Activity Analysis: The Qualitative Analysis of Stationary Points for Optimal Reasoning

Brian C. Williams Xerox Palo Alto Research Center 3333 Coyote Hill Road, Palo Alto, CA 94304 USA bwilliams@parc.xerox.com

Abstract

We present a theory of a modeler's problem decomposition skills in the context of optimal reasoning — the use of qualitative modeling to strategically guide numerical explorations of objective space. Our technique, called activity analysis, applies to the pervasive family of linear and non-linear, constrained optimization problems, and easily integrates with any existing numerical approach. Activity analysis draws from the power of two seemingly divergent perspectives the global conflict-based approaches of combinatorial satisficing search, and the local gradientbased approaches of continuous optimization combined with the underlying insights of engineering monotonicity analysis. The result is an approach that strategically cuts away subspaces that it can quickly rule out as suboptimal, and then guides the numerical methods to the remaining subspaces.

Introduction and Example

Our goal is to capture a modeler's tacit skill at decomposing physical models and its application to focusing reasoning. This work is ultimately directed towards the contruction of "self modeling" systems, operating in embedded, real time situations. This article explores the modeler's decompositional skills (Williams & Raiman 1994) in the context of optimal reasoning ---the use of qualitative modeling to strategically guide gradient-based and other numerical explorations of objective spaces. Optimal reasoning is crucial for embedded systems, where numerical methods are key to such areas as estimation, control, inductive learning and vision. The technique we present, called activity analysis, applies to the pervasive family of linear and nonlinear, constrained optimization problems, and easily integrates with any existing numerical approaches.

Activity analysis is striking in the way it merges together two styles of search that are traditionally viewed as quite disparate: first is the more strategic, conflict-based approaches used in combinatorial, satisficing search to eliminate finite, inconsistent subspaces (e.g., (de Kleer & Williams 1987)). The second is the Jonathan Cagan Department of Mechanical Engineering Carnegie Mellon University Pittsburgh, PA 15213 USA cagan+@cmu.edu

rich suite of more tactical, numeric methods(Vanderplaats 1984) used in continuous optimizing search to climb locally but monotonically towards the optimum. Activity analysis draws from the power of both perspectives, strategically cutting away subspaces that it can quickly rule out as suboptimal, and then guiding the numerical methods to the remaining subspaces.

The power of activity analysis to eliminate large suboptimal subspaces is derived from *Qualitative KT*, an abstraction in *qualitative vector algebra* of the foundational Kuhn-Tucker (KT) condition of optimization theory. The underlying algorithm achieves simplicity and completeness, by introducing the concept of generating *prime implicating assignments* of linear, qualitatice vector equations. This process of ruling out feasible, but suboptimal subspaces in a continuous domain, nicely parallels the use of conflicts and prime implicant generation for combinatorial, satisficing search. The end result is a method that achieves parsimonious descriptions, guarantees correctness, and maximizes the filtering achieved from QKT.

Finally, activity analysis can be thought of as automating the underlying principle about monotonicity used by the simplex method to examine only the vertices of the linear feasible space. It then generalizes and automatically applies this principle to nonlinear programming problems.



Figure 1: Hydraulic Cylinder

To demonstrate the task consider the design of a hydraulic cylinder, a classic optimization problem, introduced by Wilde (Wilde 1975) to demonstrate the related technique of monotonicity analysis. The cylinder (figure 1) delivers force f, through input pressure p. Weight is modeled as inside diameter (i) plus twice the cylinder thickness (t), force (f) as pressure (p) times cylinder area, and hoop stress (s) as pressure times diameter acting across the thickness. The task is to find a parametric solution that minimizes cylinder weight, while satisfying constraints including positivity of variables (i, s, t, p, f > 0), maximum pressure (P) and stress (S), and minimum force (F) and thickness (T) (design variables are in lowercase, fixed parameters in uppercase, and equality and inequality constraints are labeled h_i and g_i , respectively): Minimize i + 2t,

subject to:

$s - \frac{pi}{2t}$		0,	$(h_1 = 0):$	T - t	\leq	0,	$(g_2 \leq 0)$
$f = \frac{\pi i^3}{4}p$		0,	$(h_2 = 0):$	p - P	\leq	0,	$(g_3 \leq 0)$
F - f	\leq	0,	$(g_1 \leq 0):$	s - S	\leq	0,	$(g_4 \leq 0)$

Given this symbolic formulation, activity analysis uses qualitative arguments to classify regions of the design space where optima might lie and where they cannot. After eliminating suboptimal regions, each remaining region identifies the solution as possibly lying on the intersection of one or more constraint boundaries. Each region reduces the dimensionality of the problem by the number of intersecting boundaries, thus significantly increasing the ease with which a solution can be found. In particular, for the cylinder problem activity analysis concludes there are two subspaces of the design space that could contain the optima, one subspace in which g_1 and g_4 become strict equalities, and a second in which all but g_4 become strict equalities. The new problem formulation finds the optima of the two spaces and combines the results as follows (where "arg min" returns a set of optima):

Given: vector $\mathbf{x} = (istpf)^T$,

1. Let $\mathbf{Y} = \arg \min_{\mathbf{X}} (i + 2t)$, subject to:

$$(h_1=0) \quad (g_1=0) \quad (g_3 \le 0) \ (h_2=0) \quad (g_2 \le 0) \quad (g_4=0).$$

