

Optimization and homotopy methods for the Gibbs free energy of magmatic mixtures

A.Cassioli*, L.Consolini† M.Locatelli‡ A.Longo‡

Abstract

In this paper we consider a mathematical model for magmatic mixtures based on the Gibbs free energy. Different reformulations of the problem are presented and some theoretical results about the existence and number of solutions are derived. Finally, two homotopy methods and a global optimization one are introduced and computationally tested. One of the homotopy methods returns a single solution of the problem, while the other is able to return optimal solutions (often all of them). The global optimization method is a branch-and-reduce one with a theoretical guarantee of detecting *all* the solutions, although numerical difficulties may result in a loss of a few of them.

Keywords: Magmatic mixtures, Gibbs free energy, homotopy methods, branch-and-reduce.

1 Introduction

The ability of magma to move in underground magmatic reservoirs, to open and flow along fractures in rock, to rise upwards to the Earth's surface and give eruptions changes drastically with changes in magma physical properties, which are mainly controlled by crystallization and magmatic volatiles saturation [3, 24]. Volcanic phenomena on the Earth's surface range from effusive eruptions with quiet lava flows, low energy Strombolian explosions or Hawaiian lava fountains, to highly explosive Plinian events. Eruptions may represent from low to extremely high hazard to the human community, so that scientists are urged to make reliable eruptive predictions [18, 26]. Short-time hazard evaluation is critically based on recognition and modeling of eruptive precursors, such as volcano displacements, gravity anomalies and degassing, and of their sources, that are caused or tightly related to the subsurface motion of magma [18].

The physical-mathematical modeling of intrinsically complex natural systems such volcanoes requires advanced and performing numerical solutions that avoid oversimplifications, but still need reasonable computing efforts [31]. Magma is a multiphase and multicomponent mixture of silicatic melt and volatiles (H_2O , CO_2 , SO_2 , HCl mainly) subject to growth and resorption of crystals from the oxides in the melt, exsolution and dissolution of volatiles, depending on the

*DSI - Università degli Studi di Firenze

†Università degli Studi di Parma

‡INGV - Pisa

local conditions of pressure, temperature, and composition. Crystals, dissolved volatiles and gas bubbles control magma physical properties, such as density and viscosity, which in turn influence magma fluid-dynamics. Magma dynamics implies variations in the local conditions of pressure, temperature and composition, with consequent changes in the chemical equilibrium of solid, liquid and gaseous phases. The physical and numerical modeling of magma motion and chemical evolution thus requires the solution of the system of conservation equations of mass, momentum and energy of the magmatic mixture, closed with the constitutive relations for its physical properties, and the chemical laws for phase changes and reactions [7].

In this work we consider a mathematical model for magmatic mixtures based on the Gibbs free energy and we will propose several formulations and solution methods for its minimization. The paper is structured as follows. In Section 2 we derive the mathematical model of the problem; in Section 3 we will propose different equivalent reformulations of our problem as a nonlinear system, a fixed point problem, a complementarity problem, a global optimization problem; in Section 4 we discuss the simple two-dimensional case; in Section 5 we derive some theoretical results about the existence and number of solutions; in Section 6 we discuss homotopy methods for the detection of one or more solutions of the problem; in Section 7 we discuss a branch-and-reduce method for the detection of all the solutions; finally, in Section 8 we present some computational experiments.

2 Mathematical model of the problem

The equilibrium state of multiphase multicomponent mixtures such as the magmatic mixture, with melt, minerals, and volatiles is described by the classical thermodynamics theory of mixtures by means of the Gibbs free energy in terms of local conditions of pressure, temperature and the phase distribution of components [1]. It is applied to petrologic and geochemical reactions of solidification/readsorption of crystals, and exsolution/dissolution of gas bubbles [2, 9, 20, 21]. The mixture Gibbs free energy of the magmatic mixture (silicatic melt, crystal assemblages and gas phase) span from ideal to strongly non ideal with complex dependences on temperature, pressure, and composition. The number of independent components ranges from 2 to 15, depending on the crystals that form.

The equilibrium state of magma at a given pressure, temperature and total composition corresponds to the surface of minima of the Gibbs free energy with respect to the distribution of components among phases, or, in other words, to the locus of points of equality between the chemical potentials of components among phases [1].

Classical petrologic and geochemical studies follow two progressive steps. The Gibbs free energy of magmatic mixtures is written in terms of theoretical thermodynamic mixture models with some interaction parameters, that are computed from least squares fitting with the experimental data [9]. Substituting the best parameters in the mixture Gibbs free energy expression, the minimum surface can be computed and the equilibrium magma assemblages can be found for pressure, temperature and composition within the intervals of the fitting, typically quite narrow. Both the fitting and the minimization of the Gibbs

free energy are challenging tasks. However, they allow to derive semi-empirical phase diagrams for crystals and the melt representing the shape of the Gibbs free energy surface. Leaving out volatiles and experimental petrology that are fully volcanological subjects themselves, this work develops a method to find the equilibrium composition between solid crystals and the silicate liquid for natural silicate melts.

2.1 Numerical Modeling and Simulation Feasibility

Numerical codes aimed at volcanic hazard solve the complex non-linear conservation equations of the magmatic mixture, compute its equilibrium composition and physical properties. Hazard evaluation needs as much as possible realistic and reliable numerical solutions to the underground magma dynamics and compositional evolution. The efficient accurate computation of magma equilibrium state is thus the necessary condition to develop a fast light internal algorithms.

The computation of magmatic equilibrium assemblages, however, implies a huge increase in computing expenses. As an example, in finite element methods it should be computed at least once in each mesh element for each time step. Otherwise lookup tables may be predefined and interpolated during simulations, or changes in the equilibrium assemblages can be neglected if the thermodynamic independent local variables (pressure, temperature and composition) vary below a certain threshold.

In the thermodynamics theory of magmatic mixtures the Gibbs free energy surface is characterized by one global minimum and several multiple local minima. The global minimum represents the stable equilibrium liquid-solid assemblage. The physical meaning of each local minimum is that it represents the locus of metastable equilibria for a certain range of temperatures or liquid composition. Changing temperature and/or amount of components, the liquid-solid assemblage moves from one minimum to another. It may finally reach the global minimum and stably rest there [1].

The algorithms proposed in this work find multiple global minimum solutions, all to be retained. The more reliable solution should be chosen at simulation run time implicitly based on the petrological semi-empirical phase diagrams described in Section 2.

2.2 Gibbs Free Energy of Magmatic Mixtures

The mathematical model for the problem developed in this paper is based on the expression for the Gibbs free energy of the multicomponent multiphase (solid, liquid and gas) magmatic mixtures. It is derived following the approach in [13], and the solid-liquid equilibrium is the only reaction considered.

The solid phase has n components represented by the number of moles $(c_1^{sol}, c_2^{sol}, \dots, c_n^{sol})$. The liquid phase has p components represented by the number of moles $(c_1^{liq}, c_2^{liq}, \dots, c_p^{liq})$. Dissolution of c_i^{sol} moles of a solid i into the liquid corresponds to the generic stoichiometric mass balance

$$c_i^{sol} = \sum_{j=1}^p \nu_{ij} c_j^{liq} \quad i = 1, \dots, n \quad (1)$$

where the sum is extended over all the liquid components $j = 1, \dots, p$, and ν_{ij} are the stoichiometric coefficients of the dissolution reaction, taking null values if a liquid component does not form from the reaction.

Thermodynamics usually works with intensive variables, especially for multi-component multiphase mixtures, using the chemical potential μ_k of a component k , defined as

$$\mu_k = \left. \frac{\partial G}{\partial c_k} \right|_{p, T, j \neq k} \quad (2)$$

where p is pressure, and T is temperature. In terms of the chemical potential at standard state μ_k^0 and the activity a_k , μ_k is

$$\mu_k = \mu_k^0 + RT \ln a_k, \quad (3)$$

where R is the gas constant.

For a generic reaction, the extensive change in Gibbs free energy is

$$\Delta G^{reaction} = \Delta G^{products} - \Delta G^{reactants}$$

Writing in terms of the intensive variables μ_k , the change is the difference between the Gibbs free energy of the dissolved liquid parcel and of the solid i before dissolution:

$$-\mathcal{D}G_i^{reaction} = \sum_{j=1}^p \nu_{ij} \mu_j^{liq} - \mu_i^{sol} \quad i = 1, \dots, n \quad (4)$$

Following [13], $\mathcal{D}G_i^{reaction}$ can be replaced by the chemical affinity A , the electronic property by which dissimilar chemical species are capable of forming chemical compounds. It can be seen as an energetic measure of the extent of nearness to equilibrium of the solid in the liquid from under- or super-saturation or saturation itself (respectively, $|\mathcal{D}G_i^{reaction}| \gg 0$ and $\mathcal{D}G_i^{reaction} = 0$), yielding the form:

$$\sum_{j=1}^p \nu_{ij} \mu_j^{liq} - \mu_i^{sol} + A = 0 \quad , \quad i = 1, \dots, n \quad (5)$$

by expressing the chemical potential μ_i^{sol} of the solid component i according to (3).

The activity a_i^{sol} accounts for the deviations from the ideal behavior of a mixture (regular, isometric, non-isometric, and other extremely complex cases), i.e., the complexities of the solid-liquid mixture formed from the dissolution of the solid, or the mixture of more solids and their coexisting equilibrium liquid. In this work, temperature T and pressure P are fixed, and a symmetric regular solid solution among the $i = 1, \dots, n$ solid components is assumed (see (9)). Substitution in (5) yields

$$\sum_{j=1}^p \nu_{ij} \mu_j^{liq} - \mu_i^{0, sol} - RT \ln a_i^{sol} + A = 0 \quad , \quad i = 1, \dots, n \quad (6)$$

The latter can be simplified using for convenience the quantity $\Delta\mu_i^{sol}$, defined as the difference between the chemical potential of the solid at standard state and the Gibbs free energy of the liquid:

$$\Delta\mu_i^{sol} = - \sum_{j=1}^p \nu_{ij} \mu_j^{liq} + \mu_i^{0,sol} \quad , \quad i = 1, \dots, n \quad (7)$$

which yields the formulation

$$0 = \Delta\mu_i^{sol} + RT \ln a_i^{sol} - A \quad , \quad i = 1, \dots, n \quad (8)$$

Thus, the resulting system of n equations contains $n + 1$ unknowns, that is A and a_i^{sol} , $i = 1, \dots, n$. In order to obtain the same number of equations and variables we explicit the expressions for a_i^{sol} , $i = 1, \dots, n$.

