathbf{u}_i)$ . The whole set of parameters can be written as  $\mathbf{u} = \{\mathbf{u}_0, \mathbf{u}_1, \dots, \mathbf{u}_{N-1}\}$ . A loss function of the form  $J(\mathbf{x}_N^k)$  is defined, with an implicit dependency over the parameters  $\mathbf{u}$ , so that we can in fact label the loss as  $J(\mathbf{u})$ .

Then, posed in the framework of steepest descent, performing a single GD step for this NN over the parameters  $\mathbf{u}$  can be seen as a Bolza problem ([Bliss 1968](#)) that minimizes the following objective for the current input  $\mathbf{x}_0^k$  and current parameters  $\mathbf{u}^k$ :

$$\min_{\Delta \mathbf{u}} \nabla J(\mathbf{u}^k)^T(\Delta \mathbf{u}) + \frac{1}{2}(\Delta \mathbf{u})^T(\Delta \mathbf{u}) = \min_{\mathbf{u}} \nabla J(\mathbf{u}^k)^T(\mathbf{u}^k - \mathbf{u}) + \frac{1}{2} \sum_{i=0}^{N-1} (\mathbf{u}_i^k - \mathbf{u}_i)^T(\mathbf{u}_i^k - \mathbf{u}_i) \quad (3.1.1)$$

The above problem is quadratic and thus convex, a unique optimal solution exists given by  $\Delta \mathbf{u}^*$  (or  $\mathbf{u}^*$ ). The GD update to the parameters for the current input is then given by:

$$\mathbf{u}^{k+1} = \mathbf{u}^k + \Delta \mathbf{u}^* = \mathbf{u}^*$$

We will now show that (3.1.1) is in fact a proper LQR problem. Already, we can notice that a quadratic cost is present; what is lacking is a linear dynamic system describing the system evolution. But, we will see that the linearization arise naturally from the differentiation, that is from the term  $\nabla J(\mathbf{u}^k)^T(\Delta \mathbf{u})$  in (3.1.1).

Intuitively, this property can be understood easily from the multivariate chain rule. Indeed, if we think of a given system with  $\mathbf{x}_1 = \phi(\mathbf{u})$  and  $\mathbf{x}_2 = g(\mathbf{x}_1, \mathbf{u})$ , we have :

$$\begin{aligned} \mathcal{D}_{\mathbf{u}}(f(\mathbf{x}_2, \mathbf{u}) \circ g)(\mathbf{x}, \mathbf{u}) &= \mathcal{D}_{\mathbf{x}_2} f \mathcal{D}_{\mathbf{u}} g + \mathcal{D} f_{\mathbf{x}_2} \mathcal{D}_{\mathbf{x}_1} g \mathcal{D}_{\mathbf{u}} \phi + \mathcal{D}_{\mathbf{u}} f \\ &= \mathcal{D}_{\mathbf{x}_2} f \left( \mathcal{D}_{\mathbf{x}_1} g \mathcal{D}_{\mathbf{u}} \phi + \mathcal{D}_{\mathbf{u}} g \right) + \mathcal{D}_{\mathbf{u}} f \\ &= \mathbf{A}_2 (\mathbf{A}_1 \cdot \mathbf{C} + \mathbf{B}_1) + \mathbf{B}_2 \end{aligned}$$

That is, differentiation and linearization are closely related. Moreover, the above can be seen as a discrete-time linear dynamical system.

More formally, by applying the chain rule to the first term in (3.1.1) and denoting  $\nabla^T J(\mathbf{u}^k)$  by  $\nabla^T J^k$  (and  $f_i(\mathbf{x}_i^k, \mathbf{u}_i^k)$  by  $f_i^k$ ) we have that:

$$\begin{aligned} \nabla_{\mathbf{u}}^T J^k(\mathbf{u} - \mathbf{u}^k) &= \nabla_{\mathbf{x}_N}^T J^k (\nabla_{\mathbf{x}_{N-1}}^T f_{N-1}^k \nabla_{\mathbf{u}_0:N-2}^T f_{N-1}^k + \nabla_{\mathbf{u}_{N-1}}^T f_{N-1}^k) (\mathbf{u} - \mathbf{u}^k) \\ &= \nabla_{\mathbf{x}_N}^T J^k (\nabla_{\mathbf{x}_{N-1}}^T f_{N-1}^k (\nabla_{\mathbf{x}_{N-2}}^T f_{N-2}^k \nabla_{\mathbf{u}_0:N-3}^T f_{N-3}^k + \nabla_{\mathbf{u}_{N-2}}^T f_{N-2}^k) + \nabla_{\mathbf{u}_{N-1}}^T f_{N-1}^k) (\mathbf{u} - \mathbf{u}^k) \\ &= \vdots \\ &= \nabla_{\mathbf{x}_N}^T J^k \left( \dots \nabla_{\mathbf{x}_1}^T f_1^k (\nabla_{\mathbf{x}_0}^T f_0^k \nabla_{\mathbf{u}_0}^T f_0^k + \nabla_{\mathbf{u}_0}^T f_0^k) + \nabla_{\mathbf{u}_1}^T f_1^k \dots \right) (\mathbf{u} - \mathbf{u}^k) \end{aligned} \tag{3.1.2}$$

We chose to denote the Jacobians via the matrix of gradients,  $\nabla_{\mathbf{v}_i} f_i$ , for which each columns are the gradients of the elements of the vector function  $f_i$ . As such, we have that the Jacobian of  $f_i$  is the transpose of this matrix, i.e.  $\frac{\partial f_i(\mathbf{v}_i)}{\partial \mathbf{v}_i} = \nabla_{\mathbf{v}_i}^T f_i$ . Moreover, the notation  $\nabla_{\mathbf{u}_0:i}^T$  is used to indicate that the derivative is taken with respect to the parameters  $\mathbf{u}_0$  up to  $\mathbf{u}_i$ .

Now we note that  $\mathbf{x}_0$  is constant, so  $\nabla_{\mathbf{x}_0}^T f_1^k \nabla_{\mathbf{u}}^T f_0 = 0$  and also that we can distribute  $(\mathbf{u} - \mathbf{u}^k)$  over the sum, taking into account that  $f_i$  has only partial derivatives over  $\mathbf{u}_i$  and  $\mathbf{x}_i$ :

$$\nabla_{\mathbf{u}}^T J^k(\mathbf{u} - \mathbf{u}^k) = \nabla_{\mathbf{x}_N}^T J^k \left( \dots \nabla_{\mathbf{x}_2}^T f_2^k \left( \nabla_{\mathbf{x}_1}^T f_1^k \left( \nabla_{\mathbf{u}_0}^T f_0^k \right) (\mathbf{u}_0 - \mathbf{u}_0^k) + \nabla_{\mathbf{u}_1}^T f_1^k (\mathbf{u}_1 - \mathbf{u}_1^k) \right) + \nabla_{\mathbf{u}_2}^T f_2^k (\mathbf{u}_2 - \mathbf{u}_2^k) \dots \right) \quad (3.1.3)$$

Finally, by adopting the notation  $\mathbf{A}_i = \nabla_{\mathbf{x}_i}^T f_i^k$ ,  $\mathbf{B}_i = \nabla_{\mathbf{u}_i}^T f_i^k$ ;  $\tilde{\mathbf{u}}_i = \mathbf{u}_i - \mathbf{u}_i^k$  and  $\tilde{\mathbf{x}}_{i+1} = \mathbf{A}_i \tilde{\mathbf{x}}_i + \mathbf{B}_i \tilde{\mathbf{u}}_i$  we get:

$$\begin{aligned} \nabla_{\mathbf{u}}^T J^k(\mathbf{u} - \mathbf{u}^k) &= \nabla_{\mathbf{x}_N}^T J^k \left( \dots \mathbf{A}_2 \left( \mathbf{A}_1 (\mathbf{B}_0 \tilde{\mathbf{u}}_0) + \mathbf{B}_1 \tilde{\mathbf{u}}_1 \right) + \mathbf{B}_2 \tilde{\mathbf{u}}_2 \dots \right) \\ &= \nabla_{\mathbf{x}_N}^T J^k \cdot \tilde{\mathbf{x}}_N \end{aligned} \quad (3.1.4)$$

As it will be useful later, note that if we denote explicitly the implicit dependence of  $\mathbf{x}$  over  $\mathbf{u}$  via  $\mathbf{x} = \phi(\mathbf{u})$  (so that  $J(\mathbf{x}, \mathbf{u}) = J(\phi(\mathbf{u}), \mathbf{u})$ ), then (3.1.4) shows that:

$$\nabla_{\mathbf{u}}^T \phi(\mathbf{u}) \cdot (\mathbf{u} - \mathbf{u}^k) = \tilde{\mathbf{x}}_N \quad (3.1.5)$$

Since the cost  $J$  has only a direct dependency over the variable  $\mathbf{x}_N$  and  $\nabla_{\mathbf{u}}^T J^k(\mathbf{u} - \mathbf{u}^k) = \nabla_{\mathbf{x}_N}^T J^k \left( \nabla_{\mathbf{u}}^T \phi(\mathbf{u}) \cdot (\mathbf{u} - \mathbf{u}^k) \right)$ .

### 3.1.1 Equivalence with Gradient Descent

What we have in (3.1.4) is a discrete linear dynamical system. Indeed the transition layers of this system are obtained via the linearization of the original layers in the model. Thus, (3.1.1) can be rewritten as a standard OCP, allowing us to put forward a preliminary result of this thesis, obtained by construction from the above problem statement:

**Theorem 3.1.1.** *Retrieving the update for a single gradient descent step over the parameters  $\mathbf{u}^k$  for a given example  $\mathbf{x}^k$  in a neural network is equivalent to solving the following linear quadratic regulator problem:*

$$\begin{aligned} \min_{\tilde{\mathbf{u}}} \quad & \nabla_{\mathbf{x}_N}^T J(\mathbf{u}^k) \times \tilde{\mathbf{x}}_N + \frac{1}{2} \sum_{i=1}^{N-1} \tilde{\mathbf{u}}_i^T \tilde{\mathbf{u}}_i \\ \text{subject to} \quad & \tilde{\mathbf{x}}_{i+1} = \mathbf{A}_i \tilde{\mathbf{x}}_i + \mathbf{B}_i \tilde{\mathbf{u}}_i \\ \text{where} \quad & \mathbf{A}_i = \nabla_{\mathbf{x}_i}^T f_i^k, \mathbf{B}_i = \nabla_{\mathbf{u}_i}^T f_i^k \\ \text{given} \quad & \tilde{\mathbf{x}}_0 = 0 \end{aligned} \quad (3.1.6)$$

Based on the material presented in section 2.3, we can see that (3.1.6) is a proper LQR problem: the cost to minimize is quadratic and the transition functions are linear. Therefore, the optimal solution,  $\tilde{\mathbf{u}}^*$  to (3.1.6) can be obtained by solving the Riccati equation (2.3.4), starting with  $\mathbf{K}_N = 0$  and  $\boldsymbol{\lambda}_N = \nabla_{\mathbf{x}_N}^T J(\mathbf{u}^k)$

$$\begin{aligned}\mathbf{K}_i &= \mathbf{A}_i^T \mathbf{K}_{i+1} \mathbf{A}_i - (\mathbf{A}_i^T \mathbf{K}_{i+1} \mathbf{B}_i)(\mathbf{I} + \mathbf{B}_i^T \mathbf{K}_{i+1} \mathbf{B}_i)^{-1} (\mathbf{B}_i^T \mathbf{K}_{i+1} \mathbf{A}_i) \\ \boldsymbol{\lambda}_i &= \mathbf{A}_i^T \boldsymbol{\lambda}_{i+1} - (\mathbf{A}_i^T \mathbf{K}_{i+1} \mathbf{B}_i)(\mathbf{I} + \mathbf{B}_i^T \mathbf{K}_{i+1} \mathbf{B}_i)^{-1} \mathbf{B}_i^T \boldsymbol{\lambda}_{i+1}\end{aligned}\tag{3.1.7}$$

The equivalent GD update is then given by:

$$\begin{aligned}\tilde{\mathbf{u}}_0^* &= -(\mathbf{I} + \mathbf{B}_0^T \mathbf{K}_1 \mathbf{B}_0)^{-1} \mathbf{B}_0^T \boldsymbol{\lambda}_1 \\ \tilde{\mathbf{x}}_{i+1}^* &= \mathbf{A}_i \tilde{\mathbf{x}}_i^* + \mathbf{B}_i \tilde{\mathbf{u}}_i^* \\ \tilde{\mathbf{u}}_i^* &= -(\mathbf{I} + \mathbf{B}_i^T \mathbf{K}_{i+1} \mathbf{B}_i)^{-1} \left( (\mathbf{B}_i^T \mathbf{K}_{i+1} \mathbf{A}_i) \tilde{\mathbf{x}}_i^* + \mathbf{B}_i^T \boldsymbol{\lambda}_{i+1} \right)\end{aligned}\tag{3.1.8}$$

With a learning rate  $\alpha$ , the GD update can therefore be applied via:

$$\mathbf{u}_{k+1} = \mathbf{u}_k + \alpha \tilde{\mathbf{u}}^*\tag{3.1.9}$$

Note that the initial condition  $\tilde{\mathbf{x}}_0 = 0$  arises from  $\nabla_{\mathbf{x}_0}^T f_1^k \nabla_{\mathbf{u}}^T f_0 = 0$ ; defining the subsequent  $\tilde{\mathbf{x}}_i$  by  $\tilde{\mathbf{x}}_{i+1} = \mathbf{A}_i \tilde{\mathbf{x}}_i + \mathbf{B}_i \tilde{\mathbf{u}}_i$  is then a convenient choice to describe the state transitions happening in the linear system of (3.1.4).

Also, it is worth noting that the learning rate for the GD step can equivalently be incorporated in the optimization problem (3.1.6) by minimizing instead

$$\min_{\tilde{\mathbf{u}}} \nabla_{\mathbf{x}_N}^T J(\mathbf{u}^k) \times \tilde{\mathbf{x}}_N + \frac{1}{2} \alpha \sum_{i=1}^{N-1} \tilde{\mathbf{u}}_i^T \tilde{\mathbf{u}}_i$$

and the solution follows as above, but taking  $\alpha \mathbf{I}$  instead of  $\mathbf{I}$  in (3.1.7) and (3.1.8).

## 3.2 Steepest Descent Objective

The result shown in Theorem 3.1.1 was obtained by minimizing the steepest descent objective under the Euclidean norm (see section 2.2.1). But, the steepest descent procedure extends well beyond gradient descent and will lead to other optimization algorithms (such as NGD and Newton's method) when the divergence measure is chosen differently. We will now show

that an equivalent LQR can be obtained to the steepest descent problem for a wide choice of divergence measures.

### 3.2.1 General LQR Equivalence

To perform a steepest descent step under a given P-norm,  $\|\cdot\|_P$ , we need to solve:

$$\min_{\mathbf{v}} \nabla f(x)^T \mathbf{v} + \frac{1}{2} \|\mathbf{v}\|_P$$

The P-norm is often induced by a divergence measure, and for such cases, we can recover an equivalent LQR. We provide this result below, by simply extending the demonstration of the equivalence between a LQR problem and one Newton's method step stated in Bertsekas 2016. The use of the notation  $H$  in (3.2.3) refers to the *Hamiltonian* defined in the development of the Pontryagin's maximum principle (PMP) Pontriagin *et al.* 1962.

