

A Look at the Generalized Heron Problem through the Lens of Majorization-Minimization

Eric C. Chi and Kenneth Lange

Abstract

The generalized Heron problem states: on a closed convex subset of \mathbb{R}^d , find a point such that the sum of the distances from that point to k closed convex subsets of \mathbb{R}^d is minimal. In a recent issue of this journal, Mordukhovich, Nam, and Salinas pose this problem and solve it with the tools of modern convex analysis. In light of the majorization-minimization principle from computational statistics, we revisit the problem and construct a much faster solution algorithm using only rudimentary techniques from differential calculus.

1 Introduction.

In a recent article in this journal, Mordukhovich et al. [21] presented the following generalization of the classical Heron problem. Given a collection of closed convex sets $\{C_1, \dots, C_k\}$ in \mathbb{R}^d , find a point \mathbf{x} in the closed convex set $S \subset \mathbb{R}^d$ such that the sum of the Euclidean distances from \mathbf{x} to C_1 through C_k is minimal. In other words,

$$\text{minimize } D(\mathbf{x}) := \sum_{i=1}^k d(\mathbf{x}, C_i) \text{ subject to } \mathbf{x} \in S, \quad (1)$$

where $d(\mathbf{x}, \Omega) = \inf\{\|\mathbf{x} - \mathbf{y}\| : \mathbf{y} \in \Omega\}$.

A rich history of special cases motivates this problem formulation. When $k = 2$, C_1 and C_2 are singletons, and S is a line, we recover the problem originally posed by the ancient mathematician Heron of Alexandria. The special

case where $k = 3$; C_1 , C_2 , and C_3 are singletons; and $S = \mathbb{R}^2$ was suggested by Fermat nearly 400 years ago and solved by Torricelli [13]. In his *Doctrine and Application of Fluxions*, Simpson generalized the distances to weighted distances. In the 19th century, Steiner made several fundamental contributions, and his name is sometimes attached to the problem [9, 11]. At the turn of the 20th century, the German economist Weber generalized Fermat’s problem to an arbitrary number of singleton sets C_i . Weiszfeld published the first iterative algorithm¹ for solving the Fermat-Weber problem in 1937 [26, 27]. In the modern era, the Fermat-Weber problem has enjoyed a renaissance in various computational guises. Both the problem and associated algorithms serve as the starting point for many advanced models in location theory [18, 28].

In this article, we take a second look at this distinguished problem from the perspective of algorithm design. Mordukhovich et al. [21] present an iterative subgradient algorithm for numerically solving problem (1), a reasonable choice given that the objective function is convex but non-differentiable. However, it is natural to wonder if there might be better alternatives. Here we present one that generalizes Weiszfeld’s algorithm by invoking the majorization-minimization (MM) principle from computational statistics. Although the new algorithm displays the same kind of singularities that plagued Weiszfeld’s algorithm [15], the dilemmas can be resolved by perturbing the generalized Heron problem slightly. In the limit, one recovers the solution to the unperturbed problem. The new MM algorithm is vastly superior to the subgradient algorithms in computational speed.

Solving a perturbed version of the problem by the MM principle yields extra dividends as well. The convergence of MM algorithms on smooth problems is well understood theoretically. This fact enables us to show that solutions to the original problem can be characterized without appealing to the full machinery of modern convex analysis. While this body of mathematical knowledge is both beautiful and powerful, we demonstrate that solving problem (1) is well within the scope of classical differential calculus. Its resolution can be understood by undergraduate mathematics majors. Viewing equation (1) through the MM lens also makes it obvious how to weaken the key assumption of Mordukhovich et al. [21] that $S \cap C_i = \emptyset$ for all i . Relaxing this empty intersection condition is important because it brings the related Fermat-Weber problem under the same general umbrella.

¹Kuhn [15] points out that Weiszfeld’s algorithm has been rediscovered several times.

As a brief summary of things to come, we begin by recalling background material on the MM principle and convex analysis of differentiable functions. This is followed with a derivation of the MM algorithm for problem (1) and some relevant numerical examples. We end by proving convergence of the algorithm and characterizing solution points.

2 The MM principle.

Although first articulated by the numerical analysts Ortega and Rheinboldt [22], the MM principle currently enjoys its greatest vogue in computational statistics [1, 17]. The basic idea is to convert a hard optimization problem (for example, non-differentiable) into a sequence of simpler ones (for example, smooth). The MM principle requires majorizing the objective function $f(\mathbf{y})$ by a surrogate function $g(\mathbf{y} \mid \mathbf{x})$ anchored at the current point \mathbf{x} . Majorization is a combination of the tangency condition $g(\mathbf{x} \mid \mathbf{x}) = f(\mathbf{x})$ and the domination condition $g(\mathbf{y} \mid \mathbf{x}) \geq f(\mathbf{y})$ for all $\mathbf{y} \in \mathbb{R}^d$. The associated MM algorithm is defined by the iteration scheme

$$\mathbf{x}_{k+1} = \arg \min_{\mathbf{y} \in S} g(\mathbf{y} \mid \mathbf{x}_k). \quad (2)$$

Because

$$f(\mathbf{x}_{k+1}) \leq g(\mathbf{x}_{k+1} \mid \mathbf{x}_k) \leq g(\mathbf{x}_k \mid \mathbf{x}_k) = f(\mathbf{x}_k), \quad (3)$$

the MM iterates generate a descent algorithm driving the objective function downhill. Constraint satisfaction is enforced in finding \mathbf{x}_{k+1} . Under appropriate regularity conditions, an MM algorithm is guaranteed to converge to a local minimum of the original problem.

3 Background on Convex Analysis.

As a prelude to deriving an MM algorithm, we review some basic facts from convex analysis in the limited context of differentiable functions. Deeper treatment can be found in the references [3, 4, 12, 23, 24]. Recall that a differentiable function $f(\mathbf{y})$ is convex if and only if its domain S is convex and

$$f(\mathbf{y}) \geq f(\mathbf{x}) + \langle \nabla f(\mathbf{x}), \mathbf{y} - \mathbf{x} \rangle, \quad (4)$$

for all $\mathbf{x}, \mathbf{y} \in S$. Provided $f(\mathbf{x})$ is twice differentiable, it is convex when its second differential $d^2f(\mathbf{x})$ is positive semidefinite for all \mathbf{x} and strictly convex when $d^2f(\mathbf{x})$ is positive definite for all \mathbf{x} . These characterizations are a direct consequence of executing a second-order Taylor expansion of $f(\mathbf{y})$ and applying the supporting hyperplane inequality (4). The supporting hyperplane inequality (4) also leads to a succinct necessary and sufficient condition for a global minimum. A point $\mathbf{x} \in S$ is a global minimizer of $f(\mathbf{y})$ on S if and only if

$$\langle \nabla f(\mathbf{x}), \mathbf{y} - \mathbf{x} \rangle \geq 0 \quad (5)$$

for all $\mathbf{y} \in S$. Intuitively speaking, every direction pointing into S must lead uphill.