Each solid component i in equilibrium with the liquid phase of the magmatic mixture is assumed to behave as a strictly symmetric regular solution, so that its affinity a_i^{sol} in terms of the molar fractions x_i , $i = 1, \dots, n$ is [1]:

$$\left\{ \begin{array}{l} a_i^{sol} = x_i \exp \left\{ \frac{1}{RT} \left[\sum_{k=1}^n x_k W_{ki} - \sum_{h,k=1, h < k}^n x_h x_k W_{hk} \right] \right\} \quad i = 1, \dots, n \\ \sum_{i=1}^n x_i = 1 \\ x_i \geq 0 \quad \quad \quad i = 1 \dots n \end{array} \right. \quad (9)$$

Note that molar fractions must be non-negative and their sum must be equal to one. Coefficients W_{kh} represent binary interaction parameters for solid components (which are temperature and pressure independent) and such that $W_{kk} = 0$, $W_{kh} = W_{hk}$, $h, k = 1, \dots, n$, and can be grouped in a symmetric square matrix $\mathbf{W} \in \mathbb{R}^{n \times n}$. Substituting (9) in (8), we obtain the following system of nonlinear equations

$$\left\{ \begin{array}{l} 0 = \Delta\mu_i^{sol} - A + RT \ln x_i + \sum_{k=1}^n x_k W_{ik} - \sum_{h,k=1, h < k}^n x_h x_k W_{hk} \quad i = 1, \dots, n \\ \sum_{i=1}^n x_i = 1 \\ x_i \geq 0 \quad \quad \quad i = 1 \dots n \end{array} \right. \quad (10)$$

We can eliminate variable A and the bilinear terms $x_h x_k$ by subtracting pairs

of different equations, thus ending up with the following equivalent system

$$\left\{ \begin{array}{l} RT \ln \left(\frac{x_i}{x_j} \right) + \sum_{k=1}^n (W_{ik} - W_{jk}) x_k + \Delta\mu_i^{sol} - \Delta\mu_j^{sol} = 0 \quad i, j = 1, \dots, n, i \neq j \\ \sum_{i=1}^n x_i = 1 \\ x_i \geq 0 \end{array} \right. \quad i = 1 \dots n \quad (11)$$

2.3 Related Works

Several simple algorithms to find the equilibrium state of reduced problems exist, but they are limited to few sets of minerals. The most complex software MELTS was developed by M. S. Ghiorso and collaborators [2, 11, 12]. A simpler version is by A.E. Boudreau [4]. Both programs are intended to run with a user-friendly windowing system, no source files are released, and they can not be linked to any efficient fluid-dynamics simulator. Smith and Azimoth [25] provided the program `adiabat_1ph`, a simple text-menu driver for subroutine versions of the algorithms by M. S. Ghiorso and collaborators. This software runs under Linux, by means of executable (and not source) files, called by Perl scripts, eventually in batch mode. This kind of program could be called during fluid-dynamic simulations, but at the expenses of the huge loss of computing time.

In the field of global optimization, the solution of the minimization of energy potential functions has been a very challenging problem since the very early days of Mathematical Programming (see [6]). Many references on the use of global optimization techniques for these problems can be found for instance in [10]. A special case of convex formulation of the Gibbs Energy Potential in the context of geochemical science has been exploited in [17]; a heuristic approach has been used in [27], where tabu-search and differential evolution algorithms have been compared.

For what concerns the use of homotopy continuation methods to solve equilibrium problems arising in chemistry, we can cite, e.g., [16, 28].

3 Different problem reformulations

In this section we show different possible reformulations for our problem. Each of them will be useful for the subsequent development. As already remarked, problem (11) is in the form of a *nonlinear system*. We slightly modify it by: (i) setting, without loss of generality $RT = 1$; (ii) setting, for the sake of compactness, $c_i = \Delta\mu_i^{sol}$; (iii) parameterizing it as follows

$$\begin{cases} \ln(x_i) - \ln(x_j) + \lambda[(\mathbf{W}_i - \mathbf{W}_j)^T \mathbf{x} + c_i - c_j] = 0 & i, j \in \{1, \dots, n\}, i \neq j \\ \mathbf{e}^T \mathbf{x} - 1 = 0 \\ \mathbf{x} \geq \mathbf{0} \end{cases} \quad (12)$$

where:

- $\lambda \geq 0$ is a parameter;
- \mathbf{W}_i denotes the i -th row of matrix \mathbf{W} ;
- \mathbf{e} denotes the n -dimensional vector whose components are all equal to 1.

Of course, we are only interested in the solution for $\lambda = 1$ but, as we will see, it will turn out to be useful to consider the parameterized problem. Note that the unique solution for $\lambda = 0$ is $\frac{1}{n}\mathbf{e}$. It is easy to see that many equations in (12) are redundant and the system can be rewritten as follows

$$\begin{cases} F_i(\mathbf{x}; \lambda) := \ln(x_i) - \ln(x_n) + \lambda[(\mathbf{W}_i - \mathbf{W}_n)^T \mathbf{x} + c_i - c_n] = 0 & i = 1, \dots, n-1 \\ F_n(\mathbf{x}; \lambda) := \mathbf{e}^T \mathbf{x} - 1 = 0 \\ \mathbf{x} \geq \mathbf{0} \end{cases} \quad (13)$$

By introducing the n -dimensional vector of functions $\mathbf{F}(\mathbf{x}; \lambda)$, whose components are defined in (13), we can also present the problem in the more compact form

$$\mathbf{F}(\mathbf{x}; \lambda) = \mathbf{0}, \quad \mathbf{x} \geq \mathbf{0}.$$

The nonlinear system can also be written in the equivalent form

$$\begin{cases} G_i(\mathbf{x}, \lambda) := x_i - x_n \exp\{-\lambda[(\mathbf{W}_i - \mathbf{W}_n)^T \mathbf{x} + c_i - c_n]\} = 0 & i = 1, \dots, n-1 \\ G_n(\mathbf{x}, \lambda) := \mathbf{e}^T \mathbf{x} - 1 = 0 \\ \mathbf{x} \geq \mathbf{0} \end{cases} \quad (14)$$

or, in compact form

$$\mathbf{G}(\mathbf{x}; \lambda) = \mathbf{0}, \quad \mathbf{x} \geq \mathbf{0}. \quad (15)$$

A further reformulation is the following.

$$\begin{cases} x_i = x_n \exp\{-\lambda[(\mathbf{W}_i - \mathbf{W}_n)^T \mathbf{x} + c_i - c_n]\} & i = 1, \dots, n-1 \\ x_n = \sum_{i=1}^{n-1} x_i - 1 \\ \mathbf{x} \geq \mathbf{0} \end{cases}$$

After setting

$$\begin{cases} H_i(\mathbf{x}; \lambda) := x_n \exp\{-\lambda[(\mathbf{W}_i - \mathbf{W}_n)^T \mathbf{x} + c_i - c_n]\} & i = 1, \dots, n-1 \\ H_n(\mathbf{x}; \lambda) := \sum_{i=1}^{n-1} x_i - 1 \end{cases}$$

we can rewrite the problem in the compact form

$$\mathbf{x} = H(\mathbf{x}; \lambda), \quad \mathbf{x} \geq \mathbf{0}, \quad (16)$$

which is a *Fixed Point Problem* (FPP in what follows).

Next, problem (15) can also be viewed as a (parametric) *Nonlinear Complementarity Problem* (NCP in what follows)

$$\mathbf{x}^T \mathbf{G}(\mathbf{x}; \lambda) = 0, \quad \mathbf{x}, \mathbf{G}(\mathbf{x}; \lambda) \geq \mathbf{0}. \quad (17)$$

Indeed, the following observation holds.

Observation 3.1 *A point $\mathbf{x}^*(\lambda)$ solves (15) if and only if it also solves (17).*

Proof. The fact that a solution of (15) is also a solution of (17) is trivial. But also the vice versa holds true. Indeed, we notice that for finite λ values

$$x_i^*(\lambda) = 0 \quad \text{for some } i \in \{1, \dots, n\}, \quad \mathbf{x}^*(\lambda) \geq \mathbf{0} \quad \Rightarrow \quad \mathbf{x}^*(\lambda) = \mathbf{0}$$

which violates the inequality $\mathbf{e}^T \mathbf{x} \geq 1$. Therefore, all solutions of (17) satisfy $\mathbf{x}^*(\lambda) > \mathbf{0}$ and, consequently, they are also solutions of (15). \square

The following observation, which will turn out to be useful in what follows, states that for $\lambda = \infty$, the problem reduces to a Linear Complementarity one (LCP in what follows).

Observation 3.2 *All limit solutions of problem (12) as $\lambda \rightarrow \infty$ can be obtained by solving a LCP.*

Proof. Let us denote by $\mathbf{x}^*(\lambda)$ a solution of (12) for some value λ and let \mathbf{x}^* be the limit of such solution as $\lambda \rightarrow \infty$. In order to detect all possible limit solutions \mathbf{x}^* at infinity, we can proceed as follows. Let us fix a subset $T \subseteq \{1, \dots, n\}$ and set

$$x_i^* = 0 \quad \forall i \in T, \quad x_i^* > 0 \quad \forall i \notin T.$$

Now, we notice that dividing by λ the equation

$$\ln(x_i^*(\lambda)) - \ln(x_j^*(\lambda)) + \lambda[(\mathbf{W}_i - \mathbf{W}_j)^T \mathbf{x}^*(\lambda) + c_i - c_j] = 0$$

for $i, j \notin T$, we end up with

$$\frac{\ln(x_i^*(\lambda))}{\lambda} - \frac{\ln(x_j^*(\lambda))}{\lambda} + (\mathbf{W}_i - \mathbf{W}_j) \mathbf{x}^*(\lambda) + c_i - c_j = 0$$

Since $\ln(x_i^*), \ln(x_j^*) > -\infty$, we will have that

$$\frac{\ln(x_i^*(\lambda))}{\lambda}, \frac{\ln(x_j^*(\lambda))}{\lambda} \rightarrow 0 \quad \text{as } \lambda \rightarrow \infty.$$

Therefore, for $i, j \notin T$ we must have

$$(\mathbf{W}_i - \mathbf{W}_j) \mathbf{x}^* + c_i - c_j = 0$$

For $i \in T$ e $j \notin T$, we observe that

$$\frac{\ln(x_i^*(\lambda)) - \ln(x_j^*(\lambda))}{\lambda} < 0 \quad \forall \lambda \text{ large enough} \quad (18)$$

Since we must have

$$\frac{\ln(x_i^*(\lambda))}{\lambda} - \frac{\ln(x_j^*(\lambda))}{\lambda} + (\mathbf{W}_i - \mathbf{W}_j)\mathbf{x}^*(\lambda) + c_i - c_j = 0$$

from (18) we can conclude that

$$(\mathbf{W}_i - \mathbf{W}_j)\mathbf{x}^*(\lambda) + c_i - c_j > 0 \quad \forall \lambda \text{ large enough}$$

and, thus

$$(\mathbf{W}_i - \mathbf{W}_j)\mathbf{x}^* + c_i - c_j \geq 0$$

Then, the solution \mathbf{x}^* at infinity, related to subset T , is a feasible solution (if it exists!) of the following system of linear equalities and inequalities.

$$\begin{cases} x_i^* = 0 & i \in T \\ (\mathbf{W}_i - \mathbf{W}_j)\mathbf{x}^* + c_i - c_j = 0 & \forall i, j \notin T \\ (\mathbf{W}_i - \mathbf{W}_j)\mathbf{x}^* + c_i - c_j \geq 0 & \forall i \in T, j \notin T \\ \mathbf{e}^T \mathbf{x}^* = 1 \\ \mathbf{x}^* \geq \mathbf{0} \end{cases}$$

Then, the problem of identifying *all* the solutions at infinity can be reformulated as follows

$$\begin{cases} \mathbf{x}^*(\mathbf{W}\mathbf{x}^* + \mathbf{c} - z\mathbf{e}) = \mathbf{0} \\ \mathbf{W}\mathbf{x}^* + \mathbf{c} - z\mathbf{e} \geq \mathbf{0} \\ \mathbf{x}^* \geq \mathbf{0} \\ \mathbf{e}^T \mathbf{x}^* = 1 \end{cases} \quad (19)$$

where z is an added variable. This is not a LCP yet, but it can be easily transformed into the following LCP

$$\begin{cases} \mathbf{x}^*(\mathbf{W}\mathbf{x}^* + \mathbf{c} - (z_1 - z_2)\mathbf{e}) = \mathbf{0} \\ z_1(\mathbf{e}^T \mathbf{x}^* - 1) = 0 \\ z_2(\mathbf{e}^T \mathbf{x}^* - 1) = 0 \\ w(z_1 + z_2 - 1) = 0 \\ \mathbf{W}\mathbf{x}^* + \mathbf{c} - (z_1 - z_2)\mathbf{e} \geq \mathbf{0} \\ \mathbf{e}^T \mathbf{x}^* \geq 1 \\ z_1 + z_2 - 1 \geq 0 \\ \mathbf{x}^* \geq \mathbf{0} \\ z_1, z_2, w \geq 0 \end{cases}$$

□

We remark that solutions for $\lambda = \infty$ will be further characterized in Section 6.2.

