

INSTITUTO TECNOLÓGICO AUTÓNOMO DE MÉXICO



Ecuaciones lineales diofantinas
aplicadas a
programas lineales enteros

TESIS

PARA OBTENER EL TÍTULO DE

LICENCIADO EN MATEMÁTICAS APLICADAS

PRESENTA

IÑAKI SEBASTIÁN LIENDO INFANTE

ASESOR

DR. ANDREAS WACHTEL

Agradecimientos

Resumen

El algoritmo de Ramificación y Acotamiento (R&A) es el estándar para resolver programas lineales enteros. Este método se basa en el famoso paradigma de división y conquista, el cual combina las soluciones de subproblemas más pequeños a fin de que se obtenga una solución del problema original. Los problemas se estructuran de manera que generen un árbol: el problema original genera una colección de subproblemas, y luego cada subproblema genera su propia colección de subsubproblemas, etcétera.

En la sección 1.1.2 describimos cómo es que Ramificación y Acotamiento cuenta con reglas o políticas de poda para evitar resolver todos los subproblemas, pues de manera contraria el número de nodos a recorrer crece exponencialmente. No obstante, R&A jamás podará subárboles que contengan una posible solución. Por lo tanto, las políticas de poda operan de manera subóptima siempre que exista una gran cantidad de posibles soluciones distribuidas en subárboles disjuntos. Si, en el peor de los casos, cada hoja del árbol contiene una solución, entonces R&A deberá resolver todos los subproblemas.

Se ha observado que esto último ocurre siempre que el vector objetivo es ortogonal a una de las restricciones del programa lineal entero. En efecto, el programa relajado cuenta con una infinidad de soluciones, por lo que todo subproblema tendrá al menos una solución, lo cual implica que las políticas resultan ser ineficientes. La instancia más simple que exhibe estas

ineficiencias en las políticas de poda es

$$\max_{\mathbf{x} \in \mathbb{Z}^n} \{\mathbf{p}^T \mathbf{x} : \mathbf{p}^T \mathbf{x} \leq u, \mathbf{x} \geq \mathbf{0}\}. \quad (0.1)$$

En este caso, la restricción $\mathbf{p}^T \mathbf{x} \leq u$ es ortogonal al vector objetivo.

En esta tesis analizamos a profundidad el problema anterior para desarrollar un nuevo método que nos permita resolver de manera más eficiente este tipo de instancias. Por la misma estructura de este problema, supondremos sin pérdida de generalidad que \mathbf{p} no tiene entradas nulas. Ciertamente, nos desharemos de esta suposición una vez que introduzcamos el caso de múltiples restricciones en el capítulo 4.

De manera resumida, mostramos que existe una equivalencia entre resolver problemas del tipo (0.1) y resolver ecuaciones lineales diofantinas en n incógnitas. Los coeficientes de las ecuaciones lineales diofantinas estarán dadas por las entradas de un vector \mathbf{q} asociado a \mathbf{p} . Así también, el número de ecuaciones que deberemos resolver depende, en gran medida, de los signos en las entradas de \mathbf{q} . El teorema 1.2.9 muestra que si alguna entrada q_i es negativa, entonces es necesario resolver una sola ecuación y, en caso de que todas las entradas de \mathbf{q} sean positivas, el número de ecuaciones a resolver es finito.

El capítulo 1 presenta los prerrequisitos necesarios para obtener los resultados que se encuentran a lo largo de esta tesis. Definimos una clase de vectores a la cual supondremos que \mathbf{p} pertenece y obtendremos varias de sus propiedades. Sin duda, esta clase de vectores contiene cualquier vector representable en aritmética finita por lo que, en la práctica, este supuesto es razonable. Los resultados que más destacan en esta parte de la tesis son, en opinión del autor, los teoremas 1.2.9 y 1.2.15, así como el corolario 1.2.16.

El capítulo 2 analiza el caso en el que el vector \mathbf{q} asociado a \mathbf{p} tiene una entrada negativa. Bajo esta hipótesis adicional, la solución del problema (0.1) se obtiene al resolver una sola ecuación lineal diofantina. Por un lado,

mostramos que el valor objetivo del problema (0.1) se puede determinar de manera inmediata sin tener conocimiento de la solución óptima. Por el otro lado, presentamos un algoritmo que construye la solución óptima y cuya complejidad es polinomial. Finalmente, realizamos una serie de experimentos numéricos que permiten comparar los tiempos de terminación de nuestro algoritmo con los de Ramificación y Acotamiento.

El capítulo 3 analiza el caso en el que el vector \mathbf{q} asociado a \mathbf{p} tiene entradas estrictamente positivas. Bajo esta hipótesis solamente podemos asegurar la finitud del número de ecuaciones lineales diofantinas que debemos resolver para encontrar la solución de (0.1). No obstante, mostramos que si el lado derecho de la restricción $\mathbf{p}^T \mathbf{x} \leq u$ es suficientemente grande, entonces sí basta con resolver una sola ecuación lineal diofantina para obtener el óptimo. De manera incidental, a través de las herramientas desarrolladas en este capítulo, encontramos también nuevas cotas superiores para el Problema de la Moneda, también conocido como el problema de Frobenius. Además, presentamos un algoritmo que construye la solución óptima de (0.1) bajo el supuesto $\mathbf{q} > \mathbf{0}$. Al igual que el capítulo 2, realizamos una serie de experimentos numéricos que permiten comparar los tiempos de terminación de nuestro algoritmo con los de Ramificación y Acotamiento.

