

Optimization

Jacob Denson

December 3, 2016

Table Of Contents

I	Linear Programming	2
1	Introduction	3
1.1	Piecewise Linear Convex Functions	9
2	Polyhedra	11
3	Duality	13
4	Applications to Combinatorial Optimization	16
5	The Ellipsoid Method	19
II	Convex Optimization	20
6	Convex Functions	22
7	Mirror Descent	25

Part I

Linear Programming

Chapter 1

Introduction

Most problems in computing science can be described in a specific format. The problem, given an input, has a certain set of feasible outputs. On this output we can place a partial ordering. The problem is then phrased as an optimization – find the maximum or minimum feasible solution according to the ordering placed on the solutions. Examples of these problems are finding the ‘shortest’ path, the ‘longest’ common subsequence. Since we only have a finite set of solutions, the maximum can always be obtained in these problems.

Linear Programming attempts to provide a general solution to optimization problems with a ‘linear’ structure. We attempt to optimize a linear reward function, subject to linear constraints, where the parameters we attempt to optimize take values in a linear space. We give an example to show this idea in action.

Example. *Suppose we are stocking items at a store. Each item needs to be created from a specific number of ingredients we need to buy, each with a specific stock. For simplicity, suppose we can only stock two items, apple pie, and chocolate bars. Let a be the number of apple pies we make, and c the number of chocolate bars.*

For each apple pie we make, we must buy 3 apples, and use $1/2$ bags of sugar, and for every three chocolate bars, 1 bag of sugar must be used. Each apple is \$1, and each bag of sugar is \$3. Then we can model the cost in dollars of the process by the equation

$$9a/2 + 1/3c$$

However, we are not only buying items, but selling them as well. Each apple pie is sold for \$6, and each chocolate bar for \$1. The problem is to find the optimal number of apple pies and chocolate bars to make, maximizing the equation

$$(6 - 9/2)a + (1 - 1/3)c$$

The problem with the current description is that theoretically we could make as many apple pies and chocolate bars as we wanted, making theoretically an unbounded amount of money. Thus we add additional constraints. There is only a limited supply of sugar: we can only buy 50 bags. This means that, if we let $s = 1/2a + 1/3c$ be the number of bags of sugar we buy

$$1.2a + 1.3c \leq 50$$

We can additionally only buy 100 apples

$$a/3 \leq 100$$

What's more, it takes time to produce these items. In the time allotted, we can only produce 40 items.

$$a + c \leq 40$$

With these constraints, the problem of maximizing revenue becomes much more involved. This is an example of a linear program. If we require the number of apple pies and chocolate bars we make to be integers, then this is an integer linear programming, which is computationally much harder.

A linear program is specified by a reward (objective) functional $r \in V^*$, constraint functionals $c_1, \dots, c_n \in V^*$, and a constraint vector $b \in \mathbf{F}^n$ (where \mathbf{F} is either the rational or real numbers). The aim of the linear program is to maximize $r(v)$, subject to the constraints that $c_i(v) \leq b_i$ (the set of all v such that $c_i(v) \leq b_i$ is known as the solution space, and forms a polyhedron, a finite intersection of half planes). For computational purposes, we can choose a m dimensional basis for V , and specify the problem with a constraint matrix $M \in M_{n,m}(\mathbf{F})$ and reward functional $r \in \mathbf{R}^m$. Our job is then to maximize $\langle r, v \rangle$, subject to the constraints

$$\sum M_{ij}v_j \leq b_i$$

which we write in condensed form as $Mv \leq b$ (which makes sense if we make \mathbf{R}^n into a Riesz space, viewing it as functions over a discrete set

of n points). We say a linear program is **infeasible** if the solution space is empty, and unbounded if there is x_i in the solution space such that $r(x_i) \rightarrow \infty$.

In a more general form, we allow more fluid constraints for the problem, which can be simplified to the existing problem with little asymptotic cost. For instance, we may allow constraints of the form $c_i(v) \geq b_i$. This can be reduced to the form $-c_i(v) \leq -b_i$, which fits into our current assumptions nicely. We can also allow equalities of the form $c_i(v) = b_i$, which corresponds to two inequalities of the form $c_i(v) \leq b_i$ and $c_i(v) \geq b_i$. Finally, we shall assume there is a basis x_1, \dots, x_n such that, relative to this basis, the constraints $x_i^*(v) \geq 0$ are satisfied for any v in the solution space. In general, if we do not have these constraints on a certain basis x_1, \dots, x_n , we may introduce a new linear program on $V \times V$, with constraints (and objective function) induced from V by the map $T : (y, x) \mapsto y - x$. Any maximum solution on this new linear program induces a solution on the original programming by considering the image under T , and since each element of V can be written as $y - x$, where $y_i, x_i \geq 0$ for all i , we can introduce these new constraints at no fault.

Linear programming can be considered a specialization of convex optimization, since the solution space is always the intersection of finitely many closed hyperplanes, and is therefore convex. We recall that a point in a convex set is an **extreme point** if it does not lie between any two points in the convex set. An alternate criteria for x on a set X is that for any y , either $x + y$ or $x - y$ is not in X . Furthermore, a point on a polyhedron is an extreme point if it is the maximum of some linear functional over the polyhedron, which indicates their importance in linear optimization. The extreme points of a linear constraint problem are particularly easily to quantify. If we consider a particular constraint matrix $Ax \leq b$, we shall let A^x denote the rows the matrix obtained from A for which the inequality is tight, i.e. $(Ax)_i = b_i$. Similarly, we let b^x denote the subvector of b such that $A^x x = b^x$.