As a final reformulation for our problem, we notice that, as usual for non-linear systems, it can also be viewed as a *Global Optimization* (GO in what follows) problem

$$\begin{aligned} \min_{\mathbf{x}} \quad & \mathbf{F}(\mathbf{x}; \lambda)^T \mathbf{F}(\mathbf{x}; \lambda) \quad (\text{or } \mathbf{G}(\mathbf{x}; \lambda)^T \mathbf{G}(\mathbf{x}; \lambda)) \\ & \mathbf{e}^T \mathbf{x} = 1 \\ & \mathbf{x} \geq \mathbf{0} \end{aligned} \tag{20}$$

4 The two-dimensional case

Before discussing the general n -dimensional case, we develop the case $n = 2$. As we will see, some results obtained for the two-dimensional case can be generalized to the n -dimensional one. By eliminating variable x_2 , the two-dimensional problem can be written as follows.

$$\begin{aligned} f(x_1; \lambda) &:= \ln(x_1) - \ln(1 - x_1) + \lambda W(1 - 2x_1) + \lambda(c_1 - c_2) = 0 \\ 0 &\leq x_1 \leq 1 \end{aligned}$$

W.l.o.g., we can assume that $c_1 \geq c_2$. In fact, we only deal with the case $c_1 > c_2$ (the case $c_1 = c_2$ can be dealt with in a completely analogous way). A simple analysis of function f shows that

- f is continuous;
-

$$\lim_{x_1 \rightarrow 0} f(x_1; \lambda) = -\infty, \quad \lim_{x_1 \rightarrow 1} f(x_1; \lambda) = +\infty;$$

- for $W \leq \frac{2}{\lambda}$, f is increasing; while for $W > \frac{2}{\lambda}$ the function has the following local maximizer and local minimizer

$$x_{max}(\lambda) = \frac{1 - \sqrt{1 - \frac{2}{\lambda W}}}{2} < \frac{1}{2}, \quad x_{min}(\lambda) = \frac{1 + \sqrt{1 - \frac{2}{\lambda W}}}{2} > \frac{1}{2};$$

- f is concave for $x \leq \frac{1}{2}$ and convex for $x > \frac{1}{2}$;
- it holds that

$$f(x_{max}(\lambda); \lambda) \geq f\left(\frac{1}{2}; \lambda\right) > 0.$$

Therefore, we notice that there always exists a solution $x_1'(\lambda) < x_{max}(\lambda)$ of the problem. Moreover, if

$$f(x_{min}(\lambda); \lambda) = 0 \tag{21}$$

then we have the further solution $x_1''(\lambda) = x_{min}(\lambda)$, while if

$$f(x_{min}(\lambda); \lambda) < 0$$

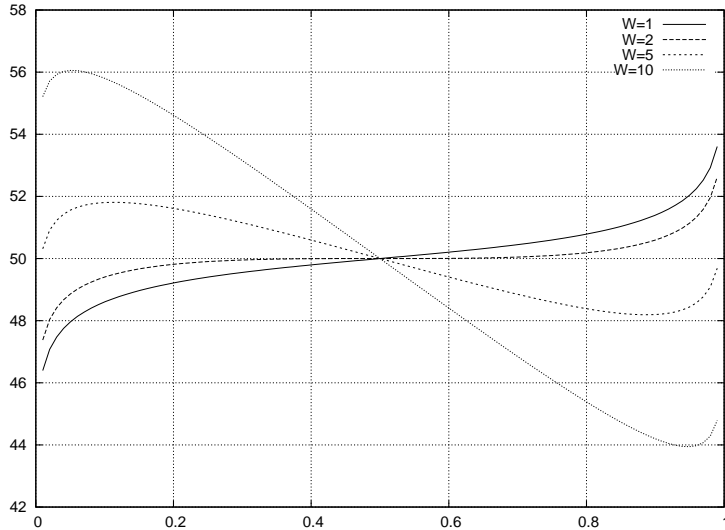


Figure 1: The energy function for the two-dimensional case with $\lambda = 1$ and different values of W .

we have two further solutions $\frac{1}{2} < x_1''(\lambda) < x_{min}(\lambda)$ and $x_1'''(\lambda) > x_{min}(\lambda)$. Therefore, we always have at least one solution and, except when λ satisfies (21), we have either one or three solutions. In fact, noticing that when λ satisfies (21) $x_{min}(\lambda)$ is a solution with multiplicity two, we can conclude that, by counting multiplicities, we always have either one or three solutions.

5 Existence and number of solutions

The question now is how the results for $n = 2$ can be generalized to larger n values. In particular, can we always guarantee the existence of at least one solution? Is it always true that the number of solutions (counting also multiplicities) is always odd? For what concerns the first question, in Section 3 we observed that our problem can be reformulated as the FPP (16). Function \mathbf{H} is defined over the n -dimensional unit simplex

$$\Sigma_n = \{\mathbf{x} \in \mathbb{R}^n : \mathbf{e}^T \mathbf{x} = 1, \mathbf{x} \geq \mathbf{0}\},$$

and returns values in Σ_n itself. We can exploit this formulation to prove existence of a solution.

Observation 5.1 *The FPP (16) always admits a solution.*

Proof. Existence of a solution is guaranteed by a variant of Brouwer's fixed point theorem [5], stating that every continuous function \mathbf{H} from a convex compact subset K of an Euclidean space to K itself ($K = \Sigma_n$ in this case) has a fixed point. \square

Later, in observation 6.5, we will prove that the number of solutions of (16) is odd for almost any choice of vector \mathbf{c} .

Finally, we observe that for any finite λ value, solutions of our problem always have strictly positive coordinate values.

Observation 5.2 *For each solution $\mathbf{x}^*(\lambda)$, $\lambda < \infty$, of our problem, it holds that*

$$x_i^*(\lambda) \geq \frac{e^{-\lambda \max_{\mathbf{x} \in \Sigma_n} [(\mathbf{W}_i - \mathbf{W}_n)^T \mathbf{x} + c_i - c_n]}}{n \sum_{j=1}^n e^{-\lambda \min_{\mathbf{x} \in \Sigma_n} [(\mathbf{W}_j - \mathbf{W}_n)^T \mathbf{x} + c_j - c_n]}} > 0.$$

Proof. In view of reformulation (16) we have that, for each $i = 1, \dots, n$

$$x_i^*(\lambda) = x_n^*(\lambda) e^{-\lambda [(\mathbf{W}_i - \mathbf{W}_n)^T \mathbf{x} + c_i - c_n]}.$$

Taking into account $\mathbf{e}^T \mathbf{x}^*(\lambda) = 1$, we have

$$x_n^*(\lambda) = \frac{1}{\sum_{j=1}^n e^{-\lambda [(\mathbf{W}_j - \mathbf{W}_n)^T \mathbf{x} + c_j - c_n]}}$$

and, consequently, also

$$x_i^*(\lambda) = \frac{e^{-\lambda [(\mathbf{W}_i - \mathbf{W}_n)^T \mathbf{x} + c_i - c_n]}}{\sum_{j=1}^n e^{-\lambda [(\mathbf{W}_j - \mathbf{W}_n)^T \mathbf{x} + c_j - c_n]}}$$

for $i = 1, \dots, n$. Then, the result immediately follows. \square

6 Homotopy methods

In this section we present continuation methods that allow finding solutions of (15). To this end, we first note that (15) can be rewritten as follows

$$\mathbf{K}(\mathbf{x}, \lambda) = \mathbf{x} - \mathbf{L}(\mathbf{x}, \lambda), \quad (22)$$

with $\mathbf{L}(\mathbf{x}, \lambda) = \mathbf{g}(\lambda(\mathbf{W}\mathbf{x} + \mathbf{c}))$ and

$$\mathbf{g}(\mathbf{x}) = \frac{e^{\mathbf{x}}}{\|e^{\mathbf{x}}\|_1}, \quad (23)$$

(here $e^{\mathbf{x}}$ denotes vector $(e^{x_1}, e^{x_2}, \dots, e^{x_n})$).

6.1 Forward homotopy from $\lambda = 0$ to $\lambda = 1$

We prove the following observation.

Observation 6.1 *There exists a continuous function $\{(\mathbf{x}(s), \lambda(s)), 0 \leq s \leq 1\}$ such that $\mathbf{K}(\mathbf{x}(s), \lambda(s)) = 0$.*

Proof. It follows from Leray-Schauder continuation theorem (Corollary 1 of [19]). In fact function $\mathbf{L} : \Sigma_n \times [0, 1] \rightarrow \Sigma_n$ is completely continuous and satisfies the following two conditions

- i) $\deg[\mathbf{x} - \mathbf{L}(\mathbf{x}; 0), \Sigma_n, 0] = 1$ (here \deg denotes the Leray-Schauder degree, see for instance [19]): this follows from the properties of Leray-Schauder degree, since $\mathbf{L}(\mathbf{x}; 0)$ is constant, being $\mathbf{L}(\mathbf{x}, 0) = \frac{\mathbf{e}}{n}, \forall \mathbf{x} \in \Sigma_n$;
- ii) the set $\Sigma_n \times [0, 1]$ is bounded. □

Observation 6.1 allows to use a continuation method to find a solution of $\mathbf{K}(\mathbf{x}, 1) = 0$. We first need a further observation.

Observation 6.2 *The Jacobian of \mathbf{g} is given by*

$$\mathbf{g}'(\mathbf{x}) = \text{diag}(\mathbf{g}(\mathbf{x}))(\mathbf{I} - \mathbf{e}\mathbf{g}(\mathbf{x})^T), \quad (24)$$

where $\text{diag}(\mathbf{g}(\mathbf{x}))$ denotes the diagonal matrix whose diagonal entries are equal to $g_i(\mathbf{x}), i = 1, \dots, n$.