El capítulo 4 introduce el caso de múltiples restricciones. Observamos en este capítulo que la división en casos del teorema 1.2.9 deja de ser vigente. Desarrollamos, un nuevo método que permite resolver este tipo de problemas. Es decir, exhibimos una manera de resolver programas lineales enteros bajo la perspectiva de la búsqueda de soluciones de sistemas de ecuaciones lineales diofantinas. Por lo tanto, es necesario introducir nueva maquinaria para resolver este tipo de sistemas de ecuaciones. Puesto que este nuevo análisis requerido sería demasiado grande para añadirlo a la tesis, el autor prefirió ser más informal en su exposición. Ciertamente, este capítulo sirve como directriz inicial para la realización de futuras investigaciones.

Índice general

| | |
|---|------------|
| 1. Aspectos Teóricos | 1 |
| 1.1. Prerrequisitos | 2 |
| 1.2. Fundamentos | 13 |
| 2. El caso infinito | 39 |
| 2.1. Experimentos numéricos | 46 |
| 3. El caso finito | 49 |
| 3.1. Análisis de capas enteras | 50 |
| 3.2. Construcción de soluciones | 71 |
| 3.3. Experimentos numéricos | 78 |
| 4. Múltiples restricciones | 86 |
| A. Algoritmo de Ramificación y Acotamiento | 101 |

Capítulo 1

Aspectos Teóricos

Este capítulo presenta los prerequisites necesarios para obtener los resultados que se encuentran a lo largo de esta tesis. En primer lugar, la sección 1.1 recopila resultados básicos de teoría de números y de programación lineal. Estas herramientas forman parte de la literatura tradicional y constituyen lo mínimo necesario para la derivación de nuestros propios resultados. En segundo lugar, la sección 1.2 presenta enunciados y definiciones obtenidos de [BH09], los cuales utilizaremos para continuar con la construcción de nuestros propios resultados que, en pleno conocimiento del autor, son originales.

Como mencionamos en la motivación de esta tesis, nos concentraremos casi exclusivamente en problemas de programación lineal entera del tipo

$$\max_{\mathbf{x} \in \mathbb{Z}^n} \mathbf{p}^T \mathbf{x}, \quad (1.1a)$$

$$\text{s.a. } \mathbf{p}^T \mathbf{x} \leq u, \quad (1.1b)$$

$$\mathbf{x} \geq \mathbf{0}.$$

En la sección 1.2 analizaremos a profundidad este problema, cuyo punto de culminación será el teorema 1.2.9. De manera resumida, el análisis de este problema se divide en dos casos, dependiendo de los signos de las entradas de un vector \mathbf{q} asociado a \mathbf{p} .

Observemos que, para este problema, podemos suponer sin pérdida de generalidad que todas las entradas de \mathbf{p} son distintas de cero. En efecto, si alguna entrada p_i es cero, podremos decir que $x_i = 0$. Cuando introduzcamos en el capítulo 4 múltiples restricciones nos desharemos de este supuesto.

1.1. Prerrequisitos

El autor consideró pertinente no incluir demostraciones en esta sección, pues los enunciados son mostrados en cualquier clase de álgebra superior, programación lineal, o investigación de operaciones. Las referencias principales para las subsecciones de teoría de números y de programación lineal son [Lav14] y [Oli17], respectivamente.

1.1.1. Teoría de Números

Máximo común divisor y mínimo común múltiplo

Definición 1.1.1. Dados dos enteros a y b , decimos que a **divide a** b (y escribimos $a \mid b$) si existe un entero k tal que $a = k \cdot b$. También denotamos por $D(a)$ al **conjunto de divisores de** a , es decir, definimos

$$D(a) := \{b \in \mathbb{Z} : b \mid a\}.$$

Observación. Si a es distinto de cero, entonces $D(a)$ es finito. En efecto, si $b \mid a$ y $a \neq 0$, es posible mostrar que $|b| \leq |a|$, lo cual implica que $|D(a)| \leq 2|a|$. En caso de que a sea nulo, se sigue que $D(a) = \mathbb{Z}$.

Observación. Para cualquier entero a se satisface $\{-1, 1\} \subseteq D(a)$, pues $a = a \cdot 1$ y también $a = (-a) \cdot (-1)$.

Definición 1.1.2. Sean a_1, \dots, a_n enteros no todos iguales a cero, entonces definimos su **máximo común divisor** d como el elemento maximal del conjunto $\bigcap_{i=1}^n D(a_i)$, y escribimos $d = \text{mcd}\{a_1, \dots, a_n\}$. Si $d = 1$, entonces decimos que a_1, \dots, a_n son **coprimos**.

Puesto que alguna entrada a_i es distinta de cero en la definición anterior, encontramos que el conjunto $\bigcap_{i=1}^n D(a_i)$ es finito y también es no vacío (ver observación anterior), por lo que este conjunto tiene un elemento maximal. Es decir, el máximo común divisor d siempre está bien definido.

Observación. El máximo común divisor siempre es positivo, pues se cumple que $1 \in D(a)$ para todo entero a , lo que implica por maximalidad que $1 \leq \text{mcd}\{a_1, \dots, a_n\}$ para cualquier colección de enteros no todos nulos.

