# 1. Parameter Optimization Problems:

## 1.1. Unconstrained Parameter Optimization

©Erik I. Verriest*

November 19, 2008

# 1   Central Theme of Optimization Problems

There are five central questions for any optimization problem:

1. *Existence*: Do there exist solutions to the problem? If the answer is negative, no further effort should be spent. The question of existence is therefore important, and not just a technicality for mathematicians.

2. *Uniqueness*: If a solution exists, is it the onbly one, or are there more solutions?

3. *Sufficient Conditions*: When are we certain that a candidate solution is indeed a solution?

4. *Necessary Conditions*: What clues do we have to find a solution? If we know what conditions are necessary, then we now that if a candidate does not satisfy this condition, it cannot be a solution to the problem.

5. *Computational Methods*: What are efficient algorithms to solve a problem?

**Remark**: Computational methods are always *finite* methods, involving algebraic operations. It is well known that typically such methods will yield solutions. Think for instance about solving ordinary differential equations via discretization. Such a "solution" can always be constructed even if the original differential equation can be *proven* to have no solution! Therefore, the questions of solvability, stability, and necessary and sufficient conditions are not just mathematical oddities or purely academic problems, but are indeed very important, and should form an integral part of every analysis and design. Too often a blind belief in

---

*School of ECE, Georgia Tech, Atlanta, GA, USA. Presently with SCD, ESAT, KULeuven.

computer power has led to catastrophic failures!

The simplest optimization problems are the problems of extremizing smooth functions of several variables. Several techniques are available.

i) Geometric Methods. These involve the optimization of the Euclidean distances. Recall that the shortest distance between two points is the length of a straight line segment between these points, many problems can be reduced to this. The parametrization may enter in the choice of intermediate points.

ii) Algebraic Methods. This involves in particular the use of inequalities. If one can show that the objective function is bounded by a simpler expression, and if it can be shown that the bound is actually achievable, then the optimum is readily found.

iii) Analytic Methods. If the objective function is a smooth function of the parameters, then the extrema are determined by setting the partial derivatives equal to zero. Special care must be taken if the parameter domain is bounded.

iv) Numerical Methods. Basic methods are the gradient method (steepest descent) and Newton's method.

# 2   Geometric Methods

Euclidean geometry has a wealth of simple (and not so simple) distance minimization problems. A basic one is the following:

**Problem 1:** *Let A and B be two points on the same side of a line $\ell$. Find a point $C \in \ell$ such that the sum of the distances, $d(A, C) + d(C, B)$, is minimal.*

This problem was first posed by Heron of Alexandria, first century A.D. Mirror symmetry is the technique that leads to the quick solution of this problem.

Here are some additional problems:

**Problem 2:** *Let C be a given point in the interior of a given angle aOb. Find a point A on Oa and B on Ob such that the perimeter of the triangle ABC is a minimum.*
Be sure to look at all possible cases for the angle aOb from 0 to $2\pi$.

**Problem 3:** *Given an angle aOb and two points, C and D, in its interior, determine a point A on Oa and B on Ob such that $d(A, B) + d(B, C) + d(C, D) + d(D, A)$ is a minimum.*

**Problem 4:** *Given three points A,B,C, find the point D such that the sum of the distances $d(D, A) + d(D, B) + d(D, C)$ is a minimum.* This problem is known as *Steiner's problem.*

Its solution is known as the Torricelli point of the triangle

**Problem 5:** *Solve Problem 4 for four points.*

Note that all these problems have practical relevance. For instance Problems 4 and 5 ask for the minimal length of freeway (cables, pipelines or canals) linking three or four towns.

**Problem 6:** *Will the answer to Problem 5 change if one does not require that the highways intersect in one point?*

**Problem 7:** This is a "dual" to Steiner's problem: *Find a point A from which the sum of the distances to three given lines, a, b, and c, is a minimum.*

We do not have to stick to problems in the plane: One may ask for instance for the shortest distance between two points on a cube, a pyramid, etc. These are easily solved by "unfolding and flattening". A typical application is the following problem.

**Problem 8:** *Find the path for the minimal length of cable between a provider P and user U, for a given mountainous terrain.*

Here is a problem from the first textbook known to mankind: Euclid's *Elements*, written in the 4-th century B.C.

