

QPLayer: Leveraging augmented-Lagrangian techniques for differentiating over infeasible quadratic programs

A. Bambade^{1,2} F. Schramm¹ A. Taylor² J. Carpentier¹

¹Willow team
Inria and ENS Paris

²Sierra team
Inria and ENS Paris

PGMO days, 29 November 2023

Table of Contents

- 1 Introduction
- 2 Solution outline
- 3 Solving the closest feasible problem
- 4 Differentiating closest feasible problem solutions
- 5 Numerical results
- 6 Conclusion

Table of Contents

- 1 Introduction
- 2 Solution outline
- 3 Solving the closest feasible problem
- 4 Differentiating closest feasible problem solutions
- 5 Numerical results
- 6 Conclusion

Introduction: QP layers v.s. neural networks

Outputs of current learning pipelines are explicit function of the inputs.

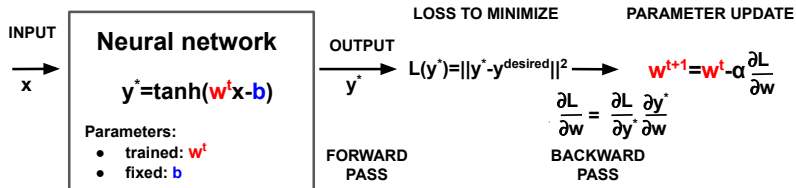


Figure: Example of a feedforward neural network.

Introduction: QP layers v.s. neural networks

More recent literature considers differentiable optimization problems as layers.

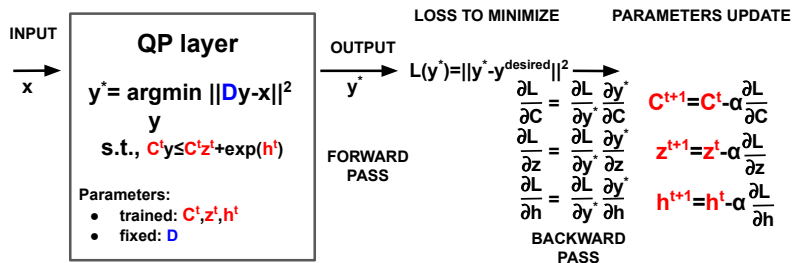


Figure: Example of a Quadratic Programming layer (with D nonsingular).

Introduction: QP layers pros (representative power)

Why using such more complex architecture ?

- For some class of layers, the representative power is similar to those of standard Neural Networks,

Theorem (B. Amos & Z. Kolter (2021))

Let $f : \mathbb{R}^n \rightarrow \mathbb{R}^n$ be an elementwise piecewise linear function with k linear regions. Then function can be represented as an OptNet layer using $O(nk)$ parameters. Additionally, the layer $z_{i+1} = \max(Wz_i + b, 0)$ for $W \in \mathbb{R}^{n \times m}$, $b \in \mathbb{R}^m$ can be represented by an OptNet layer with $O(mn)$ parameters.

Introduction: QP layers pros (representative power)

Why using such more complex architecture ?

- For some class of layers, the representative power is similar to those of standard Neural Networks,
- For some optimization based problem, it performs better.

QP layers in machine learning: example 1

Using a convex QP as a deep learning layer performs better than a ConvNet for solving Sudokus.

			3
1			
		4	
4			1

2	4	1	3
1	3	2	4
3	1	4	2
4	2	3	1

Figure: Example of Sudoku.

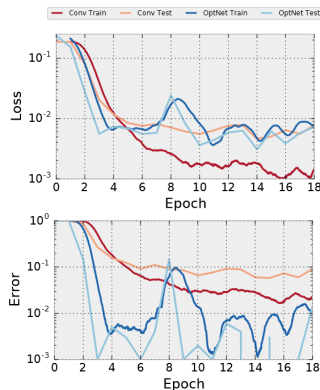


Figure: Sudoku Training plots¹.

¹Brandon Amos and J. Zico Kolter (2021). OptNet:Differentiable Optimization as a Layer in Neural Networks.

Introduction: QP layers pros (practical speed)

Why using such more complex architecture ?

- For some class of layers, the representative power is similar to those of standard Neural Networks,
- For some optimization based problem, it performs better,
- Achievable practical use.

Introduction: QP layers pros (practical speed)

Trick: use the Implicit Function Theorem.

$$\begin{aligned} \min_{x \in \mathbb{R}^n} & \left\{ f(x; \theta) \triangleq \frac{1}{2} x^\top H(\theta) x + x^\top g(\theta) \right\} \\ \text{s.t.} & C(\theta) x \leq u(\theta), \end{aligned} \quad (\text{QP}(\theta))$$

Introduction: QP layers pros (practical speed)

Trick: use the Implicit Function Theorem.

$$\begin{aligned} \min_{x \in \mathbb{R}^n} & \left\{ f(x; \theta) \triangleq \frac{1}{2} x^\top H(\theta) x + x^\top g(\theta) \right\} \\ \text{s.t.} & C(\theta) x \leq u(\theta), \end{aligned} \quad (\text{QP}(\theta))$$

Noting $(x^*(\theta), z^*(\theta))$ a solution to $(\text{QP}(\theta))$, applying the IFT outputs:

$$\begin{aligned} & \begin{bmatrix} H & C^\top \\ D(z^*)C & D(Cx^* - u) \end{bmatrix} \begin{bmatrix} \frac{\partial x^*}{\partial \theta} \\ \frac{\partial z^*}{\partial \theta} \end{bmatrix} \\ &= - \begin{bmatrix} \frac{\partial H}{\partial \theta} x^* + \frac{\partial g}{\partial \theta} + \frac{\partial C^\top}{\partial \theta} z^* \\ D(z^*) \left(\frac{\partial C}{\partial \theta} x^* - \frac{\partial u}{\partial \theta} \right) \end{bmatrix} \end{aligned}$$

