

# Some popular nonlinear programming algorithms

Mohamed Tarek<sup>1,2</sup>

<sup>1</sup>PumasAI Inc., USA

<sup>2</sup>UNSW Canberra, Australia

October 2021

Many design and decision-making tasks can be formulated as a constrained nonlinear optimization problem with an objective, equality constraints, inequality constraints and bound constraints on the decision variables. The problem can then be solved with any of the many constrained nonlinear optimization (aka nonlinear programming) algorithms available in literature and implemented in software. In this document, I attempt to explain three common mathematical optimization algorithms for solving constrained nonlinear optimization problems. These algorithms are:

- The original method of moving asymptotes (Svanberg, 1987).
- The primal-dual interior point optimization algorithm available in IPOPT as explained in Wächter and Biegler (2006).
- The augmented Lagrangian family of algorithms.

These are all extremely popular optimization algorithms with one thing in common: there is no sufficiently general yet somewhat accessible explanation of these algorithms (to my knowledge) beyond their original papers or some large textbook which might be too far out of reach for the general user base of these algorithms. For the benefit of the readers who are not familiar with nonlinear programming (NLP) concepts and terminology, I also present a short primer on nonlinear programming concepts and formulations in this document.

All of the algorithms discussed here are available in the `Nonconvex.jl`<sup>1</sup> software package in the Julia programming language. I maintain `Nonconvex.jl` and this document might grow to become a technical companion to the software. For code examples, please refer to the documentation of `Nonconvex.jl`.

---

<sup>1</sup><https://github.com/JuliaNonconvex/Nonconvex.jl>

# Contents

<b>1</b>	<b>Target audience and disclaimers</b>	<b>3</b>
<b>2</b>	<b>Background</b>	<b>4</b>
2.1	Convex vs. non-convex . . . . .	4
2.1.1	Convexity . . . . .	4
2.1.2	Conic constraints . . . . .	5
2.1.3	Non-convexity . . . . .	6
2.2	Formulations . . . . .	7
2.2.1	A linear objective is all you need . . . . .	8
2.2.2	Slack variables . . . . .	8
2.2.3	Equality constraints and empty interiors . . . . .	9
2.3	Regularity conditions . . . . .	10
2.4	Sufficient optimality conditions for regular points . . . . .	11
<b>3</b>	<b>Method of moving asymptotes</b>	<b>13</b>
3.1	Separable convex approximation . . . . .	13
3.2	Efficient primal-dual algorithm . . . . .	14
3.3	Algorithm strengths . . . . .	15
<b>4</b>	<b>Primal-dual interior point method</b>	<b>17</b>
4.1	Original problem . . . . .	17
4.2	Barrier problem . . . . .	17
4.3	Damped Newton's method . . . . .	20
4.4	Line search and filter . . . . .	22
4.4.1	Maximum step size . . . . .	22
4.4.2	Line search . . . . .	23
4.4.3	Filter set . . . . .	23
4.4.4	Second order correction . . . . .	24
4.4.5	Accelerating heuristics . . . . .	25
4.4.6	Clamping . . . . .	25
<b>5</b>	<b>Augmented Lagrangian algorithm</b>	<b>26</b>

# 1 Target audience and disclaimers

The target audience of this document is people with a bit of a mathematical background but with little to no constrained optimization background, and people who are (future) users or (future) developers of nonlinear optimization software. This document was originally part of my PhD thesis on topology optimization but I decided to present it separately here because then it might be useful to a larger audience. For this reason, there will be references to challenges in topology optimization problems which the general readers are free to ignore. For the benefit of the readers who are not familiar with nonlinear programming (NLP) concepts and terminology, I also present a short primer on nonlinear programming concepts and formulations. For readers familiar with basic concepts of NLP, feel free to jump straight to section 3.

In this document, I attempt to make my explanations:

1. Accessible but without sacrificing generality as much as I can, and
2. Correct and detailed but without too many derailing proofs and obscure mathematical language.

Correctness and accessibility are tough to combine in one document. This is my attempt.

Completeness is not necessarily a goal here if it comes at the cost of accessibility (in my subjective opinion) and if I find that there is a sufficiently clear reference I can point the interested readers to. Finally, in this document I focus on "what the algorithms are" more than the reasons "why they are what they are", although I also give reasons and proofs whenever the reasons/proofs can be concisely described (mostly) in words.

In this document, the terms *linear* and *affine* will be used inter-changeably in the context of classifying a function/constraint as either a linear/affine or a nonlinear function/constraint. Even though there is a distinction between the 2 terms, optimization literature commonly uses these 2 terms inter-changeably for the same purpose.

## 2 Background

### 2.1 Convex vs. non-convex

Optimization is a field filled with seemingly cryptic but fundamentally simple and useful concepts. Two such concepts are:

1. Convexity
2. Conic constraints

#### 2.1.1 Convexity

A set of points  $S = \{\mathbf{x}\}$  is said to be convex if every weighted average of some points inside the set is also inside the set. For sets in 2D and 3D, one can visually check for the convexity of a set  $S$ . If you can never draw a line of finite length with a start point and an end point inside the set  $S$  but part of the line lies outside the set, this set is convex. Your inability to draw such a line implies that no such line exists which implies that the set is convex.

However, set convexity generalizes to more dimensions than 2 or 3. More generally, one can formally define convexity as follows. A set  $S$  is convex if for every  $\mathbf{x}_1, \mathbf{x}_2 \in S$ ,  $\alpha\mathbf{x}_1 + (1 - \alpha)\mathbf{x}_2 \in S$  for all  $0 < \alpha < 1$ . This is a straightforward generalization of the line test in 2D/3D. It is also simple to show that the intersection of any number of convex sets is either the empty set or another convex set.

That's set convexity. Function convexity is a slightly different but related concept. A function  $f(\mathbf{x})$  is called convex if the set of points satisfying the constraint  $f(\mathbf{x}) \leq c$  is either an empty or a convex set for every  $c \in \mathbb{R}$ . Note that the set of such points doesn't have to be a compact set, that is it can extend to  $\pm\infty$  along one or more directions. Showing that a function is convex can be done by proving the above property. This however is only one way to check if a function is convex. Another common check for twice differentiable functions is to check if the function's Hessian is positive semidefinite everywhere in the domain of the function. Identifying and building convex functions in a disciplined way is a well-studied field, often termed "disciplined convex programming" (DCP). For more on DCP, the readers are advised to check the excellent online resources available in <https://dcp.stanford.edu>.

In optimization, a constraint is said to be a convex constraint if the set of points satisfying this constraint is a convex set. One common way to construct a convex constraint is using a convex function  $f(\mathbf{x})$ . If  $f$  is a convex function, the following constraint is convex:

$$f(\mathbf{x}) \leq c \tag{1}$$

for any constant  $c$ , assuming there exists at least 1 point satisfying this constraint. Note that it has to be a  $\leq$  constraint. Other common convex constraints are linear constraints:

$$f(\mathbf{x}) = 0 \tag{2}$$

where  $f$  is an affine function of  $\mathbf{x}$ , e.g.  $f(\mathbf{x}) = \mathbf{b}^T \mathbf{x} + a$  for some constants  $\mathbf{b}$  and  $a$ .

A constrained optimization problem with a convex *minimization* objective function and convex constraints is called a "convex program". Formulating decision problems as convex programs and solving them with convex optimization algorithms is often called "convex programming" which is a sub-field of "nonlinear programming". Convex programs are theoretically appealing because under mild assumptions, they have efficient algorithms that can find the so-called *global optimal solution*, which is a (not necessarily unique) feasible point  $\mathbf{x}$  that absolutely minimizes the objective function, i.e. no other solution can do strictly better.

Note that it has to be a minimization problem. Maximizing a general convex function is NP-hard! If your objective  $f(\mathbf{x})$  is a quantity you want to maximize, that's equivalent to minimizing  $-f(\mathbf{x})$  instead. If  $-f(\mathbf{x})$  is a convex function,  $f$  is known as a concave function. Therefore for a maximization optimization problem to be convex, the objective function must be concave and the constraints must be convex. This might be counter-intuitive that the maximization objective function must be *concave* for the optimization problem to be *convex*, but hey no one said optimization jargon was intuitive!

Similarly, the following constraint:

$$f(\mathbf{x}) \geq c \tag{3}$$

is a convex constraint if  $f$  was a concave function since this constraint is equivalent to:

$$-f(\mathbf{x}) \leq -c \tag{4}$$

where  $-f(\mathbf{x})$  is convex which has the correct  $\leq$  form we used above. Again, it's slightly counter-intuitive but just go with it.

### 2.1.2 Conic constraints

Among convex constraints, there are some special classes of constraints that are generally considered nicer than "generic" convex constraints. These constraints have a specific structure which can be exploited in more efficient algorithms. These are broadly known as "*structured constraints*". A simple class of structured constraints are the so-called linear constraints, e.g.:

$$\mathbf{b}^T \mathbf{x} + a = 0 \tag{5}$$

There are a number of algorithms that can handle linear constraints more efficiently than generic inequality ( $\leq/\geq$ ) or equality ( $=$ ) constraints. In fact, linear equality constraints are the only convex equality constraints! Any other equality constraint where the function in the constraint is not linear can never be convex. Linear constraints are therefore among those special classes of so-called *structured constraints*.

Another broad class of structured constraints is conic constraints. A conic constraint is a convex constraint whose feasible set is a special convex set known as a cone. Consider almost every child's dream, an infinitely tall ice cream cone. The set of points inside this cone comprise a convex set, recall the line drawing test. This however is only one type of cone. Just like with convex sets, mathematicians generalized the concept of a cone to multiple dimensions.

A convex set  $S$  is considered a cone if for every  $\mathbf{x}$  inside the set  $S$ ,  $c\mathbf{x}$  was also inside the set  $S$  for all  $c \geq 0$ . Note how a cone must be a non-compact set since  $c$  is allowed to go to  $\infty$ . A nice way to visualize this is to think of the tip of the ice cream cone as the origin. Then extend an arrow from the origin to any point  $\mathbf{x}$  inside the cone. Stretching or contracting this arrow is equivalent to changing the value of  $c$ , where the tip of the new arrow is  $c\mathbf{x}$ . For  $S$  to be a cone, the tip of every such stretched/contracted arrow must be inside the set  $S$  for any point  $\mathbf{x} \in S$ .

Considering the new generalized definition of a cone, one can come up with various types of cones beside the famous ice cream cone. The technical term of the ice cream cone is a second-order cone by the way (leave it to mathematicians to ruin a perfectly delicious *coneept*). A structured conic constraint can be written as:

$$\mathbf{x} \in C \tag{6}$$

where  $C$  is a cone. There are various efficient algorithms that can handle conic constraints more efficiently than regular convex constraints. It is therefore often appealing to convert regular mathematical constraints that are known to be convex to their conic equivalents to make use of these efficient algorithms. For instance, the following convex constraint:

$$\mathbf{x}^T \mathbf{A} \mathbf{x} \leq c \tag{7}$$

where  $\mathbf{A}$  is a positive semi-definite matrix is equivalent to a second-order conic constraint. For more on conic constraints and transforming regular convex constraints to their conic structured equivalent, the readers are referred to the excellent online book titled: "MOSEK Modeling Cookbook" by MOSEK ApS (ApS, 2018).

