

N/40-11809  
N/40-11809  
N/40-11809

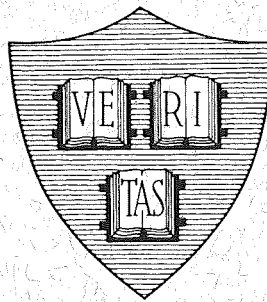
Office of Naval Research

Contract N00014-67-A-0298-0006 NR-372-012

NATIONAL AERONAUTICS AND SPACE ADMINISTRATION

Grant NGR 22-007-143

**MATHEMATICAL PROGRAMMING METHODS  
OF PATTERN CLASSIFICATION**



**CASE FILE  
COPY**

By

**Richard C. Grinold**

**June 1969**

**Technical Report No. 591**

This document has been approved for public release and sale; its distribution is unlimited. Reproduction in whole or in part is permitted by the U. S. Government.

**Division of Engineering and Applied Physics  
Harvard University • Cambridge, Massachusetts**

Office of Naval Research

Contract N00014-67-A-0298-0006

NR-372-012

National Aeronautics and Space Administration

Grant NGR 22-007-143

MATHEMATICAL PROGRAMMING METHODS  
OF PATTERN CLASSIFICATION

By

Richard C. Grinold

Technical Report No. 591

This document has been approved for public release and sale; its distribution is unlimited. Reproduction in whole or in part is permitted by the U. S. Government.

June 1969

The research reported in this document was made possible through support extended the Division of Engineering and Applied Physics, Harvard University by the U. S. Army Research Office, the U. S. Air Force Office of Scientific Research and the U. S. Office of Naval Research under the Joint Services Electronics Program by Contracts N00014-67-A-0298-0006, 0005, and 0008 and by the National Aeronautics and Space Administration under Grant NGR 22-007-143.

Division of Engineering and Applied Physics  
Harvard University · Cambridge, Massachusetts

MATHEMATICAL PROGRAMMING METHODS  
OF PATTERN CLASSIFICATION

By

Richard C. Grinold

Division of Engineering and Applied Physics  
Harvard University Cambridge, Massachusetts

ABSTRACT

This paper studies four mathematical programming methods which are useful in pattern classification. Two of the models are for linearly separable problems, while the others work without separability.

## INTRODUCTION

This report is designed to supplement "On Pattern Classification-Introduction and Survey,"<sup>\*</sup> [7], by describing several mathematical programming approaches to the classification problem. We'll assume that the reader is familiar with the Ho and Agrawala paper (at least sections I, II, and IV) and draw on the motivation, notation, and definitions used there.

Four mathematical programming models are described in detail, and two more are mentioned briefly. Others exist, and are referenced in the publications cited here. The four models were selected for their computational and conceptual properties.

Before describing the contents of the paper, we'll expand upon and change some of the notation adopted in [7].

Definitions:

- (i). Instead of  $x^1(i)^T$  and  $x^0(j)^T$ , the training samples from classes one and zero will be denoted by  $m$  component row vectors  $A_i^1$  and  $A_j^0$ .
- (ii). For  $k = 0, 1$ ;  $A^k$  is the  $n_k \times m$  matrix whose rows are  $A_i^k$ ,  $i = 1, 2, \dots, n_k$
- (iii).  $h, \ell, e, q$ , and  $f$  are vectors of ones. Their dimensions are given below.  
 $h \sim n_1 \times 1$ ;  $\ell \sim n_0 \times 1$ ;  $e \sim m \times 1$ ;  $q \sim n_1 \cdot n_0 \times 1$  and  
 $f \sim n \times 1$ , where  $n = n_1 + n_0$ .

---

<sup>\*</sup>Y. C. Ho and A. K. Agrawala, Technical Report No. 557, Division of Engineering and Applied Physics, Harvard University March (1968), also published in IEEE Trans. on Auto. Cont. Vol. 13, No. 6, December 1968 and Proceedings of the IEEE, Vol. 56, No. 12, December 1968.

(iv). A linear decision function is an  $(m+1) \times 1$  vector  $w = \begin{bmatrix} \lambda \\ u \end{bmatrix}$ .

Using these definitions we note that

$$A = \begin{bmatrix} h & A^1 \\ -\ell & -A^0 \end{bmatrix}$$

Suppose the patterns are described by an  $m$  component row vector  $x$ , then the decision function defined by  $w$  is

$$f(x) = \lambda + xu = (1, x)w.$$

(v). A linear decision function  $w$  is a separator if

$$Aw > 0$$

Problems are specified by their range of attributes.