2. Let $\mathbf{Z} = \arg \min_{\mathbf{X}} (i + 2t)$, subject to:

3. Return $\arg \min_{\mathbf{x}}(i+2t)$, subject to:

$$\mathbf{x} \in \mathbf{Y} \cup \mathbf{Z}.$$

Originally, the problem has a 3 dimensional space to be explored (3 degrees of freedom – DOF) resulting from 5 variables, 2 equality constraints. The reformulated problem rules out the interior and boundaries, except some intersections. The first remaining subspace corresponds to a *line* (1 DOF) produced by the intersection of the g_1 and g_4 constraint boundaries with the h_i . The second remaining space is a point (0 DOF) produced by the intersection of g_1 , g_2 , g_3 and the h_i . Thus finding a solution to the first problem involves a single, one dimensional line search, and the second involves solving the system of equalities to find the unique solution. Using parameter values F=1000 lbf, T=.05 in, S=30000 psi, T=1000 in, applying matlab to the original problem took 46.3 seconds. The optimal solution lies in **Z**, which took only 8.1 seconds to run; no feasible solution exists in Y for these parameter values.

Activity analysis draws inspiration from monotonicity analysis (MA) (Papalambros & Wilde 1979; Papalambros 1982). Monotonicity analysis began as a set of principles and methods used by modelers to identify ill-posed problems and to partially solve them, based on monotonic arguments alone. These principles were encoded in several rule-based implementations (Azram & Papalambros 1984; Choy & Agogino 1986; Rao & Papalambros 1987; Hansen, Jaumard, & Lu 1989), presented informally as heuristic methods.

The problem activity analysis addresses is similar in spirit to that of MA; nevertheless, the approach is quite different. First, activity analysis operates directly on an abstraction (QKT) of the Kuhn-Tucker (KT) conditions of optimization theory. While much easier to apply, QKT and KT are equivalent for the task, given only knowledge of monotonicities. Second, activity analysis provides a precise formulation of the problem in terms of *minimal pstationary coverings*, that guarantees the solution is parsimonious, maximizes the filtering derived from QKT, and insures correctness. Finally, a mapping to *prime assignments* and the introduction of a simple but complete prime assignment engine guarantees that these three properties are achieved.

Stationary Points and Kuhn-Tucker

For a point $\mathbf{x} *$ to be an optimum it is necessary that the point be *stationary*, that is any "down hill" direction is blocked by the constraints. Activity analysis exploits this fact to eliminate sets of points that can quickly be proven to be *nonstationary*, using a condition we call *Qualitative Kuhn-Tucker* (QKT). This section introduces the optimization problem, the concept of stationary point, and the traditional algebraic (Kuhn-Tucker) condition for testing stationary points. Activity analysis applies to the pervasive family of linear and non-linear, constrained optimization problems $OP = \langle \mathbf{x}, f, \mathbf{g}, \mathbf{h} \rangle$:

Find
$$\mathbf{x} * = \arg \min \quad f(\mathbf{x})$$

subject to: $\mathbf{g}(\mathbf{x}) \leq \mathbf{0}$
 $\mathbf{h}(\mathbf{x}) = \mathbf{0}$,

where column vectors are denoted in bold (e.g., \mathbf{x}, \mathbf{x}^* , $\mathbf{g}(\mathbf{x})$ and $\mathbf{h}(\mathbf{x})$), $f(\mathbf{x})$ is the objective function, $\mathbf{g}(\mathbf{x})$ is a vector of *inequality constraints* and $\mathbf{h}(\mathbf{x})$ is a vector of equality constraints. A point $\mathbf{x} \in \mathbb{R}^n$ is feasible if it satisfies the constraints, and feasible space $\mathcal{F} \subseteq \mathbb{R}^n$

denotes all feasible points (represented $\mathcal{F} = \langle \mathbf{g}, \mathbf{h} \rangle$). A feasible direction \vec{s} from a feasible point is one through which a non-zero distance can be moved before hitting a constraint boundary. $f(\mathbf{x})$ is decreasing at \mathbf{x} in direction \vec{s} if $\nabla f(\mathbf{x}) \cdot \vec{s} < 0$. Finally, a point is stationary (denoted $\mathbf{x}*$) if any direction that decreases the objective is infeasible. The Kuhn-Tucker (KT) conditions (Kuhn & Tucker 1951) provide a set of vector equations that are satisfied for a feasible point $\mathbf{x}*$ exactly when that point is stationary:

$$\nabla f(\mathbf{x}^*) + \lambda^T \nabla \mathbf{h}(\mathbf{x}^*) + \mu^T \nabla \mathbf{g}(\mathbf{x}^*) = \mathbf{0}^T \quad (\mathrm{KT1})$$

subject to

$$\mu^{T} \mathbf{g}(\mathbf{x}^{*}) = \mathbf{0}^{T}, \qquad (KT2)$$

$$\mu \geq \mathbf{0}. \qquad (KT3)$$

 μ^T transposes column vector μ to a row. Gradients ∇f , $\nabla \mathbf{g}$ and $\nabla \mathbf{h}$ denote Jacobian matrices. ∇f is a row vector $\left(\frac{\partial f}{\partial x_1} \dots \frac{\partial f}{\partial x_n}\right)$. $\nabla \mathbf{g}$ and $\nabla \mathbf{h}$ are matrices $\left(\frac{\partial g_i}{\partial x_j}\right)$ and $\left(\frac{\partial h_i}{\partial x_j}\right)$, respectively, where (a_{ij}) denotes a matrix whose element in the ith row and jth column is a_{ij} , for all i and j. For example, KT1 and KT2 are equivalences between row vectors, and KT3 is a relation between column vectors.