**Theorem 3.2.1.** *If the matrix  $P(\mathbf{x}, \mathbf{u})$  inducing the quadratic norm  $\|\cdot\|_P$  can be written as  $P(\mathbf{x}, \mathbf{u}) = \nabla_{\mathbf{u}\mathbf{u}}^2 D(\mathbf{x}, \mathbf{u})$  and  $D$  is a function that can be expressed layer-wise (i.e.  $D(\mathbf{x}, \mathbf{u}) = D_N(\mathbf{x}_N) + \sum_i^{N-1} D_i(\mathbf{x}_i, \mathbf{u}_i)$ ), then the minimization problem*

$$\begin{aligned} \min_{\mathbf{v}} \nabla^T J(F(\mathbf{x})) \mathbf{v} + \frac{1}{2} \|\mathbf{v}\|_P \\ = \min_{\tilde{\mathbf{u}}} \nabla^T J(\mathbf{u}) \tilde{\mathbf{u}} + \frac{1}{2} \tilde{\mathbf{u}}^T (\nabla_{\mathbf{u}\mathbf{u}}^2 D(\mathbf{x}, \mathbf{u})) \tilde{\mathbf{u}} \end{aligned} \quad (3.2.1)$$

is equivalent to the following LQR when  $F$  is the dynamical system of an OCP (for example, when  $F$  is a NN) and  $J$  the associated cost :

$$\begin{aligned} \min_{\tilde{\mathbf{u}}} \nabla_{\mathbf{x}_N}^T J(\mathbf{u}^k) \times \tilde{\mathbf{x}}_N + \frac{1}{2} \{ \tilde{\mathbf{x}}^T \nabla_{\mathbf{x}\mathbf{x}}^2 H(\phi(\mathbf{u}), \mathbf{u}, p(\mathbf{u})) \tilde{\mathbf{x}} + \tilde{\mathbf{x}}^T \nabla_{\mathbf{x}\mathbf{u}}^2 H(\phi(\mathbf{u}), \mathbf{u}, p(\mathbf{u})) \tilde{\mathbf{u}} + \\ \tilde{\mathbf{u}}^T \nabla_{\mathbf{u}\mathbf{x}}^2 H(\phi(\mathbf{u}), \mathbf{u}, p(\mathbf{u})) \tilde{\mathbf{x}} + \tilde{\mathbf{u}}^T \nabla_{\mathbf{u}\mathbf{u}}^2 H(\phi(\mathbf{u}), \mathbf{u}, p(\mathbf{u})) \tilde{\mathbf{u}} \} \end{aligned} \quad (3.2.2)$$

subject to  $\tilde{\mathbf{x}}_{i+1} = \mathbf{A}_i \tilde{\mathbf{x}}_i + \mathbf{B}_i \tilde{\mathbf{u}}_i$

where  $\mathbf{A}_i = \nabla_{\mathbf{x}_i}^T f_i^k$ ,  $\mathbf{B}_i = \nabla_{\mathbf{u}_i}^T f_i^k$

given  $\tilde{\mathbf{x}}_0 = 0$

Where:

$$\phi(\mathbf{u}) = [\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N]^T$$



$$H(\phi(\mathbf{u}), \mathbf{u}, p(\mathbf{u})) = D(\phi(\mathbf{u}), \mathbf{u}) + h(\phi(\mathbf{u}), \mathbf{u})^T p(\mathbf{u}) \quad (3.2.3)$$

$$h(\phi(\mathbf{u}), \mathbf{u}) = \begin{bmatrix} f_0(\mathbf{u}_0, \mathbf{x}_0) - \mathbf{x}_1 \\ \vdots \\ f_{N-1}(\mathbf{u}_{N-1}, \mathbf{x}_{N-1}) - \mathbf{x}_N \end{bmatrix} = \vec{0}$$

And

$$p(\mathbf{u}) = \begin{bmatrix} \nabla_{\mathbf{x}_1} f_1 \cdot \mathbf{p}_2 + \nabla_{\mathbf{x}_1} D \\ \vdots \\ \nabla_{\mathbf{x}_{N-1}} f_{N-1} \cdot \mathbf{p}_N + \nabla_{\mathbf{x}_{N-1}} D \\ \nabla_{\mathbf{x}_N} D \end{bmatrix}$$

**Proof.** When going over the proof, it is important to realize that the objective to minimize in (3.2.1) is *not* a valid cost function for LQR problems, because it is not decoupled with respect to the controls and the states. This is why introducing the Hamiltonian is so helpful: the Hamiltonian is such that the controls are minimizers at each transition step (i.e. at every layer) (see [Lenhart and Workman, 2007](#)). It is this property that allows us to cast the optimization objective as *a sum of decoupled terms*, thus making it a proper objective for a LQR problem.

We already know that the term  $\nabla^T J(\mathbf{u}) \tilde{\mathbf{u}}$  gives rise to the "skeleton" of the LQR, so to complete the proof we focus on showing that

$$\begin{aligned} \tilde{\mathbf{u}}^T (\nabla_{\mathbf{u}\mathbf{u}}^2 D(\mathbf{x}, \mathbf{u})) \tilde{\mathbf{u}} &= \tilde{\mathbf{x}}^T \nabla_{\mathbf{x}\mathbf{x}}^2 H(\phi(\mathbf{u}), \mathbf{u}, p(\mathbf{u})) \tilde{\mathbf{x}} + \tilde{\mathbf{x}}^T \nabla_{\mathbf{x}\mathbf{u}}^2 H(\phi(\mathbf{u}), \mathbf{u}, p(\mathbf{u})) \tilde{\mathbf{u}} \\ &\quad + \tilde{\mathbf{u}}^T \nabla_{\mathbf{u}\mathbf{x}}^2 H(\phi(\mathbf{u}), \mathbf{u}, p(\mathbf{u})) \tilde{\mathbf{x}} + \tilde{\mathbf{u}}^T \nabla_{\mathbf{u}\mathbf{u}}^2 H(\phi(\mathbf{u}), \mathbf{u}, p(\mathbf{u})) \tilde{\mathbf{u}} \end{aligned}$$

To see this, we need to view  $D$  as an "alternate loss term" for which we can develop the appropriate adjoint terms; an approach inspired by the PMP.

Since  $h$  is a zero valued vector, we have that  $\forall \mathbf{p}$ ,  $D(\phi(\mathbf{u}), \mathbf{u}) = H(\phi(\mathbf{u}), \mathbf{u}, p(\mathbf{u}))$  when  $H$  is defined as:

$$H(\phi(\mathbf{u}), \mathbf{u}, p(\mathbf{u})) = D(\phi(\mathbf{u}), \mathbf{u}) + h(\phi(\mathbf{u}), \mathbf{u})^T p(\mathbf{u})$$

By taking the first derivative with respect to  $\mathbf{u}$  and applying the chain rule, we get

$$\nabla_{\mathbf{u}}H = \nabla_{\mathbf{u}}\phi(\mathbf{u})\nabla_{\mathbf{x}}H(\phi(\mathbf{u}),\mathbf{u},p(\mathbf{u})) + \nabla_{\mathbf{u}}H(\phi(\mathbf{u}),\mathbf{u},p(\mathbf{u}))$$

This is true for any  $\mathbf{p}$ , and in particular for the given  $\mathbf{p}$  s.t.  $\nabla_{\mathbf{x}}H(\phi(\mathbf{u}),\mathbf{u},p(\mathbf{u})) = 0$ . This means that we can rewrite  $\mathbf{p}_i$  as being

$$\begin{aligned}\mathbf{p}_i &= \nabla_{x_i}f_i \cdot \mathbf{p}_{i+1} + \nabla_{x_i}D \\ \mathbf{p}_N &= \nabla_{x_N}D\end{aligned}$$

Taking the derivative of  $\nabla_{\mathbf{u}}H$  once more, we now get:

$$\begin{aligned}\nabla_{\mathbf{u}}^2L &= \nabla_{\mathbf{u}}^2\phi(\mathbf{u})\nabla_{\mathbf{x}}H(\phi(\mathbf{u}),\mathbf{u},p(\mathbf{u})) + \nabla_{\mathbf{u}}\phi(\mathbf{u})\nabla_{\mathbf{x}\mathbf{x}}H(\phi(\mathbf{u}),\mathbf{u},p(\mathbf{u}))\nabla_{\mathbf{u}}^T\phi(\mathbf{u}) \\ &\quad + \nabla_{\mathbf{u}}\phi(\mathbf{u})\nabla_{\mathbf{x}\mathbf{u}}H(\phi(\mathbf{u}),\mathbf{u},p(\mathbf{u})) + \nabla_{\mathbf{u}\mathbf{x}}H(\phi(\mathbf{u}),\mathbf{u},p(\mathbf{u}))\nabla_{\mathbf{u}}^T\phi(\mathbf{u}) \\ &\quad + \nabla_{\mathbf{u}\mathbf{u}}H(\phi(\mathbf{u}),\mathbf{u},p(\mathbf{u}))\end{aligned}$$

Remembering that  $\tilde{\mathbf{u}}^T\nabla_{\mathbf{u}}\phi(\mathbf{u}) = \tilde{\mathbf{x}}^T$  (3.1.5) and  $\nabla_{\mathbf{u}}^2\phi(\mathbf{u})\nabla_{\mathbf{x}}H(\phi(\mathbf{u}),\mathbf{u},p(\mathbf{u})) = 0$  because of our choice of  $p(\mathbf{u})$ , we get the desired result. □

The conditions imposed on  $D$  in Theorem 3.2.1 may seem quite restrictive, but we argue that it should not be the case in practice. Indeed, in most scenario, the divergence is only measured over the final state, i.e.:  $D(\mathbf{x},\mathbf{u}) = D_N(\mathbf{x}_N)$ . Moreover, the divergence measure can often be approximated (sometimes exactly) via their second-order Taylor's expansion (again, [Nocedal and Wright, 2006](#)):

**Theorem 3.2.2.** *For a divergence measure  $D \in \mathbf{C}^2$  having the following properties:*

$$D(\mathbf{x},\mathbf{x}) = 0 ; D(\mathbf{x},\mathbf{y}) \geq 0$$

*We have that:*

$$D(\mathbf{u}) \approx \mathbf{u}^T \left( \nabla_{\mathbf{u}\mathbf{u}}^2 D(\mathbf{u}') \right) \mathbf{u} \tag{3.2.4}$$

**Proof.** Taking the second order Taylor's expansion of a divergence measure  $D(\mathbf{u}) = D(\mathbf{u}',\mathbf{u})$  around  $\mathbf{u} = \mathbf{u}'$ :

$$D(\mathbf{u}) \approx D(\mathbf{u}') + \nabla_{\mathbf{u}} D(\mathbf{u}') \mathbf{u} + \mathbf{u}^T (\nabla_{\mathbf{u}\mathbf{u}}^2 D(\mathbf{u}')) \mathbf{u}$$

By definition  $D(\mathbf{u}') = D(\mathbf{u}', \mathbf{u}') = 0$ , and  $\nabla_{\mathbf{u}} D(\mathbf{u}') = 0$  as well since  $D(\mathbf{u})$  is minimal at  $\mathbf{u}'$ .  $\square$

### 3.2.2 Riccati's Solution to the Steepest Descent Objective

Since (3.2.2) is a LQR problem, just as for the GD case, we can recover the steepest descent update by solving the LQR with Riccati's equation. By defining  $H_i$  as:

$$H_i(\mathbf{x}_i, \mathbf{u}, \mathbf{p}_i) = D_i(\mathbf{x}_i, \mathbf{u}_i) + h_i(\mathbf{x}_i, \mathbf{u}_i)^T \mathbf{p}_i$$

And adopting the notation  $\mathbf{Q}_N = \nabla^2 D_N(\mathbf{x}_N)$ ,  $\mathbf{a}_N = \nabla_{\mathbf{x}_N} J(\mathbf{u}^k)$ ,  $\mathbf{Q}_i = \nabla_{\mathbf{x}_i \mathbf{x}_i}^2 H_i$ ,  $\mathbf{R}_i = \nabla_{\mathbf{u}_i \mathbf{u}_i}^2 H_i$ ,  $\mathbf{M}_i = \nabla_{\mathbf{u}_i \mathbf{x}_i}^2 H_i$ , we can rewrite (3.2.2) more succinctly as:

$$\begin{aligned} \min_{\tilde{\mathbf{u}}} \quad & \mathbf{a}_N^T \tilde{\mathbf{x}}_N + \frac{1}{2} \tilde{\mathbf{x}}_N^T \mathbf{Q}_N \tilde{\mathbf{x}}_N + \sum_{i=0}^{N-1} \left( \frac{1}{2} \tilde{\mathbf{x}}_i^T \mathbf{Q}_i \tilde{\mathbf{x}}_i + \frac{1}{2} \tilde{\mathbf{u}}_i^T \mathbf{R}_i \tilde{\mathbf{u}}_i + \tilde{\mathbf{u}}_i^T \mathbf{M}_i \tilde{\mathbf{x}}_i \right) \\ \text{subject to} \quad & \tilde{\mathbf{x}}_{i+1} = \mathbf{A}_i \tilde{\mathbf{x}}_i + \mathbf{B}_i \tilde{\mathbf{u}}_i \\ \text{where} \quad & \mathbf{A}_i = \nabla_{\mathbf{x}_i}^T f_i^k, \mathbf{B}_i = \nabla_{\mathbf{u}_i}^T f_i^k \\ \text{given} \quad & \tilde{\mathbf{x}}_0 = 0 \end{aligned} \tag{3.2.5}$$

a problem for which the optimal solution is obtained by solving backward:

$$\mathbf{K}_i = \mathbf{A}_i^T \mathbf{K}_{i+1} \mathbf{A}_i + \mathbf{Q}_i - (\mathbf{A}_i^T \mathbf{K}_{i+1} \mathbf{B}_i + \mathbf{M}_i^T) (\mathbf{R}_i + \mathbf{B}_i^T \mathbf{K}_{i+1} \mathbf{B}_i)^{-1} (\mathbf{M}_i + \mathbf{B}_i^T \mathbf{K}_{i+1} \mathbf{A}_i) \tag{3.2.6}$$

$$\boldsymbol{\lambda}_i = \mathbf{A}_i^T \boldsymbol{\lambda}_{i+1} - (\mathbf{A}_i^T \mathbf{K}_{i+1} \mathbf{B}_i + \mathbf{M}_i^T) (\mathbf{R}_i + \mathbf{B}_i^T \mathbf{K}_{i+1} \mathbf{B}_i)^{-1} \mathbf{B}_i^T \boldsymbol{\lambda}_{i+1} , \tag{3.2.7}$$

with  $\mathbf{K}_N = \mathbf{Q}_N$  and  $\boldsymbol{\lambda}_N = \mathbf{a}_N$ .