The ranges are defined below.

$$(vi). \quad S^1 = \{x \mid P(x \mid H^1) > 0\}$$

$$S^0 = \{x \mid P(x \mid H^0) > 0\}$$

(vii). The operator  $C$  will denote convex closure. Thus  $C(S^1)$  is the closed convex hull of  $S^1$ .

(viii). The problem is separable if  $C(S^0)$  and  $C(S^1)$  are disjoint. If they intersect the problem is nonseparable.

(ix). The problem is decidable if  $S^0$  and  $S^1$  are disjoint.

(x). We shall define  $C(A^k)$  as the convex hull of the rows of  $A^k$ ,  $k = 0, 1$ .

Thus

$$C(A^0) = \{b \mid b = zA^0, z\ell = 1, z \geq 0\}$$

$$C(A^1) = \{b \mid b = yA^1, yh = 1, y \geq 0\}$$

Five sections and an appendix follow. Sections one and two describe models used in the separable and nonseparable cases. The third section remarks on the model's flexibility in terms of accommodating new data and use in judging new features. Section four considers the model's generalization properties, while the last section describes an application. The appendix is a brief introduction to linear and quadratic programming.

Of the four models, two have been described in the published literature. One of the unpublished models is due to Canon and Cullum [3], the other is the author's responsibility [6].

# I. THE SEPARABLE CASE

Charnes [4] and Mangasarian [9] independently proposed a linear programming model for separating disjoint polyhedrons. Their distinct approaches illustrate the duality principle of linear programming. Charnes asks if the sets  $C(A^0)$  and  $C(A^1)$  are disjoint, while Mangasarian looks directly for a separating hyperplane.

By definition, the  $C(A^k)$  will intersect if and only if the system (1) has a feasible solution.

$$zA^1 - yA^0 = 0, \quad y_h = 1, \quad z_l = 1, \quad y \geq 0, \quad z \geq 0 \quad (1)$$

We can discover a solution of (1) by adding artificial variables to the system and minimizing the infeasibility. This gives us a linear program.

$$\text{Minimize } (r + s)e \quad (2)$$

Subj. to

$$zA^0 - yA^1 + rI - sI = 0$$

$$z_l = 1$$

$$y_h = 1$$

$$z \geq 0, \quad y \geq 0, \quad r \geq 0, \quad s \geq 0$$

This problem has  $m + 2$  equality constraints with  $n + 2m$  nonnegative variables. The value,  $(r + s)e$ , is nonnegative since  $r$  and  $s$  are nonnegative. Finally, we can easily construct a first basic feasible solution of (2).

The alternate approach involves the decision function directly. Suppose the  $m$  vector  $u$  and scalars  $(\gamma, \delta)$  satisfy the following conditions:

$$\begin{aligned} \delta - \gamma &> 0 \\ -A^1 u - h\gamma &\leq 0 \\ A^0 u + \ell\delta &\leq 0 \end{aligned} \tag{3}$$

then  $w = (\frac{\gamma + \delta}{2}, u)$  is a separator. A solution of (3) can be discovered by solving (4).

$$\begin{aligned} &\text{Maximize } \delta - \gamma \\ &\text{Subj. to} \\ &-A^1 u - h\gamma \leq 0 \\ &+A^0 u + \ell\delta \leq 0 \\ &-e \leq u \leq e \end{aligned} \tag{4}$$

Problem (4) has  $n$  inequality constraints, two free variables  $(\gamma, \delta)$ , and  $m$  variables with upper and lower bounds. The bounds-rule out infinite solutions. Evidently,  $(\gamma, \delta, u) = (0, 0, 0)$  is a feasible solution of (4).

Appealing to the results in the appendix we can state that problem (4) is the dual of problem (2), and the duality theorem applies. This guarantees the existence of optimal solutions  $(\bar{z}, \bar{y}, \bar{r}, \bar{s})$  and  $(\bar{\gamma}, \bar{\delta}, \bar{u})$  such that:

$$(\bar{r} + \bar{s})e = \bar{\delta} - \bar{\gamma} \geq 0.$$

There are two possibilities. If  $\delta - \gamma > 0$ , then  $(\frac{\bar{\gamma} + \bar{\delta}}{2}, \bar{u})$  is a separator. If  $e(\bar{r} + \bar{s}) = 0$ , then  $(\bar{z}, \bar{y})$  solves (1), and the convex hulls intersect. These facts are summarized below.