**Problem 9** (Euclid's Problem). *In a given triangle ABC, inscribe a parallelogram ADEF (EF//AB, DE//AC) of maximal area.*

Finally, we state also the classical jewel: *Dido's problem*, which while not a minimum distance problem, it still is amenable to solution by geometric reasoning:

**Problem 10:** *Find the closed curve of fixed length for which the enclosed area is maximal.* This was the first extremal problem discussed in the scientific literature. It appeared in the ninth century B.C., in the Aeneid of Vergil, although stated somewhat differently:

> They bought as much land - and called it Birsa -as could be encircled by a bull's hide.

Its solution was surely well known to Aristotle (4-th century B.C.) The geometric solution, due to Steiner is sketched below:

0. The tacit assumption is that there exists such a curve. (This is justified).

1. The extremal curve is convex. Prove by contradiction.

2. If points A and B halve the length of the extremal curve, then the chord AB halves the area it encloses. Prove by contradiction.

3. Suppose that points A and B halve the extremal curve. If C is any point on the curve, then the angle ACB is a right angle.

# 3   Fundamental Inequalities

As mentioned, many extremal problems are concealed in various inequalities. An inequality is called *exact* if equality is actually attained. (e.g. $2 \leq 3$ is not exact, whereas $a^2 \geq 0$ is, since equality holds if $a = 0$).

## 3.1   Arithmetic-Geometric Means

Let $a$ and $b$ be nonnegative numbers. Their arithmetic mean is $\frac{a+b}{2}$, and the geometric mean $\sqrt{ab}$. The following inequality holds:

$$\sqrt{ab} \leq \frac{a+b}{2}. \tag{1}$$

This inequality is exact (equality holds when $a = b$.)

## 3.2   Arithmetic-Geometric Means (General Case)

For nonnegative numbers $x_1, x_2, \ldots, x_n$ we have the inequality

$$[x_1 \cdots x_n]^{1/n} \leq \frac{x_1 + \cdots x_n}{n} \tag{2}$$

Equality holds if all numbers are equal.

## 3.3   Inequality of Arithmetic-Quadratic Means

The quadratic mean of the numbers $x_1, \ldots, x_n$ is defined by $\sqrt{\frac{x_1^2 + \cdots + x_n^2}{n}}$. The following holds:

$$\frac{x_1 + \ldots x_n}{n} \leq \sqrt{\frac{x_1^2 + \cdots + x_n^2}{n}} \tag{3}$$

It is an equality if all numbers are equal. It follows from this and the previous inequality that

$$[x_1 \cdots x_n]^{1/n} \leq \sqrt{\frac{x_1^2 + \cdots + x_n^2}{n}} \tag{4}$$

## 3.4   Cauchy-Bunyakovskii Inequality

For arbitrary numbers $x_1, \ldots x_n, y_1, \ldots, y_n$, we have the inequality

$$x_1 y_1 + \ldots + x_n y_n \le (x_1^2 + \ldots + x_n^2)^{1/2}(y_1^2 + \ldots + y_n^2)^{1/2}. \tag{5}$$

This inequality is again exact: Equality holds for $x_1 = y_1$, $\ldots$, $x_n = y_n$.

## 3.5   Hölder Inequality

This is an important generalization of the Cauchy-Bunyakowskii inequality:
For nonnegative numbers $x_1, \ldots x_n, y_1, \ldots, y_n$, and for $p > 1$ and $q$ satisfying

$$\frac{1}{p} + \frac{1}{q} = 1$$

we have:

$$x_1 y_1 + \ldots + x_n y_n \le (x_1^p + \ldots + x_n^p)^{1/p}(y_1^q + \ldots + y_n^q)^{1/q}. \tag{6}$$

## 3.6   Matrix Inequalities

Some concepts from matrix theory, which have plenty of applications in systems theory, optimization, and numerical linear algebra, are introduced. The first section deals with the eigen decomposition of symmetric matrices and some related optimization problems. Some important matrix groups are discovered along the way.