Introduction: QP layers pros (practical speed)

Trick: use the Implicit Function Theorem.

$$\min_{x \in \mathbb{R}^n} \left\{ f(x; \theta) \triangleq \frac{1}{2} x^\top H(\theta) x + x^\top g(\theta) \right\} \quad (\text{QP}(\theta))$$
$$\text{s.t. } C(\theta) x \leq u(\theta),$$

Noting $(x^*(\theta), z^*(\theta))$ a solution to $(\text{QP}(\theta))$, applying the IFT outputs:

$$\begin{aligned} & \begin{bmatrix} H & C^\top \\ D(z^*)C & D(Cx^* - u) \end{bmatrix} \begin{bmatrix} \frac{\partial x^*}{\partial \theta} \\ \frac{\partial z^*}{\partial \theta} \end{bmatrix} \\ &= - \begin{bmatrix} \frac{\partial H}{\partial \theta} x^* + \frac{\partial g}{\partial \theta} + \frac{\partial C}{\partial \theta}^\top z^* \\ D(z^*) \left(\frac{\partial C}{\partial \theta} x^* - \frac{\partial u}{\partial \theta} \right) \end{bmatrix} \end{aligned}$$
$$\begin{aligned} \frac{\partial \mathcal{L}}{\partial \theta} &= \begin{bmatrix} \frac{\partial \mathcal{L}}{\partial x} \\ \frac{\partial \mathcal{L}}{\partial z} \end{bmatrix}^\top \begin{bmatrix} \frac{\partial x}{\partial \theta} \\ \frac{\partial z}{\partial \theta} \end{bmatrix} \\ &= \left(- \begin{bmatrix} H & C^\top D(z^*) \\ C & D(Cx^* - u) \end{bmatrix} \begin{bmatrix} d_x \\ d_z \end{bmatrix} \right)^\top \begin{bmatrix} \frac{\partial x}{\partial \theta} \\ \frac{\partial z}{\partial \theta} \end{bmatrix} \end{aligned}$$

Introduction: cons (limited trainable architectures)

Two issues

- IFT assumptions (non singularity of KKT matrix etc.),
- Structural feasibility at training and test time.

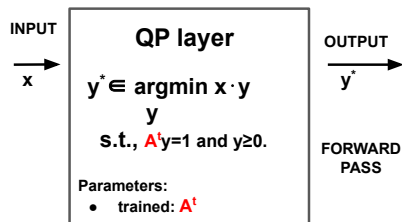


Figure: A LP layer. Nothing guarantees during training that the vector of 1 lies in the range space of the A^t .

Table of Contents

- 1 Introduction
- 2 Solution outline**
- 3 Solving the closest feasible problem
- 4 Differentiating closest feasible problem solutions
- 5 Numerical results
- 6 Conclusion

Solution outline: ideal pipeline

General idea: consider a "broader" problem.

Solution outline: ideal pipeline

General idea: consider a "broader" problem.

- Forward pass: solve the closest feasible QP.

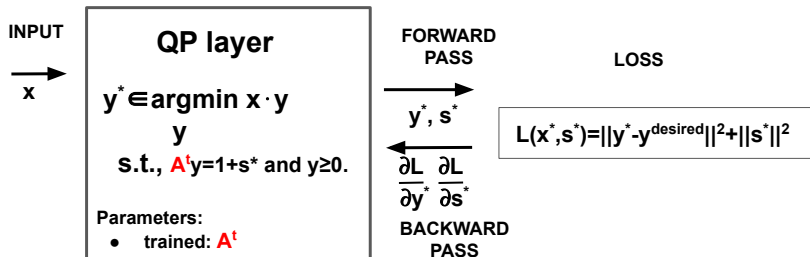


Figure: General solution outline.

Solution outline: ideal pipeline

General idea: consider a "broader" problem.

- Forward pass: solve the closest feasible QP.
- Backward pass: differentiate through one solution.

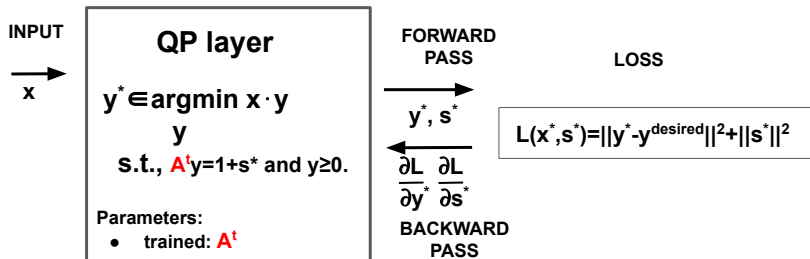


Figure: General solution outline.

Solution outline: plan

General plan:

- Forward pass: solve the closest feasible QP.

Solution outline: plan

General plan:

- Forward pass: solve the closest feasible QP.
 - Define the closest feasible problem.
 - Propose an efficient algorithm to solve it.

General plan:

- Forward pass: solve the closest feasible QP.
 - Define the closest feasible problem.
 - Propose an efficient algorithm to solve it.
- Backward pass: differentiate through one solution.

General plan:

- Forward pass: solve the closest feasible QP.
 - Define the closest feasible problem.
 - Propose an efficient algorithm to solve it.
- Backward pass: differentiate through one solution.
 - Prove applying the IFT makes sense.
 - Provide tractable algorithms.
 - In non differentiable case, provide alternatives.

Table of Contents

- 1 Introduction
- 2 Solution outline
- 3 Solving the closest feasible problem**
- 4 Differentiating closest feasible problem solutions
- 5 Numerical results
- 6 Conclusion

The closest feasible QP problem: definition

$$x^*(\theta) \in \operatorname{argmin}_{x \in \mathbb{R}^n} \left\{ f(x; \theta) \triangleq \frac{1}{2} x^\top H(\theta) x + x^\top g(\theta) \right\} \quad (\text{QP}(\theta))$$

s.t. $C(\theta) x \leq u(\theta)$,

where $H(\theta) \in \mathcal{S}_+^n(\mathbb{R})$, $g(\theta) \in \mathbb{R}^n$, $C(\theta) \in \mathbb{R}^{n_i \times n}$ and $u(\theta) \in \mathbb{R}^{n_i}$.