We conclude this section by reviewing projection operators [16]. Denote the projection of \mathbf{x} onto a set $\Omega \subset \mathbb{R}^d$ by $P_\Omega(\mathbf{x})$. By definition $P_\Omega(\mathbf{x})$ satisfies

$$P_\Omega(\mathbf{x}) := \arg \min_{\mathbf{y} \in \Omega} \|\mathbf{x} - \mathbf{y}\|.$$

If Ω is a closed convex set in \mathbb{R}^d , then $P_\Omega(\mathbf{x})$ exists and is unique. Furthermore, the projection operator is non-expansive in the sense that

$$\|P_\Omega(\mathbf{x}) - P_\Omega(\mathbf{y})\| \leq \|\mathbf{x} - \mathbf{y}\|$$

for all $\mathbf{x}, \mathbf{y} \in \mathbb{R}^d$. Non-expansion clearly entails continuity. Explicit formulas for the projection operator $P_\Omega(\mathbf{x})$ exist when Ω is a box, Euclidean ball, hyperplane, or halfspace. Fast algorithms for computing $P_\Omega(\mathbf{x})$ exist for the unit simplex, the ℓ_1 ball, and the cone of positive semidefinite matrices [10, 19].

The projection operator and the distance function are intimately related through the gradient identity $\nabla d(\mathbf{x}, C)^2 = 2[\mathbf{x} - P_C(\mathbf{x})]$. For the sake of completeness, we repeat the standard proof of this fact [12, p. 181]. Let Δ denote the difference $d(\mathbf{x} + \mathbf{y}, C)^2 - d(\mathbf{x}, C)^2$. In view of the inequality $d(\mathbf{x}, C)^2 \leq \|\mathbf{x} - P_C(\mathbf{x} + \mathbf{y})\|^2$, one can construct the lower bound

$$\begin{aligned} \Delta &\geq \|\mathbf{x} + \mathbf{y} - P_C(\mathbf{x} + \mathbf{y})\|^2 - \|\mathbf{x} - P_C(\mathbf{x} + \mathbf{y})\|^2 \\ &= \|\mathbf{y}\|^2 + 2\langle \mathbf{y}, \mathbf{x} - P_C(\mathbf{x} + \mathbf{y}) \rangle \\ &= \|\mathbf{y}\|^2 + 2\langle \mathbf{y}, \mathbf{x} - P_C(\mathbf{x}) \rangle + 2\langle \mathbf{y}, P_C(\mathbf{x}) - P_C(\mathbf{x} + \mathbf{y}) \rangle, \\ &\geq \|\mathbf{y}\|^2 + 2\langle \mathbf{y}, \mathbf{x} - P_C(\mathbf{x}) \rangle - 2\|\mathbf{y}\|\|P_C(\mathbf{x}) - P_C(\mathbf{x} + \mathbf{y})\| \\ &\geq \|\mathbf{y}\|^2 + 2\langle \mathbf{y}, \mathbf{x} - P_C(\mathbf{x}) \rangle - 2\|\mathbf{y}\|^2 \end{aligned} \quad (6)$$

The penultimate inequality here reflects the Cauchy-Schwartz inequality. The final inequality is a consequence of the non-expansiveness of the projection operator. The analogous inequality $d(\mathbf{x} + \mathbf{y}, C)^2 \leq \|\mathbf{x} + \mathbf{y} - P_C(\mathbf{x})\|^2$ gives the upper bound

$$\Delta \leq \|\mathbf{x} + \mathbf{y} - P_C(\mathbf{x})\|^2 - \|\mathbf{x} - P_C(\mathbf{x})\|^2 = \|\mathbf{y}\|^2 + 2\langle \mathbf{y}, \mathbf{x} - P_C(\mathbf{x}) \rangle. \quad (7)$$

The two bounds (6) and (7) together imply that $\Delta = 2\langle \mathbf{y}, \mathbf{x} - P_C(\mathbf{x}) \rangle + o(\|\mathbf{y}\|)$ and consequently that $\nabla d(\mathbf{x}, C)^2 = 2[\mathbf{x} - P_C(\mathbf{x})]$ according to Fréchet's definition of the differential. If $d(\mathbf{x}, C)^2 > 0$, then the chain rule gives

$$\nabla d(\mathbf{x}, C) = \nabla \sqrt{d(\mathbf{x}, C)^2} = \frac{\mathbf{x} - P_C(\mathbf{x})}{d(\mathbf{x}, C)}.$$

On the interior of C , it is obvious that $\nabla d(\mathbf{x}, C) = \mathbf{0}$. In contrast, differentiability of $d(\mathbf{x}, C)$ at boundary points of C is not guaranteed.

4 An MM Algorithm for the Heron Problem.

Since it adds little additional overhead, we recast problem (1) in the Simpson form

$$\text{minimize } D(\mathbf{x}) := \sum_{i=1}^k \gamma_i d(\mathbf{x}, C_i) \text{ subject to } \mathbf{x} \in S \quad (8)$$

involving a convex combination of the distances $d(\mathbf{x}, C_i)$ with positive weights γ_i . We first derive an MM algorithm for solving problem (8) when $S \cap C_i = \emptyset$ for all i . This exercise will set the stage for attacking the more general case where S intersects one or more of the C_i . In practice quadratic majorization is desirable because it promotes exact solution of the minimization step of the MM algorithm. It takes two successive majorizations to achieve quadratic majorization in our setting. The first is the simple majorization

$$d(\mathbf{x}, C_i) \leq \|\mathbf{x} - P_{C_i}(\mathbf{x}_m)\|$$

flowing directly from the definition of the distance function. The second is the majorization

$$\sqrt{u} \leq \sqrt{u_m} + \frac{1}{2\sqrt{u_m}}(u - u_m), \quad (9)$$

of the concave function \sqrt{u} on the interval $(0, \infty)$. The combination of these two majorizations yields the quadratic majorization

$$d(\mathbf{x}, C_i) \leq \|\mathbf{x}_m - P_{C_i}(\mathbf{x}_m)\| + \frac{\|\mathbf{x} - P_{C_i}(\mathbf{x}_m)\|^2 - \|\mathbf{x}_m - P_{C_i}(\mathbf{x}_m)\|^2}{2\|\mathbf{x}_m - P_{C_i}(\mathbf{x}_m)\|}. \quad (10)$$

