# Optimization of Static and Dynamic Systems

**Summary** Fabian Damken November 8, 2023



# Contents

1	Einfi	ührung	9
	1.1	Beispiele	9
	1.2	Fragestellungen	9
	1.3	Allgemeine Formulierung eines Optimierungsproblems	9
	1.4	Statische vs. Dynamische Optimierung	9
	1.5	Klassifizierung von Optimierungsverfahren	9
	1.6	Typische Struktur	9
2	0	diant Deced Optimization without Constraints	•••
Z		dient-Based Optimization without Constraints	
	2.1	Solution Characterization	.0
		2.1.1 One-Dimensional Optimization	.0
	<u>.</u>	2.1.2 Multi-Dimensional Optimization	.U
	2.2	Numerical Gradient-Based Methods	. I
		2.2.1 Starting Point	.1
		2.2.2 Steepest Descent	.2
		2.2.3 Conjugate Gradient	.2
		2.2.4 Newton Method	.3
		2.2.5 Quasi-Newton Methods	.5
		2.2.6 Comparison	.6
		2.2.7 Notes and Discussion	.6
	2.3	Step Size Rules, Line Search	.7
		2.3.1 Inexact Line Search	.8
		2.3.2 Notes	.8
	2.4	Trust Region Methods    1	.9
	2.5	Rate of Convergence	.9
		2.5.1 Gradient-Based Methods	20
3	Grad	dient-Free Optimization without Constraints	21
•	3.1		21
		3.1.1 Simulation-Based Optimization	21
		3.1.2 Black-Box Optimization	22
	3.2	Metaheuristics	22
	0.2	3.2.1 Evolutionary Algorithm (EA)	22
		3.2.2 Genetic Algorithms (GA)	22
		3 2 3 Further Metabeuristics	בי. גע
	33	Deterministic Sampling Methods (Pattern Search Methods)	,0 )2
	5.5	3 3 1 Nelder-Mead Simpley Method	,0 )2
		3.3.2 Multidirectional Search Methods	.5 76
		3.3.2 Multurectional Search Methods	20 26
		2.3  Implicit Filtering	10 76
		$5.5.4$ miphen finering $\ldots$ $2$	0

	3.4	Surrogate Optimization	27
		3.4.1 Approximation Methods	27
		3.4.2 Select of the Sampling Points	28
		3.4.3 Minimizing the Surrogate Function	29
		3.4.4 Discussion	30
	3.5	Comparison	30
		3.5.1 Magnetic Bearing Design	30
		3.5.2 Walking Optimization of a Humanoid Robot	30
	3.6	Discussion	30
4	Grad	lient-Based Optimization with Constraints	31
	4.1	Solution Characterization	31
		4.1.1 First-Order Necessary Optimality Conditions (Karush-Kuhn-Tucker Conditions, KKT) .	32
		4.1.2 Second-Order Necessary Optimality Conditions	32
		4.1.3 Example	32
	4.2	Simple Bounds, Box Constraints	32
	4.3	Penalty Function	33
		4.3.1 Exterior Penalty Functions	33
		4.3.2 Interior Penalty Functions	34
		4.3.3 Exact Penalty Functions	34
		4.3.4 Augmented Lagrangian	35
	4.4	Constraint Elimination	35
	4.5	Sequential Quadratic Programming (SQP)	36
		4.5.1 Finding the Search Direction	37
		4.5.2 Step Size Rules	38
		4.5.3 Approximation of the Lagrange Multipliers	38
		4.5.4 Termination Criteria	39
		4.5.5 Hessian Approximation	39
		4.5.6 SQP Method (Algorithm)	42
		4.5.7 Notes	42
		4.5.8 Examples	44
		4.5.9 Wrap-Up	44
5	Calc	culation of Derivatives	45
	5.1	Finite Difference Approximation (numerical Differentiation)	45
		5.1.1 Forward Difference Approximation	45
		5.1.2 Central-Difference Approximation	47
	5.2	Numerical Differentiation of Simulation Models	48
		5.2.1 Derivative of ODE-Simulation Models	48
		5.2.2 External Numerical Differentiation	49
		5.2.3 Internal Numerical Differentiation	50
	5.3	Symbol Differentiation	50
	5.4	Automatic Differentiation	50
6	Para	ameter Estimation	52
	6.1	Objective Functions	52
	6.2	Linear Least Squares	53
		•	

	6.3	Optimality Conditions and Special Methods
		6.3.1 Gauss-Quasi-Newton Method
		6.3.2 Levenberg-Marquardt Methods
		6.3.3 Notes
	6.4	Conditioning of Normal Equations
	6.5	Result Interpretation
		6.5.1 Common Problems
		6.5.2 The Covariance Matrix
	66	Ontimal Experimental Design
	67	Fyamples
	0.7	6.7.1 Deremeter Dependent Vehicle Dynamics
		6.7.2 Deremeter Estimation for "PioPiped"
		0.7.2 Parameter Estimation for Biobiped
7	Mini	imization of Functionals 57
	7.1	Euler-Lagrange Equation
		7.1.1 Example
		7.1.2 Notes
		713 Derivation
8	Opti	mal Control 59
	8.1	Necessary Optimality Conditions for the Basis Problem
		8.1.1 Boundary Conditions
		8.1.2 First-Order Necessary Optimality Conditions (Maximum Principle)
		8.1.3 Second-Order Necessary Optimality Condition (Legendre-Clebsch Condition) 6
		8.1.4 Example
		8.1.5 Application: Optimal Robot Control
	8.2	Bang-Bang Singular Control
	0.2	8 2 1 Singular Control
		8.2.2 Application: Time-Minimal Robot Control
		8.2.3 Notes
	83	Value Function and Hamilton Jacobi Bellman Equation
	0.5	9.2.1 Derivation
	0.4	
	8.4	
		8.4.1 Mixed Inequality Constraints
		8.4.2 State Inequality Constraints
		8.4.3 Examples
		8.4.4 Summary
9	Calc	culating Optimal Trajectories 69
-	9.1	First Computation Methods
		9.1.1 Dynamic Programming
		9.1.2 Gradient Methods (Min-H Methods)
	92	Indirect Methods
	0.2 0.2	Direct Methods
	7.5	0.3.1 Direct Collocation Methods 7/
		7.5.1 Direct Conocation Methods
	0.4	9.3.2   Direct Shooting Methods   /4
	9.4	Notes

10	Optimal Feedback Control	76
	10.1 Classical Feedback Control (Position Control)	76
	10.2 Optimal Feedback Control	76
	10.3 Linear Quadratic Regulator (LQR)	77
	10.3.1 Derivation	77
	10.4 Neighboring Extremals	77
	10.4.1 Indirect Methods	78
	10.4.2 Direct Methods	78
	10.4.3 Nonlinear Model Predictive Control (NMPC)	78
	10.5 Numerical Synthesis of the Nonlinear Feedback Control	79
11	Further Topics on Optimal Control	80
	11.1 Inverse Optimal Control	80
	11.2 Differential/Dynamic Games	80
	11.2.1 Non-Cooperative Two-Player Zero-Sum Differential Games	80
	11.3 Learning Methods and Optimization	81
	11.3.1 Foundations	81
	11.3.2 Reinforcement Learning	81

# List of Figures

10.1 Feedforward Control				 	•									•	•								•	 76
10.2 Feedback Control	•		•	 	•	•	•	•	•	•	•	 •	•	•	•	•	•	•			•	•	•	 77

# **List of Tables**

8.1 Possibilities for active state-only constraints in optimal control given the order of the constraint. 68

# List of Algorithms

1	Algorithmic structure of a gradient-based optimization algorithms.	12
2	Conjugate Gradient for nonlinear Objective Function.	13
3	Newton Method	14
4	Quasi-Newton Method with BFGS Update	16
5	Implicit Filtering	26
6	Sequential Quadratic Programming	43
7	Direct Collocation Algorithm.	74

# 1 Einführung

## 1.1 Beispiele

## 1.2 Fragestellungen

## **1.3 Allgemeine Formulierung eines Optimierungsproblems**

## 1.4 Statische vs. Dynamische Optimierung

## 1.5 Klassifizierung von Optimierungsverfahren

## 1.6 Typische Struktur

## **2** Gradient-Based Optimization without Constraints

#### 2.1 Solution Characterization

This section covers the theoretical results for solving a nonlinear optimization problem using calculus.

#### 2.1.1 One-Dimensional Optimization

For a one-dimensional function  $\varphi(p) : \mathbb{R} \to \mathbb{R}$  the first-order necessary condition for a minimum is that the derivative of  $\varphi(p)$  w.r.t. the parameter p vanishes:

$$\frac{\mathrm{d}\varphi(p^*)}{\mathrm{d}p} = 0$$

Where  $p^*$  denotes the optimal solution, i.e. the minimum.

All solutions that fulfill this condition are *candidates* for a minimum. If  $\varphi$  is twice continuous differentiable, the sufficient condition for a minimum is that the second-order derivative is positive:

$$\frac{\mathrm{d}\varphi(p^*)}{\mathrm{d}p} > 0$$

Then  $p^*$  is called a *strict minimum*. This condition is sufficient, but not necessary! The second-order necessary condition for a minimum is that the second-order derivative is non-negative, i.e.  $\frac{d\varphi(p^*)}{dp} \ge 0$ .

#### **Possibilities for a Minimum**

There are three cases for a minimum:

- $\varphi(p)$  is twice continuously differentiable everywhere
- $\varphi'(p)$  is not continuous everywhere but at  $p^*$
- $\varphi'(p)$  is not continuous everywhere, not even at  $p^*$

While the latter case is common, it is problematic as the solution can typically not be determined analytically (if a function is not continuous at one point, it is rarely invertible).

#### 2.1.2 Multi-Dimensional Optimization

For multi-dimensional objective functions  $\varphi : \mathbb{R}^{n_p} \to \mathbb{R}$ , where  $n_p$  is the dimensionality of the parameters, the first-order necessary condition is that the gradient vanishes:

$$oldsymbol{
abla} oldsymbol{
abla} oldsymbol{
abla} oldsymbol{
abla} oldsymbol{
bla} oldsymbol{
bla}$$

If  $\varphi(\mathbf{p})$  is twice continuously differentiable, the second-order sufficient condition is that the Hessian of  $\varphi(\mathbf{p})$  is positive definite. Analogous to the one-dimensional case, the second-order necessary condition is that the Hessian is positive semi-definite, i.e.:

$$oldsymbol{H}_{arphi}(oldsymbol{p}^{*}) = egin{bmatrix} rac{\partial^{2}arphi}{\partial p_{1}^{2}} & \cdots & rac{\partialarphi^{2}}{\partial p_{n_{p}}p_{1}} \ dots & \ddots & dots \ rac{\partialarphi^{2}}{\partial p_{1}p_{n_{p}}} & \cdots & rac{\partial^{2}arphi}{\partial p_{n_{p}}^{2}} \end{bmatrix} > 0 \quad ext{or respectively} \quad oldsymbol{H}_{arphi}(oldsymbol{p}^{*}) \geq 0$$

Example

#### 2.2 Numerical Gradient-Based Methods

#### 2.2.1 Starting Point

#### **Structure of Gradient-Based Methods**

Given a initial approximation  $p^{(0)}$ , an approximation of the minimum  $p^*$  is wanted. Gradient-based methods are iteration methods based on the iteration rule

$$p^{(k+1)} = p^{(k)} + \alpha^{(k)} d^{(k)}, \quad k = 0, 1, 2, \cdots$$

where

- $d^{(k)}$  is the search direction found as the solution of a linear sub problem and
- $\alpha^{(k)}$  is the step size found by a one-dimensional *line search*.

The iteration terminates once  $p^{(k+1)}$  is "close to"  $p^*$ , e.g. when the gradient nearly vanishes.

#### **Descent Direction**

Gradient-based methods have to ensure the local search direction  $d^{(k)}$  really is a descent direction (the algorithm shall not "run up the hill"). This property is ensured iff the angle  $\delta$  between the search direction and the gradient  $\nabla \varphi^{(k)}$  greater than 90°, i.e.

$$\cos \delta = \frac{\left(\boldsymbol{d}^{(k)}\right)^{T} \left(\boldsymbol{\nabla}\varphi(\boldsymbol{p}^{(k)})\right)}{\left\|\boldsymbol{d}^{(k)}\right\| \cdot \left\|\boldsymbol{\nabla}\varphi(\boldsymbol{p}^{(k)})\right\|} < 0 \quad \iff \quad \left(\boldsymbol{d}^{(k)}\right)^{T} \left(\boldsymbol{\nabla}\varphi(\boldsymbol{p}^{(k)})\right) < 0 \tag{2.1}$$

This is called the "necessary descent condition".

#### **Algorithmic Structure**

The algorithm 1 shows the basic structure of any gradient-based optimization algorithm.

Algorithm 1: Algorithmic structure of a gradient-based optimization algorithms.

1 Initialization: Choose an initial approximation  $p^{(0)}$ , set  $k \leftarrow 0$ 

2 while not converged do

- 3 Determine new search direction:  $oldsymbol{d}^{(k)} \in \mathbb{R}^{n_p}$
- 4 Determine new step size:  $\alpha^{(k)} \in \mathbb{R}^+$
- 5 Update the approximation:  $\boldsymbol{p}^{(k+1)} \leftarrow \boldsymbol{p}^{(k)} + \alpha^{(k)} \boldsymbol{d}^{(k)}$
- $\mathbf{6} \quad | \quad k \leftarrow k+1$

#### 2.2.2 Steepest Descent

Steepest descent is the straightforward way for getting a search direction. The search direction is just set to the negative of the gradient:

$$oldsymbol{d}^{(k)} = - oldsymbol{
abla} arphi(oldsymbol{p}^{(k)})$$

- Advantages:
  - Often quickly reaches areas around the local minimum.
  - No second derivatives needed.
- Disadvantages:
  - Very slow in areas around the local minimum compared to (Quasi-) Newton Methods.

#### 2.2.3 Conjugate Gradient

Basic approach for conjugate gradient: r(0) = -r(0)

$$d^{(0)} = -\nabla \varphi(p^{(0)})$$
  
 $d^{(k)} = \text{Component of } -\nabla \varphi(p^{(k)}) \text{ that is conjugate to } d^{(0)}, d^{(1)}, \cdots, d^{(k-1)}$ 

For a quadratic objective function

$$arphi(oldsymbol{p}) = rac{1}{2}oldsymbol{p}^Toldsymbol{A}oldsymbol{p} - oldsymbol{b}^Toldsymbol{p}$$

with a positive semi-definite matrix A and constant A, b, the search direction is given as the solution of:

$$(\boldsymbol{d}^{(k)})^T \boldsymbol{A} \boldsymbol{d}^{(j)} = 0, \quad j = 1, \cdots, k-1$$

With an optimal step size  $\alpha^{(k)}$ , i.e.

$$\alpha^{(k)} = \arg\min_{\alpha} \varphi(\boldsymbol{p}^{(k)} + \alpha \boldsymbol{d}^{(k)}) \implies \alpha^{(k)} = -\frac{1}{\left(\boldsymbol{d}^{(k)}\right)^T \boldsymbol{A} \boldsymbol{d}^{(k)}} \left(\boldsymbol{\nabla} \varphi(\boldsymbol{p}^{(k)})\right)^T (\boldsymbol{d}^{(k)})$$

the minimum of  $\varphi$  is reached in  $n_p$  steps.

The extension for nonlinear objective functions is given in algorithm 2.

• Exact line search necessary.

Algorithm 2: Conjugate Gradient for nonlinear Objective Function.

1 Initialization: Choose an initial approximation  $p^{(0)}$ , set  $d^{(0)} \leftarrow -\nabla \varphi(p^{(0)})$  and  $k \leftarrow 0$ 2 while not converged do  $\alpha^{(k)} \leftarrow \arg \min_{\alpha} \varphi(p^{(k)} + \alpha d^{(k)})$  $p^{(k+1)} \leftarrow p^{(k)} + \alpha^{(k)} d^{(k)}$  $\beta^{(k+1)} \leftarrow \frac{(\nabla \varphi(p^{(k+1)}))^T (\nabla \varphi(p^{(k+1)}))}{(\nabla \varphi(p^{(k)}))^T (\nabla \varphi(p^{(k)}))}$  $d^{(k+1)} \leftarrow -\nabla \varphi(p^{(k+1)}) + \beta^{(k+1)} d^{(k)}$  $k \leftarrow k + 1$ 

• Different variants of GC-algorithms mainly distinguish in the choice of of  $\beta^{(k)}$ .

- Advantages:
  - Faster then steepest descent.
  - No explicit storing of the Hessian  $H_{\omega}(p^{(k)})$  necessary.
  - No explicit matrix-vector multiplication.
  - Useful even for extreme high dimensions  $n_p$ .
  - Exact for quadratic objectives  $\varphi(\mathbf{p})$ .
- Disadvantages:
  - A lot slower then (Quasi-) Newton Methods.
  - In general not useful for optimizing simulation models.

#### 2.2.4 Newton Method

Assuming the approximation of each iteration,  $p^{(k)}$ , is close to the minimum  $p^*$ , the gradient  $\nabla \varphi(p^*)$  can be taylor-expanded around  $p^{(k)}$ :

$$\boldsymbol{\nabla} \varphi(\boldsymbol{p}^*) \stackrel{T(\boldsymbol{p}^{(k)})}{=} \boldsymbol{\nabla} \varphi(\boldsymbol{p}^{(k)}) + \boldsymbol{H}_{\varphi}(\boldsymbol{p}^{(k)})(\boldsymbol{p}^* - \boldsymbol{p}^{(k)}) + \cdots \stackrel{!}{=} \boldsymbol{0}$$

By leaving our the higher order terms the search direction  $d^{(k)} := p^* - p^{(k)}$  is given by the solution of the system of linear equations

$$oldsymbol{H}_arphi(oldsymbol{p}^{(k)})oldsymbol{d}^{(k)} = -oldsymbol{
abla}arphi(oldsymbol{p}^{(k)})$$

The realization is shown in algorithm 3.

When plugging the search direction into the necessary descent condition (2.1)

$$(\boldsymbol{d}^{(k)})^T \Big( \boldsymbol{\nabla} \varphi(\boldsymbol{p}^{(k)}) \Big) = - \Big( \boldsymbol{\nabla} \varphi(\boldsymbol{p}^{(k)})^T \Big) \Big( \boldsymbol{H}_{\varphi}(\boldsymbol{p}^{(k)}) \Big)^{-1} \Big( \boldsymbol{\nabla} \varphi(\boldsymbol{p}^{(k)}) \Big) < 0$$

it is clear that this is only fulfilled iff the Hessian is positive definite. But this is only the case in a region around the minimum! If the approximation is far away from the minimum, the search direction might also be an ascent direction causing the Newton method to diverge. There are two main solutions to this problem:

Algorithm 3: Newton Method

1 Initialization: Choose an initial approximation  $p^{(0)}$ , set  $k \leftarrow 0$ 2 while not converged do 3 Solve  $H_{\varphi}(p^{(k)})d^{(k)} = -\nabla\varphi(p^{(k)})$  for  $d^{(k)}$ 4  $\alpha^{(k)} \leftarrow \arg\min_{\alpha}\varphi(p^{(k)} + \alpha d^{(k)})$ 5  $p^{(k+1)} \leftarrow p^{(k)} + \alpha^{(k)}d^{(k)}$ 6  $k \leftarrow k + 1$ 

- 1. If the Hessian is not positive definite, replace it by an identity matrix. That is, set the search direction to the steepest descent.
- 2. Regularize the equation system with a weight  $\nu > 0$  such that the new matrix is positive definite (this "rotates" the matrix in the direction of the steepest descent such that the new search direction always fulfills the descent condition):

$$\left(oldsymbol{H}_{arphi}oldsymbol{\left(p^{(k)}
ight)}+
uoldsymbol{I}oldsymbol{d}^{(k)}=-oldsymbol{
abla}arphioldsymbol{\left(p^{(k)}
ight)}$$

- Advantages:
  - Near to strong local minima of twice continuous differentiable objective, the Newton method is quadratic convergent.
- Disadvantages:
  - Computationally expensive as a linear system has to be solved in every iteration.
  - Not only first, but also second-order derivatives have to be available. This is a big disadvantage:
    - \* In practice, the first derivative is rarely and the second derivative is never available.
    - \* Even a single wrong component in the gradient or the Hessian destroys the quadratic convergence.

#### **Availability of Second-Order Derivatives**

The obvious idea is to approximate the Hessian using finite differences. The approximated Hessian is then given as

$$oldsymbol{H}_arphiig(oldsymbol{p}^{(k)}ig) = rac{1}{2}ig( ilde{oldsymbol{H}}+ ilde{oldsymbol{H}}^T)$$

where  $\tilde{H}$  is given by

$$\tilde{\boldsymbol{H}}_{i} = \frac{\partial}{\partial p_{i}} \Big( \boldsymbol{\nabla} \varphi(\boldsymbol{p}^{(k)}) \Big) \approx \frac{1}{h_{i}} \Big( \boldsymbol{\nabla} \varphi(\boldsymbol{p}^{(k)} + h_{i}\boldsymbol{e}_{i}) - \boldsymbol{\nabla} \varphi(\boldsymbol{p}^{(k)}) \Big)$$

where  $\tilde{H}_i$  is the *i*-th column of  $\tilde{H}$ . The equation (KKT.i) is used to force the Hessian to be symmetric. Problems:

• The Hessian  $H_{\varphi}(p^{(k)}) = \frac{1}{2}(\tilde{H} + \tilde{H}^T)$  is not necessarily positive definite.

- In every iteration the Gradient has to be evaluated  $n_p$  times more.
- The linear system still needs to be solved.
- Only useful for high-dimensional problems with sparse gradients!

Another possibility are Quasi-Newton Methods.

#### 2.2.5 Quasi-Newton Methods

Quasi-Newton methods are equivalent to the Newton method, however, the Hessian (or its inverse) is approximated by a positive definite matrix

$$\hat{oldsymbol{H}}^{(k)}pproxoldsymbol{H}_arphi(oldsymbol{p}^{(k)})$$

that is updated in every iteration. This yields a lot of advantages over the classic Newton method:

- Only first-order derivatives needed.
- As  $\hat{H}$  constructed positive definite, the descent condition is fulfilled anytime.
- If the inverse Hessian is directly approximated, only  $\mathcal{O}(n_p^2)$  multiplications instead of  $\mathcal{O}(n_p^3)$  for solving the linear system.

But how to do the Hessian update? By Taylor-expanding the gradient  $oldsymbol{
abla} arphi(oldsymbol{p}^{(k)})$  around  $oldsymbol{p}^{(k+1)}$ 

$$oldsymbol{
abla} oldsymbol{
abla} oldsymbol{p} arphi(oldsymbol{p}^{(k+1)}) \stackrel{Tig(oldsymbol{p}^{(k+1)}ig)}{=} oldsymbol{
abla} arphi(oldsymbol{p}^{(k+1)}ig) + oldsymbol{H}_arphi(oldsymbol{p}^{(k+1)}ig) ig(oldsymbol{p}^{(k+1)}ig) + oldsymbol{e} oldsymbol{e} oldsymbol{0} = oldsymbol{0}$$

and cutting off the higher-order terms, the following approximation holds:

$$oldsymbol{H}_arphi(oldsymbol{p}^{(k+1)})oldsymbol{d}^{(k)}pproxoldsymbol{
abla}arphi(oldsymbol{p}^{(k+1)})-oldsymbol{
abla}arphi(oldsymbol{p}^{(k)})$$

The approximation of the Hessian must therefore fulfill the secant condition

$$ilde{oldsymbol{H}}^{(k+1)}oldsymbol{d}^{(k)} = oldsymbol{
abla} arphioldsymbol{p}^{(k+1)}ilde{oldsymbol{--}} - oldsymbol{
abla} arphioldsymbol{(p^{(k)})}$$

There exist a lot of different approaches for doing the Hessian updates  $\tilde{H}^{(k+1)} = \tilde{H}^{(k)} + U(k)$  for rank-1 or rank-2 matrices  $U^{(k)}$ :

- Approach for rank-1 corrections:  $\tilde{H}^{(k+1)} = \tilde{H}^{(k)} + \beta_1 u u^T$
- Approach for rank-2 corrections:  $\tilde{H}^{(k+1)} = \tilde{H}^{(k)} + \beta_1 u u^T + \beta_2 v v^T$

The vectors  $u, v \in \mathbb{R}^{n_p}$  and scalars  $\beta_1, \beta_2 \in \mathbb{R}$  must have to be chosen such that  $\tilde{H}^{(k+1)}$  is

- positive definite,
- symmetric,
- fulfills the secant condition and
- adding up the matrices is efficient and robust.

#### **BFGS-Update**

The most known rank-2 update for the Hessian is the BFGS-Update<sup>1</sup>

$$egin{aligned} m{u} &= ilde{m{H}}^{(k)} m{d}^{(k)} & & eta_1 &= -rac{1}{ig(m{d}^{(k)}ig)^T ilde{m{H}}^{(k)} m{d}^{(k)}} \ m{v} &= m{g}^{(k)} & & eta_2 &= rac{1}{ig(m{g}^{(k)}ig)^T m{d}^{(k)}} \end{aligned}$$

where  $g^{(k)} = \nabla \varphi(p^{(k+1)}) - \nabla \varphi(p^{(k)})$ . Plugging that into the general approach for rank-2 updates yields the update rule for BFGS-approximations:

$$\tilde{\bm{H}}^{(k+1)} = \tilde{\bm{H}}^{(k)} - \frac{1}{\left(\bm{d}^{(k)}\right)^T \tilde{\bm{H}}^{(k)} \bm{d}^{(k)}} \tilde{\bm{H}}^{(k)} \bm{d}^{(k)} \left(\tilde{\bm{H}}^{(k)} \bm{d}^{(k)}\right)^T + \frac{1}{\left(\bm{g}^{(k)}\right)^T \bm{d}^{(k)}} \bm{g}^{(k)} \left(\bm{g}^{(k)}\right)^T$$

- The direct approximation of the Hessian inverse is not really robust (e.g. for a non-optimal step size rule).
- A better alternative is to directly approximate a useful factorization, e.g. the Cholesky decomposition. This is more robust and equally efficient  $(\mathcal{O}(n_p^2))$ .

The pseudo code for the BFGS update is shown in algorithm 4.

Algorithm 4: Quasi-Newton Method with BFGS Update.

1 Initialization: Choose an initial approximation  $p^{(0)}$ , set  $\tilde{H}^{(0)} = I$  and  $k \leftarrow 0$ 2 while not converged do 3 Solve  $\tilde{H}^{(k)}d^{(k)} = -\nabla\varphi(p^{(k)})$  for  $d^{(k)}$  $\alpha^{(k)} \leftarrow \arg\min_{\alpha}\varphi(p^{(k)} + \alpha d^{(k)})$  $p^{(k+1)} \leftarrow p^{(k)} + \alpha^{(k)}d^{(k)}$  $g^{(k)} \leftarrow \nabla\varphi(p^{(k+1)}) - \nabla\varphi(p^{(k)})$  $\tilde{H}^{(k+1)} \leftarrow \tilde{H}^{(k)} - \frac{1}{(d^{(k)})^T\tilde{H}^{(k)}d^{(k)}}\tilde{H}^{(k)}d^{(k)}(\tilde{H}^{(k)}d^{(k)})^T + \frac{1}{(g^{(k)})^Td^{(k)}}g^{(k)}(g^{(k)})^T$  $k \leftarrow k+1$ 

#### 2.2.6 Comparison

#### 2.2.7 Notes and Discussion

- The convergence of gradient-based methods can be shown under weak preconditions.
- As the search direction is only a local descent direction, gradient-based algorithms only yields local minima.
- There is no algorithm that can guarantee to find the global minimum!
- Some approaches for determining a global minimum:

<sup>1</sup>"BFGS" stands for the authors Broyden, Fletcher, Goldfarb and Shanno.