The closest feasible QP problem: definition

$$x^*(\theta) \in \operatorname{argmin}_{x \in \mathbb{R}^n} \left\{ f(x; \theta) \triangleq \frac{1}{2} x^\top H(\theta) x + x^\top g(\theta) \right\} \quad (\text{QP}(\theta))$$

s.t. $C(\theta) x \leq u(\theta)$,

where $H(\theta) \in \mathcal{S}_+^n(\mathbb{R})$, $g(\theta) \in \mathbb{R}^n$, $C(\theta) \in \mathbb{R}^{n_i \times n}$ and $u(\theta) \in \mathbb{R}^{n_i}$.

Assumption

$H(\theta)$ is symmetric positive definite in the direction of $g(\theta)$ or $g(\theta)$ is orthogonal to the recession cone of $\text{QP}(\theta)$, i.e.,
 $g(\theta) \perp C^\infty(\theta) \triangleq \{y \in \mathbb{R}^n \mid C(\theta)[x + \tau y] \leq u(\theta) \text{ s.t. } C(\theta)x \leq u(\theta), \tau \geq 0\}$.

The closest feasible QP problem: definition

Under Assumption (1) the closest feasible problem (QP-H(θ)) is well-posed:

$$\begin{aligned} s^*(\theta) = \arg \min_{s \in \mathbb{R}^{n_i}} \frac{1}{2} \|s\|_2^2 \\ \text{s.t. } x^*(\theta), z^*(\theta) \in \arg \min_{x \in \mathbb{R}^n} \max_{z \in \mathbb{R}_+^{n_i}} L(x, z, s; \theta), \end{aligned} \quad (\text{QP-H}(\theta))$$

with $L(x, z, s; \theta) \triangleq \frac{1}{2} x^\top H(\theta) x + x^\top g(\theta) + z^\top (C(\theta)x - u(\theta) - s)$

The closest feasible QP problem: definition

Under Assumption (1) the closest feasible problem (QP-H(θ)) is well-posed:

$$s^*(\theta) = \arg \min_{s \in \mathbb{R}^{n_i}} \frac{1}{2} \|s\|_2^2 \quad \text{(QP-H}(\theta)\text{)}$$
$$\text{s.t. } x^*(\theta), z^*(\theta) \in \arg \min_{x \in \mathbb{R}^n} \max_{z \in \mathbb{R}_+^{n_i}} L(x, z, s; \theta),$$

with $L(x, z, s; \theta) \triangleq \frac{1}{2} x^\top H(\theta) x + x^\top g(\theta) + z^\top (C(\theta)x - u(\theta) - s)$

Remark

$$s^* \triangleq \arg \min_{s \in \mathcal{R}(C) +]-\infty, u]} \|s\|_2^2$$

The closest feasible QP problem: solution method

The Augmented Lagrangian

$$\mathcal{L}_A(x, z; \mu) \triangleq f(x) + \frac{1}{2\mu} (\| [Cx - u + \mu z]_+ \|_2^2 - \|\mu z\|_2^2)$$

The AL method

$$x^{k+1} = \arg \min_x \mathcal{L}_A(x, z; \mu)$$
$$z^{k+1} = \left[\frac{1}{\mu} (Cx^{k+1} - u) + z^k \right]_+$$

Theorem (A. Chiche, JC Gilbert (2016)¹)

Under Assumption (1), the revised AL algorithm does not terminate with a direction of unboundedness and generates a sequence $\{(x^k, z^k)\}$ converging towards a solution to (QP-H(θ)).

Remark

(i) revised AL algorithm = modified stopping criterion; (ii) exact AL method!

¹Chiche, A., Gilbert, J. C. (2016). How the augmented Lagrangian algorithm can deal with an infeasible convex quadratic optimization problem.

The closest feasible QP problem: solution method in practice

The ProxSuite library

ProxSuite

ADVANCED OPTIMIZERS FOR ROBOTICS AND BEYOND

- ✓ **fast:** C++ implementation, with homemade linear Cholesky solver
- ✓ **scalable:** various backends for dense, sparse and matrix-free optimization
- ✓ **easy-to-use:** API closed to OSQP, Python and Julia bindings
- ✓ **open-source:** BSD-license, easily installable



github.com/Simple-Robotics/proxsuite



```
conda install -c conda-forge proxsuite
```



Table of Contents

- 1 Introduction
- 2 Solution outline
- 3 Solving the closest feasible problem
- 4 Differentiating closest feasible problem solutions**
- 5 Numerical results
- 6 Conclusion

Differentiating (QP-H(θ)) solutions: methodology

- Analyze the properties of KKT map for (QP-H(θ)),
- Extended Conservative Jacobian definition.
- Practical algorithms.

Differentiating (QP-H(θ)) solutions: KKT conditions of (QP-H(θ))

$$G(x, z, t; \theta) \triangleq \begin{bmatrix} H(\theta)x + g(\theta) + C(\theta)^\top z \\ C(\theta)x - u(\theta) - t \\ [[t]_- + z]_+ - z \\ C(\theta)^\top [t]_+ \end{bmatrix} \quad (\text{G})$$

Lemma

It holds that (x^, z^*, s^*) solves QP-H(θ) iff there exists $t^* \in \mathbb{R}^{n_i}$ s.t. $G(x^*, z^*, t^*; \theta) = 0$ and $s^* = [t^*]_+$.*

Differentiating (QP-H(θ)) solutions: path differentiability of G

Lemma

G is path differentiable w.r.t. x^* , z^* and t^* . Furthermore, if $H(\theta)$, $g(\theta)$, $C(\theta)$ and $u(\theta)$ are differentiable w.r.t. θ , then G is path differentiable w.r.t. θ .