Summing these majorizations over i leads to quadratic majorization of $D(\mathbf{x})$ and ultimately to the MM algorithm map

$$\psi(\mathbf{x}) = \arg \min_{\mathbf{z} \in S} \left\{ \frac{1}{2} \sum_{i=1}^k w_i \|\mathbf{z} - P_{C_i}(\mathbf{x})\|^2 \right\}$$

with weights $w_i = \gamma_i \|\mathbf{x}_m - P_{C_i}(\mathbf{x}_m)\|^{-1}$. When the C_i are singletons and $S = \mathbb{R}^d$, the map $\psi(\mathbf{x})$ implements Weiszfeld's algorithm for solving the Fermat-Weber problem [26, 27].

The quadratic majorization of $D(\mathbf{x})$ just derived can be written as

$$g(\mathbf{x} \mid \mathbf{x}_m) = \frac{1}{2} \left(\sum_{i=1}^k w_i \right) \left\| \mathbf{x} - \sum_i \alpha_i P_{C_i}(\mathbf{x}_m) \right\|^2 + c,$$

where

$$\alpha_i = \frac{w_i}{\sum_{i=1}^k w_i},$$

and c is a constant that does not depend on \mathbf{x} . Thus, the MM update boils down to projection onto S of a convex combination of the projections of the previous iterate onto the sets C_i ; in symbols

$$\mathbf{x}_{m+1} = P_S \left[\sum_i \alpha_i P_{C_i}(\mathbf{x}_m) \right]. \quad (11)$$

The majorization (10) involves dividing by 0 when \mathbf{x}_m belongs to C_i . This singularity also bedevils Weiszfeld's algorithm. Fortunately, perturbation of the objective function salvages the situation. One simply replaces the function $D(\mathbf{x})$ by the related function

$$D_\epsilon(\mathbf{x}) = \sum_{j=1}^k \gamma_j \sqrt{d(\mathbf{x}, C_j)^2 + \epsilon}$$

for ϵ small and positive. Ben-Tal and Teboulle [2] cover further examples of this perturbation strategy. In any case observe that the smooth function $f_\epsilon(u) = \sqrt{u^2 + \epsilon}$ has derivatives

$$f'_\epsilon(u) = \frac{u}{\sqrt{u^2 + \epsilon}}, \quad f''_\epsilon(u) = \frac{\epsilon}{(u^2 + \epsilon)^{3/2}}$$

and is therefore strictly increasing and strictly convex on the interval $[0, \infty)$. Hence, the function $D_\epsilon(\mathbf{x})$ is also convex. Because $\sqrt{u^2 + \epsilon} - \sqrt{\epsilon}$ is a good approximation to $u \geq 0$, the solutions of the two problems should be close. In fact, we will show later that the minimum point of $D_\epsilon(\mathbf{x})$ tends to the minimum point of $D(\mathbf{x})$ as ϵ tends to 0. In the presence of multiple minima, this claim must be rephrased in terms of cluster points.

The majorization $d(\mathbf{x}, C_j) \leq \|\mathbf{x} - P_{C_j}(\mathbf{x}_m)\|$ around the current iterate \mathbf{x}_m yields the majorization

$$\sqrt{d(\mathbf{x}, C_j)^2 + \epsilon} \leq \sqrt{\|\mathbf{x} - P_{C_j}(\mathbf{x}_m)\|^2 + \epsilon}.$$

Application of the majorization (9) implies the further majorization

$$D_\epsilon(\mathbf{x}) \leq \frac{1}{2} \sum_{j=1}^k \gamma_j \frac{\|\mathbf{x} - P_{C_j}(\mathbf{x}_m)\|^2}{\sqrt{\|\mathbf{x}_m - P_{C_j}(\mathbf{x}_m)\|^2 + \epsilon}} + c,$$

where c is an irrelevant constant. The corresponding MM update \mathbf{x}_{m+1} is identical to the previous MM update (11) except for one difference. The weights w_i are now defined by the benign formula

$$w_i = \frac{\gamma_i}{\sqrt{\|\mathbf{x}_m - P_{C_i}(\mathbf{x}_m)\|^2 + \epsilon}}$$

involving no singularity.

5 Examples.

We now consider four examples illustrating the performance of the MM algorithm and framing our expectations for convergence. The projected sub-gradient algorithm [21] serves as a benchmark for comparison in the first example. This algorithm relies on the updates

$$\mathbf{x}_{m+1} = P_S \left[\mathbf{x}_m - \eta_m \sum_{i=1}^k \mathbf{v}_{im} \right],$$

where

$$\mathbf{v}_{im} = \begin{cases} \frac{\mathbf{x}_m - P_{C_i}(\mathbf{x}_m)}{d(\mathbf{x}_m, C_i)} & \text{if } \mathbf{x}_m \notin C_i \\ 0 & \text{if } \mathbf{x}_m \in C_i, \end{cases}$$

and the nonnegative constants η_m satisfy $\sum_{m=1}^{\infty} \eta_m = \infty$ and $\sum_{m=1}^{\infty} \eta_m^2 < \infty$.

Iteration	x_1	x_2	x_3
1	0.000000000000000	2.000000000000000	0.000000000000000
2	-0.93546738305698	1.66164748416805	0.10207032020482
3	-0.92881282698649	1.63915389878166	0.08424264751830
4	-0.92645373003448	1.63220797263449	0.08007815377225
5	-0.92567602259658	1.63004821970935	0.07911751670489
6	-0.92542515217106	1.62937435413374	0.07889815178685
7	-0.92534495711879	1.62916364685109	0.07884864943702
8	-0.92531944712805	1.62909766226627	0.07883765997470
9	-0.92531135783449	1.62907697582185	0.07883527888603
10	-0.92530879826106	1.62907048520349	0.07883478238381
20	-0.92530761702316	1.62906751412014	0.07883466748783
30	-0.92530761701184	1.62906751409212	0.07883466748878
50	-0.92530761701184	1.62906751409212	0.07883466748878

Table 1: Cubes and ball example in \mathbb{R}^3 : MM Algorithm.