Theorem 1.1. *The extremities to $Ax \leq b$ are those where A^x has full rank.*

Proof. Let A^x have full rank, and suppose $Ax \leq b - Ay$, $Ax \leq b + Ay$. Then $A^x y = 0$, so $y = 0$ since A^x has full rank. Conversely, suppose that A^x does not have full rank. Then there is $y \neq 0$ for which $A^x y = 0$. If $\varepsilon > 0$ is made small enough, we have $Ax \leq b - Ay$ and $Ax \leq b + Ay$. \square

Lemma 1.2. *If a linear program has an optimal solution, it has an optimal solution which is an extreme point.*

Proof. Let X be the solution space of the linear program. We take some optimal solution x , and slide it along $\ker(r)$. If y is an extreme point of $Y = (x + \ker(r)) \cap X$ (which must exist because of the $x_i \geq 0$ condition), then we claim $x + y$ is also an extreme point of X . If $x + y + z \in X$ and $x + y - z \in X$, then we cannot have $r(z) = 0$, since y is an extreme point of Y . But then since $r(x + y) = r(x)$ maximizes r , we cannot have both $x + y + z$ and $x + y - z$ be elements of X , because either $r(x + y + z)$ and $r(x + y - z)$ is greater than $r(x + y) = r(x)$. \square

Note that the last theorem is not true if we do not assume $x_i \geq 0$ for all i , which is why we prefer discussing linear programs in standard form. Given m constraints on an n -dimensional linear program, there can only be $\binom{m}{n}$ extreme points on the solution space, because a point is uniquely specified by its particular value on an invertible matrix, and every extreme point is determined by its values on some invertible submatrix of the constraint matrix. This gives us a particular algorithm for finding the solution to a linear program – just find the value of the function at the extreme points. Unfortunately, this method can be exponential in the description of the problem. However, we can represent these points in polynomial time.

Theorem 1.3. *If a polyhedron is defined by an inequality of the form $Ax \leq b$, where $A \in M_{n,m}(\mathbf{Q})$, and all coefficients involved can be expressed in K bits, then the extreme points of the polyhedron can be expressed in coefficients using only $2m \log(m)K$ bits.*

Proof. Each extreme point x can be defined as the unique solution to $Bx = c$, where B is an invertible $m \times m$ submatrix of A . Then, using Cramer's rule, we may write

$$x_i = \frac{\det(A_i)}{\det(A)}$$

where A_i is obtained from A by replacing the i 'th column of A with c . Now the determinant can be expressed as

$$\det(X) = \sum_{\sigma \in S_m} \text{sgn}(\sigma) \prod_{i=1}^m X_{i\sigma(i)}$$

The product of m numbers representable in K bits can be written in mK bits, and the sum of these numbers can be expressed in $m\log(m)K$ bits. Similarly, the number of bits to express a/b is bounded by the sum of the number of bits to express a and b . Thus each coefficient of x can be expressed in $2m\log(m)K$ bits. \square

Since A and c take $(n+1)mK$ bits to describe, this gives us a polynomial bound in the size of solutions. Thus, if we can choose extreme points wisely, then we may well achieve a polynomial solution.

Theorem 1.4. *A feasible, bounded linear program has an optimal solution.*

Proof. If X is the feasible solution space of a linear program, then it is closed and convex. Thus $r(X)$ is also closed and convex. If we are working over \mathbf{R} , we may apply completeness and conclude that $r(X) = (-\infty, b]$ or $r(X) = [a, b]$. In either case, we conclude there is $v \in X$ for which $r(x) = b$ maximizes the functional. To obtain results in \mathbf{Q} , we first maximize over \mathbf{R} , and then justify why we can obtain estimates over \mathbf{Q} . The last lemma implies that if $x \in \mathbf{R}^n$ maximizes $r(x)$, and we take x to be an extreme point, then A^x has full rank. Since $A^x x = b^x$, we may throw away rows from A^x to make an invertible matrix, and then express x as the product of b^x with an inverse, which has rational entries, and thus $x \in \mathbf{Q}^n$. \square

The main goal of this report is to find ways of optimizing linear programs in an algorithmically simple way. In the combinatorial method (where we maximize over the finite number of extreme points) the linear optimization problem is often easier to solve if almost all constraints are equalities. That is, a linear program is in **slack form** if we attempt to maximize $c(v)$ such that $Av = b$, and $v \geq 0$. We can make an inequality of the form $\lambda(v) \leq b$ into an equality by introducing a new variable x , defining $\lambda(v, x) = \lambda(v) + x$, and enforcing that $\lambda(v, x) = b$, and $x \geq 0$. This only increases the dimensionality of the problem by a linear amount in the constraints of the problem.

We finish this chapter by discussing some applications of linear programming, mainly to combinatorial problems in computing science (since they are the most interesting).