It means G has a **conservative**¹ Jacobian. It is equivalent to having a chain rule for the Clarke subdifferential

$$\text{Jac}^c G(w) \triangleq \text{conv} \left\{ \lim_{k \rightarrow \infty} \text{Jac} G(w_k) \in \text{diff}_F, w_k \rightarrow w \right\},$$

with diff_G the full measure set where G is differentiable, and $\text{Jac} G$ the standard Jacobian of G .

¹Bolte, J., Le, T., Pauwels, E., Silveti-Falls, T. (2021). Nonsmooth implicit differentiation for machine-learning and optimization.

Differentiating (QP-H(θ)) solutions: Extended Conservative Jacobian

Let $v^* = (x^*, z^*, t^*) \in \mathbb{R}^n \times \mathbb{R}_+^{n_i} \times \mathbb{R}^{n_i}$ s.t. $G(x^*, z^*, t^*; \theta) = 0$.

$$\left(\frac{\partial x^*}{\partial \theta}, \frac{\partial z^*}{\partial \theta}, \frac{\partial t^*}{\partial \theta} \right) \in \arg \min_w \left\| \frac{\partial G(x^*, z^*, t^*; \theta)}{\partial v^*} w + \frac{\partial G(x^*, z^*, t^*; \theta)}{\partial \theta} \right\|_2^2.$$

$$\Pi \frac{\partial t^*}{\partial \theta} \in \frac{\partial s^*}{\partial \theta}, \text{ with } \Pi \in \partial([\cdot]_+)(t^*).$$

Differentiating (QP-H(θ)) solutions: Extended Conservative Jacobian

Why do we use a least-square ?

- It makes sense in "standard" cases:
 - If (QP-H(θ)) reduces to (QP(θ)) (i.e., feasible problem), and if the IFT applies, we recover standard Jacobian.

Differentiating (QP-H(θ)) solutions: Extended Conservative Jacobian

Why do we use a least-square ?

- It makes sense in "standard" cases:
 - If (QP-H(θ)) reduces to (QP(θ)) (i.e., feasible problem), and if the IFT applies, we recover standard Jacobian.
 - If (QP-H(θ)) is infeasible and satisfies non-singularity conditions for a "conservative" IFT, we recover conservative Jacobians

Lemma

Let $C(\theta)$, $u(\theta)$ be differentiable w.r.t. θ , and $H = 0$, and g be fixed w.r.t. θ and satisfying Assumption 1. If $s^ > 0$ and $z^* = 0$ (i.e., it does not satisfy strict complementarity) and C is full row rank, then the ECJs of x^* , z^* , t^* , and s^* correspond to conservative Jacobians.*

Differentiating (QP-H(θ)) solutions: Extended Conservative Jacobian

Why do we use a least-square ?

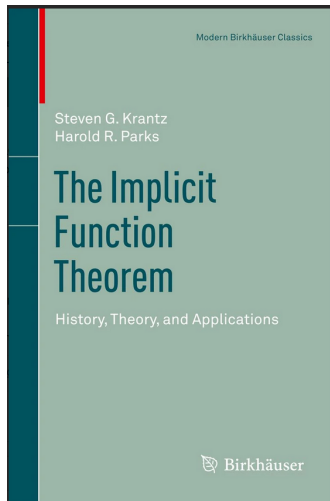
- It makes sense in "standard" cases:
 - If (QP-H(θ)) reduces to (QP(θ)) (i.e., feasible problem), and if the IFT applies, we recover standard Jacobian.
 - If (QP-H(θ)) is infeasible and satisfies non-singularity conditions for a "conservative" IFT, we recover conservative Jacobians
- In non standard cases (i.e., IFT does not apply):
 - Solutions may still have a notion of derivatives

Differentiating (QP-H(θ)) solutions: Extended Conservative Jacobian

Why do we use a least-square ?

- It makes sense in "standard" cases:
 - If (QP-H(θ)) reduces to (QP(θ)) (i.e., feasible problem), and if the IFT applies, we recover standard Jacobian.
 - If (QP-H(θ)) is infeasible and satisfies non-singularity conditions for a "conservative" IFT, we recover conservative Jacobians
- In non standard cases (i.e., IFT does not apply):
 - Solutions may still have a notion of derivatives
 - If not, it is a form of smoothing, hopefully, the "ECJ" is informative.

Differentiating (QP-H(θ)) solutions: other examples



Differentiating (QP-H(θ)) solutions: a toy example

Consider the feasible LP parameterized by $\theta \in (0, 1)$

$$x^*(\theta) \in \arg \min_{x_1, x_2 \in \mathbb{R}^2} x_1 + x_2$$

$$\text{s.t. } \theta \leq x_1 + x_2,$$

$$0 \leq x_1 \leq 1,$$

$$0 \leq x_2 \leq 1.$$

Differentiating (QP-H(θ)) solutions: a toy example

Some differential calculus:

$$\frac{\partial H}{\partial \theta} = 0, \quad \frac{\partial \mathbf{g}}{\partial \theta} = 0, \quad \frac{\partial \mathbf{C}}{\partial \theta} = 0, \quad \frac{\partial \mathbf{u}}{\partial \theta} = (-1 \ 0 \ 0 \ 0 \ 0)^\top.$$

Differentiating (QP-H(θ)) solutions: a toy example

Some differential calculus:

$$\frac{\partial H}{\partial \theta} = 0, \quad \frac{\partial \mathbf{g}}{\partial \theta} = 0, \quad \frac{\partial \mathbf{C}}{\partial \theta} = 0, \quad \frac{\partial \mathbf{u}}{\partial \theta} = (-1 \ 0 \ 0 \ 0 \ 0)^\top.$$

$$\begin{bmatrix} H & \mathbf{C}^\top & 0 \\ \mathbf{C} & 0 & -I \\ 0 & \Pi_1 - I & \Pi_1 \Pi_2 \\ 0 & 0 & \mathbf{C}^\top \Pi_3 \end{bmatrix} \in \frac{\partial G(x^*, z^*, t^*; \theta)}{\partial \mathbf{v}^*},$$

for some $\Pi_1 \in \partial[\cdot]_+([t^*]_- + z^*)$, $\Pi_2 \in \partial[\cdot]_-(t^*)$ and $\Pi_3 \in \partial[\cdot]_+(t^*)$.