- Choose the initialization well, i.e. close to the global minimum.
- Execute the algorithm multiple times with different starting points.
- Validate the solution against properties of the original problem.
- Execute direct search methods beforehand to find promising regions for the local minimum search.
- Advantages:
  - If gradient-based algorithms converge, they converge utterly fast.
  - Efficient also for high-dimensional problems, i.e. a large  $n_p$ .
- Disadvantages:
  - Only applicable for functions that are differentiable almost everywhere.
  - Require gradient information exact up to four to eight decimal points.
  - Convergence to a local minimum near the initialization  $p^{(0)}$ .
  - Require some expert knowledge.

#### 2.3 Step Size Rules, Line Search

In every iteration of gradient-based algorithms, the step size has to be determined by minimizing the onedimensional function:

$$\psi(\alpha) = \varphi(\boldsymbol{p}^{(k)} + \alpha \boldsymbol{d}^{(k)})$$

As the necessary first-order condition for a minimum, the derivative w.r.t.  $\alpha$  has to vanish:

$$\frac{\mathrm{d}\psi(\alpha^{(k)})}{\mathrm{d}\alpha} = \frac{\mathrm{d}}{\mathrm{d}\alpha}\varphi(\boldsymbol{p}^{(k)} + \alpha^{(k)}\boldsymbol{d}^{(k)}) = \left(\boldsymbol{\nabla}\varphi(\boldsymbol{p}^{(k)} + \alpha^{(k)}\boldsymbol{d}^{(k)})\right)^{T}\boldsymbol{d}^{(k)} \stackrel{!}{=} 0$$

Thus the gradient of  $\varphi$  at the minimum  $\alpha^{(k)}$  as to be orthogonal to the search direction  $d^{(k)}$ . Intuitively, the optimal step size has to be chosen such that the iteration step cannot go any further without ascending again ("hitting an ascending contour line").

Goal of the line search is to reach the minimum of  $\psi$  with least invocations of  $\psi$  as possible. Most of the existing search methods can be classified into

- Polynomial approximation, e.g. quadratic or cubic interpolation
- Direct search methods, e.g. Fibonacci-search, golden ratio search
- Optimal vs. non-optimal search methods, e.g. by finding an improvement but not the minimum
- Usage of the gradient information  $\psi'$  or not.

**Requirements:** 

- Finding the  $\alpha^{(k)}$  with a minimal value of  $\psi$ .
- Do not waste too much computation time on the line search.

In general, an exact line search requires lots of  $\psi$ -evaluations. But how far does  $\psi$  need to be reduced in order to guarantee convergence? In general, the condition  $\varphi(p^{(k)} + \alpha^{(k)}d^{(k)}) < \varphi(p^{(k)})$  is not enough!

#### 2.3.1 Inexact Line Search

Procedure: Generate and inspect a series of candidates for  $\alpha^{(k)}$  and terminate once one of the candidates fulfills specific criteria, e.g. the Armijo rule or Wolfe conditions.

#### Armijo Rule

The *Armijo rule* guarantees a sufficient reduction in  $\varphi$ :

$$\varphi(\boldsymbol{p}^{(k)} + \alpha^{(k)}\boldsymbol{d}^{(k)}) \leq \varphi(\boldsymbol{p}^{(k)}) + c_1\alpha^{(k)} \left(\boldsymbol{\nabla}\varphi(\boldsymbol{p}^{(k)})\right)^T \boldsymbol{d}^{(k)} = \varphi(\boldsymbol{p}^{(k)}) + c_1\alpha^{(k)}\psi'(0)$$

Where  $0 < c_1 < 1$  is any constant, e.g.  $c_1 = 10^{-4}$ .

Hence, the minimal reduction has to be proportional to  $\alpha^{(k)}$  and the derivative  $\psi'(0)$ .

#### **Curvature Condition**

But a sufficient descent condition is not enough as the step sizes must not be too small (otherwise progress would stop). Thus a second condition has to be employed, the *curvature condition* that requires a minimum curvature on  $\psi$ :

$$\left(\boldsymbol{\nabla}\varphi(\boldsymbol{p}^{(k)}+\alpha^{(k)}\boldsymbol{d}^{(k)})\right)^{T}\boldsymbol{d}^{(k)} \geq c_{2}\cdot\left(\boldsymbol{\nabla}\varphi(\boldsymbol{p}^{(k)})\right)^{T}\boldsymbol{d}^{(k)} = c_{2}\psi'(0) \quad \iff \quad \psi'(\boldsymbol{p}^{(k)}) \geq c_{2}\psi'(0)$$

Where  $c_1 < c_2 < 1$  is any constant, e.g.  $c_2 = 0.9$ .

#### **Wolfe and Goldstein Conditions**

Combining the Armijo rule and the curvature condition yields the Wolfe conditions that guarantee both a minimal reduction and a minimal curvature. They are especially useful for Quasi-Newton methods as the Wolfe conditions are scale invariant, i.e. independent of

- multiplying  $\varphi$  with any constant and
- affine transformations of *p*.

There are other possibilities are, e.g. the Goldstein conditions

$$\varphi(\boldsymbol{p}^{(k)}) + (1-c)\alpha^{(k)} \left(\boldsymbol{\nabla}\varphi(\boldsymbol{p}^{(k)})\right)^T \boldsymbol{d}^{(k)} \le \varphi(\boldsymbol{p}^{(k)} + \alpha^{(k)}\boldsymbol{d}^{(k)}) \le \varphi(\boldsymbol{p}^{(k)}) + c\alpha^{(k)} \left(\boldsymbol{\nabla}\varphi(\boldsymbol{p}^{(k)})\right)^T \boldsymbol{d}^{(k)}$$

with any 0 < c < 1/2, that are useful for Newton, but not for Quasi-Newton methods.

#### 2.3.2 Notes

- For gradient-based methods a step size  $\alpha^{(k)} > 1$  is in general not useful because:
  - The search direction  $d^{(k)}$  is determined using a linear or quadratic Taylor approximation.
  - The Taylor approximation is only valid in a small region around the current approximation  $p^{(k)}$ , i.e. for  $0 < \alpha^{(k)} \le 1$ .
- The local quadratic or super-linear convergence of Newton-type methods is visible in practice as the last step can be executed with full step size  $\alpha^{(k)} = 1$ .

#### 2.4 Trust Region Methods

Gradient-based methods with line search determine a fixed search direction and adjust the step size  $\alpha^{(k)}$  according to that search direction to reach global convergence.

Another approach is to determine both the length and direction of  $d^{(k)}$ . The iteration step then becomes

$$p^{(k+1)} = p^{(k)} + d^{(k)}$$

without any explicit step size. The length and direction of  $d^{(k)}$  are then determined as a solution of the quadratic sub-problem

$$\min_{\boldsymbol{d} \in \mathbb{R}^{n_p}} \left( \boldsymbol{\nabla} \varphi(\boldsymbol{p}^{(k)}) \right)^T \boldsymbol{d} + \frac{1}{2} \boldsymbol{d}^T \boldsymbol{H}_{\varphi}(\boldsymbol{p}^{(k)}) \boldsymbol{d}$$
subject to  $\|\boldsymbol{d}\|_2 \leq \delta$ 

where  $\delta$  describes the area around the current approximation  $p^{(k)}$  where the quadratic approximation of  $\varphi(p^{(k)} + d)$  makes sense, i.e. the *trust region*.

The value of  $\delta$  is extremely important for the efficiency of the method.

- If  $\delta$  is too small, opportunities for large steps are missed.
- If δ is too large, the minimum of the quadratic approximation might be far off the minimum of the objective if the Hessian is indefinite or negative definite.

It is possible to add a regularization parameter  $\beta \ge 0$  to the quadratic approximation

$$\left( oldsymbol{H}_{arphi}ig(oldsymbol{p}^{(k)}ig)+etaoldsymbol{I}ig)oldsymbol{d}=-oldsymbol{
abla}arphiig(oldsymbol{p}^{(k)}ig)$$

such that the matrix is positive semidefinite. The solution of this regularized problem also solves the trust region problem if either  $\beta = 0$ ,  $\|d\| \le \delta$  or  $\beta \ge 0$ ,  $\|d\| = \delta$ .

#### 2.5 Rate of Convergence

A common criteria to measure the performance of a gradient method are *rates of convergence*. These provide information on how fast an algorithm converges, i.e. how fast  $p^{(k)} \rightarrow p^*$  or  $||p^{(k)} - p^*|| \rightarrow 0$ .

**Definition:** Let  $\{\{p^{(k)}\}\)$  be the series of approximations produced by an optimization method. Then this series has rate of convergence r if r is the greatest positive number such that the limit

$$0 \leq \lim_{k o \infty} rac{\left\| oldsymbol{p}^{(k+1)} - oldsymbol{p}^{st} 
ight\|}{\left\| oldsymbol{p}^{(k)} - oldsymbol{p}^{st} 
ight\|^r} = \gamma < \infty$$

converges (where  $p^* = \lim_{k \to \infty} p^{(k)}$ ). If r = 1, then  $\gamma < 1$  has to hold for the method to converge.

A sequence is said to converge superlinearly if

$$\lim_{k \to \infty} \frac{\left\| \boldsymbol{p}^{(k+1)} - \boldsymbol{p}^* \right\|}{\left\| \boldsymbol{p}^{(k)} - \boldsymbol{p}^* \right\|} = 0$$

holds. Even though technically this condition holds for r > 1, in practice only methods with 1 < r < 2 are said to converge superlinearly (e.g. for r = 2, the sequence is called to converge quadratically).

#### Examples

#### 2.5.1 Gradient-Based Methods

Under ideal conditions (i.e.  $\varphi$  is twice continuously differentiable and  $H_{\varphi}(p^*)$  is positive definite), all of the following hold:

- Steepest descent is (locally) linearly convergent (with exact line search).
- The Newton method is (locally) quadratically convergent.
- Quasi-Newton methods with BFGS-update are (locally) superlinearly convergent (for inexact line search using the Wolfe conditions).

But: Even a single wrong component in the Hessian reduced the quadratic convergence of the Newton method to linear convergence!

# **3 Gradient-Free Optimization without Constraints**

This chapter covers different types of sampling methods (direct search methods):

- 1. Metaheuristics (random search),
- 2. Deterministic Sampling Methods (pattern search) and
- 3. Surrogate optimization.

These optimization methods only use evaluations of the objective  $\varphi$  and do not user gradient information (neither analytically nor using numerical differentiation). That is,  $\varphi$  is used as a "Black Box" for function evaluations.

Even for  $\varphi$  that are not differentiable and have lots of local minima, gradient-free algorithms are remarkably robust, but can also fail fast.

## 3.1 Introduction

The objective  $\varphi$  that is to be optimized often has suboptimal properties, e.g.:

- the evaluation is noisy
- not differentiable
- high computation time for evaluation (e.g. when a simulation has to be run for each evaluation)

But gradient-free techniques have a wide range of applications, e.g. in automobile, aerospace industry, robotics, financial, etc. The goal is to reduce the objective function as much as possible and find regions in which the objective is differently sensitive w.r.t. changes in the optimization variables.

#### 3.1.1 Simulation-Based Optimization

In a simulation-based setting, the objective is evaluated by running a simulation (e.g. by solving a set of differential equations or running a real experiment). This yields a suboptimal setting for optimization, yet a common one:

- Only function values are computable and not gradient information is available.
- Source code of the optimization is often not available. Hence, no information about how the simulation is computed.
- Discontinuous/non-differentiable systems and model properties, e.g. due to collisions.
- Non-differentiable structure inside the simulation (e.g. if-else).
- Discontinuities caused by subprograms, heuristics, table data, ...

- Function evaluations are computationally expensive.
- Numerical "noise" overlay the actual system properties.
- Non-deterministic simulations (e.g. due to complex friction).

Hence, simulation-based optimization is a black box problem and can be solved (or approximated) using gradient-free methods!

#### 3.1.2 Black-Box Optimization

Naturally, gradient-based optimization methods are not well-suited for problems with "low" differentiability and computationally expensive function evaluations. These problems may be solved using gradient-free optimization methods. But...

- Gradient-based techniques are really slow in comparison to gradient-based ones for differentiable optimization problems.
- They need a lot of function evaluations for high-dimensional problems and are thus practically only applicable for problems with dimensions  $n_p < 100$ , better  $n_p < 20$ .
- Have lots of problems with nonlinear equality constraints!
- The theory of gradient-free methods is not as mature as the theory for gradient-based methods.

## 3.2 Metaheuristics

#### 3.2.1 Evolutionary Algorithm (EA)

A evolutionary algorithm mimics the biological evolutionary strategy with random search methods.

- 1. Initialization: Choose one "Parent"  $p^{(0)}$  and a number of descendants  $\ell$ .
- 2. Iteration: Create  $\ell$  descendants via "mutation"

$$\boldsymbol{p}^{(k,i)} = \boldsymbol{p}^{(k)} + \alpha_i^{(k)} \boldsymbol{d}_i, \quad i = 1, \cdots, \ell$$

where  $d_i \in \mathbb{R}^{n_p}$  are vectors of Gaussian distributed variables and  $\alpha_i \in \mathbb{R}$  art suitable "mutation step sizes" Then select the descent with the lower  $\varphi$  value and repeat.

#### 3.2.2 Genetic Algorithms (GA)

*Genetic algorithms* are inspired by biological evolutionary strategies the positive properties caused by mutation are kept through natural selection. GAs are applicable for both real and discrete optimization variables p.

- 1. Initialization: Choose a suitable set of different "individuals" (first generation).
- 2. Evaluation: Determine the "fitness" of each candidate using the objective/fitness function.
- 3. Selection: Randomly select candidates of the current generation (the higher the fitness, the higher the probability to be chosen).

- 4. Recombination: Combine values (genomes) of the selected individuals and create new individuals.
- 5. Mutation: Randomly change the genomes.
- 6. New Generation: Select new individuals as the new generation and continue with step 2.

#### Example

#### 3.2.3 Further Metaheuristics

Further metaheuristics based on real representations of p like in evolutionary algorithms:

- Particle Swam: Population method, uses idea of combining local and swarm knowledge, direction and speed for particles are adjusted
- ...

Further metaheuristics based on binary representations of p like in genetic algorithms:

- Tabu Search: A list of possible manipulations is given, another lists dynamically the inverses of them, these cannot be applied any more
- ...

## 3.3 Deterministic Sampling Methods (Pattern Search Methods)

Deterministic sampling methods can be further categorized into

- Qualitative methods: Only ranking w.r.t. to the function value.
  - Simplex Methods
  - Coordinate- or compass-search
  - Multidirectional search
  - Pattern search methods
  - ...
- Quantitative methods: Consideration of the real function values.
  - Implicit filtering (based on the simplex method)
  - DIRECT (dividing rectangles)
  - ...

## 3.3.1 Nelder-Mead Simplex Method

A simplex is a simple object that consists of  $n_p + 1$  points  $p^{(i)}$  in the parameter space (which is  $n_p$ -dimensional). In a 2D space, the three points form a triangle. In the iteration phase the values of  $\varphi$  at the corners are compared and the simplex is transformed according to specific rules (see subsubsection 3.3.1). The algorithm terminates once the simplex contracts onto a single point.

#### **Iteration Phase**

The iteration phase starts by sorting the edges according to its  $\varphi$ -values:

$$\varphi(\boldsymbol{p}^{(1)}) \leq \varphi(\boldsymbol{p^{(2)}}) \leq \cdots \leq \varphi(\boldsymbol{p}^{(n_p+1)})$$

where  $p^{(1)}$  is called the *best* point and  $p^{(n_p+1)}$  is called the *worst*. The algorithm now tried to replace the worst point  $p^{(n_p+1)}$  with another point of the form

$$\boldsymbol{p}(\mu) = (1+\mu)\bar{\boldsymbol{p}} - \mu \boldsymbol{p}^{(n_p+1)}$$

Where  $\bar{p}$  is the centroid of the of all points *except the worst*, i.e.:

$$ar{oldsymbol{p}} = rac{1}{n_p}\sum_{i=1}^{n_p}oldsymbol{p}^{(i)}$$

This corresponds to a reflection of the worst point over the centroid with a weight  $\mu$  that specifies "how far the point point gets pushed out", i.e. the ratio of the original distance of the worst point to the centroid that is preserved while reflecting. If  $\mu = 1$ , the point is mirrored.

In every iteration, the value  $\mu$  is chosen of a set of four values

$$-1 < \mu_{ic} < 0 < \mu_{oc} < \mu_r < \mu_e$$

for example  $(\mu_{ic}, \mu_{oc}, \mu_r, \mu_e) = (-0.5, 0.5, 1, 2).$ 

#### Algorithm

Some termination criteria are for example:

- Exactness in the objective:  $\varphi(p^{(n_p+1)}) \varphi(p^{(1)}) \leq \varepsilon$
- Maximum number of function evaluations:  $k = k_{\text{max}}$
- Sufficient small distance on the simplex corners.
- Initialization: Choose a start simplex, evaluate the objective at the corners, sort them and set  $k = n_p + 1$ .
- Iteration: While  $\varphi(p^{(n_p+1)}) \varphi(p^{(1)}) > \varepsilon$  and  $k < k_{\max}$ , do:
  - (a) Calculate the centroid  $\bar{p} = \frac{1}{n_n} \sum_{i=1}^{n_p} p^{(i)}$ .
  - (b) *Reflection:* If  $\varphi(\mathbf{p}^{(1)}) \leq \varphi(\mathbf{p}(\mu_r)) < \varphi(\mathbf{p}^{(n_p)})$ , replace  $\mathbf{p}^{(n_p+1)}$  with  $\mathbf{p}(\mu_r)$ ; go to (g). "Use the reflected point if it is better than the second worst, but not better than the best."
  - (c) *Expansion*: If  $\varphi(\mathbf{p}(\mu_r)) < \varphi(\mathbf{p}^{(1)})$ , then: "If the reflected point is better than the best, ..."
    - If  $\varphi(\boldsymbol{p}(\mu_e)) < \varphi(\boldsymbol{p}(\mu_r))$ , replace  $\boldsymbol{p}^{(n_p+1)}$  with  $\boldsymbol{p}(\mu_e)$ ; go to (g).
      - "... and the expanded point is better than the reflected point, use the expanded point."
    - If  $\varphi(\boldsymbol{p}(\mu_r)) < \varphi(\boldsymbol{p}(\mu_e))$ , replace  $\boldsymbol{p}^{(n_p+1)}$  with  $\boldsymbol{p}(\mu_r)$ ; go to (g).
      - "... and the expanded point is worst than the reflected point, use the reflected point."

(d) *Outer Contraction*: If  $\varphi(\mathbf{p}^{(n_p)}) \leq \varphi(\mathbf{p}(\mu_r)) < \varphi(\mathbf{p}^{(n_p+1)})$ , then: "If the reflected point is better than the worst, but worst than second worst, ..."

- If φ(p(μ<sub>oc</sub>)) < φ(p(μ<sub>r</sub>)), replace p<sup>(n<sub>p</sub>+1)</sup> with p(μ<sub>oc</sub>); go to (g).
  "... and the outer contraction point is better than the reflected point, use the outer contraction point."
- Else, go to (f).
- (e) *Inner Contraction:* If  $\varphi(\mathbf{p}^{(n_p+1)}) \leq \varphi(\mathbf{p}(\mu_r))$ , then: "If the reflected point is worst than the worst point, ..."
  - If φ(p(μ<sub>ic</sub>)) < φ(p<sup>(n<sub>p</sub>+1)</sup>), replace p<sup>(n<sub>p</sub>+1)</sup> with p(μ<sub>ic</sub>); go to (g).
     "...and the inner contraction point is better than the worst, use the inner contraction point."
  - Else, go to (f).
- (f) Shrink: For all  $2 \le i \le n_p + 1$ , set  $p^{(i)} = p^{(1)} \frac{1}{2} (p^{(i)} p^{(1)})$ .
- (g) Sort: Sorting the current simplex corners such that  $\varphi(p^{(1)}) \leq \varphi(p^{(2)}) \leq \cdots \leq \varphi(p^{(n_p+1)})$  holds again. Repeat.

The Nelder-Mead method therefore always tries to create a big simplex and only shrink if every other action would yield a worst corner/simplex.

#### Notes

- The method is not guaranteed to converge. But in practice, it yields good results.
- Can get stuck on a suboptimal point such that the algorithm has to be restarted with other initial simplex corners.

**Simplex-Gradient** It is possible to detect stagnation using a *Simplex-Gradient*  $D^{(k)} \in \mathbb{R}^{n_p}$ ,  $D^{(k)} = (V^{(k)})^{-T} \delta^{(k)}$  where  $V^{(k)}$  is the matrix of the simplex directions

$$\boldsymbol{V}^{(k)} = \begin{bmatrix} \boldsymbol{p}^{(2)} - \boldsymbol{p}^{(1)} & \boldsymbol{p}^{(3)} - \boldsymbol{p}^{(1)} & \cdots & \boldsymbol{p}^{(n_p+1)} - \boldsymbol{p}^{(1)} \end{bmatrix} \eqqcolon \begin{bmatrix} \boldsymbol{v}^{(1)} & \boldsymbol{v}^{(2)} & \cdots & \boldsymbol{v}^{(n_p)} \end{bmatrix} \in \mathbb{R}^{n_p \times n_p}$$

and  $\delta^{(k)}$  is the vector of the objective differences:

$$\boldsymbol{\delta}^{(k)} = \begin{bmatrix} \varphi(\boldsymbol{p}^{(2)}) - \varphi(\boldsymbol{p}^{(1)}) \\ \varphi(\boldsymbol{p}^{(3)}) - \varphi(\boldsymbol{p}^{(1)}) \\ \vdots \\ \varphi(\boldsymbol{p}^{(n_p+1)}) - \varphi(\boldsymbol{p}^{(1)}) \end{bmatrix} \in \mathbb{R}^{n_p}$$

Analogous to a gradient-based method, this yields a condition for minimum progress

$$\hat{\varphi}^{(k+1)} - \hat{\varphi}^{(k)} < -\alpha \| \boldsymbol{D}^{(k)} \|^2, \qquad \hat{\varphi} = \frac{1}{n_p + 1} \sum_{i=1}^{n_p + 1} \varphi(\boldsymbol{p}^{(i)})$$

with a small  $\alpha > 0$ . One approach for a condition on when to restart is to restart if both

$$\hat{\varphi}^{(k+1)} - \hat{\varphi}^{(k)} > -\alpha \| \boldsymbol{D}^{(k)} \|^2 \text{ and } \hat{\varphi}^{(k+1)} - \hat{\varphi}^{(k)} < 0$$

hold.

#### 3.3.2 Multidirectional Search Methods

- In the Nelder-Mead method a bad conditioning of the simplex, i.e. the matrix  $V^{(k)}$ , leads to problems that cannot be avoided.
- In multidirectional search methods this problem is avoided by making every simplex congruent to its predecessors.
- The algorithm uses similar steps for reflection, expansion and contraction, but possibly needs a lot more function evaluations.

#### 3.3.3 Asynchronous Parallel Pattern Search (APPS)

- Asynchronous Parallel Pattern Search is a pattern-based search method on a grid.
- The direction of the pattern determines the descent direction.
- Patterns can be varied while maintaining their mathematical properties.
- It is "naturally" parallelizable.

#### 3.3.4 Implicit Filtering

*Implicit filtering* is a descent method using "smooth" approximations of the gradients. It uses a central approximation of the simplex gradient

$$oldsymbol{D}_C^{(k)} = rac{1}{2} ig( oldsymbol{D}^{(k)} + oldsymbol{D}_R^{(k)} ig)$$

where  $D_R^{(k)}$  is the gradient of the simplex that is reflected around  $p^{(k)}$ .

The structure of implicit filtering is sketched inalgorithm 5.

Algorithm 5: Implicit Filtering.

1 Initialization: Choose  $\alpha, \beta \in (0, 1)$  and set H = I2 while not converged do Calculate  $\varphi(p^{(k)})$ ,  $D_C^{(k)}$  and the search direction  $d^{(k)} = -H^{-1}D_C^{(k)}$ 3 Inexact line search for  $j = 1, \dots, j_{\text{max}}, \lambda \coloneqq \beta^j$  until the following holds: 4  $\varphi(\boldsymbol{p} + \lambda \boldsymbol{d}) - \varphi(\boldsymbol{p}) \leq \alpha \lambda \boldsymbol{\nabla}_{\Delta \boldsymbol{p}_k} \varphi^T(\boldsymbol{p}) \boldsymbol{d}^{(k)}$  $\boldsymbol{p}^{(k+1)} \leftarrow \boldsymbol{p}^{(k)} + \lambda \boldsymbol{d}^{(k)}$ 5 if line search successful then 6 Quasi-Newton update of the Hessian H with  $p^{(k+1)} - p^{(k)}$  and  $D_C^{(k+1)} - D_C^{(k)}$ 7 else 8  $H \leftarrow I$ 9 Shrink the simplex 10

#### 3.4 Surrogate Optimization

In *surrogate optimization methods*, the (complex) objective is replaced with a simpler approximation that maintains the key properties of the objective (e.g. a noisy measurement might be replaced by a simpler regression model). This surrogate function is then minimized and adjusted in order to find a good approximation of the solution of the original problem (this can also be applied for constraint functions).

Requirements for the surrogate function  $\hat{\varphi} : \mathbb{R}^{n_p} \to \mathbb{R}$ : For all function evaluation (or "sampling") points  $(p^{(i)}, \varphi(p^{(i)}))$ ,  $i = 1, \dots, m$  it must hold that

$$\varphi(\boldsymbol{p}^{(i)}) = \hat{\varphi}(\boldsymbol{p}^{(i)}) + \epsilon$$

where  $\epsilon \in \mathbb{R}$  is some "slack" constant that determines how exact the surrogate function shall be.

- For  $\epsilon = 0$ , the problem is the same as interpolation.
- For  $\epsilon > 0$ , the surrogate function does not perfectly reproduce the objective, but might be smoother.

Further requirements are that  $\hat{\varphi}$  should be fast to compute and the gradients of  $\hat{\varphi}$  should be available in closed form.

This raises some questions:

- 1.  $\varphi$  might be too complex for a simple approximation  $\implies$  which approximation method should be used?
- 2. How to generate the data basis of the function evaluations? Generating all in one point will not yields good results...
- 3. Which method is feasible to minimize  $\hat{\varphi}$ ?

#### 3.4.1 Approximation Methods

#### **Response Surface Methods (RSMs)**

*Response surface methods* use simple polynomials of a low degree as the model function  $\hat{\varphi}$ , e.g.:

- Degree one (linear):  $\hat{\varphi}(\boldsymbol{p}) = \beta_0 + \beta_1^T \boldsymbol{p}$
- Degree one with mixed terms:  $\hat{\varphi}(\boldsymbol{p}) = \beta_0 + \boldsymbol{\beta}_1^T \boldsymbol{p} + \sum_i \sum_{\substack{j \neq i \\ j \neq i}} \beta_2^{i,j} p_i p_j$
- Degree two (quadratic):  $\hat{\varphi}(\boldsymbol{p}) = \beta_0 + \beta_1 \boldsymbol{p} + \boldsymbol{p}^T \beta_2 \boldsymbol{p}$
- Higher degree: ...

The unknown parameters  $\beta_1 \in \mathbb{R}$ ,  $\beta_2 \in \mathbb{R}^{n_p}$  and  $\beta_2 \in \mathbb{R}^{n_p \times n_p}$  can be approximated using least squares:

$$\min_{eta_1,m{eta}_2,m{eta}_2}\sum_i \left(arphi(m{p}^{(i)}) - \hat{arphi}(m{p}^{(i)})
ight)$$

- Advantage: Simple and the approximations are easy to compute.
- Disadvantage: The RSMs cover only the global behavior and not local accuracy.

