# Significant Generalization of the Convergence Proof for the Direct Transcription Method for Constrained Optimal Control Problems

Martin P. Neuenhofen

July 23, 2018

### Abstract

In the arXiv paper [arXiv:1712.07761] from December 2017 we presented a convergent direct transcription method for optimal control problems. In the present paper we present a significantly generalized convergence theory in succinct form. Therein, we replace strong assumptions that we had formerly made on local uniqueness of the solution, and on differentiability of a particular functional. These assumptions are removed now.

# Contents

1

# 1 Introduction

## 1.1 Problem Statement

In Optimal Control or Dynamic Optimization one seeks properties that are not constant, but that change dynamically over a time-interval

$$\Omega := (t_0, t_E) \subset \mathbb{R}.$$

For the scope of this work, we seek two functions on this interval:

$$y : \Omega \to \mathbb{R}^{n_y}, t \mapsto y(t)$$
$$z : \Omega \to \mathbb{R}^{n_z}, t \mapsto z(t)$$

For a compact writing we also use $x := (y, z)$ and $n_x := n_y + n_z \in \mathbb{N}$.

The component $y$ is continuous, while $z$ may have discontinuities. According to this, we define the solution space for $x$ as

$$\mathcal{X} := \left( H^1(\Omega) \right)^{n_y} \times \left( L^2(\Omega) \right)^{n_z}$$

with the scalar product

$$\langle (y, z), (v, w) \rangle_{\mathcal{X}} := \sum_{j=1}^{n_y} \langle y_{[j]}, v_{[j]} \rangle_{H^1(\Omega)} + \sum_{j=1}^{n_z} \langle z_{[j]}, w_{[j]} \rangle_{L^2(\Omega)}$$

and induced norm $\| \cdot \|_{\mathcal{X}}$. $\mathcal{X}$ is a Hilbert space.

Dynamic optimization problems impose differential-algebraic constraints and point-constraints on $x$. Therefore we use the notation

$$\dot{y} := \frac{\mathrm{d}y}{\mathrm{d}t}$$

and introduce $M \in \mathbb{N}$ points

$$t^{(k)} \in \overline{\Omega}, \quad \forall\, k \in \{1, 2, ..., M\}.$$

We can then write the *optimal control problem* in general form as

$$\left.\begin{cases} \min_{x \in \mathcal{X}} & F(x) \\ \text{subject to} & b\big( y(t^{(1)}), y(t^{(2)}), ..., y(t^{(M)}) \big) = 0, \\ & c\big( \dot{y}(t), y(t), z(t), t \big) = 0, \quad \text{f.a.e. } t \in \Omega, \\ & z(t) \geq 0, \quad \text{f.a.e. } t \in \Omega. \end{cases}\right\} \text{(OCP)}$$

with a global minimizer $x^\star$. The abbreviation "f.a.e." means "for almost every", and is meant in the Lebesgue sense, since $z$ is only defined in the Lebesque sense (i.e., not pointwise). In contrast to that, pointwise values of $y$ over $\overline{\Omega}$ are well-defined because $H^1(\Omega) \hookrightarrow \mathcal{C}^0(\overline{\Omega})$. $\mathcal{C}^0(\overline{\Omega})$ is the space of continuous functions over $\overline{\Omega}$.

In (OCP), there appear functions

$$f : \mathbb{R}^{n_y} \times \mathbb{R}^{n_y} \times \mathbb{R}^{n_z} \times \Omega \to \mathbb{R},$$
$$c : \mathbb{R}^{n_y} \times \mathbb{R}^{n_y} \times \mathbb{R}^{n_z} \times \Omega \to \mathbb{R}^m,$$
$$b : \underbrace{\mathbb{R}^{n_y} \times \mathbb{R}^{n_y} \times ... \times \mathbb{R}^{n_y}}_{M \text{ times}} \to \mathbb{R}^p,$$

and the functional

$$F(x) := \int_\Omega f\big(\dot{y}(t), y(t), z(t), t\big)\, \mathrm{d}t.$$

## 1.2  Assumptions

We make a couple of assumptions on (OCP):

1. **Feasibility:** $x^\star$ exists.

2. **Boundedness:** $|f|$, $\|c\|_\infty$, $\|b\|_\infty$ are globally bounded.

3. **L-continuity:** $f$, $c$, $b$ are globally Lipschitz-continuous in all arguments except $t$.

**Remarks**

- $x^\star$ does not need to be unique.

- We do not require differentiability of $f$, $c$, $b$.

- Assumptions 2 and 3 can be forced by simply bounding the outputs of $f$, $c$, $b$ as well as their sensitivities.