The updates are then:

$$\tilde{\mathbf{u}}_0^* = -(\mathbf{R}_0 + \mathbf{B}_0^T \mathbf{K}_1 \mathbf{B}_0)^{-1} \mathbf{B}_0^T \lambda_1 \quad (3.2.8)$$

$$\tilde{\mathbf{x}}_{i+1}^* = \mathbf{A}_i \tilde{\mathbf{x}}_i^* + \mathbf{B}_i \tilde{\mathbf{u}}_i^* \quad (3.2.9)$$

$$\tilde{\mathbf{u}}_i^* = -(\mathbf{R}_i + \mathbf{B}_i^T \mathbf{K}_{i+1} \mathbf{B}_i)^{-1} \left( (\mathbf{M}_i + \mathbf{B}_i^T \mathbf{K}_{i+1} \mathbf{A}_i) \tilde{\mathbf{x}}_i^* + \mathbf{B}_i^T \lambda_{i+1} \right) . \quad (3.2.10)$$

### 3.2.3 Equivalence with Natural Gradients

Natural gradient descent [Amari 1998](#) is an optimization algorithm well suited to the modeling of probability measures and has known advantages over GD in the context of deep learning (e.g. [Pascanu and Bengio 2014](#)). NGD can be stated in the language of steepest descent as the search for the optimal unitary step in a space scaled by the Fisher’s information matrix ( $\mathbf{F}$ ):

$$\begin{aligned} \Delta \mathbf{v} &= \arg \min_{\mathbf{v}} \nabla f(\mathbf{u})^T \mathbf{v} + \frac{1}{2} \|\mathbf{v}\|_F^2 \\ &= \arg \min_{\mathbf{v}} \nabla f(\mathbf{u})^T \mathbf{v} + \frac{1}{2} \mathbf{v}^T \mathbf{F} \mathbf{v} \end{aligned} \quad (3.2.11)$$

leading to the NGD update rule for some parameters  $\mathbf{u}$ :

$$\mathbf{u}^{k+1} = \mathbf{u}^k - \mathbf{F}^{-1} \nabla f(\mathbf{u})$$

Furthermore, there is a well-known connection between the Hessian of the Kullback–Leibler (KL) divergence and the Fisher’s information matrix, allowing to develop the second-order Taylor’s approximation of the KL divergence as:

$$\begin{aligned} D_{\text{KL}}(p_{\mathbf{u}} \| p_{\mathbf{u}+\mathbf{v}}) &\approx D_{\text{KL}}(p_{\mathbf{u}} \| p_{\mathbf{u}}) + \nabla_{\mathbf{u}'}^T D_{\text{KL}}(p_{\mathbf{u}} \| p_{\mathbf{u}'}) \Big|_{\mathbf{u}'=\mathbf{u}} \cdot \mathbf{v} + \frac{1}{2} \mathbf{v}^T \left( \nabla_{\mathbf{u}'\mathbf{u}'}^2 D_{\text{KL}}(p_{\mathbf{u}} \| p_{\mathbf{u}'}) \Big|_{\mathbf{u}'=\mathbf{u}} \right) \mathbf{v} \\ &= D_{\text{KL}}(p_{\mathbf{u}} \| p_{\mathbf{u}}) + \nabla_{\mathbf{u}'}^T D_{\text{KL}}(p_{\mathbf{u}} \| p_{\mathbf{u}'}) \Big|_{\mathbf{u}'=\mathbf{u}} \cdot \mathbf{v} + \frac{1}{2} \mathbf{v}^T \mathbf{F} \mathbf{v} \end{aligned}$$

The first two terms are zero (see Theorem 3.2.2), which allows us to equivalently state the steepest descent step above as (see [Amari and Nagaoka 2007](#)):

$$\Delta \mathbf{v} = \arg \min_{\mathbf{v}} \nabla f(\mathbf{u})^T \mathbf{v} + \frac{1}{2} \mathbf{v}^T \left( \nabla_{\mathbf{u}'\mathbf{u}'}^2 D_{\text{KL}}(p_{\mathbf{u}} \| p_{\mathbf{u}'}) \Big|_{\mathbf{u}'=\mathbf{u}} \right) \mathbf{v} \quad (3.2.12)$$

This form is an exact fit to the formalism developed in Theorem 3.2.1. Thus, one can recover the NGD update step by optimizing the following LQR:

$$\begin{aligned}
& \min_{\tilde{\mathbf{u}}} \mathbf{a}_N^T \tilde{\mathbf{x}}_N + \frac{1}{2} \tilde{\mathbf{x}}_N^T \mathbf{Q}_N \tilde{\mathbf{x}}_N + \sum_{i=0}^{N-1} \left( \frac{1}{2} \tilde{\mathbf{x}}_i^T \mathbf{Q}_i \tilde{\mathbf{x}}_i + \frac{1}{2} \tilde{\mathbf{u}}_i^T \mathbf{R}_i \tilde{\mathbf{u}}_i + \tilde{\mathbf{u}}_i^T \mathbf{M}_i \tilde{\mathbf{x}}_i \right) \\
& \text{subject to } \tilde{\mathbf{x}}_{i+1} = \mathbf{A}_i \tilde{\mathbf{x}}_i + \mathbf{B}_i \tilde{\mathbf{u}}_i \\
& \text{where } \mathbf{A}_i = \nabla_{\mathbf{x}_i}^T f_i^k, \mathbf{B}_i = \nabla_{\mathbf{u}_i}^T f_i^k \\
& \text{given } \tilde{\mathbf{x}}_0 = 0
\end{aligned} \tag{3.2.13}$$

where

$$\begin{aligned}
\mathbf{Q}_N &= \nabla_{\mathbf{x}'_N}^2 D_{\text{KL}}(p_{\mathbf{x}_N} || p_{\mathbf{x}'_N}) \Big|_{\mathbf{x}'_N = \mathbf{x}_N}; \mathbf{a}_N = \nabla_{\mathbf{x}_N} J(\mathbf{u}^k), \\
\mathbf{Q}_i &= \nabla_{\mathbf{x}_i \mathbf{x}_i}^2 H_i; \mathbf{R}_i = \nabla_{\mathbf{u}_i \mathbf{u}_i}^2 H_i; \mathbf{M}_i = \nabla_{\mathbf{u}_i \mathbf{x}_i}^2 H_i
\end{aligned}$$

and with  $H_i = f_i(\mathbf{x}^k, \mathbf{u}^k)^T \mathbf{p}_i$  simply here since the KL-divergence has a direct dependency only over  $\mathbf{x}_N$ , the final state of the dynamic system that is a representation of the distribution we are trying to model.

The solution can be recovered via Riccati equation directly with (3.2.6) and (3.2.8). It is worth noticing that the solution is obtained *without* retrieving the Fisher's information matrix; indeed one only need to recover the Hessian of the KL divergence with respect to the final state (much simpler than recovering the Hessian with respect to the whole system parameters) before applying the Riccati equation to get the gain matrix  $\mathbf{K}$ .

Another advantage of the method is that we do not need to invert the Fisher to recover the NGD update, but instead, we perform  $N - 1$  matrix inversion of the quantity  $(\mathbf{R}_i + \mathbf{B}_i^T \mathbf{K}_{i+1} \mathbf{B}_i)$ . When the dynamic system is a NN and with  $n_l$  parameters per layer (the total number of parameters is  $n = \sum_l n_l$ ), this means we must invert  $N - 1$  matrices of size  $n_l \times n_l$  instead of inverting the Fisher of size  $n \times n$ .

### 3.2.4 Equivalence with Newton's method

As mention earlier, the motivation for this thesis came from the realization that the connection between the use of Newton's method for OCP optimization and LQR formalism generalizes to many steepest descent problems. For completeness, we develop this equivalence (originally shown in Bertsekas 2016; Dunn and Bertsekas 1989; De O. Pantoja and

Mayne 1989) below.

Newton's method can be easily stated as a steepest descent problem by optimizing the second-order Taylor's approximation of our function around the current parameters ( $\mathbf{u}^k$ ):

$$\begin{aligned} J(\mathbf{u}) &\approx J^k + \nabla_{\mathbf{u}}^T J^k (\mathbf{u} - \mathbf{u}^k) + \frac{1}{2} (\mathbf{u} - \mathbf{u}^k)^T \nabla_{\mathbf{u}}^2 J^k (\mathbf{u} - \mathbf{u}^k) \\ &= J^k + \nabla_{\mathbf{u}}^T J^k (\mathbf{u} - \mathbf{u}^k) + \frac{1}{2} (\mathbf{u} - \mathbf{u}^k)^T \mathbf{H} (\mathbf{u} - \mathbf{u}^k) \end{aligned} \quad , \quad (3.2.14)$$

where we used  $\mathbf{H}$  to denote the Hessian. If  $\mathbf{H}$  is positive semi-definite, then we are facing a convex problem and (3.2.14) minimum can be obtained by setting the derivative to 0. We get:

$$\begin{aligned} 0 &= \frac{\partial}{\partial \mathbf{u}} \left( J^k + \nabla_{\mathbf{u}}^T J^k (\mathbf{u} - \mathbf{u}^k) + \frac{1}{2} (\mathbf{u} - \mathbf{u}^k)^T \mathbf{H} (\mathbf{u} - \mathbf{u}^k) \right) \Big|_{\mathbf{u}=\mathbf{u}^*} \\ &= \nabla_{\mathbf{u}}^T J^k + \mathbf{H} (\mathbf{u}^* - \mathbf{u}^k) \\ &\implies \mathbf{u}^* = -\mathbf{H}^{-1} \nabla_{\mathbf{u}}^T J^k \end{aligned} .$$

Thus, to recover the Newton's update, we can solve the following steepest descent problem:

$$\Delta \mathbf{u} = \arg \min_{\mathbf{u}} \nabla_{\mathbf{u}}^T J^k (\mathbf{u} - \mathbf{u}^k) + \frac{1}{2} (\mathbf{u} - \mathbf{u}^k)^T \mathbf{H} (\mathbf{u} - \mathbf{u}^k) \quad (3.2.15)$$

If we are trying to optimize a NN, the equivalence with LQR problems developed in 2.3 allows us to express (3.2.15) as :

$$\begin{aligned} \min_{\tilde{\mathbf{u}}} & \mathbf{a}_N^T \tilde{\mathbf{x}}_N + \frac{1}{2} \tilde{\mathbf{x}}_N^T \mathbf{Q}_N \tilde{\mathbf{x}}_N + \sum_{i=0}^{N-1} \left( \frac{1}{2} \tilde{\mathbf{x}}_i^T \mathbf{Q}_i \tilde{\mathbf{x}}_i + \frac{1}{2} \tilde{\mathbf{u}}_i^T \mathbf{R}_i \tilde{\mathbf{u}}_i + \tilde{\mathbf{u}}_i^T \mathbf{M} \tilde{\mathbf{x}} \right) \\ \text{subject to} & \tilde{\mathbf{x}}_{i+1} = \mathbf{A}_i \tilde{\mathbf{x}}_i + \mathbf{B}_i \tilde{\mathbf{u}}_i \quad , \quad (3.2.16) \\ \text{where} & \mathbf{A}_i = \nabla_{\mathbf{x}_i}^T f_i^k, \mathbf{B}_i = \nabla_{\mathbf{u}_i}^T f_i^k \\ \text{given} & \tilde{\mathbf{x}}_0 = 0 \end{aligned}$$

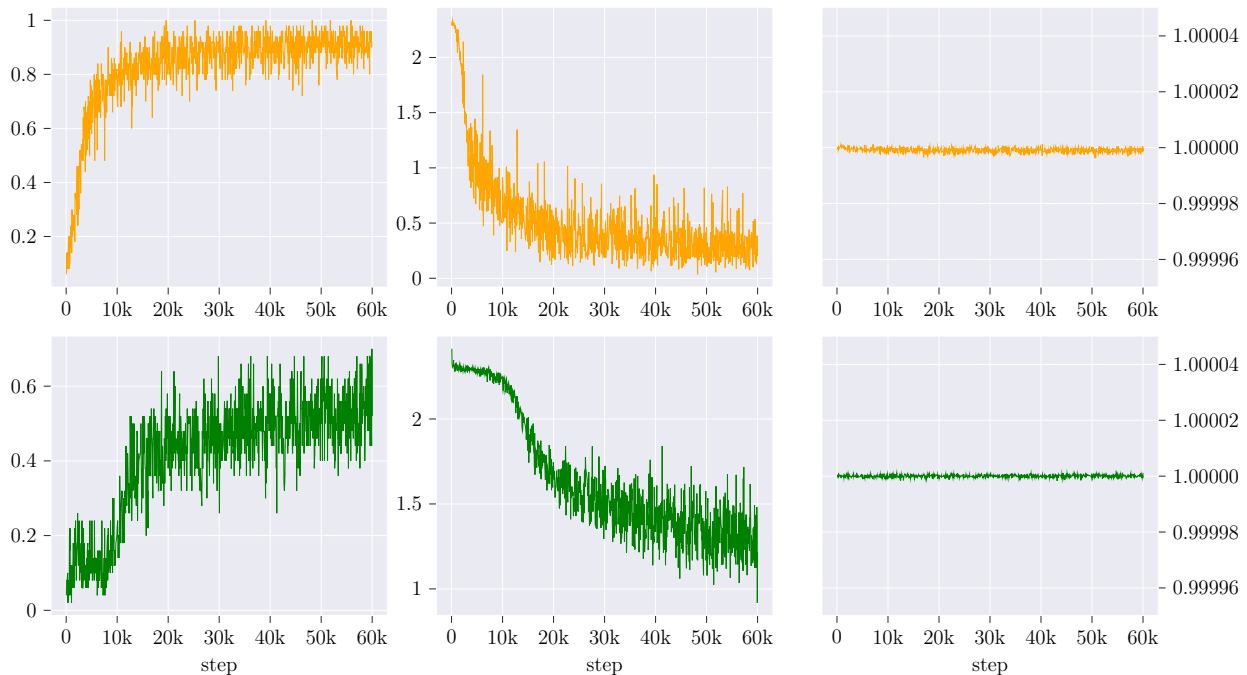
where

$$\begin{aligned}
\mathbf{Q}_N &= \nabla_{\mathbf{x}_N}^2 J(\mathbf{u}^k); \quad \mathbf{a}_N = \nabla_{\mathbf{x}_N} J(\mathbf{u}^k), \\
\mathbf{Q}_i &= \nabla_{\mathbf{x}_i \mathbf{x}_i}^2 H_i; \quad \mathbf{R}_i = \nabla_{\mathbf{u}_i \mathbf{u}_i}^2 H_i; \quad \mathbf{M}_i = \nabla_{\mathbf{u}_i \mathbf{x}_i}^2 H_i \\
H_i &= f_i(\mathbf{x}^k, \mathbf{u}^k)^T \mathbf{p}_i .
\end{aligned}$$

The above is a problem we can solve via the Riccati equation (see section 3.2.2).

### 3.3 Experimental validation

To validate the theory developed in the previous section, we set up a simple experiment: a NN trained with SGD and NGD, via the LQR algorithm, for one epoch. At every time step, before updating, we measured the cosine similarity between the update proposed by the LQR algorithm and the one obtained by plain automatic differentiation.



**Figure 3.1. top:** A NN composed of 3 convolution layers followed by 2 dense layers trained on MNIST database with SGD (5 examples per minibatch) but implemented with the LQR approach. **left:** Accuracy over one epoch of training. **center:** Loss over one epoch. **right:** The cosine similarity, measured at every step, between the gradients calculated via the LQR approach and by traditional automatic differentiation. **bottom:** Same as above, but for NGD.

As we can see in figure 3.1, the cosine similarity between the two schemes always remains close to 1, up to a factor of  $10^{-5}$ , meaning they are pointing in the same direction. The very

small difference measured between the updates can thus safely be attributed to numerical precision.

Yet, the results presented here are mostly theoretical by now, as we are still lacking an efficient implementation that could justify a wide deployment of this approach. Those limitations will be discussed further in the conclusion, and we will touch upon possible improvement directions in chapter 5. But first, we will turn our attention toward a question that arose quickly in the development of our results: could the equivalence be extended to infinite horizon problems?



# Chapter 4

---

## Extension to the Infinite Horizon Case

The previous developments were made with the discrete horizon setting in mind. Yet, Linear Quadratic Regulation problems are also well-defined and understood in the infinite horizon setting, and we would like to extend the equivalence developed in the finite setting to this problem family as well.

In the deep learning setting, we can think of the DEQ (2.4.1) as an example of an infinite horizon discrete problem. The implicit block (usually made of multiple individual layers) can be thought of as the sole layer repeated indefinitely until a fixed point is reached. The loss is then usually taken at this fixed point.

Beginning with the equivalent LQR we have for steepest descent in finite horizon setting, we will try to manipulate the obtained solution to make it compatible with the LQR formulation for infinite horizon problems.

One benefit that arises from such a correspondence is that the optimal solution of a discrete LQR problem with infinite horizon is obtained by solving the associated DARE problem. Fortunately for us, the DARE is a well-studied problem for which a variety of powerful solvers are available.