La definición más común del máximo común es dada de manera inductiva. Decimos que d es el máximo común divisor de dos enteros a_1, a_2 , no ambos iguales a cero, si se satisface

1. $d \mid a_1$ y $d \mid a_2$, y también,
2. si $d' \mid a_1$ y $d' \mid a_2$, entonces $d' \mid d$.

Luego, para un conjunto de enteros $a_1, a_2 \dots a_n$, no todos iguales a cero, definimos el máximo común divisor entre ellos a partir de

$$\text{mcd}\{a_1, a_2, \dots, a_{n-1}, a_n\} := \text{mcd}\{a_1, \text{mcd}\{a_2, \dots, \text{mcd}\{a_{n-1}, a_n\}\}\}.$$

Sin embargo, debemos ser cuidadosos con esta manera de definir las cosas, pues puede ser el caso, por ejemplo, que a_{n-1} y a_n sean ambos cero y entonces $\text{mcd}\{a_{n-1}, a_n\}$ no está bien definido.

Para que esta manera de definir el máximo común divisor sea equivalente a la definición 1.1.2, deberemos presuponer o bien que $a_{n-1} \neq 0$ o bien que $a_n \neq 0$. A partir de este punto usaremos ambas definiciones de manera indistinta. Independientemente de qué definición usemos, la manera

de calcular el máximo común divisor siempre es a través del Algoritmo de Euclides.

Observación. No porque una colección de enteros sea coprima se sigue que estos enteros son coprimos a pares. Por ejemplo, los enteros 1, 3 y 3 son coprimos pero evidentemente 3 y 3 no lo son.

Definición 1.1.3. Decimos que $c \in \mathbb{Z}$ es una **combinación lineal entera** de un conjunto de enteros a_1, \dots, a_n si existen enteros x_1, \dots, x_n tales que $c = a_1x_1 + \dots + a_nx_n$. Si c es positivo, también decimos que esto último es una **combinación lineal entera positiva**.

Teorema 1.1.4. Sea d un entero y sean a_1, \dots, a_n una colección de enteros no todos iguales a cero. Entonces $d = \text{mcd}\{a_1, \dots, a_n\}$ si y solo si d es la mínima combinación lineal entera positiva de a_1, \dots, a_n .

Ejemplo 1.1.5. El máximo común divisor d de los enteros $a_1 := 2$, $a_2 := 3$ y $a_3 := 5$ es 1 y además se cumple que $-3a_1 - a_2 + 2a_3 = 1 = d$.

Lema 1.1.6. Si $d = \text{mcd}\{a_1, \dots, a_n\}$, entonces $\text{mcd}\{\frac{a_1}{d}, \dots, \frac{a_n}{d}\} = 1$.

Además del máximo común divisor, haremos uso del mínimo común múltiplo, aunque será en menor medida.

Definición 1.1.7. Definimos el **conjunto de múltiplos** de un entero a como

$$M(a) := \{x \in \mathbb{Z} : a \mid x\}.$$

También definimos el **mínimo común múltiplo** de un conjunto de enteros a_1, \dots, a_n , no todos iguales a cero, como el elemento minimal del conjunto $\mathbb{Z}_{\geq 0} \cap \bigcap_{i=1}^n M(a_i)$. Escribimos $\text{mcm}\{a_1, \dots, a_n\}$ para denotar a este mínimo común múltiplo.

Observación. Si a es nulo, entonces $M(a) = \{0\}$. En caso contrario encontramos que $M(a)$ es un conjunto infinito.

Para mostrar que el múltiplo común múltiplo está bien definido, basta observar que el producto $|a_1 \cdots a_n|$ es no negativo y también es un elemento de $M(a_i)$ para toda $i \in \{1, \dots, n\}$.

Ecuaciones lineales diofantinas

Sea $c \in \mathbb{Z}$ y sean a_1, \dots, a_n enteros. Una ecuación lineal diofantina es una ecuación donde deseamos determinar enteros x_1, \dots, x_n que satisfagan

$$a_1x_1 + \cdots + a_nx_n = c.$$

Será de nuestro interés en las siguientes secciones resolver iterativamente este tipo de ecuaciones. Por el momento basta mencionar que podemos enfocarnos en el caso $n = 2$ sin ninguna pérdida de generalidad. No obstante, los resultados se mantienen para cualquier $n \in \mathbb{N}$.

Los siguientes enunciados abordan el problema de determinar la existencia de soluciones para este tipo de ecuaciones, así como de construir estas soluciones.

Teorema 1.1.8 (Existencia). *Sean a y b enteros, no ambos iguales a cero. Entonces la ecuación lineal diofantina $ax + by = c$ tiene solución entera si y solo si $\text{mcd}\{a, b\} \mid c$.*

Para construir el conjunto de soluciones a una ecuación lineal diofantina, primero encontramos una solución particular.

Definición 1.1.9. Sea $d := \text{mcd}\{a, b\}$ y sean x', y' enteros tales que $ax' + by' = d$ (c.f. teorema 1.1.4). Decimos entonces que x', y' son **coeficientes de Bézout** asociados a a y b , respectivamente¹.

¹Los coeficientes de Bézout se pueden calcular a través del Algoritmo Extendido de Euclides. Véase https://en.wikipedia.org/wiki/Extended_Euclidean_algorithm.