## 1.3  Numerical Goal

We want to find a numerical solution $x_h^\star \in \mathcal{X}$ that solves (OCP) in a *tolerance-accurate* sense. To describe what we mean by that, we define the following measure of feasibility for elements $x \in \mathcal{X}$:

$$r(x) := \int_\Omega \|c\big(\dot{y}(t), y(t), z(t), t\big)\|_2^2\, \mathrm{d}t + \|b\big(y(t^{(1)}), y(t^{(2)}), ..., y(t^{(M)})\big)\|_2^2$$

We define the optimality gap and feasibility residual of $x_h^\star$:

$$g_{\mathsf{opt}} := \max\{\, F(x_h^\star) - F(x^\star)\,,\, 0\,\},$$
$$r_{\mathsf{feas}} := r(x_h^\star).$$

The goal for our numerical method is the construction of $x_h^\star$ such that this gap and this residual are driven below an arbitrary prescribed positive tolerance. This is what we mean with a tolerance-accurate solution.

Beyond that, the numerical solution satisfies

$$z_h^\star(t) > 0 \qquad \forall\, t \in \overline{\Omega}\,.$$

Notice that due to our assumptions, both $F$ and $r$ are bounded. Also, both $F$ and $r$ are Lipschitz-continuous with respect to the norm $\|\cdot\|_{\mathcal{X}}$. We bound their Lipschitz-constants with $L_F, L_r \geq 1$.

## 1.4 Outline

We review our numerical method from December 2017, that shows a way for computationally constructing the numerical solution $x_h^\star$. We also review —here in an improved form— the full proof of convergence for this method.

In Section 2 we introduce a reformulation of (OCP) into an unconstrained problem. In Section 3 we apply a Finite Element Method to compute a numerical minimizer of this unconstrained problem.

Eventually, we show that the numerical unconstrained minimizer is a tolerance-accurate solution to (OCP) in the aforementioned way.

# 2 Reformulation

The reformulation into an unconstrained problem works in two steps. We first introduce quadratic penalties and then add logarithmic barriers.

## 2.1 Penalty Form

We introduce a penalty parameter $\omega \in (0, 0.5]$. We define

$$F_\omega(x) := F(x) + \frac{1}{2 \cdot \omega} \cdot r(x)\,.$$

Notice that $F_\omega$ is Lipschitz-continuous with constant

$$L_\omega := L_F + \frac{1}{2 \cdot \omega} \cdot L_r\,.$$

Using $F_\omega$, we introduce the *penalty problem*

$$\left\{ \begin{aligned} &\min_{x \in \mathcal{X}} & & F_\omega(x) \\ &\text{subject to} & & z(t) \geq 0 \qquad \text{f.a.e. } t \in \Omega\,. \end{aligned} \right\} \tag{PP}$$

with a global minimizer $x_\omega^\star$.

The following theorem shows that $\varepsilon$-optimal solutions of (PP) solve (OCP) in a tolerance-accurate way.

**Theorem 2.1** (Penalty Solution). *Let $\varepsilon \geq 0$. Consider an $\varepsilon$-optimal solution $x_\omega^\varepsilon$ to* (PP)*, i.e.*

$$F_\omega(x_\omega^\varepsilon) \leq F_\omega(x_\omega^\star) + \varepsilon \quad \text{and} \quad z_\omega^\varepsilon(t) \geq 0 \quad \text{f.a.e. } t \in \Omega \,.$$

*Define*

$$C_r := 2 \cdot \sup_{x \in \mathcal{X}} |F(x)| \,,$$

*which is bounded due to boundedness of $|f|$. Then it holds:*

$$F(x_\omega^\varepsilon) \leq F(x^\star) + \varepsilon \,,$$
$$r(x_\omega^\varepsilon) \leq 2 \cdot \omega \cdot (C_r + \varepsilon) \,.$$

<u>Proof:</u> $x^\star, x_\omega^\star, x_\omega^\varepsilon$ are all feasible for (PP), but $x_\omega^\star$ is optimal and $x_\omega^\varepsilon$ is $\varepsilon$-optimal. Thus

$$\underbrace{F_\omega(x_\omega^\varepsilon)}_{\geq F(x_\omega^\varepsilon)} \quad \leq \quad \underbrace{F_\omega(x^\star)}_{=F(x^\star)} \quad +\varepsilon \,.$$

From this follows

$$F(x_\omega^\varepsilon) \leq F(x^\star) + \varepsilon$$

because $r(x_\omega^\varepsilon) \geq 0$ and $r(x^\star) = 0$. Besides, it also follows

$$r(x_\omega^\varepsilon) \leq 2 \cdot \omega \cdot \big( \underbrace{|F(x_\omega^\star)| + |F(x^\star)|}_{\leq C_r} + \varepsilon \big) \,.$$

$\square$

## 2.2 Penalty-Barrier Form

In this subsection we reformulate (PP) once more in order to remove the inequality constraints. We do so, using logarithmic barriers.