A square matrix $A \in \mathbf{R}^{n \times n}$ is called symmetric if $A = A'$, and anti (or skew-) symmetric if $A = -A'$. Any matrix $P \in \mathbf{R}^{n \times n}$ can be decomposed in a symmetric and an anti (or skew-) symmetric part:

$$P = \frac{1}{2}(P + P') + \frac{1}{2}(P - P') = P_{\text{Sym}} + P_{\text{Asym}}. \tag{7}$$

Consider now the *quadratic form*, $x'Px = \sum_{i=1}^{n} x_i x_j P_{ij}$, where $x \in \mathbf{R}^n$, and $P \in \mathbf{R}^{n \times n}$. First, we show that the *symmetric part*, $P_{Sym}$, is the only quantity that matters. Indeed, since $x'Px$ is a scalar quantity, we have $x'Px = (x'Px)' = x'P'x$. Hence: $x'Px = x'P_{\text{Sym}}x$. Next we show that the eigenvalues of symmetric matrices have remarkable properties:

### 3.6.1   Eigen Decomposition

**Proposition:**   *Any real symmetric matrix $P \in \mathbf{R}^{n \times n}$ has $n$ real eigenvalues, and one can always find a set of $n$ mutually orthogonal eigenvectors.*

*Proof:* See class notes.

It follows from the proposition that a symmetric matrix can always be diagonalized by a *real* similarity transformation. Namely let the transformation matrix $Y$ be formed by putting $n$ mutually orthogonal eigenvectors (such a set can always be found) 'shoulder to shoulder'.

Now the eigenvectors can always be *normalized*. A matrix $X$ formed by putting $n$ mutually orthonormal vectors shoulder to shoulder has the property that

$$(X'X)_{ij} = x_i'x_j = \delta_{ij}, \qquad \text{or,} \qquad X'X = I \tag{8}$$

It then follows at once that $XX' = I$ as well, and hence $\det(X) = \pm 1$. The subset of matrices in $\mathbf{R}^{n \times n}$ with the property $X'X = I$ forms a *group* under composition (matrix product). It is a subgroup of the general linear group $Gl_n(\mathbf{R})$, called the *orthogonal* group, and denoted by $O_n$. Clearly, it has two disjoint components. (In order to talk about disjointness, one needs to introduce a distance function; it is easily checked that a distance, induced by an inner product $d(X, Y) = \|X - Y\| = \sqrt{\text{Tr}\,(X-Y)'(X-Y)}$ satisfies all needs.) The two components are the set of orthogonal matrices with unit determinant, and the set of orthogonal matrices with determinant $-1$. The first set forms itself a group under composition, and is therefore a *sub*group of $O_n$, called the *special orthogonal group*, and denoted by $SO_n$. The group $SO_n$ is also a subgroup of the *special linear group*, denoted $Sl_n(\mathbf{R})$, which is itself a subgroup of $Gl_n(\mathbf{R})$, consisting of the matrices with determinant equal to 1. The group $SO_n$ correponds to the pure rotations. Orthogonal matrices with determinant $-1$ are also called improper rotations. A determinant $-1$ is associated geometrically with a *reflection*. We summarize the above in the very important

**Theorem 3.6.1:** Eigen decomposition of a symmetric matrix
*Any symmetric matrix $S \in \mathbf{R}^{n \times n}$ can be decomposed in the product $U'\Lambda U$, where the columns of $U \in O_n$ are the orthonormal eigenvectors of $S$, and $\Lambda$ is the diagonal matrix of the (real) eigenvalues of $S$.*

This eigen decomposition, also called *spectral decomposition*, is closely related to the singular value decomposition, to be discussed further. Quadratic forms play also an important role in optimization problems. The smallest and largest eigenvalues of a symmetric matrix are solutions to a constraint optimization problem, as discovered by Rayleigh and Ritz.

**Therorem 3.6.2:** Rayleigh-Ritz
*Let $S$ be a real symmetric matrix, then the smallest and the largest eigenvalue are respectively the solutions:*

$$\lambda_{\min} = \min_{x'x=1} x'Sx$$
$$\lambda_{\max} = \max_{x'x=1} x'Sx$$

*Proof:* Let $S$ have the eigen decomposition $U'\Lambda U$, with the eigenvalues ordered as $\lambda_{\min} = $

$\lambda_1 \le \lambda_2 \le \cdots \le \lambda_n = \lambda_{\max}$, so that

$$x'Sx = x'U'\Lambda U x = \sum_{i=1}^{n} \lambda_i (Ux)_i^2 = \sum_{i=1}^{n} \lambda_i y_i^2,$$

where we have set $y = Ux$. Since $U$ is an orthogonal matrix, the constraint $x'x = 1$ is equivalent to $y'y = 1$. Hence the optimization problems are solved if one chooses for $y$ respectively the vectors $e_1$ and $e_n$. QED