5.1 Five Cubes and a Ball in \mathbb{R}^3 .

Our first example is taken from the reference [21]. This three-dimensional example involves five cubes C_i with side lengths equal to 2 and centers $(0, -4, 0)$, $(-4, 2, -3)$, $(-3, -4, 2)$, $(-5, 4, 4)$, and $(-1, 8, 1)$. The set S is a ball with center $(0, 2, 0)$ and radius 1. Iteration commences at the point $\mathbf{x}_1 = (0, 2, 0) \in S$ and takes projected subgradient steps with $\eta_m = 1/m$. Table 1 shows the MM iterates with $\epsilon = 0$. Convergence to machine precision occurs within 30 iterations. In contrast Table 2 shows that parameter values (x_1, x_2, x_3) are still changing after 10^6 projected subgradient iterates. For brevity we omit a second example of four squares and a disk in \mathbb{R}^2 from the same source [21]. In this example the superiority of the MM algorithm over the projected subgradient algorithm is equally evident.

Iteration	x_1	x_2	x_3
1	0.000000000000000	2.000000000000000	0.000000000000000
10	-0.92583298353433	1.63051788239768	0.07947484741743
100	-0.92531325048300	1.62908232435160	0.07883822912883
1000	-0.92530767419684	1.62906766065418	0.07883468589312
10000	-0.92530761758555	1.62906751554109	0.07883466757273
100000	-0.92530761701755	1.62906751410641	0.07883466748904
1000000	-0.92530761701233	1.62906751409334	0.07883466748881
1500000	-0.92530761701231	1.62906751409328	0.07883466748881
2000000	-0.92530761701229	1.62906751409324	0.07883466748881

Table 2: Cubes and ball example in \mathbb{R}^3 : Projected Subgradient Algorithm.

5.2 The Closest Point to Three Disks in \mathbb{R}^2 .

This example from the reference [20] illustrates the advantage of minimizing a sequence of approximating functions $D_{\epsilon_m}(\mathbf{x})$. The sets C_i are three unit balls in \mathbb{R}^2 centered at $(0, 2)$, $(2, 0)$, and $(-2, 0)$. The set S equals \mathbb{R}^2 . By inspection the minimum distance occurs at $(0, 1)$. Figure 1 displays the iteration paths for 50 different starting values (dots) and their corresponding fixed point (the square). Along the m th leg of the path we set ϵ_m to be $\max\{10^m, 10^{-16}\}$. The solution to the current problem is taken as the initial point for the next problem. All solution paths initially converge to a point just below $(0, 1)$ and then march collectively upwards to $(0, 1)$. The passage of the MM iterates through the unit balls is facilitated by our strategy of systematically reducing ϵ .

5.3 Three Collinear Disks in \mathbb{R}^2 .

Here we illustrate the behavior of the MM algorithm when there is more than one solution. Consider two unit balls in \mathbb{R}^2 centered at $(2, 0)$, and $(-2, 0)$, and take S to be the unit ball centered at the origin. By inspection there is a continuum of solutions extending along the line segment from $(-1, 0)$ to $(1, 0)$. Figure 2 shows the iteration paths for 100 different initial values (dots) and their corresponding fixed points (squares). In this example we take $\epsilon = 0$. Although the iterates are not guaranteed to converge and may in principle cycle among multiple cluster points, this behavior is not observed

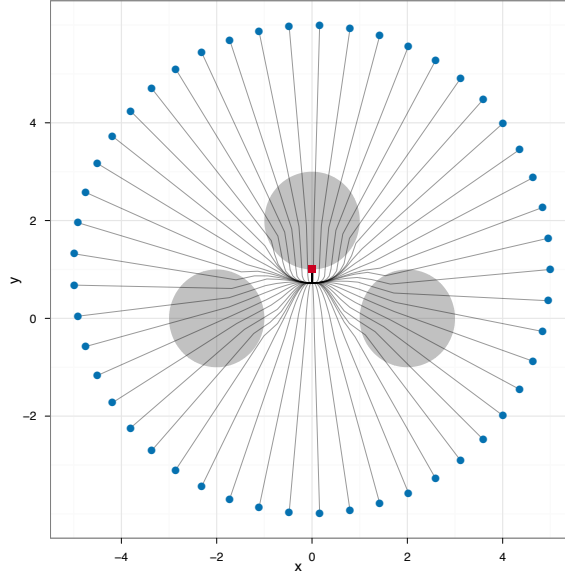


Figure 1: Finding the closest point to three disks in \mathbb{R}^2 .

in practice. The iterates simply converge to different fixed points depending on where they start.

5.4 Kuhn's Problem.

Our last example was originally concocted by Kuhn [14] to illustrate how Weiszfeld's algorithm can stall when its iterates enter one of the sets C_i . Although this event rarely occurs in practice, characterizing the initial conditions under which it happens has been a subject of intense scrutiny [5, 6, 7, 8, 15]. The occasional failure of Weiszfeld's algorithm prompted Vardi and Zhang [25] to redesign it. Their version preserves the descent property but differs substantially from ours. In any event the example shown in Figure 3 involves two points with weights γ_i proportional to 5 placed at (59,0) and (20,0) and two more points with weights proportional to 13 placed at (-20, 48) and (-20, -48). The optimal point is the origin. Starting at (44,0), Weiszfeld's algorithm stalls at (20,0) after one iteration. Our MM iterates (dots) with ϵ decreasing from 0.1 to 0, in contrast, move across (20,0) and correctly converge to (0,0) in 99 steps.