Observación. Los coeficientes de Bézout asociados a un par de enteros no son únicos. En efecto, si x', y' son coeficientes de Bézout de a y b , entonces $x' + b, y' - a$ también lo son:

$$a(x' + b) + b(y' - a) = ax' + by' + ab - ab = ax' + by' = d.$$

Para fines de esta tesis basta la existencia de estos coeficientes, por lo que decimos de manera indistinta “los coeficientes de Bézout” y “una elección de coeficientes de Bézout”.

Definamos $d := \text{mcd}\{a, b\}$ y supongamos que la ecuación $ax + by = c$ tiene solución. Por el teorema 1.1.8, se sigue que $d \mid c$, y entonces existe $c' \in \mathbb{Z}$ tal que $c = c' \cdot d$. Sean x', y' los coeficientes de Bézout asociados a a, b respectivamente. Así,

$$a(c' \cdot x') + b(c' \cdot y') = c'(ax' + by') = c'd = c,$$

por lo que $(c' \cdot x', c' \cdot y')$ es una solución particular de la ecuación $ax + by = c$.

Teorema 1.1.10 (Construcción). *Sea (x_0, y_0) una solución particular de la ecuación lineal diofantina $ax + by = c$. Entonces todas las soluciones enteras de aquella ecuación están dadas por*

$$\begin{cases} x = x_0 + \frac{b}{d}t, \\ y = y_0 - \frac{a}{d}t, \end{cases} \quad (1.2)$$

donde $d := \text{mcd}\{a, b\}$ y $t \in \mathbb{Z}$ es una variable libre.

Ejemplo 1.1.11. Consideremos la ecuación lineal $2x + 3y = 5$. Los coeficientes de Bézout asociados a 2 y 3 son, respectivamente, -1 y 1. Luego, una solución particular para la ecuación es $(x_0, y_0) = (-5, 5)$. Por el teorema

anterior encontramos que todas las soluciones están dadas por

$$\begin{cases} x = -5 + 3t, \\ y = 5 - 2t, \end{cases}$$

donde $t \in \mathbb{Z}$ es una variable libre. En efecto, sustituyendo obtenemos

$$2(-5 + 3t) + 3(5 - 2t) = -10 + 15 + 6t - 6t = 5.$$

1.1.2. Programación lineal

La programación lineal se encarga de resolver problemas de optimización de la forma

$$\max_{\mathbf{x}} \{\mathbf{c}^T \mathbf{x} : \mathbf{x} \in P\}, \quad (1.3)$$

donde P es un poliedro.

En esta sección repasamos brevemente propiedades del poliedro P al cual llamamos **región factible**. Así también, indicamos dónde se encuentra el óptimo del problema (1.3) y hacemos mención rápida sobre cómo obtenerlo. Finalmente, nos enfocamos en programas lineales enteros y, más importantemente, describimos cómo funciona el algoritmo de Ramificación y Acotamiento para encontrar soluciones de estos programas enteros.

Definición 1.1.12. Sea $\mathbf{a} \in \mathbb{R}^n$ un vector no nulo y sea $b \in \mathbb{R}$ un escalar. Llamamos **hiperplano afino** al conjunto de vectores $\mathbf{x} \in \mathbb{R}^n$ que satisfacen $\mathbf{a}^T \mathbf{x} = b$. Así también, llamamos **semi-espacios afinos** a los conjuntos de vectores $\mathbf{x}, \mathbf{y} \in \mathbb{R}^n$ que satisfacen $\mathbf{a}^T \mathbf{x} \geq b$ y $\mathbf{a}^T \mathbf{y} \leq b$.

Definición 1.1.13. Sea $A \in \mathbb{R}^{m \times n}$ una matriz con renglones linealmente independientes y $\mathbf{b} \in \mathbb{R}^m$ un vector. Entonces al conjunto definido por

$$P := \{\mathbf{x} \in \mathbb{R}^n : A\mathbf{x} \geq \mathbf{b}\} \quad (1.4)$$

lo llamamos **poliedro**. Si, además, P es acotado, entonces decimos que P es un **politopo**.

Observación. Todo poliedro P definido de esta manera representa la intersección de m semi-espacios afines. Esto se debe a que $A\mathbf{x} \geq \mathbf{b}$ si y solo si $\mathbf{a}_i^T \mathbf{x} \geq b_i$ para toda $1 \leq i \leq m$ y donde \mathbf{a}_i^T representa el i -ésimo renglón de la matriz A . En la Figura 1.1 se muestra visualmente esta relación entre hiperplanos afines y poliedros.

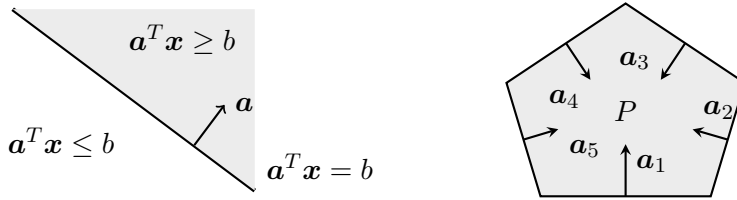


Figura 1.1.: *Izquierda:* Un hiperplano afino $\{\mathbf{x}: \mathbf{a}^T \mathbf{x} = b\}$ junto con los dos semi-espacios que induce. *Derecha:* Un politopo P .