In KT1 the $- \bigtriangledown f$ term denotes directions of decreasing objective from \mathbf{x}_* , the term $(\lambda^T \bigtriangledown \mathbf{h}(\mathbf{x}_*) +$ $\mu^T \bigtriangledown \mathbf{g}(\mathbf{x}^*)$ denotes infeasible directions from \mathbf{x}^* , and the equality says the decreasing directions are all infeasible; hence, \mathbf{x} * is stationary. More specifically, \vec{s} decreases the objective if it has a component in the $- \bigtriangledown f$ direction ($\vec{s} \cdot \bigtriangledown f < 0$). A direction is infeasible with respect to inequality constraint $g_i(\mathbf{x}^*)$ if \mathbf{x}^* lies on the constraint boundary $(g_i(\mathbf{x}*) = 0)$ and it has a component in the $+ \bigtriangledown g_i(\mathbf{x}^*)$ direction. A direction is infeasible with respect to equality constraint $h_i(\mathbf{x}^*)$ if it has a component in either the $-\nabla h_i(\mathbf{x}^*)$ or $+\nabla h_i(\mathbf{x}^*)$ direction. Most importantly, if x* lies on multiple constraint boundaries, then an infeasible direction has a component which is a linear, weighted combination of the above gradients for these constraints. The weights are μ and λ , (called Lagrange multipliers), and the combination is $\mu^T \bigtriangledown \mathbf{g} + \lambda^T \bigtriangledown \mathbf{h}$ subject to KT2 and KT3. Hence all decreasing directions are infeasible when $- \bigtriangledown f$ equals one of these linear combinations (KT1). Figure 2 shows an example of ∇f and ∇g gradient vectors, and the combined weighted vector, which exactly cancels $\bigtriangledown f$.

A key property of KT is that it identifies active inequality constraints. Intuitively, a constraint $[g_i]$ is active at a point **x** when **x** is on the constraint boundary and the direction of decreasing objective, ∇f , is pointing into the boundary. When this is true μ_i is positive. The basis of our approach is to conclude, by looking at signs of μ , that the stationary points lie at the intersection of the constraint boundaries. One or more constraints have been identified as active, hence the name activity analysis.



Figure 2: Example gradient vector diagram for KT.

Qualitative KT Conditions

Qualitative KT (QKT) is an abstraction of KT that is a necessary, but insufficient, condition for a point being stationary. It is the means by which activity analysis quickly rules out suboptimal subspaces. Qualitative properties used by QKT to test a point \mathbf{x} include whether each constraint is active at \mathbf{x} , and the quadrant of the coordinate axes each gradient ∇f , $\nabla \mathbf{g}$ and $\nabla \mathbf{h}$ lies within. These properties can be extracted quickly and hold uniformly for large subsets of the feasible space, and parameterized families of optimization problems. QKT, its proof (see (Williams 1994)), and manipulations by activity analysis rely on a matrix version of SR1 - a hybrid algebra combining signs and reals. This algebra behaves as one expects given a familiarity with (scalar) sign algebra and traditional matrix algebra (see (Williams 1994; 1991)). Derived from KT, QKT states that a feasible point $\mathbf{x} *$ is stationary only if (QKT1):

$$\left[\bigtriangledown f(\mathbf{x}^*) \right] + \left[\lambda \right]^T \left[\bigtriangledown \mathbf{h}(\mathbf{x}^*) \right] + \left[\mu \right]^T \left[\bigtriangledown \mathbf{g}(\mathbf{x}^*) \right] \supseteq \mathbf{0}^T,$$

subject to

$$\begin{aligned} [\mu]^T[\mathbf{g}(\mathbf{x}*)] &= \mathbf{0}^T, \text{ and } (\text{QKT2}) \\ [\mu_i] &\neq \hat{-}, \quad (\text{QKT3}) \end{aligned}$$

where $[\mathbf{v}]$, called a sign vector, denotes the signs of the elements of \mathbf{v} , such that $[v_i] \in \{\hat{-}, 0, \hat{+}\}$. Recall KT said that to be stationary there must exist a weighted sum $(\vec{\mathbf{w}})$ of $\bigtriangledown \mathbf{g}$ and $\bigtriangledown \mathbf{h}$ that exactly cancels $\bigtriangledown f$ (note $\vec{\mathbf{w}}$ is a row vector). QKT says a point is nonstationary unless there exists a $\vec{\mathbf{w}}$ that lies in the quadrant diagonal from that which contains $\bigtriangledown f$. For example, in figure $2 \bigtriangledown f$ lies in the upper left quadrant; thus, a $\vec{\mathbf{w}}$ must exist that lies in the lower right. The sign vector $[\mathbf{v}]$ denotes the quadrant containing a vector \mathbf{v} , and each component $[v_i]$ describes where \mathbf{v} lies relative to the $v_i = 0$ plane. For example, $[\vec{\mathbf{w}}] = (\hat{+} -)$ indicates that $\vec{\mathbf{w}}$ is in the lower right. Using this algebraic representation, the condition on diagonal quadrants becomes $-[\bigtriangledown f] = [\vec{\mathbf{w}}]$.