Note that, when one goes back to the original problem, the optimizing vectors are respectively $x_* = Ue_1 = u_1$ and $x^* = Ue_n = u_n$. The Rayleigh-Ritz theorem can be extended further to:

$$\lambda_k = \min_{\substack{x'x = 1 \\ x \perp \mathrm{span}\{u_1, \ldots, u_{k-1}\}}} x'Sx,$$

where the $u_i$ are the eigenvectors of $S$, i.e. the columns of $U$. We cite a related result, not involving the exact knowledge of the eigenvectors of $S$, known as the "min-max theorem". For a proof, we refer to [2, p. 179]:

**Theorem 3.6.3:** (Courant-Fisher)
*Let $S$ be a symmetric matrix, with eigenvalues ordered by $\lambda_{\min} = \lambda_1 \le \lambda_2 \le \cdots \le \lambda_n = \lambda_{\max}$, and let $k$ be any integer $1 \le k \le n$. Then*

$$\lambda_k = \min_{w_1, w_2, \ldots, w_{n-k} \in \mathbf{R}^n} \quad \max_{\substack{x'x = 1 \\ x \perp \mathrm{span}\{w_1, \ldots, w_{n-k}\}}} x'Sx$$

$$\lambda_k = \max_{w_1, w_2, \ldots, w_{k-1} \in \mathbf{R}^n} \quad \min_{\substack{x'x = 1 \\ x \perp \mathrm{span}\{w_1, \ldots, w_{k-1}\}}} x'Sx$$

Note that if $k = 1$ or $k = n$, then the outer optimization is over an empty set, and does therefore not take place. In these cases this theorem coincides with the Rayleigh-Ritz results. An important consequence of the Courant-Fisher theorem is the following property of the eigenvalues of a bordered symmetric matrix.

**Definition**: Let $A$ be a symmetric matrix, then a *bordered* matrix extension of $A$ is the matrix

$$\overline{A}(a, \alpha) = \begin{bmatrix} A & a \\ a' & \alpha \end{bmatrix} \in \mathbf{R}^{n \times n}$$

with $\alpha \in \mathbf{R}$ and $a \in \mathbf{R}^n$

**Theorem 3.6.4:** *Let the eigenvalues of the bordered matrix $\overline{A}$ be $\mu_1 \ge \mu_2 \ge \cdots \ge \mu_{n+1}$, and the eigenvalues of $A$ be $\lambda_1 \ge \lambda_2 \ge \cdots \ge \lambda_n$, then the interlacing inequalities*

$$\mu_1 \ge \lambda_1 \ge \mu_2 \ge \lambda_2 \cdots \ge \lambda_n \ge \mu_{n+1},$$

*hold.*

*Proof:* Uses the Courant-Fisher variational characterization [2, p. 179]. QED
For a nice exposition of other applications of the min-max theorem, or the variational characterization of the eigenvalues, we refer to the book [2].

### 3.6.2  Definiteness

A *symmetric* matrix $P$ is called *positive semi-definite* if, $\forall x \in \mathbf{R}^n$, the quadratic form $x'Px \geq 0$. It is called *positive definite* if, $\forall x \neq 0 \in \mathbf{R}^n, x'Px > 0$. The following properties of positive definite matrices are easily shown:

1. Any principal submatrix of a positive definite matrix is positive definite.

2. Any nonnegative linear combination of positive definite matrices is positive definite. One can also express this property by saying that the set of positive definite matrices forms a *positive cone* in the vector space $\mathbf{R}^{n \times n}$.

3. Each eigenvalue of a positive definite matrix is real and positive.

4. If $P \in \mathbf{R}^{n \times n}$ is positive definite, then for arbitrary $M \in \mathbf{R}^{p \times n}$, $MPM'$ is positive semi-definite. Since $\rho(MPM') = \rho(M)$, the matrix $MPM'$ is positive definite iff $\rho(M) = p$. ($\rho(M)$ is the rank of $M$.)

Obviously, for all $x \neq 0$ the quadratic form $x'Ax > 0$ iff $A_{\text{Sym}}$ is positive definite.