Definición 1.1.14. Sea P un poliedro. Decimos que el vector $\mathbf{x} \in P$ es un **vértice** de P si existe $\mathbf{c} \in \mathbb{R}^n$ de manera que $\mathbf{c}^T \mathbf{x} < \mathbf{c}^T \mathbf{y}$ para todo $\mathbf{y} \in P \setminus \{\mathbf{x}\}$.

En términos gráficos, decimos que \mathbf{x} es un vértice si se satisfacen dos condiciones: en primer lugar, existe un hiperplano afino que pasa por \mathbf{x} y uno de sus semi-espacios inducidos contiene completamente al poliedro P ; en segundo lugar, ningún otro punto de P se encuentra sobre este hiperplano.

Definición 1.1.15. Sea P un poliedro y sea $\mathbf{c} \in \mathbb{R}^n$ un vector. Todo problema de optimización de la forma (1.3) entra en una de las siguientes tres categorías:

1. El valor óptimo no existe: ningún vector $\mathbf{x} \in \mathbb{R}^n$ satisface el sistema de desigualdades $A\mathbf{x} \geq \mathbf{b}$. Es decir, la región factible es vacía.

2. El valor óptimo existe y es infinito: el poliedro P no es acotado y somos capaces de encontrar una sucesión de vectores $\{\mathbf{x}_k\}_{k \in \mathbb{N}}$ en el poliedro P que satisface $\mathbf{c}^T \mathbf{x}_{k+1} > \mathbf{c}^T \mathbf{x}_k$ para todo $k \in \mathbb{N}$.
3. El valor óptimo existe y es finito: este caso es la negación de los dos casos anteriores, pero cabe recalcar que esto no significa que el poliedro P es acotado.

En el primer caso decimos que **el problema es infactible**, mientras que en los últimos dos decimos que **el problema es factible**. También diremos comúnmente del segundo caso que **el problema es no acotado**.

Es posible mostrar que todo poliedro $P := \{\mathbf{x} \in \mathbb{R}^n : A\mathbf{x} \geq \mathbf{b}\}$ puede ser transformado a la forma estándar

$$\{(\mathbf{x}^+, \mathbf{x}^-, \mathbf{s}) \in \mathbb{R}^{n+n+m} : A(\mathbf{x}^+ - \mathbf{x}^-) - \mathbf{s} = \mathbf{b}, (\mathbf{x}^+, \mathbf{x}^-, \mathbf{s}) \geq \mathbf{0}\},$$

de manera que todo problema de optimización de la forma (1.3) puede ser escrito sin pérdida de generalidad como

$$\max_{\mathbf{x} \in \mathbb{R}^n} \quad \mathbf{c}^T \mathbf{x}, \tag{1.5a}$$

$$\text{s.a.} \quad A\mathbf{x} = \mathbf{b}, \tag{1.5b}$$

$$\mathbf{x} \geq \mathbf{0}.$$

De ahora en adelante nuestro análisis se concentrará exclusivamente en problemas lineales de este tipo. Es decir, supondremos, sin pérdida de generalidad, que todo problema lineal se encuentra en esta forma estándar.

Teorema 1.1.16. *Sea P un poliedro que tiene al menos un vértice, consideremos el problema (1.5), y supongamos que el valor óptimo z^* existe y es finito. Entonces el conjunto de soluciones óptimas contiene al menos un vértice de P .*

Este teorema fundamental constituye el primer paso para la construcción de varios algoritmos que encuentran soluciones del problema (1.5). Ciertamente, el más famoso de todos es el algoritmo simplex, el cual “salta” de vértice en vértice hasta llegar a uno con valor óptimo. Otros, más modernos y conocidos como métodos de puntos interiores, comienzan en el interior del poliedro P y son “atraídos” como imanes a uno de los vértices con valor óptimo. No es el objetivo de esta tesis exponer la maquinaria matemática detrás de estos algoritmos².

Ahora describimos brevemente los programas lineales enteros y pasamos a explicar el método de Ramificación y Acotamiento. Por ello, lo que se encuentra a continuación supone que contamos con un algoritmo para resolver problemas del tipo (1.5).

Definición 1.1.17. Sea $A \in \mathbb{R}^{m \times n}$ una matriz con renglones linealmente independientes y sea $\mathbf{b} \in \mathbb{R}^m$ un vector. Al problema de optimización lineal (1.5) lo llamamos **problema relajado** del programa lineal entero

$$\max_{\mathbf{x} \in \mathbb{Z}^n} \mathbf{c}^T \mathbf{x}, \quad (1.6a)$$

$$\text{s.a. } A\mathbf{x} = \mathbf{b}, \quad (1.6b)$$

$$\mathbf{x} \geq \mathbf{0}.$$

Resalta el hecho de que la formulación de un programa lineal entero es idéntico a su formulación relajada, solamente agregamos la restricción de que nuestra solución óptima \mathbf{x}^* sea entera. Es decir, lo único que cambia es la región de factibilidad. De hecho, si definimos el poliedro

$$P := \{\mathbf{x} \in \mathbb{R}^n : A\mathbf{x} = \mathbf{b}, \mathbf{x} \geq \mathbf{0}\},$$

entonces tenemos que $P \cap \mathbb{Z}^n$ corresponde a la región factible de (1.6),

²Sin embargo, la literatura para explicar estos métodos es abundante. Véase, por ejemplo, [NW06].

mientras que P corresponde a la región factible de su problema relajado.