We introduce a barrier parameter $\tau \in (0, 0.5]$. We define

$$\Gamma(x) := -\sum_{j=1}^{n_z} \int_\Omega \log \big( z_{[j]}(t) \big) \, \mathrm{d}t$$

$$F_{\omega,\tau}(x) := F_\omega(x) + \tau \cdot \Gamma(x)$$

Using $F_{\omega,\tau}$, we introduce the *penalty-barrier problem*

$$\left\{ \quad \min_{x \in \mathcal{X}} \quad F_{\omega,\tau}(x) \quad \right\} \tag{PBP}$$

with a global minimizer $x_{\omega,\tau}^\star$.

We have introduced the logarithmic barriers to keep $z_{\omega,\tau}^\star$ strictly positive. But we know that $L^2(\Omega)$ contains functions that are pointwise unbounded in terms of poles or other singularities. So one may ask whether these logarithmic barriers will actually hinder the components of $z_{\omega,\tau}^\star$ from becoming negative. This question motivates the following theorem.

5

**Theorem 2.2** (Strict Interiorness)**.**

$$z^\star_{\omega,\tau}(t) \geq \frac{\tau}{L_\omega} \qquad \text{f.a.e. } t \in \Omega \,.$$

<u>Proof:</u> We consider a worst-case example, that forces $z$ closest possible to 0 in every component and almost every time-value.

Since $F_\omega$ is Lipschitz-continuous, this worst-case example would be obtained when

$$F_\omega(x) := L_\omega \cdot \|z\|_{L^1(\Omega)} \,.$$

Consequently,

$$F_{\omega,\tau}(x) = L_\omega \cdot \|z\|_{L^1(\Omega)} + \tau \cdot \Gamma(x) \,.$$

The analytic minimizer of $F_{\omega,\tau}$ is

$$z_{[j]}(t) = \frac{\tau}{L_\omega} \qquad \text{f.a.e. } t \in \Omega \,, \quad \forall\, j \in \{1, 2, ..., n_y\} \,.$$

Hence, in general $\tau/L_\omega$ is an essential lower bound. $\qquad\square$

The following theorem shows that $x^\star_{\omega,\tau}$ is $\varepsilon$-optimal for (PP).

**Theorem 2.3** (Penalty-Barrier Solution)**.** *Let* $\|z^\star_\omega\|_{L^\infty(\Omega)}, \|z^\star_{\omega,\tau}\|_{L^\infty(\Omega)} = \mathcal{O}(1)$. *Then:*

$$F_\omega(x^\star_{\omega,\tau}) - F_\omega(x^\star_\omega) = \mathcal{O}\big( \tau \cdot |\log(\tau) + \log(\omega)| \big)$$

<u>Proof:</u> We use the following bar-operator:

> For $x \in \mathcal{X}$, we write $\bar{x}$ to denote a modified version of $x$, where $z_{[j]}(t)$ is replaced pointwise with
>
> $$\bar{z}_{[j]}(t) := \max\left\{ z_{[j]}(t) \,, \frac{\tau}{L_\omega} \right\} \qquad \forall\, j \in \{1, 2, ..., n_z\}, \ \forall\, t \in \Omega \,.$$
>
> Notice that $\bar{x} \in \mathcal{X}$.

Notice also that, according to Theorem 2.2, it holds $x^\star_{\omega,\tau} \equiv \bar{x}^\star_{\omega,\tau}$.

We further notice the general algebraic result

$$\left| \tau \cdot \log\left( \frac{\tau}{L_\omega} \right) \right| = \mathcal{O}\big( \tau \cdot |\log(\tau) + \log(\omega)| \big) \,.$$

Using this relation and the bar-operator from above, we find $\forall\, x \in \mathcal{X}$:

$$
\begin{aligned}
|\tau \cdot \Gamma(\bar{x})| &\leq \left| \tau \cdot \sum_{j=1}^{n_z} \int_\Omega \log\big( \bar{z}_{[j]}(t) \big) \, \mathrm{d}t \right| \\
&\leq n_z \cdot |\Omega| \cdot \max_{1 \leq j \leq n_z} \|\tau \cdot \log(\bar{z}_{[j]})\|_{L^\infty(\Omega)} \\
&\leq n_z \cdot |\Omega| \cdot \Big( \underbrace{\mathcal{O}\big( \tau \cdot |\log(\tau) + \log(\omega)| \big)}_{\text{estimate for } \bar{z}_{[j]} < 1} + \underbrace{\mathcal{O}(\tau)}_{\text{estimate for } \bar{z}_{[j]} \geq 1} \Big) \quad (1)
\end{aligned}
$$

In the last line of this bound, we have distinguished two cases: Namely,

$$\left| \log \left( \bar{z}_{[j]}(t) \right) \right|$$

attains its largest value at a $t \in \overline{\Omega}$ where either $\bar{z}_{[j]} < 1$ (case 1) or where $\bar{z}_{[j]} \geq 1$ (case 2). In the first case, we can use the above algebraic result together with Theorem 2.2 and arrive at the left term in the big brackets of (1). In the second case, we simply bound the logarithm using the extremal value of $|\bar{z}_{[j]}(t)| \leq |z_{[j]}(t)| = \mathcal{O}(1)$.