#### **Radial Basis Functions (RBFs)**

Now the surrogate function  $\hat{\varphi}$  uses a linear combination of *radial basis functions*:

$$\hat{\varphi}(\boldsymbol{p}) = \sum_{i=1}^{m} \gamma_i h\Big( \left\| \boldsymbol{p} - \boldsymbol{p}^{(i)} \right\| \Big)$$

with basis function  $h(\cdot)$  based only on the euclidean distance from the interpolation point, e.g.

- Linear:  $h(r_i) = r_i$
- Cubic:  $h(r_i) = r_i^3$
- Thin-Plate:  $h(r_i) = r_i^2 \log r$

where  $r_i = \| p - p^{(i)} \|$ .

A suitable combination of RSM and RBF yield cubic spline-approximation:

$$\hat{\varphi}(\boldsymbol{p}) = eta_0 + eta_1 \boldsymbol{p} + \sum_i \gamma_i h(\boldsymbol{p}), \quad eta_1 \in \mathbb{R}, eta_2 \in \mathbb{R}^{n_p}, \gamma_i \in \mathbb{R}$$

- Univariate:  $h(\mathbf{p}) = \frac{1}{12} \sum_{i} r_i^3$
- Bivariate:  $h(\mathbf{p}) = \frac{1}{16\pi} \sum_{i} r_i^2 \log r_i$

#### **Design and Analysis of Computer Experiments (DACE)**

Assuming the model function is a realization of a stochastic process

$$\hat{\varphi}(\boldsymbol{p}) = \boldsymbol{v}^T(\boldsymbol{p})\boldsymbol{\beta} + Z(\boldsymbol{p})$$

where v(p) is a vector of basis functions (e.g. RSM, RBF) and Z(p) is a stationary random variable that is Gaussian distributed with zero mean. The covariance between two points  $p^{(l)}$  and  $p^{(k)}$  is given as

$$\operatorname{Cov}\left[Z(\boldsymbol{p}^{(l)}), Z(\boldsymbol{p}^{(k)})\right] = \sigma^2 R(\boldsymbol{p}^{(l)}, \boldsymbol{p}^{(k)}) \quad \text{with} \quad R(\boldsymbol{p}^{(l)}, \boldsymbol{p}^{(k)}) = \prod_{i=1}^{n_p} e^{-\theta_i d_i^2}, \quad d_i = \|\boldsymbol{p}_i^{(l)} - \boldsymbol{p}_i^{(k)}\|_2$$

Die unknown parameters  $\beta$ ,  $\theta$  and  $\sigma^2$  are then estimated using statistical estimators, e.g. maximum likelihood.

#### 3.4.2 Select of the Sampling Points

#### Design of Experiments (DoE)

The classical strategy for selection sampling points, *design of experiments* is mainly designed for physical experiments, not for deterministic ones! Typical approach:

- Classical selection
  - orthogonal arrays
  - latin hypercubes
  - combinations

- Metric-bases methods
  - MiniMax: minimizing the maximal distance between the sampling points
  - MaxiMin: maximizing the minimal distance
- Stochastic selection for Gaussian processes
  - Entropy Design or D-Opt: maximizing the determinant of the covariance matrix
  - A-Opt: depends on the trace of the covariance matrix
  - G-Opt: minimize the maximum mean squared error

#### 3.4.3 Minimizing the Surrogate Function

To successfully minimize the original objective function, the data basis of the surrogate function has to be expanded sequentially. The following sections describe two methods for this, the Strawman and the Shoemaker method.

#### Strawman

- 1. Calculate the current minimum of the surrogate function (e.g. using gradient descent).
- 2. Add the minimum of the surrogate function as a sampling point.
- 3. Calculate the new surrogate function. Repeat.

#### Shoemaker

- 1. Calculate the minimum with a minimal distance to all sampling points.
- 2. Extend the sampling points by this point.
- 3. Determine the new surrogate function. Repeat.

#### **DACE-Based, Sequential Update Strategy**

- The the expected mean error as a criteria for the quality of the surrogate function.
- Weigh small function values and uncertainties in the approximations.
- There exist different strategies following this basic idea.
- Termination criteria:
  - Number of function evaluations
  - No more improvements in the objective function
  - $\varphi(p)$  close to the minimal possible function value
- Sequential methods are better on normal computers, for parallel computers a special scheme should be used.

#### 3.4.4 Discussion

- Independent of the approximation method, the surrogate function has to be minimized one or more times per iteration (depending on the actual method).
- But the effort for the minimization is negligible as one simulation run for the evaluation of  $\varphi$  often takes a lot longer.
- Every method of **??** can be used for minimizing  $\hat{\varphi}$  as the gradients are available by design.
- Advantages:
  - $\hat{\varphi}$  is given in closed form as well as the gradientd.
  - Easy to compute, Newton-type methods applicable!
  - "Smooth" surrogate function.
- Disadvantages:
  - Approximation accuracy is limited.
  - The number of sample points rises a lot for high-dimensional problems (curse of dimensionality).

## 3.5 Comparison

#### 3.5.1 Magnetic Bearing Design

#### 3.5.2 Walking Optimization of a Humanoid Robot

#### 3.6 Discussion

- Advantages:
  - Application is easy, no or little prior knowledge required.
  - Robust toward discontinuities of  $\varphi$  or  $\nabla \varphi$ .
  - No calculation of gradients necessary.
  - No need to start "close to" a solution.
  - Some methods (e.g. evolutionary algorithms) are highly parallelizable, some are not (e.g. Nelder-Mead).
- Disadvantages:
  - Slow convergence, high number of steps needed and high computation time due to many  $\varphi$ -evaluations.
  - Inefficient for large  $n_p$ .
  - Major difficulties for nonlinear constraints in *p*.

## **4** Gradient-Based Optimization with Constraints

This chapter covers the optimization problems with nonlinear equality- and inequality-constraints:

$$egin{aligned} \min_{oldsymbol{p}\in\mathbb{R}^{n_p}}arphi(oldsymbol{p})\ \mathrm{subject \ to} \quad oldsymbol{a}(oldsymbol{p}) = oldsymbol{0}, \quad oldsymbol{a}:\mathbb{R}^{n_p}
ightarrow\mathbb{R}^{n_a}\ oldsymbol{b}(oldsymbol{p}) \geq oldsymbol{0}, \quad oldsymbol{b}:\mathbb{R}^{n_p}
ightarrow\mathbb{R}^{n_b} \end{aligned}$$

Such an optimization problem is called a nonlinear programming problem (NLP).

- A point  $p \in \mathbb{R}^{n_p}$  that fulfills the constraints is called a *feasible point*.
- The set of all feasibly points is called the *feasible region*.
- For a local minimum  $p^*$  it holds that  $\varphi(p^*) \leq \varphi(p)$  for all feasible points p in a neighborhood of  $p^*$ .
- All constraints that are *active* at a point *p* make up the *active set*. This set includes all equality constraints and the active inequality constraints A(*p*) := { *j* ∈ N : 1 ≤ *j* ≤ n<sub>b</sub>, b<sub>j</sub>(*p*) = 0 }.

#### 4.1 Solution Characterization

By a geometric view it is clear that the gradient of the active constraints must be parallel to the gradient of the objective function, i.e. there has to be a constant  $\mu_i^*$  for each active constraint  $a_i$  such that

$$\boldsymbol{\nabla}\varphi(\boldsymbol{p}^*) = \mu_i^* \boldsymbol{\nabla} a_i(\boldsymbol{p}^*) \quad \iff \quad \boldsymbol{\nabla}\varphi(\boldsymbol{p}^*) - \mu_i^* \boldsymbol{\nabla} a_i(\boldsymbol{p}^*) = \boldsymbol{0}$$

This leads directly leads to the definition of the Lagrangian

$$L(\boldsymbol{p}, \boldsymbol{\mu}, \boldsymbol{\sigma}) = \varphi(\boldsymbol{p}) - \boldsymbol{\mu}^T \boldsymbol{a}(\boldsymbol{p}) - \boldsymbol{\sigma}^T \boldsymbol{b}(\boldsymbol{p})$$

with the Lagrange multiplier  $\mu$ ,  $\sigma$  which encodes the insight above.

For formulating the *first-order optimality conditions* analogous to the ones for unconstrained optimization (see subsection 2.1.2), a *constraint qualification* must hold: The gradients

 $\nabla a_1(p^*), \cdots, \nabla a_{n_p}(p^*)$  and  $\nabla b_j(p^*), j \in A(p^*)$ 

of the active constraints have to linearly independent.

#### 4.1.1 First-Order Necessary Optimality Conditions (Karush-Kuhn-Tucker Conditions, KKT)

Let  $\varphi : \mathbb{R}^{n_p} \to \mathbb{R}$ ,  $a : \mathbb{R}^{n_p} \to \mathbb{R}^{n_a}$  and  $b : \mathbb{R}^{n_p} \to \mathbb{R}^{n_b}$  be continuously differentiable and let the constraint qualification be fulfilled. If p is a feasible local minimum of the NLP, then there exist Lagrange multiplier  $\mu \in \mathbb{R}^{n_a}$ ,  $\sigma \in \mathbb{R}^{n_b}$  such that the *Karush-Kuhn-Tucker* conditions hold:

$$\nabla_{\boldsymbol{p}} L(\boldsymbol{p}, \boldsymbol{\mu}, \boldsymbol{\sigma}) = \boldsymbol{0} \tag{KKT.i}$$

$$u^{r} a(p) = 0 \tag{KK1.11a}$$

$$\boldsymbol{\sigma}^T \boldsymbol{b}(\boldsymbol{p}) = \mathbf{0}$$
 (KKT.iib)

$$\sigma \ge 0$$
 (KKT.iii)

Here,  $L(\boldsymbol{p}, \boldsymbol{\mu}, \boldsymbol{\sigma})$  is the Lagrangian

$$L(\boldsymbol{p}, \boldsymbol{\mu}, \boldsymbol{\sigma}) \coloneqq \varphi(\boldsymbol{p}) - \boldsymbol{\mu}^T \boldsymbol{a}(\boldsymbol{p}) - \boldsymbol{\sigma}^T \boldsymbol{b}(\boldsymbol{p})$$

and the inequality in (KKT.iii) is element-wise.

#### 4.1.2 Second-Order Necessary Optimality Conditions

The second-order necessary optimality condition is fulfilled if the first-order condition is fulfilled and the Hessian of the Lagrangian has a positive curvature *along the feasible directions*:

$$oldsymbol{z}^T oldsymbol{H}^{oldsymbol{p}}_L(oldsymbol{p},oldsymbol{\mu},\sigma)oldsymbol{z} \geq 0$$
  
for all  $oldsymbol{z} \in \mathbb{R}^{n_p} \setminus \{\mathbf{0}\}$  with  $ig(oldsymbol{z}^T \cdot oldsymbol{
abla} a_i(oldsymbol{p}) = oldsymbol{0}ig)_{i=1,\cdots,n_a}$ , and  $ig(oldsymbol{z}^T \cdot oldsymbol{
abla} b_j(oldsymbol{p}) = oldsymbol{0}ig)_{j\in A(oldsymbol{p})}$ 

#### 4.1.3 Example

## 4.2 Simple Bounds, Box Constraints

The most common type of inequality-constraints are simple lower and upper bounds on the optimization variables:

$$\min_{oldsymbol{p} \in \mathbb{R}^{n_p}} arphi(oldsymbol{p})$$
  
subject to  $oldsymbol{p}_{i,\min} \leq oldsymbol{p}_i \leq oldsymbol{p}_{i,\max}$ 

In practice, these are not only used in nearly all constrained, but also in unconstrained optimization problems as they might increase the efficiency by constraining the search space. They also increase robustness of optimization techniques that ensure the fulfilling of the constrains in every iteration.

Box constraints can also be used to prohibit "insecure values", e.g. negative variables when the objective function takes square roots.

For gradient-free methods, the constraints can be ensured by clipping the new iteration value to the bounds:

 $p_i^{(k+1)} = \max\left\{p_{i,\min}, \min\left\{p_{i,\max}, p_i^{(k+1)}\right\}\right\}$ 

For gradient-based methods, box constraints can be considered like every other linear and nonlinear constraint (more on that later).

## 4.3 Penalty Function

The approach of *penalty functions* is to transform the constrained problem to an unconstrained problem by punishing violations of the constraints. As this arises big problems with lots of and highly nonlinear constraints, penalty functions are rarely used in practice anymore in favor of sequential quadratic programming (see section 4.5).

As of today, penalty functions are mainly used for step size determination for nonlinear programs and in the application of robust, gradient-free methods in the unconstrained optimization for solving NLPs.

#### 4.3.1 Exterior Penalty Functions

The original nonlinear problem is replaced by an unconstrained problem with a penalty function  $\Phi$ 

$$\min_{oldsymbol{p}\in\mathbb{R}^{n_p}}\Phi(oldsymbol{p},
ho),\quad \Phi(oldsymbol{p},
ho)=arphi(oldsymbol{p})+
ho\sum_j\pi_j(oldsymbol{p}),\quad
ho\in\mathbb{R}^+$$

with a function  $\pi_i$  per constraint that is positive if the constraint is violated and zero otherwise.

Often a series of unconstrained optimization problems  $\Phi(\mathbf{p}, \rho)$  is solved with an increasing  $\rho$ , such that the solution  $\mathbf{p}^*(\rho)$  gets pushed into the feasible region of the NLP step by step in in the hope of

$$\lim_{\rho \to \infty} \boldsymbol{p}^*(\rho) = \boldsymbol{p}^*$$

where  $p^*$  is the real minimum.

This approach is called *exterior penalty functions* as the penalty term is only relevant if p violates the according constraints.

#### **Quadratic Penalty Function**

The *quadratic penalty function* is the most common exterior penalty function:

$$\Phi_Q(\boldsymbol{p}, \rho) = \varphi(\boldsymbol{p}) + \rho \cdot \frac{1}{2} \left( \sum_{k=1}^{n_a} \left( a_k(\boldsymbol{p}) \right)^2 + \sum_{jk=1}^{n_b} \left( \hat{b}_k(\boldsymbol{p}) \right)^2 \right)$$

Here,  $\hat{b}_i(\mathbf{p})$  denotes the violation of the inequality constraints:

$$\hat{b}_j(oldsymbol{p}) = egin{cases} 0 & ext{iff} \ b_j(oldsymbol{p}) \geq 0 \ -b_j(oldsymbol{p}) & ext{iff} \ b_j(oldsymbol{p}) < 0 \ \end{cases} = ig|\min\{0, \ b_j(oldsymbol{p})\}ig|$$

But in the case of  $\rho \to \infty$ , the Hessian of  $\Phi_C$  gets more and more ill-condition and may even be singular in the limit.

#### Example

#### 4.3.2 Interior Penalty Functions

When using exterior penalty functions, it is not guaranteed that every solution of the iterating unconstrained problems fulfills the constraints of the original NLP. Therefore, exterior penalty functions are not applicable if fulfilling the constraints in every iteration is required. *Interior penalty functions* guarantee exactly that: Every subproblem yields a feasible solution.

Given a NLP with only inequality constraints, interior penalty functions use *barrier functions* that have the following properties:

- Value of infinity everywhere except inside the feasible region.
- Continuously differentiable inside the feasible region.
- The values go to  $+\infty$  as p gets closer to the edge of the feasible region.

#### **Logarithmic Barrier Function**

The *logarithmic barrier function* is the most used interior penalty function:

$$\Phi_B(\boldsymbol{p},r) = \varphi(\boldsymbol{p}) - r \sum_{k=1}^{n_b} \ln b_i(\boldsymbol{p})$$

The barrier parameter r > 0 will be decreased step by step and the solution should converge to the minimum as  $r \to 0$ .

But in the case of  $r \to 0$ , the Hessian of  $\Phi_B$  gets more and more ill-condition and may even be singular in the limit.

#### Example

#### 4.3.3 Exact Penalty Functions

Neither the quadratic penalty function nor the logarithmic barrier function are "exact" penalty function so that, with an appropriate  $\rho$  or respectively r, the solution of the unconstrained NLP yields the exact solution.

One important penalty function of this class is the *exact*  $\ell_1$ -*penalty function*:

$$\Phi_{\ell_1}(\boldsymbol{p},\rho) = \varphi(\boldsymbol{p}) + \rho \sum_{k=1}^{n_a} \left| a_k(\boldsymbol{p}) \right| + \rho \sum_{k=1}^{n_b} \left| \min\{0, b_j(\boldsymbol{p})\} \right|$$

However, the  $\ell_1$ -penalty function is not differentiable! This make the numerical solution difficult. But the minimization yields, for adequate big  $\rho$ , the minimum of the NLP!

#### Example 1

#### Example 2

#### 4.3.4 Augmented Lagrangian

The downsides of exterior, interior and the  $\ell_1$ -penalty function can be avoided by not using the objective directly, but by using the Lagrangian:

$$L_Q(\boldsymbol{p}, \boldsymbol{\mu}, \boldsymbol{\sigma}, \rho) = \varphi(\boldsymbol{p}) - \boldsymbol{\mu}^T \boldsymbol{a}(\boldsymbol{p}) - \boldsymbol{\sigma}^T \boldsymbol{b}(\boldsymbol{p}) + \rho \cdot \frac{1}{2} \left( \sum_{k=1}^{n_a} \left( a_k(\boldsymbol{p}) \right)^2 + \sum_{jk=1}^{n_b} \left( \min\{0, b_j(\boldsymbol{p})\} \right)^2 \right)$$

But this requires a good approximation of the Lagrange multiplier  $\mu$  and  $\sigma$ .

#### Example 1

#### Example 2

#### Notes

- In practice, the augmented Lagrangian  $L_Q$  is used in the following fashion:
  - 1. Choose Lagrange multipliers  $\mu$ ,  $\sigma$  and the parameter  $\rho > 0$ .
  - 2. Calculate a local minimum  $p^*(\rho)$  of  $L_Q$  using methods of the unconstrained optimization.
  - 3. Update the Lagrange multipliers and  $\rho$ . Repeat with 2.
- The update of the Lagrange multipliers according to μ(ρ) = −ρa(p\*(ρ)) is based on the convergence properties of the quadratic penalty function.
- If
- $p^*$  is a minimum of the NLP,
- the constraint qualification hold and
- the KKT-conditions and the second-order sufficient optimality condition is fulfilled for  $\mu^*, \sigma^*,$

then

- exists a threshold  $\hat{\rho}$  for which it holds that for all  $\rho \geq \hat{\rho}$  it holds that
- $p^*$  is a strict local minimum of the (quadratic) augmented Lagrangian  $L_Q(p, \mu^*, \sigma^*, \rho)$ .

#### 4.4 Constraint Elimination

For NLPs with active constraints it is often tempting to transform the NLP to remove some of the constraints and to reduce the number of optimization variables (the degrees of freedom). But this might yields wrong results! It has to be assured that:

- The original minimum is not eliminated.
- The nonlinearity of the problem is not increased too much (causing problems in finite difference approximations of derivatives).
- No singularities arise in the transformed objective.
- No new discontinuities and non-differentiabilities are added to the objective.

- The Hessian of the surrogate problem is not singular or ill-conditioned near the minimum.
- The transformed problem does not have additional local minima or stationary points.
- And a lot more.

In general: Keep more constraints and keep the function as linear as possible instead of applying transformations that behave badly.

#### Example 1

Example 2

#### Example 3

## 4.5 Sequential Quadratic Programming (SQP)

To take care of highly nonlinear equality and inequality constraints, information about the trajectory of these have to be considered, i.e. gradient and Hessian information. Assuming that the constraints that are active on the solution are known, the optimization problem is given as:

$$\begin{split} \min_{\boldsymbol{p} \in \mathbb{R}^{n_p}} \varphi(\boldsymbol{p}) \\ \text{subject to} \quad \boldsymbol{a}(\boldsymbol{p}) = \boldsymbol{0}, \quad \boldsymbol{a} : \mathbb{R}^{n_p} \to \mathbb{R}^{n_a} \end{split}$$

The KKT-conditions are then equivalent to the simpler formulation

$$\boldsymbol{\nabla} L(\boldsymbol{p}, \boldsymbol{\mu}) \coloneqq \begin{bmatrix} \boldsymbol{\nabla}_{\boldsymbol{p}} L(\boldsymbol{p}, \boldsymbol{\mu}) \\ \boldsymbol{\nabla}_{\boldsymbol{\mu}} L(\boldsymbol{p}, \boldsymbol{\mu}) \end{bmatrix} = \begin{bmatrix} \boldsymbol{\nabla} \varphi(\boldsymbol{p}) - \sum_{k=1}^{n_a} \mu_k \cdot \boldsymbol{\nabla} a_k(\boldsymbol{p}) \\ -\boldsymbol{a}(\boldsymbol{p}) \end{bmatrix} = \boldsymbol{0}$$

with the Lagrangian

$$L(\boldsymbol{p}, \boldsymbol{\mu}) = \varphi(\boldsymbol{p}) - \boldsymbol{\mu}^T \boldsymbol{a}(\boldsymbol{p})$$

This yields a system of  $n_p + n_a$  nonlinear equations for  $n_p + n_a$  unknowns p,  $\mu$ .

Taylor-expanding the gradient of the Lagrangian around  $(m{p}^{(k)},m{\mu}^{(k)})$  yields

$$\boldsymbol{\nabla}L(\boldsymbol{p}^*,\boldsymbol{\mu}^*) \stackrel{T\left(\boldsymbol{p}^{(k)},\boldsymbol{\mu}^{(k)}\right)}{=} \boldsymbol{\nabla}L\left(\boldsymbol{p}^{(k)},\boldsymbol{\mu}^{(k)}\right) + \boldsymbol{H}_L\left(\boldsymbol{p}^{(k)},\boldsymbol{\mu}^{(k)}\right) \begin{bmatrix} \boldsymbol{d}_p^{(k)} \\ \boldsymbol{d}_\mu^{(k)} \end{bmatrix} + \cdots \stackrel{!}{=} 0$$

with  $d_p^{(k)} \coloneqq p^* - p^{(k)}$  and  $d_{\mu}^{(k)} \coloneqq \mu^* - \mu^{(k)}$ . Cutting off the higher-order terms yields the *Lagrange-Newton* method where the search direction  $d_p^{(k)}$ ,  $d_{\mu}^{(k)}$  is given as the solution of

$$\boldsymbol{H}_{L}(\boldsymbol{p}^{(k)},\boldsymbol{\mu}^{(k)}) \begin{bmatrix} \boldsymbol{d}_{p}^{(k)} \\ \boldsymbol{d}_{\mu}^{(k)} \end{bmatrix} = -\boldsymbol{\nabla}L(\boldsymbol{p}^{(k)},\boldsymbol{\mu}^{(k)})$$
(4.1)

The Iteration equation  $k \rightarrow k + 1$  is analogous to the Newton method:

$$\begin{bmatrix} \boldsymbol{p}^{(k+1)} \\ \boldsymbol{\mu}^{(k+1)} \end{bmatrix} = \begin{bmatrix} \boldsymbol{p}^{(k)} \\ \boldsymbol{\mu}^{(k)} \end{bmatrix} + \begin{bmatrix} \boldsymbol{d}^{(k)}_p \\ \boldsymbol{d}^{(k)}_\mu \end{bmatrix}$$

But there is one catch: The set of active constraints at the minimum is generally not known and can change in every iteration.
## 4.5.1 Finding the Search Direction

Further analysis of the linear system (4.1) with the Lagrangian

$$L(\boldsymbol{p}, \boldsymbol{\mu}) = \varphi(\boldsymbol{p}) - \boldsymbol{\mu}^T \boldsymbol{a}(\boldsymbol{p})$$

yields that the linear system has the following structure:

$$\begin{bmatrix} \boldsymbol{H}_{L}^{\boldsymbol{p}}(\boldsymbol{p}^{(k)},\boldsymbol{\mu}^{(k)}) & -\boldsymbol{J}_{\boldsymbol{a}}(\boldsymbol{p}^{(k)}) \\ -\boldsymbol{J}_{\boldsymbol{a}}^{T}(\boldsymbol{p}^{(k)}) & \boldsymbol{0} \end{bmatrix} \begin{bmatrix} \boldsymbol{d}_{p}^{(k)} \\ \boldsymbol{d}_{\mu}^{(k)} \end{bmatrix} = \begin{bmatrix} -\boldsymbol{\nabla}\varphi(\boldsymbol{p}^{(k)}) + \boldsymbol{J}_{\boldsymbol{a}}(\boldsymbol{p}^{(k)})\boldsymbol{\mu}^{(k)} \\ \boldsymbol{a}(\boldsymbol{p}^{(k)}) \end{bmatrix}$$

Where  $H_L^p(p^{(k)})$  is the Hessian of the Lagrangian w.r.t. p and  $J_a(p^{(k)})$  is the Jacobian

$$oldsymbol{J}_{oldsymbol{a}}oldsymbol{p}oldsymbol{a}^{(k)}ig) = egin{bmatrix} rac{\partial a_1}{\partial p_1} & \cdots & rac{\partial a_{n_a}}{\partial p_1} \ dots & \ddots & dots \ rac{\partial a_1}{\partial p_{n_p}} & \cdots & rac{\partial a_{n_a}}{\partial p_{n_p}} \end{bmatrix}$$

Note that this definition of the Jacobian differs from the usual definition! This one has the gradients of the function has columns!

This linear system can be transformed to

$$\begin{bmatrix} \boldsymbol{H}_{L}^{\boldsymbol{p}}(\boldsymbol{p}^{(k)},\boldsymbol{\mu}^{(k)}) & -\boldsymbol{J}_{\boldsymbol{a}}(\boldsymbol{p}^{(k)}) \\ -\boldsymbol{J}_{\boldsymbol{a}}^{T}(\boldsymbol{p}^{(k)}) & \boldsymbol{0} \end{bmatrix} \begin{bmatrix} \boldsymbol{d}_{p}^{(k)} \\ \boldsymbol{d}_{\mu}^{(k)} + \boldsymbol{\mu}^{(k)} \\ \boldsymbol{\mu}^{(k+1)} \end{bmatrix} = \begin{bmatrix} -\boldsymbol{\nabla}\varphi(\boldsymbol{p}^{(k)}) \\ \boldsymbol{a}(\boldsymbol{p}^{(k)}) \end{bmatrix}$$

where  $d_p^{(k)}$  can be viewed as the solution of a quadratic minimization problem!

#### **Quadratic Problem (QP)**

The said quadratic problem is given as:

$$\min_{oldsymbol{d}_p \in \mathbb{R}^{n_p}} arphi(oldsymbol{p}^{(k)}) + \left( oldsymbol{
abla} arphi(oldsymbol{p}^{(k)}) 
ight)^T oldsymbol{d}_p + rac{1}{2} oldsymbol{d}_p^T oldsymbol{H}_L^{oldsymbol{p}}(oldsymbol{p}^{(k)},oldsymbol{\mu}^{(k)}) oldsymbol{d}_p$$
  
subject to  $oldsymbol{a}(oldsymbol{p}^{(k)}) + oldsymbol{J}_a^T(oldsymbol{p}^{(k)}) oldsymbol{d}_p = oldsymbol{0}$ 

Where the quadratic objective function of the QP consists of a quadratic Taylor-approximation of the NLP objective plus a weighted curvature condition via the Hessian of the activate conditions:

$$oldsymbol{H}_L^{oldsymbol{p}}oldsymbol{\left(p^{(k)},oldsymbol{\mu}^{(k)}
ight)} = oldsymbol{H}_arphioldsymbol{\left(p^{(k)}
ight)} - \sum_{i=1}^{n_a} \mu_i^Toldsymbol{H}_{a_i}oldsymbol{\left(p^{(k)}
ight)}$$

The linear constraints of the QP also consist of Taylor-approximations of the active NLP constraints.