A partir de lo anterior, deducimos inmediatamente que el valor óptimo z_{PE}^* de un programa entero es una cota inferior del valor óptimo z^* de su problema relajado, pues ambos son problemas de maximización y es cierto que $P \cap \mathbb{Z}^n \subseteq P$. De aquí se sigue que si $z_{\text{PE}}^* = z^*$, entonces la solución óptima \mathbf{x}^* del problema relajado también es la solución óptima del programa lineal entero si $\mathbf{x}^* \in \mathbb{Z}^n$.

El algoritmo estándar para encontrar soluciones de programas lineales enteros es Ramificación y Acotamiento. Este método consiste en generar un árbol binario donde cada nodo representa un subproblema lineal a resolver. En la raíz del árbol resolvemos el problema relajado (1.5) y, si la solución óptima $\mathbf{x}^* \in \mathbb{R}^n$ no es entera, entonces para alguna entrada x_i^* no entera agregamos la restricción $x_i \leq \lfloor x_i^* \rfloor$ para crear un subproblema, y también añadimos la restricción $x_i \geq \lceil x_i^* \rceil$ para crear otro subproblema. Este procedimiento se realiza de manera recursiva.

Observemos que, si decidimos recorrer todos los nodos del árbol binario, entonces tendremos que resolver al menos 2^n subproblemas, donde n es la dimensión del problema lineal. Por esta razón, el algoritmo cuenta con políticas para deshacerse de subárboles que nunca proveerán la solución óptima. El autor considera que es mejor ilustrar estas políticas a partir de un ejemplo. El Algoritmo 6 en el Apéndice A presenta una versión rudimentaria del método de Ramificación y Acotamiento.

Ejemplo 1.1.18 ([Oli17]). Consideremos el programa lineal entero

$$\begin{aligned} \max_{\mathbf{x} \in \mathbb{Z}^2} \quad & 4x_1 - x_2, \\ \text{s. a.} \quad & 7x_1 - 2x_2 \leq 14, \\ & 2x_1 - 2x_2 \leq 3, \\ & x_2 \leq 3, \end{aligned}$$

$$x_1, x_2 \geq 0.$$

La región factible de este problema se muestra en la Figura 1.2. La solución al problema relajado, cuya región factible denotamos por S_0 , está dada por $\mathbf{x}^0 := (20/7, 3)^T$. Como $x_1^0 = 20/7$ no es entero, generamos dos nuevos subproblemas con regiones factibles

$$S_{00} := S_0 \cup \{x_1 \leq \lfloor 20/7 \rfloor = 2\},$$

$$S_{01} := S_0 \cup \{x_1 \geq \lceil 20/7 \rceil = 3\}.$$

De la Figura 1.2, observamos que S_{01} es vacío y por lo tanto de este problema no podemos generar otros subproblemas. En este caso, decimos que **podamos S_{01} por infactibilidad**.

Ahora bien, la solución al problema S_{00} está dada por $\mathbf{x}^1 := (2, 1/2)^T$. Encontramos que $x_2^1 = 1/2$ no es entero y por lo tanto generamos dos nuevos subproblemas:

$$S_{000} := S_{00} \cup \{x_2 \leq \lfloor 1/2 \rfloor = 0\},$$

$$S_{001} := S_{00} \cup \{x_2 \geq \lceil 1/2 \rceil = 1\}.$$

Observemos que la solución \mathbf{x}^2 de S_{001} es $(2, 1)^T$, la cual es entera y tiene valor objetivo $z_2^* := 7$. No generamos otros subproblemas a partir de este problema porque sus regiones factibles estarán contenidas en S_{001} y por lo tanto sus valores objetivos serán menores o iguales al de S_{001} . Así pues, decimos que **podamos S_{001} por integralidad**.

La solución de S_{000} , en cambio, es $\mathbf{x}^3 := (3/2, 0)^T$ y tendríamos que ramificar de nuevo en otros dos subproblemas. No obstante, observemos que el valor objetivo de este subproblema es $z_3^* := 6$, el cual es menor que $z_2^* = 7$. Como la región factible de cualquier subproblema generado a partir de este último problema está contenido en S_{000} , se sigue que su valor objetivo

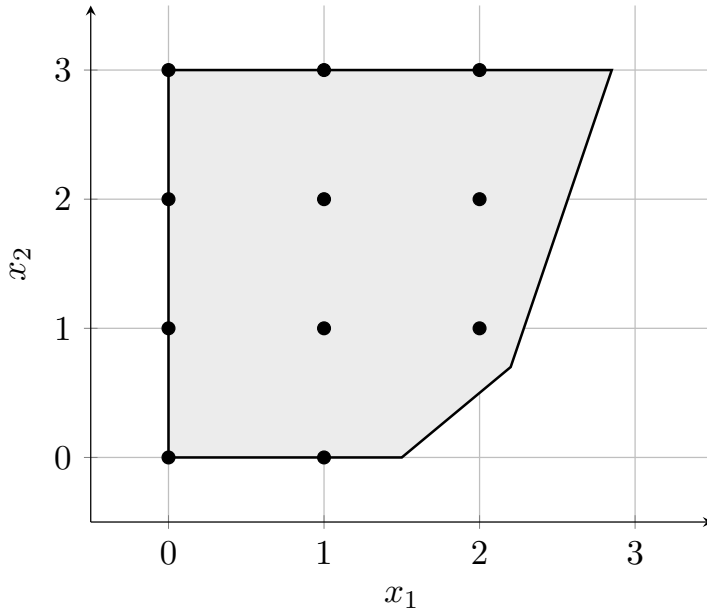


Figura 1.2.: Los puntos negros forman la región factible del programa lineal entero del Ejemplo 1.1.18, mientras que la región sombreada es la región factible de su problema relajado.

será menor o igual al de S_{000} . Es decir, no encontraremos otro vector cuyo valor objetivo sea mayor que z_2^* . Decimos entonces que **podamos** S_{000} **por cota**.