Example. *The classical computing science problem to be explained in linear programming is the max flow problem. Given a directed graph, where every edge e has a certain capacity c , we want to maximize the amount of ‘flow’*

between two vertices s and t , such that flow conservation is preserved at every vertex other than s and t . That is, we wish to find a function $f : E \rightarrow \mathbf{R}^+$, such that

$$\sum_{sv \in E} f(sv)$$

is maximized, and for any vertex $v \neq s, t$ in the graph,

$$\sum_{wv \in E} f(wv) = \sum_{vw \in E} f(vw)$$

which is the flow conservation law. After the formalization of linear programming, we see this is just a linear constraint problem over \mathbf{R}^E .

Example. The shortest path problem on an undirected graph with edge costs $c : E \rightarrow \mathbf{R}^+$ asks us to find the path between two fixed vertices s and t which minimizes the sum of the costs of the edges in the path. To obtain a linear program describing the program, we would like to associate each subset of edges as an element $f \in \{0, 1\}^E$. The only problem with this is that no system of linear constraints can enforce that $f(e) \in \{0, 1\}$ for all e , because the set formed is not convex. The problem is we are discussing an **integer linear program**, where we attempt to optimize over a \mathbf{Z} -module rather than an \mathbf{R} -vector space (a mathematically well posed problem, though normally exponentially more complex than a normal linear program, though this problem, an exception of the rule, can be computed using linear programming in subexponential time). In addition to the 0/1 constraint, we desire that for each vertex v ,

$$\sum_{wv \in E} f(wv) \leq 1 \quad \sum_{vw \in E} f(vw) - \sum_{wv \in E} f(wv) = 0 \quad \sum_{vw \in E} f(vw) \leq 1$$

and

$$\sum_{ws \in E} f(ws) = \sum_{tw \in E} f(tw) = 1$$

Once the constraints are specified, we just need to minimize $\sum c(e)f(e)$.

Example. Suppose that we are classifying objects into two categories, based on information contained in a feature vectors v_1, \dots, v_n . To prevent overfitting, we attempt to find a linear classifier – an linear functional λ together with a scalar b such that the object corresponding to v_i is in category A if $\lambda(v_i) \geq b$, and in category B if $\lambda(v_i) < b$. These constraints describe exactly a convex region corresponding to a linear problem, so this problem can be reduced to a finding a **feasible solution** with respect to the linear constraints.

Example. Consider a set X described by constraints of the form $\langle y_i, x \rangle \leq b_i$ for some $y_i, x \in \mathbf{R}^n$ and $b_i \in \mathbf{R}$. The **Chebyshev center** is the largest circle contained within X . A circle can be specified by a center x and a radius r . We of course require that $r > 0$. To verify the circle is contained within X , we require $\langle y_i, x + tv \rangle \leq b_i$ is true for all $t \leq r$ and $\|v\|^2 = 1$. This can be reexpressed as $t\langle y_i, v \rangle \leq b_i - \langle y_i, x \rangle$. Since $t\langle y_i, v \rangle \leq t\|y_i\|$. The Cauchy Schwarz gives us bounds of the form $\langle y_i, tv \rangle \leq t\|y_i\|$, and this bound is tight since we can choose v to be a scalar multiple of y_i . Thus $B_r(x) \subset X$ if and only if $\langle y_i, x \rangle \leq b_i - r\|y_i\|$ for each i . This gives us a linear constraint problem to find the Chebyshev center.

1.1 Piecewise Linear Convex Functions

A function f is **convex** if $f(\lambda a + (1-\lambda)b) \leq \lambda f(a) + (1-\lambda)f(b)$ for every a, b in the domain of f . All linear functions are convex (as well as linear functions plus a constant, known as affine functions), but not all convex functions are linear (for instance, the exponential function is convex, and any differentiable function which satisfies $f'' \geq 0$). If f_1, \dots, f_n are convex functions, then $\max(f_1, \dots, f_n)$ is also a convex function. Thus the maximum of affine functions is convex. This tells us that the absolute value function $|x| = \max(x, -x)$ is convex. A more general optimization problem, given affine functionals $\lambda_i(\cdot) + a_i$, would be to maximize $\max(\lambda_i(v) + a_i)$ subject to some constraint function $Av \leq b$. This turns out to be no more general than linear programming. Note that $\max(\lambda_i(v) + a_i)$ is equal to the minimum value of z such that $z \geq \lambda_i(v) + a_i$ for all i . Thus we just need to add a linear number of constraints of this form, in addition to the new variable z , and then minimize the value of z over the constraint region.

What about if we are just discussing absolute value. For example, consider the problem of maximizing $\sum c_i |x_i|$, such that $Ax \geq b$, and all c_i are non-negative. We could perform the trick in the last problem, but it is quite involved. Instead, we note that $|x_i|$ is the smallest value of z_i such that $x_i \leq z_i$ and $-x_i \leq z_i$. Thus we are minimizing $\sum c_i z_i$, subject to these new constraints. Alternatively, we replace each x_i with non-negative values x_i^+ and x_i^- for each x_i , such that $x_i = x_i^+ - x_i^-$.