Solving this quadratic optimization problem is more robust and possibly faster than solving the linear system of equations. There exist special algorithms for solving QPs. But even though they are not active, the inequality constraints must be fulfilled. Therefore, the QP is extended to also determine the Lagrange multipliers  $\mu^{(k)}$ ,  $\sigma^{(k)}$  along with the search direction  $d_p^{(k)}$ :

$$\min_{\boldsymbol{d}_{p} \in \mathbb{R}^{n_{p}}} \varphi(\boldsymbol{p}^{(k)}) + \left(\boldsymbol{\nabla}\varphi(\boldsymbol{p}^{(k)})\right)^{T} \boldsymbol{d}_{p} + \frac{1}{2} \boldsymbol{d}_{p}^{T} \boldsymbol{H}_{L}^{\boldsymbol{p}}(\boldsymbol{p}^{(k)}, \boldsymbol{\mu}^{(k)}, \boldsymbol{\sigma}^{(k)}) \boldsymbol{d}_{p}$$
subject to
$$\boldsymbol{a}(\boldsymbol{p}^{(k)}) + \boldsymbol{J}_{\boldsymbol{a}}^{T}(\boldsymbol{p}^{(k)}) \boldsymbol{d}_{p} = \boldsymbol{0}$$

$$\boldsymbol{b}(\boldsymbol{p}^{(k)}) + \boldsymbol{J}_{\boldsymbol{b}}^{T}(\boldsymbol{p}^{(k)}) \boldsymbol{d}_{p} \geq \boldsymbol{0}$$
(4.2)

Iteration step  $k \to k + 1$ :  $p^{(k+1)} = p^{(k)} + d_p^{(k)}$ 

#### Notes

- Similar to the Newton method, it can be shown under some assumptions that the SQP method converges to a local minimum.
- Algorithms for solving general QPs can be used to determine the active constraints in each iteration by solving the QP.
   Even though it is not needed, it is useful for efficiency to use as much information as possible from the last iteration ("hot start").
- Other approaches determine the active constraints (working set) outside of the QP solution to only solve QPs with equality constraints which is especially efficient.

### 4.5.2 Step Size Rules

If the initialization  $p^{(0)}$  is "far" away from the minimum (or the QP is a bad, local approximation of the NLP), the convergence can be improved by determining the optimal step size for

$$\boldsymbol{p}^{(k+1)} = \boldsymbol{p}^{(k)} + \alpha^{(k)} \boldsymbol{d}_p^{(k)} \quad \text{or respectively} \quad \begin{bmatrix} \boldsymbol{p}^{(k+1)} \\ \boldsymbol{\mu}^{(k+1)} \end{bmatrix} = \begin{bmatrix} \boldsymbol{p}^{(k)} \\ \boldsymbol{\mu}^{(k)} \end{bmatrix} + \alpha^{(k)} \begin{vmatrix} \boldsymbol{d}_p^{(k)} \\ \boldsymbol{\mu}_{QP}^{(k+1)} - \boldsymbol{\mu}^{(k)} \end{vmatrix}$$

Common methods for determining the step size are

- the (quadratic) augmented Lagrangian  $L_Q$  or
- the exact  $\ell_1$ -penalty function  $\Phi_{\ell_1}$

with step size rules like the Armijo rule similar as in the unconstrained optimization (see section 2.3).

## 4.5.3 Approximation of the Lagrange Multipliers

If the approximation  $p^{(k)}$  is far away from the NLP solution, it is not useful to use the Lagrange multipliers of the QP for the NLP (as the linearization is only valid locally). Another method is to approximate the Lagrange multipliers after calculating  $p^{(k+1)}$  (i.e. solving the QP) using minimum least squares:

$$\min_{\boldsymbol{\mu}\in\mathbb{R}^{n_a}} \left\| \boldsymbol{\nabla}\varphi(\boldsymbol{p}^{(k+1)}) - \sum_{i=1}^{n_a} \mu_i^{(k+1)} \cdot \boldsymbol{\nabla} a_i(\boldsymbol{p}^{(k+1)}) \right\|_2^2$$

or equivalently model the active inequality constraints explicitly:

$$\min_{\boldsymbol{\mu}\in\mathbb{R}^{n_a}} \left\|\boldsymbol{\nabla}\varphi(\boldsymbol{p}^{(k+1)}) - \sum_{i=1}^{n_a} \mu_i^{(k+1)} \cdot \boldsymbol{\nabla}a_i(\boldsymbol{p}^{(k+1)}) - \sum_{i\in A(\boldsymbol{p}^{(k+1)})} \sigma_i^{(k+1)} \cdot \boldsymbol{\nabla}a_i(\boldsymbol{p}^{(k+1)})\right\|_2^2$$

Solving least squares problem will be discussed in detail in chapter 6.

## 4.5.4 Termination Criteria

An obvious termination criteria would be to check whether the necessary KKT-conditions are sufficiently fulfilled, i.a.  $\nabla L(p, \mu) \approx 0$ .

The commonly used SQP method NPSOL terminates if all of the following criteria are fulfilled:

1. Old and new approximation do not change any more:

$$\left\|\boldsymbol{p}^{(k+1)} - \boldsymbol{p}^{(k)}\right\| = \alpha^{(k)} \cdot \left\|\boldsymbol{d}_{p}^{(k)}\right\|_{2} \leq \sqrt{\varepsilon_{\text{opt}}} \left(1 + \left\|\boldsymbol{p}^{(k+1)}\right\|_{2}\right)$$

2. Gradient of the objective function that is projected onto the active constraints vanishes:

$$\left\| \boldsymbol{Z}^{T} \cdot \boldsymbol{\nabla} \varphi(\boldsymbol{p}^{(k+1)}) \right\|_{2} \leq \sqrt{\varepsilon_{\text{opt}}} \left( 1 + \max\left\{ 1 + \left| \varphi(\boldsymbol{p}^{(k+1)}) \right|, \left\| \boldsymbol{\nabla} \varphi(\boldsymbol{p}^{(k+1)}) \right\|_{2} \right\} \right)$$

3. Constraints are sufficiently fulfilled:

$$\begin{aligned} \left| a_i(\boldsymbol{p}^{(k+1)}) \right| &\leq \varepsilon_{\text{ft}}, \quad i = 1, \cdots, n_a \\ b_j(\boldsymbol{p}^{(k+1)}) &\geq -\varepsilon_{\text{ft}}, \quad i = 1, \cdots, n_b \end{aligned}$$

Where the constraints  $\varepsilon_{\rm opt},\,\varepsilon_{\rm ft}$  have to be chosen by the user.

In the more modern method SNOPT, all of the following two criteria have to be fulfilled:

$$\frac{\max_{i=1,\cdots,n_{p}}\left\{\left|\frac{\partial\varphi\left(\boldsymbol{p}^{(k)}\right)}{\partial p_{i}}-\sum_{j=1}^{n_{a}}\mu_{j}^{(k)}\cdot\frac{\partial a_{j}\left(\boldsymbol{p}^{(k)}\right)}{\partial p_{i}}-\sum_{l\in A\left(\boldsymbol{p}^{(k)}\right)}\sigma_{j}^{(k)}\cdot\frac{\partial b_{l}\left(\boldsymbol{p}^{(k)}\right)}{\partial p_{i}}\right|\right\}}{\sqrt{\sum_{j=1}^{n_{a}}\left(\mu_{j}^{(k)}\right)^{2}+\sum_{l\in A\left(\boldsymbol{p}^{(k)}\right)}\left(\sigma_{j}^{(k)}\right)^{2}}}$$
$$\frac{\max\left\{\left|a_{j}\left(\boldsymbol{p}^{(k)}\right)\right|:i=1,\cdots,n_{p}\right\}\cup\left\{\left|\min\left\{0,\ b_{l}\left(\boldsymbol{p}^{(k)}\right)\right\}\right|:l\in A\left(\boldsymbol{p}^{(k)}\right)\right\}\right\}}{\left\|\boldsymbol{p}^{(k)}\right\|_{2}}\leq\varepsilon_{\mathrm{ft}}$$

A common criteria to detect failures is a maximum number of iterations  $k_{\text{max}}$ .

#### 4.5.5 Hessian Approximation

The quadratic problem (4.2) needs the  $(n_p \times n_p)$ -dimensional Hessian of the Lagrangian

$$\boldsymbol{H}_{L}^{\boldsymbol{p}}(\boldsymbol{p}^{(k)}, \boldsymbol{\mu}^{(k)}, \boldsymbol{\sigma}^{(k)}) = \boldsymbol{H}_{\varphi}(\boldsymbol{p}^{(k)}) - \sum_{i=1}^{n_{a}} \mu_{i}^{(k)} \boldsymbol{H}_{a_{i}}(\boldsymbol{p}^{(k)}) - \sum_{i \in A(\boldsymbol{p}^{(k)})} \sigma_{i}^{(k)} \boldsymbol{H}_{b_{i}}(\boldsymbol{p}^{(k)})$$

for (theoretical) quadratic convergence of the SQP method. But in practice, the second-order derivatives (the Hessian) is often not available! Additionally it is assumed that the Hessian has a positive curvature (i.e. is positive definite) along all feasible directions, which is fulfilled near a strict minimum. However, this cannot be assumed in every iteration.

Hence, approximations/modifications of the Hessian are required.

#### Naïve Approach: BFGS Approximation

It is tempting to use the BFGS update for the Hessian that is really successful in the unconstrained optimization. The update rule is given as:

$$\tilde{\boldsymbol{H}}^{(k+1)} = \tilde{\boldsymbol{H}}^{(k)} - \frac{1}{\left(\boldsymbol{d}^{(k)}\right)^T \tilde{\boldsymbol{H}}^{(k)} \boldsymbol{d}^{(k)}} \tilde{\boldsymbol{H}}^{(k)} \boldsymbol{d}^{(k)} \left(\tilde{\boldsymbol{H}}^{(k)} \boldsymbol{d}^{(k)}\right)^T + \frac{1}{\left(\boldsymbol{g}^{(k)}\right)^T \boldsymbol{d}^{(k)}} \boldsymbol{g}^{(k)} \left(\boldsymbol{g}^{(k)}\right)^T \\ \boldsymbol{g}^{(k)} = \boldsymbol{\nabla}_{\boldsymbol{p}} L(\boldsymbol{p}^{(k+1)}, \boldsymbol{\mu}^{(k+1)}, \boldsymbol{\sigma}^{(k+1)}) - \boldsymbol{\nabla}_{\boldsymbol{p}} L(\boldsymbol{p}^{(k)}, \boldsymbol{\mu}^{(k+1)}, \boldsymbol{\sigma}^{(k+1)})$$

But it is not necessary that the full Hessian of the Lagrangian is positive definite! Additionally, it is very inefficient to calculate the full Hessian if the NLP has lots of active constraints. Hence, the approximation has to be modified in order to be useful for SQP methods.

#### **Reduced Hessian**

Every active constraints reduces the degrees of freedom by one Hence the degrees of freedom are the number of optimization variables  $n_p$  minus the number of active and linearly independent constraints. In all of the following it is assumed that it is known which constraints are active at the solution, yielding the following, simpler, view of the NLP:

$$\min_{\boldsymbol{d}_{p} \in \mathbb{R}^{n_{p}}} \varphi(\boldsymbol{p}^{(k)}) + \left(\boldsymbol{\nabla}\varphi(\boldsymbol{p}^{(k)})\right)^{T} \boldsymbol{d}_{p} + \frac{1}{2} \boldsymbol{d}_{p}^{T} \boldsymbol{H}_{L}^{\boldsymbol{p}}(\boldsymbol{p}^{(k)}, \boldsymbol{\mu}^{(k)}, \boldsymbol{\sigma}^{(k)}) \boldsymbol{d}_{p}$$
subject to  $\boldsymbol{a}(\boldsymbol{p}^{(k)}) + \boldsymbol{J}_{\boldsymbol{a}}^{T}(\boldsymbol{p}^{(k)}) \boldsymbol{d}_{p} = \boldsymbol{0}$ 
(4.3)

where the last Jacobian  $J_a^T$  might also contain active inequality constraints which are left out here for brevity. But they are handled the same as regular equality constraints and can be considered to be part of it knowing which are active (as assumed). The degrees of freedom of this NLP are therefore  $n_p - n_a$ .

Assuming the current approximation  $p^{(k)}$  fulfills the constraints,  $a(p^{(k)}) = 0$ , the next approximation also has to fulfill the constraints:

$$oldsymbol{a}(oldsymbol{p}^{(k)}+oldsymbol{d}_p)=oldsymbol{0}$$

By Taylor-expanding this equation around  $p^{(k)}$ 

$$oldsymbol{a}ig(oldsymbol{p}^{(k)}+oldsymbol{d}_pig)\stackrel{Tig(oldsymbol{p}^{(k)}ig)}{=}oldsymbol{a}ig(oldsymbol{p}^{(k)}ig)+oldsymbol{J}_{oldsymbol{a}}^Tig(oldsymbol{p}^{(k)}ig)oldsymbol{d}_p+\dots=0$$

a linear system for the search direction  $d_p$  can be found such that the new iteration is also feasible:

$$oldsymbol{J}_{oldsymbol{a}}^T(oldsymbol{p}^{(k)})oldsymbol{d}_p = oldsymbol{0}$$

Thus,  $d_p$  has to lie in the kernel of  $J_a^T$  and the dimensionality of the kernel is  $n_p - n_a$ .

Hence, the kernel is spanned by  $n_p - n_a$  basis vectors (which are not uniquely determined). Let  $Z^{(k)}$  be the matrix that contains these basis vectors as columns, then

$$\boldsymbol{J}_{\boldsymbol{a}}^{T}(\boldsymbol{p}^{(k)})\boldsymbol{Z}^{(k)}=\boldsymbol{0}$$

holds and the search direction  $d_p \in \mathbb{R}^{n_p}$  can be represented as

$$oldsymbol{d}_p = oldsymbol{Y}^{(k)}oldsymbol{d}_p^y + oldsymbol{Z}^{(k)}oldsymbol{d}_p^z$$

where  $Z^{(k)} \in \mathbb{R}^{n_p \times (n_p - n_a)}$  are the basis vectors of the kernel and  $Y^{(k)} \in \mathbb{R}^{n_p \times n_a}$  are the basis vectors of the image space of  $J_a^T(p^{(k)})$ . The vectors  $d_p^y \in \mathbb{R}^{n_a}$  and  $d_p^z \in \mathbb{R}^{n_p - n_a}$  are unknown and have to be computed to solve the QP. Plugging this formulation of  $d_p$  into the constraints of (4.3):

$$0 = a(p^{(k)}) + J_{a}^{T}(p^{(k)})d_{p}$$

$$\iff -a(p^{(k)}) = J_{a}^{T}(p^{(k)})d_{p}$$

$$= J_{a}^{T}(p^{(k)})(Y^{(k)}d_{p}^{y} + Z^{(k)}d_{p}^{z})$$

$$= J_{a}^{T}(p^{(k)})Y^{(k)}d_{p}^{y} + \underbrace{J_{a}^{T}(p^{(k)})Z^{(k)}}_{=0}d_{p}^{z}$$

$$= J_{a}^{T}(p^{(k)})Y^{(k)}d_{p}^{y} \qquad (4.4)$$

The vector  $d_p^y$  is therefore determined by the active constraints via the linear system (4.4). Now there remain  $n_p - n_a$  degrees of freedom for the actual optimization.

Plugging the formulation of  $d_p$  into the objective of the QP<sup>1</sup>

$$\begin{split} \varphi(\boldsymbol{p}^{(k)}) &+ \left(\boldsymbol{\nabla}\varphi(\boldsymbol{p}^{(k)})\right)^T \boldsymbol{d}_p + \frac{1}{2} \boldsymbol{d}_p^T \boldsymbol{H}_L^{\boldsymbol{p}} \boldsymbol{d}_p \\ &= \varphi(\boldsymbol{p}^{(k)}) + \left(\boldsymbol{\nabla}\varphi(\boldsymbol{p}^{(k)})\right)^T \left(\boldsymbol{Y}^{(k)} \boldsymbol{d}_p^y + \boldsymbol{Z}^{(k)} \boldsymbol{d}_p^z\right) + \frac{1}{2} \left(\boldsymbol{Y}^{(k)} \boldsymbol{d}_p^y + \boldsymbol{Z}^{(k)} \boldsymbol{d}_p^z\right)^T \boldsymbol{H}_L^{\boldsymbol{p}} \left(\boldsymbol{Y}^{(k)} \boldsymbol{d}_p^y + \boldsymbol{Z}^{(k)} \boldsymbol{d}_p^z\right) \end{split}$$

yields the following objective when leaving out all constant parts w.r.t.  $d_p^z$ , as that is the optimization variable:

$$\varphi(\boldsymbol{p}^{(k)}) + \left(\boldsymbol{\nabla}\varphi(\boldsymbol{p}^{(k)})\right)^{T} \boldsymbol{Z}^{(k)} \boldsymbol{d}_{p}^{z} + \left(\boldsymbol{d}_{p}^{y}\right)^{T} \left(\boldsymbol{Z}^{(k)}\right)^{T} \boldsymbol{H}_{L}^{\boldsymbol{p}} \left(\boldsymbol{Z}^{(k)} \boldsymbol{d}_{p}^{z}\right) + \frac{1}{2} \left(\boldsymbol{d}_{p}^{z}\right)^{T} \left(\boldsymbol{Z}^{(k)}\right)^{T} \boldsymbol{H}_{L}^{\boldsymbol{p}} \boldsymbol{Z}^{(k)} \boldsymbol{d}_{p}^{z}$$

The solution of this optimization problem can be computed by solving the following linear system (if the reduced Hessian is positive definite, which it is close to a strict local minimum):

$$\left(\left(\boldsymbol{Z}^{(k)}\right)^{T}\boldsymbol{H}_{L}^{(k)}\boldsymbol{Z}^{(k)}
ight)\boldsymbol{d}_{p}^{z}=-\left(\boldsymbol{Z}^{(k)}
ight)^{T}\boldsymbol{H}_{L}^{(k)}\boldsymbol{Y}^{(k)}\boldsymbol{d}_{p}^{y}-\left(\boldsymbol{Z}^{(k)}
ight)^{T}\cdot\boldsymbol{
abla}arphi(\boldsymbol{p}^{(k)})$$

**Summary** Using the representation  $d_p = Y^{(k)}d_p^y + Z^{(k)}d_p^z$ , the search direction as the solution of the QP (4.3) can be computed by solving two staggered linear systems:

$$egin{aligned} oldsymbol{J}_{oldsymbol{a}}^Toldsymbol{p}^{(k)}oldsymbol{Y}^{(k)}oldsymbol{d}_p^y &= -oldsymbol{a}oldsymbol{p}^{(k)}ig) \ igg(oldsymbol{(Z^{(k)})}^Toldsymbol{H}_L^{(k)}oldsymbol{Z^{(k)}}igg)oldsymbol{d}_p^z &= -oldsymbol{(Z^{(k)})}^Toldsymbol{H}_L^{(k)}oldsymbol{Y}^{(k)}oldsymbol{d}_p^y -oldsymbol{(Z^{(k)})}^T\cdotoldsymbol{
abla}arphioldsymbol{p}^{(k)}igg) \ oldsymbol{d}_p^z &= -oldsymbol{(Z^{(k)})}^Toldsymbol{H}_L^{(k)}oldsymbol{Y}^{(k)}oldsymbol{d}_p^y -oldsymbol{(Z^{(k)})}^T\cdotoldsymbol{
abla}arphioldsymbol{p}^{(k)}oldsymbol{d}_p^z &= -oldsymbol{(Z^{(k)})}^Toldsymbol{H}_L^{(k)}oldsymbol{Y}^{(k)}oldsymbol{d}_p^y -oldsymbol{(Z^{(k)})}^T\cdotoldsymbol{
abla}arphioldsymbol{p}^{(k)}oldsymbol{d}_p^y &= -oldsymbol{(Z^{(k)})}^Toldsymbol{H}_L^{(k)}oldsymbol{Y}^{(k)}oldsymbol{d}_p^y -oldsymbol{(Z^{(k)})}^Toldsymbol{D}_p^y oldsymbol{e}^{(k)}oldsymbol{e}^{(k)}oldsymbol{(Z^{(k)})}^Toldsymbol{
abla} oldsymbol{D}_p^y oldsymbol{d}_p^y &= -oldsymbol{(Z^{(k)})}^Toldsymbol{H}_L^{(k)}oldsymbol{\nabla}arphioldsymbol{d}_p^y -oldsymbol{(Z^{(k)})}^Toldsymbol{D}_p^y oldsymbol{e}^{(k)}oldsymbol{D}_p^y oldsymbol{d}_p^y oldsymbol{\nabla}_p^y oldsymbol{D}_p^y o$$

- If  $a(p^{(k)}) = 0$ , i.e. the constraints are linear,  $d_p^y = 0$ .
- If additionally  $(\mathbf{Z}^{(k)})^T \cdot \nabla \varphi(\mathbf{p}^{(k)}) = \mathbf{0}$ , then  $d_p^z = \mathbf{0}$ .

**First- and Second-Order Conditions** Using the matrix Z some of the necessary and sufficient conditions can be formulated equivalent:

• First-order necessary condition:

<sup>&</sup>lt;sup>1</sup>Note that the parameters of the Hessian are kept implicitly for brevity.

$$\boldsymbol{\nabla}\varphi(\boldsymbol{p}^*) - \sum_{i=1}^{n_a} \mu_i \cdot \boldsymbol{\nabla} a_i(\boldsymbol{p}^*) - \sum_{i \in A(\boldsymbol{p}^*)} = \boldsymbol{0} \quad \Longleftrightarrow \quad \boldsymbol{Z}^T(\boldsymbol{p}^*) \cdot \boldsymbol{\nabla}\varphi(\boldsymbol{p}^*) = \boldsymbol{0}$$

- Second-order necessary condition: The reduced Hessian of the Lagrangian Z<sup>T</sup>(p<sup>\*</sup>) · H<sup>p</sup><sub>L</sub>(p<sup>\*</sup>, μ<sup>\*</sup>, σ<sup>\*</sup>) · Z(p<sup>\*</sup>) is positive semidefinite.
- Second-order sufficient condition: The reduced Hessian of the Lagrangian is positive definite.

### Example

#### Approximation of the Reduced Hessian

The reduced Hessian contains all information that is needed to compute the QP solution! Hence, SQP methods can be built on this reduced Hessian. A Quasi-Newton approximation, e.g. BFGS, can be used to update the reduced Hessian:

$$\begin{split} \tilde{\boldsymbol{H}}^{(k+1)} &= \tilde{\boldsymbol{H}}^{(k)} - \frac{1}{\left(\boldsymbol{d}^{(k)}\right)^{T} \tilde{\boldsymbol{H}}^{(k)} \boldsymbol{d}^{(k)}} \tilde{\boldsymbol{H}}^{(k)} \boldsymbol{d}^{(k)} \left(\tilde{\boldsymbol{H}}^{(k)} \boldsymbol{d}^{(k)}\right)^{T} + \frac{1}{\left(\boldsymbol{g}^{(k)}\right)^{T} \boldsymbol{d}^{(k)}} \boldsymbol{g}^{(k)} \left(\boldsymbol{g}^{(k)}\right)^{T} \\ \boldsymbol{d}^{(k)} &= \boldsymbol{d}_{p}^{z} \\ \boldsymbol{g}^{(k)} &= \left(\boldsymbol{Z}^{(k+1)}\right)^{T} \cdot \boldsymbol{\nabla}_{p} L \left(\boldsymbol{p}^{(k+1)}, \boldsymbol{\mu}^{(k+1)}, \boldsymbol{\sigma}^{(k+1)}\right) - \left(\boldsymbol{Z}^{(k)}\right)^{T} \cdot \boldsymbol{\nabla}_{p} L \left(\boldsymbol{p}^{(k)}, \boldsymbol{\mu}^{(k+1)}, \boldsymbol{\sigma}^{(k+1)}\right) \end{split}$$

Analogous to unconstrained optimization, it is helpful to replace the rank-2 update with a rank-1 update and directly approximate a Cholesky decomposition.

## 4.5.6 SQP Method (Algorithm)

A sketch of the implementation of an SQP method is given in algorithm 6.

## 4.5.7 Notes

- When using the reduced Hessian, the matrix  $Z^{(k)}$  has to be updated in every iteration.
- Especially for high-dimensional problems, the gradients and Jacobians are sparse, i.e. only a few entries are nonzero. This behavior can be exploited in order to optimize implementations a lot.
- To rise efficiency and robustness further, sophisticated implementations of SQP methods (e.g. NPSOL, SNOPT) differentiate between constraints like
  - upper and lower bounds
  - linear constraints
  - nonlinear constraints

## Algorithm 6: Sequential Quadratic Programming

3

4

5

6

7

8

1 Initialization: Choose an initial approximation  $p^{(0)}$ , set  $k \leftarrow 0$ 2 while not converged do Calculate the Lagrange multipliers  $\mu^{(k)}$ ,  $\sigma^{(k)}$  using least squares:  $\min_{\boldsymbol{\mu}\in\mathbb{R}^{n_a}} \left\|\boldsymbol{\nabla}\varphi(\boldsymbol{p}^{(k+1)}) - \sum_{i=1}^{n_a} \mu_i^{(k+1)} \cdot \boldsymbol{\nabla}a_i(\boldsymbol{p}^{(k+1)}) - \sum_{i\in A(\boldsymbol{p}^{(k+1)})} \sigma_i^{(k+1)} \cdot \boldsymbol{\nabla}a_i(\boldsymbol{p}^{(k+1)})\right\|_2^2$ if termination criteria fulfilled then return Calculate new search direction  $d_p^{(k)}$  by solving the quadratic problem:  $\min_{\boldsymbol{d}_p \in \mathbb{R}^{n_p}} \varphi\big(\boldsymbol{p}^{(k)}\big) + \Big(\boldsymbol{\nabla}\varphi\big(\boldsymbol{p}^{(k)}\big)\Big)^T \boldsymbol{d}_p + \frac{1}{2} \boldsymbol{d}_p^T \boldsymbol{H}_L^{\boldsymbol{p}}\big(\boldsymbol{p}^{(k)}, \boldsymbol{\mu}^{(k)}, \boldsymbol{\sigma}^{(k)}\big) \boldsymbol{d}_p$ subject to  $\boldsymbol{a}(\boldsymbol{p}^{(k)}) + \boldsymbol{J}_{\boldsymbol{a}}^T(\boldsymbol{p}^{(k)})\boldsymbol{d}_p = \boldsymbol{0}$  $oldsymbol{b}(oldsymbol{p}^{(k)})+oldsymbol{J}_{oldsymbol{b}}^{T}(oldsymbol{p}^{(k)})oldsymbol{d}_{p}\geq oldsymbol{0}$ Calculate the step size  $\alpha^{(k)}$  using by minimizing a merit function, e.g. the augmented Lagrangian  $\min_{\alpha \in \mathbb{R}^+} L_Q \big( \boldsymbol{p}^{(k)} + \alpha \boldsymbol{d}_p^{(k)}, \, \boldsymbol{\mu}^{(k)} + \alpha \boldsymbol{d}_\mu^{(k)}, \, \boldsymbol{\sigma}^{(k)} + \alpha \boldsymbol{d}_\sigma^{(k)}, \, \rho^{(k)} \big)$ or the exact  $\ell_1$ -penalty function  $\min_{\alpha \in \mathbb{R}^+} \Phi_{\ell_1}(\boldsymbol{p}^{(k)} + \alpha \boldsymbol{d}_p^{(k)}, \rho^{(k)})$ 

Calculate the new solution approximation:

$$\boldsymbol{p}^{(k+1)} \leftarrow \boldsymbol{p}^{(k)} + \alpha^{(k)} \boldsymbol{d}_p^{(k)}$$

## 4.5.8 Examples

## **Optimal Control of a 6-DoF Industry Robot**

### **Car Drive**

## 4.5.9 Wrap-Up

- Motivation:
  - The Newton method is applied to determine a root ("zero point") of the gradient of the Lagrangian.
  - This yields a linear equation system to determine a search direction.
  - Requires first- and second-order derivatives of the objective and the constraints.
  - Preliminaries for the derivative where differentiability and knowing which constraints are active.
- Basic structure:
  - The linear system derived as a first step is equivalent to the solution of a quadratic problem (QP).
  - Solving the QP is faster and more robust than solving the linear system directly.
  - This yields a sequence of quadratic problems, thus sequential quadratic programming (SQP).
- Improvement of the basic structure:
  - For globalizing the method, a step size determination was introduced using line search on an appropriate test function (penalty function).
  - As the second-order derivatives are commonly not available, Quasi-Newton approximations of the Lagrangian are used.
- SQP for high-dimensional NLPs:
  - If the NLP is high-dimensional, the reduced Hessian shall be used.
  - The QP can be solved faster and the Hessian of the Lagrangian can be approximation using Quasi-Newton approaches!
  - High-dimensional NLPs often have sparse gradients as Jacobians which can be used for further performance improvements.
- Outlook:
  - There are lots of SQP methods, e.g. trust-region SQP that can work with indefinite or negative definite Hessians or methods that allow only feasible approximations in every iteration (feasible SQP methods), while the "classic" SQP only fulfills the constraints at the end (infeasible SQP method).
  - Other numerical methods for solving nonlinear, constrained optimization problems exist, e.g. inner point methods.