Proof. Since $(e^{\mathbf{x}})' = \text{diag}(e^{\mathbf{x}})$, the derivative of (23) is given by

$$\mathbf{g}'(\mathbf{x}) = \frac{\text{diag}(e^{\mathbf{x}})}{\|e^{\mathbf{x}}\|_1} - \frac{e^{\mathbf{x}}}{\|e^{\mathbf{x}}\|_1^2} \mathbf{e}^T \text{diag}(e^{\mathbf{x}}) = \text{diag}(\mathbf{g}(\mathbf{x})) - \mathbf{g}(\mathbf{x})\mathbf{e}^T \text{diag}(\mathbf{g}(\mathbf{x})).$$

Then (24) follows using the identities $\mathbf{g}(\mathbf{x}) = \text{diag}(\mathbf{g}(\mathbf{x}))\mathbf{e}$ and $\mathbf{e}^T \text{diag}(\mathbf{g}(\mathbf{x})) = \mathbf{g}(\mathbf{x})^T$. □

To compute function $\gamma(s) = (\mathbf{x}(s), \lambda(s))$, it is convenient to require that 0 is a regular value of \mathbf{K} , according to the following definition.

Definition 6.1 *0 is a regular value of \mathbf{K} if the Jacobian of \mathbf{K} , i.e., $\mathbf{K}'(\mathbf{x}, \lambda) = [\partial_{\mathbf{x}}\mathbf{K}(\mathbf{x}, \lambda), \partial_{\lambda}\mathbf{K}(\mathbf{x}, \lambda)]$ is full rank for all $(\mathbf{x}, \lambda) \in \mathbf{K}^{-1}(0) \subset \Sigma_n \times [0, 1]$.*

To introduce the continuation method, following the notation of Section 2.1 of [8], for a full rank matrix $\mathbf{A} \in \mathbb{R}^{n \times (n+1)}$, we denote with $\mathbf{t}(\mathbf{A})$ the unique vector of \mathbb{R}^n such that

- 1) $\mathbf{A}\mathbf{t} = \mathbf{0}$,
- 2) $\|\mathbf{t}\| = 1$,
- 3) $\det \begin{pmatrix} \mathbf{A} \\ \mathbf{t} \end{pmatrix} > 0$.

Observation 6.3 *If 0 is a regular value of \mathbf{K} , then the solution of*

$$\begin{cases} \dot{\gamma}(s) = \mathbf{t}(\mathbf{K}'(\gamma(s))) \\ \gamma(0) = \frac{\mathbf{e}}{n} \end{cases} \quad (25)$$

satisfies $\mathbf{K}(\gamma(s)) = 0$ and intersects the set $\Sigma_n \times 1$. The intersection value $\gamma(\bar{s}) = \begin{pmatrix} \bar{\mathbf{x}} \\ 1 \end{pmatrix}$ satisfies $\mathbf{K}(\bar{\mathbf{x}}, 1) = 0$ and is therefore a solution of (15).

Ideally, curve $\gamma(s)$ can be computed by simply integrating (25). Numerically it is much more convenient to use an algorithm specifically designed for continuation procedures, such as one of the predictor-corrector methods.

One possible problem arises from the fact that 0 may not be a regular value of \mathbf{K} . This can be overcome by replacing homotopy (22) with

$$\mathbf{K}_2(\mathbf{x}, \lambda, \mathbf{p}) = \mathbf{x} - \mathbf{g}(\lambda(\mathbf{W}\mathbf{x} + \mathbf{c} + \mathbf{p})), \quad (26)$$

where $\mathbf{p} \in \mathbb{R}^n$ is a (small) perturbation term.

Observation 6.4 *For almost all $\mathbf{p} \in \mathbb{R}^n$, 0 is a regular value of \mathbf{K}_2 .*

Proof. By Observation 6.2

$$\partial_{\mathbf{p}}\mathbf{K}_2 = -\mathbf{g}'(\lambda(\mathbf{W}\mathbf{x} + \mathbf{c} + \mathbf{p}))\lambda$$

is full rank for any $\lambda \in (0, 1]$ and $\partial_{\mathbf{x}}\mathbf{K}_2$ is full rank for $\lambda = 0$ (here we mean full rank with respect to the dimension $n - 1$ of the codomain of function g , i.e., the unit simplex). Therefore the Jacobian of \mathbf{K}_2 is full rank for any $\lambda \in [0, 1]$ and, by Sard's Theorem (see 11.2.2 of [8]), 0 is a regular value of map $\mathbf{K}_2(\cdot, \cdot, \mathbf{p})$ for almost all perturbations $\mathbf{p} \in \mathbb{R}^n$. \square

The proof of the last observation leads also to the following observation.

Observation 6.5 *The number of solutions of (16) is odd for almost all \mathbf{c} .*

Proof. By the same reasoning of the previous proof, Sard's theorem implies that 0 is a regular value of map $\mathbf{K}(\mathbf{x}, 1)$ for almost all \mathbf{c} . By the properties of the degree, it follows that

$$\deg[\mathbf{K}(\cdot, 1), \Sigma_n, 0] = \sum_{x \in \mathbf{K}(\cdot, 1)^{-1}(0)} \text{sign}(\partial_x K(x, 1)),$$

which implies that set $\mathbf{K}(\cdot, 1)^{-1}(0)$ has a odd number of elements. \square

6.2 Backward homotopy from $\lambda = +\infty$ to $\lambda = 1$

To find multiple solutions of $\mathbf{K}(\mathbf{x}, 1) = 0$, we apply a similar continuation procedure starting from each solution of the asymptotic problem (19) for $\lambda \rightarrow \infty$.

To this end, we present some properties of map $g(\lambda(\mathbf{W}\mathbf{x} + \mathbf{c}))$ for $\lambda \rightarrow \infty$. First of all consider the following definition.

Definition 6.2 *For any $\mathbf{x} = (x_1, \dots, x_n) \in \mathbb{R}^n$, let $\max(\mathbf{x})$ be the maximum of \mathbf{x} , $l_{\mathbf{x}}$ be the cardinality of the set $\mathcal{M}_{\mathbf{x}} = \{i = 1, \dots, n \mid x_i = \max \mathbf{x}\}$, and define function $\text{sat}(\mathbf{x}) = (s_1, \dots, s_n)$ as*

$$s_i = \begin{cases} 0 & \text{if } x_i \notin \mathcal{M}_{\mathbf{x}} \\ 1/l_{\mathbf{x}} & \text{if } x_i \in \mathcal{M}_{\mathbf{x}}. \end{cases}$$

Observation 6.6 *The following properties hold, $\forall \mathbf{x} \in \mathbb{R}^n, \forall p \in \mathbb{N}$*

- i) $\lim_{\lambda \rightarrow \infty} \mathbf{g}(\lambda \mathbf{x}) = \text{sat}(\mathbf{x})$,
- ii) $\lim_{\lambda \rightarrow \infty} \lambda^p (\text{sat}(\mathbf{x}) - \mathbf{g}(\lambda \mathbf{x})) = \mathbf{0}$.
- iii) *for any Lipschitz vector-valued function \mathbf{h}*

$$\lim_{\lambda \rightarrow \infty} \lambda^p \mathbf{h}(\text{sat}(\mathbf{x})) = \lim_{\lambda \rightarrow \infty} \lambda^p \mathbf{h}(\mathbf{g}(\lambda \mathbf{x})).$$

Proof. i) follows from the definition of \mathbf{g} . To prove ii) let $\beta_i(\mathbf{x}, \lambda) = (\text{sat}(\mathbf{x}) - \mathbf{g}(\lambda\mathbf{x}))_i$, where notation $(\mathbf{z})_i$ denotes the i -th component of a vector \mathbf{z} . Assume first that $i \in \mathcal{M}_{\mathbf{x}}$, then $(\text{sat}(\mathbf{x}))_i = \frac{1}{l_{\mathbf{x}}}$ and

$$\beta_i(\mathbf{x}, \lambda) = \frac{1}{l_{\mathbf{x}}} - \frac{e^{\lambda x_i}}{l_{\mathbf{x}} e^{\lambda x_i} + \sum_{k \notin \mathcal{M}_{\mathbf{x}}} e^{\lambda x_k}} = \frac{1}{l_{\mathbf{x}}} \frac{\sum_{k \notin \mathcal{M}_{\mathbf{x}}} e^{\lambda x_k}}{e^{\lambda x_i} + l_{\mathbf{x}}^{-1} \sum_{k \notin \mathcal{M}_{\mathbf{x}}} e^{\lambda x_k}},$$

therefore

$$|\beta_i(\mathbf{x}, \lambda)| \leq \frac{1}{l_{\mathbf{x}}} \sum_{k \notin \mathcal{M}_{\mathbf{x}}} e^{\lambda(x_k - x_i)},$$

since $x_k - x_i < 0, \forall k \notin \mathcal{M}_{\mathbf{x}}$, $\beta_i(\mathbf{x}, \lambda)$ decreases exponentially with λ and ii) follows. Similarly, if $i \notin \mathcal{M}_{\mathbf{x}}$, let $x_m = \max \mathbf{x}$, then $(\text{sat}(\mathbf{x}))_i = 0$ and

$$|\beta_i(\mathbf{x}, \lambda)| \leq \frac{e^{\lambda x_i}}{e^{\lambda x_m}} = e^{\lambda(x_i - x_m)},$$

since $x_i - x_m < 0, \forall k \notin \mathcal{M}_{\mathbf{x}}$, $\beta_i(\mathbf{x}, \lambda)$ decreases exponentially with λ . iii) Follows directly from ii) since

$$\begin{aligned} \left\| \lim_{\lambda \rightarrow \infty} \lambda^p \mathbf{h}(\text{sat}(\mathbf{x})) - \lim_{\lambda \rightarrow \infty} \lambda^p \mathbf{h}(\mathbf{g}(\lambda\mathbf{x})) \right\| &= \left\| \lim_{\lambda \rightarrow \infty} \lambda^p (\mathbf{h}(\text{sat}(\mathbf{x})) - \mathbf{h}(\mathbf{g}(\lambda\mathbf{x}))) \right\| \\ &\leq \lim_{\lambda \rightarrow \infty} \lambda^p L (\|\text{sat}(\mathbf{x}) - \mathbf{g}(\lambda\mathbf{x})\|) = 0, \end{aligned}$$

where $L > 0$ is the Lipschitz constant for function \mathbf{h} . □

Observation 6.7

$$\lim_{\lambda \rightarrow \infty} \|(\partial_{\mathbf{x}} \mathbf{K}(\mathbf{x}, \lambda))^{-1}\| \leq 1. \quad (27)$$

Proof.