Como hemos agotado todos los subproblemas que podríamos generar, entonces concluimos que la solución óptima de este programa lineal entero es $\mathbf{x}^2 = (2, 1)^T$ y tiene valor objetivo $z_2^* = 7$.

1.2. Fundamentos

Esta sección se divide en cuatro partes. En primer lugar, damos a conocer las definiciones y enunciados provistos por [BH09]. Es importante aclarar que el autor tradujo libremente algunos términos a falta de encontrar

fuentes en español que hicieran uso de ellos. En particular, el autor decidió nombrar “vectores esencialmente enteros” a los *projectively rational vectors* (ver definición 1.2.1) y “capas enteras” a los *c-layers* (ver definición 1.2.3).

En segundo lugar, mostramos la equivalencia entre resolver problemas del tipo (1.1) y resolver ecuaciones lineales diofantinas. Nos basamos en los teoremas 1.1.8 y 1.1.10 para construir inductivamente el conjunto de soluciones de una ecuación lineal diofantina en n incógnitas. Las soluciones de estas ecuaciones serán definidas a partir de una relación de recurrencia y también dependerán de una colección de $n - 1$ parámetros libres.

En tercer lugar, establecemos una transformación lineal entre el conjunto de soluciones de una ecuación lineal diofantina en n incógnitas y el conjunto de $n - 1$ parámetros libres que determinan estas soluciones. Luego, investigamos las propiedades de esta transformación lineal que serán de gran utilidad teórica para los siguientes capítulos, en especial para el capítulo 4.

Finalmente, mostramos que el vector objetivo \mathbf{p} del problema original (1.1) induce una descomposición interesante de \mathbb{Z}^n y analizamos cómo se relacionan estas descomposiciones al considerar distintos vectores \mathbf{p} . Esto nos permitirá mostrar equivalencias entre distintas instancias del problema (1.1).

1.2.1. Capas enteras

Definición 1.2.1. Decimos que un vector $\mathbf{v} \in \mathbb{R}^n \setminus \{\mathbf{0}\}$ es **esencialmente entero** si existen un vector $\mathbf{w} \in \mathbb{Z}^n$ y un escalar $m \neq 0$ tales que $\mathbf{v} = m\mathbf{w}$. Además, decimos que \mathbf{w} es el **múltiplo coprimo** de \mathbf{v} si sus entradas son coprimas y si su primera entrada no nula es positiva.

Ejemplo 1.2.2. El vector $(-\sqrt{2}, 1/\sqrt{2})^T = 2\sqrt{2}(-2, 1)^T$ es esencialmente entero y $(2, -1)^T$ es su múltiplo coprimo. En contraste, el vector $(\sqrt{2}, \sqrt{3})^T$ no es esencialmente entero.

Observación. Todo vector racional $\mathbf{v} \in \mathbb{Q}^n \setminus \{\mathbf{0}\}$ es esencialmente entero. En efecto, para cada $i \in \{1, \dots, n\}$ existen $p_i, q_i \in \mathbb{Z}$ con $q_i \neq 0$ tales que $v_i = p_i/q_i$. Luego, si definimos $m := \text{mcm}\{q_1, \dots, q_n\} \neq 0$ y también $\mathbf{w} := m\mathbf{v} \in \mathbb{Z}^n$, encontramos que $\mathbf{v} = \frac{1}{m}\mathbf{w}$.

Observación. Todo vector \mathbf{v} esencialmente entero tiene a lo más dos vectores coprimos asociados. Sean $m \in \mathbb{R}$ y $\mathbf{w} \in \mathbb{Z}^n$ tales que $\mathbf{v} = m\mathbf{w}$. Entonces

$$\pm \frac{1}{\text{mcd}\{w_1, \dots, w_n\}} \mathbf{w}$$

son dos vectores cuyas entradas son coprimas, de acuerdo al lema 1.1.6. Como la primera entrada no nula w_i también debe ser positiva, se sigue que solo uno de estos dos vectores es el múltiplo coprimo de \mathbf{v} . Así, el múltiplo coprimo de un vector esencialmente entero es único.

Puesto que todo número representable en cualquier sistema de aritmética finita es necesariamente racional, decidimos enfocar nuestro análisis en vectores esencialmente enteros. Desde el punto de vista puramente teórico, esta condición reduce de manera drástica el tipo de programas lineales que podemos resolver. No obstante, esta clase de vectores es un poco más general que los considerados en otros textos de programación lineal. Por ejemplo, [MT90] y [Sch98] toman en cuenta vectores puramente racionales.