# **5** Calculation of Derivatives

To use efficient gradient-based algorithms, the first-order derivatives

$\partial arphi$	$\partial a_k$	$\partial b_l$
$\overline{\partial p_i}$	$\overline{\partial p_i}$	$\overline{\partial p_i}$

of the objective and the constraints are needed. But these are typically not available directly! And even one wrong derivative could destroy the fast convergence properties...

# 5.1 Finite Difference Approximation (numerical Differentiation)

## 5.1.1 Forward Difference Approximation

The most known approximation of the first derivative is the forward difference approximation

$$rac{\partial arphi(oldsymbol{p})}{\partial p_i} pprox D_{V,i} arphi(oldsymbol{p}) = rac{1}{\delta_i} ig( arphi(oldsymbol{p}+oldsymbol{e}_i \delta_i) - arphi(oldsymbol{p}) ig)$$

where  $e_i \in \mathbb{R}^{n_p}$  is the *i*-th unit vector and  $\delta_i$  is an appropriate step size. The complete error is:

#### Error

The complete error of the approximation is composed of the following:

- Approximation error (theoretical error)
- Function precision
- Rounding error

**Approximation Error** The (theoretical) approximation error is given as the neglected terms of the Taylor approximation (i.e. the Lagrangian remainder):

$$\varphi(\boldsymbol{p} + \boldsymbol{e}_{i}\delta_{i})^{T(\boldsymbol{p})} \varphi(\boldsymbol{p}) + \frac{\partial\varphi(\boldsymbol{p})}{\partial p_{i}}\delta_{i} + \frac{1}{2}\frac{\partial^{2}\varphi(\tilde{\boldsymbol{p}})}{\partial p_{i}^{2}}\delta_{i}^{2}, \quad \tilde{\boldsymbol{p}} \in [\boldsymbol{p}, \boldsymbol{p} + \boldsymbol{e}_{i}\delta_{i}]$$

$$\iff \qquad \varphi(\boldsymbol{p} + \boldsymbol{e}_{i}\delta_{i}) - \varphi(\boldsymbol{p}) - \frac{\partial\varphi(\boldsymbol{p})}{\partial p_{i}}\delta_{i} = \frac{1}{2}\frac{\partial^{2}\varphi(\tilde{\boldsymbol{p}})}{\partial p_{i}^{2}}\delta_{i}^{2}$$

$$\iff \qquad \underbrace{\frac{1}{\delta_{i}}\left(\varphi(\boldsymbol{p} + \boldsymbol{e}_{i}\delta_{i}) - \varphi(\boldsymbol{p})\right)}_{=D_{V,i}\varphi(\boldsymbol{p})} - \frac{\partial\varphi(\boldsymbol{p})}{\partial p_{i}} = \frac{1}{2}\frac{\partial^{2}\varphi(\tilde{\boldsymbol{p}})}{\partial p_{i}^{2}}\delta_{i}$$

$$\iff \qquad D_{V,i}\varphi(\boldsymbol{p}) - \frac{\partial\varphi(\boldsymbol{p})}{\partial p_{i}} = \frac{1}{2}\frac{\partial^{2}\varphi(\tilde{\boldsymbol{p}})}{\partial p_{i}^{2}}\delta_{i} =: T_{V,i}(\varphi;\delta_{i})$$

In theory, the error should decrease with the step size. But today's computers have finite arithmetic! Other error factors have a serious role demolishing this theoretical result.

**Function Precision** The function precision takes into account that the target function  $\varphi$  cannot be calculated with machine precision, e.g. because the evaluation depends on other methods or because rounding errors have summed up due to cancellation or ill-conditioning.

This can be taken into account with the absolute errors  $\varepsilon$ ,  $\varepsilon_{\delta_i}$  of the function evaluations:

Where the absolute error can be expressed in terms of a relative error  $\varepsilon_R = 10^{-n_d}$  as  $\varepsilon = \varepsilon_R \varphi(\mathbf{p})$  where  $n_d$  are the number of decimal places that are correct. Plugging  $\hat{\varphi}$  into the forward approximation yields the function precision error  $C(D_{V,i}\varphi; \delta_i)$ :

$$D_{V,i}\hat{\varphi}(\boldsymbol{p}) = \frac{1}{\delta_i} \left( \hat{\varphi}(\boldsymbol{p} + \boldsymbol{e}_i \delta_i) - \hat{\varphi}(\boldsymbol{p}) \right) = \frac{1}{\delta_i} \left( \varphi(\boldsymbol{p} + \boldsymbol{e}_i \delta_i) - \varphi(\boldsymbol{p}) \right) + \frac{\varepsilon_{\delta_i} - \varepsilon}{\delta_i} =: D_{V,i}\varphi(\boldsymbol{p}) + C(D_{V,i}\varphi;\delta_i)$$

**Rounding Error** Additionally to the function precision and the theoretical error, rounding errors are produced by the subtraction and the division. But if  $\delta_i$  does not get "too small", these are negligible compared to the approximation error and the function precision.

**Total Error** Hence, the total error is given as:

$$T_{V,i}(\varphi;\delta_i) + C(D_{V,i}\varphi;\delta_i) = \frac{1}{2} \frac{\partial^2 \varphi(\tilde{\boldsymbol{p}})}{\partial p_i^2} \delta_i + \frac{\varepsilon_{\delta_i} - \varepsilon}{\delta_i}$$
(5.1)

#### Choosing the Step Size

Ideally, the step size should be chooses such that the error is minimal. As the error term (5.1) contains second-order derivatives that are not available, an upper bound has to be drawn on the error that more or less independent of the derivatives:

$$\begin{aligned} \left| \frac{1}{2} \frac{\partial^2 \varphi(\tilde{\boldsymbol{p}})}{\partial p_i^2} \delta_i + \frac{\varepsilon_{\delta_i} - \varepsilon}{\delta_i} \right| &\leq \frac{1}{2} \delta_i \left| \frac{\partial^2 \varphi(\tilde{\boldsymbol{p}})}{\partial p_i^2} \right| + \frac{1}{\delta_i} |\varepsilon_{\delta_i} - \varepsilon| \leq \frac{1}{2} \delta_i L_{\varphi'',i} + \frac{2}{\delta_i} \varepsilon_R L_{\varphi} \\ L_{\varphi'',i} &\coloneqq \max \left\{ \left| \frac{\partial^2 \varphi(\tilde{\boldsymbol{p}})}{\partial p_i^2} \right| : \tilde{\boldsymbol{p}} \in [\boldsymbol{p}, \boldsymbol{p} + \boldsymbol{e}_i \delta_i] \right\} \\ L_{\varphi} &\coloneqq \max \left\{ \left| \varphi(\boldsymbol{p}) \right|, \left| \varphi(\boldsymbol{p} + \boldsymbol{e}_i \delta_i) \right| \right\} \end{aligned}$$

Minimizing this w.r.r. to the step size yields:

$$\min_{\delta_i \in \mathbb{R}^+} \Phi(\delta_i), \quad \Phi(\delta_i) = \frac{1}{2} \delta_i L_{\varphi'',i} + \frac{2}{\delta_i} \varepsilon_R L_{\varphi}$$
$$\implies \quad \Phi'(\delta_i) = \frac{1}{2} L_{\varphi'',i} - \frac{2}{\delta_i^2} \varepsilon L_{\varphi} \stackrel{!}{=} 0 \quad \stackrel{L_{\varphi'',i} \neq 0}{\longleftrightarrow} \quad \delta_i = \sqrt{\frac{4\varepsilon_R L_{\varphi}}{L_{\varphi'',i}}}$$

If  $L_{\varphi}/L_{\varphi'',i} \approx$ , then  $\delta_i \approx 2\sqrt{\varepsilon_R}$ . If additionally it is possible to evaluate the function with machine precision, i.e.  $\varepsilon_R = \varepsilon_{\text{mach}}$ , the optimal step size is simply

$$\delta_i \approx 2\sqrt{\varepsilon_{\text{mach}}}$$

This step size often is a good choice and thus set as the default in most implementations. It also explains the rule of thumb that forward-differences can approximately evaluate half of decimal places correctly.

#### Notes

- For evaluating the gradient  $\nabla \varphi(\mathbf{p})$  it is necessary to
  - Determine  $n_p$  step sized  $\delta_i$  and to evaluate  $\varphi$  at least  $2n_p$  times to approximate  $L_{\varphi'',i}$ .
  - Every iteration of a gradient-based method needs the gradients causing high computation times.
  - It is better to use a one-time approximation of *relative step sizes*  $\varepsilon_i$  with  $\delta_i = \varepsilon_i (1 + |p_i|)$  at the initialization  $p^{(0)}$ .
  - The value  $\varepsilon_i = 5\sqrt{\varepsilon_{\text{mach}}}$  can be used as an initial relative step size.
  - If the optimization fails, restart with a new initialization and re-calculate the step sizes.

#### 5.1.2 Central-Difference Approximation

For forward difference approximation often yields results that are good enough, except the gradients are too small. Additionally, the forward approximation is not sufficient if the optimization step size  $\alpha^{(k)}$  such that the changes in p are less than the step size  $\delta$  or the differences in the function values are "too small" relative to  $\delta$ .

A potentially more exact approximation are central differences:

$$\frac{\partial \varphi(\boldsymbol{p})}{\partial p_i} \approx D_{Z,i} \varphi(\boldsymbol{p}) = \frac{1}{2\delta_i} \big( \varphi(\boldsymbol{p} + \boldsymbol{e}_i \delta_i) - \varphi(\boldsymbol{p} - \boldsymbol{e}_i \delta_i) \big)$$

Analogous to the forward differences, Taylor-expand this formula yields the insight that the order of the approximation error is  $\mathcal{O}(\delta_i^2)$  while the order of the forward differences is  $\mathcal{O}(\delta_i)$ . Analogous to the forward differences, the optimal step size  $\delta_i^Z$  can be calculated by minimizing an upper bound on the total error, yielding the optimal step size

$$\delta_i^Z = \sqrt[3]{\frac{3\varepsilon_R L_{\varphi}}{L_{\varphi''',i}}}, \quad L_{\varphi''',i} = \max\left\{ \left| \frac{\partial^3 \varphi(\tilde{\boldsymbol{p}})}{\partial p_i^3} \right| : \tilde{\boldsymbol{p}} \in \left[ \boldsymbol{p} - \boldsymbol{e}_i \delta_i^Z, \boldsymbol{p} + \boldsymbol{e}_i \delta_i^Z \right] \right\}$$

For functions with  $L_{\varphi}/L_{\varphi''',i} \approx 1$  it follows  $\delta_i^Z \approx \sqrt[3]{3\varepsilon_R} \approx \delta_i^{2/3}$ . If additionally  $L_{\varphi} \approx 1$  and  $L_{\varphi''',i} \approx 1$ , it follows:

$$\left| D_{Z,i} \hat{\varphi}(\boldsymbol{p}) - \frac{\partial \varphi(\boldsymbol{p})}{\partial p_i} 
ight| \leq \varepsilon_R^{2/3}$$

Hence, the rule of thumb that central difference can approximately evaluate two third of decimal places correctly.

But central differences produce much more computational overhead compared to forward/backward differences! Thus they should only be used if needed (by switching from forward to central differences).

To approximate the second derivative, the following schema can be used

$$rac{\partial^2 \varphi(oldsymbol{p})}{\partial p_i^2} pprox rac{1}{\delta_i^2} (\varphi(oldsymbol{p} + oldsymbol{e}_i \delta_i) - 2\varphi(oldsymbol{p}) + \varphi(oldsymbol{p} - oldsymbol{e}_i \delta_i))$$

that is a combination of forward and backward differences to approximate the second-order derivative. This directly gives the *i*-th diagonal entry of the Hessian of  $\varphi$ , which can reduce the number of iterations needed.

## 5.2 Numerical Differentiation of Simulation Models

An important class in optimization is simulation-based optimization where the system state x(t) is given as the numerical solution of ODEs or PDEs. In this setting, the objective  $\varphi$  and the constraints a, b are dependent on the state variables x of an ODE/PDE system. Hence, the calculation of  $\varphi(x(p))$  requires solving the ODE/PDE numerically:

- Every calculation of  $\varphi(\boldsymbol{x}(\boldsymbol{p}))$ ,  $\boldsymbol{a}$  and  $\boldsymbol{b}$  is computationally expensive.
- The calculation is only possible with simulation errors (i.e. approximation error in the ODE/PDE solver and accumulated rounding errors).
- The gradients (see below) are generally not available and have to me approximated.

$$oldsymbol{
abla} oldsymbol{arphi}oldsymbol{\left(x(p)
ight)} = rac{\partial arphiig(x(p)ig)}{\partial oldsymbol{p}} = rac{\partial arphiig)}{\partial oldsymbol{x}}rac{\partial arphi}{\partial oldsymbol{x}}$$

## 5.2.1 Derivative of ODE-Simulation Models

Given an IVP (initial value problem)  $\dot{x} = f(t, x; p)$ ,  $x(0) = x_0$ ,  $0 \le t \le t_f$ , the derivatives of the (numerical) solution x(t; p) w.r.t. to the parameters p are required. Formally, the parameters p can be transformed to initial values  $x_0$  and thus the derivatives w.r.t. the parameters to derivatives w.r.t. the initial values. The derivative w.r.t. the initial values is called the *sensitivity matrix*:

$$rac{\partial oldsymbol{x}(t;oldsymbol{x}_0)}{\partial oldsymbol{x}_0}$$

The IVP with the parameters p transformed to initial values is given as

$$\begin{bmatrix} \dot{x}_1 \\ \vdots \\ \dot{x}_{n_x} \end{bmatrix} = \dot{\boldsymbol{x}} = \boldsymbol{f}(t, \boldsymbol{x}, x_{n_x+1}, \cdots, x_{n_x+n_p})$$
$$\begin{bmatrix} \dot{x}_{n_x+1} \\ \vdots \\ \dot{x}_{n_x+n_p} \end{bmatrix} = \begin{bmatrix} 0 \\ \vdots \\ 0 \end{bmatrix}$$

with  $x_{n_x+1} \coloneqq p_1, \dots, x_{n_x+n_p} \coloneqq p_{n_p}$  and the initial values

$$\boldsymbol{x}(0) = \boldsymbol{x}_0 \qquad \begin{bmatrix} x_{n_x+1}(0) \\ \vdots \\ x_{n_x+n_p}(0) \end{bmatrix} = \begin{bmatrix} p_1 \\ \vdots \\ p_{n_p} \end{bmatrix}$$

That is, the "parameter states" are added as time invariant states always equaling the parameters.

## 5.2.2 External Numerical Differentiation

#### Naïve Approach

Just use forward differences:

$$rac{\partial oldsymbol{x}(t;oldsymbol{p})}{\partial p_i}pprox rac{1}{\delta_i}ig(oldsymbol{x}(t;oldsymbol{p}+oldsymbol{e}_i\delta_i)-oldsymbol{x}(t;oldsymbol{p})ig)$$

- This approach requires solving  $n_p$  additional IVP solutions for calculating the tweaked evaluations.
- To solve the IVP as best as possible, it may use variable step sizes.
- This causes gradient-based algorithms to have extremely big problems in finding the minimum!
- But this is not only caused by the optimization method...

**Runge-Kutta Methods** When using Runge-Kutta methods of order m, the numerical solutions depends on the integration tolerance and the internal step width of the integration. By just looking at the initial values, the error of RK methods of order m is

$$\tilde{\boldsymbol{x}}(t; \boldsymbol{x}_{0}, h_{1}) = \boldsymbol{x}(t; \boldsymbol{x}_{0}) + \sum_{j=m}^{N} \tilde{c}_{j}(t, \boldsymbol{x}_{0}) h_{1}^{j} + \mathcal{O}(h_{1}^{N+1})$$
$$\tilde{\boldsymbol{x}}(t; \boldsymbol{x}_{0} + \boldsymbol{e}_{i}\delta_{i}, h_{2}) = \boldsymbol{x}(t; \boldsymbol{x}_{0} + \boldsymbol{e}_{i}\delta_{i}) + \sum_{j=m}^{N} \tilde{c}_{j}(t, \boldsymbol{x}_{0} + \boldsymbol{e}_{i}\delta_{i}) h_{2}^{j} + \mathcal{O}(h_{2}^{N+1})$$

with differentiable functions  $\tilde{c}_i(t, x_0)$ . Plugging this into the forward difference scheme yields:

$$\frac{1}{\delta_i} \left( \tilde{\boldsymbol{x}}(t; \boldsymbol{x}_0 + \boldsymbol{e}_i \delta_i, h_2) - \tilde{\boldsymbol{x}}(t; \boldsymbol{x}_0, h_1) \right)$$

$$= \frac{\partial \boldsymbol{x}(t; \boldsymbol{x}_0)}{\partial x_{0,i}} + \mathcal{O}(\delta_i) + \sum_{j=m}^N \tilde{c}_j(t, \boldsymbol{x}_0 + \boldsymbol{e}_i \delta_i) \cdot \underbrace{\frac{h_2^j - h_1^j}{\delta_i}}_{\rightarrow \infty} + \sum_{j=m}^N \left( \frac{\partial \tilde{c}_j(t, \boldsymbol{x}_0)}{\partial x_{0,i}} + \mathcal{O}(\delta_i) \right) h_1^j + \mathcal{O}\left( \frac{h_1^{N+1}}{\delta_i} \right) + \mathcal{O}\left( \frac{h_2^{N+1}}{\delta_i} \right)$$

This is the problem! If the integration steps  $h_1 \neq h_2$  are not equal, the error term becomes dominant as  $\delta_i$  usually is "small".

A "solution" would be to set  $h_1 = h_2$ , causing bad integration results.

#### **Coupled Forward Differences Approximation**

It is possible to simultaneously integrate an  $(n_p + 1)$ -times big IVP for each tweaked value guaranteeing  $h_1 = h_2$ :

$$\begin{aligned} \dot{\boldsymbol{x}} &= \boldsymbol{f}(t, \boldsymbol{x}), \quad \boldsymbol{x}(0) = \boldsymbol{x}_0 \\ \dot{\boldsymbol{x}} &= \boldsymbol{f}(t, \boldsymbol{x}) \\ &\vdots \\ \dot{\boldsymbol{x}} &= \boldsymbol{f}(t, \boldsymbol{x}) \end{aligned}$$

# 5.2.3 Internal Numerical Differentiation

The external numerical differentiation costs a lot of time! One insight: The sensitivity matrix can be expresses as the solution of a matrix-ODE:

$$\frac{\mathrm{d}}{\mathrm{d}t}\frac{\partial \boldsymbol{x}(t;\boldsymbol{x}_0)}{\partial \boldsymbol{x}_0} = \frac{\partial \boldsymbol{f}\big(t,\boldsymbol{x}(t;\boldsymbol{x}_0)\big)}{\partial \boldsymbol{x}}\frac{\partial \boldsymbol{x}(t;\boldsymbol{x}_0)}{\partial \boldsymbol{x}_0}, \quad \frac{\partial \boldsymbol{x}(0;\boldsymbol{x}_0)}{\partial \boldsymbol{x}_0} = \boldsymbol{I}$$

- Variant 1: Simultaneously integrate the ODE and the matrix-ODE with a standard integrator.
- Variant 2 (better): Differentiated integrate method calculates x(t; p),  $\frac{\partial x(t; p)}{\partial n}$ ,  $\frac{\partial \dot{x}(t; p)}{\partial n}$ 
  - Advantage: Really efficient.
  - Disadvantage: Implementation complexity; especially complication for switching points.

# 5.3 Symbol Differentiation

If the functions  $\varphi$ , a, b are given as explicit formulas, the derivatives  $\nabla \varphi(p)$ ,  $J_a(p)$ ,  $J_b(p)$  could be evaluated using a computer algebra system (e.g. Maple, Mathematica, SymPy, ...) in closed form. This is based on a systematic application of chain, product, ...rules and often requires a special input format.

- Advantage: No approximation error, only rounding errors.
- Disadvantages: Can lead to complex functions that are computationally expensive to evaluate.

# 5.4 Automatic Differentiation

Automatic differentiation refers to a technique to generate first and possible second-order derivatives of an existing program that calculates  $\varphi$ , a, b. This is based on the insight that every every so complex function can be composed of as a sequence of elementary functions with one or two arguments:

- 1-argument functions: Sine, Cosine, Exponential, Logarithm, ...
- 2-argument functions: Addition, Subtraction, Multiplication, Division, Exponentiation

AD-methods are based on an analysis of the evaluation sequence as a *computation graph* of elementary functions.

There are two main design decisions:

- Pre-compile a program to get the derivative function to be able to evaluation the function and the derivatives simultaneously (forward mode).
- Build the computation graph after the function has been evaluated and evaluate the derivative afterwards (backward mode).

The computational complexity of the gradient in forward mode on a scalar function with multiply variables  $\varphi(\mathbf{p})$  can be as high as for symbolic differentiation:  $\mathcal{O}(n_p)$ . But for the backward mode, a complexity of  $\mathcal{O}(5)$  can be reached independent of  $n_p$ ! But the memory complexity can rise a lot...

• Advantages:

- No approximation error, only rounding errors.
- AD is continuously improving and even today's algorithms are capable of a lot of calculations.
- Disadvantages:
  - It is problematic to handle piecewise functions (if-then-else), approximation of tabular data, approximations of Sine, Cosine, ... by rationale functions.
     But this is problematic even for analytical derivatives...

Even ODE- and PDE-simulations can be handled using AD!

Popular implementations, commonly used in machine learning, are libraries like TensorFlow (static computation graph) and PyTorch (dynamic computation graph).

# 6 Parameter Estimation

This chapter covers a special type of objective function: the sum of squares, called *least squares* optimization problems:

$$arphi(m{p}) = rac{1}{2} \sum_{i=1}^{n_r} r_i^2(m{p}) = rac{1}{2} \|m{r}(m{p})\|_2^2$$

This objective function arises in a lot of optimization problems, e.g.: Given some detected corner points  $(x_i, y_i)_{i=1,\dots,n_r}$  of a ball, what is the radius  $R_K$  and the position  $(x_K, y_K)$  if the ball? The residual (error) function r can be determined from the circle equation:

$$R_K^2 = (x_i - x_K)^2 + (y_i - y_K)^2 \implies r_i(x_K, y_K, R_K) = \sqrt{(x_i - x_K)^2 + (y_i - y_K)^2} - R_K$$

If, in practice, some parametric model is used, this almost always leads to a parameter fitting problem to measure/minimize the differences between the model and measurements. By minimizing the differences, the model with that most correspond to the measurements can be found (e.g. the parameters of a friction model or inertia of a robot).

## 6.1 Objective Functions

There are various different objective functions that could be used:

- Absolute sum:  $\varphi_1(p) = \|r(p)\|_1 = \sum_{i=1}^{n_r} |r_i(p)|$
- Sum of squares:  $\varphi_2(p) = \frac{1}{2} \|r(p)\|_2^2 = \frac{1}{2} \sum_{i=1}^{n_r} r_i^2(p)$
- Maximum difference:  $\varphi_{\infty}(\boldsymbol{p}) = \|\boldsymbol{r}(\boldsymbol{p})\|_{\infty} = \max\{|r_i(\boldsymbol{p})|: i = 1, \dots, n_r\}$

But even if r is differentiable,  $\varphi_1$  and  $\varphi_{\infty}$  are generally not. But  $\varphi_{\infty}$  can be transformed into a NLP with differentiable functions:

$$\min_{\boldsymbol{p}, p_{n_r+1}} p_{n_r+1}$$
  
subject to  $-p_{n_r+1} \leq r_i(\boldsymbol{p}) \leq p_{n_r+1}, \quad i = 1, \cdots, n_r$ 

Anyway,  $\varphi_{\infty}$  is extremely sensitive towards outliers. On the other hand,  $\varphi_2$  is less sensitive towards outliers and is differentiable (if *r* is differentiable)!

Hence, in most cases the least squares objective  $\varphi_2$  is used. Besides the differentiability, it is also statistically appealing: If the measurement noise  $\varepsilon_{ij}$  are statistically independent and Gaussian distributed with mean 0 and variance  $\sigma^2$ , the solution of  $\varphi_2 \rightarrow \min$  is a maximum likelihood estimator!

Additionally, by exploiting the structure of  $\varphi_2$ , the efficiency of gradient-based algorithms can be increased (e.g. due to sparse Jacobians).

# 6.2 Linear Least Squares

In the special case of a linear residual function  $r(p) = J^T p + f$  and the least squares objective

$$arphi(m{p}) = arphi_2(m{p}) = rac{1}{2} \|m{r}(m{p})\|_2^2$$

the optimization problem is quadratic in *p*:

$$\varphi(\boldsymbol{p}) = \frac{1}{2} \left\| \boldsymbol{J}^T \boldsymbol{p} + \boldsymbol{f}_r \right\| = \frac{1}{2} \left( \boldsymbol{p}^T \boldsymbol{J} + \boldsymbol{f}_r^T \right) \left( \boldsymbol{J}^T \boldsymbol{p} + \boldsymbol{f}_r \right)$$

Zeroing the gradient  $\nabla \varphi(\mathbf{p}) \stackrel{!}{=} \mathbf{0}$  yields the *normal equations* 

$$\boldsymbol{J}\boldsymbol{J}^T\boldsymbol{p}^* = -\boldsymbol{J}\boldsymbol{f}_r$$

that a solution must fulfill (where  $JJ^T$  is symmetric and positive definite). This linear system should not be solved as a normal linear system as it is ill-conditioned! Better use special methods like orthogonalization or single value decomposition.