### 4.1 Problem Statement

The extension of the equivalence developed in (3.2.2) to the infinite-horizon setting leads to an OCP in which the fixed point  $\mathbf{x}_j$  is reached after  $j$  iterations of the dynamic system.

$$\begin{aligned}
& \min_{\tilde{\mathbf{u}}} \nabla_{\mathbf{x}_N}^T J(\mathbf{u}^k) \times \tilde{\mathbf{x}}_N + g(\tilde{\mathbf{x}}^2, \tilde{\mathbf{u}}^2) \\
& \text{where } N \rightarrow \infty \\
& \text{subject to } \tilde{\mathbf{x}}_{i+1} = \mathbf{A}\tilde{\mathbf{x}}_i + \mathbf{B}\tilde{\mathbf{u}}_i \text{ if } i \geq j \\
& \text{subject to } \tilde{\mathbf{x}}_{i+1} = \mathbf{A}_i\tilde{\mathbf{x}}_i + \mathbf{B}_i\tilde{\mathbf{u}}_i \text{ if } i < j \quad , \quad (4.1.1) \\
& \text{where } \mathbf{A} = \nabla_{\mathbf{x}}^T f(\mathbf{x}_j, \mathbf{u}^k), \mathbf{B} = \nabla_{\mathbf{u}}^T f(\mathbf{x}_j, \mathbf{u}^k) \\
& \text{where } \mathbf{A}_i = \nabla_{\mathbf{x}}^T f(\mathbf{x}_i, \mathbf{u}^k), \mathbf{B}_i = \nabla_{\mathbf{u}}^T f(\mathbf{x}_i, \mathbf{u}^k) \\
& \text{given } \tilde{\mathbf{x}}_0 = 0
\end{aligned}$$

where we used the notation  $g(\tilde{\mathbf{x}}^2, \tilde{\mathbf{u}}^2)$  to denote the quadratic term in the LQR objective. Note that the above system is equivalent to a linear OCP with transition layers having *frozen weights*  $\mathbf{u}$  across all layers. Moreover, the linear transition layers become static after reaching the fixed point (i.e.  $\mathbf{A}_i$  and  $\mathbf{B}_i$  stabilize).

Just as it is the case for a NN with frozen weights, the total gradient over  $\mathbf{u}$  is given by the sums of the gradients recovered at each layer. This means that if the optimal solution to the above problem is  $\tilde{\mathbf{u}}_i^*$ , then the update over the frozen weights  $\mathbf{u}^k$  must be:

$$\Delta \mathbf{u}^k = \sum_i^{N \rightarrow \infty} \tilde{\mathbf{u}}_i^* . \quad (4.1.2)$$

We should not expect that  $\tilde{\mathbf{u}}_i \rightarrow 0$  as  $N \rightarrow \infty$ , thus the infinite summation above is obviously divergent.

To overcome that problem, we should not worry about the total gradient value, but instead its *direction*. In that case, the overall direction of the update stabilizes only if the linear system in (4.1.1) obtained after reaching the fixed point also leads to a fixed point. That is, the simplified dynamic system denoted by  $\tilde{\mathbf{x}}_{i+1} = \mathbf{A}\tilde{\mathbf{x}}_i + \mathbf{B}\tilde{\mathbf{u}}_i$  must reach a fixed point, that we can denote by the pair  $(\tilde{\mathbf{x}}_g, \tilde{\mathbf{u}}_g)$ .

We can convince ourselves that this will always be the case by using some notions from dynamical system theory. It is known that to have a fixed point, the dynamic system space must be contractive (i.e. spectral radius smaller than one) in the vicinity of the fixed point (see section 2.2 in [Nayfeh and Balachandran 1995](#)). This means that locally around the fixed point, the linearization of the dynamic system must also be contractive. Thus, since  $\mathbf{A}\tilde{\mathbf{x}}_i + \mathbf{B}\tilde{\mathbf{u}}_i$  is a linearization around a fixed point of our system, it itself must have a fixed

point (see proposition 1, [Gilbert 1992](#)).

Now, if the linearized system has reached its fixed point after  $j'$  iterations, then we can rewrite (4.1.1) equivalently as :

$$\begin{aligned}
& \min_{\tilde{\mathbf{u}}} \nabla_{\mathbf{x}_N}^T J(\mathbf{u}^k) \times \tilde{\mathbf{x}}_N + g(\tilde{\mathbf{x}}^2, \tilde{\mathbf{u}}^2) \\
& \text{where } N \rightarrow \infty \\
& \text{subject to } \tilde{\mathbf{x}}_{i+1} = \tilde{\mathbf{x}}_g = \mathbf{A}\tilde{\mathbf{x}}_g + \mathbf{B}\tilde{\mathbf{u}}_g \text{ if } i \geq g = j + j' \\
& \text{subject to } \tilde{\mathbf{x}}_{i+1} = \mathbf{A}\tilde{\mathbf{x}}_i + \mathbf{B}\tilde{\mathbf{u}}_i \text{ if } j + j' > i \geq j \\
& \text{subject to } \tilde{\mathbf{x}}_{i+1} = \mathbf{A}_i\tilde{\mathbf{x}}_i + \mathbf{B}_i\tilde{\mathbf{u}}_i \text{ if } i < j \\
& \text{where } \mathbf{A} = \nabla_{\mathbf{x}}^T f(\mathbf{x}_j, \mathbf{u}^k), \mathbf{B} = \nabla_{\mathbf{u}}^T f(\mathbf{x}_j, \mathbf{u}^k) \\
& \text{where } \mathbf{A}_i = \nabla_{\mathbf{x}}^T f(\mathbf{x}_i, \mathbf{u}^k), \mathbf{B}_i = \nabla_{\mathbf{u}}^T f(\mathbf{x}_i, \mathbf{u}^k) \\
& \text{given } \tilde{\mathbf{x}}_0 = 0
\end{aligned} \tag{4.1.3}$$

Therefore, it means that as  $N \rightarrow \infty$ , we will have a stabilization of the optimal update  $\tilde{\mathbf{u}}_i^*$  around  $\tilde{\mathbf{u}}_g^*$  and  $\Delta \mathbf{u} \rightarrow N\tilde{\mathbf{u}}_g^*$ . That is, the optimal update goes toward infinity, but it does so by being parallel to  $\tilde{\mathbf{u}}_g^*$ .

This means that to recover the *optimal direction* we only need to solve the system after it has reached its fixed point and the resulting problem is then equivalent to solving (4.1.1):

$$\begin{aligned}
& \min_{\tilde{\mathbf{u}}} \nabla_{\mathbf{x}_N}^T J(\mathbf{u}^k) \times \tilde{\mathbf{x}}_N + g(\tilde{\mathbf{x}}^2, \tilde{\mathbf{u}}^2) \\
& \text{where } N \rightarrow \infty \\
& \text{subject to } \tilde{\mathbf{x}}_{i+1} = \mathbf{A}\tilde{\mathbf{x}}_i + \mathbf{B}\tilde{\mathbf{u}}_i \\
& \text{where } \mathbf{A} = \nabla_{\mathbf{x}}^T f(\mathbf{x}_j, \mathbf{u}^k), \mathbf{B} = \nabla_{\mathbf{u}}^T f(\mathbf{x}_j, \mathbf{u}^k) \\
& \text{given } \tilde{\mathbf{x}}_0 = 0
\end{aligned} \tag{4.1.4}$$

Keeping the initialization at  $\tilde{\mathbf{x}}_0 = 0$  makes no difference here since we know that the linear system given by  $\tilde{\mathbf{x}}_{i+1} = \mathbf{A}\tilde{\mathbf{x}}_i + \mathbf{B}\tilde{\mathbf{u}}_i$  is contractive and will eventually reach the same fixed point,  $\tilde{\mathbf{x}}_g$  as (4.1.3).

The above OCP is almost a proper LQR in the infinite horizon setting, except that the terms  $\nabla_{\mathbf{x}_N}^T J(\mathbf{u}^k) \times \tilde{\mathbf{x}}_N$  makes no sense in such a setting. But since  $\tilde{\mathbf{x}}_i$  eventually reaches a

fixed point, we can replace it by  $\sum_{i=0}^{N \rightarrow \infty} \nabla_{\mathbf{x}_j}^T J(\mathbf{u}^k) \times \tilde{\mathbf{x}}_i$ . Indeed, as  $N \rightarrow \infty$ , this sum will tend to  $N \left( \nabla_{\mathbf{x}_j}^T J(\mathbf{u}^k) \times \tilde{\mathbf{x}}_g \right)$  and is thus a good proxy to the minimization of  $\nabla_{\mathbf{x}_N}^T J(\mathbf{u}^k) \times \tilde{\mathbf{x}}_N$ .

By construction, we can then state the following equivalence between steepest descent and LQR problems in a discrete infinite horizon setting:

**Theorem 4.1.1.** *The steepest descent update over the parameters  $\mathbf{u}^k$  for a given example (input)  $\mathbf{x}^k$  in a discrete infinite horizon OCP (for example, a DEQ) is equivalent to solving the following LQR problem:*

$$\begin{aligned} \min_{\tilde{\mathbf{u}}} \quad & \sum_{i=0}^{\infty} \left[ \nabla_{\mathbf{x}_j}^T J(\mathbf{u}^k) \times \tilde{\mathbf{x}}_i + g(\tilde{\mathbf{x}}_i^2, \tilde{\mathbf{u}}_i^2) \right] \\ \text{subject to} \quad & \tilde{\mathbf{x}}_{i+1} = \mathbf{A}\tilde{\mathbf{x}}_i + \mathbf{B}\tilde{\mathbf{u}}_i \quad , \\ \text{where} \quad & \mathbf{A} = \nabla_{\mathbf{x}}^T f(\mathbf{x}_j, \mathbf{u}), \mathbf{B} = \nabla_{\mathbf{u}}^T f(\mathbf{x}_j, \mathbf{u}) \\ \text{given} \quad & \tilde{\mathbf{x}}_0 = 0 \end{aligned} \quad (4.1.5)$$

where  $\mathbf{x}_j$  is the fixed point reached by the underlying dynamical system and  $g(\tilde{\mathbf{x}}_i^2, \tilde{\mathbf{u}}_i^2)$  denotes the quadratic cost associated with the chosen divergence function used to measure the minimal steepest descent step.

### 4.1.1 Riccati's Solution in Infinite Horizon Setting

Just as for their finite counterpart, an exact solution exists for the problem of the form (4.1.5) and it can be obtained by solving the discrete algebraic Riccati equation (DARE). For a problem like (4.1.5), we need to solve the following Riccati system:

$$\mathbf{K} = \mathbf{A}^T \mathbf{K} \mathbf{A} + \mathbf{Q} - (\mathbf{M} + \mathbf{B}^T \mathbf{K} \mathbf{A})^T (\mathbf{R} + \mathbf{B}^T \mathbf{K} \mathbf{B})^{-1} (\mathbf{M} + \mathbf{B}^T \mathbf{K} \mathbf{A}) \quad (4.1.6)$$

$$\boldsymbol{\lambda} = \mathbf{a} + \mathbf{A}^T \boldsymbol{\lambda} - (\mathbf{M} + \mathbf{B}^T \mathbf{K} \mathbf{A})^T (\mathbf{R} + \mathbf{B}^T \mathbf{K} \mathbf{B})^{-1} (\mathbf{B}^T \boldsymbol{\lambda}) \quad , \quad (4.1.7)$$

with

$$\mathbf{a} = \nabla_{\mathbf{x}_j} J(\mathbf{u}^k) \quad (4.1.8)$$

$$\mathbf{Q} = \nabla_{\mathbf{x}\mathbf{x}}^2 H ; \mathbf{R} = \nabla_{\mathbf{u}\mathbf{u}}^2 H ; \mathbf{M} = \nabla_{\mathbf{u}\mathbf{x}}^2 H \quad (4.1.9)$$

$$H = D(\mathbf{x}_j) + f(\mathbf{x}_j, \mathbf{u}^k)^T \mathbf{p} \quad , \quad (4.1.10)$$

where  $D(\cdot)$  is the divergence measure we are using to perform the steepest descent. The adjoint vector is recovered by solving the linear system:

$$\begin{aligned}
\mathbf{p} &= \nabla_{\mathbf{x}} f(\mathbf{x}_j, \mathbf{u}^k) \cdot \mathbf{p} + \nabla_{\mathbf{x}} D(\mathbf{x}_j) \\
\implies \mathbf{p} &= (\mathbf{I} - \mathbf{A}^T)^{-1} \nabla_{\mathbf{x}} D(\mathbf{x}_j)
\end{aligned} \tag{4.1.11}$$

The optimal update is then obtained by solving the following dynamic linear system:

$$(\tilde{\mathbf{x}}^*, \tilde{\mathbf{u}}^*) = \left( \mathbf{A}\tilde{\mathbf{x}}^* + \mathbf{B}\tilde{\mathbf{u}}^*, -(\mathbf{R} + \mathbf{B}^T \mathbf{K} \mathbf{B})^{-1} ((\mathbf{M} + \mathbf{B}^T \mathbf{K} \mathbf{A})\tilde{\mathbf{x}}^* + \mathbf{B}^T \boldsymbol{\lambda}) \right) \tag{4.1.12}$$

## 4.2 Equivalence with Implicit Differentiation

To our surprise, it appears that the LQR approach in infinite horizon leads to the exact same algorithm as the one obtained when calculating the gradient via the Implicit Function Theorem (see section 2.4). We summarize this result in the following theorem:

**Theorem 4.2.1.** *Taking a gradient step in a NN with infinite depth (such as a DEQ) via the Implicit Function Theorem is equivalent to finding the solution to a discrete-time LQR with an infinite horizon.*

**Proof.** We know from Theorem 4.1.1 that to retrieve the GD update for a NN with infinite depth that reaches a fixed point, we need to solve the following LQR problem:

$$\begin{aligned}
&\min_{\tilde{\mathbf{u}}} \sum_{i=0}^{\infty} \left[ \nabla_{\mathbf{x}_j}^T J(\mathbf{u}^k) \times \tilde{\mathbf{x}}_i + \tilde{\mathbf{u}}_i^T \tilde{\mathbf{u}}_i \right] \\
&\text{subject to } \tilde{\mathbf{x}}_{i+1} = \mathbf{A}\tilde{\mathbf{x}}_i + \mathbf{B}\tilde{\mathbf{u}}_i \\
&\text{where } \mathbf{A} = \nabla_{\mathbf{x}}^T f(\mathbf{x}_j, \mathbf{u}), \mathbf{B} = \nabla_{\mathbf{u}}^T f(\mathbf{x}_j, \mathbf{u}) \\
&\text{given } \tilde{\mathbf{x}}_0 = 0
\end{aligned} \tag{4.2.1}$$

Leveraging Riccati's equation (4.1.6), we need to solve a matrix  $\mathbf{K}$  that satisfies:

$$\mathbf{K} = \mathbf{A}^T \mathbf{K} \mathbf{A} - (\mathbf{B}^T \mathbf{K} \mathbf{A})^T (\mathbf{I} + \mathbf{B}^T \mathbf{K} \mathbf{B})^{-1} (\mathbf{B}^T \mathbf{K} \mathbf{A})$$

to find the optimal update. But, the above problem is in fact a trivial case, for which the obvious solution is  $\mathbf{K} = \mathbf{0}$ . The optimal update is therefore given by (4.1.12):

$$\tilde{\mathbf{u}}^* = -\mathbf{B}^T \boldsymbol{\lambda}$$

With  $\boldsymbol{\lambda}$  obtained by solving the linear dynamic system:

$$\boldsymbol{\lambda} = \mathbf{a} + \mathbf{A}^T \boldsymbol{\lambda} = \nabla_{\mathbf{x}_j} J(\mathbf{u}^k) + \nabla_{\mathbf{x}} f(\mathbf{x}_j, \mathbf{u}^k) \boldsymbol{\lambda}$$

Or equivalently:

$$\begin{aligned} \mathbf{0} &= (\nabla_{\mathbf{x}} f(\mathbf{x}_j, \mathbf{u}^k) - \mathbf{I}) \boldsymbol{\lambda} + \nabla_{\mathbf{x}_j} J(\mathbf{u}^k) \\ \implies \boldsymbol{\lambda} &= -(\nabla_{\mathbf{x}} f(\mathbf{x}_j, \mathbf{u}^k) - \mathbf{I})^{-1} \nabla_{\mathbf{x}_j} J(\mathbf{u}^k) \end{aligned}$$

And thus, the parameters update can be rewritten as:

$$\tilde{\mathbf{u}}^* = \nabla_{\mathbf{u}} f(\mathbf{x}_j, \mathbf{u}) (\nabla_{\mathbf{x}} f(\mathbf{x}_j, \mathbf{u}^k) - \mathbf{I})^{-1} \nabla_{\mathbf{x}_j} J(\mathbf{u}^k)$$

Remembering that the matrix of gradients is the transpose of the Jacobian, we can recognize that the parameter update recovered by solving the LQR is the same one as when we retrieve the gradient via the Implicit Function Theorem (2.4.7)  $\square$

Although this result does not improve or lead to a new approach to recovering the gradient, it provides a novel view on the training process behind the optimization of NN with infinite depth.