Theorem: (5)

- (i). Problems (2) and (4) have optimal solutions with equal, nonnegative values.



- (ii). If the optimal value is zero, the patterns are not linearly separable.
- (iii). If the optimal value is positive, then  $\bar{w} = (\frac{\bar{y} + \bar{\delta}}{2}, \bar{u})$  defines a separator which maximizes

$$\text{Min}[A_i w \mid i = 1, 2, \dots, n]$$

Subj. to

$$-1 \leq w_j \leq 1, \quad \text{for } j = 1, 2, \dots, m$$

Statements (i), (ii), and the first part of (iii) are established above. The final statement can be established by contradiction.

The linear programs will be solved using some variant of Dantzig's, [2], simplex method. This is a rapidly convergent combinatorial procedure, while the adaption algorithms, see [7] Table I, are gradient descent techniques which converge slowly. If the patterns are not separable, slow and nonconvergence can be confused. See [9], pg. 451, for a more detailed comment along this line. The adaption algorithms do have the advantage of simplicity, but this is largely offset by the wide availability of professionally written linear programming codes. Either (2) or (4) can be solved, but the simplex algorithm is more efficient with fewer nontrivial constraints. It is not very sensitive to the number of variables. Since  $m + 2 \ll n$  it is reasonable to solve (2).

Canon and Cullum [3] have proposed a quadratic programming method for the separable case. Although it is generally more difficult to solve quadratic programs, the authors take advantage of the problem's special structure and claim their method is competitive with the linear programming model.

For each  $i$  and  $j$ ,  $i = 1, 2, \dots, n_1$ ,  $j = 1, 2, \dots, n_0$ ; we can define a difference vector

$$D_k = A_i^1 - A_j^0 \quad \text{for } k = 1, 2, \dots, n_1 n_0.$$

The vectors  $D_k$  are the rows of the  $n_1 n_0 \times m$  matrix  $D$ . Recall  $C(D) = \{u \mid u = yD, yg = 1, y \geq 0\}$ . It is easy to establish that  $C(A^0)$  and  $C(A^1)$  will be separable if and only if the origin is not contained in  $C(D)$ .

This suggests a test for separability: find the vector in  $C(D)$  with minimum norm. The problem can be written in two ways:

$$\text{Minimize } \frac{uIu'}{2} \tag{6}$$

Subj. to

$$yD - uI = 0$$

$$yg = 1$$

$$y \geq 0$$

$$\text{Minimize } \frac{yDD'y'}{2} \tag{7}$$

Subj. to

$$yg = 1$$

$$y \geq 0$$

The following facts about (6) and (7) should be clear: they are equivalent, the objectives are convex and quadratic, they have optimal solutions with nonnegative values, and the sets are separable if and only if the optimal value is positive.

It is well known that any point in  $C(D)$  can be expressed as a convex combination of at most  $m + 1$  rows of  $D$ . This fact is used to reduce the problem's size. The algorithm solves a modified version of (6), restricting

attention to a subset of  $m + 1$  rows. A test sees if the restricted solution is optimal with all rows considered. If so, (6) is solved. If not, a new row is added, an old row dropped, and the algorithm proceeds, finding an optimal solution in a finite number of steps. The optimal solution of (6) defines the linear decision surface.

Suppose  $(\bar{y}, \bar{u})$  solves (6),  $\bar{u} \neq 0$ , and

$$\gamma = \text{Min} \{A_i^1 \bar{u}^1 \mid i = 1, 2, \dots, n_1\}$$

$$\delta = \text{Max} \{A_j^0 \bar{u}^1 \mid j = 1, 2, \dots, n_0\}$$

then  $(\frac{\gamma + \delta}{2}, \bar{u})$  is a separator. If  $\bar{u} = 0$ , no separator exists. This is demonstrated in the appendix using the Kuhn-Tucker theorem. Canon and Cullum do the same by showing problem (6) is equivalent to:

$$\text{Max} [\text{Min} \{uz \mid u \in C(D)\}] \quad (8)$$

Subj. to

$$z^1 I z \leq 1$$

## II. NONSEPARABLE

One approach to the nonseparable case was taken in [6]. A description will require two definitions.