# 6.3 Optimality Conditions and Special Methods

Let  $\varphi(\mathbf{p}) = \varphi_2(\mathbf{p})$  be two times continuously differentiable. Then the following first- and second-order necessary conditions can be formulated:

1. The gradient has to vanish:

$$abla arphi_2({m p}^*) = {m J}_{m r}({m p}^*) \cdot {m r}({m p}^*) = {m 0}$$

2. The Hessian has to be positive semidefinite:

$$H_{\varphi_2}(p^*) = J_r(p^*) J_r^T(p^*) + \sum_{i=1}^{n_r} \left( r_i(p^*) \cdot H_{r_i}(p^*) \right)$$
(6.1)

#### 6.3.1 Gauss-Quasi-Newton Method

In the Quasi-Newton method, the search direction is determined as a solution of the linear system

$$oldsymbol{H}_arphiig(oldsymbol{p}^{(k)}ig)oldsymbol{d}^{(k)}=-oldsymbol{
abla}arphiig(oldsymbol{p}^{(k)}ig)$$

where the Hessian is approximated, e.g. using a BFGS update. But in the case of a least squares optimization problem, the only part of the Hessian (6.1) depending on second-order derivatives is:

$$\sum_{i=1}^{n_r} \left( r_i(\boldsymbol{p}^*) \cdot \boldsymbol{H}_{r_i}(\boldsymbol{p}^*) \right)$$

But as the objective gets smaller and smaller, this term also vanishes. Hence, the Hessian can be approximated using first-order derivatives only

$$oldsymbol{H}_{arphi_2}(oldsymbol{p}^*) pprox oldsymbol{J}_{oldsymbol{r}}(oldsymbol{p}^*) oldsymbol{J}_{oldsymbol{r}}^T(oldsymbol{p}^*)$$

if the residuals are "small". This leads to the *Gauss-Newton Method* where the search direction is given as the solution of the normal equations

$$oldsymbol{J}_r(oldsymbol{p}^{(k)})oldsymbol{J}_r^T(oldsymbol{p}^{(k)})oldsymbol{d}^{(k)} = -oldsymbol{J}_r(oldsymbol{p}^{(k)})\cdotoldsymbol{r}(oldsymbol{p}^{(k)})$$

or as the solution of a (better conditioned) linear least squares problem:

$$\min_{oldsymbol{d} \in \mathbb{R}^{n_p}} rac{1}{2} \Big\| oldsymbol{J}_r^Toldsymbol{p}^{(k)}oldsymbol{d} + oldsymbol{r}oldsymbol{p}^{(k)}ildsymbol{d} + oldsymbol{r}oldsymbol{p}^{(k)}ildsymbol{d} \Big\|_2^2$$

If residuals are big or the problem is ill-conditioned, it can be modified with a suitable matrix  $B^{(k)}$ :

$$\left( m{J}_{m{r}}(m{p}^{(k)})m{J}_{m{r}}^T(m{p}^{(k)}) + m{B}^{(k)} 
ight) m{d}^{(k)} = -m{J}_{m{r}}(m{p}^{(k)}) \cdot m{r}(m{p}^{(k)})$$

Additionally a step size rule should be used. This method is implemented, e.g. in *NLSCON* which also allows additional nonlinear inequality constraints.

#### 6.3.2 Levenberg-Marquardt Methods

Instead of a Newton-approach, trust region methods can be used to determine the search direction. In the *Levenberg-Marquardt Method*, the search direction  $d^{(k)}$  is given as the solution of:

$$\left(\boldsymbol{J}_{\boldsymbol{r}}(\boldsymbol{p}^{(k)})\boldsymbol{J}_{\boldsymbol{r}}^{T}(\boldsymbol{p}^{(k)}) + \gamma^{(k)}\boldsymbol{I}\right)\boldsymbol{d}^{(k)} = -\boldsymbol{J}_{\boldsymbol{r}}(\boldsymbol{p}^{(k)})\cdot\boldsymbol{r}(\boldsymbol{p}^{(k)}), \quad \gamma^{(k)} \ge 0$$

That is, the search direction is a mixture of the Gauss-Newton direction and steepest descent. The same search direction can be received by solving a constrained linear least squares problem:

$$\min_{\boldsymbol{d} \in \mathbb{R}^{n_p}} \frac{1}{2} \left\| \boldsymbol{J}_{\boldsymbol{r}}^T(\boldsymbol{p}^{(k)}) \boldsymbol{d} + \boldsymbol{r}(\boldsymbol{p}^{(k)}) \right\|_2^2$$
subject to  $\|\boldsymbol{d}\|_2 \le \delta^{(k)}$ 

Where  $\delta$  and  $\gamma$  have a connection.

This is implemented, e.g. in LMDER, LMJAC of MINPACK.

#### 6.3.3 Notes

- The Gauss-Newton method needs an appropriate test function that also exploits the structure of the objective function to determine the step size to ensure global convergence.
- In the Levenberg-Marquardt method, globalization happens by choosing an appropriate trust region.
- Levenberg-Marquardt and Gauss-Newton methods are in general a lot more efficient (faster and more precise) than solving least squares problems using general purpose optimization techniques (e.g. Quasi-Newton or SQP).
- When optimization by simulation, the sensitivity matrix  $\frac{\partial x(t;p)}{\partial p}$  is needed (seesubsection 5.2.1).

# 6.4 Conditioning of Normal Equations

Often, the Jacobian of the objective function of linear least squares J is ill-conditioned, i.e. cond J is high, causing round error to be increased. In extreme cases, J does not has full rank, i.e. cond  $J = \infty$ . In this case, the solution d of the normal equations

$$JJ^Td = -Jr$$

is not unique! Common reasons for ill-conditioning are:

- "too few" measurements E.g. the measurements are not "dense enough".
- not the "correct" measurements E.g. the measurements to not depend or weakly depend on the optimization variables.
- System model does not fit to the measurements (it is incompatible). Then either the measurements or the model is wrong.

# 6.5 Result Interpretation

## 6.5.1 Common Problems

Common problems if the residuals are still high after the optimization terminates:

- Derivatives are not precise enough.
- The optimization method is not suitable for the problem, e.g. if the method cannot handle inexact function evaluations or has problems with local minima of  $\varphi_2$ .
- In parameter estimation settings,
  - the system model and measurements might be incompatible (wrong measurements) or
  - the system model and the physical might be incompatible (wrong model).

Common problems for small residuals:

- Some optimization variables might not be uniquely determined.
- Too less measurements, variables are not unique in general (e.g. linearly dependent in the model).

Practical aspects:

- Scaling/balancing of the variables: By the transformation p<sub>i</sub> → p̂<sub>i</sub> = s<sub>i</sub>p<sub>i</sub>, s<sub>i</sub> = const > 0, the derivative changes to ∂φ/∂p̂<sub>i</sub> = ∂φ/∂p̂<sub>i</sub> 1/s<sub>i</sub>.
- Scaling of the residuals by weights  $w_i > 0$ :  $\varphi(\mathbf{p}) = \frac{1}{2} \sum_{i=1}^{n_r} w_i r_i^2(\mathbf{p})$

## 6.5.2 The Covariance Matrix

Assuming the measurement errors  $\varepsilon_i$  (in the measurements  $r_i$ ) are normally distributed with zero mean and constant variance  $\sigma^2$ . Then the solution of the linear least squares problem is a maximum likelihood estimator for the parameters! The covariance matrix V is then given by  $V = \sigma^2 (JJ^T)^{-1}$  where the variance can be approximated by

$$\sigma^2 \approx \frac{\|\boldsymbol{r}(\boldsymbol{p}^*)\|^2}{n_r - n_p}$$

and the mean of the residual squares is given as

$$\mathbb{E}\big[\|\boldsymbol{r}(\boldsymbol{p}^*)\|^2\big] = (n_r - n_p)\sigma^2$$

Hence, a bad conditioning of J implies high variance!

# 6.6 Optimal Experimental Design

Goal of *optimal experimental design* is a good conditioning of the parameter identification problem, i.e. the optimal experimental parameters  $s^*$  is a solution of the optimization problem

$$\min_{\boldsymbol{s}} \phi_{\exp} \big( \boldsymbol{V}(\boldsymbol{s}) \big)$$

where V is the covariance matrix. Some objective  $\phi_{exp}$  are:

- 1. Determinant of V.
- 2. Trace mean, i.e.  $\operatorname{tr} \boldsymbol{V}/n_p$
- 3. Biggest eigenvalue of V
- 4. Absolute length of the biggest confidence interval.
- 5. Conditional number.

# 6.7 Examples

#### 6.7.1 Parameter-Dependent Vehicle Dynamics

#### Simulated Measurements

**Real Measurements** 

Comparison

## 6.7.2 Parameter Estimation for "BioBiped"

# 7 Minimization of Functionals

In the setting of *variational problems*, the unknown x is a function of t (where t is the independent variable) and the objective is a functional J[x] of integral type:

$$J[\boldsymbol{x}] = \int L(\boldsymbol{x}(t), \dot{\boldsymbol{x}}(t), t) \,\mathrm{d}t$$

the solution  $x^*$  must fulfill given constraints, e.g.:

- Initial and end conditions (boundary conditions)
- Inequality constraints
- Integral-type constraints
- Differential equations

If differential equations are given, the problem is called an optimal control problem.

## 7.1 Euler-Lagrange Equation

Given an optimization problem

$$\min_{\boldsymbol{x}} J[\boldsymbol{x}], \quad J[\boldsymbol{x}] = \int_{a}^{b} L(\boldsymbol{x}(t), \dot{\boldsymbol{x}}, t) dt$$
  
subject to  
$$\boldsymbol{x}(a) = \boldsymbol{x}_{a}$$
  
$$\boldsymbol{x}(b) = \boldsymbol{x}_{b}$$

where  $\boldsymbol{x}: \mathbb{R} \to \mathbb{R}^{n_x}$ ,  $L: \mathbb{R}^{n_x} \times \mathbb{R}^{n_x} \times R \to R$  are two times continuously differentiable.

A stationary solution  $x^*$ , i.e. a solution that is a minimum candidate, has to fulfill the *Euler-Lagrange Equation*:

$$\frac{\partial L}{\partial x_i} - \frac{\mathrm{d}}{\mathrm{d}t} \frac{\partial L}{\partial \dot{x}_i} = 0, \quad i = 1, \cdots, n_x$$

Additionally, the boundary conditions  $x(a) = x_a$ ,  $x(b) = x_b$  must be fulfilled, yielding a second-order ordinary boundary value problem.

## 7.1.1 Example

As of the Hamilton's principle x(t) between  $t_1$  and  $t_2$  is a stationary point of the action functional

$$\int_{t_1}^{t_2} L\big(\boldsymbol{x}(t), \dot{\boldsymbol{x}}(t), t\big) \,\mathrm{d}t$$

where L is the Lagrangian of the system, i.e. the difference of kinetic and potential energy:

$$L = T - V$$

For a ball with mass m and gravity acceleration g that is thrown into the air in a straight line (i.e. onedimensional), the kinetic and potential energy are given as:

$$T = \frac{m}{2}\dot{x}^2 \qquad \qquad V = mgx$$

The Lagrangian then is  $L = T - V = \frac{m}{2}\dot{x}^2 - mgx$ . Plugging that in the Euler-Lagrange equations yields a second-order differential equation for the movement of the ball:

$$\frac{\partial L}{\partial x} = -mg \qquad \frac{\partial L}{\partial \dot{x}} = mx \qquad \frac{\mathrm{d}}{\mathrm{d}t} \frac{\partial L}{\partial \dot{x}} = m\ddot{x} \implies -mg - m\ddot{x} = 0 \quad \Longleftrightarrow \quad \ddot{x} = -g$$

Solving this differential equation with the initial values x(0) = 0,  $\dot{x}(0) = \dot{x}_0$  yields the equation of movement of the ball:

$$x(t) = -\frac{g}{2}t^2 + \dot{x}_0t$$

As expected, this equation is a perfect parabola w.r.t. time! The same can be done in two dimensions.

## 7.1.2 Notes

- The Euler-Lagrange equations are in general not tractable. Hence, numerical methods have to be used.
- The equations can be extended to solutions with non-differentiabilities, (in-) equality constraints, integral-type constraints, ...
- Also, second-order necessary conditions can be formulated.

### 7.1.3 Derivation

# 8 Optimal Control

A dynamical system is described by the ODE

$$\dot{x} = f(x, u, p, z), \quad x(0) = x_0$$

where x is the state, u are the control variables, p are (constant) system parameters and z is noise. An optimal control problem has given  $x_0$ , p and  $x_i(t_f)$  and seeks for an optimal u and x.

A complete optimal control problem is given as:

$$\min J[\boldsymbol{u}], \quad J[\boldsymbol{u}] = \phi(\boldsymbol{x}(t_f), t_f) + \int_0^{t_f} L(\boldsymbol{x}(t), \boldsymbol{u}(t)) dt$$
  
subject to  
$$\dot{\boldsymbol{x}}(t) = \boldsymbol{f}(\boldsymbol{x}(t), \boldsymbol{u}(t))$$
$$x_i(0) = x_{i,0} = \text{const}, \quad i \in \{1, \dots, n_x\}$$
$$\boldsymbol{r}(\boldsymbol{x}(t_f), t_f) = \boldsymbol{0}, \quad \text{e.g. } x_j(t_f) = x_{j,f}, \ j \in \{1, \dots, n_x\}$$
$$\boldsymbol{g}(\boldsymbol{x}(t), \boldsymbol{u}(t)) \ge \boldsymbol{0}$$
$$\boldsymbol{g}(\boldsymbol{x}(t)) \ge \boldsymbol{0}$$
$$\boldsymbol{r}^i(\boldsymbol{x}(t_s - 0), \boldsymbol{x}(t_s + 0), t_s) = \boldsymbol{0}$$

The constraints (from to bottom) are called:

- Equations of movement, these are the defining property of an optimal control problem over a regular variational problem.
- Initial conditions (optional).
- Final conditions (optional).
- Control constraints (optional, often box constraints).
- State constraints (optional).
- Interior point constraints (optional).

Additionally, the final time  $t_f$  might either fixed or free. The objective (called the *Bolza functional*) is split into an endpoint cost (also called *Mayer term*) and a Lagrangian:

$$J[\boldsymbol{u}] = \underbrace{\phi(\boldsymbol{x}(t_f), t_f)}_{\text{Endpoint Cost}} + \int_0^{t_f} \underbrace{L(\boldsymbol{x}(t), \boldsymbol{u}(t))}_{\text{Lagrangian}} dt$$

Note that the endpoint cost is more general than the Lagrangian in terms that the Lagrangian term can be transformed to an endpoint cost by introducing a new state

$$\dot{x}_{n_x+1} = L(\boldsymbol{x}(t), \boldsymbol{u}(t)), \quad x_{n_x+1}(0) = 0$$

and changing the endpoint cost to:

$$\tilde{\phi}(\tilde{\boldsymbol{x}}(t_f), t_f) \coloneqq \phi(\boldsymbol{x}(t_f), t_f) + x_{n_x+1}(t_f)$$

Additionally, non-autonomous problems can be transformed to autonomous problems by introducing a "clock state"  $\dot{x}_{n_x+1} = 1$ ,  $x_{n_x+1}(0) = 0$  and changing the ODE accordingly:

$$\dot{\boldsymbol{x}}(t) = \boldsymbol{f}(t, \boldsymbol{x}(t), \boldsymbol{u}(t)) \quad \rightarrow \quad \dot{\boldsymbol{x}}(t) = \boldsymbol{f}(x_{n_x+1}, \boldsymbol{x}(t), \boldsymbol{u}(t))$$

# 8.1 Necessary Optimality Conditions for the Basis Problem

The basis optimal control problem does not have control, state or interior constraints:

 $\min J[\boldsymbol{u}], \quad J[\boldsymbol{u}] = \phi(\boldsymbol{x}(t_f), t_f) + \int_0^{t_f} L(\boldsymbol{x}(t), \boldsymbol{u}(t)) dt$ subject to  $\dot{\boldsymbol{x}}(t) = \boldsymbol{f}(\boldsymbol{x}(t), \boldsymbol{u}(t))$  $x_i(0) = x_{i,0} = \text{const}, \quad i \in \{1, \dots, n_x\}$  $\boldsymbol{r}(\boldsymbol{x}(t_f), t_f) = \boldsymbol{0}, \quad \text{e.g. } x_j(t_f) = x_{j,f}, j \in \{1, \dots, n_x\}$ 

For formulating the necessary optimality conditions, the following auxiliary functions are needed (the latter one is called *Hamiltonian*):

$$\Phi(\boldsymbol{x}, t, \boldsymbol{\nu}) \coloneqq \phi(\boldsymbol{x}, t) + \boldsymbol{\nu}^T \boldsymbol{r}(\boldsymbol{x}, t)$$
$$H(\boldsymbol{x}, \boldsymbol{u}, \boldsymbol{\lambda}) \coloneqq L(\boldsymbol{x}, \boldsymbol{u}) + \boldsymbol{\lambda}^T \boldsymbol{f}(\boldsymbol{x}, \boldsymbol{u})$$

Note that, while kept implicitly, the auxiliary variables  $\nu$  and  $\lambda$  are time-dependent! The  $\lambda$  are also called the *adjunct variables* of the optimal control problem.

## 8.1.1 Boundary Conditions

To solve the optimal control problem,  $2n_x$  boundary conditions are needed (the optimality conditions will yield one first-order ODE for every state and adjunct variable). Additional to the given boundary conditions of the problem formulation, enough conditions have to be found to get  $2n_x$  conditions. Common cases:

- (i) Fixed initial conditions  $x_i(0) = x_{i,0} = \text{const:}$ Either  $x_k(0)$  is given or, if not, set  $\lambda_i(0) = -\frac{\partial \phi(\boldsymbol{x}(0), \boldsymbol{x}(t_f), t_f)}{\partial x_i(0)}$
- (ii) Fixed final conditions  $x_i(t_f) = x_{i,f} = \text{const:}$ Either  $x_i(t_f)$  is given or, if not and  $x_i$  is not part of  $r(\cdots)$ , set  $\lambda_i(t_f) = \frac{\partial \phi(\boldsymbol{x}(t_f), t_f)}{\partial x_i(t_f)}$
- (iii) Mixed boundary conditions:
  - a) General boundary conditions  $r(x(0), x(t_f), t_f) = 0$ : If  $x_i(0)$  is free, set  $\lambda_i(0) + \frac{\partial \Phi}{\partial x_i(0)}\Big|_{t=0} = 0$ If  $x_i(t_f)$  is free, set  $\lambda_i(t_f) - \frac{\partial \Phi}{\partial x_i(t_f)}\Big|_{t=t_f} = 0$
  - b) Periodic boundary conditions  $x_i(0) x_j(t_f) = 0$ : Set  $\lambda_i(0) - \lambda_j(t_f) + \frac{\partial \phi}{\partial x_i(0)} + \frac{\partial \phi}{\partial x_j(t_f)} = 0$

If the final time  $t_f$  is free, an additional condition has to be employed:

$$H(\boldsymbol{x}(t_f), \boldsymbol{u}(t_f), \boldsymbol{\lambda}(t_f)) = -\frac{\partial \Phi}{\partial t_f}$$

## 8.1.2 First-Order Necessary Optimality Conditions (Maximum Principle)

Functions  $x^*$ ,  $u^*$  and  $\lambda^* \neq 0$  are the optimal solution of the basis problem iff the *canonical differential equations* 

$$\dot{oldsymbol{x}} = rac{\partial H}{\partial oldsymbol{\lambda}} \qquad \dot{oldsymbol{\lambda}} = -rac{\partial H}{\partial oldsymbol{x}}$$

and the boundary conditions are fulfilled and the optimal control  $u^*$  minimizes the Hamilton function:

$$H(\boldsymbol{x}^{*}(t), \boldsymbol{u}^{*}(t), \boldsymbol{\lambda}^{*}(t)) = \min_{\tilde{\boldsymbol{u}} \in U} H(\boldsymbol{x}^{*}(t), \tilde{\boldsymbol{u}}, \boldsymbol{\lambda}^{*}(t)), \quad \forall t \in [0, t_{f}]$$
(8.1)

Where  $U \subseteq \mathbb{R}^{n_u}$  is the set of feasible controls. These conditions are called the *maximum principle*. **Definition:** The Hamiltonian is called *regular* if and only if it has an unique minimum.

If the Hamiltonian is regular and the control u appears nonlinear in the Hamiltonian, condition (8.1) can be formulated as

$$\frac{\partial H}{\partial \boldsymbol{u}} = \boldsymbol{0}$$

#### Derivation

#### 8.1.3 Second-Order Necessary Optimality Condition (Legendre-Clebsch Condition)

For stationary  $x^*$ ,  $u^*$ ,  $\lambda^*$ , the second-order necessary optimality condition, also called the *Legendre-Clebsch Condition* is that the Hessian of the Hamiltonian is positive semidefinite along the optimal state and control:

$$\boldsymbol{H}_{H}^{\boldsymbol{p}}(\boldsymbol{x}^{*}(t),\boldsymbol{u}^{*}(t),\boldsymbol{\lambda}^{*}(t)) \geq 0, \quad \forall t \in [0,t_{f}]$$

#### 8.1.4 Example

#### 8.1.5 Application: Optimal Robot Control

**The Robot Dynamics** 

#### **Optimal Control**

**Objective Functionals** 

**Necessary Conditions** 

# 8.2 Bang-Bang Singular Control

This chapter covers optimal control functions where the control function appears only linear in the Lagrangian and the motion equations:

$$L(\boldsymbol{x}(t),\boldsymbol{u}(t)) = L_0(\boldsymbol{x}(t)) + \boldsymbol{L}_1^T(\boldsymbol{x}(t))\boldsymbol{u}(t), \quad L_0: \mathbb{R}^{n_x} \to R, \, \boldsymbol{L}_1: \mathbb{R}^{n_x} \to \mathbb{R}^{n_u}$$
$$\boldsymbol{f}(\boldsymbol{x}(t),\boldsymbol{u}(t)) = \boldsymbol{f}_0(\boldsymbol{x}(t)) + \boldsymbol{f}_1^T(\boldsymbol{x}(t))\boldsymbol{u}, \quad \boldsymbol{f}_0: \mathbb{R}^{n_x} \to \mathbb{R}^{n_x}, \, \boldsymbol{f}_1: \mathbb{R}^{n_x} \to \mathbb{R}^{n_u \times n_x}$$

Also, it is assumed to have box constraints for the controls:  $u_{\min} \leq u(t) \leq u_{\max}$ .

The according Hamiltonian is

$$H(\boldsymbol{x}, \boldsymbol{u}, \boldsymbol{\lambda}) = L_0(\boldsymbol{x}(t)) + \boldsymbol{L}_1^T(\boldsymbol{x}(t))\boldsymbol{u}(t) + \boldsymbol{\lambda}^T (\boldsymbol{f}_0(\boldsymbol{x}(t)) + \boldsymbol{f}_1^T(\boldsymbol{x}(t))\boldsymbol{u})$$
  
=  $\underbrace{(\boldsymbol{L}_1^T(\boldsymbol{x}(t)) + \boldsymbol{\lambda}^T \boldsymbol{f}_1^T(\boldsymbol{x}(t)))}_{\boldsymbol{s}(\boldsymbol{x}, \boldsymbol{\lambda}) \coloneqq} \boldsymbol{u}(t) + L_0(\boldsymbol{x}(t)) + \boldsymbol{\lambda}^T \boldsymbol{f}_0(\boldsymbol{x}(t))$ 

where  $s(x, \lambda)$  is called the *switching function*. As of the maximum principle, the control u is optimal iff it minimizes the Hamiltonian. Hence, the optimal control function u is given as a piecewise function depending on the switching function:

$$u_i(t) = \begin{cases} u_{i,\min} & \text{iff } s_i\big(\boldsymbol{x}(t), \boldsymbol{\lambda}(t)\big) > 0\\ u_{i,\max} & \text{iff } s_i\big(\boldsymbol{x}(t), \boldsymbol{\lambda}(t)\big) < 0\\ u_{i,\text{sing}} & \text{iff } s_i\big(\boldsymbol{x}(t), \boldsymbol{\lambda}(t)\big) \equiv 0 \end{cases}$$

Here,  $u_{i,sing}$  is the *singular control* that is only needed if the switching function is zero over an interval. The other two cases are called *bang-bang control*.

All of the following will only cover a single control  $u_i(t)$ , but the methods can be applied to more than one control variable in a setting. Of course it is possible that some control variables are nonlinear in the Lagrangian, hence a bang-bang or singular control is not needed in that case.

## 8.2.1 Singular Control

Let  $t_1$  and  $t_2$ ,  $t_1 < t_2$  denote the start and end point of a singular control interval, respectively, i.e. the switching function is zero in that interval:

$$s_i(\boldsymbol{x}, \boldsymbol{\lambda}) = 0, \quad \forall t \in [t_1, t_2]$$

With the insight that a system has to keep itself in the singular control, also the time derivative  $s_i^{(1)} \coloneqq \frac{\mathrm{d}s_i}{\mathrm{d}t}$  of the switching function has to vanish. And also the second-order time derivative  $s_i^{(2)}$ . In fact, every time derivative has to vanish, yielding a recursive process for calculating the *m*-th time derivative:

$$s_i^{(m)} = \frac{\mathrm{d}}{\mathrm{d}t} s_i^{(m-1)} (\boldsymbol{x}(t), \boldsymbol{\lambda}(t)) \equiv 0$$

Let *m* be the smallest number of which the control  $u_i$  appears explicitly in the time derivative,  $\frac{\partial s_i^{(m)}}{\partial u_1} \neq 0$ . This yields a formula to determine  $u_{i,\text{sing}}$ :

$$s_i^{(m)}(\boldsymbol{x}(t), \boldsymbol{\lambda}(t), \boldsymbol{u}(t)) \equiv 0$$

Where  $u_i$  is part of the vector  $\boldsymbol{u}$ .

## Second-Order Necessary Optimality Condition for Singular Control

**Theorem:** If  $m < \infty$ , then m is even. Let m = 2p, where p is called the *order* of the singular control.

The second-order necessary optimality condition for singular control, i.e. the generalized Legendre-Clebsch is:

$$(-1)^{p} \frac{\partial}{\partial u_{i}} s_{i}^{(m)} \left( \boldsymbol{x}^{*}(t), \boldsymbol{\lambda}^{*}(t), \boldsymbol{u}^{*}(t) \right) \geq 0, \quad \forall t \in [t_{1}, t_{2}]$$

The condition can be expanded to a sufficient condition by replacing the inequality with a strict inequality, i.e. the left side has to be positive.

If the sufficient condition is fulfilled, two conclusions can be proven:

- If p is odd,  $u_i$  is either discontinuous or continuously differentiable in  $t_1$ .
- If p is even,  $u_i$  is continuous in  $t_1$ , but *chattering* is possible. If *chattering* occurs, the control  $u_i$  "rings" around zero until reaching  $t_1$ , i.e.  $s_i$  has infinitely many root before entering the singular section.

## 8.2.2 Application: Time-Minimal Robot Control

### 8.2.3 Notes

- If all degrees of freedom would become singular simultaneously, all adjunct variables would be zero (contradicting the maximum principle). Hence, this is not possible! For bang-bang control this means that at least one control input has to operate at its maximum or minimum.
- With the given properties, a numerically calculated solution can be checked for sanity (e.g. if all controls are singular on a section, the solution cannot be optimal).
- There are approaches to eliminate singular control beforehand so that only the number and order of switching points has to be calculated. But often a control problem can be found that necessarily has singular parts, so such methods cannot determine the optimal solution.
- The necessary conditions lead to a fully determined boundary value problem (BVP) for x and  $\lambda$  that can be solved using numerical methods (see section 9.2), but the order of bang-bang and singular control has to be known.