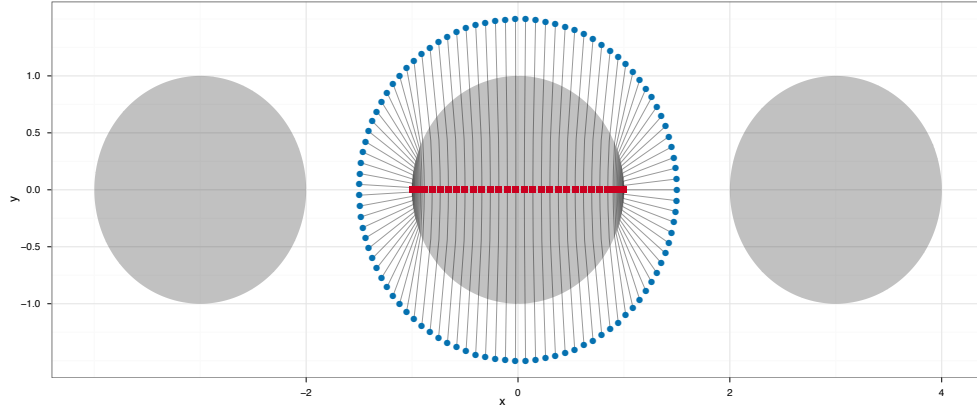


Figure 2: An example with a continuum of solutions.

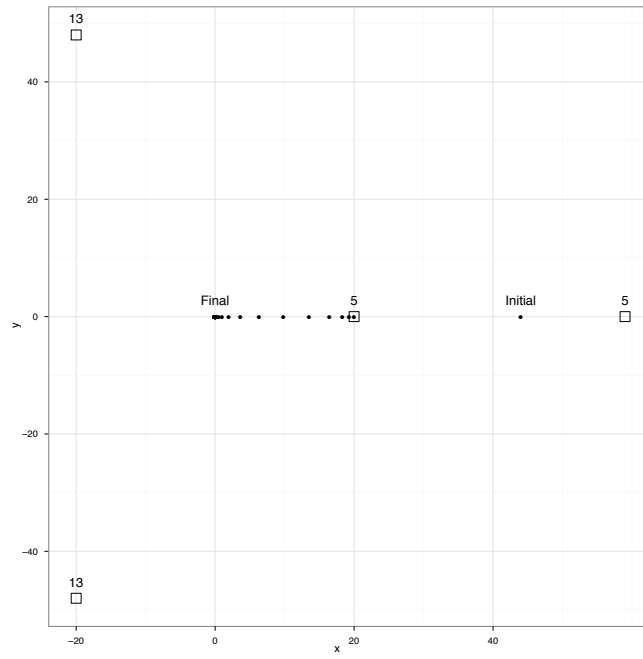


Figure 3: A problem where Weiszfeld's algorithm fails to converge.

6 Convergence Theory.

Before embarking on a proof of convergence, it is prudent to discuss whether a minimum point exists and is unique. Recall that a continuous function attains its minimum on a compact set. Thus, problem (8) possesses a minimum whenever S is bounded. If S is unbounded, then one can substitute boundedness of one or more of the sets C_i . In this circumstance $D(\mathbf{x})$ is coercive in the sense that $\lim_{\|\mathbf{x}\| \rightarrow \infty} D(\mathbf{x}) = \infty$. As pointed out in Proposition 3.1 of the reference [21], coerciveness is sufficient to guarantee existence. Because $D(\mathbf{x}) \leq D_\epsilon(\mathbf{x})$, the perturbed criterion $D_\epsilon(\mathbf{x})$ is coercive whenever the original criterion $D(\mathbf{x})$ is coercive. Henceforth, we will assume that S or at least one of the C_i is bounded.

A strictly convex function possesses at most one minimum point on a convex set. The function $|x|$ shows that this sufficient condition for uniqueness is hardly necessary. In the Fermat-Weber problem, where the closed convex sets $C_i = \{\mathbf{x}_i\}$ are singletons, the function $D(\mathbf{x})$ is strictly convex if and only if the points \mathbf{x}_i are non-collinear. To generalize this result, we require the sets C_i to be non-collinear. Geometrically this says that it is impossible to draw a straight line that passes through all of the C_i . Non-collinearity can only be achieved when $k > 2$ and $\cap_{i=1}^k C_i = \emptyset$. We also require the C_i to be strictly convex. A set C is said to be strictly convex if the interior of the line segment $[\mathbf{x}, \mathbf{y}]$ connecting two different points \mathbf{x} and \mathbf{y} of C lies in the interior of C . Put another way, the boundary of C can contain no line segments. A singleton or a closed ball is strictly convex, but a closed box is not.

Proposition 6.1. *If the closed convex sets C_1, \dots, C_k are strictly convex but not collinear, then $D(\mathbf{x})$ is strictly convex.*

Proof. Suppose the contrary, and choose $\mathbf{x} \neq \mathbf{y}$ and α strictly between 0 and 1 so that

$$D[\alpha\mathbf{x} + (1 - \alpha)\mathbf{y}] = \alpha D(\mathbf{x}) + (1 - \alpha)D(\mathbf{y}). \quad (12)$$

Let L be the line $\{s\mathbf{x} + (1 - s)\mathbf{y} : s \in \mathbb{R}\}$ passing through the points \mathbf{x} and \mathbf{y} . Then there exists at least one C_j such that $L \cap C_j = \emptyset$. In particular, \mathbf{x} , \mathbf{y} , and $\alpha\mathbf{x} + (1 - \alpha)\mathbf{y}$ all fall outside this C_j . Equality (12) implies that

$$\begin{aligned} & \alpha\|\mathbf{x} - P_{C_j}(\mathbf{x})\| + (1 - \alpha)\|\mathbf{y} - P_{C_j}(\mathbf{y})\| \\ &= \|\alpha\mathbf{x} + (1 - \alpha)\mathbf{y} - P_{C_j}[\alpha\mathbf{x} + (1 - \alpha)\mathbf{y}]\| \\ &\leq \|\alpha\mathbf{x} + (1 - \alpha)\mathbf{y} - \alpha P_{C_j}(\mathbf{x}) - (1 - \alpha)P_{C_j}(\mathbf{y})\| \\ &\leq \alpha\|\mathbf{x} - P_{C_j}(\mathbf{x})\| + (1 - \alpha)\|\mathbf{y} - P_{C_j}(\mathbf{y})\|. \end{aligned}$$

Since the projection of a point onto C_j is unique, these sandwich inequalities entail

$$P_{C_j}[\alpha \mathbf{x} + (1 - \alpha) \mathbf{y}] = \alpha P_{C_j}(\mathbf{x}) + (1 - \alpha) P_{C_j}(\mathbf{y}).$$