We can use bound (1) to show the proposition:

$$0 \leq F_\omega(x^\star_{\omega,\tau}) - F_\omega(x^\star_\omega) \leq F_\omega(\bar{x}^\star_{\omega,\tau}) - F_\omega(\bar{x}^\star_\omega) + L_\omega \cdot \|x^\star_\omega - \bar{x}^\star_\omega\|_\mathcal{X}$$

We used Lipschitz-continuity of $F_\omega$ and $x^\star_{\omega,\tau} \equiv \bar{x}^\star_{\omega,\tau}$. We bound this further to

$$
\begin{aligned}
F_\omega(x^\star_{\omega,\tau}) - F_\omega(x^\star_\omega) \leq \quad & \underbrace{F_\omega(\bar{x}^\star_{\omega,\tau}) - F_{\omega,\tau}(\bar{x}^\star_{\omega,\tau})}_{= -\tau \cdot \Gamma(\bar{x}^\star_{\omega,\tau})} + F_{\omega,\tau}(\bar{x}^\star_{\omega,\tau}) \\
& - \left( \underbrace{F_\omega(\bar{x}^\star_\omega) - F_{\omega,\tau}(\bar{x}^\star_\omega)}_{= -\tau \cdot \Gamma(\bar{x}^\star_\omega)} + F_{\omega,\tau}(\bar{x}^\star_\omega) \right) \\
& + L_\omega \cdot n_z \cdot |\Omega| \cdot \frac{\tau}{L_\omega} \, .
\end{aligned}
$$

Therein, we just added two zeros. Also, we used the bound

$$\|x^\star_\omega - \bar{x}^\star_\omega\|_\mathcal{X} \leq n_z \cdot |\Omega| \cdot \frac{\tau}{L_\omega} \leq n_z \cdot |\Omega| \cdot \tau \, ,$$

which follows easily from the definition of the bar-operator.

For the terms $|\tau \cdot \Gamma(\bar{x}^\star_{\omega,\tau})|$ and $|\tau \cdot \Gamma(\bar{x}^\star_\omega)|$ we can use the bound (1), hence obtaining:

$$
\begin{aligned}
F_\omega(x^\star_{\omega,\tau}) - F_\omega(x^\star_\omega) \leq \quad & \underbrace{F_{\omega,\tau}(\bar{x}^\star_{\omega,\tau}) - F_{\omega,\tau}(\bar{x}^\star_\omega)}_{\leq 0} \\
& + n_z \cdot |\Omega| \cdot \mathcal{O}\big( \tau \cdot |\log(\tau) + \log(\omega)| \big) \\
& + n_z \cdot |\Omega| \cdot \tau
\end{aligned}
$$

The under-braced term is bounded above by zero because $\bar{x}^\star_{\omega,\tau} \equiv x^\star_{\omega,\tau}$ is a global minimizer of $F_{\omega,\tau}$. $\qquad \square$

From the proof of the theorem follows a bound for $\Gamma$.

**Lemma 2.4** (Bound for $\Gamma$). *Let $x \in \mathcal{X}$ with $\|z\|_{L^\infty(\Omega)} = \mathcal{O}(1)$. Consider $\bar{x}$, according to the above bar-operator. Then:*

$$\left| \tau \cdot \Gamma\left( \frac{x + \bar{x}}{2} \right) \right| = \mathcal{O}\big( \tau \cdot |\log(\tau) + \log(\omega)| \big)$$

<u>Proof:</u> Notice that

$$\tilde{x} := \frac{x + \bar{x}}{2}$$

is so to say a milder version of $\bar{x}$. Namely, $\tilde{z}$ is only pushed half as much into the interior. In particular, it follows

$$\tilde{z}(t) \geq \frac{\tau}{2 \cdot L_\omega} \qquad \text{f.a.e. } t \in \Omega \,.$$

The proposition now follows in the same way as we showed the bound (1). $\quad\square$