Finally, a matrix $P$ is called negative definite (semi-definite) iff the matrix $-P$ is positive definite (semi-definite). A matrix that is neither positive nor negative definite is called *indefinite.*

### 3.6.3  Sums of Hermitian Matrices

Consider the triple of Hermitian matrices $A, B$ and $C = A + B$. Let their eigenvalues respectively be enumerated as $\alpha_1 \geq \alpha_2 \geq \cdots \geq \alpha_n$; $\beta_1 \geq \beta_2 \geq \cdots \geq \beta_n$; $\gamma_1 \geq \gamma_2 \geq \cdots \geq \gamma_n$. It follows from the previous section that

$$\gamma_1 \leq \alpha_1 + \beta_1$$
$$\gamma_n \geq \alpha_n + \beta_n.$$

Weyl proved the following more general inequalities:

$$\gamma_{i+j-1} \leq \alpha_i + \beta_j \quad \text{for} \quad i + j - 1 \leq n. \tag{9}$$

In 1949 Ky Fan proved that

$$\sum_{j=1}^{k} \gamma_j \leq \sum_{j=1}^{k} \alpha_j + \sum_{j=1}^{k} \beta_j \qquad 1 \leq k \leq n. \tag{10}$$

In turn, this was further generalized by Lidskii (1950) and Wielandt (final proof). Let $1 \le k \le n$ and let $1 \le i_1 \le i_2 \cdots \le i_k \le n$. Then

$$\sum_{j=1}^{k} \gamma_{i_j} \le \sum_{j=1}^{k} \alpha_{i_j} + \sum_{j=1}^{k} \beta_j. \tag{11}$$

# 4   Calculus: the Great Extremal Problem Solver

## 4.1   Weierstrass's Theorem

### 4.1.1   One dimension

i) Necessary conditions for a minimum at $x^* \in [a, b]$:

   a) $f'(x^*) = 0$, and $f''(x^*) \ge 0$ if $a < x^* < b$

   b) $f'(x^*) \ge 0$, if $x^* = a$

   c) $f'(x^*) \le 0$, if $x^* = b$

In b) and c), the one sided derivatives are understood.

ii) Sufficient conditions for a *local* minimum at $x^* \in [a, b]$:
If we have the strict inequalities:
$$f'(x^*) > 0 \quad \text{for b)}$$
$$f'(x^*) < 0 \quad \text{for c)}$$
$$f''(x^*) > 0 \quad \text{for a)},$$
then there exists a neighborhood, $\mathcal{N}(x^*)$, of $x^*$, such that

$$f(x^*) < f(x), \forall x \in \mathcal{N}(x^*) \cap [a, b] \setminus x^*.$$

iii) Existence of a minimum:
If $f$ is *continuous* on [a,b], then $f$ has a minimum there. (In fact it suffices that $f$ is lower-semicontinuous.)

Uniqueness of the minimum:
If $f$ is *strictly convex* on [a,b], it has a minimum at a unique point $x^* \in [a, b]$. A sufficient condition for $f$ to be strictly convex on $[a, b]$ is

$$f''(x) > 0 \quad \text{on } [a, b].$$

### 4.1.2 Multi Parameter Space

Weierstrass Theorem: Every function which is continuous in a closed domain $\mathcal{D}$ of the variables possesses a largest and a smallest value in the interior or on the boundary of the region.

A necessary condition for the differentiable function $J = L(u_1, \ldots, u_n)$ to have an extremum at an interior point is

$$\frac{\partial L}{\partial u_1} = \frac{\partial L}{\partial u_2} = \cdots = \frac{\partial L}{\partial u_n} = 0.$$

## 4.2 Gradients - Notation

If $u$ is the column vector $[u_1, \ldots, u_n]'$, then we shall denote the *gradient* of $L$ as the *row* vector

$$\frac{\partial L}{\partial u} = \left[ \frac{\partial L}{\partial u_1}, \frac{\partial L}{\partial u_2}, \ldots, \frac{\partial L}{\partial u_n} \right].$$

Thus we may write the total differential

$$dL = \sum_{i=1}^{n} \frac{\partial L}{\partial u_i} \, du_i$$

simply as

$$dL = \frac{\partial L}{\partial u} \, du.$$

## 4.3 Matrixfunctions and their Gradients

More generally, if $X$ is an $n \times m$ matrix of parameters, then we shall define the gradient of the scalar matrix function $f(X)$ as

$$\left( \frac{\partial f(X)}{\partial X} \right)_{\alpha,\beta} = \frac{\partial f(X)}{\partial X_{\beta,\alpha}}$$

Note the reordering of the indices, consistent with considering the gradient with respect to a column vector as a row vector.