### 2.1.3 Non-convexity

The vast majority of research literature on continuous mathematical optimization seems to be published on different variations and flavours of convex and conic programming, with a few exceptions. This is perhaps not surprising given the attractive theoretical properties of these classes of problems that enable mathematicians to reason about them and prove interesting, non-obvious results. There is however another extremely important class of continuous optimization problems that doesn't nearly get as much attention from mathematicians and is often treated as a "lost cause". These are the non-convex continuous optimization problems, i.e. problems where:

- The decision variables are continuous, and
- At least one constraint is not a convex constraint or the minimization (maximization) objective is not a convex (concave) function.

This is perhaps not surprising given that even some simple, non-convex, continuous optimization problems are known to be "NP-hard", where the existence of "efficient" algorithms to "solve" this class of problems is extremely unlikely by some definitions of "efficient" and "solve". This often means that efficiently solving an instance of this class requires many heuristics and is often almost as much an art as it is a science.

That said, "efficient" and "solve" are contextual terms. Often mathematicians consider the so-called "worst-case complexity" of a class of problems when talking about the computational efficiency of solving this class of problems. Worst-case complexity means that if one instance of this class is extremely difficult to "solve" for 100s of variables, we assume the whole class is extremely difficult to solve for 100s of variables. There are various classes of *problem difficulty*, aka "*complexity classes*", that are well studied in a field known as "*complexity theory*". For more on complexity theory and problem complexity analysis, the readers are referred to the excellent online book titled "Computational Complexity: A Modern Approach" by Arora and Barak (2007).

Worst-case complexity aside, if one chooses to define "efficient" and "solve" as finding a good (not necessarily the best) solution to a specific non-convex optimization problem with a reasonable number of variables and constraints in a reasonable amount of time, one can find many problem-algorithm pairings where a real-life problem can be efficiently solved, even though there would be little to no theoretical guarantees. This is generally the goal of non-convex optimization algorithms, to solve as many real-life, non-convex optimization problems as possible, as fast as possible.

In this document, a few such algorithms will be covered, those algorithms targeting generic non-convex optimization problems where no structure can be assumed. Special structured (e.g. conic) optimization algorithms on the other hand will not be covered here. That said, things aren't so black and white and there are some non-convex optimization algorithms that can efficiently handle special structured constraints when they exist while still handling generic non-convex objectives and/or constraints. Some examples of such algorithms are the Frank-Wolfe family of algorithms (Besaçon et al., 2021) and the proximal family of algorithms (Parikh and Boyd, 2013). These however are also not going to be discussed in this document.

## 2.2 Formulations

There are generally a number of ways to formulate the same optimization or decision problem using mathematical functions and decision variables. For instance, consider the following inequality constrained optimization problem:

$$\begin{aligned}
 & \underset{\mathbf{x}}{\text{minimize}} && f_0(\mathbf{x}) \\
 & \text{subject to} && \\
 & f_i(\mathbf{x}) \leq 0 && \forall i = 1..I, \\
 & l_j \leq x_j \leq u_j && \forall j = 1..V
 \end{aligned} \tag{8}$$

where  $I$  is the number of constraints,  $V$  is the number of decision variables,  $f_i$  are potentially nonlinear scalar-valued function of  $\mathbf{x}$  and  $\mathbf{l}$  and  $\mathbf{u}$  are finite vectors of lower and upper bounds on the decision vector  $\mathbf{x}$  respectively. Two formulations are considered equivalent for optimization purposes if every optimal solution of one formulation maps to an optimal solution in the other formulation. In the following sub-sections, 2 relevant re-formulation tricks will be presented. For more re-formulation techniques, the readers are referred to the excellent online book titled: "MOSEK Modeling Cookbook" by MOSEK ApS (ApS, 2018).

### 2.2.1 A linear objective is all you need

For instance, consider the alternative formulation:

$$\begin{aligned}
 & \underset{\mathbf{x}, c}{\text{minimize}} && c \\
 & \text{subject to} && \\
 & f_0(\mathbf{x}) \leq c, && (9) \\
 & f_i(\mathbf{x}) \leq 0 \quad \forall i = 1..I, \\
 & l_j \leq x_j \leq u_j \quad \forall j = 1..V
 \end{aligned}$$

where we added a new decision variable  $c$  and a new constraint  $f_0(\mathbf{x}) \leq c$  to the previous formulation and changed the objective to minimizing  $c$  which is a linear function of the decision variables.

It is clear that if  $\mathbf{x}^*$  is an optimal solution of formulation 8, then  $(\mathbf{x}, c) = (\mathbf{x}^*, f_0(\mathbf{x}^*))$  is optimal in formulation 9. This is because  $f_0(\mathbf{x}^*)$  is the lowest value  $f_0(\mathbf{x})$  can take for any feasible  $\mathbf{x}$  according to the remaining inequality constraints, and  $c = f_0(\mathbf{x}^*)$  is the lowest value the objective  $c$  can take without violating the constraint on  $c$ .

Conversely, if  $(\mathbf{x}^*, c^*)$  is an optimal solution of formulation 9,  $\mathbf{x}^*$  must be optimal in formulation 8. Since  $c$  shows up in only one constraint:

$$f_0(\mathbf{x}) \leq c \tag{10}$$

in formulation 9, for  $(\mathbf{x}^*, c^*)$  to be optimal,  $c^*$  must be equal to  $f_0(\mathbf{x}^*)$ . And since this is lowest value  $c$  can take, it must be the lowest value  $f_0(\mathbf{x})$  can take without violating any of the other constraints on  $\mathbf{x}$ , which in turn makes  $\mathbf{x}^*$  an optimal solution to formulation 8. This completes the proof.

### 2.2.2 Slack variables

Another common formulation transformation is changing all the inequality constraints to equality constraints except the variable bounds. This can be done by introducing additional slack variables. Let  $\mathbf{s}$  be a vector of so-called slack variables with length  $I$ , where  $s_i$  is the  $i^{\text{th}}$  element of the vector associated with constraint  $i$  for  $i \in 1..I$ . Formulation 8 can be re-written as:

$$\begin{aligned}
 & \underset{\mathbf{x}, \mathbf{s}}{\text{minimize}} && f_0(\mathbf{x}) \\
 & \text{subject to} && \\
 & f_i(\mathbf{x}) + s_i = 0 \quad \forall i = 1..I, && (11) \\
 & l_j \leq x_j \leq u_j \quad \forall j = 1..V, \\
 & s_i \geq 0 \quad \forall i = 1..I
 \end{aligned}$$

Since  $\mathbf{s}$  doesn't show up in the objective, in order to prove that formulations 8 and 11 are equivalent for optimization purposes, it suffices to show that every feasible solution  $\mathbf{x}$  to formulation 8 can be mapped a feasible solution  $(\mathbf{x}, \mathbf{s})$  in formulation 11 and vice versa.

It is simple to show that if  $\mathbf{x}$  is a feasible solution in formulation 8, that  $(\mathbf{x}, \mathbf{s})$  where  $\mathbf{s}$  such that  $s_i = -f_i(\mathbf{x}) \forall i \in 1..I$  is a feasible solution in formulation 11. This is because for  $\mathbf{x}$  to be feasible in formulation 8,  $f_i(\mathbf{x})$  must be non-positive for all  $i$  which makes  $s_i$  non-negative. Conversely, if  $(\mathbf{x}, \mathbf{s})$  is a feasible solution to formulation 11,  $\mathbf{x}$  is clearly feasible in formulation 8 because  $s_i$  will be non-negative which implies that  $f_i(\mathbf{x}) \leq 0$  for all  $i$ . This completes the proof.

### 2.2.3 Equality constraints and empty interiors

While the following equality constraint:

$$f(\mathbf{x}) = 0 \tag{12}$$

is in theory equivalent to the following 2 inequality constraints:

$$0 \leq f(\mathbf{x}) \leq 0 \tag{13}$$

this transformation does not change the shape or nature of the feasible domain. Not all NLP algorithms that can handle inequality constraints can also handle problems with 2 inequality constraints derived from an equality constraint like this. For instance, the convergence proof of the globally convergent method of moving asymptotes (MMA) algorithm (Svanberg, 2002) assumes that the feasible domain must have a non-empty interior. Equality constraints introduce an "empty interior". A domain  $D \subseteq \mathbb{R}^V$  is said to have a non-empty interior if  $\exists(\mathbf{x}_c \in D, r > 0)$  such that  $\mathbf{x}_c + \mathbf{u} \in D$  for all  $\mathbf{u}$  in  $\{\mathbf{u} : \mathbf{u} \in \mathbb{R}^V \wedge \|\mathbf{u}\|_2 \leq r\}$ . Since equality constraints change the feasible domain to a lower dimensional manifold embedded in  $\mathbb{R}^V$ , there exist no such  $(\mathbf{x}_c, r)$  in equality constrained problems. Therefore, the MMA algorithm will fail if the interior of the feasible domain is empty.

Linear equality constraints are generally an exception though. While technically still introducing an empty interior, most NLP algorithms can either natively handle linear equality constraints or the linear equality constrained NLP can be re-parameterized such that the interior of the feasible domain is no longer empty. The primal-dual interior point optimizer (IPOPT) discussed in this document and the augmented Lagrangian algorithm (Bertsekas, 1996) can handle linear and even nonlinear equality constraints natively. Sequential quadratic programming (SQP) methods (Boyd and Vandenberghe, 2009) can only handle linear equality constraints natively whereas nonlinear ones must get locally approximated by linear ones in an iterative process. MMA however doesn't handle either linear or nonlinear equality constraints natively. But it can be made to support linear ones with a simple nullspace re-parameterization trick.