**Remark: On Boundedness**   Notice that $\|z_\omega^\star\|_{L^\infty(\Omega)}, \|z_{\omega,\tau}^\star\|_{L^\infty(\Omega)} = \mathcal{O}(1)$ can be forced easily. E.g., the path constraints

$$z_{[1]}(t) \geq 0 \,, \quad z_{[2]}(t) \geq 0 \,, \quad z_{[1]}(t) + z_{[2]}(t) = const$$

lead to

$$\|z_{[j]}\|_{L^\infty(\Omega)} \leq const \qquad \forall\, j \in \{1, 2\} \,.$$

## 3   Finite Element Method

The approach of our method is to solve the unconstrained problem (PBP) computationally in a finite-dimensional subspace of $\mathcal{X}$, using nonlinear optimization methods. The subspace is constructed using the Finite Element method.

In this section we introduce a suitable finite-dimensional space. We then show a stability result. Eventually, we prove convergence of the Finite Element solution for (PBP) and (OCP).

### 3.1   Spaces

We use a mesh parameter $h \in \mathbb{R}_+ \setminus \{0\}$.

The set $\mathcal{T}_h := \{T\}$ is called *triangulation*, consisting of open intervals $T \subset \Omega$. These intervals satisfy:

$$
\begin{array}{lll}
(i) & \text{Disjunction:} & T_1 \cap T_2 = \emptyset \qquad \forall\, T_1 \neq T_2 \in \mathcal{T}_h \\[2mm]
(ii) & \text{Coverage:} & \bigcup_{T \in \mathcal{T}_h} \overline{T} = \overline{\Omega} \\[2mm]
(iii) & \text{Resolution:} & \max_{T \in \mathcal{T}_h} |T| \leq h \\[2mm]
(iv) & \text{Quasi-uniformity:} & \min_{T_1, T_2 \in \mathcal{T}_h} \dfrac{|T_1|}{|T_2|} \geq \sigma \in \mathbb{R}_+ \setminus \{0\}
\end{array}
$$

In this, $\sigma > 0$ is a constant that must not depend on $h$, such that $1/\sigma = \mathcal{O}(1)$.

We write $\mathcal{P}_p(\overline{T})$ for the space of functions that are polynomials of degree $\leq p \in \mathbb{N}_0$ on interval $\overline{T}$. Our Finite Element space is then given conventionally as

$$\mathcal{X}_{h,p} := \left\{ x : \overline{\Omega} \to \mathbb{R}^{n_x} \ \middle|\ x \in \mathcal{P}_p(\overline{T}) \ \forall\, T \in \mathcal{T}_h \text{ and } y \in \mathcal{C}^0(\overline{\Omega}) \right\}.$$

$\mathcal{X}_{h,p}$ is a Hilbert-space. It holds $\mathcal{X}_{h,p} \subset \mathcal{X}$.

## 3.2 Stability

The following theorem shows that two particular Lebesgue-norms are equivalent in this Finite Element space.

**Theorem 3.1** (Norm equivalence)**.** *Let $p \leq 10^3$. Then:*

$$\|x_{[j]}\|_{L^\infty(\Omega)} \leq \frac{p+1}{\sigma \cdot h} \cdot \|x\|_{\mathcal{X}} \qquad \forall\, x \in \mathcal{X}_{h,p}, \quad \forall\, j \in \{1, 2, ..., n_x\}.$$

<u>Proof:</u> Let $T \subset \mathbb{R}$, $T \neq \emptyset$. We find empirically from the optimality conditions of a convex quadratic program, that

$$\min_{0 \neq u \in \mathcal{P}_p(\overline{T})} \frac{\|u\|_{L^2(T)}^2}{\|u\|_{L^\infty(T)}^2} = \frac{|T|^2}{(p+1)^2}$$

holds $\forall\, p \in \{1, 2, ..., 10^3\}$. This implies

$$\|u\|_{L^\infty(T)} \leq \frac{p+1}{|T|} \cdot \|u\|_{L^2(T)} \qquad \forall\, u \in \mathcal{P}_p(\overline{T}). \tag{2}$$

Using the above bound, we can show the proposition:

$$\|x_{[j]}\|_{L^\infty(\Omega)} \leq \max_{T \in \mathcal{T}_h} \|x_{[j]}\|_{L^\infty(T)}$$

$$\overset{(*)}{\leq} \max_{T \in \mathcal{T}_h} \frac{p+1}{|T|} \cdot \|x_{[j]}\|_{L^2(T)}$$

$$\leq \frac{p+1}{\sigma \cdot h} \cdot \|x_{[j]}\|_{L^2(\Omega)} \leq \frac{p+1}{\sigma \cdot h} \cdot \|x\|_{\mathcal{X}}$$

We used the bound (2) for inequality $(*)$. $\qquad\qquad\square$

The following theorem gives a bound on the growth of $F_{\omega,\tau}$ in a neighborhood of $x_{\omega,\tau}^\star$ for elements of $\mathcal{X}_{h,p}$.