(xi). Let  $a = \sum_{i=1}^n \frac{A_i}{n}$  be the average of the rows of A.

(xii). For any decision function  $w$ , let the quality of  $w$  be defined as

$$\text{Min} [A_i w \mid i = 1, 2, \dots, n]$$

If  $w$  is a separator, the quality is positive. If  $w$  is not a separator, the negative of the quality (a nonnegative number) measures the largest error the decision surface makes. To obtain a decision surface of highest quality we solve

$$\begin{array}{ll} \text{Maximize } \{ \text{Min}[A_i w \mid i = 1, 2, \dots, n] \} \\ \text{Subj. to} & aw = 1 \end{array}$$

The constraint is a normalization.

This problem can be transformed into a linear program by introducing a new variable  $\rho$  and requiring  $\rho \leq A_i w$  for  $i = 1, 2, \dots, n$ . The new problem and its dual are given below.

$$\begin{array}{ll} \text{Maximize } \rho & (9) \\ \text{Subj. to} & \\ Aw - f\rho \geq 0 & \\ aw = 1 & \end{array}$$

$$\begin{array}{ll} \text{Minimize } \gamma & (10) \\ \text{Subj. to} & \\ yA - \gamma a = 0 & \\ yf = 1 & \\ y \geq 0 & \end{array}$$

Problem (10) has  $m + 2$  equality constraints,  $n$  nonnegative variables, and one free variable.

The main result of [6] is:

Theorem (11)

- (i). Problems (9) and (10) have optimal solutions with equal objective values iff and  $a \neq 0$ . When  $a = 0$ , (9) is infeasible.

- (ii). If  $(w, \rho)$  solves (9) and  $\rho > 0$ , then  $w$  defines a separator of maximum quality.
- (iii). If  $(w, \rho)$  solves (9) and  $\rho \leq 0$ , then the patterns are not separable and  $w$  defines a decision surface that minimizes the maximum error.

This is equivalent to (5) in the separable case. In addition, a meaningful decision surface is generated if the patterns are not separable.

Smith, [13], has another approach. Note that  $Aw > 0$  has a solution iff  $Aw \geq f$  has a solution. In this spirit, we can solve

$$\begin{aligned} & \text{Minimize } f'v \\ & \text{Subj. to} \\ & Aw + Iv \geq f \\ & v \geq 0 \end{aligned}$$

The  $v_i$ 's measure the size of any error in the classification of the  $i$ th sample. Thus if  $A_i w \geq 1$ , there is no error and  $v_i = 0$ . If  $A_i w < 1$ ,  $v_i$  is positive. There is some difficulty if  $0 < A_i w < 1$ . In this case the pattern is correctly classified, but an error is counted. This behavior is observed in optimal solutions.

The dual, (12), is a linear program with  $m + 1$  equality constraints and  $n$  nonnegative variables with upper bounds. It is relatively easy to solve, [5].

$$\begin{aligned} & \text{Maximize } yf & (12) \\ & \text{Subj. to} \\ & yA = 0 \\ & 0 \leq y \leq f' \end{aligned}$$

This model suggests several conceptually interesting but computationally difficult variations. For instance, we could minimize the sum of squared errors. This leads to a quadratic program:

$$\begin{aligned} & \text{Minimize } v' Iv \\ \text{Subj. to} & \\ & Aw + Iv \cong f \\ & v \cong 0 \end{aligned}$$

Another variant maximizes the number of correctly classified samples:

$$\begin{aligned} & \text{Maximize } \sum_{i=1}^n \delta(A_i w) \\ \text{Subj. to } & -1 \cong w_i \cong 1 \quad i = 0, 1, 2, \dots, m \end{aligned}$$

$\delta(\cdot)$  is the step function; one if its argument is positive, zero otherwise. This problem can be reformulated as an integer program, [12] pp. 194-8.

Another method of treating the nonseparable case was proposed by Mangasarian, [10]. The approach is similar to Arkadev and Braveman [1], i. e. a piecewise linear decision surface is created which decides correctly about all the data. Mangasarian uses mathematical programming to construct the decision surface. We will not examine that algorithm in detail, but we do comment on its generalization properties in section 4.