Let the following be the linear equality constrained NLP with a nonlinear objective and inequality constraints:

$$\begin{aligned} & \underset{\mathbf{x}}{\text{minimize}} && f_0(\mathbf{x}) \\ & \text{subject to} && \\ & && f_i(\mathbf{x}) \leq 0 \quad \forall i = 1..I, \\ & && l_j \leq x_j \leq u_j \quad \forall j = 1..V, \\ & && \mathbf{Ax} = \mathbf{b} \end{aligned} \tag{14}$$

where  $A$  is a constant matrix of size  $(M \times V)$  and  $\mathbf{b}$  is a constant vector. Let  $\mathbf{x}_0$  be an arbitrary point that satisfies the  $A\mathbf{x} = \mathbf{b}$  constraint, e.g.  $\mathbf{x}_0 = A^+\mathbf{b}$  where  $A^+$  is the Moore-Penrose pseudoinverse of  $A$  (Golub and Loan, 1996). Let  $N$  be the nullspace matrix of  $A$  such that  $A \times N = \mathbf{0}$ , where  $N$  is of size  $V \times U$ . Every feasible solution  $\mathbf{x}$  to the constraints  $A\mathbf{x} = \mathbf{b}$  can now be written as:

$$\mathbf{x} = \mathbf{x}_0 + N\mathbf{y} \quad (15)$$

for some  $\mathbf{y} \in \mathbb{R}^U$ , where  $U < V$ . Substituting for  $\mathbf{x}$  in formulation 14, we get:

$$\begin{aligned} & \underset{\mathbf{y}}{\text{minimize}} && f_0(\mathbf{x}_0 + N\mathbf{y}) \\ & \text{subject to} && \\ & && f_i(\mathbf{x}_0 + N\mathbf{y}) \leq 0 \quad \forall i = 1..I, \\ & && \mathbf{l} \leq \mathbf{x}_0 + N\mathbf{y} \leq \mathbf{u} \end{aligned} \quad (16)$$

The re-parameterization essentially limits the feasible domain to the nullspace of the linear constraints since the "interior" is non-empty once we limit ourselves to the nullspace. A similar re-parameterization trick can also be used for some nonlinear manifolds, e.g. optimization on the surface of a hyper-sphere.

Beside the iterative linear approximation of nonlinear equality constraints, MMA can also be made to approximately support nonlinear equality constraints by relaxing it as follows:

$$-\epsilon \leq f(\mathbf{x}) \leq \epsilon \quad (17)$$

where  $\epsilon > 0$  which creates a non-empty interior.

### 2.3 Regularity conditions

There are a number of optimization algorithms that can solve a general NLP without assuming convexity. Most NLP algorithms tend to converge to a locally optimal solution without global optimality guarantees and this suffices in many applications. For simplicity, assume all the inequality constraints have been converted to equality constraints using slack variables. Let the NLP with equality constraints be:

$$\underset{\mathbf{x}}{\text{minimize}} \quad f(\mathbf{x}) \quad (18a)$$

subject to

$$\mathbf{c}(\mathbf{x}) = \mathbf{0}, \quad (18b)$$

$$\mathbf{l} \leq \mathbf{x} \leq \mathbf{u} \quad (18c)$$

where  $\mathbf{c}(\mathbf{x})$  is of length  $E$ . Additionally, let  $f$  and  $\mathbf{c}$  be continuous and twice differentiable functions.

A point  $\mathbf{x}$  is called a regular point if it satisfies one of the so-called constraint qualification conditions. Three common constraint qualifications for NLPs are:

1. Linear constraint qualification (LCQ):  $\mathbf{c}(\mathbf{x})$  is an affine function.

2. Linear independence constraint qualification (LICQ): the rows of the Jacobian  $\nabla \mathbf{c}(\mathbf{x})$  and the gradients of the active (i.e. satisfied at equality) bound constraints are linearly independent at  $\mathbf{x}$ .
3. Mangasarian-Fromovitz constraint qualification (MFCQ): the rows of  $\nabla \mathbf{c}(\mathbf{x})$  are linearly independent at  $\mathbf{x}$  and there exists a direction vector  $\mathbf{d} \in \mathbb{R}^V$  where  $d_i > 0$  if  $x_i = l_i$ ,  $d_i < 0$  if  $x_i = u_i$ , and  $\nabla \mathbf{c}(\mathbf{x})^T \mathbf{d} = \mathbf{0}$ .

Alternatively, if the problem is convex and the nonlinear inequality constraints are not converted to equality constraints, every feasible point is regular if  $\exists \mathbf{x}$  such that all the inequality constraints are satisfied but not active (i.e. satisfied at a strict inequality) and all the linear equality constraints are satisfied. This condition is known as Slater's condition (Boyd and Vandenberghe, 2009).

## 2.4 Sufficient optimality conditions for regular points

Let:

$$L(\mathbf{x}, \boldsymbol{\lambda}, \mathbf{z}_+, \mathbf{z}_-) = f(\mathbf{x}) + \mathbf{c}(\mathbf{x})^T \boldsymbol{\lambda} + (\mathbf{x} - \mathbf{u})^T \mathbf{z}_+ - (\mathbf{x} - \mathbf{l})^T \mathbf{z}_- \quad (19)$$

where  $\boldsymbol{\lambda} \in \mathbb{R}^E$  is the vector of Lagrangian multipliers of the nonlinear constraints,  $\mathbf{z}_- \in \mathbb{R}_+^V$  is the vector of Lagrangian multipliers of the  $\geq$  bound constraints, and  $\mathbf{z}_+ \in \mathbb{R}_+^V$  is the vector of Lagrangian multipliers of the  $\leq$  bound constraints.

If  $\mathbf{x}$  is regular and is a local optimum of the NLP, then  $\exists (\boldsymbol{\lambda}, \mathbf{z}_+, \mathbf{z}_-)$  such that:

$$\nabla_{\mathbf{x}} L(\mathbf{x}, \boldsymbol{\lambda}, \mathbf{z}_+, \mathbf{z}_-) = \mathbf{0} \quad (20)$$

$$\mathbf{c}(\mathbf{x}) = \mathbf{0} \quad (21)$$

$$\mathbf{l} \leq \mathbf{x} \leq \mathbf{u} \quad (22)$$

$$\mathbf{z}_+ \geq \mathbf{0} \quad (23)$$

$$\mathbf{z}_- \geq \mathbf{0} \quad (24)$$

$$(\mathbf{x} - \mathbf{u})^T \mathbf{z}_+ = 0 \quad (25)$$

$$(\mathbf{x} - \mathbf{l})^T \mathbf{z}_- = 0 \quad (26)$$

$$\nabla \mathbf{c}(\mathbf{x})^T \nabla_{\mathbf{x}\mathbf{x}}^2 L(\mathbf{x}, \boldsymbol{\lambda}, \mathbf{z}_+, \mathbf{z}_-) \nabla \mathbf{c}(\mathbf{x}) \succcurlyeq \mathbf{0} \quad (27)$$

These conditions are known as the KKT sufficient conditions for optimality and  $\mathbf{x}$  would be called a KKT point (Boyd and Vandenberghe, 2009). Condition 20 is known as the stationarity condition which generalizes the 0 gradient condition for unconstrained NLPs. Conditions 21 and 22 are known as primal feasibility conditions. Conditions 23 and 24 are known as dual feasibility conditions. Finally, condition 27 is known as the second order KKT optimality condition which generalizes the second order optimality condition for unconstrained NLPs, where  $\mathbf{A} \succcurlyeq \mathbf{0}$  when  $\mathbf{A}$  is a matrix means that  $\mathbf{A}$  must be positive semi-definite.

Note that not every local optimal solution to the NLP must be a regular point or by consequence a KKT point for that matter. Therefore, these conditions are not necessary conditions for optimality for general NLPs. However for convex problems,

if Slater's condition is satisfied, these conditions are both necessary and sufficient and the local/global optimum is guaranteed to be a regular and a KKT point (Boyd and Vandenberghe, 2009). All the optimization algorithms presented next seek to find a locally optimal KKT point.

### 3 Method of moving asymptotes

#### 3.1 Separable convex approximation

Consider the following inequality constrained optimization problem:

$$\begin{aligned}
 & \underset{\mathbf{x}}{\text{minimize}} && f_0(\mathbf{x}) \\
 & \text{subject to} && \\
 & f_i(\mathbf{x}) \leq 0 && \forall i = 1..I, \\
 & l_j \leq x_j \leq u_j && \forall j = 1..V
 \end{aligned} \tag{28}$$

where  $I$  is the number of inequality constraints,  $V$  is the number of decision variables,  $f_i$  is a potentially nonlinear scalar-valued function of  $\mathbf{x}$  and  $\mathbf{l}$  and  $\mathbf{u}$  are finite vectors of lower and upper bounds on the decision vector  $\mathbf{x}$  respectively.

The first MMA algorithm proposed by Svanberg (1987) relied on a separable convex approximation of the objective and constraint functions. A function  $f(\mathbf{x})$  is linearly approximated with respect to either  $t_{l,j}$  or  $t_{u,j}$ , where:

$$t_{l,j} = t_l(x_j; L_j) = \frac{1}{x_j - L_j} \tag{29}$$

$$t_{u,j} = t_u(x_j; U_j) = \frac{1}{U_j - x_j} \tag{30}$$

for all  $j \in 1..V$  for some constants  $L_j$  and  $U_j$ , where  $l_j \leq L_j < U_j \leq u_j$ , and  $L_j$  and  $U_j$  are known as the asymptotes of the approximation. The choice of which function,  $t_{l,j}$  or  $t_{u,j}$ , to approximate  $f$  with respect to, for each variable  $x_j$ , depends on the sign of the partial derivative  $\frac{\partial f}{\partial x_j}$  such that the approximation is convex.

More specifically, let the MMA approximation of  $f(\mathbf{x})$  around  $\bar{\mathbf{x}}$  be:

$$\bar{f}(\mathbf{x}; \bar{\mathbf{x}}) = f(\bar{\mathbf{x}}) + \sum_j \bar{f}_j(x_j; \bar{\mathbf{x}}) \tag{31}$$

Let  $\bar{f}_j(x_j; \bar{\mathbf{x}})$  be:

$$\bar{f}_j(x_j; \bar{\mathbf{x}}) = \begin{cases} \left( t_l(x_j; L_j) - t_l(\bar{x}_j; L_j) \right) \frac{\partial f}{\partial t_{l,j}}(\bar{\mathbf{x}}), & \frac{\partial f}{\partial x_j}(\bar{\mathbf{x}}) < 0 \\ \left( t_u(x_j; U_j) - t_u(\bar{x}_j; U_j) \right) \frac{\partial f}{\partial t_{u,j}}(\bar{\mathbf{x}}), & \frac{\partial f}{\partial x_j}(\bar{\mathbf{x}}) \geq 0 \end{cases} \tag{32}$$

where  $\bar{x}_j$  is the  $j^{\text{th}}$  element of  $\bar{\mathbf{x}}$  and  $\frac{\partial f}{\partial t_{l,j}}$  and  $\frac{\partial f}{\partial t_{u,j}}$  are:

$$\frac{\partial f}{\partial t_{l,j}} = \frac{\partial f}{\partial x_j} / \frac{dt_{l,j}}{dx_j} = -(x_j - L_j)^2 \frac{\partial f}{\partial x_j} \tag{33}$$

$$\frac{\partial f}{\partial t_{u,j}} = \frac{\partial f}{\partial x_j} / \frac{dt_{u,j}}{dx_j} = (U_j - x_j)^2 \frac{\partial f}{\partial x_j} \tag{34}$$