### 4.3 Cyclic Solution

So far, we have treated the implicit block of layers involved in the transformation  $f(\mathbf{x}, \mathbf{u})$  as a whole. However, in practice, this transformation is usually made up of individual layers. If the implicit block is composed of  $N$  layers, we can rewrite  $f(\mathbf{x}, \mathbf{u})$  as:

$$f(\mathbf{x}, \mathbf{u}) = f_N(f_{N-1}(\dots(f_2(f_1(\mathbf{x}, \mathbf{u}))))))$$

Taking into account this decomposition, we can recast the LQR problem (4.1.5) as:

$$\begin{aligned} \min_{\tilde{\mathbf{u}}} \quad & \sum_{i=0}^{\infty} \sum_{n=1}^N \left[ \nabla_{\mathbf{x}_j}^T J(\mathbf{u}^k) \times \tilde{\mathbf{x}}_n^i + g((\tilde{\mathbf{x}}_n^i)^2, (\tilde{\mathbf{u}}_n^i)^2) \right] \\ \text{subject to} \quad & \tilde{\mathbf{x}}_{n+1}^i = \mathbf{A}_n \tilde{\mathbf{x}}_n^i + \mathbf{B}_n \tilde{\mathbf{u}}_n^i \\ \text{with} \quad & \tilde{\mathbf{x}}_{N+1}^i = \tilde{\mathbf{x}}_0^{i+1} \\ \text{where} \quad & \mathbf{A}_n = \nabla_{\mathbf{x}_n}^T f_n(\mathbf{x}_n^j, \mathbf{u}_n^k), \mathbf{B}_n = \nabla_{\mathbf{u}_n}^T f_n(\mathbf{x}_n^j, \mathbf{u}_n^k) \\ \text{given} \quad & \tilde{\mathbf{x}}_0^0 = 0 \end{aligned} \tag{4.3.1}$$

where  $\mathbf{x}_n^j$  denotes the fixed point inside the layer  $f_n$ . Note that we moved the iteration number  $i$  to a superscript to make room for the needed additional index,  $n$ , denoting the layer. Overall, (4.3.1) denotes the exact same system as the one in (4.1.5); we simply exposed here the underlying structure of the transformation  $f(\mathbf{x}, \mathbf{u})$  that is repeated infinitely.

Nevertheless, a problem with a form such as (4.3.1) is known as a *cyclic LQR* and the solution for such problem can be obtained by solving the cyclic equations system over  $\mathbf{K}_i$ :

$$\begin{aligned}
\mathbf{K}_{N-1} &= \mathbf{A}_N^T \mathbf{K}_N \mathbf{A}_N + \mathbf{Q}_N - (\mathbf{M}_N + \mathbf{B}_N^T \mathbf{K}_N \mathbf{A}_N)^T (\mathbf{R}_N + \mathbf{B}_N^T \mathbf{K}_N \mathbf{B}_N)^{-1} (\mathbf{M}_N + \mathbf{B}_N^T \mathbf{K}_N \mathbf{A}_N) \\
&\vdots \\
\mathbf{K}_1 &= \mathbf{A}_2^T \mathbf{K}_2 \mathbf{A}_2 + \mathbf{Q}_2 - (\mathbf{M}_2 + \mathbf{B}_2^T \mathbf{K}_2 \mathbf{A}_2)^T (\mathbf{R}_2 + \mathbf{B}_2^T \mathbf{K}_2 \mathbf{B}_2)^{-1} (\mathbf{M}_2 + \mathbf{B}_2^T \mathbf{K}_2 \mathbf{A}_2) \\
\mathbf{K}_N &= \mathbf{A}_1^T \mathbf{K}_1 \mathbf{A}_1 + \mathbf{Q}_1 - (\mathbf{M}_1 + \mathbf{B}_1^T \mathbf{K}_1 \mathbf{A}_1)^T (\mathbf{R}_1 + \mathbf{B}_1^T \mathbf{K}_1 \mathbf{B}_1)^{-1} (\mathbf{M}_1 + \mathbf{B}_1^T \mathbf{K}_1 \mathbf{A}_1)
\end{aligned} \tag{4.3.2}$$

The same cyclic structure arises as well for the  $\lambda_i$ :

$$\begin{aligned}
\lambda_{N-1} &= \mathbf{a} + \mathbf{A}_N^T \lambda_N - (\mathbf{M}_N + \mathbf{B}_N^T \mathbf{K}_N \mathbf{A}_N)^T (\mathbf{R}_N + \mathbf{B}_N^T \mathbf{K}_N \mathbf{B}_N)^{-1} (\mathbf{B}_N^T \lambda_N) \\
&\vdots \\
\lambda_1 &= \mathbf{a} + \mathbf{A}_2^T \lambda_2 - (\mathbf{M}_2 + \mathbf{B}_2^T \mathbf{K}_2 \mathbf{A}_2)^T (\mathbf{R}_2 + \mathbf{B}_2^T \mathbf{K}_2 \mathbf{B}_2)^{-1} (\mathbf{B}_2^T \lambda_2) \\
\lambda_N &= \mathbf{a} + \mathbf{A}_1^T \lambda_1 - (\mathbf{M}_1 + \mathbf{B}_1^T \mathbf{K}_1 \mathbf{A}_1)^T (\mathbf{R}_1 + \mathbf{B}_1^T \mathbf{K}_1 \mathbf{B}_1)^{-1} (\mathbf{B}_1^T \lambda_1)
\end{aligned} \tag{4.3.3}$$

and the optimal solution is obtained by bringing solving the following system:

$$(\tilde{\mathbf{x}}_{n+1}^*, \tilde{\mathbf{u}}_{n+1}^*) = \left( \mathbf{A}_n \tilde{\mathbf{x}}_n^* + \mathbf{B}_n \tilde{\mathbf{u}}_n^*, -(\mathbf{R}_n + \mathbf{B}_n^T \mathbf{K}_n \mathbf{B}_n)^{-1} ((\mathbf{M}_n + \mathbf{B}_n^T \mathbf{K}_n \mathbf{A}_n) \tilde{\mathbf{x}}_n^* + \mathbf{B}_n^T \lambda_n) \right) \tag{4.3.4}$$

Although it may seem that this formulation only complicates things, we will see later that it allows us to leverage a very useful approximation in the infinite setting. Moreover, we again benefit from the maturity of the LQR field since cyclic counterparts already exist for most of the DARE solvers.

## 4.4 LQR-based Natural Gradient Method for DEQs

To the best of our knowledge, second-order methods have not been deployed for the optimization of infinite depth NN such as DEQ. A possible explanation is the computational and engineering challenge necessary to recover the second derivative via implicit differentiation.

What is interesting about the LQR equivalence is that it allows us to follow a simple recipe to compute the update for second-order methods optimization. As an example, we detail below the algorithm that arises from the LQR formulation when we want to find the NGD update for a system such as DEQ, given the implicit block  $f$  and the final transformation  $h$  (usually a dense layer followed by a softmax). We assume we have access to a fixed-point solver, that we denote `fx_point_solve`, taking as input an initial value and a mapping. The solver can simply be the repeated application of the implicit block until convergence or leverage an acceleration scheme such as Anderson iteration (Anderson 1965). We also assume we have access to a DARE solver; for a variety of such see Bini *et al.* 2011.

---

**Algorithm 1** One epoch of training with natural gradient descent over an implicit model

---

```

procedure IMPLICIT-NGD( $f, \mathbf{x}_0$ )
  for  $\mathbf{x}_0^k \in \mathcal{D}$  do
     $\mathbf{x}_j \leftarrow \text{fx\_point\_solve}(f, \mathbf{x}_0^k)$ 
     $\mathbf{A} \leftarrow \nabla_{\mathbf{x}} f(\mathbf{x}_j)$ 
     $\mathbf{B} \leftarrow \nabla_{\mathbf{u}} f(\mathbf{u}^k)$ 
     $\mathbf{p} \leftarrow \text{fx\_point\_solve}(\mathbf{A}^T \mathbf{p} + \nabla_{\mathbf{x}} D_{\text{KL}}(h(\mathbf{x}_j) || h(\mathbf{x}')) \Big|_{\mathbf{x}'=\mathbf{x}_j}, \mathbf{p} = 0)$ 
    ▷  $\mathbf{p}$  is the adjoint

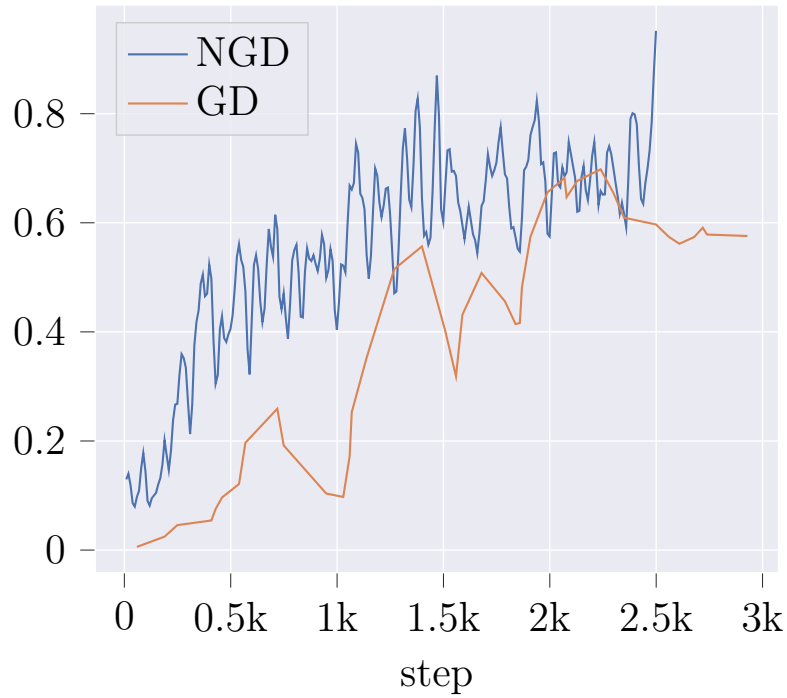
     $H(\mathbf{x}', \mathbf{u}') \leftarrow D_{\text{KL}}(h(\mathbf{x}_j) || h(\mathbf{x}')) + f(\mathbf{x}', \mathbf{u}')^T \mathbf{p}$ 
     $\mathbf{Q} \leftarrow \nabla_{\mathbf{x}\mathbf{x}} H(\mathbf{x}_j, \mathbf{u}^k)$ 
     $\mathbf{R} \leftarrow \nabla_{\mathbf{u}\mathbf{u}} H(\mathbf{x}_j, \mathbf{u}^k)$ 
     $\mathbf{Q} \leftarrow \nabla_{\mathbf{u}\mathbf{x}} H(\mathbf{x}_j, \mathbf{u}^k)$ 
     $\mathbf{a} \leftarrow \nabla_{\mathbf{x}_j} J(\mathbf{u}^k)$ 
     $\mathbf{K} \leftarrow \text{DARE}(\mathbf{a}, \mathbf{A}, \mathbf{B}, \mathbf{Q}, \mathbf{R}, \mathbf{M})$ 
     $\boldsymbol{\lambda} \leftarrow \text{fx\_point\_solve}(\mathbf{a} + \mathbf{A}^T \boldsymbol{\lambda} - (\mathbf{M} + \mathbf{B}^T \mathbf{K} \mathbf{A})^T (\mathbf{R} + \mathbf{B}^T \mathbf{K} \mathbf{B})^{-1} \mathbf{B}^T \boldsymbol{\lambda}, \boldsymbol{\lambda} = \mathbf{a})$ 
     $\mathbf{u}^{k+1} \leftarrow \text{fx\_point\_solve}(-(\mathbf{R} + \mathbf{B}^T \mathbf{K} \mathbf{B})^{-1} ((\mathbf{M} + \mathbf{B}^T \mathbf{K} \mathbf{A}) \tilde{\mathbf{x}} + \mathbf{B}^T \boldsymbol{\lambda}), \mathbf{u} = \mathbf{u}^k)$ 
    ▷  $\tilde{\mathbf{x}}$  updated simultaneously to  $\tilde{\mathbf{x}} \leftarrow \mathbf{A} \tilde{\mathbf{x}} + \mathbf{B} \mathbf{u}$ 
  end for
  return  $\mathbf{u}^N$ 
end procedure

```

---

We tried this algorithm ourselves and noticed that the number of iteration needed to converge was smaller on our toy network trained over MNIST, as showcased in figure 4.1. Sadly, we encountered stability issues and were not able to scale this algorithm to larger DEQs.





**Figure 4.1.** Zoom on the accuracy achieved (first 2500 iterations) during the training of a toy DEQ (a Resnet block with 12 inner channels) over MNIST with the NGD algorithm arising from the LQR approach. NGD starts to converge in less iteration, but it is 10x slower than GD.

# Chapter 5

---

## Computational Perspective

Computationally speaking, when thinking about second-order method optimizers, the approach described above already has a slight advantage over the direct implementation. Indeed, the difficulties with those methods come from the fact that we need to invert a matrix (Hessian for Newton's method, the Fisher for NGD); a procedure that typically requires  $\mathcal{O}(m^{2.3})$  iterations for a matrix of size  $m \times m$  when solved with an algorithm such as conjugate gradients (Hestenes and Stiefel 1952). This cost often proves prohibitive with typical architectures in DL, the dimension of the parameters being simply too big.

Fortunately for us, a benefit of building the equivalent LQR is to reduce (slightly) the computation complexity. To see it, let us consider that we want to perform NGD to optimize a NN. Then, as seen with the system of equations shown in (3.2.6) and (4.3.2), the inversion of the Fisher matrix (of size  $m \times m$ ,  $m$  being the total number of parameters) is replaced by the inversion of  $N$  matrices scaling with the number of parameters per layer ( $N$  being the number of layers). Thus, we instead need to invert the  $N$  matrices  $(\mathbf{R} + \mathbf{B}^T \mathbf{K} \mathbf{B})^{-1}$  of size  $m_1, m_2, \dots, m_N$ ; an operation that has a computation complexity of the form  $\mathcal{O}(m_1^{2.3} + m_2^{2.3} + \dots + m_N^{2.3})$ . Knowing that  $m = m_1 + m_2 + \dots + m_N$ , the second procedure has thus an advantage.

Yet, this gain is not enough to justify the use of the LQR method since the number of parameters inside an individual layer can still get pretty big with modern NN architectures. To alleviate, this problem we now develop a useful approach when facing a NN composed of layers with activation functions that have *at most* a linear dependency over their respective layer parameters.

## 5.1 Leveraging Neumann Series

Let consider the term we need to invert when using the LQR approach to make a steepest descent step:

$$(\bar{\mathbf{R}}_i + \mathbf{B}_i^T \mathbf{K}_i \mathbf{B}_i)^{-1} , \quad (5.1.1)$$

where we have  $\bar{\mathbf{R}}_i = \mathbf{R}_i + \delta_i \mathbf{I}$ ; that is a penalty term  $\delta_i$  is added over the diagonal of the matrix  $\mathbf{R}_i$ . This decision corresponds to the addition of a damping term and can be justified as a trust-region method. We will discuss this term in more detail later, so let us just accept it as is for now.