If $P_{C_j}(\mathbf{x}) \neq P_{C_j}(\mathbf{y})$, then the strict convexity of C_j implies the convex combination $\alpha P_{C_j}(\mathbf{x}) + (1 - \alpha) P_{C_j}(\mathbf{y})$ is interior to C_j . Hence, this point cannot be the closest point to the external point $\alpha \mathbf{x} + (1 - \alpha) \mathbf{y}$. Therefore, consider the possibility $P_{C_j}(\mathbf{x}) = P_{C_j}(\mathbf{y}) = \mathbf{z}$. Equality can occur in the inequality

$$\|\alpha \mathbf{x} + (1 - \alpha) \mathbf{y} - \mathbf{z}\| \leq \alpha \|\mathbf{x} - \mathbf{z}\| + (1 - \alpha) \|\mathbf{y} - \mathbf{z}\|$$

only when $\mathbf{x} - \mathbf{z} = t(\mathbf{y} - \mathbf{z})$ for some $t \neq 1$. This relation shows that

$$\mathbf{z} = \frac{1}{1 - t} \mathbf{x} - \frac{t}{1 - t} \mathbf{y}$$

belongs to $L \cap C_j$, contradicting our hypothesis. Thus, $D(\mathbf{x})$ is strictly convex. \square

The next result shows that the function $D_\epsilon(\mathbf{x})$ inherits strict convexity from $D(\mathbf{x})$. Therefore, when $D(\mathbf{x})$ is strictly convex, $D_\epsilon(\mathbf{x})$ possesses a unique minimum point.

Proposition 6.2. *If $D(\mathbf{x})$ is strictly convex, then $D_\epsilon(\mathbf{x})$ is also strictly convex.*

Proof. Fix arbitrary $\mathbf{x} \neq \mathbf{y}$ and α strictly between 0 and 1. The strict convexity of $D(\mathbf{x})$ implies that there is at least one j such that

$$d(\alpha \mathbf{x} + (1 - \alpha) \mathbf{y}, C_j) < \alpha d(\mathbf{x}, C_j) + (1 - \alpha) d(\mathbf{y}, C_j).$$

The strict inequality

$$\begin{aligned} \sqrt{d(\alpha \mathbf{x} + (1 - \alpha) \mathbf{y}, C_j)^2 + \epsilon} &< \sqrt{[\alpha d(\mathbf{x}, C_j) + (1 - \alpha) d(\mathbf{y}, C_j)]^2 + \epsilon}, \\ &\leq \alpha \sqrt{d(\mathbf{x}, C_j)^2 + \epsilon} + (1 - \alpha) \sqrt{d(\mathbf{y}, C_j)^2 + \epsilon}, \end{aligned}$$

follows because the function $f_\epsilon(u) = \sqrt{u^2 + \epsilon}$ is a strictly increasing and convex. Summing over j gives the desired result. \square

We now clarify the relationship between the minima of the $D_\epsilon(\mathbf{x})$ and $D(\mathbf{x})$.

Proposition 6.3. *For a sequence of constants ϵ_m tending to 0, let \mathbf{y}_m be a corresponding sequence minimizing $D_{\epsilon_m}(\mathbf{x})$. If \mathbf{y} is the unique minimum point of $D(\mathbf{x})$, then \mathbf{y}_m tends to \mathbf{y} . If $D(\mathbf{x})$ has multiple minima, then every cluster point of the sequence \mathbf{y}_m minimizes $D(\mathbf{x})$.*

Proof. To prove the assertion, consider the inequalities

$$D(\mathbf{y}_m) \leq D_{\epsilon_m}(\mathbf{y}_m) \leq D_{\epsilon_m}(\mathbf{x}) \leq D_1(\mathbf{x})$$

for any $\mathbf{x} \in S$ and $\epsilon_m \leq 1$. Taking limits along the appropriate subsequences proves that the cluster points of the sequence \mathbf{y}_m minimize $D(\mathbf{x})$. Convergence to a unique minimum point \mathbf{y} occurs provided the sequence \mathbf{y}_m is bounded. If S is bounded, then \mathbf{y}_m is bounded by definition. On the other hand, if any C_j is bounded, then $D(\mathbf{x})$ is coercive, and the inequality $D(\mathbf{y}_m) \leq D_1(\mathbf{x})$ forces \mathbf{y}_m to be bounded. \square

The convergence theory of MM algorithms hinges on the properties of the algorithm map $\psi(\mathbf{x}) \equiv \arg \min_{\mathbf{y}} g(\mathbf{y} | \mathbf{x})$. For easy reference, we state a simple version of Zangwill's convergence theorem [29, p. 91] instrumental in proving convergence in our setting.

Proposition 6.4. *Let $f(\mathbf{x})$ be a continuous function on a domain S and $\psi(\mathbf{x})$ be a continuous algorithm map from S into S satisfying $f(\psi(\mathbf{x})) < f(\mathbf{x})$ for all $\mathbf{x} \in S$ with $\psi(\mathbf{x}) \neq \mathbf{x}$. Suppose for some initial point \mathbf{x}_0 that the set $\mathcal{L}_f(\mathbf{x}_0) \equiv \{\mathbf{x} \in S : f(\mathbf{x}) \leq f(\mathbf{x}_0)\}$ is compact. Then (a) the sequence of iterates $\mathbf{x}_{m+1} = \psi(\mathbf{x}_m)$ has at least one cluster point, (b) all cluster points are fixed points of $\psi(\mathbf{x})$, and (c) $\lim_{m \rightarrow \infty} \|\mathbf{x}_{m+1} - \mathbf{x}_m\| = 0$.*

The convergence of the MM iterates (11) to a stationary point of $f(\mathbf{x})$ follows immediately from Proposition 6.4 provided the fixed points of $\psi(\mathbf{x})$ are stationary points of $f(\mathbf{x})$ and $\psi(\mathbf{x})$ possesses only finitely many fixed points.