Using only knowledge of the quadrant each constraint's gradient lies within and whether each constraint is active (indicated by the signs of the lagrange multipliers $[\mu]$ and $[\lambda]$), we know from KT that the quadrants $\vec{\mathbf{w}}$ may lie within are a subspace of those described by $[\mu]^T[\bigtriangledown \mathbf{g}] + [\lambda]^T[\bigtriangledown \mathbf{h}]$. Thus, $-[\bigtriangledown f] =$ $[\vec{\mathbf{w}}] \subseteq [\mu]^T[\bigtriangledown \mathbf{g}] + [\lambda]^T[\bigtriangledown \mathbf{h}]$ (i.e., QKT1). For example, in figure 2 since $\bigtriangledown g_1 (= (\hat{+} + \hat{+}))$ lies in the upper right and $\bigtriangledown g_2 (= (\hat{-} - \hat{-}))$ lies in the lower left, it is possible for a $\vec{\mathbf{w}}$ to lie in the lower right; thus, any \mathbf{x} satisfying these conditions may be stationary. But suppose $\bigtriangledown g_1$ is replaced with $\bigtriangledown g'_1$, which lies in the upper left for points in some subspace $\mathcal{F}1 \subseteq \mathcal{F}$. Then $\vec{\mathbf{w}}$ may lie in the upper or lower left, but not the lower right; thus, all points in $\mathcal{F}1$ must be nonstationary. That is, evaluating $-[\bigtriangledown f] = [\mu]^T[\bigtriangledown g]$ for $\bigtriangledown g_1$ and then $\bigtriangledown g'_1$:

$$\begin{pmatrix} \hat{+} & \hat{-} \end{pmatrix} \subseteq \begin{pmatrix} \hat{?} & \hat{?} \end{pmatrix} = \begin{pmatrix} \hat{+} & \hat{+} \end{pmatrix} \begin{pmatrix} \hat{+} & \hat{+} \\ \hat{-} & \hat{-} \end{pmatrix}$$
but
$$\begin{pmatrix} \hat{+} & \hat{-} \end{pmatrix} \not\subseteq \begin{pmatrix} \hat{-} & \hat{?} \end{pmatrix} = \begin{pmatrix} \hat{+} & \hat{+} \end{pmatrix} \begin{pmatrix} \hat{-} & \hat{+} \\ \hat{-} & \hat{-} \end{pmatrix}$$

It is this second type of conclusion, made from only qualitative properties, that activity analysis uses to eliminate feasible subspaces of nonstationary points.

Next, to instantiate QKT1 on optimization problem $OP \equiv \langle \mathbf{x}, f, \mathbf{g}, \mathbf{h} \rangle$:

- 1. Compute Jacobians $\bigtriangledown f$, $\bigtriangledown g$ and $\bigtriangledown h$ by symbolic differentiation.
- 2. Compute signs of Jacobians. For each element,
- (a) replace real operators with sign operators, using properties $[a+b] \subseteq [a] + [b]$, [ab] = [a][b], [a/b] = [a]/[b] and [-a] = -[a].
- (b) Substitute for sign variables [a] using positivity conditions ([a] = +̂), and perform sign arithmetic (e.g., [5] ⇒ +̂, (-̂) + (-̂) ⇒ -̂).
- 3. Expand QKT1 by expanding matrix sums and products.

Returning to the hydraulic cylinder problem from the introduction, recall that \mathbf{x} is the vector $(itfsp)^T$, the objective $f(\mathbf{x})$ is i + 2t, and the constraint vectors are:

$$\mathbf{h} = \left(\begin{array}{cc} s - \frac{pi}{2t} & f - \frac{\pi i^2}{4}p \end{array} \right)^T, \\ \mathbf{g} = \left(\begin{array}{cc} F - f & T - t & p - P & s - S \end{array} \right).^T$$

The following shows $[\nabla h]$ after steps 2a (middle) and 2b (right):

$$[\nabla \mathbf{h}] = \begin{pmatrix} \frac{-[p]}{[2][t]} & \frac{[p][i]}{[2][t]^2} & 0 & [1] & \frac{-[i]}{[2][t]} \\ \frac{-[\pi][i]}{[2]}[p] & 0 & [1] & 0 & \frac{-[\pi][i]^2}{[4]} \end{pmatrix} \\ = \begin{pmatrix} \hat{-} & \hat{+} & 0 & \hat{+} & \hat{-} \\ \hat{-} & 0 & \hat{+} & 0 & \hat{-} \end{pmatrix}.$$

Repeating for $[\nabla f]$ and $[\nabla g]$, and inserting into QKT:

$$\mathbf{0}^{T} \subseteq \begin{pmatrix} \hat{+} \\ \hat{+} \\ 0 \\ 0 \\ 0 \end{pmatrix}^{T} + \lambda^{T} \begin{pmatrix} \hat{-} & \hat{+} & 0 & \hat{+} & \hat{-} \\ \hat{-} & 0 & \hat{+} & 0 & \hat{-} \end{pmatrix}$$

$$+\mu^{T} \left(\begin{array}{ccccc} 0 & 0 & \dot{-} & 0 & 0 \\ 0 & \dot{-} & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & \dot{+} \\ 0 & 0 & 0 & \dot{+} & 0 \end{array} \right)$$

Expanding matrix operations for step 3 results in equations QKT1(1)-(5):

Г	0	⊆	$(+) - [\lambda_1] - [\lambda_2]$	(1)	0	⊆	$[\mu_4] + [\lambda_1]$	(4)
	0	\subseteq	$(\dot{+}) - [\mu_2] + [\lambda_1]$	(2)	0	\subseteq	$[\mu_3] - [\lambda_1] - [\lambda_2]$	(5)
L	0	Ē	$-[\mu_1] + [\lambda_2]$	(3)				

Note that the computation of sign matrices in step 2 is extremely simple, but suprisingly adequate for many problems. The symbolic algebra system Minima (Williams 1991) provides a general tool for deducing the signs of sensitivities (e.g., $\left[\frac{\partial f(\mathbf{X})}{\partial x_i}\right]$) subject to \mathbf{x} satisfying the equality and inequality constraints. Having achieved an easily evaluable condition that is sufficient for testing the suboptimality of infinite subspaces, we turn to its use for strategically focussing optimization.