$\bar{f}_j(x_j; \bar{\mathbf{x}})$  can therefore be written as:

$$\bar{f}_j(x_j; \bar{\mathbf{x}}) = \begin{cases} -\left(t_l(x_j; L_j) - t_l(\bar{x}_j; L_j)\right)(\bar{x}_j - L_j)^2 \frac{\partial f}{\partial x_j}(\bar{\mathbf{x}}), & \frac{\partial f}{\partial x_j}(\bar{\mathbf{x}}) < 0 \\ \left(t_u(x_j; U_j) - t_u(\bar{x}_j; U_j)\right)(U_j - \bar{x}_j)^2 \frac{\partial f}{\partial x_j}(\bar{\mathbf{x}}), & \frac{\partial f}{\partial x_j}(\bar{\mathbf{x}}) \geq 0 \end{cases} \quad (35)$$

The signs of the gradient and Hessian of  $\bar{f}(\mathbf{x}; \bar{\mathbf{x}})$  therefore depend mostly on  $t_l(x_j; L_j)$  and  $t_u(x_j; U_j)$ . Since the approximation is separable in  $\mathbf{x}$ , the Hessian of  $\bar{f}(\mathbf{x}; \bar{\mathbf{x}})$  wrt  $\mathbf{x}$  is a diagonal matrix where the  $j^{\text{th}}$  diagonal element is:

$$\frac{\partial^2 \bar{f}}{\partial x_j^2} = \begin{cases} -\frac{d^2 t_{l,j}}{dx_j^2} (\bar{x}_j - L_j)^2 \frac{\partial f}{\partial x_j}(\bar{\mathbf{x}}), & \frac{\partial f}{\partial x_j}(\bar{\mathbf{x}}) < 0 \\ \frac{d^2 t_{u,j}}{dx_j^2} (U_j - \bar{x}_j)^2 \frac{\partial f}{\partial x_j}(\bar{\mathbf{x}}), & \frac{\partial f}{\partial x_j}(\bar{\mathbf{x}}) \geq 0 \end{cases} \quad (36)$$

Since  $\frac{d^2 t_{l,j}}{dx_j^2}$  and  $\frac{d^2 t_{u,j}}{dx_j^2}$  are both positive for all  $x_j$ , it is clear that the MMA approximation is always going to be convex.

In each iteration of the MMA algorithm, convex approximations of the objective and all the constraint functions are formed around the current solution and the approximate problem is solved to optimality while restricting the decision variables to be between the asymptotes  $\mathbf{L}$  and  $\mathbf{U}$  instead of the original bounds  $\mathbf{l}$  and  $\mathbf{u}$ . This is similar to the trust region approach. Let the restricted convex approximation of the original nonlinear program be:

$$\begin{aligned} & \underset{\mathbf{x}}{\text{minimize}} && \bar{f}_0(\mathbf{x}; \bar{\mathbf{x}}) \\ & \text{subject to} && \\ & \bar{f}_i(\mathbf{x}; \bar{\mathbf{x}}) \leq 0 && \forall i = 1..I, \\ & \alpha_j \leq x_j \leq \beta_j && \forall j = 1..V \end{aligned}$$

where each decision variable  $x_j$  is restricted to be between  $\alpha_j = \max(l_j, L_j)$  and  $\beta_j = \min(u_j, U_j)$  and  $\bar{f}_i(\mathbf{x}; \bar{\mathbf{x}})$  is the MMA approximation of  $f_i(\mathbf{x})$  around  $\bar{\mathbf{x}}$ . In order to form the approximation above, the gradient of the objective and the full Jacobian of the constraint functions need to be computed first. This typically limits the MMA algorithm's scalability in handling many constraints where the full Jacobian can be expensive to form, e.g. in stress-constrained topology optimization. However if only a few constraints exist, the MMA algorithm is usually quite robust when starting from a feasible solution. Most importantly is that once the approximate problem is formed, no more calls to the objective or constraint functions are required by the primal-dual algorithm to solve the approximate problem to optimality.

### 3.2 Efficient primal-dual algorithm

Once the approximation is formed, the separable nature of the MMA convex approximation then allows the approximate nonlinear program to be solved to optimality using an efficient primal-dual Lagrangian optimization algorithm (Svanberg, 1987). The dual

of the convex approximation problem above is the following lower bound constrained nested optimization problem:

$$\text{maximize}_{\lambda \geq \mathbf{0}} \quad \min_{l \leq x \leq u} \mathcal{L}(x, \lambda) = r_0 + \sum_{j=1}^V \mathcal{L}_j(x_j, \lambda)$$

where

$$\mathcal{L}_j(x_j; \lambda) = \frac{p_{0,j} + \sum_i \lambda_i p_{i,j}}{U_j - x_j} + \frac{q_{0,j} + \sum_i \lambda_i q_{i,j}}{x_j - L_j} \quad (37)$$

$$r_0 = f_0(\bar{x}) - \sum_{j=1}^V \frac{p_{0,j}}{U_j - \bar{x}_j} + \frac{q_{0,j}}{\bar{x}_j - L_j} \quad (38)$$

$$p_{i,j} = \begin{cases} (U_j - \bar{x}_j)^2 \frac{\partial f_i}{\partial x_j}(\bar{x}) & \frac{\partial f_i}{\partial x_j}(\bar{x}) > 0 \\ 0 & \frac{\partial f_i}{\partial x_j}(\bar{x}) \leq 0 \end{cases} \quad (39)$$

$$q_{i,j} = \begin{cases} 0 & \frac{\partial f_i}{\partial x_j}(\bar{x}) \geq 0 \\ -(\bar{x}_j - L_j)^2 \frac{\partial f_i}{\partial x_j}(\bar{x}) & \frac{\partial f_i}{\partial x_j}(\bar{x}) < 0 \end{cases} \quad (40)$$

The primal optimal solution  $\mathbf{x}^*$  of the convex approximation has an analytic form as a function of the dual solution  $\lambda$ :

$$x_j^*(\lambda) = \begin{cases} \alpha_j & \frac{\partial \mathcal{L}_j}{\partial x_j}(\alpha_j; \lambda) \geq 0 \\ \beta_j & \frac{\partial \mathcal{L}_j}{\partial x_j}(\beta_j; \lambda) \leq 0 \\ \frac{(p_{0,j} + \sum_i \lambda_i p_{i,j})^{1/2} L_j + (q_{0,j} + \sum_i \lambda_i q_{i,j})^{1/2} U_j}{(p_{0,j} + \sum_i \lambda_i p_{i,j})^{1/2} + (q_{0,j} + \sum_i \lambda_i q_{i,j})^{1/2}} & \text{otherwise.} \end{cases} \quad (41)$$

### 3.3 Algorithm strengths

Since the optimization of the approximate nonlinear program requires no calls to the objective or constraint functions, its performance is independent of the computational complexity of computing the objective value, its gradient, the constraint values and their Jacobian.

One of the biggest numerical challenges in topology optimization is that the scale of the objective and constraints can often be vastly different. For example in the volume constrained compliance minimization problem, the objective value scales up with the volume of the design and the Young's modulus while the volume fraction constraint does not. Moreover, the  $\ell_\infty$ -norm of the gradient of the volume fraction constraint scales down as the number of elements increases while that of the gradient of the compliance function does not. This usually results in problems where the objective and constraints span multiple orders of magnitude. This difference in scale usually means that more iterations are needed to converge to sufficiently large or sufficiently small elements of the dual solution  $\lambda$  that satisfies the first order KKT optimality conditions. More specifically in the MMA algorithm, this mostly translates to requiring more iterations to solve the approximate problem to optimality and usually not many more subproblems to be solved. However as mentioned earlier, the additional iterations needed to solve

the approximate subproblem do not require any additional calls to the objective or constraint functions. This is an especially attractive feature of the MMA algorithm in topology optimization since scaling issues are common and the computational time is usually dominated by the time it takes to compute the objective, the constraints and their gradients.

Another attractive feature of the MMA algorithm is the low sensitivity of the algorithm's performance to most of the parameters. In `Nonconvex.jl`, the lower bound constrained dual of the convex approximation is solved using a log-barrier method with the nonlinear conjugate gradient algorithm (Nocedal and Wright, 2006). The parameters of the subproblem optimizer generally have little effect on the performance of the MMA algorithm for the reason mentioned above. The main parameters which tend to affect the performance of the algorithm by changing the number of subproblems that need to be solved are:  $s_{init}$  used to specify the initial asymptotes,  $s_{incr}$  used to widen the asymptotes of each variable and  $s_{decr}$  used to tighten the asymptotes for each variable. These parameters are described in details by Svanberg (1987). Finally, a tolerance  $tol$  must be picked to terminate the algorithm.

## 4 Primal-dual interior point method

In this section, the primal-dual interior point optimizer (IPOPT) as described in Wächter and Biegler (2006) and implemented in the IPOPT software will be described. This algorithm and its implementation are well tested in practice. However, the main disadvantage of the IPOPT algorithm is its complexity and the large number of hyper-parameters as will be seen next.

### 4.1 Original problem

Consider the following equality constrained optimization problem:

$$\begin{aligned}
 & \underset{\mathbf{x}}{\text{minimize}} && f(\mathbf{x}) \\
 & \text{subject to} && \\
 & \mathbf{c}(\mathbf{x}) = \mathbf{0}, && \\
 & \mathbf{l} \leq \mathbf{x} \leq \mathbf{u} &&
 \end{aligned} \tag{42}$$

where the length of  $\mathbf{x}$  is  $V$ ,  $E$  is the number of elements in the output of the vector-valued function  $\mathbf{c}(\mathbf{x})$ ,  $f$  is a potentially nonlinear scalar-valued function of  $\mathbf{x}$  and  $\mathbf{l}$  and  $\mathbf{u}$  are the lower and upper bounds on the variables. If there are inequality constraints, they can be converted to equality constraints by adding slack variables. Additionally, let  $\mathcal{I}_l$  be the set of variable indices with a finite lower bound and  $\mathcal{I}_u$  be the set of variable indices with a finite upper bound. Additionally, assume that the lower and upper bounds are not equal for any variable. If a variable's lower and upper bounds are equal, it can be fixed and removed from the optimization.

### 4.2 Barrier problem

In the IPOPT algorithm by Wächter and Biegler (2006), a so-called log barrier method (Boyd and Vandenberghe, 2009) is used to guarantee that  $\mathbf{l} \leq \mathbf{x} \leq \mathbf{u}$  remains satisfied at every intermediate solution if the initial solution is within the bounds. This is achieved using a so-called barrier function such as:

$$-\mu \left( \sum_{i \in \mathcal{I}_l} \log(x_i - l_i) + \sum_{i \in \mathcal{I}_u} \log(u_i - x_i) \right) \tag{43}$$

for some  $\mu > 0$  which would go to  $\infty$  if any of the decision variables approaches one of its finite bounds. This creates a barrier stopping the optimizer from ever reaching the finite bound. It can be shown that if a decreasing geometric series of values  $\mu$  are used that the solutions of the sub-problems will follow a so-called critical path converging to a KKT point in convex problems satisfying Slater's condition (Boyd and Vandenberghe, 2009). However, if the optimal solution is exactly on a boundary, the problem can become numerically unstable near the optimal solution. Therefore, the variable bounds are typically relaxed slightly. In particular, the lower bound  $l_i$  where  $i \in \mathcal{I}_l$  is relaxed by  $tol \times \max(1, |l_i|)$ , and the upper bound  $u_i$  where  $i \in \mathcal{I}_u$  is relaxed

by  $tol \times \max(1, |u_i|)$ .  $\mathbf{l}$  and  $\mathbf{u}$  will refer to the relaxed lower and upper bounds from now on.