Example. Consider the problem of minimizing $2|x| + |y|$, such that $x + y \geq 4$. We introduce new variables a and b , and minimize $2a + b$, such that $-x, x \leq a$, and $-y, y \leq b$. The solution is seen to be unbounded, because we can set $x = 4 - y$, and then as we let $x \rightarrow \infty$, $y \rightarrow -\infty$ and so $2|x| + |y| \rightarrow \infty$. An alternate

and equivalent linear program to optimize is minimize $2x_+ + 2x_- + y_+ + y_-$, such that $x_+ - x_- + y_+ + y_- \geq 4$, $x_+, x_-, y_+, y_- \geq 0$.

Example. Returning to linear optimization, we may have to perform a linear regression on a certain regression sample $(x_1, y_1), \dots, (x_n, y_n)$, with $x_i \in \mathbf{R}^n$ and $y_i \in \mathbf{R}^m$. We wish to find a linear relation of the form $y = Mx + b$, which fits the data best. For each data point (x_i, y_i) , we consider the prediction error $|y_i - (Mx_i + b)|$. The total error is the sum of these errors

$$\sum_i |y_i - (Mx_i + b)|$$

We have discussed how to formulate this problem as a linear program. The quadratic error $\sum_i |y_i - (Mx_i + b)|^2$ can also be discussed, which can be solved by methods of elementary calculus.

Example. Consider a linear discrete dynamical system defined by

$$x_{n+1} = Ax_n + Bz_n \quad y_n = c_n x_n$$

Where $x_n \in \mathbf{R}^n$ is the state of the system, $y_n \in \mathbf{R}^n$ is the output, and x' is a sequence of values we get to choose to control how the system evolves. We may enforce that $Cx_n \leq a$ for some constraint matrix C . Given initial conditions $x_0^i = a_i$, we may want to drive the state of the system to a state $b \in \mathbf{R}^n$ at time N . If this is possible, we might want to minimize the total perturbation of the system. This turns into an optimization problem, trying to minimize $\sum |c_n x_n|$ subject to the constraints described above.

Chapter 2

Polyhedra

We shall now describe a more optimal way to check the extreme points of a polyhedron against a constraint matrix, to find optimal solutions. Given two extreme points on a polyhedron, we say they are adjacent

Definition. Let v and w be two basic solutions in a polyhedron. We say the two vertices are adjacent if we can find $n - 1$ linearly independent constraints that are active at both of them. The line segment joining the two vertices is called an edge.

We now focus on polyhedra in slack form, that is, vectors such that $Mv = b$. How do we find basic solutions in this polyhedron? Choose n linearly independent rows in M . Find solutions to the vector $Mv = b$, where $v_i = 0$ for a row i not in the n linearly independent columns. The elements forced to be zero are non-basic, whereas the others are basic. Any vector solution is a basic solution.

Definition. A polyhedron contains a line if there is a vector v and a non-zero vector w such that the polyhedron contains $v + \lambda w$ for all scalar values λ .

Theorem 2.1. *The following are equivalent:*

1. *The polyhedron does not contain a line*
2. *The polyhedron has at least one extreme point.*
3. *There are n rows in the matrix defining the polyhedron that are linearly independant.*

Proof. We prove the implications in the order they appear:

- (1) \implies (2): We prove that if a polyhedron does not contain a line, it has a basic feasible solution, and therefore must have an extreme point. Let x be a point in the polyhedron...

□

Chapter 3

Duality

The method of Lagrangian multipliers allows us to find the optima of a function f over the level set of some function $g : \mathbf{R}^n \rightarrow \mathbf{R}$ – i.e. $g(x) = c$. For instance, we may wish to find the maximum $x^2 + y^2$, such that $x + y = 1$. If we could find a maximum over all of \mathbf{R}^n , and it happened to lie on the level set of g , then we'd be done, but this almost never happens. However, suppose we can modify f to a new function f_1 , such that it takes on the same values on the level set of g . Then, if the maximum of f_1 lies on the level set of g , we'd be done, since this would also be the maximum of f on the level set of g . The easiest method is to define $f_1(x) = f(x) + p(g(x) - c)$, for some $p \in \mathbf{R}$. We then express the values of $\partial_i f_1(x) = 0$ as functions of p , and then solve for p to guarantee these solutions lie on $g(x) = c$. For instance, we take $f_1(x, y) = x^2 + y^2 + p(x + y - 1)$. Then

$$\partial_x f_1(x, y) = 2x + p \quad \partial_y f_1(x, y) = 2y + p$$

so the optima here occur at $x = y = -p/2$. If we want $x + y = 1$, then we must have $p = -1$, and then we find the maximum occurs at $(1/2, 1/2)$, with maximum value $1/2$.

Let us apply a similar method to linear programming. Consider a problem of maximizing $c^t x$, such that $Ax \leq b$, and $x \geq 0$. We call this the **primal** linear program. We switch to the dual space, and attempt to find y which minimizes $b^t y$, such that $A^t y \geq c$, and $y \geq 0$. This is the **dual** linear program. Note that this really reflects a duality, because the double dual of a linear program is the original linear program.