Let us verify the conditions of Proposition 6.4 for minimizing $D_\epsilon(\mathbf{x})$. The function $D_\epsilon(\mathbf{x})$ is continuous on its domain S , and the set $\mathcal{L}_{D_\epsilon}(\mathbf{x}_0)$ is compact for any initial point \mathbf{x}_0 since either S is compact or $D_\epsilon(\mathbf{x})$ is coercive. The continuity of the algorithm map follows immediately from the continuity of the projection mapping. Finally, we need to establish that $D_\epsilon(\psi(\mathbf{x})) < D_\epsilon(\mathbf{x})$ whenever $\mathbf{x} \neq \psi(\mathbf{x})$. First observe that $\psi(\mathbf{x}) = \mathbf{x}$ if and only if the MM surrogate function satisfies $g_\epsilon(\mathbf{x} | \mathbf{x}) = \min_{\mathbf{y}} g_\epsilon(\mathbf{y} | \mathbf{x})$. Therefore, we have the strict inequality $g_\epsilon(\psi(\mathbf{x}) | \mathbf{x}) < g_\epsilon(\mathbf{x} | \mathbf{x})$ whenever \mathbf{x} is not a fixed point

of ψ . This forces a decrease in the objective function $D_\epsilon(\mathbf{x})$ and makes the MM algorithm strictly monotone outside the set of stationary points.

We now argue that the fixed points of the algorithm map $\psi(\mathbf{x})$ are stationary points of $D_\epsilon(\mathbf{x})$. We will show, in fact, that the two sets of points coincide. To accomplish this, we need to determine the gradients of $D_\epsilon(\mathbf{x})$ and $g_\epsilon(\mathbf{x} \mid \mathbf{y})$. Recall that $f_\epsilon(u)$ is strictly increasing and strictly convex. As a consequence the functions $f_\epsilon(\|\mathbf{x}\|)$ and $f_\epsilon[d(\mathbf{x}, C_j)]$ are convex. Even more remarkable is the fact that both functions are continuously differentiable. When $\mathbf{x} \neq \mathbf{0}$, the function $\|\mathbf{x}\|$ is differentiable. Likewise, when $\mathbf{x} \notin C_j$, the function $d(\mathbf{x}, C_j)$ is differentiable. Therefore, the chain rule implies

$$\nabla f_\epsilon(\|\mathbf{x}\|) = \frac{\|\mathbf{x}\|}{\sqrt{\|\mathbf{x}\|^2 + \epsilon}} \frac{\mathbf{x}}{\|\mathbf{x}\|} = \frac{\mathbf{x}}{\|\mathbf{x}\|^2 + \epsilon} \quad (13)$$

$$\nabla f_\epsilon[d(\mathbf{x}, C_j)] = \frac{d(\mathbf{x}, C_j)}{\sqrt{d(\mathbf{x}, C_j)^2 + \epsilon}} \frac{\mathbf{x} - P_{C_j}(\mathbf{x})}{d(\mathbf{x}, C_j)} = \frac{\mathbf{x} - P_{C_j}(\mathbf{x})}{\sqrt{d(\mathbf{x}, C_j)^2 + \epsilon}}, \quad (14)$$

respectively.

By continuity one expects the gradients to be defined for $\mathbf{x} = \mathbf{0}$ and $\mathbf{x} \in C_j$ by the corresponding limit of $\mathbf{0}$. In the former case the expansion

$$\sqrt{\|\mathbf{x}\|^2 + \epsilon} - \sqrt{\epsilon} = \sqrt{\epsilon} \sqrt{1 + \frac{\|\mathbf{x}\|^2}{\epsilon}} - \sqrt{\epsilon} = \frac{1}{2} \frac{\|\mathbf{x}\|^2}{\sqrt{\epsilon}} + \sqrt{\epsilon} o\left(\frac{\|\mathbf{x}\|^2}{\epsilon}\right).$$

shows that $\nabla f_\epsilon(\|\mathbf{0}\|) = \mathbf{0}$. In the latter case the expansion

$$\sqrt{d(\mathbf{y}, C_j)^2 + \epsilon} - \sqrt{\epsilon} = \frac{1}{2} \frac{d(\mathbf{y}, C_j)^2}{\sqrt{\epsilon}} + \sqrt{\epsilon} o\left[\frac{d(\mathbf{y}, C_j)^2}{\epsilon}\right]$$

and the bound $d(\mathbf{y}, C_j) = |d(\mathbf{y}, C_j) - d(\mathbf{x}, C_j)| \leq \|\mathbf{y} - \mathbf{x}\|$ for $\mathbf{x} \in C_j$ likewise show that $\nabla f_\epsilon[d(\mathbf{x}, C_j)] = \mathbf{0}$. Consequently, equations (13) and (14) hold for all $\mathbf{x} \in \mathbb{R}^d$. It follows that both $D_\epsilon(\mathbf{x})$ and $g_\epsilon(\mathbf{x} \mid \mathbf{y})$ are differentiable on \mathbb{R}^d , with gradients

$$\nabla D_\epsilon(\mathbf{x}) = \sum_{j=1}^k \gamma_j \frac{\mathbf{x} - P_{C_j}(\mathbf{x})}{\sqrt{d(\mathbf{x}, C_j)^2 + \epsilon}}, \quad (15)$$

and

$$\nabla g_\epsilon(\mathbf{x} \mid \mathbf{y}) = \sum_{j=1}^k \gamma_j \frac{\mathbf{x} - P_{C_j}(\mathbf{y})}{\sqrt{d(\mathbf{y}, C_j)^2 + \epsilon}}, \quad (16)$$

respectively. Note that $\mathbf{y} \in S$ minimizes $D_\epsilon(\mathbf{x})$ over S if and only if

$$\sum_{j=1}^k \gamma_j \frac{\langle \mathbf{y} - P_{C_j}(\mathbf{y}), \mathbf{x} - \mathbf{y} \rangle}{\sqrt{d(\mathbf{y}, C_j)^2 + \epsilon}} \geq 0,$$

for all $\mathbf{x} \in S$. This inequality, however, is equivalent to the inequality $\langle \nabla g_\epsilon(\mathbf{y} \mid \mathbf{y}), \mathbf{x} - \mathbf{y} \rangle \geq 0$, for all $\mathbf{x} \in S$, which in turn holds if and only if \mathbf{y} is a fixed point of $\psi(\mathbf{x})$. If $D(\mathbf{x})$ is strictly convex, then $D_\epsilon(\mathbf{x})$ has a unique minimum point, and $\psi(\mathbf{x})$ has exactly one fixed point.