**Theorem 3.2** (Lipschitz-continuity)**.** *Define*

$$\delta_{\omega,\tau,h} := \frac{\tau}{2 \cdot L_\omega} \cdot \frac{\sigma \cdot h}{p+1},$$

$$L_{\omega,\tau,h} := L_\omega + n_z \cdot 2 \cdot L_\omega \cdot \frac{p+1}{\sigma \cdot h}.$$

*Consider the spherical neighbourhood*

$$\mathcal{B} := \left\{ x \in \mathcal{X} \; \middle| \; \|x - x^\star_{\omega,\tau}\|_{\mathcal{X}} \leq \delta_{\omega,\tau} \right\}.$$

*The following holds:*

$$|F_{\omega,\tau}(x^A) - F_{\omega,\tau}(x^B)| \leq L_{\omega,\tau,h} \cdot \|x^A - x^B\|_{\mathcal{X}} \qquad \forall\, x^A, x^B \in \mathcal{B} \cap \mathcal{X}_{h,p} \,.$$

<u>Proof:</u> From Theorem 2.2 and $x^A, x^B \in \mathcal{B}$ follows

$$\min_{1 \leq j \leq n_z} \; \min_{\mathcal{Q} \in \{A,B\}} \; \operatorname*{ess\,inf}_{t \in \Omega} \; z^{\mathcal{Q}}_{[j]}(t) \geq \frac{\tau}{L_\omega} - \frac{p+1}{\sigma \cdot h} \cdot \delta_{\omega,\tau,h} = \frac{\tau}{2 \cdot L_\omega} \,. \tag{3}$$

With this bound we can show the proposition:

$$|F_{\omega,\tau}(x^A) - F_{\omega,\tau}(x^B)|$$

$$\leq |F_\omega(x^A) - F_\omega(x^B)| + \tau \cdot \sum_{j=1}^{n_z} \int_\Omega \left| \log\left(z^A_{[j]}(t)\right) - \log\left(z^B_{[j]}(t)\right) \right| \, \mathrm{d}t$$

$$\leq L_\omega \cdot \|x^A - x^B\|_{\mathcal{X}} + \tau \cdot n_z \cdot |\Omega| \cdot \max_{1 \leq j \leq n_z} \; \operatorname*{ess\,sup}_{t \in \Omega} \left| \log\left(z^A_{[j]}(t)\right) - \log\left(z^B_{[j]}(t)\right) \right|$$

The essential supremum term can be bounded with a Lipschitz-result for the logarithm, because we know lower bounds for the arguments of the logarithm from (3). We obtain

$$\max_{1 \leq j \leq n_z} \; \operatorname*{ess\,sup}_{t \in \Omega} \left| \log\left(z^A_{[j]}(t)\right) - \log\left(z^B_{[j]}(t)\right) \right|$$

$$\leq \max_{1 \leq j \leq n_z} \; \frac{1}{\frac{\tau}{2 \cdot L_\omega}} \cdot \|z^A_{[j]} - z^B_{[j]}\|_{L^\infty(\Omega)}$$

$$\leq \frac{2 \cdot L_\omega}{\tau} \cdot \frac{p+1}{\sigma \cdot h} \cdot \|x^A - x^B\|_{\mathcal{X}} \,,$$

wherein the latter inequality is obtained using Theorem 3.1 . $\qquad\square$

## 3.3  Optimality

We state the *discrete penalty-barrier problem.*

$$\left\{ \begin{array}{ll} \displaystyle\min_{x \in \mathcal{X}_{h,p}} & F_{\omega,\tau}(x) \\[2ex] \text{subject to} & z(t) \geq \dfrac{\tau}{2 \cdot L_\omega} \qquad \text{f.a.e. } t \in \Omega \end{array} \right\} \tag{PBP$_h$}$$

Using the space

$$\mathcal{X}^{\omega,\tau}_{h,p} := \left\{ x \in \mathcal{X}_{h,p} \; \middle| \; z(t) \geq \frac{\tau}{2 \cdot L_\omega} \quad \text{f.a.e. } t \in \Omega \right\},$$

10

we can write $(\mathrm{PBP}_h)$ as unconstrained problem

$$\left\{ \min_{x \in \mathcal{X}_{h,p}^{\omega,\tau}} \quad F_{\omega,\tau}(x) \right\},$$

with a global minimizer $x_h^\star$.

We show that the Finite Element solution $x_h^\star$ is $\varepsilon$-optimal for (PBP).