## 8.3 Value Function and Hamilton-Jacobi-Bellman Equation

There is a useful interpretation of the adjunct variables  $\lambda$ ! Note that by varying the initial conditions  $t_0$  and  $x(t_0) = x_0$ , different trajectories  $x^*$ ,  $u^*$ ,  $\lambda^*$  are generated by solving the optimal control problem. The *Value Function* expresses the value (of the objective) for a given initial condition  $(t_0, x_0)$ :

$$V(t_0, \boldsymbol{x}_0) = \min_{\boldsymbol{u} \in U} \left( \phi \left( \boldsymbol{x}(t_f), t_f \right) + \int_{t_0}^{t_f} L \left( \boldsymbol{x}(\tau), \boldsymbol{u}(\tau) \right) d\tau \right)$$
$$V(t_f, \boldsymbol{x}) = \phi(\boldsymbol{x}, t_f)$$

The Hamilton-Jacobi-Bellman Equation is a partial differential equation for the value function:

$$-\frac{\partial V(t, \boldsymbol{x})}{\partial t} = \min_{\tilde{\boldsymbol{u}}} \left( L(t, \boldsymbol{x}, \tilde{\boldsymbol{u}}) + \left(\frac{\partial V(t, \boldsymbol{x})}{\partial \boldsymbol{x}}\right)^T \boldsymbol{f}(t, \boldsymbol{x}, \tilde{\boldsymbol{u}}) \right)$$
$$V(t_f, \boldsymbol{x}) = \phi(\boldsymbol{x}, t_f)$$

If V(t, x) is a solution of the Hamilton-Jacobi-Bellman equation, then the minimization

$$\min_{\tilde{\boldsymbol{u}}} \left( L(t, \boldsymbol{x}, \tilde{\boldsymbol{u}}) + \left( \frac{\partial V(t, \boldsymbol{x})}{\partial \boldsymbol{x}} \right)^T \boldsymbol{f}(t, \boldsymbol{x}, \tilde{\boldsymbol{u}}) \right)$$

generates the optimal control  $u^*$  of the basis problem with  $t_0 \le t \le t_f$  with the adjunct variables:

$$\boldsymbol{\lambda}^{*}(t) = rac{\partial V(t, \boldsymbol{x}^{*}(t))}{\partial \boldsymbol{x}}$$

Hence, the adjunct variables are the gradient of the minimal objective value w.r.t. the state.

## 8.3.1 Derivation

#### 8.3.2 Notes

- If  $\lambda_i^*(t) \equiv 0$  for  $t \in [t_1, t_2]$ , the value of the objective does not depend on  $x_i^*(t)$  for the same interval.
- The value function contains all information needed for the optimal feedback control  $u^*(t, x)$ . If the value function would be available, the optimal control could be calculate for arbitrary initial conditions. But most of the time it is not available...
- It is nearly impossible to solve the HJB equation numerically (!) for relevant problems, even without constraints.

# 8.4 Constraints

It is possible to add state and control constraints to the optimal control problem in the form

$$g(x(t), u(t)) \ge 0$$

where the inequality is element-wise. Such a constraint is called a *control constraint* if u appears explicitly in g, i.e.  $\frac{\partial g}{\partial u} \neq 0$ .

#### 8.4.1 Mixed Inequality Constraints

Mixed inequality constraints are both dependent on the state and the control:

$$\boldsymbol{g}(\boldsymbol{x}(t), \boldsymbol{u}(t)) \geq \boldsymbol{0}, \quad \boldsymbol{g}: \mathbb{R}^{n_x} \times \mathbb{R}^{n_u} \to \mathbb{R}^{n_g}, \quad t \in [0, t_f]$$

This constraint is added to the Hamiltonian with multipliers  $\eta$  yielding the augmented Hamiltonian, similar to the Lagrangian of static optimization:

$$H(\boldsymbol{x}, \boldsymbol{u}, \boldsymbol{\lambda}, \boldsymbol{\eta}) = L(\boldsymbol{x}, \boldsymbol{u}) + \lambda^T \boldsymbol{f}(\boldsymbol{x}, \boldsymbol{u}) + \boldsymbol{\eta}^T \boldsymbol{g}(\boldsymbol{x}, \boldsymbol{u})$$

Analogous to the normal Hamiltonian, the augmented Hamiltonian is called *regular* along the optimal solution  $x^*(t)$ ,  $\lambda^*(t)$ ,  $\eta^*(t)$  iff  $H(x^*(t), u, \lambda^*(t), \eta^{(t)})$  has a unique minimum  $u = u^*(t)$  for all  $t \in [0, t_f]$ .

#### **Necessary Conditions (Maximum Principle)**

Let the problem be autonomous, f, g,  $\phi$  be continuously differentiable, U the set of feasible optimal controls and let the inequality constraint  $g \ge 0$  have an active component  $g_i$  with  $g_i(t) = 0$  for  $t \in [t_1, t_2]$ , i.e. let it be active on  $[t_1, t_2]$ .

If not all boundary conditions are given, they have to be determined analogous tosubsection 8.1.1.

Then the necessary optimality conditions can be defined analogous to the ones of the basis problem: Iff  $x^*$ ,  $u^*$ ,  $\lambda^*$ ,  $\eta^*$  are optimal, the canonical differential equations

$$\dot{oldsymbol{x}} = rac{\partial H}{\partial oldsymbol{\lambda}} \qquad \dot{oldsymbol{\lambda}} = -rac{\partial H}{\partial oldsymbol{x}}$$

are fulfilled and

$$\eta_j \begin{cases} = 0 & \text{iff } g_j \big( \boldsymbol{x}(t), \boldsymbol{u}(t) \big) > 0 \\ \leq 0 & \text{iff } g_j \big( \boldsymbol{x}(t), \boldsymbol{u}(t) \big) = 0 \end{cases}$$

which is, for the special constraint  $g_i$ , equivalent to:

$$\eta_j \begin{cases} = 0 & t \in [0, t_1) \cup (t_2, t_f] \\ \le 0 & t \in [t_1, t_2] \end{cases}$$

Also, the following condition have to be fulfilled on the switching points of the constraints:

$$\lim_{\delta \to 0^+} \boldsymbol{\lambda}(t_{\Delta} + \delta) = \left(\lim_{\delta \to 0^-} \boldsymbol{\lambda}(t_{\Delta} + \delta)\right) - \nu_{\Delta} \cdot \frac{\partial g_i(\boldsymbol{x}(t_{\Delta}), \boldsymbol{u}(t_{\Delta}))}{\partial \boldsymbol{x}}$$

Where  $\eta_{\Delta} = \text{const} \leq 0$  and  $t_{\Delta}$  is the entry/exit point, i.e.  $\Delta = 1$  or  $\Delta = 2$ . This implies that the adjunct variables are in general not continuous at the switching points.

The optimal control than has to minimize the augmented Hamiltonian over U:

$$H(\boldsymbol{x}^{*}(t), \boldsymbol{u}^{*}(t), \boldsymbol{\lambda}^{*}(t), \boldsymbol{\eta}^{*}(t)) = \min_{\tilde{\boldsymbol{u}} \in U} H(\boldsymbol{x}^{*}(t), \tilde{\boldsymbol{u}}, \boldsymbol{\lambda}^{*}(t), \boldsymbol{\eta}^{*}(t)), \quad \forall t \in [0, t_{f}]$$

$$(8.2)$$

Along the active constraint,

$$g_i(\boldsymbol{x}(t), \boldsymbol{u}(t)) = 0$$
 and  $\frac{\partial g_i(\boldsymbol{x}, \boldsymbol{u})}{\partial u_j} \neq 0$ 

yields a condition for determining  $u_j$  along the active constraint. All other controls as well as  $u_j$  outside of the active section have to fulfill (8.2) and the Legendre-Clebsch condition (iff *H* is regular) or the corresponding bang-bang/singular control, respectively.

As the problem is autonomous, the Hamiltonian is constant sectionally constant w.r.t. time t and it can change when constraints become active/inactive.

#### 8.4.2 State Inequality Constraints

Simple state inequality constraints

$$\boldsymbol{g}(\boldsymbol{x}(t)) \geq \boldsymbol{0}$$

that does not contain u, but at least one  $x_i$  are handled like a mixture of mixed constraints and singular control.

Let a single constraint  $g_i$  be active in the interval  $[t_1, t_2]$ , i.e.  $g_i(\boldsymbol{x}(t), \boldsymbol{u}(t)) = 0$  for  $t \in [t_1, t_2]$ . Then also the time derivative has to vanish as well as the second-order time derivative<sup>1</sup>:

$$g_i^{(m_g)} = \frac{\mathrm{d}}{\mathrm{d}t} g_i^{(m_g-1)} \left( \boldsymbol{x}(t) \right) \equiv 0$$

<sup>1</sup>Let  $g_i^{(0)} \coloneqq g_i$ .

Let  $m_g$  be the smallest number such that  $g_i^{(m_g)}$  explicitly contains the control  $u_j$ , then

$$g^{(m_g)}(\boldsymbol{x}, \boldsymbol{u}) = 0$$

yields an equation for determining  $u_j$  in the interval  $[t_1, t_2]$  and  $m_g$  is called the *order* of the state constraint. Hence, a single active constraints determined a single control variable!

#### **Augmented Hamiltonian**

Let  $\boldsymbol{g}(\boldsymbol{x}) = g(\boldsymbol{x})$ , i.e. only a single state constraint.

Mathematically, the active constraint  $g(\boldsymbol{x}(t)) = 0$  for  $t \in [t_1, t_2]$  is equivalent to

$$g^{(k)}(\boldsymbol{x}(t), \boldsymbol{u}(t)) = 0, \quad t \in [t_1, t_2]$$

for every  $k = 1, \dots, m_g$ . Hence, there are multiple, equivalent formulations of the maximum principle for active state constraints by viewing at the corresponding augmented Hamiltonian, e.g.:

$$\begin{aligned} H^{0}(\boldsymbol{x},\boldsymbol{u},\boldsymbol{\lambda}^{0},\eta^{0}) &= L(\boldsymbol{x},\boldsymbol{u}) + \left(\boldsymbol{\lambda}^{0}\right)^{T}\boldsymbol{f}(\boldsymbol{x},\boldsymbol{u}) + \eta^{0}g^{(0)}(\boldsymbol{x}) \\ &\vdots \\ H^{k}(\boldsymbol{x},\boldsymbol{u},\boldsymbol{\lambda}^{k},\eta^{k}) &= L(\boldsymbol{x},\boldsymbol{u}) + \left(\boldsymbol{\lambda}^{k}\right)^{T}\boldsymbol{f}(\boldsymbol{x},\boldsymbol{u}) + \eta^{k}g^{(k)}(\boldsymbol{x}) \\ &\vdots \\ H^{k}(\boldsymbol{x},\boldsymbol{u},\boldsymbol{\lambda}^{m_{g}},\eta^{m_{g}}) &= L(\boldsymbol{x},\boldsymbol{u}) + \left(\boldsymbol{\lambda}^{m_{g}}\right)^{T}\boldsymbol{f}(\boldsymbol{x},\boldsymbol{u}) + \eta^{m_{g}}g^{(m_{g})}(\boldsymbol{x}) \end{aligned}$$

#### **Maximum Principle**

As in the previous section, a single state constraint g is assumed.

If not all boundary conditions are given, they have to be determined analogous tosubsection 8.1.1.

The maximum principle for state constraints is similar to the one for mixed constraints. Let g be active in the interval  $[t_1, t_2]$ . Iff  $x^*$ ,  $u^*$ ,  $\lambda^{k*}$ ,  $\eta^{k*}$  are optimal, the canonical differential equations

$$\dot{oldsymbol{x}} = rac{\partial H^k}{\partial oldsymbol{\lambda}} \qquad \dot{oldsymbol{\lambda}}^k = -rac{\partial H^k}{\partial oldsymbol{x}}$$

are fulfilled and

$$\eta \begin{cases} = 0 & \text{iff } g(\boldsymbol{x}(t), \boldsymbol{u}(t)) > 0 \\ \leq 0 & \text{iff } g(\boldsymbol{x}(t), \boldsymbol{u}(t)) = 0 \end{cases} \iff \eta \begin{cases} = 0 & t \in [0, t_1) \cup (t_2, t_f] \\ \leq 0 & t \in [t_1, t_2] \end{cases}$$

Additionally  $\eta^k$  have to fulfill the sign conditions,

$$(-1)^{j} \frac{\mathrm{d}^{j}}{\mathrm{d}t^{j}} \eta^{k}(t) = \eta^{k-j}(t) \le 0, \quad j = 1, \cdots, k$$

and the final conditions

$$\lim_{\delta \to 0^{-}} \frac{\mathrm{d}^{j}}{\mathrm{d}t^{j}} \eta^{k}(t_{2} + \delta) = 0, \quad j = 0, \dots, k - 2 \text{ if } k \ge 2$$

For  $k = 1, \dots, m_g$  the active constraints have to fulfill the following at the entry point:

$$\lim_{\delta \to 0^+} \boldsymbol{\lambda}^k(t_1 + \delta) = \left(\lim_{\delta \to 0^-} \boldsymbol{\lambda}^k(t_1 + \delta)\right) - \sum_{j=1}^k \beta^j \cdot \frac{\partial g^{(j-1)}(\boldsymbol{x}(t_1))}{\partial \boldsymbol{x}}$$

with  $\beta^j \leq 0$ . And on the exit point the condition

$$\lim_{\delta \to 0^+} \lambda^k(t_2 + \delta) = \lim_{\delta \to 0^-} \lambda^k(t_2 + \delta)$$

has to be fulfilled for  $k = 1, \dots, m_g$ . That is, the adjunct variables of the  $x_i$  that appear in  $g^{(j-1)}$  are discontinuous at the entry point and continuous at the exit point of the active constraint.

The optimal control than has to minimize the augmented Hamiltonian over U:

$$H(\boldsymbol{x}^*(t), \boldsymbol{u}^*(t), \boldsymbol{\lambda}^*(t), \boldsymbol{\eta}^*(t)) = \min_{\tilde{\boldsymbol{u}} \in U} H(\boldsymbol{x}^*(t), \tilde{\boldsymbol{u}}, \boldsymbol{\lambda}^*(t), \boldsymbol{\eta}^*(t)), \quad \forall t \in [0, t_f]$$

As the problem is autonomous, the Hamiltonian is constant sectionally constant w.r.t. time t and it can change when constraints become active/inactive.

To get back and forth in different formulations for k, the following recursion can be defined. For  $m_g \ge 1$ :

$$\eta^{1}(t) = \begin{cases} \nu_{2} + \int_{t}^{t_{2}} \eta^{0}(\tau) \, \mathrm{d}\tau & \text{iff } t \in [t_{1}, t_{2}] \\ 0 & \text{else} \end{cases}$$
$$\beta^{1} = \nu_{1} + \eta^{1}(t_{1})$$

And for  $m_g \ge 2$  and  $k \ge 2$ :

$$\eta^{k}(t) = \begin{cases} \int_{t}^{t_{2}} \eta^{k-1}(\tau) \, \mathrm{d}\tau & \text{iff } t \in [t_{1}, t_{2}] \\ 0 & \text{else} \end{cases}$$
$$\beta^{k} = \eta^{k}(t_{1})$$
$$\boldsymbol{\lambda}^{k}(t) = \boldsymbol{\lambda}^{0}(t) - \sum_{j=1}^{k} \eta^{j}(t) \cdot \frac{\partial g^{(j-1)}(\boldsymbol{x}(t))}{\partial \boldsymbol{x}}$$

Where  $\nu_1, \nu_2 \leq 0$  are given by

$$\lim_{\delta \to 0^+} \boldsymbol{\lambda}^0(t_\Delta + \delta) = \left(\lim_{\delta \to 0^-} \boldsymbol{\lambda}^0(t_\Delta + \delta)\right) - \nu_\Delta \cdot \frac{\partial g_i(\boldsymbol{x}(t_\Delta), \boldsymbol{u}(t_\Delta))}{\partial \boldsymbol{x}}$$

for  $\Delta = 1$  and  $\Delta = 2$ .

#### Notes

- There are three different types of active constraints possible:
  - 1. Contact point

The states control directly "into" the infeasible region and have to do an immediate turn when having contact with the constraint (like bumping a car into a wall and bouncing off).

2. Touch point

The states just merely touch the border of the infeasible region and no changes have to be made to the control (like a a go-around with merely touching the ground).

Order $m_g$	Touch Point	Boundary Arc
1	no	yes
even	yes	yes
odd	yes	no

Table 8.1: Possibilities for active state-only constraints in optimal control given the order of the constraint.

3. Boundary arc

The states have to remain a while on the border (the "active section") until going away again (like sliding along the edge on an ice rink).

- If the Hamiltonian is regular, there are multiple possibilities shown in Table 8.1.
- Typically when working with the necessary conditions directly (analytically),  $H^0$  or  $H^{m_g}$  are used.
- The conditions and relations for  $H^k$  are primarily useful to analyze for the solutions can be transformed one into another and are as such equivalent.
- The solutions for x and u are the same for all  $H^k$ , but the  $\lambda$  and  $\eta$  might be different on the active sections (and only on these).

## 8.4.3 Examples

## **Optimal Robot Control**

## **Energy Minimization Problem**

## 8.4.4 Summary

- The conditions derivable from the maximum principle for yields a multi-point boundary value problem for  $x^*$  and  $\lambda^*$ .
- The interior switching points  $t_{S,i}$  result of:
  - Switching points at bang-bang and singular control.
  - The entry and exit into/of active sections of mixed or state constraints.
  - Other constraints at interior time steps.
  - A dynamic that is only defined piecewise.
- The switching structure (i.e. the number of order of switching points) in normally unknown!
- The switching structure has to be determined otherwise, e.g. by solving the unconstrained problem and successively tighten the constraints (continuation/homotopy technique).
- Other things like discontinuities in the system dynamics and state constraints at interior points are not covered here...
- Optimal control is hell more complex than it seems in this summary!

# 9 Calculating Optimal Trajectories

This chapter covers the numerical calculation of optimal trajectories exploiting the theoretical results of chapter 8.

# 9.1 First Computation Methods

## 9.1.1 Dynamic Programming

The technique of *dynamic programming* was introduced by Richard Bellman for numerically solving optimal control problems:

- Start from the final state  $x_f$  at  $t_f$ , the discretized Bellman Equation is used to successively iterate forward in time until t = 0.
- The effort raises exponentially with the dimensions  $n_x$  of x (the "curse of dimensionality").
- A lot of current research is based on dynamic programming, especially in the field of reinforcement learning.

## 9.1.2 Gradient Methods (Min-H Methods)

- 1. Starting with an approximation for x,  $\lambda$ ,  $0 \le t \le t_f$ , calculate an approximation for u by minimizing the Hamiltonian at discrete time steps  $t_i$ , e.g. using gradient methods.
- 2. Use this approximation for u to solve the dynamics numerically using forward integration and calculate the adjunct variables using backward iteration.
- 3. Repeat with the new approximations for x,  $\lambda$ .

# 9.2 Indirect Methods

Indirect methods are based on the necessary conditions and try to solve the arising boundary value problem.

- Advantages:
  - As it depends on the necessary conditions, the found solution is likely to be optimal.
  - By using highly precise integration methods, the solution can be determined with a high precision.
     This is especially useful, e.g. for satellite missions.
  - Increasingly powerful tools like automatic differentiation ease this method.
- Disadvantages:

- A lot of expert knowledge is needed to derive the optimality conditions and formulate the multi-point boundary value problem.
- The explicit derivation of the necessary conditions can be really costly.
- The correct, optimal switching structure is not known a-prior.
- Numerical methods need initial approximations for x and  $\lambda$  which are, especially for  $\lambda$ , hard to get.

# 9.3 Direct Methods

*Direct methods* avoid the explicit formulation of the necessary conditions by discretizing the control and possible the state variables. They transform the optimal control problem to a nonlinear optimization problem and employ methods of the nonlinear optimization.

- The user does not have to handle adjunct variables, switching structures, ...
- The efficiency of the method highly depends on the discretization and the employed nonlinear optimization problem.
- Progress in nonlinear optimization (especially SQP) yield highly efficient direct methods.
- Most commonly used.

All of the following assumes an optimal control problem given as:

$$\min_{\boldsymbol{u}} J[\boldsymbol{u}] = \phi(\boldsymbol{x}(t_0), \boldsymbol{x}(t_f), t_f)$$
  
subject to  
$$\dot{\boldsymbol{x}}(t) = \boldsymbol{f}(\boldsymbol{x}(t), \boldsymbol{u}(t)), \quad t \in [t_0, t_f]$$
$$x_i(0) = x_{i,0}, \quad i \in I_0 \subseteq \{1, \dots, n_x\}$$
$$x_j(t_f) = x_{j,f}, \quad j \in I_f \subseteq \{1, \dots, n_x\}$$
$$\boldsymbol{g}(\boldsymbol{x}(t)) \ge \boldsymbol{0}$$
$$\boldsymbol{g}(\boldsymbol{x}(t), \boldsymbol{u}(t)) \ge \boldsymbol{0}$$

The biggest difference is that this problem does not have a Lagrangian term, but as the Lagrangian can be transformed to a Mayer term, this is no loss of generality.

### 9.3.1 Direct Collocation Methods

Direct collocation methods approximate the control state variables u, x using piecewise polynomial functions. First of all, the interval  $[t_0, t_f]$  is split into  $n_s$  segments  $[t_i, t_{i+1}]$  where

$$t_i = t_0 + \tau_i \cdot (t_f - t_0)$$
 with  $0 = \tau_1 < \tau_2 < \dots < \tau_{n_s+1} = 1$ 

with all  $\tau_i$  given. The resulting segmented space is called the *mesh*. In every of these segments, u and x and approximated using a polynomial. As u is part of the integrand for x, it is reasonable to choose a higher polynomial degree for x than u.

Also let  $t_{i+1/2} \coloneqq t_i + h_i/2$  with  $h_i \coloneqq t_{i+1} - t_i$ .

#### **First Discretization, Constant Control**

A first idea is to approximate the control using constant functions

$$\tilde{u}_i(t) = u_i(t_{j+1/2}), \quad t_j \le t < t_{j+1}, \quad j = 1, \cdots, n_s, \quad i = 1, \cdots, n_i$$

and to approximate the state using linear functions:

$$\tilde{x}_i(t) = x_i(t_j) + \frac{t - t_j}{h_j} \left( x_i(t_{j+1}) - x_i(t_j) \right), \quad t_j \le t < t_{j+1}, \quad j = 1, \dots, n_s, \quad i = 1, \dots, n_x$$

The optimization variables are given as the vector

$$\boldsymbol{p} = \begin{bmatrix} \boldsymbol{x}(t_1) & \boldsymbol{u}(t_{1+1/2}) & \boldsymbol{x}(t_2) & \boldsymbol{u}(t_{2+1/2}) & \cdots & \boldsymbol{x}(t_{n_s}) & \boldsymbol{u}(t_{n_s+1/2}) & \boldsymbol{x}(t_{n_s+1}) & t_f \end{bmatrix}^T \\ = \begin{bmatrix} \boldsymbol{p}_1^x & \boldsymbol{p}_1^u & \boldsymbol{p}_2^x & \boldsymbol{p}_2^u & \cdots & \boldsymbol{p}_{n_s}^x & \boldsymbol{p}_{n_s}^u & \boldsymbol{p}_{n_s+1}^x & t_f \end{bmatrix}^T$$

where the final time  $t_f$  might be left out if it is given. The second row specifies a shorthand notation for the parameter of the approximation of x/u. The approximations shall now be determined subject to:

- Minimize the objective functional.
- Fulfill the differential equations and inequality constraints at the center of a segment.
- Comply with the initial and final conditions.

This yields the following finite-dimensional nonlinear optimization problem with  $n_p = n_s \cdot (n_x + n_u) + n_x + 1$  optimization variables (or one less if the final time is given):

$$\min_{\boldsymbol{p} \in \mathbb{R}^{n_p}} \varphi(\boldsymbol{p}), \quad \varphi(\boldsymbol{p}) = \phi(\boldsymbol{p}_1^x, \boldsymbol{p}_{n_s+1}^x, t_f)$$
  
subject to  $\boldsymbol{f}(\tilde{\boldsymbol{x}}(t_{j+1/2}), \tilde{\boldsymbol{u}}(t_{j+1/2})) - \dot{\tilde{\boldsymbol{x}}}(t_{j+1/2}) = \boldsymbol{0}, \quad j = 1, \cdots, n_s$   
 $\tilde{x}_i(t_0) = x_{i,0}, \quad i \in I_0 \subseteq \{1, \cdots, n_x\}$   
 $\tilde{x}_j(t_f) = x_{j,f}, \quad j \in I_f \subseteq \{1, \cdots, n_x\}$   
 $\boldsymbol{g}(\tilde{\boldsymbol{x}}(t_{j+1/2}), \tilde{\boldsymbol{u}}(t_{j+1/2})) \geq \boldsymbol{0}, \quad j = 1, \cdots, n_s$ 

With  $\tilde{\boldsymbol{u}}(t_{j+1/2}) = \boldsymbol{p}_j^u$ ,  $\tilde{\boldsymbol{x}}(t_{j+1/2}) = \frac{1}{2} (\boldsymbol{p}_j^x + \boldsymbol{p}_{j+1}^x)$  and  $\dot{\tilde{\boldsymbol{x}}}(t_{j+1/2}) = \frac{1}{h_j} (\boldsymbol{p}_{j+1}^x - \boldsymbol{p}_j^x)$ .

But low approximations like this have to use lots of mesh points to get good approximations. Hence, approximations of higher degree might be useful.

#### Second Discretization, Linear Control

Approximate the control using linear functions:

$$\tilde{u}_i(t) = u_i(t_j) + \frac{t - t_j}{h_j} (u_i(t_{j+1}) - u_i(t_j)), \quad t_j \le t < t_{j+1}, \quad j = 1, \dots, n_s, \quad i = 1, \dots, n_u$$

And approximate the state using linear functions:

$$\tilde{x}_i(t) = \sum_{k=0}^{3} c_{i,j,k} \left( \frac{t - t_j}{h_j} \right)^k, \quad t_j \le t < t_{j+1}, \quad j = 1, \dots, n_s, \quad i = 1, \dots, n_t$$

Here,  $c_{i,j,k}$  is the *k*-th coefficient for the *i*-th component in the *j*-th segment  $[t_j, t_{j+1}]$ , yielding four unknown parameters per component and segment. One of the parameters is determined by requiring continuity an the left side of the segment. The other three parameters are determined by requiring the fulfilling of the ODEs at three time steps in the *j*-th segment:

- Gauss Points:  $t_{j+1/2} \sqrt{3/5}h_j$ ,  $t_{j+1/2}$ ,  $t_{j+1/2} + \sqrt{3/5}h_j$ 
  - Theoretically provide the best approximation.
  - Do not provide differentiable transitions between the segments.
  - Need  $3n_s$  evaluations of the ODEs.
- Lobatto Points:  $t_j$ ,  $t_{j+1/2}$ ,  $t_{j+1}$ 
  - Provide continuously differentiable transitions between the segments.
  - Need  $2n_s + 1$  evaluations of the ODEs.

When using Lobatto points, the NLP parameters per segment can be reduced from four to two, thus reducing the collocation constraints from three to one per segment. But the remaining constraints are more nonlinear...Hence, the dimension of the NLP is similar to the constraint approximation, bit given a much better approximation of the solution.