Activity Analysis and Prime Assignments

Activity analysis reduces an optimization problem to a set of simpler subproblems by "cutting" out feasible subspaces that are suboptimal. These subspaces contain all and only those points that are provably nonstationary by QKT (see (Williams 1994)). The output of activity analysis is a concise description of the remainder, called a minimal pstationary covering ("p-" stands for "possible" according to QKT). It is a set of feasible subspaces (and corresponding optimization problems), at least one of which is guaranteed to contain the true optimum. What is key is that the descriptions are parsimonious, they maximize the "filtering" achievable from QKT, and are always correct (these three properties are theorems, stated precisely in (Williams 1994)). This section states and demonstrates the activity analysis problem, and a sound and complete solution algorithm. The core is a mapping between minimal pstationary subspaces and prime assignments, and a general prime assignment engine for arbitrary systems of linear sign equations.

To start we say a point is pnonstationary if it follows from QKT that it is nonstationary; otherwise, it is pstationary. A feasible subspace is pstationary if all its points are pstationary, and pnonstationary if all its points are pnonstationary. Activity analysis maximizes its use of QKT while preserving correctness by eliminating exactly the phonstationary subspaces from its description of the feasible space. This description is built from a set Σ whose elements result from strengthening one or more of the inequality constraints $g_i \leq 0$ to strict equalities $g_i = 0$; that is, Σ is the powerset of constraint boundary intersections. The description (called a minimal pstationary covering), covers the pstationary points by collecting all pstationary subspaces that are maximal under superset. These cover every pstationary subspace. The activity analysis problem is then: given optimization problem $OP = \langle \mathbf{x}, f, \mathbf{g}, \mathbf{h} \rangle$ and instantiation of QKT (=QKT(OP)), construct the minimal pstationary covering C.

Mapping QKT(OP) to C relies on two observations: First, from QKT2 ($\equiv [\mu_i(\mathbf{x})][g_i(\mathbf{x})] = 0$) it follows that $[\mu_i(\mathbf{x})] = \hat{+} - g_i(\mathbf{x}) = 0$ (denoted R1). That is, any point where $[\mu_i] = \hat{+}$ must be on the $g_i = 0$ constraint boundary. Thus, when activity analysis shows that a subspace of pstationary points makes $[\mu_i] = \hat{+}$ for one or more g_i 's, it concludes that these points lie along the intersection of the g_i boundaries. Second, a particular set of variable assignments for QKT1, called *prime (implicating) assignments*, directly maps to the minimal pstationary covering by applying the first observation. The key here is that achieving parsimony, maximum filtering and correctness reduces to generating complete prime assignments.

The following properties, stated informally here, are given as definitions and theorems in (Williams 1994). First, a (partial) assignment to $[\mathbf{x}]$ is a set α which assigns each $[x_i]$ at most one value, $\alpha \subseteq \{[x_i] = s \mid [x_i] \in$ $\mathbf{x}, s \in \{\hat{-}, 0, \hat{+}\}\}$. We are interested in the consistent assignments to QKT1, where the $[\mathbf{x}]$ to be assigned is a vector of lagrange multipliers $([\mu]^T[\lambda]^T)^T$. Additionally, the consistent assignments must also satisfy the restriction of QKT3 $([\mu] \neq \hat{-})$. Note that each consistent assignment C has a corresponding subset S of feasible space, produced by applying R1 to the assignment and then adding the resulting active constraints to the original constraint set. S has the property that every point in S satisfies C.

Next, an *implicating assignment* γ is a consistent assignment to QKT1, such that whenever an extension to γ satisfies restriction QKT3, it also is consistent with QKT1. That is, assignment γ *implies* QKT1 under restriction QKT3. An implicating assignment has the important property that every point in its corresponding subspace S satisfies QKT. Thus S is a pstationary subspace.

Finally, a prime assignment P is an implicating assignment no proper subset of which is also an implicating assignment. Thus P's corresponding S is a maximal pstationary subspace. Conversely, every maximal pstationary subspace is the corresponding subspace of some prime assignment. Thus the set of subspaces corresponding to all prime assignments is a minimal pstationary covering.

To produce all primes for QKT1, our prime assignment engine first computes the primes P_i of each scalar equation in QKT1, then combines them using minimal set covering. Pulling this all together, the activity analysis algorithm is:

Given problem $OP = \langle \mathbf{x}, f, \mathbf{g}, \mathbf{h} \rangle$:

- 1. Instantiate QKT1 (given earlier) $\rightarrow QKT1(OP)$,
- 2. Compute prime assignments P_i of each $QKT1_i(OP) \in QKT1(OP)$,

- 3. Compute minimal set covering of $P_i \rightarrow P$, deleting inconsistent assignments,
- 4. Extract minimal sets of $[\mu_i] = \hat{+}$ assignments from $P \rightarrow U$,
- 5. Map each element of U to a maximal pstationary subspace by applying $[\mu_i(\mathbf{x})] = \hat{+} \rightarrow g_i(\mathbf{x}) = 0$, producing a covering.
- 6. Formulate and return a new optimization problem from this covering.