Now, let us turn our attention to  $\mathbf{R}_i$ . When the divergence measure is only taken over the final state, this quantity is the Hessian of the transition function (a layer) with respect to its parameters. But, typical layers inside a NN are composed of a linear function followed by an activation ( $\sigma$ ), that we can usually write as:

$$\mathbf{R}_i = \nabla \mathbf{u}_i \mathbf{u}_i f_i = \nabla \mathbf{u}_i \mathbf{u}_i \sigma(\mathbf{W} \mathbf{x}_i + \mathbf{b}_i) \quad (5.1.2)$$

If we were not to use any activation function,  $\mathbf{R}_i$  would be 0. What is more interesting though is that the most popular activation function, the Rectified Linear Unit (ReLU) ([Jarrett \*et al.\* 2009](#); [Glorot \*et al.\* 2011](#)), does not introduce a co-dependency in the activation between the parameters of the current layer. Indeed, the ReLU pushes the negative terms to 0 and leaves unmodified the others. For this particular case, we have:

$$\mathbf{R}_i = \nabla \mathbf{u}_i \mathbf{u}_i \text{ReLu}(\mathbf{W} \mathbf{x}_i + \mathbf{b}_i) = 0 . \quad (5.1.3)$$

This means that for such networks, we are seeking to invert matrices of the form:

$$\begin{aligned} (\bar{\mathbf{R}}_i + \mathbf{B}_i^T \mathbf{K}_i \mathbf{B}_i)^{-1} &= (\mathbf{I}_{\delta_i} + \bar{\mathbf{K}}_i)^{-1} \\ &= (\delta_i (\mathbf{I} + \frac{1}{\delta_i} \bar{\mathbf{K}}_i))^{-1} , \\ &= \frac{1}{\delta_i} + (\mathbf{I} + \frac{1}{\delta_i} \bar{\mathbf{K}}_i)^{-1} \end{aligned} \quad (5.1.4)$$

where we abbreviated  $\mathbf{B}_i^T \mathbf{K}_i \mathbf{B}_i$  by  $\bar{\mathbf{K}}_i$ . If the magnitude of the biggest eigenvalue of  $\frac{1}{\delta_i} \bar{\mathbf{K}}_i$  is smaller than one, we can use Neumann's series to approximate this inverse, knowing that it would be convergent (e.g., see [Neumann 1877](#); [Suzuki 1976](#)):

$$\begin{aligned} (\mathbf{I} + \frac{1}{\delta_i} \bar{\mathbf{K}}_i)^{-1} &= \sum_{j=0}^{\infty} (-\frac{1}{\delta_i} \bar{\mathbf{K}}_i)^j \\ &\approx \mathbf{I} - \frac{1}{\delta_i} \bar{\mathbf{K}}_i + \frac{1}{\delta_i^2} \bar{\mathbf{K}}_i^2 - \frac{1}{\delta_i^3} \bar{\mathbf{K}}_i^3 + \dots \end{aligned} \quad (5.1.5)$$

The truncated series is often a pretty good approximation, as long as the spectral radius of  $\frac{1}{\delta_i} \bar{\mathbf{K}}$  remains smaller than one. What we now propose is to use the penalty  $\delta_i$  as an adaptive layer parameter that will ensure that it is the case. This idea is in line with the use of a damping factor to rescale a linear system before solving it with Chebyshev iteration, a special case of modified Richardson iteration (e.g. see [Golub and Varga 2007](#))

We will show that this is a valid approach, but first, we need to show that the matrices  $\mathbf{K}_i$  (and by extension  $\bar{\mathbf{K}}_i$ ) are always symmetric:

**Theorem 5.1.1.** *The gain matrices  $\mathbf{K}_i$  (or  $\mathbf{K}$  in the infinite horizon setting) given by (3.2.6) (respectively (4.3.2)) are symmetric matrices, and therefore the matrices  $\mathbf{B}_i^T \mathbf{K}_i \mathbf{B}_i$  (respectively  $\mathbf{B}^T \mathbf{K} \mathbf{B}$ ) are also symmetric.*

**Proof.** The proof follows from the fact that the matrices  $\mathbf{K}_i$  are a composition (via addition and subtraction) of symmetric matrices.

By Schwartz's Theorem (see [Rudin 1976](#)) we know that the partial second-order derivatives matrices  $\mathbf{Q}_N, \mathbf{Q}_i, \mathbf{M}_i$  and  $\mathbf{R}_i$  are symmetric matrices.

Moreover, if  $\mathbf{K}_{i+1}$  is symmetric, then  $\mathbf{A}_i^T \mathbf{K}_{i+1} \mathbf{A}_i$  is also symmetric because:

$$(\mathbf{A}_i^T \mathbf{K}_{i+1} \mathbf{A}_i)^T = \mathbf{A}_i^T \mathbf{K}_{i+1}^T \mathbf{A}_i = \mathbf{A}_i^T \mathbf{K}_{i+1} \mathbf{A}_i$$

The same argument holds for  $\mathbf{B}_i^T \mathbf{K}_{i+1} \mathbf{B}_i$ , and for  $(\mathbf{M}_i + \mathbf{B}_i^T \mathbf{K}_{i+1} \mathbf{A}_i)^T (\dots) (\mathbf{M}_i + \mathbf{B}_i^T \mathbf{K}_{i+1} \mathbf{A}_i)$ . Finally, remembering that the inverse of a symmetric matrix is also symmetric, we have that  $\mathbf{K}_i$ ,

$$\mathbf{K}_i = \mathbf{A}_i^T \mathbf{K}_{i+1} \mathbf{A}_i + \mathbf{Q}_i - (\mathbf{A}_i^T \mathbf{K}_{i+1} \mathbf{B}_i + \mathbf{M}_i^T) (\mathbf{R}_i + \mathbf{B}_i^T \mathbf{K}_{i+1} \mathbf{B}_i)^{-1} (\mathbf{M}_i + \mathbf{B}_i^T \mathbf{K}_{i+1} \mathbf{A}_i)$$

is symmetric if  $\mathbf{K}_{i+1}$  is also symmetric. This will always be the case as the terminal gain matrix,  $\mathbf{K}_N = \mathbf{Q}_N$  is symmetric.  $\square$

We can now use the fact that  $\bar{\mathbf{K}}_i$  is symmetric to show that we can indeed use the penalty term  $\delta_i$  to ensure that the new matrix  $\frac{1}{\delta_i}\bar{\mathbf{K}}_i$  will have a spectral radius smaller than one. Indeed, the symmetry of  $\bar{\mathbf{K}}_i$  implies that its eigenvalues are real. It also means that its singular values, obtained via Singular Values Decomposition (SVD), are the absolute values of the eigenvalues (see corollary C.5.2 in [Lebanon 2012](#)). That is,

$$\bar{\mathbf{K}}_i = \bar{\mathbf{U}}_i \bar{\mathbf{\Sigma}}_i \bar{\mathbf{V}}_i^{-1} \quad , \quad (5.1.6)$$

where  $\bar{\mathbf{\Sigma}}_i$  is a diagonal matrix for which the entries are the magnitude of the eigenvalues of  $\bar{\mathbf{K}}_i$ . Therefore, if we denote the largest absolute singular value of  $\bar{\mathbf{K}}_i$  by  $\lambda_m$  and take  $\delta_i > |\lambda_m|$ , then we will have that the SVD of  $\frac{1}{\delta_i}\bar{\mathbf{K}}_i$  can be written as:

$$\frac{1}{\delta_i}\bar{\mathbf{K}}_i = \bar{\mathbf{U}}_i \left( \frac{1}{\delta_i} \bar{\mathbf{\Sigma}}_i \right) \bar{\mathbf{V}}_i^{-1} \quad , \quad (5.1.7)$$

which is an expression that we can use to deduce that the magnitudes of the eigenvalues of  $\frac{1}{\delta_i}\bar{\mathbf{K}}_i$  would be smaller than one and therefore the Neumann's series approximation will be convergent.

### 5.1.1 An Algorithm without Matrix Inversion

The proposed approach above means that we now need to compute the spectral radius (equivalent to its maximum absolute eigenvalues) of the matrix  $\bar{\mathbf{K}}_i$ . Although such calculation seems unusual in DL, it is not without precedent. [Yoshida and Miyato 2017](#) for example regularize the weight matrix of the NN they trained by their spectral radius and claim better generalizability from it. Using the same strategy, [Miyato et al. 2018](#) showed that it could be used to stabilize the training of the discriminator of a generative adversarial network. It is worth noting that both works use the power iteration method to estimate the spectral radius of their matrix of interest, a method that is relatively quick computation-wise.

For us, this means that we can leverage a second-order method optimization method to train a NN that has ReLu activation units without explicitly inverting a matrix. Moreover, the trick can be applied in the finite horizon setting and **in the infinite horizon one** (using the cyclic formulation for the latter).

### 5.1.2 Parallel with Damping Parameter

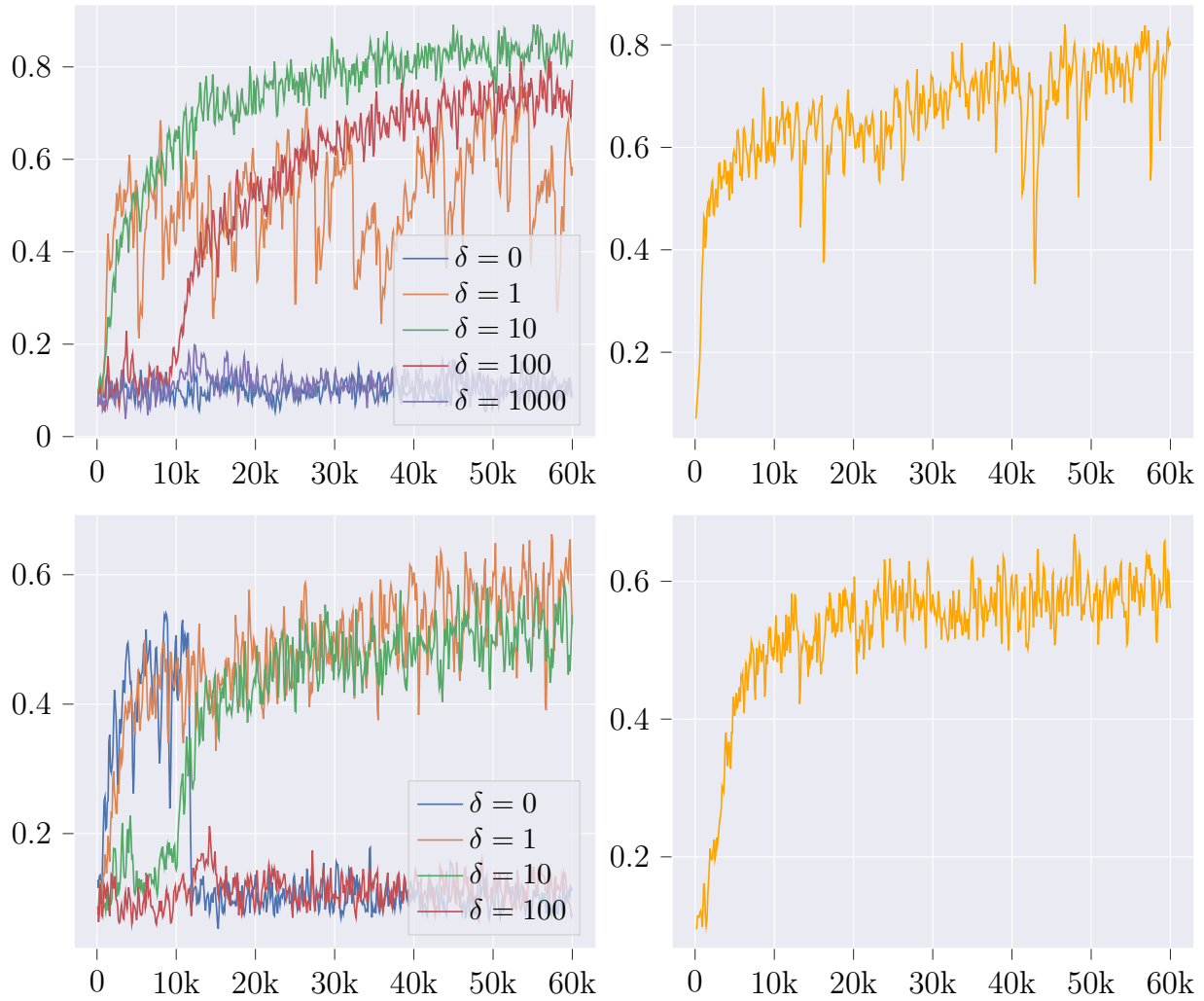
Applying the adaptive penalty  $\delta_i$  may seem arbitrary at first glance, but if we cast it as a damping parameter over the weights update, it suddenly makes much more sense. The use of a damping parameter is advocated as critical by [Martens 2016](#) and [Martens 2020](#) for second-order optimization of NN because the direction that such algorithms propose in parameter space is not aligned with the distribution space. The resulting path will therefore more often than not veers away from the target pointed to by the second-order method.

To alleviate this problem, we need to dampen the size of the updates. One strategy to do so is via trust-region methods, i.e. by penalizing the update size made in some directions. In fact, that is exactly what the penalty term  $\delta_i$  does; its introduction is equivalent to adding the following term to the cost of the LQR (3.2.5):

$$\sum_{i=0}^{N-1} \tilde{\mathbf{u}}_i^T \delta_i \tilde{\mathbf{u}}_i = \sum_{i=0}^{N-1} (\mathbf{u}_i - \mathbf{u}_i^k)^T \delta_i (\mathbf{u}_i - \mathbf{u}_i^k) .$$

Experimentally, we also used a parameter  $\delta$  to stabilize the training for second-order methods implemented as a LQR without Neumann’s series trick, in the same fashion as above. However, the term was kept constant across all layers and fine-tuned manually.

Setting  $\delta_i$  via the spectral radius of  $\bar{\mathbf{K}}_i$  allowed us to dampen adaptively the updates and offer superior performance than without its use. This finding is summarized in figure 5.1. As of now, we do not have a better explanation for this behavior than the intuition discussed above, and it will require further investigation from our side to quantify correctly the impact of  $\delta_i$  on the training.



**Figure 5.1. top:** Accuracy curves of a training experiment, the same Convnet as in figure 3.1 (5 layers deep) trained on MNIST with the LQR approach equivalent to NGD and a learning rate set to  $\alpha = 0.1$ . **left:** In addition to the learning rate, a damping parameter  $\delta$  has to be fine-tuned for the training to be stable without Neumann's series trick. **right:** On the opposite, implementing Neumann's series trick also has the added benefit of adding a layer adaptive damping term that stabilizes the training. **bottom:** Same as above, but for a learning rate set to  $\alpha = 0.01$ .

# Chapter 6

---

## Conclusion

Overall, the most important contribution of this thesis is to provide a template to express steepest descent under various divergence, both in the finite and infinite horizon cases, as a special LQR problem. Equipped with this new tool, a multitude of possible paths can now be explored. Most of our experimentation within our framework has been through second-order methods such as NGD and Newton's, but those only represent specific cases. At the moment, two major challenges are preventing us to apply our method at scale: the lack of full matrix-free implementation and stability issues (for which we may have a partial solution).

**Matrix-Free Implementation:** One of the major advantage of automatic differentiation is the fact that the method is *matrix free* (see [Baydin et al. 2017](#)). Jacobian matrices are never fully evaluated nor stored, and the computation time is much quicker. Unfortunately, the framework we have described above did not completely preserve this property, in both the finite and the infinite settings. This also means that the memory cost incurred by our algorithm is bigger than for standard backpropagation.

For the finite horizon described in chapter 3, nothing prevents us from formulating the algorithm as matrix-free; all Jacobian and Hessian matrices can be derived as matrix-vector product (MVP) functions. The problem arises from the potential combinatorial explosion in the expression of the matrices  $\mathbf{K}_i$  in (3.2.6). Hence, if we were to evaluate  $\mathbf{K}_0 \cdot \mathbf{v}$  for a NN with N layers, we would have to evaluate  $4^{N-1}$  times the product  $\mathbf{K}_N \cdot \mathbf{v}$ , and then to evaluate  $4^{N-2}$  times the product  $\mathbf{K}_{N-1} \cdot \mathbf{v}$ , etc. As such, the total number of MVP evaluations quickly becomes intractable when increasing the number of layers.