Differentiating (QP-H(θ)) solutions: a toy example

Some differential calculus:

$$\frac{\partial H}{\partial \theta} = 0, \quad \frac{\partial g}{\partial \theta} = 0, \quad \frac{\partial C}{\partial \theta} = 0, \quad \frac{\partial u}{\partial \theta} = (-1 \ 0 \ 0 \ 0 \ 0)^T.$$

$$\begin{bmatrix} H & C^T & 0 \\ C & 0 & -I \\ 0 & \Pi_1 - I & \Pi_1 \Pi_2 \\ 0 & 0 & C^T \Pi_3 \end{bmatrix} \in \frac{\partial G(x^*, z^*, t^*; \theta)}{\partial v^*},$$

for some $\Pi_1 \in \partial[\cdot]_+([t^*]_- + z^*)$, $\Pi_2 \in \partial[\cdot]_-(t^*)$ and $\Pi_3 \in \partial[\cdot]_+(t^*)$.

Feasible problem: $\Pi_2 = I$, $\Pi_3 = 0$

Differentiating (QP-H(θ)) solutions: a toy example

Some differential calculus:

$$\frac{\partial H}{\partial \theta} = 0, \quad \frac{\partial g}{\partial \theta} = 0, \quad \frac{\partial C}{\partial \theta} = 0, \quad \frac{\partial u}{\partial \theta} = (-1 \ 0 \ 0 \ 0 \ 0)^\top.$$

$$\begin{bmatrix} H & C^\top & 0 \\ C & 0 & -I \\ 0 & \Pi_1 - I & \Pi_1 \Pi_2 \\ 0 & 0 & C^\top \Pi_3 \end{bmatrix} \in \frac{\partial G(x^*, z^*, t^*; \theta)}{\partial v^*},$$

for some $\Pi_1 \in \partial[\cdot]_+([t^*]_- + z^*)$, $\Pi_2 \in \partial[\cdot]_-(t^*)$ and $\Pi_3 \in \partial[\cdot]_+(t^*)$.

Feasible problem: $\Pi_2 = I$, $\Pi_3 = 0$

Strict complementarity+unique active constraint:

$$\Pi_1 = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 0 & 0 \\ 0 & 0 & 0 \end{bmatrix}.$$

Differentiating (QP-H(θ)) solutions: a toy example

The problem simplifies to

$$\frac{\partial x^*}{\partial \theta}, \frac{\partial z^*}{\partial \theta} \in \arg \min_{b_x, b_z} \left\| \begin{bmatrix} 0 & 0 & -1 \\ 0 & 0 & -1 \\ -1 & -1 & 0 \end{bmatrix} \begin{bmatrix} (b_x)_1 \\ (b_x)_2 \\ b_z \end{bmatrix} + \begin{bmatrix} 0 \\ 0 \\ 1 \end{bmatrix} \right\|_2^2$$

with $\frac{\partial t^*}{\partial \theta} = 0$.

The IFT does not apply (degenerate constraints)!

Yet, the least square provides solutions given by the equations:

$$\begin{aligned} (b_x)_1 + (b_x)_2 &= \frac{1}{2}, \\ b_z &\in \mathbb{R}. \end{aligned}$$

Differentiating (QP-H(θ)) solutions: Backward AD algorithms (generic case)

$$\begin{aligned}\frac{\partial \mathcal{L}}{\partial \theta} &= (b_1^*)^\top \frac{\partial H}{\partial \theta} x^* + (b_1^*)^\top \frac{\partial g}{\partial \theta} + (b_2^*)^\top \frac{\partial C}{\partial \theta} x^* \\ &\quad + (z^*)^\top \frac{\partial C}{\partial \theta} b_1^* + (s^*)^\top \frac{\partial C}{\partial \theta} b_4^* - (b_2^*)^\top \frac{\partial u}{\partial \theta},\end{aligned}$$

with b_1^* , b_2^* , b_3^* and b_4^* with solutions of

$$\begin{bmatrix} H & C^\top & 0 & 0 \\ C & 0 & (I - \Pi_1) & 0 \\ 0 & -I & -\Pi_1 \Pi_2 & (1 - \Pi_2)C \end{bmatrix} \begin{bmatrix} b_1^* \\ b_2^* \\ b_3^* \\ b_4^* \end{bmatrix} = - \begin{bmatrix} \frac{\delta \mathcal{L}}{\delta x^*} \\ \frac{\delta \mathcal{L}}{\delta z^*} \\ \frac{\delta \mathcal{L}}{\delta s^*} \end{bmatrix}.$$

Differentiating (QP-H(θ)) solutions: Backward AD algorithms (feasible case)

$$\begin{aligned}\frac{\partial \mathcal{L}}{\partial \theta} &= (b_x^*)^\top \frac{\partial H}{\partial \theta} x^* + (b_x^*)^\top \frac{\partial g}{\partial \theta} + (\Pi_1 b_z^*)^\top \frac{\partial C}{\partial \theta} x^* \\ &\quad + (z^*)^\top \frac{\partial C}{\partial \theta} b_x^* - (\Pi_1 b_z^*)^\top \frac{\partial u}{\partial \theta},\end{aligned}$$

with b_x^* , b_z^* , the solution of the following linear system

$$\begin{bmatrix} H & C^\top \Pi_1 \\ C & -(I - \Pi_1) \end{bmatrix} \begin{bmatrix} b_x \\ b_z \end{bmatrix} = - \begin{bmatrix} \frac{\delta \mathcal{L}}{\delta x^*} \\ \frac{\delta \mathcal{L}}{\delta z^*} \end{bmatrix}.$$

Differentiating (QP-H(θ)) solutions: toy example (continued!)