One final catch in the barrier problem formulation is that if the set of optimal points to the original problem does not consist of isolated points, but contains an unbounded connected subspace, e.g. a line to  $\infty$  or  $-\infty$  for one or more variables, the log-barrier term may become  $-\infty$ . This is only possible if a variable is bounded from one side only where it's allowed to go to  $\pm\infty$  from the other side. If this happens, the minimum objective value of the barrier problem becomes  $-\infty$  even when the original problem's objective value is bounded. For this reason, the following additional term is added to the barrier objective:

$$\kappa_d \mu \sum_{i \in \bar{\mathcal{I}}_l \setminus \bar{\mathcal{I}}_u} \log(x_i - l_i) + \kappa_d \mu \sum_{i \in \bar{\mathcal{I}}_u \setminus \bar{\mathcal{I}}_l} \log(u_i - x_i) \quad (44)$$

for some small constant  $\kappa_d \in (0, 1)$ . This way, divergence of variables having only one bound is penalized. The effect of this additional term is reduced as  $\mu$  decreases and can be shown to not affect the local convergence proof (Wächter and Biegler, 2006). The barrier sub-problem is therefore defined as:

$$\begin{aligned} \underset{\mathbf{x}}{\text{minimize}} \quad & \phi_\mu(\mathbf{x}) = f(\mathbf{x}) - \mu \sum_{i \in \bar{\mathcal{I}}_l} \log(x_i - l_i) - \mu \sum_{i \in \bar{\mathcal{I}}_u} \log(u_i - x_i) \\ & + \kappa_d \mu \sum_{i \in \bar{\mathcal{I}}_l \setminus \bar{\mathcal{I}}_u} \log(x_i - l_i) + \kappa_d \mu \sum_{i \in \bar{\mathcal{I}}_u \setminus \bar{\mathcal{I}}_l} \log(u_i - x_i) \quad (45) \\ \text{subject to} \quad & \\ \mathbf{c}(\mathbf{x}) = \mathbf{0} \quad & \end{aligned}$$

The KKT stationarity condition is therefore:

$$\nabla f(\mathbf{x}) + \nabla \mathbf{c}(\mathbf{x})^T \boldsymbol{\lambda} - \mathbf{z}_l - \mathbf{z}_u = \mathbf{0} \quad (46)$$

where  $\boldsymbol{\lambda}$  is the vector Lagrangian multipliers associated with the equality constraint  $\mathbf{c}(\mathbf{x}) = \mathbf{0}$ ,  $\mathbf{z}_l$  is a vector whose  $i^{th}$  element is:

$$\mathbf{z}_l = \begin{cases} \frac{\mu}{x_i - l_i} & i \in \bar{\mathcal{I}}_l \cap \bar{\mathcal{I}}_u \\ \frac{(1 - \kappa_d)\mu}{x_i - l_i} & i \in \bar{\mathcal{I}}_l \setminus \bar{\mathcal{I}}_u \\ 0 & \text{otherwise,} \end{cases} \quad (47)$$

$\mathbf{z}_u$  is a vector whose  $i^{th}$  element is:

$$\mathbf{z}_u = \begin{cases} \frac{\mu}{u_i - x_i} & i \in \bar{\mathcal{I}}_l \cap \bar{\mathcal{I}}_u \\ \frac{(1 - \kappa_d)\mu}{u_i - x_i} & i \in \bar{\mathcal{I}}_u \setminus \bar{\mathcal{I}}_l \\ 0 & \text{otherwise,} \end{cases} \quad (48)$$

$\mathbf{1}_l$  is a vector whose  $i^{th}$  element is:

$$\begin{cases} 1 & i \in \bar{\mathcal{I}}_l \\ 0 & \text{otherwise,} \end{cases} \quad (49)$$

and  $\mathbf{1}_u$  is a vector whose  $i^{th}$  element is:

$$\begin{cases} 1 & i \in \mathcal{I}_u \\ 0 & \text{otherwise.} \end{cases} \quad (50)$$

Additionally, let  $\mathbf{Z}_l$  be the diagonal matrix whose diagonal is  $\mathbf{z}_l$ ,  $\mathbf{Z}_u$  be the diagonal matrix whose diagonal is  $\mathbf{z}_u$ ,  $\mathbf{I}_l$  be the diagonal matrix whose diagonal is  $\mathbf{1}_l$ ,  $\mathbf{I}_u$  be the diagonal matrix whose diagonal is  $\mathbf{1}_u$ ,  $\mathbf{X}_l$  be the diagonal matrix whose  $i^{th}$  diagonal element is:

$$\begin{cases} x_i - l_i & i \in \mathcal{I}_l \cap \mathcal{I}_u \\ \frac{x_i - l_i}{1 - \kappa_d} & i \in \mathcal{I}_l \setminus \mathcal{I}_u \\ 0 & \text{otherwise,} \end{cases} \quad (51)$$

and  $\mathbf{X}_u$  be the diagonal matrix whose  $i^{th}$  diagonal element is:

$$\begin{cases} u_i - x_i & i \in \mathcal{I}_l \cap \mathcal{I}_u \\ \frac{u_i - x_i}{1 - \kappa_d} & i \in \mathcal{I}_u \setminus \mathcal{I}_l \\ 0 & \text{otherwise.} \end{cases} \quad (52)$$

The first order KKT sufficient conditions for optimality assuming the constraint qualifications are satisfied can therefore be written as:

$$\nabla f(\mathbf{x}) + \nabla \mathbf{c}(\mathbf{x})^T \lambda - \mathbf{z}_l - \mathbf{z}_u = \mathbf{0} \quad (53)$$

$$\mathbf{c}(\mathbf{x}) = \mathbf{0} \quad (54)$$

$$\mathbf{X}_l \mathbf{Z}_l \mathbf{1}_l - \mu \mathbf{1}_l = \mathbf{0} \quad (55)$$

$$\mathbf{X}_u \mathbf{Z}_u \mathbf{1}_u - \mu \mathbf{1}_u = \mathbf{0} \quad (56)$$

$$l \leq \mathbf{x} \leq u \quad (57)$$

$$\mathbf{z}_l \geq \mathbf{0} \quad (58)$$

$$\mathbf{z}_u \geq \mathbf{0} \quad (59)$$

Conditions 55 and 56 ensure that the relationship between  $\mathbf{X}_l$ ,  $\mathbf{Z}_l$ ,  $\mathbf{X}_u$  and  $\mathbf{Z}_u$  is maintained according to the definitions of  $\mathbf{z}_l$  and  $\mathbf{z}_u$ . Assume  $l_i$  and  $u_i$  are both finite for some index  $i$ . Given that  $\mu$  is positive, condition 55 guarantees that  $x_i - l_i$  and  $z_{li}$  will either be both positive or both negative. However even if we start from a value for  $x_i > l_i$  and  $z_{li} > 0$ , during the intermediate line search steps, there is a non-zero chance that both  $x_i - l_i$  and  $z_i$  may become negative simultaneously for some  $i$ , hence the need for the explicit non-negativity constraints on  $\mathbf{z}_l$  and  $\mathbf{z}_u$ , and the bounds constraints on  $\mathbf{x}$ .

The primal-dual interior point optimizer presented in Wächter and Biegler (2006) computes an approximate solution to the barrier problem for a fixed value of  $\mu$  then decreases  $\mu$  and continues the solution of the next barrier problem from the approximate

solution of the previous one. The optimality error is defined as:

$$E_\mu(\mathbf{x}, \boldsymbol{\lambda}, \mathbf{z}_l, \mathbf{z}_u) := \max \left\{ \frac{\|\nabla f(\mathbf{x}) + \nabla \mathbf{c}(\mathbf{x})^T \boldsymbol{\lambda} - \mathbf{z}_l - \mathbf{z}_u\|}{s_d}, \|\mathbf{c}(\mathbf{x})\|_1, \frac{\|\mathbf{X}_l \mathbf{Z}_l \mathbf{1}_l - \mu \mathbf{1}_l\|_\infty}{s_c}, \frac{\|\mathbf{X}_u \mathbf{Z}_u \mathbf{1}_u - \mu \mathbf{1}_u\|_\infty}{s_c} \right\} \quad (60)$$

with scaling parameters  $s_d, s_c \geq 1$  defined below.  $E_0(\mathbf{x}, \boldsymbol{\lambda}, \mathbf{z}_l, \mathbf{z}_u)$  is the optimality error for the original problem. Let  $(\tilde{\mathbf{x}}_*, \tilde{\boldsymbol{\lambda}}_*, \tilde{\mathbf{z}}_{l,*}, \tilde{\mathbf{z}}_{u,*})$  be the approximate solution of the barrier problem. The algorithm terminates if  $E_0(\tilde{\mathbf{x}}_*, \tilde{\boldsymbol{\lambda}}_*, \tilde{\mathbf{z}}_{l,*}, \tilde{\mathbf{z}}_{u,*}) \leq \text{tol}$  for some tolerance  $\text{tol}$ .

It is possible that the magnitudes of  $\boldsymbol{\lambda}$ ,  $\mathbf{z}_l$  and  $\mathbf{z}_u$  might become very large, e.g. if the gradients of the active constraints are nearly linearly dependent in a solution to the original problem. This is why the scaling factors  $s_d$  and  $s_c$  are used to make it easier to satisfy the terminating condition around these solutions.  $s_d$  and  $s_c$  are therefore chosen as:

$$s_d = \max \left\{ s_{max}, \frac{\|\boldsymbol{\lambda}\|_1 + \|\mathbf{z}_l\|_1 + \|\mathbf{z}_u\|_1}{E + V} \right\} / s_{max} \quad (61)$$

$$s_c = \max \left\{ s_{max}, \frac{\|\mathbf{z}_l\|_1 + \|\mathbf{z}_u\|_1}{V} \right\} / s_{max} \quad (62)$$

Using these scaling factors, a component of the optimality error will be scaled, whenever the average value of the multipliers becomes larger than a fixed number  $s_{max} \geq 1$

Denoting with  $j$  the outer iteration counter of barrier sub-problems. Each barrier problem  $j$  is terminated when:

$$E_{\mu_{j-1}}(\tilde{\mathbf{x}}_{*,j}, \tilde{\boldsymbol{\lambda}}_{*,j}, \tilde{\mathbf{z}}_{l,*}, \tilde{\mathbf{z}}_{u,*}) \leq \kappa_\epsilon \mu_j \quad (63)$$

for a constant  $\kappa_\epsilon > 0$ . The new barrier parameter  $\mu_j$  is then obtained using:

$$\mu_j = \max \left\{ \frac{\text{tol}}{10}, \min \left\{ \kappa_\mu \mu_{j-1}, \mu_{j-1}^{\theta_\mu} \right\} \right\} \quad (64)$$

with constants  $\kappa_\mu \in (0, 1)$  and  $\theta_\mu \in (1, 2)$ . In this way,  $\mu$  is eventually decreased at a superlinear rate. On the other hand, this update rule stops  $\mu$  from becoming smaller than necessary given the desired tolerance  $\text{tol}$ .