Step one was demonstrated in the previous section. For steps two and three we note that QKT1 is an instance of a linear system of sign equations (denoted L([x])) and solve the prime assignment problem for arbitrary L([x]). That is, L([x]) in vector form is $0 \subseteq [B] +$ [A][x], with [A] and [B] being sign constant matrices, [x] an *n* vector, [A] an *n* by *m* matrix and [B] an *m* vector. The ith scalar equation of L([x]) (denoted $L_1([x])$) is of the form:

$$L_i([\mathbf{x}]) \equiv 0 \subseteq [b_i] + \sum_{j=1}^m [a_{ij}] [x_j].$$

For QKT1, \mathbf{x}^T is $(\mu^T \lambda^T)^T$, $[\mathbf{B}] = [\bigtriangledown f]$, and $[\mathbf{A}]$ is the matrix $(\bigtriangledown \mathbf{g} \bigtriangledown \mathbf{h})$. Additionally, we generalize the set of restrictions given by QKT3 (i.e., $[\mu_i] \neq \hat{-}$), to arbitrary sets of restrictions $\mathbf{R}([\mathbf{x}]) \subseteq \{[x_i] \neq s | x_i \in \mathbf{x}, s \in \{\hat{-}, 0, \hat{+}\}\}$. For the cylinder (table, end of QKT section), QKT1 has 5 $L_i([\mathbf{x}])$'s, with $\mathbf{x} \equiv (\mu_1 \mu_2 \mu_3 \mu_4 \lambda_1 \lambda_2)^T$. For ease of reading we wrote terms $\hat{+}[x_i]$ as $[x_i]$, $\hat{-}[x_i]$ as $-[x_i]$, and eliminated terms $0[x_i]$. The cylinder $\mathbf{R}([\mathbf{x}])$ is $\{[\mu_1] \neq \hat{-}, [\mu_2] \neq \hat{-}, [\mu_3] \neq \hat{-}, [\mu_4] \neq \hat{-}\}$.

For step 2, the prime assignments of each $L_i([\mathbf{x}])$ are constructed from three sets of scalar assignments, consistent with $\mathbf{R}([\mathbf{x}])$: those restricting one of the equation's terms $([a_{ij}][x_j])$ to be positive (P_i) , those making a term zero (Z_i) , and those making a term negative (N_i) , respectively:

 $P_{i} \equiv \{ [x_{j}] = [a_{ij}] \mid [a_{ij}] \neq 0, \ ([x_{j}] \neq [a_{ij}]) \notin \mathbf{R}([\mathbf{x}]) \}, \\ Z_{i} \equiv \{ [x_{j}] = 0 \mid [a_{ij}] \neq 0, \ ([x_{j}] \neq 0) \notin \mathbf{R}([\mathbf{x}]) \} \text{ and} \\ N_{i} \equiv \{ [x_{i}] = -[a_{ij}] \mid [a_{ij}] \neq 0, \ ([x_{j}] \neq -[a_{ij}]) \notin \mathbf{R}([\mathbf{x}]) \}.$

Justifying P_i , for example, we know in general that $[c] \neq 0 \rightarrow [c]^2 = \hat{+}$. Thus $[a_{ij}][x_j] = \hat{+}$ if $[x_j] = [a_{ij}]$ and $[a_{ij}] \neq 0$. The derivation of Z_i and N_i is similar. Constructing the prime assignments for the cylinder $L_i([\mathbf{x}])$ uses:

1	Ni	Z_1	Pi
1	$[\lambda_1] = +, [\lambda_2] = +$	$[\lambda_1] = 0, [\lambda_2] = 0$	$[\lambda_1] = -, [\lambda_2] = -$
2	$[\lambda_1] = -, [\mu_2] = +$	$[\lambda_1] = 0, [\mu_2] = 0$	$[\lambda_1] = \dot{+}$
3	$[\lambda_2] = -, [\mu_1] = +$	$\lambda_2 = 0, \mu_1 = 0$	$[\lambda_2] = \dot{+}$
4	$[\lambda_1] = -$	$\lambda_1 = 0, \mu_4 = 0$	$[\lambda_1] = \dot{+}, [\mu_4] = \dot{+}$
5	$[\lambda_1] = \dot{+}, [\lambda_2] = \dot{+}$	$\lambda_1 = 0, \lambda_2 = 0,$	$[\lambda_1] = -, [\lambda_2] = -,$
		$\mu_3 = 0$	$[\mu_3] = \dot{+}$

Next, recall that the prime (implicating) assignments for $L_i([\mathbf{x}])$ must imply $L_i([\mathbf{x}])$. That is, they guarantee that it holds, given $\mathbf{R}([\mathbf{x}])$, independent of

additional consistent assignments. This is true if the right hand side of $L_i([\mathbf{x}])$ is guaranteed to be a superset of 0 (i.e., it is either 0 or ?). The form of the assignments that achieve this for some $L_i([\mathbf{x}])$ depends on the value of $[b_i]$, where $[b_i] = \left[\frac{\partial f}{\partial x_i}\right]$ for QKT1. Suppose $[b_i] = \hat{+}$, then the right hand side must become ?. This holds exactly when at least one of the $[a_{ij}][x_j]$ terms is negative (since $0 \subseteq (\hat{-}) + (\hat{+}) = \hat{?}$). For example, in the cylinder QKT equation (2), $\lambda_1 = \hat{-}$ guarantees that the equation is satisfied. The only other assignment that guarantees this is $\mu_2 = \hat{+}$. Thus the prime assignments for (2) are $\{\lambda_1 = \hat{-}\}$ and $\{\mu_2 = \hat{+}\}$. The treatment of $[b_i] = \hat{-}$ is analogous.