LP parameterized by $\theta \in (0, 1)$:

$$x^*(\theta) \in \arg \min_{x_1, x_2 \in \mathbb{R}^2} x_1 + x_2$$

$$\text{s.t. } \theta \leq x_1 + x_2,$$

$$0 \leq x_1 \leq 1,$$

$$0 \leq x_2 \leq 1.$$

We choose: $x_1^* = x_2^* = \frac{\theta}{2}$.

The problem simplifies to

$$\frac{\partial x^*}{\partial \theta}, \frac{\partial z^*}{\partial \theta} \in \arg \min_{b_x, b_z} \left\| \begin{bmatrix} 0 & 0 & -1 \\ 0 & 0 & -1 \\ -1 & -1 & 0 \end{bmatrix} \begin{bmatrix} (b_x)_1 \\ (b_x)_2 \\ b_z \end{bmatrix} + \begin{bmatrix} 0 \\ 0 \\ 1 \end{bmatrix} \right\|^2$$

with $\frac{\partial t^*}{\partial \theta} = 0$.

The IFT does not apply (degenerate constraints)!

Yet, the least square provides solutions given by the equations:

$$(b_x)_1 + (b_x)_2 = \frac{1}{2},$$

$$b_z \in \mathbb{R}.$$

Differentiating (QP-H(θ)) solutions: toy example (continued!)

Let's minimize $\mathcal{L}(\theta) = x_1^*(\theta)$.

Differentiating (QP-H(θ)) solutions: toy example (continued!)

Let's minimize $\mathcal{L}(\theta) = x_1^*(\theta)$. If we try applying the IFT we get

$$\begin{bmatrix} 0 & 0 & -1 \\ 0 & 0 & -1 \\ -1 & -1 & 0 \end{bmatrix} \begin{bmatrix} (b_x)_1 \\ (b_x)_2 \\ b_z \end{bmatrix} = - \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix},$$

which is infeasible!

Differentiating (QP-H(θ)) solutions: toy example (continued!)

Let's minimize $\mathcal{L}(\theta) = x_1^*(\theta)$. If we try applying the IFT we get

$$\begin{bmatrix} 0 & 0 & -1 \\ 0 & 0 & -1 \\ -1 & -1 & 0 \end{bmatrix} \begin{bmatrix} (b_x)_1 \\ (b_x)_2 \\ b_z \end{bmatrix} = - \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix},$$

which is infeasible! The least-square solution, yet provides

$$\begin{bmatrix} 1 & 1 & 0 \\ 1 & 1 & 0 \\ 0 & 0 & 2 \end{bmatrix} \begin{bmatrix} (b_x^*)_1 \\ (b_x^*)_2 \\ b_z^* \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \\ 1 \end{bmatrix},$$

which outputs as ECJ for

$\frac{\partial \mathcal{L}}{\partial \theta} = b_z^* = \frac{1}{2}$. It is coherent with

$\nabla_{\theta}(\theta/2) = 1/2$, if we choose

$x_1^*(\theta) = \theta/2$.

Differentiating (QP-H(θ)) solutions: toy example (continued!)

Let's minimize $\mathcal{L}(\theta) = x_1^*(\theta)$. If we try applying the IFT we get

$$\begin{bmatrix} 0 & 0 & -1 \\ 0 & 0 & -1 \\ -1 & -1 & 0 \end{bmatrix} \begin{bmatrix} (b_x)_1 \\ (b_x)_2 \\ b_z \end{bmatrix} = - \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix},$$

which is infeasible! The least-square solution, yet provides

$$\begin{bmatrix} 1 & 1 & 0 \\ 1 & 1 & 0 \\ 0 & 0 & 2 \end{bmatrix} \begin{bmatrix} (b_x^*)_1 \\ (b_x^*)_2 \\ b_z^* \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \\ 1 \end{bmatrix},$$

which outputs as ECJ for

$\frac{\partial \mathcal{L}}{\partial \theta} = b_z^* = \frac{1}{2}$. It is coherent with $\nabla_{\theta}(\theta/2) = 1/2$, if we choose $x_1^*(\theta) = \theta/2$.

The problem simplifies to

$$\frac{\partial x^*}{\partial \theta}, \frac{\partial z^*}{\partial \theta} \in \arg \min_{b_x, b_z} \left\| \begin{bmatrix} 0 & 0 & -1 \\ 0 & 0 & -1 \\ -1 & -1 & 0 \end{bmatrix} \begin{bmatrix} (b_x)_1 \\ (b_x)_2 \\ b_z \end{bmatrix} + \begin{bmatrix} 0 \\ 0 \\ 1 \end{bmatrix} \right\|_2^2$$

with $\frac{\partial t^*}{\partial \theta} = 0$.

The IFT does not apply (degenerate constraints)!

Yet, the least square provides solutions given by the equations:

$$\begin{aligned} (b_x)_1 + (b_x)_2 &= \frac{1}{2}, \\ b_z &\in \mathbb{R}. \end{aligned}$$

Differentiating (QP-H(θ)) solutions: toy example (continued!)