### III. FLEXIBILITY

This section examines the ability of the different models, (2), (6), and (12) to handle new data and yield information useful in selecting new

features. Models (2), (10), and (12) can accept new data points and find a new decision surface easily. In each case, adding a new point is equivalent to introducing a new activity (column) into the linear program.

Model (6) has a similar property. For example, suppose a new point in class one,  $A_{n_1+1}^1$  is observed. This adds  $n_0$  new rows to the matrix D. If

$$A_{n_1+1}^1 u' \geq \text{Min} [A_{n_i}^1 u' \mid i = 1, 2, \dots, n_1]$$

no change is needed, the old decision surface is still optimal. If the inequality does not hold, we continue to apply the Canon-Cullum algorithm until a new optimal solution is obtained.

Introducing a new feature in (2), (10), or (12),<sup>†</sup> adds a new constraint (row) to the linear program. If several new features are being considered, we can devise a heuristic rule for choosing among them. Try the current optimal solution for each new constraint. Select the constraint which is the furthest from being satisfied. If the optimal solution satisfies all the new constraints, it is still optimal. This selects the feature which maximizes the rates of improvement of the solution. Then a new optimal solution can be obtained using the dual simplex method.

---

<sup>†</sup> There doesn't seem to be any way that new feature can be accommodated by model (6).

#### IV. GENERALIZATION

The generalization properties of the models are examined in this section. In particular, we are interested in the decision surfaces generated as the number of sample points  $n$  becomes large. For each  $n$  the models produce a decision surface defined by a nonzero  $m + 1$  vector. Without loss of generality we can uniformly bound these vectors. Thus, there will be subsequences which converge. We shall study the properties of the limiting decision surface.

For example, assume  $C(S^1)$  and  $C(S^0)$  are disjoint with one set compact, and consider model (2). Let  $(\lambda^n, u^n)$  be the normalized optimal decision surface for the  $n$  sample problem and let  $(\lambda, u)$  be a limiting surface: i. e.  $(\lambda^n, u^n) \rightarrow (\lambda, u)$  on some subsequence. The following theorem asserts  $(\lambda, u)$  is optimal for the limiting problem.

Theorem:

With probability one (wp. 1) there exists a  $\rho > 0$  such that  $(\lambda, u, \rho)$  solve:

$$\begin{array}{ll} \text{Maximize } \rho & \\ \text{Subj. to} & \\ \lambda + xu - \rho \geq 0 & x \in C(S^1) \\ -\lambda - xu - \rho \geq 0 & x \in C(S^0) \\ -e \leq u \leq e & \end{array}$$

Proof:

There exists a hyperplane which strictly separates  $C(S^1)$  and  $C(S^0)$ . Therefore the problem has an optimal solution with positive value.



Suppose  $(\lambda, u)$  and some  $\rho > 0$  are not optimal. A contradiction can be established by appealing to the facts that  $(\lambda, u)$  is (2) feasible for all  $n$ , and that  $(\lambda, u)$  is the limit of a subsequence of optimal solutions.

Three comments are in order. First it is obvious that similar results hold for models (6), (10), and (12). Secondly, if compactness is dropped a weaker,  $\rho \geq 0$ , statement is true. Finally, if separability doesn't hold, then (wp, 1) all models will indicate this for some large value of  $n$ .

Assuming decidability we could obtain a like result using the piecewise approach, [10]. Additional regularity assumptions are needed to allow a piecewise linear function defined by a finite number of hyperplanes. Without decidability, the piecewise approach would struggle in vain to produce a perfect decision function.

Model (10) will work in the separable case, but it has questionable generalization properties. It is very sensitive to the tails of the distributions. The decision surface minimizes the maximum error, therefore it will react to the worst points or perhaps to a faulty observation. Things can get worse.

Let  $a^k$  be the finite means of the distributions,  $P(x|H^k)$  for  $k = 0, 1$ . Then the row average of  $A$  will converge (wp. 1) to

$$a = \begin{pmatrix} P(H^1) a^1 - P(H^0) a^0 \\ P(H^1) - P(H^0) \end{pmatrix}$$

The following is an example of what can go wrong. Suppose  $P(H^1) > P(H^0)$ , and the sets described below have a nonvoid intersection:

$$L = \{d \mid d = \gamma a, \gamma \leq 1\}$$

$$Z = \{d \mid d = \begin{pmatrix} b \\ -1 \end{pmatrix}, -b \in C(S^0)\}$$

then the limiting optimal solution is given by  $w = (\frac{1}{P(H^1) - P(H^0)}, 0)$  i. e. the decision function is

$$f(x) = \frac{1}{P(H^1) - P(H^0)} > 0 \quad \text{for all } x$$

The fact that  $f$  is correct more than not offers little consolation. Note that this phenomenon will occur if  $S^0 = S^1 = R^m$ , and  $P(H^0) \neq P(H^1)$ : e. g. multivariate normal.