To avoid this problem, we currently compute and store the matrices  $\mathbf{K}_i$ . As those matrices are proportional to the size of the representation vector of their respective layer ( $n_i$ ), they



each require  $\mathcal{O}(n_i)$  MVP evaluations and a memory space that scale with  $\mathcal{O}(n_i^2)$ .

It should be noted that we could be strategic when designing the NN architecture that we want to optimize to mitigate this issue. In practice, only storing the  $\mathbf{K}_i$  matrices every 2 or 3 layers would not incur too much slow down. Knowing this, one could include a features bottleneck in their NN design and only store/compute the  $\mathbf{K}_i$  for those layers.

The issue is different in the infinite horizon setting, presented in chapter 4, since typical fast DARE solvers rely on doubling algorithms and therefore perform multiple matrix multiplications (see [Bini et al. 2011](#)). There is thus no gain in keeping everything as an MVP, and for this setting, we currently compute and store all Jacobians and Hessians when performing steepest descent (the matrices  $\mathbf{A}$ ,  $\mathbf{B}$ ,  $\mathbf{Q}$ ,  $\mathbf{R}$ , and  $\mathbf{M}$ ). At least, the implicit block in typical DEQ usually does not consist of many layers, so it is still feasible to do it that way.

**Stability:** The observed stability that occurs during training with the use of the damping parameter  $\delta_i$  introduced in chapter 5 came as a surprise to us; we expected a bit of tuning to be necessary. This penalty term directly reduces the spectral radius of the gain matrices  $\mathbf{K}_i$ , but does it implicitly impact the spectral radius of the weight matrices? We plan to investigate further to see if it is the case as it could allow us to understand the stable behavior through the lens of spectral norm regularization (see [Yoshida and Miyato 2017](#) and [Miyato et al. 2018](#)).

Also, as already mentioned, our formulation allowed us to train DEQ networks under a different divergence measure than the Euclidean norm; a formulation for which popular second-order methods are a specific case. Yet trying it with NGD, we did encounter stability issues during and extensive hyperparameter tuning was required. We hope that bringing the Neumann series trick in the infinite horizon setting will bring the same advantages noticed in the finite setting while allowing faster computation simultaneously.

## 6.1 Future Work

So far, we have limited our work to the field of discrete-time problems but we do believe it is worth exploring if the same connection can be made in the continuous-time setting. As a starting point, we could cast the optimization objective with mirror descent; a technique that was originally motivated in continuous time (see [Nemirovsky and Yudi 1983](#)). The solution of the mirror descent objective in continuous time under various divergence leads to continuous-time descent dynamics to follow to converge to an optimum (e.g. [Krichene et al. 2015](#)).

Moreover, the PMP (Pontriagin *et al.* 1962) was also first developed for control problems in continuous time; the optimality principle still allows us to design an objective that decouples the controls at every time by handing us a function that must be optimal at every time step. That is, the optimal trajectory will consist of the instantaneous optimal control of the Hamiltonian put end-to-end. We thus believe we will be able to cast the mirror descent objective as a continuous LQR (see Miquel 2020 for more details on the PMP and LQR in a continuous setting).

This would imply that a Continuous-Time Algebraic Riccati Equation (CARE) solver could be used to perform a mirror descent step. The discretization of the optimal trajectory obtained would then provide updates for the parameters of a given NN that preserves the convergence guarantee obtained via mirror descent. It could also be used directly to provide updates to a neural ODEs network (Chen *et al.* 2018). Again, since LQR problems have long been studied and since it exists efficient solver for them, we can hope to leverage those tools to obtain performance gain during the training of a neural net.

Another possible line of work stems from the question: does casting the steepest descent step as a LQR gives us a new tool to analyze the behavior of optimization technique? After all, LQR problems come with great convergence guarantees/conditions and analytical solutions. We do believe our framework opens a new door for the analysis of optimization algorithms, especially when a divergence term other than Euclidean distance is brought in to quantify the size of the updates.

Finally, we hope that establishing an equivalence between a steepest descent and LQR will now enable the design of new efficient optimization algorithms while providing new tools to study their theoretical properties.

## References

---

- Almeida, L. B. (1990). *A Learning Rule for Asynchronous Perceptrons with Feedback in a Combinatorial Environment*, page 102–111. IEEE Press.
- Amari, S. (1998). Natural gradient works efficiently in learning. *Neural Computation*, **10**, 251–276.
- Amari, S. and Nagaoka, H. (2007). *Methods of Information Geometry*. American Mathematical Society.
- Anderson, B. D. O. and Moore, J. B. (1990). *Optimal Control: Linear Quadratic Methods*. Prentice-Hall, Inc., USA.
- Anderson, D. G. (1965). Iterative procedures for nonlinear integral equations. *J. ACM*, **12**(4), 547–560.
- Bai, S., Kolter, J. Z., and Koltun, V. (2019). Deep equilibrium models. In *Advances in Neural Information Processing Systems 32: Annual Conference on Neural Information Processing Systems 2019, NeurIPS 2019, December 8-14, 2019, Vancouver, BC, Canada*, pages 688–699.
- Baydin, A. G., Pearlmutter, B. A., Radul, A. A., and Siskind, J. M. (2017). Automatic differentiation in machine learning: a survey. *J. Mach. Learn. Res.*, **18**, 153:1–153:43.
- Beck, A. and Teboulle, M. (2003). Mirror descent and nonlinear projected subgradient methods for convex optimization. *Oper. Res. Lett.*, **31**(3), 167–175.
- Bellman, R. (1952). On the theory of dynamic programming. *Proceedings of the National Academy of Sciences*, **38**(8), 716–719.
- Bertsekas, D. (2016). *Nonlinear Programming*. Athena scientific optimization and computation series. Athena Scientific.
- Bini, D. A., Iannazzo, B., and Meini, B. (2011). *Numerical Solution of Algebraic Riccati Equations*. Society for Industrial and Applied Mathematics.

- Bliss, G. A. (1968). *Lectures on the calculus of variations*. Chicago (Ill.) : University of Chicago press.
- Bottou, L. (1998). Online algorithms and stochastic approximations. In *Online Learning and Neural Networks*. Cambridge University Press, Cambridge, UK. revised, oct 2012.
- Boyd, S. P. and Vandenberghe, L. (2004). *Convex optimization*. Cambridge University Press, Cambridge, UK ; New York.
- Bryson, A. E. and Ho, Y.-C. (1969). *Applied optimal control: optimization, estimation, and control*. Blaisdell Publishing Co.
- Chen, T. Q., Rubanova, Y., Bettencourt, J., and Duvenaud, D. (2018). Neural ordinary differential equations. In *Advances in Neural Information Processing Systems 31: Annual Conference on Neural Information Processing Systems 2018, NeurIPS 2018, December 3-8, 2018, Montréal, Canada*, pages 6572–6583.
- De O. Pantoja, J. and Mayne, D. (1989). A sequential quadratic programming algorithm for discrete optimal control problems with control inequality constraints. In *Proceedings of the 28th IEEE Conference on Decision and Control*,, pages 353–357 vol.1.
- Dreyfus, S. E. (1990). Artificial neural networks, back propagation, and the kelley-bryson gradient procedure. *Journal of Guidance, Control, and Dynamics*, **13**(5), 926–928.
- Duchi, J., Hazan, E., and Singer, Y. (2011). Adaptive subgradient methods for online learning and stochastic optimization. *Journal of Machine Learning Research*, **12**(61), 2121–2159.
- Dunn, J. and Bertsekas, D. (1989). Efficient dynamic programming implementations of newton’s method for unconstrained optimal control problems. *Journal of Optimization Theory and Applications*, **63**, 23–38.
- Gilbert, J. C. (1992). Automatic differentiation and iterative processes. *Optimization methods and software*, **1**(1), 13–21.
- Glorot, X., Bordes, A., and Bengio, Y. (2011). Deep sparse rectifier neural networks. In *Proceedings of the Fourteenth International Conference on Artificial Intelligence and Statistics*, volume 15 of *Proceedings of Machine Learning Research*, pages 315–323, Fort Lauderdale, FL, USA. PMLR.
- Golub, G. H. and Varga, R. S. (2007). Chebyshev semi-iterative methods, successive over-relaxation iterative methods, and second-order richardson iterative methods, parts I and II. In *Milestones in Matrix Computation - Selected Works of Gene H. Golub, with Commentaries*, pages 45–67. Oxford University Press.

- Goodfellow, I. J. and Vinyals, O. (2015). Qualitatively characterizing neural network optimization problems. In Y. Bengio and Y. LeCun, editors, *3rd International Conference on Learning Representations, ICLR 2015, San Diego, CA, USA, May 7-9, 2015, Conference Track Proceedings*.
- Griewank, A. and Naumann, U. (2003). Accumulating jacobians as chained sparse matrix products. *Math. Program.*, **95**(3), 555–571.
- Haber, E. and Ruthotto, L. (2017). Stable architectures for deep neural networks. *Inverse Problems*, **34**(1), 014004.
- He, K., Zhang, X., Ren, S., and Sun, J. (2016). Deep residual learning for image recognition. In *2016 IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2016, Las Vegas, NV, USA, June 27-30, 2016*, pages 770–778. IEEE Computer Society.
- Hestenes, M. and Stiefel, E. (1952). Methods of conjugate gradients for solving linear systems. *Journal of Research of the National Bureau of Standards*, **49**(6), 409.
- Hinton, G. (2012). Neural networks for machine learning. Coursera, video lectures.
- Jarrett, K., Kavukcuoglu, K., Ranzato, M., and LeCun, Y. (2009). What is the best multi-stage architecture for object recognition? In *2009 IEEE 12th International Conference on Computer Vision*, pages 2146–2153.
- Jin, W., Wang, Z., Yang, Z., and Mou, S. (2020). Pontryagin differentiable programming: An end-to-end learning and control framework. In *Advances in Neural Information Processing Systems*, volume 33, pages 7979–7992. Curran Associates, Inc.
- Kingma, D. P. and Ba, J. (2015). Adam: A method for stochastic optimization. In *3rd International Conference on Learning Representations, ICLR 2015, San Diego, CA, USA, May 7-9, 2015, Conference Track Proceedings*.
- Krichene, W., Bayen, A., and Bartlett, P. L. (2015). Accelerated mirror descent in continuous and discrete time. In *Advances in Neural Information Processing Systems*, volume 28. Curran Associates, Inc.
- Le, Q. V., Ngiam, J., Coates, A., Lahiri, A., Prochnow, B., and Ng, A. Y. (2011). On optimization methods for deep learning. In *Proceedings of the 28th International Conference on Machine Learning, ICML 2011, Bellevue, Washington, USA, June 28 - July 2, 2011*, pages 265–272. Omnipress.
- Lebanon, G. (2012). *Probability: The Analysis of Data, Volume 1*. CreateSpace Independent Publishing Platform.

- Lecun, Y. (1988). A theoretical framework for back-propagation. In *Proceedings of the 1988 Connectionist Models Summer School, CMU, Pittsburg, PA*, pages 21–28. Morgan Kaufmann.
- Lenhart, S. and Workman, J. T. (2007). *Optimal Control Applied to Biological Models*. Chapman and Hall/CRC.
- Linnainmaa, S. (1976). Taylor expansion of the accumulated rounding error. *BIT*, **16**(2), 146–160.
- Liu, D. C. and Nocedal, J. (1989). On the limited memory bfgs method for large scale optimization. *Mathematical programming*, **45**(1), 503–528.
- Martens, J. (2016). *Second-order optimization for neural networks*. University of Toronto (Canada).
- Martens, J. (2020). New insights and perspectives on the natural gradient method. *J. Mach. Learn. Res.*, **21**, 146:1–146:76.
- Martens, J. and Grosse, R. B. (2015). Optimizing neural networks with kronecker-factored approximate curvature. In *Proceedings of the 32nd International Conference on Machine Learning, ICML 2015, Lille, France, 6-11 July 2015*, volume 37 of *JMLR Workshop and Conference Proceedings*, pages 2408–2417. JMLR.org.
- Martens, J. and Sutskever, I. (2011). Learning recurrent neural networks with hessian-free optimization. In L. Getoor and T. Scheffer, editors, *Proceedings of the 28th International Conference on Machine Learning, ICML 2011, Bellevue, Washington, USA, June 28 - July 2, 2011*, pages 1033–1040. Omnipress.
- Miquel, T. (2020). Introduction to Optimal Control. Lecture.
- Miyato, T., Kataoka, T., Koyama, M., and Yoshida, Y. (2018). Spectral normalization for generative adversarial networks. In *6th International Conference on Learning Representations, ICLR 2018, Vancouver, BC, Canada, April 30 - May 3, 2018, Conference Track Proceedings*. OpenReview.net.
- Nayfeh, A. H. and Balachandran, B. (1995). *Equilibrium Solutions*, chapter 2, pages 35–145. John Wiley & Sons, Ltd.
- Nemirovsky, A. S. and Yudi, D. B. (1983). *Problem Complexity and Method Efficiency in Optimization*. Wiley-Interscience series in discrete mathematics. Wiley.
- Nesterov, Y. (1983). A method for solving the convex programming problem with convergence rate  $O(1/k^2)$ . *Proceedings of the USSR Academy of Sciences*, **269**, 543–547.

- Neumann, C. (1877). *Untersuchungen über das logarithmische und Newton'sche Potential*. BG Teubner.
- Nocedal, J. and Wright, S. (2006). *Numerical Optimization*. Springer Series in Operations Research and Financial Engineering. Springer New York.
- Osawa, K., Tsuji, Y., Ueno, Y., Naruse, A., Yokota, R., and Matsuoka, S. (2019). Large-scale distributed second-order optimization using kronecker-factored approximate curvature for deep convolutional neural networks. In *IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2019, Long Beach, CA, USA, June 16-20, 2019*, pages 12359–12367. Computer Vision Foundation / IEEE.
- Pascanu, R. and Bengio, Y. (2014). Revisiting natural gradient for deep networks. In *2nd International Conference on Learning Representations, ICLR 2014, Banff, AB, Canada, April 14-16, 2014, Conference Track Proceedings*.
- Pearlmutter, B. A. (1994). Fast exact multiplication by the hessian. *Neural Comput.*, **6**(1), 147–160.
- Pineda, F. J. (1987). Generalization of backpropagation to recurrent and higher order neural networks. In *Proceedings of the 1987 International Conference on Neural Information Processing Systems, NIPS'87*, page 602–611, Cambridge, MA, USA. MIT Press.
- Polyak, B. (1964). Some methods of speeding up the convergence of iteration methods. *USSR Computational Mathematics and Mathematical Physics*, **4**(5), 1–17.
- Pontriagin, L., Boltyanskii, V., Collection, K. M. R., Neustadt, L., Gamkrelidze, R., Tirogoff, K., Mishchenko, E., and Brown, D. (1962). *The Mathematical Theory of Optimal Processes*. International series of monographs in pure and applied mathematics. Interscience Publishers.
- Raskutti, G. and Mukherjee, S. (2015). The information geometry of mirror descent. *IEEE Trans. Inf. Theory*, **61**(3), 1451–1457.
- Reid, W. (1972). *Riccati Differential Equations*. Mathematics in science and engineering : a series of monographs and textbooks. Academic Press.
- Robbins, H. (2007). A stochastic approximation method. *Annals of Mathematical Statistics*, **22**, 400–407.
- Ruder, S. (2016). An overview of gradient descent optimization algorithms. *CoRR*, **abs/1609.04747**.
- Rudin, W. (1976). *Principles of Mathematical Analysis*. International series in pure and applied mathematics. McGraw-Hill.