Theorem 3.1 (Weak Duality). *Any solution satisfying the dual of a linear program bounds the solutions of the original program.*

Proof. Note that if x and y are as in previous notation,

$$c^t x \leq y^t A x \leq y^t b$$

The first part follows because $x \geq 0$, and the last because $y \geq 0$. \square

We can actually prove that the optimum value of the dual is the same as the optimum value of the original program. The proof gives us much more about the relationship between the two programs. First, a fundamental theorem from the theory of convex bodies, which emerges from the Hilbert space theory of finite dimensional vector spaces.

Theorem 3.2. *If $X \subset \mathbf{R}^n$ is closed and convex, and $y \notin X$, then there is an affine functional λ such that $\lambda(X) < 0$, and $\lambda(y) > 0$.*

Proof. Let x be the point in X that is closest to y . Then if $x' \in X$,

$$\begin{aligned} \langle y - x, x' \rangle &= \langle y - x, x \rangle + \|y - x\|^2 - \langle y - x, y - x' \rangle \\ &\leq \langle y - x, x \rangle + \|y - x\| (\|y - x\| - \|y - x'\|) \\ &\leq \langle y - x, x \rangle \end{aligned}$$

and

$$\langle y - x, y \rangle = \|y - x\|^2 + \langle y - x, x \rangle > \langle y - x, x \rangle$$

Hence $a \mapsto \langle y - x, a \rangle - \langle y - x, x \rangle - \varepsilon$ is the affine functional we need, for small enough ε . \square

Geometrically, this theorem says there is a hyperplane separating y and X , not containing y or any points in X , since a hyperplane is just a level set of a linear functional.

A **cone** is a convex subset of \mathbf{R}^n closed under multiplication by non-negative scalars. In the two dimensional case, we obtain the points between some angle, hence the geometric name.

Theorem 3.3. *If X is a non-empty closed cone, and $y \notin X$ there is linear functional λ such that $\lambda(X) \geq 0$, and $\lambda(y) < 0$.*

Proof. Using the last theorem, find an affine function μ such that $\mu(X) > 0$ and $\mu(y) < 0$. If we define $\lambda = \mu - \mu(0)$, then λ is linear, and if $\lambda(y) > 0$, then we find $\lambda(y) > \mu(0) > 0$, which is impossible. \square

Theorem 3.4 (Farkas). *If $A \in M_{n,m}(\mathbf{R})$, $b \in \mathbf{R}^m$, and $c \in \mathbf{R}^n$. Then either there is $x \geq 0$ such that $Ax \leq b$, $y \geq 0$ such that $A^t y \geq 0$ and $b^t y < 0$, but not both.*

Proof. Suppose that both x and y existed. Then we find

$$0 \leq y^t(Ax) \leq y^t b < 0$$

We now show that one of the two conditions must hold. \square

Farkas' lemma gives us strong duality in linear programs, for it relates the slack in the dual program to slack in the original program.

Theorem 3.5 (Strong duality). *The optimal value of the dual of a linear program is equal to the optimal value of the original program.*

Proof. Let y^* be the optimal value of the dual of the linear program which maximizes $\lambda(x)$, such that $Ax \leq b$, and let $v^* = b^t y^*$. We apply Farkas's lemma to conclude that there is either $x \geq 0$ such that $-\lambda(x) \leq -v^*$ and $Ax \leq b$, which implies that $\lambda(x) \geq v^*$ (which would imply equality since we have already proven that $\lambda(x) \leq v^*$), or there is $(y_0, y) \geq 0$ such that $\langle b, y \rangle < \langle v^*, y_0 \rangle$, and if $M = \begin{pmatrix} -c^t \\ A \end{pmatrix}$, then $M^t(y_0, y) \geq 0$. If $y_0 = 0$, then $A^t y \geq 0$ and $b^t y < 0$, so $A^t(y^* + y) \geq c$ and $b^t(y^* + y) < v^*$, contradicting optimality. Similarly, if $y_0 > 0$, then we obtain the same conclusion by some tricky manipulation. This completes the proof. \square

We finalize our discussion of duality by talking about complementary slackness. Let x be a solution of a linear program, and y a solution of its dual. Then x and y are both optimal if for all $1 \leq j \leq n$, $x_j(A^t)_j y - c_j = 0$, and for all $1 \leq i \leq m$, $y_i(A_i x - b_i) = 0$. To see this, we recall the proof of weak duality, and calculate

$$\sum_j c_j x_j \leq \sum_j \sum_i (A^t)_{j,i} y_i x_j = \sum_i \left(\sum_j A_{i,j} x_j \right) y_i \leq \sum_i b_i y_i$$

and the last theorem shows the inequalities are equalities if and only if the condition holds.

Chapter 4

Applications to Combinatorial Optimization

Unfortunately, we normally have to perform combinatorial problems over \mathbf{Z} -modules, a much more intractable problem than over real vector spaces. However, there is a trick, known as unimodularity, which enables us to optimize over the real vector space corresponding to the \mathbf{Z} -module, and then recover results in the original space. Call a matrix unimodular if the determinant of every square submatrix is 0, 1, or -1 .

Theorem 4.1. *If a polyhedron is expressed by $Ax \leq b$ where A is totally unimodular, and $b \in \mathbf{Z}^n$, then every extreme point of the polyhedron is integral.*

Proof. We apply Cramer's rule. If we take a submatrix equality $Bx = c$ to find an extreme point x , then

$$x_i = \frac{\det B_i}{\det B}$$

Then $\det B \in \{-1, 0, 1\}$, and $\det B_i \in \mathbf{Z}$, so that $x \in \mathbf{Z}^n$. □

Thus we may apply linear programming techniques to integer linear programs, provided they are specified by a unimodular constraint.