The generalization properties of (12) seem to be the best. It is a reasonable conjecture that the limiting decision surfaces of (12) are optimal solutions to the following:

Minimize  $F(w)$

Subj. to

$$-1 \leq w_i \leq 1 \quad i = 0, 1, 2, \dots, m$$

where

$$F(w) = P(H^1) \int_{X^0} (-\lambda - ux) P(x \mid H^1) dx + P(H^0) \int_{X^1} (\lambda + ux) P(x \mid H^0) dx$$

is the expected error distance. It is also reasonable to assume that the limiting decision surfaces of the integer program mentioned in section two will minimize the probability of error among all linear decision functions. A brief attempt was made to prove these conjectures, but the proof is elusive.

## V. EXAMPLE

Models (10) and (12) were employed to design decision functions using data from a NASA biomedical experiment. Two types of electroencephalograms (brainwaves, EEG) were recorded. In one instance the subject was watching a strobe light. In the other case the light was not visible. The object is to distinguish the two cases using the EEG data.

Of a possible one hundred features K. Prahbu selected five, using a distance-dispersion technique and prepared the data for the linear programming models. The parameters were  $n_0 = 165$ ,  $n_1 = 155$ ,  $m = 5$ ,  $n = 320$ , and the problems were solved on an IBM 360-65 using the mathematical programming package, MPS 360, [11]. Results are tabulated below.

Model (10) Solution Time 0.09 min.

Errors	Number	Percentage
Type I	25	16.6
Type II	21	12.7
Total	46	14.4

Model (12) Solution Time 0.92 min.

Errors	Number	Percentage
Type I	18	11.6
Type II	24	14.5
Total	42	13.1

Notice the performance of model (12) is slightly better although the solution time is longer. Both problems had unique optimal solutions and 31 of the points were incorrectly classified by both techniques.

## APPENDIX

### Linear and Quadratic Programming

Several results from mathematical programming have been used in this report. This appendix attempts to motivate and explain these results while citing more substantial references.

A linear program is an optimization problem

$$\begin{array}{ll} \text{Min} & \sum_{j=1}^m x_j c_j \\ \text{Subj. to} & \sum_{j=1}^m x_j a_{ji} = b_i \quad i = 1, 2, \dots, n \\ & x_j \geq 0 \quad j = 1, 2, \dots, m \end{array}$$

Our vector notation is

$$\begin{array}{ll} \text{Min} & xc \\ \text{Subj. to} & \end{array}$$

$$xA = b$$

$$x \geq 0$$

$c$  is  $m \times 1$ ,  $A$   $m \times n$ , and  $b$   $1 \times n$ . We shall call this problem the primal.

There is an associated dual problem:

$$\begin{array}{ll} \text{Max} & by \\ \text{Subj. to} & \end{array}$$

$$Ay \leq c$$

Linear programs appear in many forms: maximization or minimization, equalities or inequalities, nonnegative or unrestricted variables. Any problem can be transformed into the same form as our primal, which allows us to know its dual. The dual can be found directly using the diagrams on pp. 126-7 of [2].

An efficient algorithm known as the simplex method, has been devised to solve linear programs. In a finite number of steps it finds a feasible solution (if one exists), then again in a finite number of steps it determines an optimal or an unbounded solution. An optimal dual solution is supplied as a by product of the calculations.

The principle theoretical result in linear programming relates primal and dual.