### 4.3 Damped Newton's method

In order to solve the barrier problem for a given value  $\mu = \mu_j$ , damped Newton's method is applied to the primal-dual optimality conditions. Here we use  $k$  to denote the iteration counter for the inner iterations when solving the barrier problem. Given an iterate  $(\mathbf{x}_k, \boldsymbol{\lambda}_k, \mathbf{z}_{l,k}, \mathbf{z}_{u,k})$  with  $\mathbf{l} < \mathbf{x}_k < \mathbf{u}$  and  $\mathbf{z}_{l,k}, \mathbf{z}_{u,k} > \mathbf{0}$ , some search directions

$(d_k^x, d_k^\lambda, d_k^{z_l}, d_k^{z_u})$  are obtained using the regularized linearization of the optimality conditions (excluding the inequality constraints):

$$\begin{bmatrix} W_k + \delta_w I & A_k & -I_l & -I_u \\ A_k^T & -\delta_c I & \mathbf{0} & \mathbf{0} \\ Z_{l,k} & \mathbf{0} & X_{l,k} & \mathbf{0} \\ Z_{u,k} & \mathbf{0} & \mathbf{0} & X_{u,k} \end{bmatrix} \begin{pmatrix} d_k^x \\ d_k^\lambda \\ d_k^{z_l} \\ d_k^{z_u} \end{pmatrix} = - \begin{pmatrix} \nabla f(x_k) + A_k^T \lambda_k - z_{l,k} - z_{u,k} \\ c(x_k) \\ X_{l,k} Z_{l,k} \mathbf{1}_l - \mu_j \mathbf{1}_l \\ X_{u,k} Z_{u,k} \mathbf{1}_u - \mu_j \mathbf{1}_u \end{pmatrix} \quad (65)$$

where  $\delta_w$  and  $\delta_c$  are non-negative constants and  $A_k := \nabla c(x_k)^T$  and  $W_k := \nabla_{xx}^2 \mathcal{L}(x_k, \lambda_k)$  which is the Hessian of the Lagrangian function of the original problem, where:

$$\mathcal{L}(x, \lambda) := f(x) + c(x)^T \lambda \quad (66)$$

Note that the Lagrangian terms from the bounds constraint  $\mathbf{0} \leq x \leq u$  in the original problem can be ignored since their contribution to the Hessian is 0. When the Hessian of the Lagrangian is not available, an l-BFGS approximation (Nocedal and Wright, 2006) of the Hessian can be used instead. This changes the IPOPT algorithm from a second order algorithm to a first order one. And the Newton update becomes a quasi-Newton update.

When  $\delta_w$  and  $\delta_c$  are 0s, we get the linearization of the primal-dual first order KKT conditions. However, since the Hessian of the Lagrangian function isn't necessarily a positive definite matrix when the NLP is non-convex and in order to guarantee that the search directions obtained are descent directions (Wächter and Biegler, 2006), a positive value for  $\delta_w$  is used. Additionally, if the gradients of the active constraints are (nearly) linearly dependent, the matrix on the LHS will be (nearly) singular even if the Hessian of the Lagrangian is positive definite. To stop the matrix from becoming singular in this case, a positive value for  $\delta_c$  can be used. If the matrix is so ill-conditioned even with large values for  $\delta_w$  and  $\delta_c$ , the algorithm gives up on finding a search direction and attempts a feasibility restoration step, hoping that the matrix has better properties close to feasible points. For more on the feasibility restoration phase or on the heuristic used to choose  $\delta_w$  and  $\delta_c$ , the readers are referred to section 3.1 in Wächter and Biegler (2006).

Instead of solving the non-symmetric system of equations above, one can instead change the system as such:

$$\begin{bmatrix} W_k + \delta_w I & A_k & -I_l & -I_u \\ A_k^T & -\delta_c I & \mathbf{0} & \mathbf{0} \\ X_{l,k}^+ Z_{l,k} & \mathbf{0} & I_l & \mathbf{0} \\ X_{u,k}^+ Z_{u,k} & \mathbf{0} & \mathbf{0} & I_u \end{bmatrix} \begin{pmatrix} d_k^x \\ d_k^\lambda \\ d_k^{z_l} \\ d_k^{z_u} \end{pmatrix} = - \begin{pmatrix} \nabla f(x_k) + A_k^T \lambda_k - z_{l,k} - z_{u,k} \\ c(x_k) \\ z_{l,k} - \mu_j X_{l,k}^+ \mathbf{1}_l \\ z_{u,k} - \mu_j X_{u,k}^+ \mathbf{1}_u \end{pmatrix} \quad (67)$$

where  $X^+$  is the Moore-Penrose pseudoinverse of  $X$  (Golub and Loan, 1996). Adding the third and fourth equations to the first one, the third and fourth blocks of coefficients of the first equation will be eliminated.

$$\begin{bmatrix} W_k + \delta_w I + \Sigma_k & A_k & \mathbf{0} & \mathbf{0} \\ A_k^T & -\delta_c I & \mathbf{0} & \mathbf{0} \\ X_{l,k}^+ Z_{l,k} & \mathbf{0} & I_l & \mathbf{0} \\ X_{u,k}^+ Z_{u,k} & \mathbf{0} & \mathbf{0} & I_u \end{bmatrix} \begin{pmatrix} d_k^x \\ d_k^\lambda \\ d_k^{z_l} \\ d_k^{z_u} \end{pmatrix} = - \begin{pmatrix} \nabla f(x_k) + A_k^T \lambda_k - \mu_j X_{l,k}^+ \mathbf{1}_l - \mu_j X_{u,k}^+ \mathbf{1}_u \\ c(x_k) \\ z_{l,k} - \mu_j X_{l,k}^+ \mathbf{1}_l \\ z_{u,k} - \mu_j X_{u,k}^+ \mathbf{1}_u \end{pmatrix} \quad (68)$$

where  $\Sigma_k = X_{l,k}^+ Z_{l,k} + X_{u,k}^+ Z_{u,k}$ . Therefore, one can now solve for  $d_k^x$  and  $d_k^\lambda$  by solving the following symmetric linear system:

$$\begin{bmatrix} W_k + \delta_w I + \Sigma_k & A_k \\ A_k^T & -\delta_c I \end{bmatrix} \begin{pmatrix} d_k^x \\ d_k^\lambda \end{pmatrix} = - \begin{pmatrix} \nabla f(x_k) + A_k^T \lambda_k - \mu_j X_{l,k}^+ \mathbf{1}_l - \mu_j X_{u,k}^+ \mathbf{1}_u \\ c(x_k) \end{pmatrix} \quad (69)$$

then use the value of  $d_k^x$  to find  $d_k^{z_l}$  and  $d_k^{z_u}$  using:

$$d_k^{z_l} = -z_{l,k} + \mu_j X_{l,k}^+ \mathbf{1}_l - X_{l,k}^+ Z_{l,k} d_k^x \quad (70)$$

$$d_k^{z_u} = -z_{u,k} + \mu_j X_{u,k}^+ \mathbf{1}_u - X_{u,k}^+ Z_{u,k} d_k^x \quad (71)$$

Note that the elements in  $z_l$  that are always 0 because their indices correspond to decision variables  $x$  unbounded from below can be ignored in implementation. The same can be done for the elements of  $z_u$  that are always 0 because their indices correspond to decision variables  $x$  unbounded from above.

## 4.4 Line search and filter

### 4.4.1 Maximum step size

Having computed the solution to the search directions linear system of equations to obtain  $(d_k^x, d_k^\lambda, d_k^{z_l}, d_k^{z_u})$ , now two step sizes  $\alpha_k, \alpha_k^z \in (0, 1]$  have to be determined in order to obtain the next iterate using:

$$x_{k+1} := x_k + \alpha_k d_k^x \quad (72)$$

$$\lambda_{k+1} := \lambda_k + \alpha_k d_k^\lambda \quad (73)$$

$$z_{l,k+1} := z_{l,k} + \alpha_k^z d_k^{z_l} \quad (74)$$

$$z_{u,k+1} := z_{u,k} + \alpha_k^z d_k^{z_u} \quad (75)$$

When the step sizes are chosen, the bounds on  $x$ ,  $z_l$  and  $z_u$  are enforced. Let  $\tau_j \in (0, 1)$  be the so-called fraction-to-the-boundary parameter given by:

$$\tau_j = \max\{\tau_{min}, 1 - \mu_j\} \quad (76)$$

where  $\tau_{min} \in (0, 1)$  is its minimum value. The step sizes are then chosen using the following fraction-to-the-boundary rule:

$$\alpha_k^{max} := \max \left\{ \alpha \in (0, 1] : \left( \tau_l + (1 - \tau_j)x_k \right) \leq \left( x_k + \alpha d_k^x \right) \leq \left( \tau_u + (1 - \tau_j)x_k \right) \right\} \quad (77)$$

$$\alpha_k^z := \max \left\{ \alpha \in (0, 1] : (z_{l,k} + \alpha d_k^{z_l} \geq (1 - \tau_j)z_{l,k}) \wedge (z_{u,k} + \alpha d_k^{z_u} \geq (1 - \tau_j)z_{u,k}) \right\} \quad (78)$$

where the actual step size  $\alpha_k \in (0, \alpha_k^{max}]$  is determined using a backtracking line search procedure exploring a decreasing sequence of trial step sizes:

$$\alpha_{k,l} = 2^{-l} \alpha_k^{max} \quad (79)$$

with  $l = 0, 1, 2, \dots$

#### 4.4.2 Line search

Choosing the step sizes is done using a so-called line search filter method. During line search, a step size is considered acceptable if it leads to an acceptable reduction in the objective value of the barrier problem  $\phi_\mu(\mathbf{x})$ , and/or the constraint violation  $\theta(\mathbf{x}) = \|\mathbf{c}(\mathbf{x})\|_1$  with a certain emphasis on the latter quantity. Let  $\mathbf{x}_k(\alpha_{k,l})$  be:

$$\mathbf{x}_k(\alpha_{k,l}) := \mathbf{x}_k + \alpha_{k,l} \mathbf{d}_k^x \quad (80)$$

There are 2 acceptable criteria used to accept a trial step size. The following conditions are called the "switching conditions":

$$\nabla \phi_{\mu_j}(\mathbf{x}_k)^T \mathbf{d}_k^x < 0 \quad (81)$$

$$\alpha_{k,l} [-\nabla \phi_{\mu_j}(\mathbf{x}_k)^T \mathbf{d}_k^x]^{S_\phi} > \delta [\theta(\mathbf{x}_k)]^{S_\theta} \quad (82)$$

for some constants  $\delta > 0$ ,  $S_\phi > 1$  and  $S_\theta \geq 1$ . If  $\theta(\mathbf{x}_k) \leq \theta^{min}$ , for some constant  $\theta^{min} \in (0, \infty)$ , and the following so-called "switching conditions" are satisfied, a step size  $\alpha_{k,l}$  is considered acceptable if the following so-called Armijo condition is satisfied:

$$\phi_{\mu_j}(\mathbf{x}_k(\alpha_{k,l})) \leq \phi_{\mu_j}(\mathbf{x}_k) + \eta_\phi \alpha_{k,l} \nabla \phi_{\mu_j}(\mathbf{x}_k)^T \mathbf{d}_k^x \quad (83)$$

for some constant  $\eta_\phi \in (0, 0.5)$  and  $(\theta(\mathbf{x}_k(\alpha_{k,l})), \phi_{\mu_j}(\mathbf{x}_k(\alpha_{k,l}))) \notin \mathcal{F}_k$  where  $\mathcal{F}_k$  is a filter set (defined below) used to reject a trial solution if its constraint violation is too high or if it has a close barrier objective value and constraint violation value as certain previous iterates  $\mathbf{x}_k$ . This filter method ensures that the algorithm cannot cycle between 2 points that alternately decrease the constraint violation and barrier objective value. If the trial step is rejected, a so-called second order correction (SOC) step is performed. The SOC step is explained below.