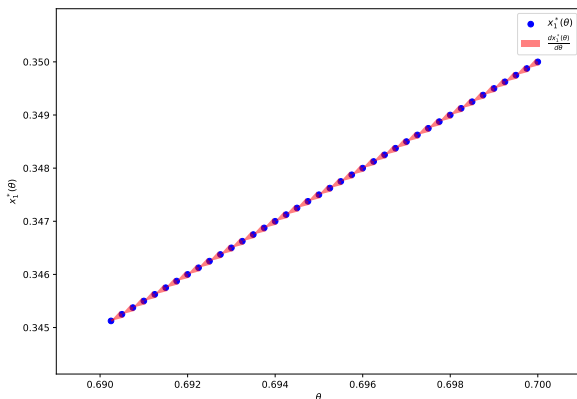


Figure: 40 iterations of gradient descent for minimizing $x_1^*(\theta)$ using ECJs.

Table of Contents

- 1 Introduction
- 2 Solution outline
- 3 Solving the closest feasible problem
- 4 Differentiating closest feasible problem solutions
- 5 Numerical results**
- 6 Conclusion

Back to the Sudoku problem: setup

Three architectures compared:

Back to the Sudoku problem: setup

Three architectures compared:

- structurally feasible with strictly convex loss (standard model of OptNet);

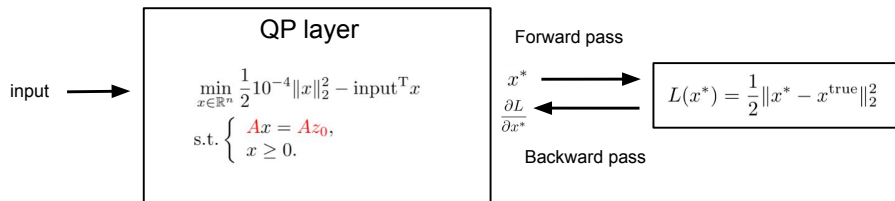


Figure: OptNet layer (structurally feasible at training time).

Back to the Sudoku problem: setup

Three architectures compared:

- structurally feasible with strictly convex loss (standard of OptNet);
- training infeasible LP towards feasibility (with QPLayer, ours).

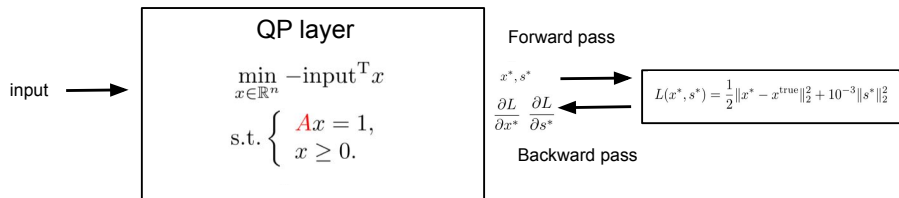


Figure: QPLayer training an infeasible LP.

Back to the Sudoku problem: setup

Three architectures compared:

- structurally feasible with strictly convex loss (standard model of OptNet);
- training infeasible LP towards feasibility (with QPLayer, ours).
- Reformulation of $(QP-H(\theta))$ as a convex QP (non standard)

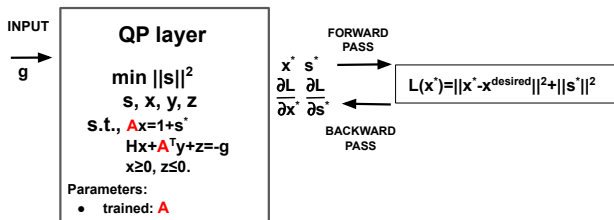


Figure: OptNet layer used for trying to learn A via a reformulation of $(QP-H(\theta))$.

Back to the sudoku: results

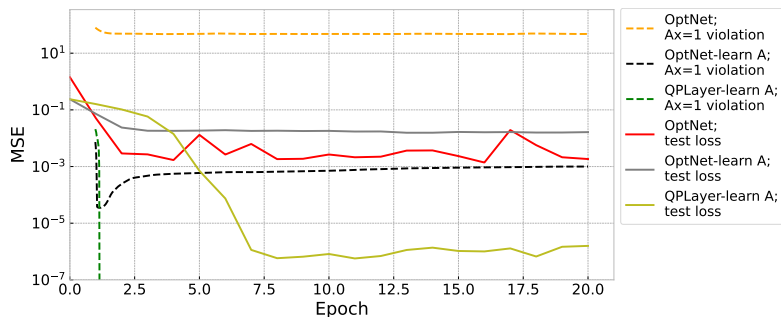


Figure: Test MSE loss of QPLeayr, OptNet, QPLeayr-learn A, and OptNet-learn A specialized for learning A. It includes Sudoku $Ax = 1$ violation.

Back to the sudoku: results

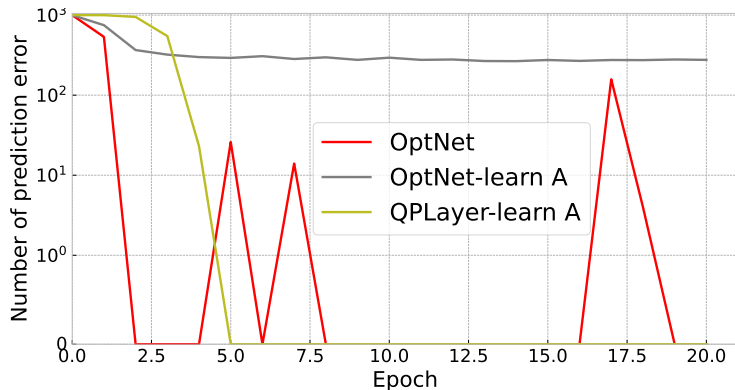


Figure: Test prediction errors over 1000 puzzles of OptNet, QPlayer, QPlayer-learn A and OptNet-learn A specialized for learning A.