Next, suppose $[b_i] = 0$, then to imply $L_i([\mathbf{x}])$ the prime assignment can make the right hand side either 0 or $\hat{?}$. The first holds exactly when all terms are 0. The second holds when at least one term is positive and the other is negative. For example, $\left[\frac{\partial f(\mathbf{X})}{\partial x_i}\right] = 0$ in cylinder QKT1(3): $0 \subseteq -[\mu_1] + [\lambda_2]$. Thus, the prime assignments are $\{\lambda_2 = 0, \mu_1 = 0\}$ and $\{\lambda_2 = \hat{+}, \mu_1 = \hat{+}\}$. Note that $\{\lambda_2 = \hat{-}, \mu_1 = \hat{-}\}$ is not acceptable, since by restriction $[\mu_i] \neq \hat{-}$. To summarize, the prime assignments of $L_i([\mathbf{X}])$ are 1) N_i if $[b_i] = \hat{+}, 2$) P_i if $[b_i] = \hat{-}, and 3$ $\{Z_i\} \cup \{\{p, n\}| p \in P_i, n \in N_i\}$ if $[b_i] = 0$ (where p and n in $\{p, n\}$ do not contradict each other). Completing step two for the table of cylinder equations QKT1(1) - (5) produces:

$$\{ \lambda_1 = \hat{+} \}, \{ \lambda_2 = \hat{+} \}$$
 P(1)

$$\{ \lambda_1 = \hat{-} \}, \{ \mu_2 = \hat{+} \}$$
 P(2)

$$\{ \lambda_2 = 0, \mu_1 = 0 \} \{ \lambda_2 = \hat{+}, \mu_1 = \hat{+} \}$$
 P(3)

$$\{ \lambda_1 = 0, \mu_4 = 0 \}, \{ \lambda_1 = \hat{-}, \mu_4 = \hat{+} \}$$
 P(4)

$$\{ \lambda_1 = 0, \lambda_2 = 0, \mu_3 = 0 \},$$

$$\{ \lambda_1 = \hat{+}, \lambda_2 = \hat{-}, \} \{ \lambda_1 = \hat{+}, \mu_3 = \hat{+} \}$$

$$\{ \lambda_1 = \hat{-}, \lambda_2 = \hat{+} \}, \{ \lambda_2 = \hat{+}, \mu_3 = \hat{+} \}$$
 P(5)

The third step, constructing the composite primes for L([x]), is based on:

$$\bigvee_{p \in P(\mathbf{L}([\mathbf{x}]))} \left(\bigwedge_{a \in p} a \right) \equiv \bigwedge_{i=1}^{n} \left(\bigvee_{p \in P(\mathbf{L}_{i}([\mathbf{x}]))} \left(\bigwedge_{a \in p} a \right) \right)$$

The left hand side is a disjunction of the $L([\mathbf{x}])$ prime assignments, and the right hand side is an expression in terms of the primes of $L_i([\mathbf{x}])$, just computed. Thus, the desired primes result from reducing the expression on the right to minimal, disjunctive normal form. For this specialized case, this step is equivalent to computing minimal set covering of the $P(L_i([\mathbf{x}]))$ and then removing inconsistent assignments (see a standard algorithm text, or (Williams 1994) for our algorithm). For the cylinder, the minimal covering of P(1) - (5) produces just two prime assignments,

$$\{ \{ [\lambda_1] = \hat{-}, [\lambda_2] = \hat{+}, [\mu_1] = \hat{+}, [\mu_4] = \hat{+} \}, \\ \{ [\lambda_1] = 0, [\lambda_2] = \hat{+}, [\mu_1] = \hat{+}, [\mu_2] = \hat{+}, [\mu_3] = \hat{+}, [\mu_4] = 0 \} \}.$$

The fourth step, extracting the minimal sets of $[\mu_i] = \hat{+}$ assignments results in $\{[\mu_1] = \hat{+}, [\mu_4] = \hat{+}\}$ and $\{[\mu_1] = \hat{+}, [\mu_2] = \hat{+}, [\mu_3] = \hat{+}\}$. The fifth step uses $[\mu_i] = \hat{+} \rightarrow g_i(\mathbf{x}) = 0$ to map these sets to the equivalent minimal pstationary covering. The sets tell us that g_1 and g_4 must be active, or g_1 , g_2 and g_3 . The resulting cover is:

$$egin{aligned} \mathcal{F}_1 &\equiv \langle \{g_2,g_3\}, \{h_1,h_2,g_1,g_4\}
angle \ \mathcal{F}_2 &\equiv \langle \{g_4\}, \{h_1,h_2,g_1,g_2,g_3\}
angle, \end{aligned}$$
 and

where (\mathbf{g}, \mathbf{h}) is a space defined by inequality \mathbf{g} and equality \mathbf{h} constraints. \mathcal{F}_1 and \mathcal{F}_2 denote the line and point highlighted in the introduction to the cylinder example. The final step, formulating a new optimization problem, produces:

Given: $S \equiv \{\mathbf{x} * \mid \mathbf{x} * = \arg\min_{\mathbf{x} \in \mathcal{F}} f(\mathbf{x}), \mathcal{F} \in \{\mathcal{F}_1, \mathcal{F}_2\}\},$ Find: $\min_{\mathbf{x} \in \mathcal{S}} f(\mathbf{x}).$

The first part finds the minimum of each subspace in the covering. The second part selects from these the global minimum. A more expanded form was given in the introduction. Thus through this example we have demonstrated activity analysis' capability of partially solving constrained optimization problems from monotonicity constraints, and for synthesizing special purpose optimization codes.

Discussion

As we mentioned in the introduction, activity analysis builds upon a large body of work from the mechanical engineering community on monotonicity analysis(Wilde 1975; Papalambros & Wilde 1979; Papalambros 1982), a method that uses derivative information to address the boundedness and global optimality of optimization problems. Monotonicity analysis provides two rules that test the boundedness of a formulation:

Rule 1: If the objective function is monotonic with respect to a variable, then there exists at least one active constraint that bounds the variable in the direction opposite of the objective function.