Remark that if for a n -dimensional square matrix \mathbf{A} ,

$$\|\mathbf{A}\mathbf{w}\| \geq \|\mathbf{w}\| \text{ for any } \mathbf{w} \neq 0 \quad (28)$$

then $\|\mathbf{A}^{-1}\| \leq 1$. In fact, (28) implies that \mathbf{A} is nonsingular and that for any \mathbf{v} there exists \mathbf{w} such that $\mathbf{v} = \mathbf{A}\mathbf{w}$. Moreover

$$\|\mathbf{A}^{-1}\mathbf{v}\| = \|\mathbf{w}\| \leq \|\mathbf{A}\mathbf{w}\| = \|\mathbf{v}\|.$$

Therefore to prove (27) it is sufficient to show that, $\forall \mathbf{x} \in \mathbb{R}^n \setminus \{\mathbf{0}\}$

$$\lim_{\lambda \rightarrow \infty} \|(\lambda \mathbf{g}'(\lambda(\mathbf{W}\mathbf{x} + \mathbf{c}))\mathbf{W} - \mathbf{I})\mathbf{v}\| \geq \|\mathbf{v}\|.$$

By (24) it follows that

$$\begin{aligned} &\lim_{\lambda \rightarrow \infty} \|(\lambda \mathbf{g}'(\lambda(\mathbf{W}\mathbf{x} + \mathbf{c}))\mathbf{W} - \mathbf{I})\mathbf{v}\| \\ &= \lim_{\lambda \rightarrow \infty} \|(\lambda \text{diag}(\mathbf{g}(\lambda(\mathbf{W}\mathbf{x} + \mathbf{c}))) (\mathbf{I} - \mathbf{e}\mathbf{g}(\lambda(\mathbf{W}\mathbf{x} + \mathbf{c}))^T) \mathbf{W} - \mathbf{I})\mathbf{v}\| \quad (29) \\ &= \lim_{\lambda \rightarrow \infty} \|(\lambda \text{diag}(\text{sat}(\mathbf{W}\mathbf{x} + \mathbf{c})) (\mathbf{I} - \mathbf{e}\text{sat}(\mathbf{W}\mathbf{x} + \mathbf{c}))^T) \mathbf{W} + \psi(\lambda) - \mathbf{I})\mathbf{v}\| \end{aligned}$$

where we have used iii) of Observation 6.6 with

$$\mathbf{h}(\mathbf{z}) = \text{diag}(\mathbf{z})(\mathbf{I} - \mathbf{e}\mathbf{z})\mathbf{W},$$

$p = 1$, and $\psi(\lambda)$ is a function such that $\lim_{\lambda \rightarrow \infty} \psi(\lambda) = 0$. The thesis follows from the fact that either \mathbf{v} is such that $\text{diag}(\text{sat}(\mathbf{W}\mathbf{x} + \mathbf{c}))(\mathbf{I} - \mathbf{e}\text{sat}(\mathbf{W}\mathbf{x} + \mathbf{c}))^T \mathbf{W}\mathbf{v} = \mathbf{0}$ and then the value of limit (27) is $\|\mathbf{v}\|$, or $\text{diag}(\text{sat}(\mathbf{W}\mathbf{x} + \mathbf{c}))(\mathbf{I} - \mathbf{e}\text{sat}(\mathbf{W}\mathbf{x} + \mathbf{c}))^T \mathbf{W}\mathbf{v} \neq \mathbf{0}$, in which case the value of the limit is infinity. \square

The following observation shows that, for λ sufficiently high, it is possible to find a distinct solution of $\mathbf{K}(\mathbf{x}, \lambda) = 0$ by gradient descent from any of the solutions of the asymptotic problem (19).

Observation 6.8 *Let $\hat{\mathbf{x}}_i, i = 1, \dots, n$ be distinct solutions of (19). There exists $\hat{\lambda}$ sufficiently high such that the following properties are satisfied:*

i) The solutions \mathbf{x}_i of the following differential equations, for $i = 1, \dots, n$ are well defined for all $t \geq 0$

$$\begin{cases} \dot{\mathbf{x}}_i(t) = -(\partial_{\mathbf{x}}\mathbf{K}(\mathbf{x}_i, \hat{\lambda}))^{-1}\mathbf{K}(\mathbf{x}_i, \hat{\lambda}) \\ \mathbf{x}_i(0) = \hat{\mathbf{x}}_i. \end{cases} \quad (30)$$

ii) Setting $\mathbf{z}_{i,\hat{\lambda}} = \lim_{t \rightarrow \infty} \mathbf{x}_i(t)$, then $\mathbf{K}(\mathbf{z}_{i,\hat{\lambda}}, \hat{\lambda}) = 0$, for $i = 1, \dots, n$.

iii) All vectors $\mathbf{z}_{i,\hat{\lambda}}$ are distinct.

Proof. Since Σ_n is compact, by (27) there exists $\bar{\lambda}$ sufficiently high such that

$$\|(\partial_{\mathbf{x}}\mathbf{K}(\mathbf{x}, \lambda))^{-1}\| \leq 2, \forall \mathbf{x} \in \Sigma_n, \forall \lambda \geq \bar{\lambda}. \quad (31)$$

For $\lambda \geq \bar{\lambda}$, the solution of (30) is defined for all $t \geq 0$ and satisfies

$$\partial_t \mathbf{K}(\mathbf{x}_i(t), \lambda) = -\mathbf{K}(\mathbf{x}_i(t), \lambda), \quad (32)$$

therefore

$$\mathbf{K}(\mathbf{x}_i(t), \lambda) = e^{-t}\mathbf{K}(\hat{\mathbf{x}}_i, \lambda),$$

which implies that $\mathbf{z}_{i,\lambda}$ satisfies $\mathbf{K}(\mathbf{z}_{i,\lambda}, \lambda) = 0$.

It remains to prove that there exists $\hat{\lambda}$, sufficiently high, such that $\mathbf{z}_{i,\hat{\lambda}}$ are distinct. Let $\rho = \min_{i \neq j} \|\hat{\mathbf{x}}_i - \hat{\mathbf{x}}_j\|$ denote the minimum distance between each pair of distinct asymptotic solutions. For any $i = 1, \dots, n$, by (31) and (32), $\|\dot{\mathbf{x}}_i(t)\| \leq 2\|\mathbf{K}(\mathbf{x}_i(t), \lambda)\| = 2e^{-t}\|\mathbf{K}(\hat{\mathbf{x}}_i, \lambda)\|$, which implies that

$$\|\mathbf{x}_i(t) - \mathbf{x}_i(0)\| = \|\mathbf{x}_i(t) - \hat{\mathbf{x}}_i\| \leq 2 \int_0^{+\infty} e^{-t} \|\mathbf{K}(\hat{\mathbf{x}}_i, \lambda)\| dt = 2\|\mathbf{K}(\hat{\mathbf{x}}_i, \lambda)\|$$

and $\|\mathbf{z}_{\lambda,i} - \hat{\mathbf{x}}_i\| \leq 2\|\mathbf{K}(\hat{\mathbf{x}}_i, \lambda)\|$. Since for $i = 1, \dots, n$, $\lim_{\lambda \rightarrow \infty} \mathbf{K}(\hat{\mathbf{x}}_i, \lambda) = 0$, it is possible to choose $\hat{\lambda}$ sufficiently high such that

$$\|\mathbf{K}(\hat{\mathbf{x}}_i, \hat{\lambda})\| < \frac{1}{6}\rho, \forall i = 1, \dots, n,$$

then for each pair $\mathbf{z}_{i,\hat{\lambda}}, \mathbf{z}_{j,\hat{\lambda}}$

$$\begin{aligned} \|\mathbf{z}_{i,\hat{\lambda}} - \mathbf{z}_{j,\hat{\lambda}}\| &= \|(\mathbf{z}_{i,\hat{\lambda}} - \hat{\mathbf{x}}_i) - (\mathbf{z}_{j,\hat{\lambda}} - \hat{\mathbf{x}}_j) + \hat{\mathbf{x}}_i - \hat{\mathbf{x}}_j\| \\ &\geq -\|\mathbf{z}_{i,\hat{\lambda}} - \hat{\mathbf{x}}_i\| - \|\mathbf{z}_{j,\hat{\lambda}} - \hat{\mathbf{x}}_j\| + \|\hat{\mathbf{x}}_i - \hat{\mathbf{x}}_j\| \geq \frac{1}{3}\rho. \end{aligned}$$

□

From Observation 6.8 one obtains the following procedure for finding a subset of the solutions of (15).

- 1) Find all asymptotic solutions of (19) for $\lambda \rightarrow \infty$, e.g., by solving a LCP as stated in Observation 3.2 (in this work, due to the limited size of n , we solve it through a simple enumeration).
- 2) For each solution, find a $\bar{\lambda}$ sufficiently high such that (30) converges.
- 3) Continue, if possible, the solution branch from $\lambda = \bar{\lambda}$ to $\lambda = 1$ by a numerical continuation procedure.

As will be shown in the section about computational experiments, in most (though not all) cases this method allows finding all solutions of (15).

7 Branch-and-Reduce method

In this section we first derive a convex and a linear relaxation for problem (12) (with $\lambda = 1$ in what follows), and then we employ such relaxations within a branch-and-reduce scheme (B&R in what follows), which allows to detect all solutions.

7.1 Relaxations

Let us assume that a lower bound $\ell_i > 0$ and an upper bound u_i for each variable $x_i, i = 1, \dots, n$, are given. E.g., we can initially set ℓ_i equal to the limitation from below presented in Observation 5.2, while we can set $u_i = 1$.

As a first step, we substitute each equation

$$\ln(x_i) - \ln(x_j) + [(\mathbf{W}_i - \mathbf{W}_j)^T \mathbf{x} + c_i - c_j] = 0$$

with the pair of inequalities

$$\ln(x_i) - \ln(x_j) + [(\mathbf{W}_i - \mathbf{W}_j)^T \mathbf{x} + c_i - c_j] \leq 0$$

$$\ln(x_i) - \ln(x_j) + [(\mathbf{W}_i - \mathbf{W}_j)^T \mathbf{x} + c_i - c_j] \geq 0$$

Such inequalities are easily convexified by substituting $\ln(x_i)$ in the \leq inequality, and $\ln(x_j)$ in the \geq inequality, with the corresponding convex envelopes $\xi_i(x_i)$ and $\xi_j(x_j)$ over the intervals $[\ell_i, u_i]$ and $[\ell_j, u_j]$, respectively. In particular, $\xi_i(x_i)$ (similar for $\xi_j(x_j)$) is the line interpolating $\ln(x_i)$ at the extremes of the interval. Therefore, the original nonlinear and nonconvex system is relaxed into the still nonlinear, but now convex, system

$$\begin{cases} \xi_i(x_i) - \ln(x_j) + [(\mathbf{W}_i - \mathbf{W}_j)^T \mathbf{x} + c_i - c_j] \leq 0 & i, j \in \{1, \dots, n\}, i \neq j \\ \ln(x_i) - \xi_j(x_j) + [(\mathbf{W}_i - \mathbf{W}_j)^T \mathbf{x} + c_i - c_j] \geq 0 & i, j \in \{1, \dots, n\}, i \neq j \\ \mathbf{e}^T \mathbf{x} - 1 = 0 \\ \mathbf{x} \geq \mathbf{0} \end{cases}$$

A relevant observation from the practical point of view (see [30]), is that while theoretically solvable in polynomial time, nonlinear convex problems encounter

numerical problems and are solved less efficiently than linear problems. Therefore, in [30] it is suggested to transform nonlinear convex problems into (higher dimensional) linear problems. In order to do that we can limit from above the concave function $\ln(x_i)$ through a piecewise linear function, obtained by taking the minimum of the tangent lines to $\ln(x_i)$ computed at a finite set of points in $[\ell_i, u_i]$. Different choices are possible for the set of t points X_i^t at which the tangent lines are computed. The resulting piecewise linear function is

$$\min_{\bar{x} \in X_i^t} \frac{1}{\bar{x}}(x_i - \bar{x}) + \ln(\bar{x}). \quad (33)$$

By introducing auxiliary variables $z_i, i = 1 \dots n$, we end up with the linear system

$$\begin{cases} \xi_i(x_i) - z_j + [(\mathbf{W}_i - \mathbf{W}_j)^T \mathbf{x} + c_i - c_j] \leq 0 & i, j \in \{1, \dots, n\}, i \neq j \\ z_i - \xi_j(x_j) + [(\mathbf{W}_i - \mathbf{W}_j)^T \mathbf{x} + c_i - c_j] \geq 0 & i, j \in \{1, \dots, n\}, i \neq j \\ z_i \leq \frac{1}{\bar{x}}(x_i - \bar{x}) + \ln(\bar{x}) & i = 1, \dots, n, \forall \bar{x} \in X_i^t \\ \mathbf{e}^T \mathbf{x} - 1 = 0 \\ \mathbf{x} \geq \mathbf{0} \end{cases} \quad (34)$$

The set of points X_i^t can be chosen in many different ways (see, e.g., [22]). In our approach we use the following formula

$$X_i^{t,k} = \{x \in [\ell_i, u_i] : x = \ell_i + \left(\frac{i}{t}\right)^k (u_i - \ell_i), i = 0, \dots, t-1\},$$

where k is an integer larger than one, so that more points are located close to lower bound of the interval, where the slope of the logarithmic function is higher.