Table of Contents

- 1 Introduction
- 2 Solution outline
- 3 Solving the closest feasible problem
- 4 Differentiating closest feasible problem solutions
- 5 Numerical results
- 6 Conclusion**

We introduced:

- QPLayer: framework for learning new types of QP layers.
- Practical concept of "Extended" Conservative Jacobian.
- Practical algorithms making use of powerful AL properties.

Appendix: Augmented Lagrangian method

The Augmented Lagrangian

$$\mathcal{L}_A(x, z; \mu) \triangleq f(x) + \frac{1}{2\mu} (\| [Cx - u + \mu z]_+ \|_2^2 - \|\mu z\|_2^2)$$

The AL method

$$x^{k+1} \approx_{\epsilon^k} \arg \min_x \mathcal{L}_A(x, z; \mu)$$
$$z^{k+1} = \left[\frac{1}{\mu} (Cx^{k+1} - u) + z^k \right]_+$$



Magnus Hestenes



Michael J.D. Powell

Appendix: QP layers in machine learning (example 2)

Coupling visual sensing and calibration of control algorithms in robotics.

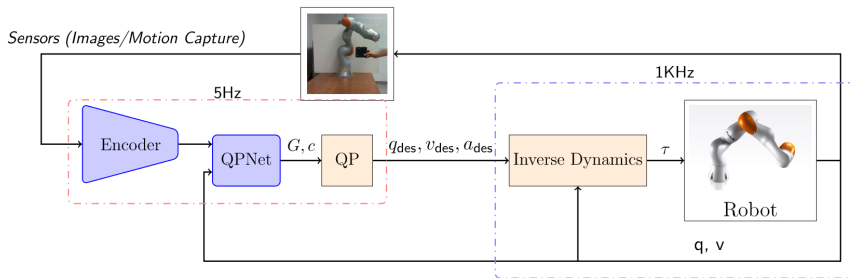


Figure: Pipeline of image measurements for learning QP cost functions (for optimal reaching tasks).

The closest feasible QP problem: solution method

Augmented Lagrangian-based methods have the property of converging towards a solution to $(\text{QP-H}(\theta))$ if $(\text{QP}(\theta))$ is primal infeasible.

Theorem (A. Chiche, JC Gilbert (2016))

Under Assumption (1), the revised AL algorithm does not terminate with a direction of unboundedness and generates a sequence $\{(x^k, z^k)\}$ converging towards a solution to $(\text{QP-H}(\theta))$.

Remark

revised AL algorithm means here two things: (i) the AL algorithm takes into account a modified stopping criterion (closest feasible optimality); (ii) the AL method is exact!

Differentiating (QP-H(θ)) solutions: Backward AD algorithms

Lemma

Let $h : \mathbb{R}^n \times (\mathbb{R}^{n_i})^2 \rightarrow \mathbb{R}$ be a differentiable function, and let $H(\theta)$, $g(\theta)$, $C(\theta)$ and $u(\theta)$ be differentiable w.r.t. θ and satisfying 1. Then, denoting $\mathcal{L}(\theta) \triangleq h(x^*(\theta), z^*(\theta), s^*(\theta))$ and under IFT assumptions, we have that $\frac{\partial \mathcal{L}}{\partial \theta}$ can be derived as follows

$$\frac{\partial \mathcal{L}}{\partial \theta} = (b_1^*)^\top \frac{\partial H}{\partial \theta} x^* + (b_1^*)^\top \frac{\partial g}{\partial \theta} + (b_2^*)^\top \frac{\partial C}{\partial \theta} x^* + (z^*)^\top \frac{\partial C}{\partial \theta} b_1^* + (s^*)^\top \frac{\partial C}{\partial \theta} b_4^* - (b_2^*)^\top \frac{\partial u}{\partial \theta},$$

where b_1^* , b_2^* , b_3^* and b_4^* are the solutions of the linear system

$$\begin{bmatrix} H & C^\top & 0 & 0 \\ C & 0 & (I - \Pi_1) & 0 \\ 0 & -I & -\Pi_1 \Pi_2 & (1 - \Pi_2)C \end{bmatrix} \begin{bmatrix} b_1^* \\ b_2^* \\ b_3^* \\ b_4^* \end{bmatrix} = - \begin{bmatrix} \frac{\delta \mathcal{L}}{\delta x^*} \\ \frac{\delta \mathcal{L}}{\delta z^*} \\ \frac{\delta \mathcal{L}}{\delta s^*} \end{bmatrix}.$$

Differentiating (QP-H(θ)) solutions: Backward AD algorithms (feasible case)

Lemma

Let $h : \mathbb{R}^n \times (\mathbb{R}^{n_i}) \rightarrow \mathbb{R}$ be a differentiable function, and let $H(\theta)$, $g(\theta)$, $C(\theta)$ and $u(\theta)$ be differentiable w.r.t. θ and satisfying Assumption 1. Then, denoting $\mathcal{L}(\theta) \triangleq h(x^*(\theta), z^*(\theta))$, we have under IFT assumptions that $\frac{\partial \mathcal{L}}{\partial \theta}$ can be derived as follows

$$\frac{\partial \mathcal{L}}{\partial \theta} = (b_x^*)^\top \frac{\partial H}{\partial \theta} x^* + (b_x^*)^\top \frac{\partial g}{\partial \theta} + (\Pi_1 b_z^*)^\top \frac{\partial C}{\partial \theta} x^* + (z^*)^\top \frac{\partial C}{\partial \theta} b_x^* - (\Pi_1 b_z^*)^\top \frac{\partial u}{\partial \theta},$$

with b_x^* , b_z^* , the solution of the following linear system

$$\begin{bmatrix} H & C^\top \Pi_1 \\ C & -(I - \Pi_1) \end{bmatrix} \begin{bmatrix} b_x \\ b_z \end{bmatrix} = - \begin{bmatrix} \frac{\delta \mathcal{L}}{\delta x^*} \\ \frac{\delta \mathcal{L}}{\delta z^*} \end{bmatrix}.$$