It follows now easily from the definition that

$$\frac{\partial \operatorname{Tr} AX}{\partial X} = A$$

and for a symmetric matrix $P$,

$$\frac{\partial \operatorname{Tr} X'PX}{\partial X} = 2X'P$$

In computing matrix gradients, the following identities are helpful:

$$\operatorname{Tr} P = \operatorname{Tr} P'$$

and the *cyclic* permutation property:

$$\operatorname{Tr} ABC = \operatorname{Tr} BCA = \operatorname{Tr} CAB.$$

The following general gradient rule is easily established for a matrix function $\mathcal{A}(\cdot)$. It provides a 'product' rule for the matrix calculus.

$$\frac{\partial \operatorname{Tr} \mathcal{A}(X)X}{\partial X} = \mathcal{A}(X) + \left.\frac{\partial \operatorname{Tr} \mathcal{A}(Y)X}{\partial Y}\right|_{Y=X}.$$

**Example**: Let $X$ be a square matrix:

$$\frac{\partial \operatorname{Tr} XX'X}{\partial X} = XX' + X'X + (X')^2.$$

For matrix functions involving the *inverse matrix*, note that

$$\frac{\partial X^{-1}}{\partial X_{\alpha,\beta}} = -X^{-1}M^{\alpha,\beta},$$

where

$$M_{i,j}^{\alpha,\beta} = \delta_{\alpha j}\left(X^{-1}\right)_{\beta j}$$

**Example**: Let $X$ be nonsingular:

$$\frac{\partial \operatorname{Tr} AX^{-1}}{\partial X} = -X^{-1}AX^{-1}.$$

For matrix functions involving the *determinant*, Laplace's expansion easily shows for nonsingular $X$ that:

$$\frac{\partial \det X}{\partial X} = \operatorname{Adj} X = \det X \cdot X^{-1}.$$

Combining with the product rule we get, letting $\mathcal{A}(X)$ be a square matrix function:

$$\frac{\partial \det \mathcal{A}(X)}{\partial X} = \left.\frac{\partial \operatorname{Tr} \left[\operatorname{Adj} Y \mathcal{A}(X)\right]}{\partial X}\right|_{Y=\mathcal{A}(X)}.$$

**Example**:

$$\frac{\partial \det AX}{\partial X} = \operatorname{Adj}(AX) \cdot X.$$

**Example**:

$$\frac{\partial \det X'X}{\partial X} = 2\operatorname{Adj}(X'X) \cdot X'.$$

**Example**:

$$\frac{\partial \det XX'}{\partial X} = 2X'\operatorname{Adj}(XX').$$

As an interesting application, consider the gradient of the closed loop characteristic polynomial $\alpha_K(s) = \det[sI - A + BK]$ with respect to the feedback gain:

$$\frac{\partial \det[sI - A + BK]}{\partial K} = (sI - A + BK)^{-1}B\,\alpha_K(s)$$

# 5   Numerical Methods

## 5.1   Steepest Descent

## 5.2   Newton's Algorithm

# References

[1] G. Golub and C. Van Loan, *Matrix Computations*, Johns Hopkins University Press, Baltimore, 1983.

[2] R.A. Horn and C.R. Johnson, *Matrix Analysis*, Cambridge University Press, Cambridge, 1990.

[3] G. Strang, "The Fundamental Theorem of Linear Algebra,", *American Mathematical Monthly*, Vol. 100, No. 9, pp. 848-855, November 1993.

[4] R. Bhatin, "Linear Algebra to Quantum Cohomology: The Story of Alfred Horn's Inequalities," *American Mathematical Monthly,* Vol. 108, pp.289-318, April 2001.

[5] J.R. Magnus and H. Neudecker, *Matrix Differential Calculus with Applications in Statistics and Econometrics*, Wiley, 1988.

[6] V.M. Tikhomirov, *Stories about Maxima and Minima*, AMS-MAA, 1991.