7.2 Range reduction strategy

Once we have a relaxation, the next step is that of applying a range reduction strategy (see, e.g., [29]), i.e., we try to reduce the range $[\ell_i, u_i]$ of each variable x_i , without removing any solution of the original system.

A standard way to accomplish that is by solving the following linear problems

$$\min \setminus \max \quad x_i$$

$$\begin{cases} \xi_i(x_i) - z_j + [(\mathbf{W}_i - \mathbf{W}_j)^T \mathbf{x} + c_i - c_j] \leq 0 & i, j \in \{1, \dots, n\}, i \neq j \\ z_i - \xi_j(x_j) + [(\mathbf{W}_i - \mathbf{W}_j)^T \mathbf{x} + c_i - c_j] \geq 0 & i, j \in \{1, \dots, n\}, i \neq j \\ z_i \leq \frac{1}{\bar{x}}(x_i - \bar{x}) + \ln(\bar{x}) & i = 1, \dots, n, \forall \bar{x} \in X_i^t \\ \mathbf{e}^T \mathbf{x} - 1 = 0 \\ \mathbf{x} \geq \mathbf{0} \end{cases} \quad (35)$$

The interval defined by the optimal values of these two linear problems includes all values of x_i which can be found in all solutions of our problem. In order to

strengthen the reduction, we can define a cycle where at each iteration problems (35) are solved for all variables $x_i, i = 1, \dots, n$, and the cycle stops as soon as no (significant) reduction is observed from one iteration to the next.

7.3 Branch-and-Reduce approach

Different outcomes of such range reduction strategy are possible for a given initial box $B = \prod_{i=1}^n [\ell_i, u_i]$:

Case 1 a single point: in this case the point is also the unique solution of the problem which can be found within box B ;

Case 2 an empty set: in this case the problem has no solution within box B (if the box is defined by the original bounds, then Observation 5.1 rules out this possibility);

Case 3 a new n -dimensional box.

In the first two cases the problem over box B is solved. Therefore, we need to proceed only if we are in the third case. What we can do is to perform a branching of the box, i.e., the box is subdivided into two sub-boxes (e.g., through a bisection at the midpoint of the largest edge of the box), and the range reduction strategy is applied again over the sub-boxes. The reduction of a sub-box ends up with one of the three cases previously discussed and if Case 3 holds, we need to further subdivide the sub-box. This way we are basically performing a B&R approach (see, e.g., [23]), where nodes/sub-boxes are fathomed as soon as the range reduction strategy ends up with Case 1 (a solution has been detected), or Case 2 (the sub-box does not contain any solution). Note that, by this approach, we are able, in principle, to detect *all* the solutions of the system. Also note that B&R approaches are usually applied for the solution of GO problems. In fact, this is also our case, where we are basically solving the GO problem (20). However, with respect to typical implementations, here we do not need to compute lower and upper bounds, rather we only need to perform range reductions.

8 Computational experiments

In the following subsections we will present computational results about the B&R method, introduced in Section 7, and the homotopy-based heuristics, as in Section 6.

We performed random test sets, named R1, where we have uniformly sampled T from $[900, 1500]$ and \mathbf{c} from $[-10^5, 10^5]^n$, while \mathbf{W} has been uniformly sampled from $[-10^5, 10^5]^{n \times n}$. We generated 1000 tests for each n value. In view of the limitations on the number of components n in the solid phase, we only considered n values up to 12. In fact, we also considered a further set, named R2, where \mathbf{W} has been generated from real-life values taken from Table (A 4-3) in [14]. However, since the results for both test sets are very similar, we only report those for the R1 test set.

Data are summarized using boxplots, which allow for a compact representation of standard statistics: for each data series, a box is plot for the interval

embracing the interval from the 25th to the 75th percentile, with a red line representing the median; whiskers comprise an interval 1.5 times larger than the box, while the remaining outliers are marked as red crosses.

8.1 Experiments with Branch-and-Reduce

The B&R strategy proposed in Section 7 has been implemented in C++, using CPLEX 11.1 as a solver for the linear programming subproblems. All the tests have been performed on a Pentium i5 standard desktop machine running Linux, 2Gbyte RAM and 2.67GHz quad-core CPU (but no multi-thread or parallel computation has been carried out).

Although in principle B&R is able to return all solutions of our problem, we need to consider some numerical difficulties. The maximization/minimization of a variable x_i might not return exactly the solution of the problem but, mainly due to rounding errors and tolerances, a slightly interior solution is obtained. In particular, this is likely to occur when solutions with some variables having values close to 0 exist, because of the large negative values attained by the logarithmic function at these values. Then, due to such numerical errors, we might cut a solution of our problem. If the whole process ends up with no solution at all, which is not possible in view of Observation 5.1, we are aware an error occurred. Analogously, if we end up with an even number of solutions, we have very likely missed some solution. But even if we end up with an odd number of solutions, it might have happened that a solution has been discarded. To ease these problems, we need to add a tolerance, i.e., the solutions of problems (35) are decreased (lower limit) or increased (upper limit) by a small quantity $\rho > 0$. The latter plays an important role in the trade-off between performance and accuracy. Moreover, it turned out that accuracy could be improved by the choice to loosen CPLEX feasibility accuracy (that we will denote as ϵ_c) and to force the solver to use the simplex algorithm (interior point methods are likely to produce strictly interior solutions). Further strategies to increase the accuracy are under investigation.

We also notice that a box is reduced until the width of its largest edge falls below a given threshold and a single solution is detected within that box. In fact, it might happen that such box contains different solutions quite close to each other. In such case the approach consider such solutions as indistinguishable.

As we will see, the computation times of the B&R approach are usually not competitive with those of the homotopy methods. This is the reason why we have not performed a detailed comparison of the proposed approach with existing B&R methods and, in particular, with BARON [23]. In fact, the proposed B&R approach has been mainly employed to compare the solutions returned by such method and those returned by the homotopy methods. It is important to emphasize here an important difference between B&R and homotopy methods. The former have a theoretical guarantee of detecting *all* the solutions, while the latter do not have the same guarantee. This is obvious for the forward homotopy, which only returns a single solution, but it may also happen for the backward homotopy since not all solutions for $\lambda = 1$ can be reached through continuous trajectories starting at solutions for $\lambda = \infty$. On the other hand, as previously commented, B&R methods might miss some solutions because of numerical difficulties (a few low-dimensional experiments revealed that also BARON share the same numerical difficulties). The differences between the

# components	# solutions											no solution	
	1	2	3	4	5	6	7	8	9	10	11		
2	897	0	103	0	0	0	0	0	0	0	0	0	0
3	849	0	150	0	0	0	1	0	0	0	0	0	0
4	806	0	190	0	2	0	1	0	0	0	0	0	1
5	795	0	201	0	3	0	1	0	0	0	0	0	0
6	771	3	212	0	10	0	2	0	0	0	0	0	2
7	718	1	264	0	16	0	1	0	0	0	0	0	0
8	710	1	258	0	21	0	7	0	1	0	0	0	2
9	684	2	290	0	20	0	2	0	0	0	1	0	1
10	664	1	304	0	19	0	9	0	1	0	0	0	2
11	642	1	307	0	37	1	8	0	1	0	0	0	3
12	646	6	306	0	33	0	7	0	0	0	0	0	2

Table 1: Number of solutions found in the R1 test set (optimality threshold set to 10^{-6} and similarity error threshold set to 10^{-5}).

solutions returned by the different methods will be commented later on.

In Figure 2 we report the execution time, the number of linear subproblems solved as well as the number of nodes expanded during the B&R execution, for each dimension, i.e., number of solid components.

As a comparison, we also run the same Branch-and-Reduce algorithm, but with no range reduction (CPLEX parameters are the same as before) on test set R1. Results are reported in Figure 3. Such results clearly show the advantages brought by range reductions: the computational cost of the reduction is largely compensated by a much lower number of nodes.

The actual number of solutions found has been checked starting from them a local search on Problem (20) (using the SNOPT 7.0 SQP solver, see [15]). We have considered only solutions with an objective function value smaller than 10^{-6} , clustering together those closer than 10^{-5} . The results for test set R1 are reported in Table 1. The last column accounts for the instances for which no solution has been found. We remark that the overall number of "anomalies" (no solution or even number of solutions) is quite limited: 29 out of 11000, with the largest number of "anomalies" (8 out of 1000) for $n = 12$.

In terms of computing times, the performance of the proposed B&R approach is enhanced by decreasing ρ , as depicted in Figure 4 where we show results for the R1 test set with $\rho = 10^{-10}$. On the other hand, we also need to point out that numerical difficulties increase in such case and the number of "anomalies" is doubled.

8.2 Experiments with homotopies

The homotopy methods presented in Section 6 have been implemented with MATLAB and run on an Intel Core2Duo CPU at 3.16 GHz. In particular the LCP (19) has been solved by brute force by iterating on all possible choices of the zero entries of vector \mathbf{x} and the numerical continuation was based on a Runge-Kutta predictor step and a Newton corrector step.

# components	Homotopy only sols.(%)	B&R only sols.(%)
2	0	0
3	0.152	0
4	0.711	0.427
5	0.560	0.560
6	0.738	0.268
7	0.312	0.249
8	0.241	0.241
9	0.531	0.118
10	0.569	0.228
11	0.598	0.489
12	0.774	0.387

Table 2: Discrepancies between B&R and homotopy solutions

8.2.1 Forward homotopy from $\lambda = 0$ to $\lambda = 1$

The execution times for this homotopy method are reported in Figure 5. We notice that the times compare quite favorably with the other two methods, but we should also remark that this method allows finding a single solution. It can also be proved (see, e.g., [19]) that the single solution detected by this homotopy is always one on those also detected by the backward homotopy.

8.2.2 Backward homotopy from $\lambda = +\infty$ to $\lambda = 1$

This method allows finding a set of many solutions. The execution times are depicted in Figure 6.

This figure shows that homotopy is considerably faster than the B&R method in the higher dimensional problems, while it is slower in the problem of dimension 2 and 3. The performance of the homotopy method could be further improved using a faster method for the asymptotic LCP and implementing the numerical continuation method, e.g., in the C programming language.