By plugging the constraints/locations of the Lobatto points into the state approximation, this yields the following formulas for the coefficients:

$$\begin{bmatrix} c_{i,j,0} \\ c_{i,j,1} \\ c_{i,j,2} \\ c_{i,j,3} \end{bmatrix} = \begin{bmatrix} 1 & 0 & 0 & 0 \\ 0 & h_j & 0 & 0 \\ -3 & -2h_j & 3 & -h_j \\ 2 & h_j & -2 & h_j \end{bmatrix} \cdot \begin{bmatrix} p_{i,j}^x \\ p_{i,j}^x \\ p_{i,j+1}^x \\ p_{i,j+1}^x \end{bmatrix} = \begin{bmatrix} p_{i,j}^x \\ h_j p_{i,j}^x \\ -3p_{i,j}^x - 2h_j p_{i,j}^x + 3p_{i,j+1}^x - h_j p_{i,j+1}^x \\ 2p_{i,j}^x + h_j p_{i,j}^x - 2p_{i,j+1}^x + h_j p_{i,j+1}^x \end{bmatrix}$$

Where  $p_{i,j}^{\dot{x}} = \dot{x}_i(t_j) = f_i(\boldsymbol{x}(t_j), \boldsymbol{u}(t_j)) = f_i(\boldsymbol{p}_j^x, \boldsymbol{p}_j^u)$ . Hence, all coefficients can be computed using only the parameter vector

$$\boldsymbol{p} = \begin{bmatrix} \boldsymbol{x}(t_1) & \boldsymbol{u}(t_{1+1/2}) & \boldsymbol{x}(t_2) & \boldsymbol{u}(t_{2+1/2}) & \cdots & \boldsymbol{x}(t_{n_s}) & \boldsymbol{u}(t_{n_s+1/2}) & \boldsymbol{x}(t_{n_s+1}) & \boldsymbol{u}(t_{n_s+1+1/2}) & t_f \end{bmatrix}^T \\ = \begin{bmatrix} \boldsymbol{p}_1^x & \boldsymbol{p}_1^u & \boldsymbol{p}_2^x & \boldsymbol{p}_2^u & \cdots & \boldsymbol{p}_{n_s}^x & \boldsymbol{p}_{n_s}^u & \boldsymbol{p}_{n_s+1}^x & \boldsymbol{p}_{n_s+1} \end{bmatrix}^T$$

with one parameter more than the constant control approximation. That is,  $n_p = (n_s + 1)(n_x + n_u) + 1$  parameters (or one less if the final time is fixed). This yields the following nonlinear optimization problem:

$$\min_{\boldsymbol{p} \in \mathbb{R}^{n_p}} \varphi(\boldsymbol{p}), \quad \varphi(\boldsymbol{p}) = \phi(\boldsymbol{p}_1^x, \boldsymbol{p}_{n_s+1}^x, t_f)$$
  
subject to  $\boldsymbol{f}(\tilde{\boldsymbol{x}}(t_{j+1/2}), \tilde{\boldsymbol{u}}(t_{j+1/2})) - \dot{\tilde{\boldsymbol{x}}}(t_{j+1/2}) = \boldsymbol{0}, \quad j = 1, \cdots, n_s$   
 $\tilde{x}_i(t_0) = x_{i,0}, \quad i \in I_0 \subseteq \{1, \cdots, n_x\}$   
 $\tilde{x}_j(t_f) = x_{j,f}, \quad j \in I_f \subseteq \{1, \cdots, n_x\}$   
 $\boldsymbol{g}(\tilde{\boldsymbol{x}}(t_{j+1/2}), \tilde{\boldsymbol{u}}(t_{j+1/2})) \ge \boldsymbol{0}, \quad j = 1, \cdots, n_s$ 

With  $\tilde{\boldsymbol{u}}(t_{j+1/2}) = \boldsymbol{p}_j^u$ ,  $\tilde{\boldsymbol{x}}(t_{j+1/2}) = \frac{1}{2} (\boldsymbol{p}_j^x + \boldsymbol{p}_{j+1}^x)$  and  $\dot{\tilde{\boldsymbol{x}}}(t_{j+1/2}) = \frac{1}{h_j} (\boldsymbol{p}_{j+1}^x - \boldsymbol{p}_j^x)$ .

#### **Discretization of Inequality Constraints**

In the NLP, the inequality constraints are evaluated at the points  $t_j$ , i.e. the mesh points. But why not at other points, e.g.  $g(t_{j+1/2})$ ? As of the theoretical results, one control is determined by

$$g^{(m_g)}(\boldsymbol{x}(t), \boldsymbol{u}(t)) = \frac{\mathrm{d}^{m_g}}{\mathrm{d}t^{m_g}}g(\boldsymbol{x}(t)) \equiv 0$$

in an active segment  $t \in [t_1, t_2]$ . Hence, for consistency, the number of degrees of freedom on a potential arc have to equal the number of degrees of freedom of the discretized control constraints. This is only ensured when the inequality constraints are evaluated at the mesh points  $t_j$ .
### Sparsity

In direct collocation methods, the resulting NLPs have sparse gradients and Jacobians. For example, the memory complexity of the full Jacobian is  $O(n_s^2)$ , while the non-zero elements rise linearly. By exploiting this sparse structure, the efficiency for NLP solvers can be highly enhanced!

### **Convergence Properties and Solution Validation**

By studying the Lagrangian of the NLP, it is possible to show that a gradient method that is using the Lagrangian for determining the step size, the maximum principle is taken into account in a discretized matter. That way, using direct collocation methods, it is possible to compute approximations for the adjunct variables and the multipliers  $\eta^0$ . Hence, the theoretical optimality conditions can be validated afterwards!

For autonomous problems, the Hamiltonian has to be sectionally constant which can be validated using the calculated approximations  $\tilde{a}(t)$ ,  $\tilde{u}(t)$ ,  $\tilde{\lambda}^0(t)$ ,  $\tilde{\eta}^0(t)$ . Another thing that can be validates is whether the initial and final conditions for the adjunct variables are fulfilled.

### Derivation

### **Choosing the Mesh Points**

Direct collocation methods typically start with a rough grid and a rough approximation for  $x^*$  and  $u^*$  that is then successively adjusted until the solution is smooth enough. But how to do this mesh adjustments?

- Ideally: A segment-wise approximation error  $\|\tilde{\boldsymbol{x}}(t) \boldsymbol{x}^*(t)\|, \|\tilde{\boldsymbol{u}}(t) \boldsymbol{u}^*(t)\|$ .
- Alternatively: Use the error of in the ODE and inequality constraints:  $\|\boldsymbol{f}(\tilde{\boldsymbol{x}}(t), \tilde{\boldsymbol{u}}(t)) \dot{\tilde{\boldsymbol{x}}}(t)\|, \|\boldsymbol{g}_{-}(\tilde{\boldsymbol{x}}(t), \tilde{\boldsymbol{u}}(t))\|$ with  $g_{-} \coloneqq |g|$  as the element-wise absolute value.

And refine the mesh until  $\|\cdot\| \leq \varepsilon_{tol}$ .

**Segment-Wise Error Estimation** Assuming  $\tilde{u}(t) = u^*(t)$ , theories of collocation methods can be used to estimate the error

$$\|\tilde{\boldsymbol{x}}(t) - \boldsymbol{x}^*(t)\|$$

where  $x^*$  is the "real" solution of the BVP. But this method does not work well for direct collocation methods.

**Segment-Wise Error Estimation of Time-Continuous Constraints** Simple strategy that works good in practice: Check the fulfilling of the constraints at test points, e.g. for the second discretization:

$$\left\|\boldsymbol{f}\big(\tilde{\boldsymbol{x}}(t),\tilde{\boldsymbol{u}}(t)\big) - \dot{\tilde{\boldsymbol{x}}}(t)\right\| \quad \left\|\boldsymbol{g}_{-}\big(\tilde{\boldsymbol{x}}(t),\tilde{\boldsymbol{u}}(t)\big)\right\| \quad \text{for} \quad t = t_{j+k/4} = t_j + \frac{k}{4}h_j, \quad k = 1, 2, 3$$

Then successively split the segments with too big errors into smaller segments.

**Segment-Wise Estimation of Optimality Error** The previous error estimation methods only take the constraints into account, not that this is an optimal control problem. But it is also possible to approximate the "optimality error" using the calculated values for the adjunct variables (e.g. using quadrature rules):

$$L = \phi \left( \tilde{\boldsymbol{x}}(0), \tilde{\boldsymbol{x}}(t_f), t_f \right) + \sum_{j=1}^{n_s} \underbrace{\int_{t_j}^{t_{j+1}} \left[ \tilde{\boldsymbol{\lambda}}^T(t) \cdot \left( \boldsymbol{f} \left( \tilde{\boldsymbol{x}}(t), \tilde{\boldsymbol{u}}(t) \right) - \dot{\tilde{\boldsymbol{x}}}(t) \right) + \tilde{\boldsymbol{\eta}}^T(t) \cdot \boldsymbol{g} \left( \tilde{\boldsymbol{x}}(t), \tilde{\boldsymbol{u}}(t) \right) \right] \mathrm{d}t}_{\boldsymbol{\lambda}}$$

Optimality Error in the *j*-th Segment

### Implementation

The general approach to implement direct collocation methods is shown inalgorithm 7.

Algorithm 7: Direct Collocation Algorithm.

1 Initialization: Choose a start grid  $(\tau_j^{(0)})_{j=1}^{n_s^{(0)}+1}$  and initial values  $p_j^{x,(0)} = x(t_j)$ ,  $p_j^{u,(0)} = u(t_{j+1/2})$ 2 while not converged (error in  $\dot{x}$ , maximum iterations or maximum mesh size) do

3 Adjust the mesh: 
$$(\tau_j^{(k)})_{j=1}^{n_s^{(k)}+1}$$
  
4 Solve the resulting NLP:  

$$\min_{\boldsymbol{p} \in \mathbb{R}^{n_p}} \varphi(\boldsymbol{p})$$
subject to  $\boldsymbol{a}(\boldsymbol{p}) = \boldsymbol{0}$   
 $\boldsymbol{b}(\boldsymbol{p}) \ge \boldsymbol{0}$   
5  $k \leftarrow k+1$ 

### Notes

- For a rough mesh, the tolerances  $\varepsilon_{opt}$  and  $\varepsilon_{ft}$  do not have to be chosen too tight as the NLP solver might not ding a solution then.
- While successively refining the grid, the tolerances can be set to tighter values.
- Ideally, the most SQP iterations are needed for the first two meshes. For the high-dimensional NLPs for fine grids, only a few iterations are needed.
- System parameters can be optimized simultaneously by adding them to the variable vector *p*.

### **Application: Optimal Robot Control**

### 9.3.2 Direct Shooting Methods

Similar to the direct collocation methods, shooting methods rely on a segmentation of the the interval  $[t_0, t_f]$  into  $n_s$  segments  $[t_j, t_{j+1}]$  with  $n_s + 1$  mesh point  $t_j$ . Then the control u is approximated per section, e.g. using constant or linear functions.

Direct shooting methods then solve various initial value problems until one fulfills the boundary conditions. Informally speaking, multiple trajectories are "shot into the room" until one is feasible. More formally, the process is:

- 1. Divide the interval into  $n_s$  segments and approximate the control  $\boldsymbol{u}$ , e.g. using a piecewise linear function  $\tilde{\boldsymbol{u}}$ . Then parameterize  $\boldsymbol{u}$  by  $\boldsymbol{p} = \begin{bmatrix} \boldsymbol{u}(t_1) & \cdots & \boldsymbol{u}(t_f) & t_f \end{bmatrix}^T \in \mathbb{R}^{n_u \cdot (n_s + 1)}$ .
- 2. Simulate the movement by forward-integrating the system dynamics starting from the initial values and by using the approximated control  $\tilde{u}$ . This gives an approximation of the state,  $\tilde{x}$ .
- 3. Calculate the objective and the violations of the constraints, especially the final constraints.
- 4. Optimize p such that the objective is minimized while fulfilling the constraints. Repeat.

But the calculation of the NLP gradients and the Jacobians needs calculating the sensitivity matrix  $\frac{\partial x(t;p)}{\partial n}$ .

- Advantages:
  - The resulting dimension of the NLPs are lot less than for direct collocation methods.
  - In every iteration of the NLP solver, a solution of the ODEs is available.
- Disadvantages:
  - Solving the initial value problem may highly depend on the approximations  $\tilde{u}$ ,  $\tilde{t}_f$ . To increase the robustness, multiple shooting methods can be used.
  - Approximations for the adjunct variables are not as natural as in direct collocation methods, but possible.

# 9.4 Notes

- Some big disadvantages of indirect methods, that good approximations of the optimal state- and adjunct variables and the switching structure are needed, can be overcome by using a direct method first.
- The majority of optimal control problems are nowadays solved using direct methods (other approaches exist besides direct collocation and shooting methods).
- Important applications are for example, aerospace engineering, robots, vehicles, process engineering, economy, biology, ecology, . . .
- Direct collocation methods are also called "simultaneous simulation and optimization".
- Direct shooting methods are also called "iterative simulation and optimization".
  - Two different algorithms are used for simulation and optimization.
  - Highly efficient and stable calculation of the sensitivity matrix using special integration techniques.

# **10 Optimal Feedback Control**

Applying the computed control trajectory directly as an open loop (feedforward) control as illustrated in Figure 10.1 would cause growing divergence from the nominal state  $x_d^*(t)$ . This can be caused by e.g.:

- Errors in the model (complex models cannot be 100% accurate).
- Noise: During the run of a dynamic process noise may affect the systems behavior.

Hence, the wanted end state might not be reached when just using feedforward control.

# **10.1 Classical Feedback Control (Position Control)**

It seems obvious to just use a classical position control for the nominal stater state trajectory  $x^*(t)$  (a set point trajectory control) as illustrated in Figure 10.2. But this causes "ringing" around the nominal trajectory, e.g. using a PID control law. Additionally, using a position control has more flaws:

- An individual control for each state component is merely possible.
- The quality of the position control depends on the desired states and the control parameters.
- Returning to the nominal trajectory using control laws is not the optimal trajectory from the disturbed state as the initial state to the final state!
- Also, using position control can cause violations of the constraints, e.g. a bang-bang control always operates on the border of the control constraints. Hence, returning to the nominal trajectory might cause overshooting this constraints.

# **10.2 Optimal Feedback Control**

If the real trajectory at time step  $t_1$  differs from the nominal state  $x^*(t_1; x_0)$  that was calculated starting from  $x_0$  it would be optimal to now following a new trajectory with the initial state  $x_d(t_1)$ . But how to calculate this new trajectory? These computations have to be done anytime. But data from the old trajectory can be used for faster calculations!

These re-computations would not be necessary if the optimal control trajectory would not be calculated as a function of time  $u^*(t)$  (as a feedforward control), but as a function of the initial state  $u^*(x_0)$ , i.e. as a feedback control. Sadly, it is nearly impossible to calculate this function except for some special cases...



Figure 10.1: Feedforward Control



Figure 10.2: Feedback Control

# 10.3 Linear Quadratic Regulator (LQR)

Linear systems with quadratic objective are the only system for which the optimal feedback control  $u^*(x_0)$  can be computed in closed form! A linear system with quadratic objective is given as

$$\begin{split} \min_{\boldsymbol{u}} J[\boldsymbol{u}], \quad J[\boldsymbol{u}] &= \int_{t_0}^{t_f} \left( \boldsymbol{x}^T(t) \boldsymbol{Q} \boldsymbol{x}(t) + \boldsymbol{u}^T(t) \boldsymbol{\Gamma} \boldsymbol{u}(t) \right) \mathrm{d}t \\ \dot{\boldsymbol{x}}(t) &= \boldsymbol{A}(t) \boldsymbol{x}(t) + \boldsymbol{B}(t) \boldsymbol{u}(t) \\ \text{subject to} \qquad \boldsymbol{x}(t) &= \boldsymbol{x}_0 \\ & \boldsymbol{x}(t_f) \quad \text{free} \end{split}$$

with a symmetric and positive definite matrix Q and a diagonal Matrix  $\Gamma = \text{diag}(\gamma_1, \dots, \gamma_{n_u})$ . The optimal control law is then given as

$$oldsymbol{u}(oldsymbol{x}) = -oldsymbol{\Gamma}^{-1}oldsymbol{B}^T\!(t)oldsymbol{P}(t)oldsymbol{x}$$

where the matrix P(t) is found by solving the matrix-Riccati ODE

$$\dot{\boldsymbol{P}}(t) + \boldsymbol{Q} + \boldsymbol{A}^{T}(t)\boldsymbol{P}(t) + \boldsymbol{P}(t)\boldsymbol{A}(t) - \boldsymbol{P}(t)\boldsymbol{\Gamma}^{-1}\boldsymbol{B}(t)\boldsymbol{B}^{T}(t)\boldsymbol{P}(t) = \boldsymbol{O}$$

with the boundary condition  $P(t_f) = O$ .

If the system is time invariant, i.e.  $\dot{A} = O$ ,  $\dot{B} = O$  or  $t_f = \infty$ , it follows  $\dot{P} = O$  and hence P is given by solving the algebraic matrix-Riccati equation:

$$\boldsymbol{Q} + \boldsymbol{A}^T \boldsymbol{P} + \boldsymbol{P} \boldsymbol{A} - \boldsymbol{P} \boldsymbol{\Gamma}^{-1} \boldsymbol{B} \boldsymbol{B}^T \boldsymbol{P} = \boldsymbol{O}$$

As P is symmetric, these are  $n_x \cdot (n_x + 1)/2$  equations to determine every element of P.

These matrix-Riccati algebraic/differential equations can be solved efficiently using special numerical methods. LQR systems can approximate a lots of systems and are therefore more useful than it appears at first (e.g. active damping in cars or linearized inverted pendulum). In theory, LQR methods can be applied to any system using an infinite-dimensional nonlinear embedding ("measurement") which can be finitely approximated, still yielding good results. Another approach is iterative LQR (iLQR) that linearizes the system in every time step.

### 10.3.1 Derivation

### **10.4 Neighboring Extremals**

TO approximately the new trajectory  $x^*(t; x_a(t_1))$ , all information that have already been expended for calculating the trajectory  $x^*(t; x_0)$  can be used. These updates have to be be done continuously.

### 10.4.1 Indirect Methods

Indirect methods calculates a solution trajectory  $x^*(t)$ ,  $u^*(t)$ ,  $\lambda^*(t)$ ,  $\eta^*(t)$  as the numerical solution of the multi-point boundary value problem rising from the necessary optimality conditions. By Taylor-expanding the disturbed trajectory  $x^*(t; x_0 + \epsilon_0)$  around the nominal trajectory yields a linear-quadratic optimal control problem (called "accessory minimum problem") for

$$\delta \boldsymbol{x}(t), \quad \delta \boldsymbol{u}(t), \quad \delta \boldsymbol{\lambda}(t)$$

Different variants of this are possible, e.g. the repeated correction method:

- 1. Calculate  $x^*(t)$ ,  $u^*(t)$ ,  $\lambda^*(t)$ ,  $\eta^*(t)$  using the multi-point BVP width multiple-shooting methods.
- 2. Feedback schema for noisy nominal trajectory  $\boldsymbol{u}^*(t_0; \boldsymbol{x}_0 + \boldsymbol{\epsilon}_0) = \boldsymbol{u}^*(t_0; \boldsymbol{x}_0) + \delta \boldsymbol{u}(t_0)$  containing the nominal control and a correction term.
- 3. Calculate the control correction using:  $\Delta u(t_0) = K_1(t_0) \cdot \partial x(t_0) + K_2(t) \cdot \partial \lambda(t_0)$
- 4. Apply this repeatedly for different time steps  $t_0 = t_1$ .
- Advantages:
  - Fast, numerical calculation.
  - First-order optimal trajectories.
  - Can be used to correct noise in the model parameters of the ODE.
- Disadvantages:
  - Only locally applicable along a nominal trajectory (i.e. "little" noise that does not change the switching structure)
  - Depends on the multi-point BVP of the necessary conditions. Hence, application is time consuming as the ODEs have to be formulated.

### 10.4.2 Direct Methods

Direct methods solve the problem by discretizing the problem, resulting in a NLP that contains the initial value  $x(t_0)$  in its equality constraints. Studying the dependency on this parameters leads to sensitivity and stability analysis of nonlinear optimization problems.

For direct collocation methods it can be shown that the linear-quadratic optimal control problem is analogous to a quadratic problem (QP) for calculating the corrections  $p^*(x_0 + \epsilon_0) = p^*(x_0) + \delta p$ .

### 10.4.3 Nonlinear Model Predictive Control (NMPC)

For nonlinear model predictive control, a series of optimal control problems is solved for a varying time horizon  $P_j$ . These methods are commonly used for chemical processes with slow reaction time and aerospace engineering, but not so often in fast environments like robot movements.

# 10.5 Numerical Synthesis of the Nonlinear Feedback Control

Direct collocation methods can compute optimal control problems with different initial conditions  $x_0$  pretty fast and robust. Using an appropriate "step size control", multiple neighboring trajectories can be computed kind of automatic.

- First approach for synthesis  $u^*(x)$ : Calculate  $u^*$  using the HJB equation where  $x^*(t)$ ,  $\lambda^*(t)$  are approximated using reference trajectories for  $x_a(t_1)$ .
- Second approach: Approximate  $u^*(t)$ , e.g. an approximation that has been trained on lots of open-loop optimal trajectories.

# **11 Further Topics on Optimal Control**

## **11.1 Inverse Optimal Control**

Inverse optimal control approaches the following problem:

- Given a dynamic process x
   *x* = *f*(*x*(t), *u*(t)) and measurements of a run *x<sub>k</sub>* = *x*(t<sub>k</sub>) + ε<sub>k</sub> that is in general noisy,
- Assume the run was controlled optimally,
- Find the objective function that was used.

One possible approach is the linear combination of multiple "basis" functionals  $J_k[u]$ , e.g. one for minimum energy and one for minimum time:

$$J^*[oldsymbol{u}] = \sum_{k=1}^{n_J} \omega_k J_k[oldsymbol{u}]$$

This yields a two-level optimization problem:

- "Outer" problem: Finite-dimensional, nonlinear optimization problem:  $\Phi(\boldsymbol{\omega}) = \sum_{k=1}^{n_m} \|\boldsymbol{x}^*(t_k, \boldsymbol{\omega}) \boldsymbol{x}_k\|_2^2$
- "Inner" problem: Constraints of the outer problem; Optimal control problem with solution  $x^*$

$$\min_{\boldsymbol{u}} J^*[\boldsymbol{u}], \quad J^*[\boldsymbol{u}] = \sum_{k=1}^{n_J} \omega_k J_k[\boldsymbol{u}]$$
subject to  $\dot{\boldsymbol{x}}(t) = \boldsymbol{f}(\boldsymbol{x}(t), \boldsymbol{u}(t))$ 

### 11.2 Differential/Dynamic Games

Differential games have ODEs for the complete state x containing multiple ODEs for every player and controls  $v_1(t), \dots, v_{n_p}(t)$  for every placer. Additionally they might have state or control constraints. The players then might minimize or maximize one or more objective cooperative or non-cooperative.

### 11.2.1 Non-Cooperative Two-Player Zero-Sum Differential Games

A common class of differential games are non-cooperative two-player zero-sum differential games where:

- Non-cooperative means that one player tries to maximize the objective and the other tries to minimize it.
- Zero-sum means that the loss of one player is the reward of the other and vice versa.

Necessary conditions for these games can be derived from the minimax principle similar to the maximum principle using the ODEs for both players, adjunct differential equations and a Hamiltonian.

This also yields a multi-point BVP that can be solved numerically using indirect methods and it is also possible to use direct methods.

These games can be used, e.g. to generate robust trajectories by letting one player "be the noise" or friction or similar that tries to destroy the optimal trajectory.

### Example

# **11.3 Learning Methods and Optimization**

Learning methods (from machine learning) often rely on formulations of solving specific optimization problems. Hence, machine learning and optimization is highly coupled. The focus is a little bit different as most optimization methods try to find solutions as accurate as possible while ML algorithms try to generalize as best as possible to handle unseen scenarios.

### 11.3.1 Foundations

The underlying task often is to find a model  $f_{\theta}(x)$  that has a specific input/output behavior that is enforced using a *loss function*, e.g. the least squares loss function. These models are then optimized using basic gradient methods, gradient descent, without step size determination but using a small step size or an adaptive one (e.g. adam, adagrad, adadelta). A common function approximator are neural networks which use multiple *layers*, *activation functions* and *weights* to approximate arbitrary functions (in fact, any function can be approximated using a two-layer neural net).

Some open questions are:

- Why and when does gradient descent work and how fast?
- Why does not training not always cause overfitting (even when the number of parameters are much higher than the number of samples)?
- How to interpret the trained models (explainable AI)?

### 11.3.2 Reinforcement Learning

*Reinforcement learning* is highly related to optimal control in terms if a reward that is maximized by tweaking controls etc. However, optimal control is often discrete and models are often treated as stochastic models allowing reasoning about uncertainty. It is also possible to not have any model (e.g. in model-free reinforcement learning)! The model is then learned implicitly.

### Reward

In the RL setting, it is possible to study both finite and infinite time horizons using discounted rewards ensuring that the sum of rewards converges.

### **Value Function**

In RL, a policy  $\pi$  is searched which is commonly dependent on the state, not time:  $\pi = \pi(x)$ . The *value function* describes the discounted reward starting from the current state:

$$V^{\pi}(\boldsymbol{x}_{0}) = \sum_{k=0}^{\infty} \gamma^{k} r_{k} (\boldsymbol{x}_{k}, \pi(\boldsymbol{x}_{k}))$$

The *state-action function* describes the discounted reward of the current state if a specific action is taken next and following the policy afterwards:

$$Q^{\pi}(oldsymbol{x}_0,oldsymbol{u}_0)=r_1(oldsymbol{x}_0,oldsymbol{u}_0)+\sum_{min}^{max}$$

### **First Approach**

A first approach would be to learn the control directly by maximizing the reward J using gradient descent:

$$\boldsymbol{\nabla}_{\boldsymbol{\theta}} J = \frac{\partial J}{\partial \pi_{\boldsymbol{\theta}}} \frac{\partial \pi_{\boldsymbol{\theta}}}{\partial \boldsymbol{\theta}}$$

But this requires a gradient of the objective...The approximation may have high variances. Additionally, new gradient is computed independently of old approximations. Hence, no learning happens.

#### Learning the Value Function

It is better to solve the optimization problem

$$oldsymbol{u}^*(oldsymbol{x}) = rg\max_{oldsymbol{u}} \, Q^{\pi}(oldsymbol{x},oldsymbol{u})$$

by approximating the value function  $V(\mathbf{x}_0)$  with a function  $V_{\boldsymbol{\theta}}(\mathbf{x}_0)$ , e.g. with *Temporal Difference Learning* (TD).

This uses the Bellman equation

$$V(\boldsymbol{x}) = \max_{\boldsymbol{x}} r(\boldsymbol{x}, \boldsymbol{u}) + \gamma V(\boldsymbol{x}, \boldsymbol{u}) = \max_{\boldsymbol{u}} Q(\boldsymbol{x}, \boldsymbol{u})$$

which is a discretized version of the H-J-B equation. Temporal difference learning measures the TD-error  $\delta_n = r(\mathbf{x}, \mathbf{u}) + V_{\theta}(\mathbf{x}_n) - \gamma V_{\theta}(\mathbf{x}_{n+1})$  and uses gradient descent to update the parameters  $\theta$ . An instance of this class of algorithms is Q-Learning.

- Advantages:
  - Small variance in the approximations of the expected reward.
- Disadvantages:
  - In every state, an optimization problem has to be solved. The overhead can be reduced by discretizing the control and just trying out all possibilities.
  - No guarantee of convergence (but this holds for every algorithm proposed in this course as it may guarantee convergence but convergence to a poor local minimum).

### Actor-Critic

Another approach is to approximate the control and the value function at the same time:

- The *actor* generates control u for a given state x and
- the *critic* estimates the quality of the current control and learns the value functions that is used to tweak the control parameters.

This method combined the advantages of the previous approaches as it does not require to solve an optimization problem in each step and no control discretization. It also allows a more precise measurement of the gradients.

### Notes

- All proposed methods are based on an approximation of the value function or the control. In practice, neural networks are often used for these approximations.
- Stability proves have been made for systems with special structures (e.g. affine dynamics).
- Better convergence than actor-only methods.
- The approach can be extended to models with continuous time.