Thus, Proposition 6.3 and Proposition 6.4 together tell us that \mathbf{y} is a solution to (8) if there is a sequence of ϵ_m tending to zero and a sequence of points \mathbf{y}_m tending to \mathbf{y} that satisfy

$$\left\langle -\sum_{j=1}^k \gamma_j \frac{\mathbf{y}_m - P_{C_j}(\mathbf{y}_m)}{\sqrt{d(\mathbf{y}_m, C_j)^2 + \epsilon_m}}, \mathbf{x} - \mathbf{y}_m \right\rangle \leq 0,$$

for all $\mathbf{x} \in S$. The above sufficient condition becomes necessary as well if $D(\mathbf{x})$ is strictly convex. When the sets $S \cap C_j$ are all empty and the weights γ_j are identical, we recover the characterization of the optimal points given in Theorem 3.2 of reference [21].

Let us close with a few words of praise for the MM principle. It lit the way to an efficient numerical algorithm for solving problem (8) using only elementary principles of convex analysis. It also suggested how to derive an MM algorithm that removes the singularities of Weiszfeld's algorithm. Finally, it clarified the optimality conditions derived by Mordukhovich et al. [21]. Similar advantages accrue across a broad spectrum of optimization problems. The ability of MM algorithms to handle high-dimensional problems in imaging, genomics, statistics, and a host of other fields testifies to the potency of a simple idea consistently invoked. Mathematical scientists are well advised to be on the lookout for new applications.

Acknowledgments.

This research was supported by the United States Public Health Service grant GM53275.

References

- [1] M. P. Becker, I. Yang, and K. Lange, EM algorithms without missing data, *Statistical Methods in Medical Research* **6** (1997) 38–54.
- [2] A. Ben-Tal and M. Teboulle, A smoothing technique for nondifferentiable optimization problems, in *Optimization*, S. Dolecki, ed., Lecture Notes in Mathematics, vol. 1405, Springer Berlin / Heidelberg, 1989, 1–11.
- [3] D. P. Bertsekas, *Convex Optimization Theory*, Athena Scientific, Belmont, MA, 2009.
- [4] J. M. Borwein and A. S. Lewis, *Convex Analysis and Nonlinear Optimization: Theory and Examples*, Springer, New York, 2000.
- [5] J. Brimberg, The Fermat-Weber location problem revisited, *Mathematical Programming* **71** (1995) 71–76.
- [6] ———, Further notes on convergence of the Weiszfeld algorithm, *Yugoslav Journal of Operations Research* **13** (2003) 199–206.
- [7] L. Cánovas, A. Marín, and R. Cañflavate, On the convergence of the Weiszfeld algorithm, *Mathematical Programming* **93** (2002) 327–330.
- [8] R. Chandrasekaran and A. Tamir, Open questions concerning Weiszfeld’s algorithm for the Fermat-Weber location problem, *Mathematical Programming* **44** (1989) 293–295.
- [9] R. Courant and H. Robbins, *What is Mathematics?*, Oxford University Press, New York, 1961.
- [10] J. Duchi, S. Shalev-Shwartz, Y. Singer, and T. Chandra, Efficient projections onto the ℓ_1 -ball for learning in high dimensions, in *Proceedings of the International Conference on Machine Learning*, 2008.
- [11] S. Gueron and R. Tessler, The Fermat-Steiner problem, *The American Mathematical Monthly* **109** (2002) 443–451.
- [12] J.-B. Hiriart-Urruty and C. Lemaréchal, *Fundamentals of Convex Analysis*, Springer, 2004.

- [13] J. Krarup and S. Vajda, On Torricelli's geometrical solution to a problem of Fermat, *IMA Journal of Mathematics Applied in Business and Industry* **8** (1997) 215–224.
- [14] H. W. Kuhn, On a pair of dual nonlinear programs, in *Nonlinear Programming*, J. Abadie, ed., North-Holland Publishing, Amsterdam, The Netherlands, 1967, 37–54.
- [15] ———, A note on Fermat's problem, *Mathematical Programming* **4** (1973) 98–107.
- [16] K. Lange, *Numerical Analysis for Statisticians*, 2nd ed., Springer, New York, 2010.
- [17] K. Lange, D. R. Hunter, and I. Yang, Optimization transfer using surrogate objective functions (with discussion), *Journal of Computational and Graphical Statistics* **9** (2000) 1–20.
- [18] R. F. Love, J. G. Morris, and G. O. Wesolowsky, *Facilities Location: Models and Methods*, Appleton and Lange, North-Holland, 1988.
- [19] C. Michelot, A finite algorithm for finding the projection of a point onto the canonical simplex of \mathbb{R}^n , *Journal of Optimization Theory and Applications* **50** (1986) 195–200.
- [20] B. Mordukhovich and N. M. Nam, Applications of variational analysis to a generalized Fermat-Torricelli problem, *Journal of Optimization Theory and Applications* **148** (2011) 431–454.
- [21] B. Mordukhovich, N. M. Nam, and J. Salinas, Solving a generalized Heron problem by means of convex analysis, *The American Mathematical Monthly* **119** (2012) 87–99.
- [22] J. M. Ortega and W. C. Rheinboldt, *Iterative Solutions of Nonlinear Equations in Several Variables*, Academic, New York, 1970.
- [23] R. T. Rockafellar, *Convex Analysis*, Princeton University Press, Princeton, NJ, 1996.
- [24] A. Ruszczyński, *Nonlinear Optimization*, Princeton University Press, Princeton, NJ, 2006.

- [25] Y. Vardi and C.-H. Zhang, A modified Weiszfeld algorithm for the Fermat-Weber location problem, *Mathematical Programming, Series A* **90** (2001) 559–566.
- [26] E. Weiszfeld, Sur le point pour lequel la somme des distances de n points donnés est minimum, *Tôhoku Mathematics Journal* **43** (1937) 355–386.
- [27] E. Weiszfeld and F. Plastria, On the point for which the sum of the distances to n given points is minimum, *Annals of Operations Research* **167** (2009) 7–41.
- [28] G. O. Wesolowsky, The Weber problem: Its history and perspectives, *Location Science* **1** (1993) 5–23.
- [29] W. I. Zangwill, *Nonlinear Programming: A Unified Approach*, International Series in Management, Prentice-Hall, Englewood Cliffs, New Jersey, 1969.

Department of Human Genetics, University of California, Los Angeles, CA 90095
ecchi@ucla.edu

Departments of Human Genetics, Biomathematics, and Statistics, University of California, Los Angeles, CA 90095
klange@ucla.edu