Differently from B&R, the homotopy method does not allow finding all solutions, since it can find only those solutions which are part of a solution branch that continues to $\lambda = +\infty$. In numerical experiments, comparing the solutions found with the two methods, we discovered that this kind of solution is rare, since it represents less than 0.5% of the total number of solutions.

Experiments have shown that, in more or less the 99% of the cases, the two methods completely agree and give the same set of solutions. In some cases the homotopy method did not find some of the solutions found by B&R for the reason explained above. In other cases, for numerical reasons, B&R failed to find some solutions that were found with homotopy.

Table 8.2.2 shows the discrepancies in the solution sets found by the two methods. In particular, the first column shows the problem dimension, the second column the percentage of the solutions that have been found by homotopy and not found by B&R over the total solutions, conversely the third column shows the percentage of solutions that have been found by B&R and not by homotopy.

9 Conclusions

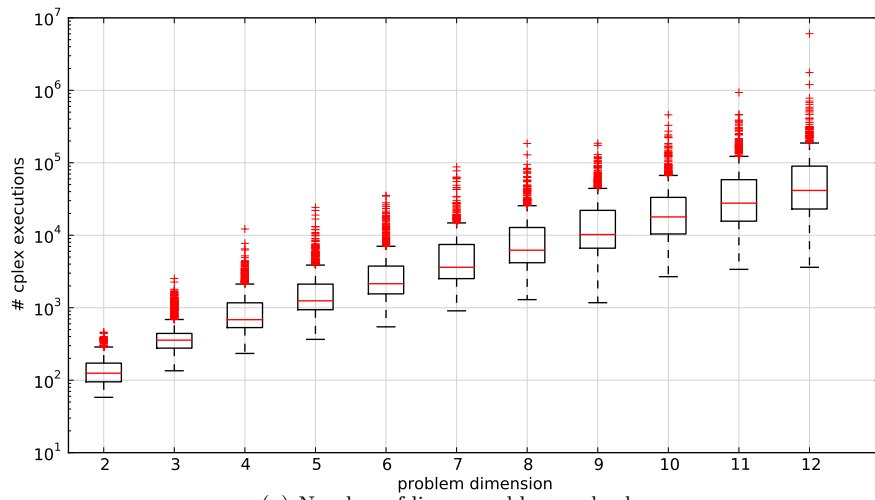
In this paper we considered a mathematical model based on the Gibbs free energy for magmatic mixtures. We have discussed different reformulations and solution methods for this problem. Some theoretical results concerning the existence and the number of solutions have been proved. Two rather fast homotopy methods have been proposed, one returning a single solution, the other returning multiple solutions (although not necessarily all of them). A slower B&R approach has also been proposed which is able, in principle, to return all solutions, but in practice might miss a few of them because of numerical difficulties. Future developments might include: new theoretical results, e.g., related to solutions which do not stem from trajectories starting at $\lambda = 0$ or $\lambda = +\infty$; new strategies to remove numerical difficulties in the B&R approach; more efficient implementations of all the proposed approaches.

References

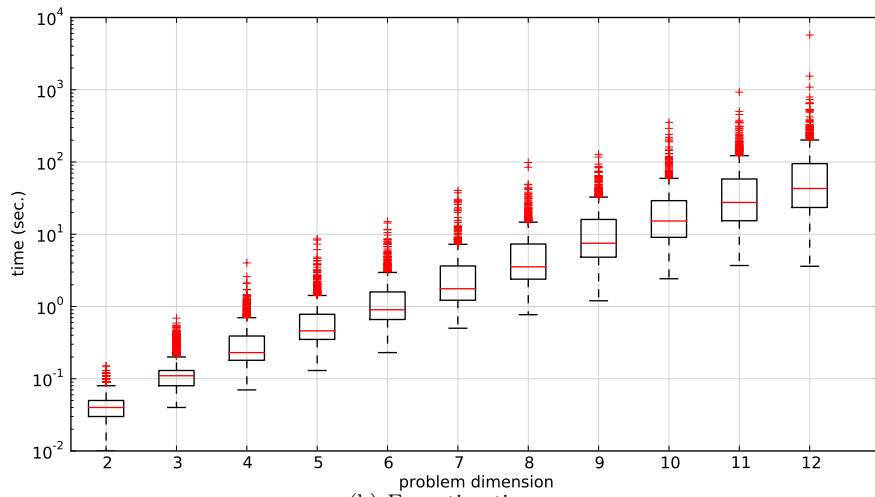
- [1] G. M. Anderson and D. A. Crerar. *Thermodynamics in Geochemistry. The Equilibrium Model*. Oxford University Press, 1993.
- [2] Ghiorso M.S. Asimow P.D. Algorithmic modifications extending MELTS to calculate subsolidus phase relations. *American Mineralogist*, 83:1127–1131, 1998.
- [3] M. G. Best and E. H. Christiansen. *Igneous Petrology*. Wiley-Blackwell, 2001.
- [4] A. E. Boudreau. PELE—a version of the MELTS software program for the PC platform. *Computers & Geosciences*, 25(2):201–203, 1999.
- [5] L.E.J. Brouwer. Über Abbildung von Mannigfaltigkeiten. *Math. Ann.*, 71:97–115, 1912.
- [6] George Dantzig, Selmer Johnson, and Wayne White. A linear programming approach to the chemical equilibrium problem. *Management Science*, 5(1):pp. 38–43, 1958.
- [7] F. Dobran and J.I. Ramos. Chapter 7 global volcanic simulation: Physical modeling, numerics, and computer implementation. In Flavio Dobran, editor, *Vesuvius - Education, Security and Prosperity*, volume 8 of *Developments in Volcanology*, pages 311–372. Elsevier, 2006.
- [8] K. Georg E. L. Allgower. *Introduction to Numerical Continuation Methods*. SIAM, 2003.
- [9] O. Fabrichnaya, S. K. Saxena, P. Richet, and E. F. Westrum. *Thermodynamic Data, Models and Phase*. Springer, 2004.
- [10] CA Floudas, IG Akrotirianakis, S. Caratzoulas, CA Meyer, and J. Kallrath. Global optimization in the 21st century: Advances and challenges. *Computer Aided Chemical Engineering*, 18:23–51, 2004.

- [11] Mark S. Ghiorso, Marc M. Hirschmann, Peter W. Reiners, and Victor C. III Kress. The pMELTS: An revision of MELTS aimed at improving calculation of phase relations and major element partitioning involved in partial melting of the mantle at pressures up to 3 GPa. *Geochemistry, Geophysics, Geosystems*, 3(5), 2002.
- [12] Mark S. Ghiorso and Richard O. Sack. Chemical Mass Transfer in Magmatic Processes. IV. A Revised and Internally Consistent Thermodynamic Model for the Interpolation and Extrapolation of Liquid-Solid Equilibria in Magmatic Systems at Elevated Temperatures and Pressures. *Contributions to Mineralogy and Petrology*, 119:197–212, 1995.
- [13] M.S. Ghiorso. Algorithms for the estimation of phase stability in heterogeneous thermodynamic systems. *Geochimica et Cosmochimica Acta*, 58(24):5489–5502, 1994.
- [14] M.S. Ghiorso, I.S.E. Carmichael, M.L. Rivers, and R.O. Sack. The gibbs free energy of mixing of natural silicate liquids; an expanded regular solution approximation for the calculation of magmatic intensive variables. *Contributions to Mineralogy and Petrology*, 84(2):107–145, 1983.
- [15] P. E. Gill, W. Murray, and Michael A. Saunders. *User’s guide for SNOPT version 7: Software for large-scale nonlinear programming*. Systems Optimization Laboratory - University of Stanford, 2007.
- [16] F. Jalali and JD Seader. Homotopy continuation method in multi-phase multi-reaction equilibrium systems. *Computers & Chemical Engineering*, 23(9):1319–1331, 1999.
- [17] I.K. Karpov, K.V. Chudnenko, D.A. Kulik, and V.A. Bychinskii. The convex programming minimization of five thermodynamic potentials other than Gibbs energy in geochemical modeling. *American Journal of Science*, 302(4):281, 2002.
- [18] W. Marzocchi, L. Sandri, and J. Selva. BET_EF: a probabilistic tool for long- and short-term eruption forecasting. *Bullettin of Volcanology*, 70:623–632, 2008.
- [19] J. Mawhin. Leray-schauder continuation theorems in the absence of a priori bounds. *Topological Methods in Nonlinear Analysis, Journal of the Juliusz Schauder Center*, 9:179–200, 1997.
- [20] Roberto Moretti and Paolo Papale. On the oxidation state and volatile behavior in multicomponent gas-melt equilibria. *Chemical Geology*, 213(1-3):265 – 280, 2004. 7th Silicate Melt Workshop.
- [21] Paolo Papale, Roberto Moretti, and David Barbato. The compositional dependence of the saturation surface of H₂O+CO₂ fluids in silicate melts. *Chemical Geology*, 229(1-3):78–95, 2006. Physics, Chemistry and Rheology of Silicate Melts and Glasses.
- [22] G. Rote. The convergence rate of the sandwich algorithm for approximating convex functions. *Computing*, 48:337–361, 1992.

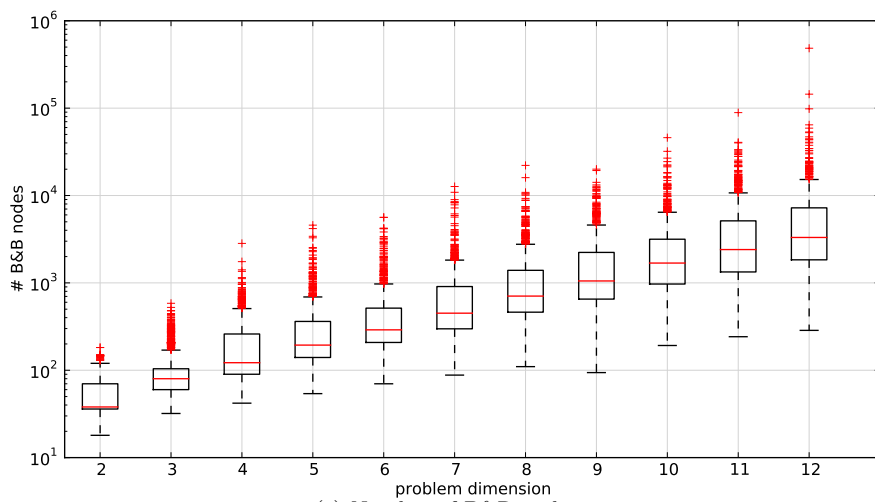
- [23] Nikolaos V. Sahinidis. *BARON Branch And Reduce Optimization Navigator*. University of Illinois at Urbana-Champaign Department of Chemical Engineering, 600 South Mathews Avenue Urbana, Illinois 61801, USA, 0.4 edition, June 2000.
- [24] Haraldur Sigurdsson. *Encyclopedia of Volcanoes*. Academic Press, 1999.
- [25] Paula M. Smith and Paul D. Asimow. Adibat - 1ph: A new public front-end to the melts, pmelts, and phmelts models. *Geochemistry Geophysics Geosystems*, 6(1), 2005.
- [26] R.S.J. Sparks. Forecasting volcanic eruptions. *Earth and Planetary Science Letters*, 210(1-2):1 – 15, 2003.
- [27] Mekapati Srinivas and G.P. Rangaiah. A study of differential evolution and tabu search for benchmark, phase equilibrium and phase stability problems. *Computers & Chemical Engineering*, 31(7):760 – 772, 2007.
- [28] A.C. Sun and W.D. Seider. Homotopy-continuation method for stability analysis in the global minimization of the Gibbs free energy. *Fluid Phase Equilibria*, 103(2):213–249, 1995.
- [29] Mohit Tawarmalani and Nikolaos V. Sahinidis. Global optimization of mixed-integer nonlinear programs: A theoretical and computational study. *Mathematical Programming*, 99(3):563–591, April 2004.
- [30] Mohit Tawarmalani and Nikolaos V. Sahinidis. A polyhedral branch-and-cut approach to global optimization. *Mathematical Programming*, 103:225–249, 2005.
- [31] C. Textor, H. Graf, A. Longo, A. Neri, T. E. Ongaro, P. Papale, C. Timmreck, and G. G. J. Ernst. Numerical simulation of explosive volcanic eruptions from the conduit flow to global atmospheric scales. *Annals of Geophysics*, 48(4-5), 2009.