Example. Let $G = (L \cup R, E)$ be a Bipartite graph. The problem of finding a maximal matching in this graph can be reduced to integer linear programming. Let $V = \mathbf{Z} \cdot E$ be the free module generated by the edges. Consider the linear function λ_v which assigns a value of one to edges containing v as an endpoint,

and assigning zero to other edges. We can also consider the functionals μ_e which maps e to one, and other edges to zero. The bipartite matching problem can then be expressed as maximizing $\sum_e \mu_e(w)$ over V , such that $0 \leq \mu_e(w) \leq 1$ for all e , and $\lambda_v(w) \leq 1$. It turns out that the matrix constraints related to this problem are unimodular, so we may solve the problem using non-integer linear programming in $\mathbf{R} \cdot E$, and then project results back into $\mathbf{Z} \cdot E$, which makes the problem tractable.

Example. The maximum $s - t$ flow problem for a graph (V, E) with capacities c can be represented as a linear program over $\mathbf{R} \cdot E$. We attempt to maximize the operator $\mu_{\text{out}} - \mu_{\text{in}}$, where $\mu_{\text{out}} = \lambda_t$ and $\mu_{\text{in}} = \lambda_s$, using the notation of the last example. The constraints are that $0 \leq \mu_e \leq c(e)$. The constraint matrix is not necessarily unimodular, but here we allow real values, so we can solve the problem in polynomial time using linear programming.

A **network matrix** is a matrix corresponding to a tree (V, E) , with edges oriented arbitrarily (and non-empty). Given a pair of vectors (v, w) , and an edge e , we define $A_{(v,w),e} = 1$ if e appears on the $v - w$ path positively, -1 if it appears negatively, and zero if it doesn't occur on the path. More generally, we can consider any submatrix of a matrix of this form to be a network matrix, for we only need to discuss these matrices to prove unimodularity.

Theorem 4.2. *Every network matrix is totally unimodular.*

Proof. We prove by induction. If we only have one edge, then our network matrix is a column of elements in $\{0, 1, -1\}$, and is therefore unimodular. Inductively, suppose a network matrix A is presented by (V, E) , and a network matrix corresponding to fewer nodes is unimodular. Let B be a submatrix. We may assume E contains more than one edge, because otherwise the matrix is just a columns of elements of $\{0, 1, -1\}$, and is therefore unimodular. Fix some edge $vw \in E$. If the row corresponding to vw and wv does not lie in B , then B is a submatrix of the network matrix corresponding to $(V, E - e)$, and is unimodular by induction. We may without loss of generality assume vw lies in B . If there is no ur such that vw is on the path from u to r , then the row corresponding to vw is equal to zero, and $\det(B) = 0$. Thus we may assume ur exists. By negating columns, we may assume that for any such ur vw occurs positively. If we subtract the column corresponding to ur with every other column where vw occurs, we obtain a new matrix B' with the same determinant as B . What's

more, the row corresponding to vw has value one only on the column corresponding to ur , so we may consider the submatrix C of B' obtained from removing the row vw and column ur , and C has the same determinant of B' , thought perhaps multiplied by -1 . Now we may apply induction to C , since it doesn't contain vw or wv , to conclude that the determinant of B is in $\{0, -1, 1\}$. \square

It is very useful in the intersection of linear programming and combinatorial optimization to be able to link intuitions about the specific problems at hand (in this case graph theory) with the intuitions about linearity from in the programming itself.

Example. *The Bipartite problem has unimodular constraints. To see this, we need only note that each row of the constraint matrix corresponds to a row of the network matrix corresponding to some tree in $G = (L \cup R, E)$. To prove this, we may assume none of the rows corresponding to the constraints $0 \leq \mu_e(w) \leq 1$ are in the matrix, since the row is just a single value of 1. Thus we take n vertices v_1, \dots, v_n and consider certain columns corresponding to edges e_1, \dots, e_n . If the edges have a cycle $e_{i_1}, e_{i_2}, \dots, e_{i_k}$, then by summing the corresponding columns by alternately adding and subtracting the columns, we end up with a zero vector, and so the determinant of the submatrix is zero. Thus the*

The Gouila-Houri theorem is also useful for verifying unimodularity.

Theorem 4.3. *A matrix $A \in M_{m,n}(\mathbf{R})$ is unimodular if and only if for every subset of rows can be partitioned into two subset v_1, \dots, v_l and w_1, \dots, w_k such that $\sum v_i - \sum w_i \in \{-1, 0, 1\}^n$.*

Chapter 5

The Ellipsoid Method

Recall that every Ellipsoid in \mathbf{R}^n at the origin is the unit ball relative to some positive definite inner product on \mathbf{R}^n . The classification theorem on inner product tells us that we can identify every ellipsoid in \mathbf{R}^n with a positive definite matrix $M > 0$ and $x \in \mathbf{R}^n$. To find a polynomial time algorithm for linear programming, we shall shrink ellipsoids around polyhedra to find optimal points.