If either  $\theta(\mathbf{x}_k) > \theta^{min}$  or the switching conditions are violated, a step size  $\alpha_{k,l}$  is considered acceptable if either:

$$\theta(\mathbf{x}_k(\alpha_{k,l})) \leq (1 - \gamma_\theta) \theta(\mathbf{x}_k) \quad (84)$$

$$(\theta(\mathbf{x}_k(\alpha_{k,l})), \phi_{\mu_j}(\mathbf{x}_k(\alpha_{k,l}))) \notin \mathcal{F}_k \quad (85)$$

or

$$\phi_{\mu_j}(\mathbf{x}_k(\alpha_{k,l})) \leq \phi_{\mu_j}(\mathbf{x}_k) - \gamma_\phi \theta(\mathbf{x}_k) \quad (86)$$

$$(\theta(\mathbf{x}_k(\alpha_{k,l})), \phi_{\mu_j}(\mathbf{x}_k(\alpha_{k,l}))) \notin \mathcal{F}_k \quad (87)$$

is satisfied for some constants  $\gamma_\theta, \gamma_\phi \in (0, 1)$ . If the trial step is rejected, the SOC step is performed.

#### 4.4.3 Filter set

The filter set  $\mathcal{F}_k$  is initialized as:

$$\mathcal{F}_0 = \{(\theta, \phi) \in \mathbb{R}^2 : \theta \geq \theta^{max}\} \quad (88)$$

Then after every iteration  $k$  if an accepted trial point  $\mathbf{x}_k$  doesn't satisfy either the switching conditions or the Armijo rule, the filter set is updated as:

$$\mathcal{F}_{k+1} = \mathcal{F}_k \cup \left\{ (\theta, \phi) \in \mathbb{R}^2 : \left( \theta \geq (1 - \gamma_\theta)\theta(\mathbf{x}_k) \right) \wedge \left( \phi \geq \phi_{\mu_j}(\mathbf{x}_k) - \gamma_\phi\theta(\mathbf{x}_k) \right) \right\} \quad (89)$$

Every time the barrier parameter  $\mu$  is decreased in the outer iterations of the algorithm, the filter set is reset to its definition at  $k = 0$ .

If the backtracking procedure (SOC steps included) fails to find an acceptable trial step  $\alpha_{k,l} \geq \alpha_k^{min}$ , where:

$$\alpha_k^{min} := \begin{cases} \min \left\{ \gamma_\theta, \frac{\gamma_\theta \theta(\mathbf{x}_k)}{-\nabla \phi_{\mu_j}(\mathbf{x}_k)^T \mathbf{d}_k^x}, \frac{\delta[\theta(\mathbf{x}_k)^{S_\theta}]}{[-\nabla \phi_{\mu_j}(\mathbf{x}_k)^T \mathbf{d}_k^x]^{S_\phi}} \right\} & \left( \nabla \phi_{\mu_j}(\mathbf{x}_k)^T \mathbf{d}_k^x < 0 \right) \wedge \left( \theta(\mathbf{x}_k) \leq \theta^{min} \right) \\ \min \left\{ \gamma_\theta, \frac{\gamma_\phi \theta(\mathbf{x}_k)}{-\nabla \phi_{\mu_j}(\mathbf{x}_k)^T \mathbf{d}_k^x} \right\} & \left( \nabla \phi_{\mu_j}(\mathbf{x}_k)^T \mathbf{d}_k^x < 0 \right) \wedge \left( \theta(\mathbf{x}_k) > \theta^{min} \right) \\ \gamma_\theta & \text{otherwise,} \end{cases} \quad (90)$$

for some safety factor  $\gamma_\alpha \in (0, 1]$ , the algorithm reverts to a feasibility restoration phase. For details on the feasibility restoration phase, please refer to section 3.3 in Wächter and Biegler (2006).

#### 4.4.4 Second order correction

If a trial step  $\alpha_{k,l}$  is rejected for some iteration  $l$  in the backtracking procedure and  $\theta(\mathbf{x}_k(\alpha_{k,l})) \geq \theta(\mathbf{x}_k)$ , before moving on to the next trial, an SOC step is performed. The SOC step aims to reduce the infeasibility by applying an additional Newton-type step for the constraints at the point  $\mathbf{x}_k + \tilde{\mathbf{d}}_k^x$  where  $\tilde{\mathbf{d}}_k^x = \alpha_{k,l} \mathbf{d}_k^x$  using the Jacobian  $A_k^T$  at  $\mathbf{x}_k$  to obtain a direction update  $\mathbf{d}_k^{x,soc}$  by solving:

$$A_k^T \mathbf{d}_k^{x,soc} + \mathbf{c}(\mathbf{x}_k + \alpha_{k,l} \mathbf{d}_k^x) = \mathbf{0} \quad (91)$$

The new corrected search direction is then obtained from:

$$\mathbf{d}_k^{x,cor} = \alpha_{k,l} \mathbf{d}_k^x + \mathbf{d}_k^{x,soc} \quad (92)$$

The details of efficiently and carefully calculating  $\mathbf{d}_k^{x,soc}$  can be found in section 2.4 in Wächter and Biegler (2006). Once the corrected search direction  $\mathbf{d}_k^{x,cor}$  has been computed, we again apply the fraction-to-the-boundary rule:

$$\alpha_k^{soc} := \max \left\{ \alpha \in (0, 1] : \left( \tau l + (1 - \tau_j) \mathbf{x}_k \right) \leq \left( \mathbf{x}_k + \alpha \mathbf{d}_k^{x,cor} \right) \leq \left( \tau u + (1 - \tau_j) \mathbf{x}_k \right) \right\} \quad (93)$$

and check if the resulting trial point:

$$\mathbf{x}_k^{soc} := \mathbf{x}_k + \alpha_k^{soc} \mathbf{d}_k^{x,cor} \quad (94)$$

is acceptable to the filter and satisfies the acceptance criteria set earlier replacing  $\mathbf{x}(\alpha_{k,l})$  with  $\mathbf{x}_k^{soc}$  while keeping  $\mathbf{d}_k^x$  as-is. If the new iterate is still rejected, the correction step is repeated (replacing  $\alpha_{k,l}$  with  $\alpha_k^{soc}$ , and  $\mathbf{d}_k^x$  with  $\mathbf{d}_k^{x,cor}$  in the SOC step) unless the correction step has not decreased the constraint violation by a fraction  $\kappa_{soc} \in (0, 1)$  or a maximum number  $p^{max}$  of SOC steps has been performed. In that case, the original search direction  $\mathbf{d}_k^x$  is reverted to and the regular backtracking line search is resumed with a shorter step size.

#### 4.4.5 Accelerating heuristics

When multiple trial steps are rejected because the filter set is too strict or the progress made in the objective or constraint violation is insufficient, 2 accelerating heuristics are further employed to relax the acceptance criteria under some conditions. More on this can be found in section 3.2 in Wächter and Biegler (2006).

#### 4.4.6 Clamping

Finally after taking a step successfully, each variable  $z_{l_i}$  in  $\mathbf{z}_l$  and each variable  $z_{u_i}$  in  $\mathbf{z}_u$  then get clamped to an interval to avoid their extreme deviation:

$$z_{l_i} \in \begin{cases} \left[ \frac{\mu_j}{\kappa_\Sigma(x_i - l_i)}, \frac{\kappa_\Sigma \mu_j}{(x_i - l_i)} \right] & i \in \bar{\mathcal{I}}_l \cap \bar{\mathcal{I}}_u \\ \left[ \frac{(1 - \kappa_d) \mu_j}{\kappa_\Sigma(x_i - l_i)}, \frac{\kappa_\Sigma(1 - \kappa_d) \mu_j}{x_i - l_i} \right] & i \in \bar{\mathcal{I}}_l \setminus \bar{\mathcal{I}}_u \\ [0, 0] & otherwise, \end{cases} \quad (95)$$

$$z_{u_i} \in \begin{cases} \left[ \frac{\mu_j}{\kappa_\Sigma(u_i - x_i)}, \frac{\kappa_\Sigma \mu_j}{u_i - x_i} \right] & i \in \bar{\mathcal{I}}_l \cap \bar{\mathcal{I}}_u \\ \left[ \frac{(1 - \kappa_d) \mu_j}{\kappa_\Sigma(u_i - x_i)}, \frac{\kappa_\Sigma(1 - \kappa_d) \mu_j}{u_i - x_i} \right] & i \in \bar{\mathcal{I}}_u \setminus \bar{\mathcal{I}}_l \\ [0, 0] & otherwise. \end{cases} \quad (96)$$

This clamping step creates a safeguard needed for the global convergence proof of the algorithm (Wächter and Biegler, 2006). A large value for  $\kappa_\Sigma = 10^{10}$  is used to only minimally interfere with the Newton or quasi-Newton update.

## 5 Augmented Lagrangian algorithm

Both the MMA and IPOPT algorithms discussed previously require the full Jacobian of the constraints. This can be computationally prohibitive in some applications. In some cases, computing the Jacobian is much more expensive than calling the operator  $\mathbf{v} \rightarrow \mathbf{J}'\mathbf{v}$  where  $\mathbf{J}$  is the Jacobian of the constraints and  $\mathbf{v}$  is a vector. The first order augmented Lagrangian algorithm only requires this operator and not the full Jacobian matrix.