(a) Number of linear problems solved

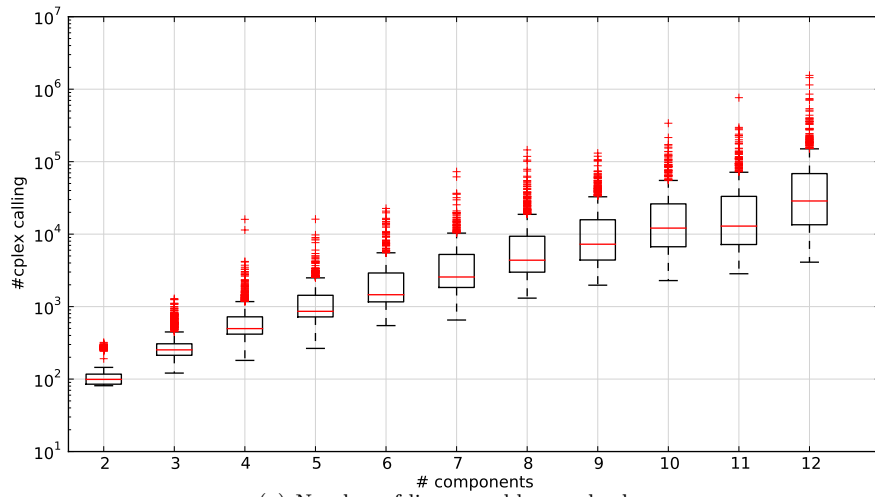


(b) Execution time.

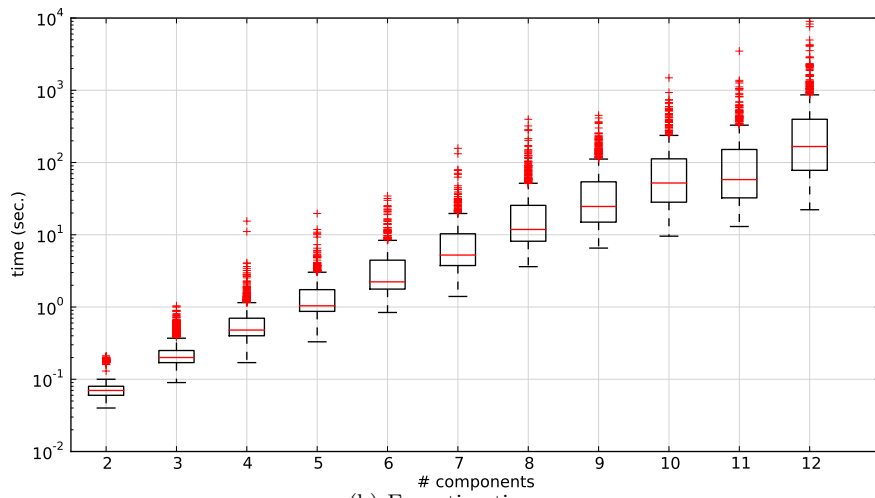


(c) Number of B&R nodes

Figure 2: Results for test set R1, $\rho = 10^{-6}$, $\epsilon_c = 10^{-7}$.

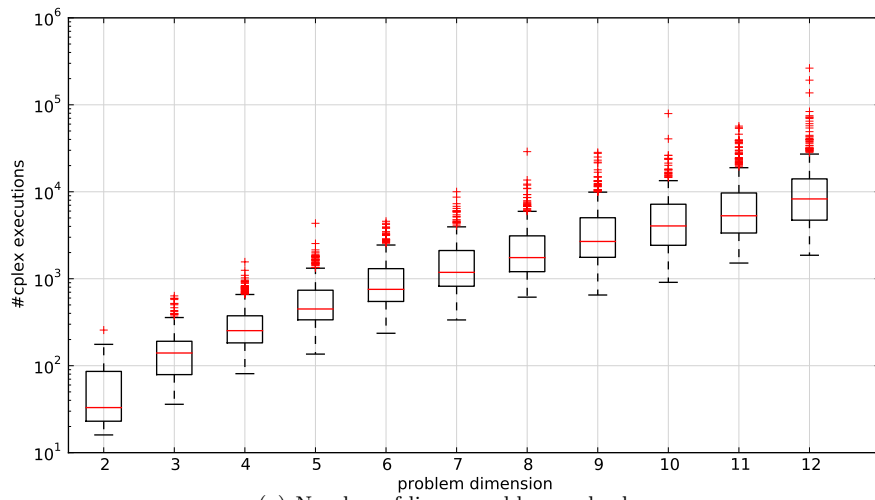


(a) Number of linear problems solved.

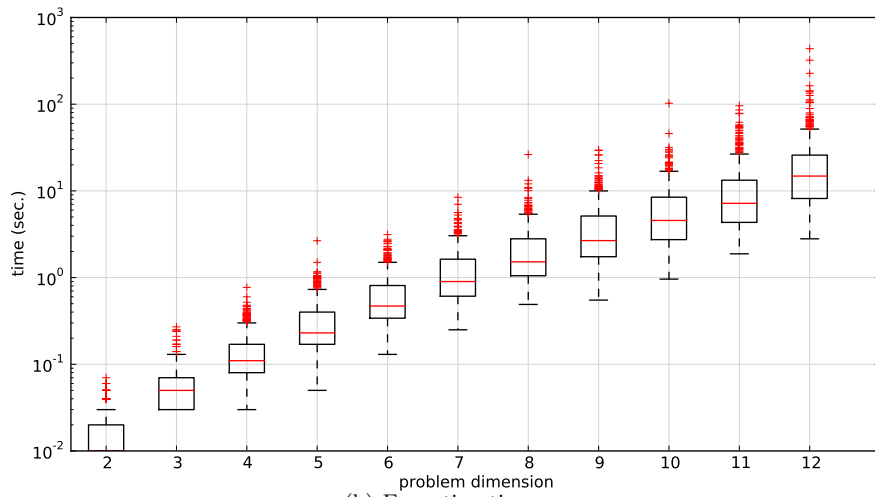


(b) Execution time.

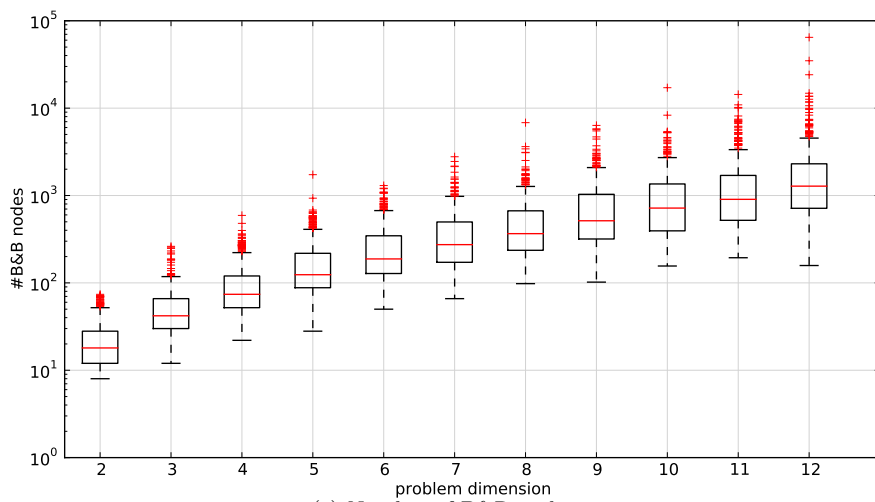
Figure 3: Execution of B&B on test set R1.



(a) Number of linear problems solved.



(b) Execution time.



(c) Number of B&B nodes.

Figure 4: Results for test set R1 ($\rho = 10^{-10}$, $\epsilon_c = 10^{-6}$).

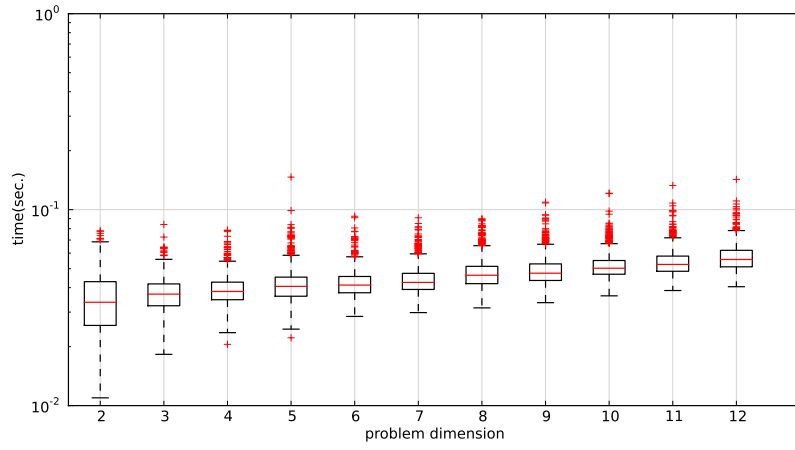


Figure 5: Results for test set R1 using the forward homotopy method for finding one solution.

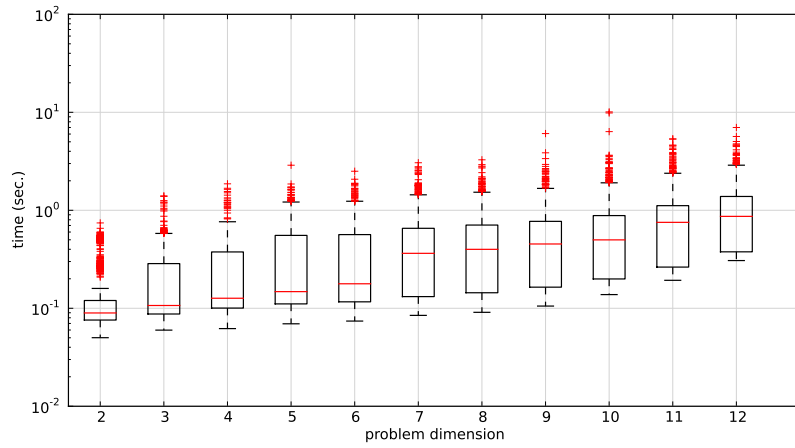


Figure 6: Results for test set R1 using the backward homotopy method for finding a set of solutions.