- Seidman, J. H., Fazlyab, M., Preciado, V. M., and Pappas, G. J. (2020). Robust deep learning as optimal control: Insights and convergence guarantees.
- Sutskever, I., Martens, J., Dahl, G., and Hinton, G. (2013). On the importance of initialization and momentum in deep learning. In *Proceedings of the 30th International Conference on Machine Learning*, volume 28 of *Proceedings of Machine Learning Research*, pages 1139–1147, Atlanta, Georgia, USA. PMLR.
- Suzuki, N. (1976). On the convergence of neumann series in banach space. *Mathematische Annalen*, **220**(2), 143–146.
- Yoshida, Y. and Miyato, T. (2017). Spectral norm regularization for improving the generalizability of deep learning. *ArXiv*, **abs/1705.10941**.
- Zhang, D., Zhang, T., Lu, Y., Zhu, Z., and Dong, B. (2019). You only propagate once: Accelerating adversarial training via maximal principle. In *Advances in Neural Information Processing Systems 32: Annual Conference on Neural Information Processing Systems 2019, NeurIPS 2019, December 8-14, 2019, Vancouver, BC, Canada*, pages 227–238.



# Appendix A

---

## Steepest descent – Minimization Without Constraints

**Theorem A.0.1.** *The steepest descent direction can be obtained via the following minimization problem.*

$$\mathbf{v}^* = \alpha \left[ \underset{\mathbf{v}}{\operatorname{argmin}} \left\{ \nabla f(\mathbf{x})^T \mathbf{v} + \frac{1}{2} \|\mathbf{v}\|^2 \right\} \right] \quad (\text{A.0.1})$$

**Proof.** First, let us suppose that  $\mathbf{w}^*$  is the solution to normalized steepest descent problem, that is :

$$\mathbf{w}^* = \underset{\mathbf{v}}{\operatorname{argmin}} \nabla f(\mathbf{x})^T \mathbf{v} \quad \text{s.t.} \quad \|\mathbf{v}\| = 1$$

Now, in (A.0.1), we can reexpress  $\mathbf{v}$  as  $\mathbf{v} = t\mathbf{w}$ ,  $t \geq 0$  being a scalar and with  $\mathbf{w} = \mathbf{v}/\|\mathbf{v}\|$ . With this substitution, we are facing a minimization problem over two variables, that we can try to optimize successively.

By first taking  $\mathbf{w}$  as fixed and minimizing for  $t$  we have:

$$\hat{t} = \underset{t}{\operatorname{argmin}} \left\{ t \nabla f(\mathbf{x})^T \mathbf{w} + \frac{t^2}{2} \right\}$$

This is a quadratic problem for which the minimum can be obtained by putting the derivative to zero:

$$\begin{aligned} 0 &= t + \nabla f(\mathbf{x})^T \mathbf{w} \\ \implies t &= -\nabla f(\mathbf{x})^T \mathbf{w} \end{aligned}$$

This in turn implies that  $\hat{t} = 0$  if we choose  $w$  s.t.  $\nabla f(\mathbf{x})^T \mathbf{w} > 0$  since we have that  $t \geq 0$ . And therefore, to minimize  $t\nabla f(\mathbf{x})^T \mathbf{w} + \frac{t^2}{2}$  we must pick a  $\mathbf{w}$  s.t.  $\nabla f(\mathbf{x})^T \mathbf{w} \leq 0$  and  $\hat{t} = -\nabla f(\mathbf{x})^T \mathbf{w}$

Turning our attention toward the minimization over  $\mathbf{w}$  we have:

$$\begin{aligned} \hat{\mathbf{w}} &= \arg \min_{\mathbf{w}, \|\mathbf{w}\|=1} \left\{ t\nabla f(\mathbf{x})^T \mathbf{w} + \frac{t^2}{2} \right\} \\ &= \arg \min_{\mathbf{w}, \|\mathbf{w}\|=1} \{ t\nabla f(\mathbf{x})^T \mathbf{w} \} \\ &= \arg \min_{\mathbf{w}, \|\mathbf{w}\|=1} \{ -(\nabla f(\mathbf{x})^T \mathbf{w})^2 \} \\ &= \arg \max_{\mathbf{w}, \|\mathbf{w}\|=1} \{ |\nabla f(\mathbf{x})^T \mathbf{w}| \} \end{aligned}$$

A problem for which we already know the solutions, that is  $\hat{\mathbf{w}} = \pm \mathbf{w}^*$ . But since we require  $\nabla f(\mathbf{x})^T \mathbf{w} \leq 0$ , we get that  $\hat{\mathbf{w}} = \mathbf{w}^*$ .

Therefore, we have that:

$$\begin{aligned} \arg \min_{\mathbf{v}} \left\{ \nabla f(\mathbf{x})^T \mathbf{v} + \frac{1}{2} \|\mathbf{v}\|^2 \right\} &= \hat{t} \hat{\mathbf{w}} \\ &= -\nabla f(\mathbf{x})^T \hat{\mathbf{w}} \hat{\mathbf{w}} \\ &= -\nabla f(\mathbf{x})^T \mathbf{w}^* \mathbf{w}^* \\ &= \alpha \mathbf{w}^* \\ &\text{with } \alpha = -\nabla f(\mathbf{x})^T \mathbf{w}^* \geq 0 \end{aligned}$$

□

# Appendix B

---

## Riccati's Equation – Solution with a Linear Term in the Final Cost

Since an exhaustive development of the discrete Riccati equations for the case of an LQR with a linear term in final cost was hard to come by, we provide one below.

We seek the Riccati equation for the following problem, which is the problem we encounter when transforming the steepest descent step into a LQR:

$$\begin{aligned} \min_{\mathbf{u}} \quad & \mathbf{a}_N \mathbf{x}_N + \mathbf{x}_N^T \mathbf{Q}_N \mathbf{x}_N \sum_{i=0}^{N-1} \left[ \frac{1}{2} \mathbf{x}_i^T \mathbf{Q} \mathbf{x}_i + \frac{1}{2} u_i^T \mathbf{R} u_i + \mathbf{u}_i^T \mathbf{M} \mathbf{x}_i \right] \\ \text{subject to} \quad & \mathbf{x}_{i+1} = \mathbf{A}_i \mathbf{x}_i + \mathbf{B}_i \mathbf{u}_i \\ \text{given} \quad & \mathbf{x}_0 \end{aligned}$$

With the Hamiltonian functions given by :

$$\begin{aligned} H_i(\mathbf{x}_i, \mathbf{u}_i) &= \mathbf{x}_i^T \mathbf{Q} \mathbf{x}_i + \frac{1}{2} u_i^T \mathbf{R} u_i + \mathbf{u}_i^T \mathbf{M} \mathbf{x}_i + \mathbf{p}_{i+1} (\mathbf{A}_i \mathbf{x}_i + \mathbf{B}_i \mathbf{u}_i) \\ H_N(\mathbf{x}_N) &= \mathbf{a}_N \mathbf{x}_N + \mathbf{x}_N^T \mathbf{Q}_N \mathbf{x}_N \end{aligned}$$

The adjoints are defined by  $\mathbf{p}_i = \nabla_{\mathbf{x}_i} H_i$ , and in particular when evaluated at the optimal pair  $(\mathbf{x}_i^*, \mathbf{u}_i^*)$ :

$$\begin{aligned} \mathbf{p}_i &= \nabla_{\mathbf{x}_i} H_i(\mathbf{x}_i^*, \mathbf{u}_i^*) \\ &= \mathbf{Q}_i \mathbf{x}_i^* + \mathbf{M}_i^T \mathbf{u}_i^* + \mathbf{A}_i^T \mathbf{p}_{i+1}^* \end{aligned} \tag{B.0.1}$$

They are also chosen such that they satisfy  $\nabla_{\mathbf{u}_i} H_i = 0$ , so that:

$$\begin{aligned}
0 &= \nabla_{\mathbf{u}_i} H_i(\mathbf{x}_i^*, \mathbf{u}_i^*) \\
&= \mathbf{R}_i \mathbf{u}_i^* + \mathbf{M}_i \mathbf{x}_i^* + \mathbf{B}_i^T \mathbf{p}_{i+1}^* \\
&\implies \mathbf{u}_i^* = -\mathbf{R}_i^{-1} (\mathbf{M}_i \mathbf{x}_i^* + \mathbf{B}_i^T \mathbf{p}_{i+1}^*)
\end{aligned}$$

We now want to prove by induction that we can rewrite  $\mathbf{p}_i^*$  as:

$$\mathbf{p}_i^* = \mathbf{K}_i \mathbf{x}_i^* + \boldsymbol{\lambda}_i \quad (\text{B.0.2})$$

$\forall i = 0, \dots, N-1$  starting with  $i = N-1$  with  $\mathbf{p}_N^* = \mathbf{Q}_N \mathbf{x}_N^* + \mathbf{a}_N$

Supposing it holds for  $i+1$ , we have that:

$$\begin{aligned}
\mathbf{u}_i^* &= -\mathbf{R}_i^{-1} (\mathbf{M}_i \mathbf{x}_i^* + \mathbf{B}_i^T \mathbf{K}_{i+1} \mathbf{x}_{i+1}^* + \mathbf{B}_i^T \boldsymbol{\lambda}_{i+1}) \\
&= -\mathbf{R}_i^{-1} (\mathbf{M}_i \mathbf{x}_i^* + \mathbf{B}_i^T \mathbf{K}_{i+1} (\mathbf{A}_i \mathbf{x}_i^* + \mathbf{B}_i \mathbf{u}_i^*) + \mathbf{B}_i^T \boldsymbol{\lambda}_{i+1}) \\
\mathbf{R}_i \mathbf{u}_i^* + \mathbf{B}_i^T \mathbf{K}_{i+1} \mathbf{B}_i \mathbf{u}_i^* &= -(\mathbf{M}_i \mathbf{x}_i^* + \mathbf{B}_i^T \mathbf{K}_{i+1} \mathbf{A}_i \mathbf{x}_i^* + \mathbf{B}_i^T \boldsymbol{\lambda}_{i+1}) \\
\mathbf{u}_i^* &= -(\mathbf{R}_i + \mathbf{B}_i^T \mathbf{K}_{i+1} \mathbf{B}_i)^{-1} ((\mathbf{M}_i + \mathbf{B}_i^T \mathbf{K}_{i+1} \mathbf{A}_i) \mathbf{x}_i^* + \mathbf{B}_i^T \boldsymbol{\lambda}_{i+1})
\end{aligned} \quad (\text{B.0.3})$$

Using this result, we can rewrite  $\mathbf{x}_{i+1}^*$  as :

$$\mathbf{x}_{i+1}^* = \mathbf{A}_i \mathbf{x}_i^* + \mathbf{B}_i \mathbf{u}_i^* = \mathbf{A}_i \mathbf{x}_i^* - \mathbf{B}_i (\mathbf{R}_i + \mathbf{B}_i^T \mathbf{K}_{i+1} \mathbf{B}_i)^{-1} ((\mathbf{M}_i + \mathbf{B}_i^T \mathbf{K}_{i+1} \mathbf{A}_i) \mathbf{x}_i^* + \mathbf{B}_i^T \boldsymbol{\lambda}_{i+1})$$

Multiplying each side by  $\mathbf{K}_{i+1}$  and then adding  $\boldsymbol{\lambda}_{i+1}$  we get from (B.0.2):

$$\mathbf{p}_{i+1}^* = \mathbf{K}_{i+1} \mathbf{A}_i \mathbf{x}_i^* - \mathbf{K}_{i+1} \mathbf{B}_i (\mathbf{R}_i + \mathbf{B}_i^T \mathbf{K}_{i+1} \mathbf{B}_i)^{-1} ((\mathbf{M}_i + \mathbf{B}_i^T \mathbf{K}_{i+1} \mathbf{A}_i) \mathbf{x}_i^* + \mathbf{B}_i^T \boldsymbol{\lambda}_{i+1}) + \boldsymbol{\lambda}_{i+1}$$

Now multiplying each side by  $\mathbf{A}_i^T$ :

$$\mathbf{A}_i^T \mathbf{p}_{i+1}^* = \mathbf{A}_i^T \mathbf{K}_{i+1} \mathbf{A}_i \mathbf{x}_i^* - \mathbf{A}_i^T \mathbf{K}_{i+1} \mathbf{B}_i (\mathbf{R}_i + \mathbf{B}_i^T \mathbf{K}_{i+1} \mathbf{B}_i)^{-1} ((\mathbf{M}_i + \mathbf{B}_i^T \mathbf{K}_{i+1} \mathbf{A}_i) \mathbf{x}_i^* + \mathbf{B}_i^T \boldsymbol{\lambda}_{i+1}) + \mathbf{A}_i^T \boldsymbol{\lambda}_{i+1}$$

And finally adding  $(\mathbf{Q}_i \mathbf{x}_i^* + \mathbf{M}_i^T \mathbf{u}_i^*)$  on each side, replacing  $\mathbf{u}_i^*$  its expression in (B.0.3):

$$\begin{aligned}
\mathbf{p}_i^* &= \mathbf{Q}_i \mathbf{x}_i^* + \mathbf{M}_i^T \mathbf{u}_i^* + \mathbf{A}_i^T \mathbf{p}_{i+1}^* \quad \text{from (B.0.1)} \\
&= \left( \mathbf{A}_i^T \mathbf{K}_{i+1} \mathbf{A}_i + \mathbf{Q}_i - (\mathbf{A}_i^T \mathbf{K}_{i+1} \mathbf{B}_i + \mathbf{M}_i^T) (\mathbf{R}_i + \mathbf{B}_i^T \mathbf{K}_{i+1} \mathbf{B}_i)^{-1} (\mathbf{M}_i + \mathbf{B}_i^T \mathbf{K}_{i+1} \mathbf{A}_i) \right) \mathbf{x}_i^* \\
&\quad + \mathbf{A}_i^T \boldsymbol{\lambda}_{i+1} - (\mathbf{A}_i^T \mathbf{K}_{i+1} \mathbf{B}_i + \mathbf{M}_i^T) (\mathbf{R}_i + \mathbf{B}_i^T \mathbf{K}_{i+1} \mathbf{B}_i)^{-1} \mathbf{B}_i^T \boldsymbol{\lambda}_{i+1}
\end{aligned}$$

Completing the induction proof while also giving convenient expressions for  $\mathbf{K}_i$  and  $\boldsymbol{\lambda}_i$ :

$$\mathbf{K}_i = \mathbf{A}_i^T \mathbf{K}_{i+1} \mathbf{A}_i + \mathbf{Q}_i - (\mathbf{A}_i^T \mathbf{K}_{i+1} \mathbf{B}_i + \mathbf{M}_i^T) (\mathbf{R}_i + \mathbf{B}_i^T \mathbf{K}_{i+1} \mathbf{B}_i)^{-1} (\mathbf{M}_i + \mathbf{B}_i^T \mathbf{K}_{i+1} \mathbf{A}_i) \tag{B.0.4}$$

$$\boldsymbol{\lambda}_i = \mathbf{A}_i^T \boldsymbol{\lambda}_{i+1} - (\mathbf{A}_i^T \mathbf{K}_{i+1} \mathbf{B}_i + \mathbf{M}_i^T) (\mathbf{R}_i + \mathbf{B}_i^T \mathbf{K}_{i+1} \mathbf{B}_i)^{-1} \mathbf{B}_i^T \boldsymbol{\lambda}_{i+1} \tag{B.0.5}$$

The optimal trajectory can then be recovered from the system dynamics and from (B.0.3):

$$\mathbf{u}_0^* = -(\mathbf{R}_0 + \mathbf{B}_0^T \mathbf{K}_1 \mathbf{B}_0)^{-1} \mathbf{B}_0^T \boldsymbol{\lambda}_1 \tag{B.0.6}$$

$$\mathbf{x}_{i+1}^* = \mathbf{A}_i \mathbf{x}_i^* + \mathbf{B}_i \mathbf{u}_i^* \tag{B.0.7}$$

$$\mathbf{u}_i^* = (\mathbf{R}_i + \mathbf{B}_i^T \mathbf{K}_{i+1} \mathbf{B}_i)^{-1} ((\mathbf{B}_i^T \mathbf{K}_{i+1} \mathbf{A}_i) \mathbf{x}_i^* + \mathbf{B}_i^T \boldsymbol{\lambda}_{i+1}) \tag{B.0.8}$$