Consider the following inequality constrained optimization problem:

$$\begin{aligned}
 & \underset{\mathbf{x}}{\text{minimize}} && f_0(\mathbf{x}) \\
 & \text{subject to} && \\
 & f_i(\mathbf{x}) \leq 0 && \forall i = 1..I, \\
 & f_i(\mathbf{x}) = 0 && \forall i = I + 1..I + E, \\
 & l_j \leq x_j \leq u_j && \forall j = 1..V
 \end{aligned} \tag{97}$$

where  $I$  is the number of inequality constraints,  $E$  is the number of equality constraints and  $V$  is the number of decision variables,  $f_i$  is a potentially nonlinear scalar-valued function of  $\mathbf{x}$  and  $\mathbf{l}$  and  $\mathbf{u}$  are finite vectors of lower and upper bounds on the decision vector  $\mathbf{x}$  respectively.

The augmented Lagrangian algorithm (AugLag) is not a single algorithm but a family of algorithms in which the main idea is to reduce the above minimization problem to the following max-min formulation:

$$\underset{\lambda \in \mathbf{Y}_\lambda}{\text{maximize}} \quad \min_{\mathbf{l} \leq \mathbf{x} \leq \mathbf{u}} \left\{ \mathcal{L}_c(\mathbf{x}, \lambda) = f_0(\mathbf{x}) + \sum_{i=1}^{I+E} \lambda_i f_i(\mathbf{x}) + c \sum_{i=1}^I \max\{f_i(\mathbf{x}), 0\}^2 + c \sum_{i=I+1}^{I+E} |f_i(\mathbf{x})|^2 \right\} \tag{98}$$

for some constant  $c \geq 0$ , where  $\mathbf{Y}_\lambda = \{\lambda : \lambda_i \geq 0 \forall i \in [1, I], \lambda_i \in \mathbb{R} \forall i \in [I + 1, I + E]\}$  and  $\mathcal{L}_c$  is known as the augmented Lagrangian function. Assume all the functions used are continuous and twice differentiable.

Note how the inequality and equality constraints were relaxed and added to the objective but the bounds constraints were not. That is the inner minimization still needs to respect the bounds constraints. The choice of which constraints to relax and which to handle directly, as well as which optimizer to use for the inner optimization problem and which one to use for the outer problem can lead to different variants of the augmented Lagrangian family of algorithms. The optimal value of the augmented Lagrangian function subject to some relaxed constraints is known to be always less than or equal to the optimal value of the original optimization problem. It is also known that the outer maximization problem is always concave (i.e. tractable) even if the objective and constraint functions were non-convex (Bertsekas, 1996).

It is known (Bertsekas, 1996) that if:

1. The original problem has an isolated set of local minima  $\mathbf{X}^*$  which is compact,
2. The constant  $c$  is increased in a sequence  $\{c_k\}$  such that  $0 < c_k < c_{k+1}$ , and

3. The Lagrangian multipliers  $\lambda$  follow an arbitrary bounded sequence  $\{\lambda_k\}$  where  $\lambda_{k,i}$  is the  $i^{\text{th}}$  element of  $\lambda_k$

then there exists a sub-sequence of points  $\{\mathbf{x}_k\}_K$  converging to a point  $\mathbf{x}^* \in \mathbf{X}^*$  such that  $\mathbf{x}_k$  is a local minimum for the following augmented Lagrangian sub-problem:

$$\underset{l \leq \mathbf{x} \leq \mathbf{u}}{\text{minimize}} \quad \left\{ \mathcal{L}_{c_k}(\mathbf{x}, \lambda_k) = f_0(\mathbf{x}) + \sum_{i=1}^{I+E} \lambda_{k,i} f_i(\mathbf{x}) + c_k \sum_{i=1}^I \max\{f_i(\mathbf{x}), 0\}^2 + c_k \sum_{i=I+1}^{I+E} f_i(\mathbf{x})^2 \right\} \quad (99)$$

If additionally each  $\mathbf{x}_k$  is a global minimum of the augmented Lagrangian sub-problem with  $c_k$  and  $\lambda_k$ , and the feasible set is compact, e.g.  $l$  and  $u$  are finite, then every limit point of the sequence  $\{\mathbf{x}_k\}$  is a global minimum of the original problem. Note that it is possible for the original problem to have a (unique) global minimum but the augmented Lagrangian sub-problem is unbounded from below. This is a weakness of the AugLag algorithm that needs to be addressed in implementation. However if the feasible domain is compact, this weakness is eliminated.

If additionally the sub-sequence  $\{\mathbf{x}_k\}_K$  converges to a point  $\mathbf{x}^* \in \mathbf{X}^*$  that satisfies one of the constraint qualification conditions, then the sequence  $\{\text{proj}_{Y_\lambda}(\lambda_k + c_k \mathbf{G}_k)\}$  converges to a point  $\lambda^*$  such that  $(\mathbf{x}^*, \lambda^*)$  satisfy the first order sufficient KKT optimality conditions, where  $\mathbf{G}_k$  is a vector of length  $I + E$  and whose  $i^{\text{th}}$  element is  $f_i(\mathbf{x}_k)$  (Bertsekas, 1996). Another weakness of the AugLag algorithm is that the sub-problem's optimizer may converge to a point that doesn't satisfy any of the constraint qualification conditions for the original problem. One case where this is guaranteed to happen is if the original problem is infeasible so the quadratic term dominates the augmented Lagrangian function as  $c \rightarrow \infty$  and every (approximate) KKT point of the sub-problem must (approximately) violate the constraint qualification conditions of the original problem.

Instead of choosing  $\lambda_k$  arbitrarily, if  $\lambda_{k+1}$  is chosen as:

$$\lambda_{k+1} = \text{proj}_{Y_\lambda}(\lambda_k + \alpha_k \mathbf{G}_k) \quad (100)$$

where  $\alpha_k = c_k$ , convergence can often be significantly accelerated and the algorithm typically converges at lower values of  $c_k$ . The dual update step size  $\alpha_k$  can also be improved to a more "optimal" choice than  $c_k$  or a Newton/quasi-Newton direction can be used instead of  $\mathbf{G}_k$ . These can lead to faster convergence but come at an increased computational cost. For more on the augmented Lagrangian algorithm, see Bertsekas (1996).

In the context of topology optimization, one challenge with AugLag is that the Hessian of the augmented Lagrangian function becomes rather ill-conditioned for large values of  $c$  where the condition number of the Hessian matrix goes to  $\infty$  as  $c \rightarrow \infty$ . Given that computing the Hessian of the augmented Lagrangian in topology optimization is generally intractable, one must rely on first order methods to optimize the sub-problem. However even quasi-Newton methods may very well encounter difficulty converging if the Hessian approximation is not good enough and/or the starting point is not near a solution. Adding to this the scaling issues that are commonly found in topology optimization and this makes AugLag extremely difficult to fine-tune in practice

when using it in topology optimization. The convergence speed and numerical stability of AugLag can be extremely sensitive to the initial value of  $c$ , its increment rate, the initial solution, and the sub-problem's optimizer's terminating conditions.

The main computational advantage of AugLag is that it can handle block constraints more efficiently. Let  $\mathbf{G}(\mathbf{x})$  be the vector-valued function whose components are the constraint functions  $f_i(\mathbf{x}) \forall i \in [1, I + E]$  and let  $\mathbf{M}(\mathbf{x}) = \text{proj}_{\mathbf{Y}_\lambda} \mathbf{G}(\mathbf{x})$  be the projection of  $\mathbf{G}(\mathbf{x})$  on  $\mathbf{Y}_\lambda$ . The augmented Lagrangian formulation can therefore be written as:

$$\underset{\lambda \in \mathbf{Y}_\lambda}{\text{maximize}} \quad \min_{l \leq \mathbf{x} \leq u} \{ \mathcal{L}_c(\mathbf{x}, \lambda) = f(\mathbf{x}) + (\lambda + c\mathbf{M}(\mathbf{x}))^T \mathbf{G}(\mathbf{x}) \} \quad (101)$$

Note that the gradient of  $\mathcal{L}_c(\mathbf{x}, \lambda)$  wrt  $\mathbf{x}$ ,  $\nabla_{\mathbf{x}} \mathcal{L}_c(\mathbf{x}, \lambda)$ , can be written in terms of the gradient of  $f(\mathbf{x})$ ,  $\nabla_{\mathbf{x}} f(\mathbf{x})$ , and the Jacobian of  $\mathbf{G}(\mathbf{x})$ ,  $\nabla_{\mathbf{x}} \mathbf{G}(\mathbf{x})$ , as:

$$\nabla_{\mathbf{x}} \mathcal{L}_c(\mathbf{x}, \lambda) = \nabla_{\mathbf{x}} f(\mathbf{x}) + \nabla_{\mathbf{x}} \mathbf{G}(\mathbf{x})^T (\lambda + 2c\mathbf{M}(\mathbf{x})) \quad (102)$$

To compute the above gradient, the full Jacobian of the constraints need not be built. Only the operator  $\mathbf{w} \rightarrow \nabla_{\mathbf{x}} \mathbf{G}(\mathbf{x})^T \mathbf{w}$  needs to be defined. For the purposes of this paper, the term *block constraint* will be used to refer to constraints on functions  $\mathbf{G}(\mathbf{x})$  where the operator  $\mathbf{w} \rightarrow \nabla_{\mathbf{x}} \mathbf{G}(\mathbf{x})^T \mathbf{w}$  can be defined more efficiently than simply calculating the entire Jacobian,  $\nabla_{\mathbf{x}} \mathbf{G}(\mathbf{x})$ , and then multiplying its transpose by  $\mathbf{w}$ .

## References

- M. ApS. Mosek modeling cookbook, 2018. URL <https://docs.mosek.com/MOSEKModelingCookbook-v2.pdf>.
- S. Arora and B. Barak. *Computational Complexity: A Modern Approach*. 2007. URL <https://theory.cs.princeton.edu/complexity/book.pdf>.
- D. P. Bertsekas. *Constrained Optimization and Lagrange Multiplier Methods*. Athena Scientific, 1996.
- M. Besançon, A. Carderera, and S. Pokutta. Frankwolfe.jl: a high-performance and flexible toolbox for frank-wolfe algorithms and conditional gradients, 2021.
- S. Boyd and L. Vandenberghe. *Convex Optimization*. 2009.
- G. H. Golub and C. F. V. Loan. *Matrix Computations*, volume 10. 1996.
- J. Nocedal and S. J. Wright. *Numerical Optimization*. Springer Sc, 2006. ISBN 354035445X. doi: 10.1002/lsm.21040.
- N. Parikh and S. Boyd. *Proximal Algorithms*. 2013.
- K. Svanberg. The method of moving asymptotes - a new method for structural optimization. *International Journal for Numerical Methods in Engineering*, 24(2):359–373, 1987.
- K. Svanberg. A Class of Globally Convergent Optimization Methods Based on Conservative Convex Separable Approximations. *SIAM Journal on Optimization*, 12(2): 555–573, 2002.
- A. Wächter and L. T. Biegler. *On the Implementation of a Primal-Dual Interior Point Filter Line Search Algorithm for Large-Scale Nonlinear Programming*, volume 106. 2006. ISBN 1010700405.