Theorem: [2] pg. 129

If both primal and dual have feasible solutions, they have optimal solutions  $(x, y)$  such that

$$xc = by$$

We shall discuss quadratic programming in the context of problem (6).

$$\text{Min } \frac{uIu^t}{2}$$

Subj. to

$$yD - uI = 0$$

$$yg = 1$$

$$y \geq 0$$

A central result in the study of these problems is the Kuhn-Tucker theorem, [8]. In our case it states:

Theorem:

$(y, u)$  is optimal for (6) if and only if there exist  $(x, z, \lambda)$  such that

$$g\lambda + Dx + z = 0$$

$$u' - x = 0$$

$$y \geq 0$$

$$z \geq 0$$

$$yg = 1$$

$$yD - uI = 0$$

$$yz = 0$$

Suppose  $u \neq 0$  is optimal in (6), then  $uIu' > 0$ . Juggling the above equations we can easily establish that

$$\lambda = -uIu' < 0$$

and

$$Du' \geq g(uIu') > 0.$$

## REFERENCES

1. Arkadev, A. G. and Braveman, E. M., *Computers and Pattern Recognition*, Washington, D. C., Thompson 1967.
2. Dantzig, G. B., *Linear Programming and Extensions*, Princeton University Press, Princeton, N. J., 1963.
3. Canon, M. D. and Cullum, C. D., "The Determination of Optimum Separating Hyperplanes, I: A Finite Step Procedure," Report (un-numbered) of IBM Watson Research Center, Yorktown Heights, New York.
4. Charnes, A., "Some Fundamental Theorems of Perception Theory and Their Geometry," in *Computer and Information Sciences*, J. T. Tou and R. H. Wilcox (eds.) pp. 67-74, Spartan Books, Washington, D. C. 1964.
5. Grinold, R. C., "A Comment on 'Pattern Classifier Design by Linear Programming,'" *IEEE Transactions on Computers*, to appear.
6. Grinold, R. C., "A Note on Pattern Separation," submitted to *Operations Research*.
7. Ho, Y. C. and A. K. Agrawala, "On Pattern Classification Algorithms- Introduction and Survey," *Proc. IEEE* Vol. 56, No. 12, December 1968.
8. Kuhn, H. W. and A. W. Tucker, "Nonlinear Programming," *Econometrica*, Vol. 19, No. 1, January 1951.
9. Mangasarian, O. L., "Linear and Non-Linear Separation of Patterns by Linear Programming," *Operations Research* 13, 1965, pp. 444-452.
10. Mangasarian, O. L., "Multi-Surface Method of Pattern Separation," *IEEE, Transactions on Information Theory*, November 1968.
11. IBM Publication H20-0476-1, *Mathematical Programming System 360, Linear and Separable Programming - Users Manual*.
12. Simmonard, M., "Linear Programming," Prentice Hall, 1966.
13. Smith, F. W., "Pattern Classifier Design by Linear Programming," *IEEE Transactions on Computers*, Vol. C-17, No. 4, April 1968, pp. 367-372.





Unclassified

Security Classification

## DOCUMENT CONTROL DATA - R &amp; D

(Security classification of title, body of abstract and indexing annotation must be entered when the overall report is classified)

1. ORIGINATING ACTIVITY (Corporate author) Division of Engineering and Applied Physics Harvard University Cambridge, Massachusetts		2a. REPORT SECURITY CLASSIFICATION Unclassified	
3. REPORT TITLE  MATHEMATICAL PROGRAMMING METHODS OF PATTERN CLASSIFICATION		2b. GROUP	
4. DESCRIPTIVE NOTES (Type of report and, inclusive dates) Interim technical report			
5. AUTHOR(S) (First name, middle initial, last name)  Richard C. Grinold			
6. REPORT DATE June 1969		7a. TOTAL NO. OF PAGES 25	7b. NO. OF REFS 13
8a. CONTRACT OR GRANT NO. N00014-67-A-0298-0006 & NASA Grant		9a. ORIGINATOR'S REPORT NUMBER(S)  Technical Report No. 591	
b. PROJECT NO. NGR 22-007-143		9b. OTHER REPORT NO(S) (Any other numbers that may be assigned this report)	
c.			
d.			
10. DISTRIBUTION STATEMENT This document has been approved for public release and sale; its distribution is unlimited. Reproduction in whole or in part is permitted by the U.S. Government.			
11. SUPPLEMENTARY NOTES		12. SPONSORING MILITARY ACTIVITY  Office of Naval Research	
13. ABSTRACT  This paper studies four mathematical programming methods which are useful in pattern classification. Two of the models are for linearly separable problems, while the others work without separability.			

4 KEY WORDS	LINK A		LINK B		LINK C	
	ROLE	WT	ROLE	WT	ROLE	WT
Linear Programming Quadratic Programming Pattern Classification Pattern Recognition Non Parametric Hypothesis Testing						