Thus, given $M > 0$, $x, a \in \mathbf{R}^n$, we wish to find a new ellipsoid $E(N, y)$, which shrinks around a , such that

$$E(N, y) \supset E(M, x) \cap \{z \in \mathbf{R}^n : a_z \geq a_x\}$$

and the volume of $E(N, y)$ is less than or equal to $\sqrt{\alpha}$ times the volume of $E(M, x)$, for small enough α (we shall find $\alpha = e^{-1/2n}$). Since the determinant measures the volume change of the respective linear transformation, the volume condition algebraically means $\det(N) \leq \alpha \det(M)$.

Define

$$b = \frac{Ma}{\|a\|_M} \quad N = \frac{n^2}{n^2 - 1} \left(M - \frac{2}{n+1} b \cdot b^t \right)$$

First, we verify that N is positive definite. Note

$$M - \frac{2}{n+1} b b^t = M - \frac{2}{n+1} \frac{M a a^t M^t}{\|a\|_M^2}$$

Part II

Convex Optimization

It is often the case that the functions we attempt to maximize are non-linear, or the set we are optimizing over is not a polyhedron. Convex optimization is the safest step to generalizing the linear programming optimization problem, for it is still a semi-tractable problem.

Chapter 6

Convex Functions

Let $D \subset \mathbf{R}^n$ be a convex subset of euclidean space. A function $f : D \rightarrow \mathbf{R}$ is **convex** if the graph of f lies below the lines between any two points on the graph. Analytically, this can be expressed by the inequality

$$f(\lambda x + (1 - \lambda)y) \leq \lambda f(x) + (1 - \lambda)f(y)$$

Or equivalently, the epigraph $\{(x, y) : f(x) \leq y\}$ of the function is convex.

The inner product on \mathbf{R}^n gives us lots of information about the convex subsets of \mathbf{R}^n . Indeed, the Banach and Hilbert space structure of \mathbf{R}^n tells us nice separation theorems about convex sets.

Theorem 6.1. *If $D \subset \mathbf{R}^n$ is convex and closed, and $x_0 \in \mathbf{R}^n - D$, then there is $y \in \mathbf{R}^n$ such that for any $x \in D$,*

$$\langle y, x \rangle < \langle y, x_0 \rangle$$

Geometrically, there is a hyperplane separating D and x_0 . If $x_0 \in \partial D$, then there is a nonzero $y \in \mathbf{R}^n$ such that for $x \in D$,

$$\langle y, x \rangle \leq \langle y, x_0 \rangle$$

So there is a hyperplane supporting the convex set in \mathbf{R}^n .

If f is a function defined on a convex set, then v is a **subgradient** of f at x if the hyperplane defined at x by v lies below the graph of f . Analytically, for any y ,

$$f(x) - f(y) \geq \langle v, x - y \rangle$$

The set of subgradients of f at x is denoted $\partial f(x)$. It is a pleasant way to obtain a form of differentiation of a convex function that need not be convex.

Theorem 6.2. *If a convex f is differentiable at x , then $\partial f(x)$ contains a single element.*

Proof. Write $f(x+h) = f(x) + \langle \nabla f(x), h \rangle + O(\|h\|^2)$, so

$$f(x+h) - f(x) = \langle \nabla f(x), h \rangle + O(\|h\|^2)$$

Now if $v \in \partial f(x)$, this implies that

$$|\langle v - \nabla f(x), h \rangle| = O(\|h\|^2)$$

But this implies $v = \nabla f(x)$, before linear functions cannot approach the origin quadratically. Indeed, find $x_0 \in S^{n-1}$ with

$$|\langle \nabla f(x) - v, x_0 \rangle| = \|\nabla f(x) - v\|$$

Then

$$\frac{\langle \nabla f(x) - v, \lambda x_0 \rangle}{|\lambda|^2} = \frac{\|\nabla f(x) - v\|}{|\lambda|}$$

which must vanish as $\lambda \rightarrow 0$, impossible unless $\|\nabla f(x) - v\| = 0$.

Conversely, we verify that $\nabla f(x)$ is a subgradient for f . We use the convexity of f to write

$$f(y) \geq \frac{f((1-\lambda)x + \lambda y) - (1-\lambda)f(x)}{\lambda} = f(x) + \frac{f((1-\lambda)x + \lambda y) - f(x)}{\lambda}$$

By the convexity in one dimension, the values on the right converge monotonically to $\langle \nabla f(x), y \rangle$, and this completes the proof. \square

Theorem 6.3. *If $\partial f(x)$ is non-empty for each x , then f is convex. Conversely, if f is convex, $\partial f(x)$ is non-empty on the interior of the domain of f .*

Proof. If $v \in \partial f((1-\lambda)x + \lambda y)$, then

$$f((1-\lambda)x + \lambda y) \leq f(x) - \lambda \langle v, y - x \rangle$$

$$f((1-\lambda)x + \lambda y) \leq f(y) + (1-\lambda) \langle v, y - x \rangle$$

If we multiply the first equation by λ , the second by $(1 - \lambda)$, and then by adding the two equations together, we obtain the convexity inequality.