**Theorem 3.3** (Optimality of FEM Solution). *Let $\mathcal{B}$ as in Theorem 3.2. If $\mathcal{B} \cap \mathcal{X}_{h,p} \neq \emptyset$, then $x_h^\star$ satisfies:*

$$F_{\omega,\tau}(x_h^\star) - F_{\omega,\tau}(x_{\omega,\tau}^\star) \leq L_{\omega,\tau,h} \cdot \inf_{x_h \in \mathcal{X}_{h,p}} \left\{ \|x_{\omega,\tau}^\star - x_h\|_{\mathcal{X}} \right\}.$$

Proof: Consider the Finite Element best-approximation

$$\tilde{x}_h := \operatorname*{argmin}_{x_h \in \mathcal{X}_{h,p}} \left\{ \|x_{\omega,\tau}^\star - x_h\|_{\mathcal{X}} \right\},$$

as well-defined by the Hilbert-space $\mathcal{X}_{h,p}$.

Since $\mathcal{B} \cap \mathcal{X}_{h,p} \neq \emptyset$, it follows

$$\exists\, x_h \in \mathcal{X}_{h,p} \;:\; \|x_{\omega,\tau}^\star - x_h\|_{\mathcal{X}} \leq \delta_{\omega,\tau,h}\,.$$

Thus, $\tilde{x}_h \in \mathcal{B} \cap \mathcal{X}_{h,p}$. Hence,

$$\tilde{x}_h \equiv \operatorname*{argmin}_{x_h \in \mathcal{B} \cap \mathcal{X}_{h,p}} \left\{ \|x_{\omega,\tau}^\star - x_h\|_{\mathcal{X}} \right\}.$$

In analogy to (3), we find $\mathcal{B} \cap \mathcal{X}_{h,p} \subset \mathcal{X}_{h,p}^{\omega,\tau}$. Thus, $\tilde{x}_h \in \mathcal{X}_{h,p}^{\omega,\tau} \subset \mathcal{X}_{h,p}$. Hence,

$$\tilde{x}_h \equiv \operatorname*{argmin}_{x_h \in \mathcal{X}_{h,p}^{\omega,\tau}} \left\{ \|x_{\omega,\tau}^\star - x_h\|_{\mathcal{X}} \right\}.$$

Since $x_{\omega,\tau}^\star$ is a global minimizer of $F_{\omega,\tau}$ in $\mathcal{X}$, whereas $x_h^\star$ is only a global minimizer of $F_{\omega,\tau}$ in the subspace $\mathcal{X}_{h,p}^{\omega,\tau} \subset \mathcal{X}$, and whereas $\tilde{x}_h$ just lives in $\mathcal{X}_{h,p}^{\omega,\tau}$, it holds

$$F_{\omega,\tau}(x_{\omega,\tau}^\star) \leq F_{\omega,\tau}(x_h^\star) \leq F_{\omega,\tau}(\tilde{x}_h)\,.$$

For the latter, using Theorem 3.2, we obtain the bound

$$F_{\omega,\tau}(\tilde{x}_h) \leq F_{\omega,\tau}(x_{\omega,\tau}^\star) + L_{\omega,\tau,h} \cdot \|x_{\omega,\tau}^\star - \tilde{x}_h\|_{\mathcal{X}}\,.$$

$\square$

**Remark: Non-emptiness** We show that usually the space $\mathcal{B} \cap \mathcal{X}_{h,p}$ is non-empty. If

$$\underbrace{\inf_{x_h \in \mathcal{X}_{h,p}} \left\{ \|x_{\omega,\tau}^\star - x_h\|_{\mathcal{X}} \right\}}_{(*)} \leq \delta_{\omega,\tau,h} \tag{4}$$

then the minimizing argument of $(*)$ lives in $\mathcal{B} \cap \mathcal{X}_{h,p}$, showing non-emptiness. Also, if we keep $\omega, \tau$ fixed, then

$$\delta_{\omega,\tau,h} \in \Theta(h) \,.$$

Thus, if the infimum term $(*)$ converges superlinearly in $h$, then for $h > 0$ sufficiently small it follows (4).

## 3.4 Convergence

We obtain a bound for the optimality gap and feasibility residual of $x_h^\star$ for (OCP).