Rule 2: If a variable is not contained in the objective function then it must be either bounded from both above and below by active constraints or not actively bounded at all (i.e., in the latter case any constraint that is monotonic with respect to that variable must be inactive or irrelevant).

Both of these rules can be derived from the Kuhn-Tucker Conditions. They also follow as an instance of QKT and are embodied within activity analysis.

The result of monotonicity analysis (exhaustive application of the rules) are several sets of constraints one of which must be active for a problem to be well bounded. Various levels of rule-based implementations of monotonicity analysis have been described in (Michelena & Agogino 1988; Rao & Papalambros 1987; Azram & Papalambros 1984; Hansen, Jaumard, & Lu 1989), which guide numerical optimization codes. Choy and Agogino (Choy & Agogino 1986) and Agogino and Almgren (Agogino & Almgren 1987) incorporate symbolic algebraic methods to aid in the evaluation of monotonicities and the solution of the optima. Cagan and Agogino (Cagan & Agogino 1987) apply monotonicity analysis to identify topological changes to designs that improve performance. While these systems address the optimal reasoning problem, they do not present algorithms proven to be sound and complete (each of these implementations has been described as "heuristic" (Rao & Papalambros 1987; Hansen, Jaumard, & Lu 1989)).

Activity analysis provides the following contributions: it formalizes the strategic way in which a modeler focuses optimization, as the process of generating minimal pstationary coverings. It introduces QKT as a powerful condition for quickly eliminating large, suboptimal subspaces. Finally, it exploits this condition through a novel problem reformulation based on the prime, implicating assignments of linear sign equations. The activity analysis algorithm is sound and complete with respect to classifying the design space into pstationary and pnonstationary subspaces. The method of pruning suboptimal subspaces provides a continuous analog to the conflict-based approaches prevalent in combinatorial satisficing search (such as those used in model-based diagnosis (de Kleer & Williams 1987)). Activity analysis automates the intuitions about monotonicity exploited by the simplex method to examine only the vertices of the linear feasible space, most importantly, extending its application to nonlinear problems.

Activity analysis has been demonstrated on several engineering problems. The implementation is in Franz Lisp running on a Sparc 2. The problem reformulation is passed to Matlab's Optimization toolbox, where a wide variety of nonlinear gradient methods are available. (Williams 1994) describes an extension to activity analysis for cases where monotonicities are only partially known. Activity analysis is currently being pursued in the context of visual 3D matching problems and other embedded, realtime problems. Activity analysis can also be extended to provide explainable optimizers, ones that use QKT to provide commonsense explanations about optimality. Activity analysis is one of several techniques being developed that capture a modeler's expertise at strategically guiding numerical codes.

References

Agogino, A. M., and Almgren, A. S. 1987. Techniques for Integrating Qualitative Reasoning and Symbolic Computation in Engineering Optimization. *Engineering Optimization* 12:117-135.

Azram, S., and Papalambros, P. 1984. An Automated Procedure for Local Monotonicity Analysis. Trans. ASME, Journal of Mechanisms, Transmissions, and Automation in Design 106:82-89. Cagan, J., and Agogino, A. M. 1987. Innovative Design of Mechanical Structures from First Principles. *AI EDAM* 1(3):169-189.

Choy, J. K., and Agogino, A. M. 1986. SYMON: Automated SYMbolic MONotonicity Analysis System for Qualitative Design Optimization. In *Proceedings* of ASME 1986 International Computers in Engineering Conference, 305-310.

de Kleer, J., and Williams, B. C. 1987. Diagnosing Multiple Faults. Artif. Intell. 32:97-130.

Hansen, P.; Jaumard, B.; and Lu, S. H. 1989. An Automated Proceedure for Globally Optimal Design. Trans. of the ASME, Journal of Mechanisms, Transmissions, and Automation in Design 361-367.

Kuhn, H. W., and Tucker, A. W. 1951. Nonlinear Programming. In Neyman, J., ed., *Proceedings of the* Second Berkeley Symposium on Mathematical Statistics and Probability. Berkeley, CA: University of California Press.

Michelena, N., and Agogino, A. M. 1988. Multiobjective Hydraulic Cylinder Design. Journal of Mechanisms, Transmission and Automation in Design 110:81-87.

Papalambros, P., and Wilde, D. J. 1979. Global Non-Iterative Design Optimization Using Monotonicity Analysis. Trans. ASME, Journal of Mechanical Design 101(4):645-649.

Papalambros, P. 1982. Monotonicity in Goal and Geometric Programming. Transactions of the ASME, Journal of Mechanical Design 104:108-113.

Rao, J. R., and Papalambros, P. 1987. Implementation of Semi-Heuristic Reasoning for Bounded Analysis of Design Optimization Models. In Advances in Design Automation, proceedings of the ASME Design Automation Conference, 59-65.

Vanderplaats, G. N. 1984. Numerical Optimization Techniques for Engineering Design With Applications. New York: McGraw-Hill.

Wilde, D. J. 1975. Monotonicity and Dominance in Optimal Hydraulic Cylinder Design. Trans of the ASME, Journal of Engineering for Industry 94(4):1390-1394.

Williams, B. C., and Raiman, O. 1994. Decompositional Modelling through Caricatural Reasoning. In AAAI.

Williams, B. C. 1991. A theory of interactions: unifying qualitative and quantitative algebraic reasoning. *Artif. Intell.* 51.

Williams, B. C. 1994. Characterizing Activity Analysis. in progress.