Now let $x \in D$. Then $(x, f(x))$ is on the boundary of the epigraph of f , so using the supporting hyperplanes theorem, there is $(a, b) \in \mathbf{R}^n \times \mathbf{R}$ such that for any (y, t) , with $t \geq f(y)$, with

$$\langle a, x \rangle + bf(x) \leq \langle a, y \rangle + bt$$

Letting $t \rightarrow \infty$, we conclude $b \leq 0$. If $b = 0$, and $x \in D^\circ$, then we can let $y = x - \varepsilon a$, for which we conclude that

$$\langle a, x \rangle \leq \langle a, x - \varepsilon a \rangle$$

so that $-\varepsilon \|a\|^2 \geq 0$, an impossibility unless $a = 0$, which is impossible due to the separating hyperplanes theorem. Thus $b \neq 0$, and we may divide out, choosing $t = f(y)$, we find

$$f(y) - f(x) \geq \langle v/b, y - x \rangle$$

So $v/b \in \partial f(x)$. □

Why are convex functions so interesting in optimization theory? Optimizing an arbitrary general convex function is *not* easy in general, but if we have a pleasant convex domain, the problem becomes much easier. First, note that if we have an arbitrary function $f : D \rightarrow \mathbf{R}$, then we can write a minimization problem $\inf_{x \in D} f(x)$ as minimizing the linear function $\inf_{(x,t) \in \text{epi}(f)} t$ over the epigraph of f , so regardless of the qualities of f , we can view a minimization of any function as a minimization of a linear functional on your set. Next, note that minimizing a function over D is the same as minimizing the function over the convex hull of D . However, obtaining a solution on D from a solution over the convex hull of D is normally very difficult, for it involves, given a set $y \in D' \supset D$, finding $x_1, \dots, x_n \in D$ for which $y = \sum a_i x_i$, with $a_i \geq 0$ and $\sum a_i = 1$. This is a very difficult problem algorithmically.

Chapter 7

Mirror Descent

Mirror descent is a way to apply Gradient descent in space with arbitrary norms, that is, on a Banach space rather than a Hilbert space. We cannot apply standard gradient descent to a function f , because for Banach spaces the gradients (Fréchet derivatives) $\nabla f(x)$ lie in the dual space of the Banach space. On a Hilbert space, there is a canonical way of identifying $\nabla f(x)$ with a vector in the original space, so we can consider the iterative procedure $x \mapsto x - \eta \nabla f(x)$. The main idea of mirror descent is to generalize this procedure to Banach spaces by mapping the Banach space into its dual, finding some way of performing gradient descent in the dual space of the Banach space, and then mapping the space back into the original space, in such a way that we truly obtain a minimizing procedure.

Suppose we are optimizing over a compact, convex subset $D \subset \mathbf{R}^n$ contained in the closure of an open convex subset U , for which $U \cap D$ is non-empty. Now consider a convex function Φ defined on U , which will enable us to map D into the dual space by the map $x \mapsto \nabla \Phi(x)$, the ‘mirror map’. Provided that $\nabla \Phi(U) = \mathbf{R}^n$, we can always reverse the mirror map. If we consider the iteration

$$x \mapsto (\nabla \Phi)^{-1}(\Phi(x) - \eta \nabla f(x))$$

where $(\nabla \Phi)^{-1}(\Phi(x) - \eta \nabla f(x))$ is a point $y \in U$ for which

$$\nabla \Phi(y) = \Phi(x) - \eta \nabla f(x)$$

Now this point may not lie in D , so we define the Bregman divergence

$$D_\Phi(x, y) = \Phi(x) - \Phi(y) - \langle \nabla \Phi(y), x - y \rangle$$

and project y onto a point x' in D which minimizes the divergence $D_\Phi(x', y)$, for $x' \in D \cap U$. This projection exists provided that

$$\lim_{x \rightarrow \partial U} \|\nabla \Phi(x)\| = \infty$$

which implies $x \mapsto D_\Phi(x, y)$ is locally increasing on the boundary of U , and if $D_\Phi(x, y)$ is strongly convex, then the existence and uniqueness of a projection is guaranteed (I don't have a proof right now, so let's just assume this is true).

Now we describe mirror descent. Let x_1 minimize $\Phi(x)$ on $U \cap D$. Then, once x_n is specified, let y_n be defined by $\nabla \Phi(y_n) = \nabla \Phi(x_n) - \eta v$, where $v \in \partial f(x_n)$ and $x_{n+1} \in \operatorname{argmin} D_\Phi(x, y_n)$.

Theorem 7.1. *If Φ is strongly convex on $D \cap U$ with respect to the norm $\|\cdot\|$, that is*

$$\Phi(x) - \Phi(y) \leq \langle \nabla \Phi(x), x - y \rangle - \frac{m}{2} \|y - x\|^2$$

Define $\lambda^2 = \sup_{x \in U \cap D} \Phi(x) - \Phi(x_1)$, let f be convex, and let L be Lipschitz with respect to $\|\cdot\|$. If we let $\eta = (\lambda/L)\sqrt{2a/K}$, then

$$f\left(\frac{1}{K} \sum_{k=1}^K x_k\right) - \min_{x \in D} f(x) \leq \lambda L \sqrt{\frac{2}{mK}}$$

Proof. To Be Continued...

□