**Theorem 3.4** (Convergence for OCP). *Let $\|z_h^\star\|_{L^\infty(\Omega)} = \mathcal{O}(1)$. Then the numerical solution $x_h^\star$ satisfies:*

$$g_{opt} = \mathcal{O}\left( \tau \cdot |\log(\tau) + \log(\omega)| + L_{\omega,\tau,h} \cdot \inf_{x_h \in \mathcal{X}_{h,p}} \left\{ \|x_{\omega,\tau}^\star - x_h\|_{\mathcal{X}} \right\} \right)$$

$$r_{feas} = \mathcal{O}\left( \omega + \omega \cdot L_{\omega,\tau,h} \cdot \underbrace{\inf_{x_h \in \mathcal{X}_{h,p}} \left\{ \|x_{\omega,\tau}^\star - x_h\|_{\mathcal{X}} \right\}}_{=: \varepsilon_{h,p}} \right) \,.$$

<u>Proof:</u> From Theorem 3.3 we know that for $\varepsilon_{h,p} > 0$ it holds:

$$F_{\omega,\tau}(x_h^\star) \leq F_{\omega,\tau}(x_{\omega,\tau}^\star) + \varepsilon_{h,p}$$

This is equivalent to

$$F_\omega(x_h^\star) + \tau \cdot \Gamma(x_h^\star) \leq F_\omega(x_{\omega,\tau}^\star) + \tau \cdot \Gamma(x_{\omega,\tau}^\star) + \varepsilon_{h,p}$$

$$\Rightarrow \qquad F_\omega(x_h^\star) \leq F_\omega(x_{\omega,\tau}^\star) + \underbrace{|\tau \cdot \Gamma(x_h^\star)| + |\tau \cdot \Gamma(x_{\omega,\tau}^\star)|}_{(*)} \,.$$

Since $x_h^\star \in \mathcal{X}_{h,p}^{\omega,\tau}$, it follows $z_h^\star \geq \frac{\tau}{2 \cdot L_\omega}$ and thus

$$x_h^\star = \frac{x_h^\star + \bar{x}_h^\star}{2} \,.$$

Similarly, since $x_{\omega,\tau}^\star \equiv \bar{x}_{\omega,\tau}^\star$, it also follows

$$x_{\omega,\tau}^\star = \frac{x_h^\star + \bar{x}_{\omega,\tau}^\star}{2} \,.$$

12

Therefore, we can apply Lemma 2.4 to bound $(*)$ in $\mathcal{O}(\tau \cdot |\log(\tau) + \log(\omega)|)$. It follows:

$$F_\omega(x_h^\star) \leq F_\omega(x_{\omega,\tau}^\star) + \mathcal{O}(\tau \cdot |\log(\tau) + \log(\omega)|) + \varepsilon_{h,p}$$

Since, according to Theorem 2.3, $x_{\omega,\tau}^\star$ is $\tilde{\varepsilon}$-optimal for (PP), where $\tilde{\varepsilon} = \mathcal{O}(\tau \cdot |\log(\tau) + \log(\omega)|)$, it follows further:

$$F_\omega(x_h^\star) \leq F_\omega(x_\omega^\star) + \underbrace{\mathcal{O}(\tau \cdot |\log(\tau) + \log(\omega)|) + \varepsilon_{h,p}}_{=:\varepsilon}$$

I.e., $x_h^\star$ is $\varepsilon$-optimal for (PP). The proposition now follows from Theorem 2.1.
$\qquad\square$

**Remark: Order of Convergence**  It holds

$$L_{\omega,\tau,h} = \mathcal{O}\left(\frac{1}{\omega \cdot h}\right).$$

If $\tau \leq \omega$ then

$$\tau \cdot |\log(\tau) + \log(\omega)| = \mathcal{O}(\sqrt{\tau}).$$

Assuming that the solution $x_{\omega,\tau}^\star$ is sufficiently smooth, such that

$$\exists\, \ell \in (0,p] \;:\; \inf_{x_h \in \mathcal{X}_{h,p}} \left\{ \|x_{\omega,\tau}^\star - x_h\|_\mathcal{X} \right\} = \mathcal{O}(h^{1+\ell}),$$

and choosing $h \in (0,1]$, $\tau = h^\ell$, and $\omega = h^{\ell/2}$, then

$$g_{\mathtt{opt}} = \mathcal{O}(h^{\ell/2}), \qquad r_{\mathtt{feas}} = \mathcal{O}(h^{\ell/2}).$$

# 4   Conclusions

We summarized the convergence result at a glance. In this, we removed some assumptions.

In the present form, basic assumptions like uniqueness of $x^\star$ and continuity of $f$, $c$, $b$ are no longer required. The convergence of the numerical solution is characterized in the measures of optimality gap and feasibility residual.

Certainly, for the numerical computation of $x_h^\star$ with optimization methods, it is beneficial when $f, c, b$ are smooth. But, for the convergence of the direct transcription scheme, this is not crucial.

**Outlook**   Typically, when solving (PBP$_h$), one uses numerical quadrature for the integrals in $F$, $r$. I.e., $F$, $r$ are replaced with perturbed functionals $F_h$, $r_h$. At this point, we have not looked into the effects that these perturbations have on the convergence of $x_h^\star$.