Definición 1.2.3. Sea $\mathbf{v} \in \mathbb{R}^n$ un vector esencialmente entero y sea $t \in \mathbb{R}$ un escalar. Decimos que su hiperplano afino asociado

$$H_{\mathbf{v},t} := \ker\{\mathbf{x} \mapsto \mathbf{v}^T \mathbf{x}\} + t\mathbf{v} = \{\mathbf{v}^\perp + t\mathbf{v} : \mathbf{v}^T \mathbf{v}^\perp = 0\} \quad (1.7)$$

es una **capa entera** si contiene al menos un punto entero.

Lema 1.2.4. Sean $\mathbf{v}, \mathbf{x} \in \mathbb{R}^n$ con \mathbf{v} distinto de cero. Entonces $\mathbf{x} \in H_{\mathbf{v},t_{\mathbf{x}}}$, donde $t_{\mathbf{x}} := \frac{\mathbf{v}^T \mathbf{x}}{\|\mathbf{v}\|^2}$.

Las capas enteras son invariantes ante reescalamientos en el vector \mathbf{v} : si

$r \neq 0$, entonces $H_{\mathbf{v},t} = H_{r\mathbf{v},t/r}$. En efecto, sea $\mathbf{x} \in H_{\mathbf{v},t}$. Luego, existe \mathbf{v}^\perp ortogonal a \mathbf{v} tal que

$$\mathbf{x} = \mathbf{v}^\perp + t\mathbf{v} = \mathbf{v}^\perp + \frac{t}{r}(r\mathbf{v}).$$

Pero si \mathbf{v}^\perp es ortogonal a \mathbf{v} , entonces también es ortogonal a $r\mathbf{v}$. Así, encontramos que $\mathbf{x} \in H_{r\mathbf{v},t/r}$. La otra contención se muestra de manera similar.

En particular, si \mathbf{w} es el múltiplo coprimo de \mathbf{v} , se cumple que

$$\{H_{\mathbf{v},t} : t \in \mathbb{R}\} = \{H_{\mathbf{v}/m,tm} : t \in \mathbb{R}\} = \{H_{\mathbf{w},t} : t \in \mathbb{R}\}, \quad (1.8)$$

donde $m \neq 0$ es el escalar que satisface $\mathbf{v} = m\mathbf{w}$. Así pues, para analizar las capas enteras, basta con fijarnos en los múltiplos coprimos que las definen en vez de sus vectores esencialmente enteros asociados.

Teorema 1.2.5. *Sea $\mathbf{v} \in \mathbb{R}^n$ un vector esencialmente entero y sea \mathbf{w} su múltiplo coprimo. Entonces la familia de capas enteras $\{H_{\mathbf{w},k\|\mathbf{w}\|^{-2}} : k \in \mathbb{Z}\}$ cubre a \mathbb{Z}^n .*

Lema 1.2.6. *Sea $\mathbf{v} \in \mathbb{R}^n$ un vector esencialmente entero y sea \mathbf{w} su múltiplo coprimo. Entonces $\mathbf{q}^T \mathbf{x} = k$ para todo $\mathbf{x} \in H_{\mathbf{w},k\|\mathbf{w}\|^{-2}}$.*

Demostración. Sea $\mathbf{x} \in H_{\mathbf{w},k\|\mathbf{w}\|^{-2}}$, por lo que existe un vector \mathbf{w}^\perp ortogonal a \mathbf{w} tal que

$$\mathbf{x} = \mathbf{w}^\perp + \frac{k}{\|\mathbf{w}\|^2} \mathbf{w}.$$

Luego,

$$\mathbf{w}^T \mathbf{x} = \mathbf{w}^T \mathbf{w}^\perp + \frac{k}{\|\mathbf{w}\|^2} \mathbf{w}^T \mathbf{w} = 0 + \frac{k}{\|\mathbf{w}\|^2} \|\mathbf{w}\|^2 = k.$$

que es lo que deseábamos obtener. \square

Consideremos el problema (1.1) y supongamos que el vector objetivo \mathbf{p} es esencialmente entero. Sea \mathbf{q} su múltiplo coprimo. Sabemos de (1.8) que $\{H_{\mathbf{p},t} : t \in \mathbb{R}\} = \{H_{\mathbf{q},t} : t \in \mathbb{R}\}$. Puesto que nos concentramos en puntos

enteros, por el teorema 1.2.5 podemos considerar exclusivamente el subconjunto de capas enteras $\{H_{\mathbf{q},k\|\mathbf{q}\|^{-2}} : k \in \mathbb{Z}\}$.

Observemos del lema 1.2.6 junto con la restricción presupuestaria (1.1b) que los puntos enteros de la k -ésima capa entera $H_{\mathbf{q},k\|\mathbf{q}\|^{-2}}$ o bien respetan todos esta restricción o bien ninguno la respeta. Nos gustaría entonces determinar el primer parámetro $\eta \in \mathbb{Z}$ que induce a que todos los puntos enteros de $H_{\mathbf{q},\eta\|\mathbf{q}\|^{-2}}$ respeten la restricción presupuestaria.

Para respetar la restricción (1.1b), debe ser el caso que todo $\mathbf{x} \in H_{\mathbf{q},k\|\mathbf{q}\|^{-2}}$ satisfaga

$$\mathbf{p}^T \mathbf{x} = m\mathbf{q}^T \mathbf{x} = mk \leq u \iff \begin{cases} k \geq u/m, & m < 0, \\ k \leq u/m, & m > 0, \end{cases} \quad (1.9)$$

donde $m \neq 0$ es el escalar que satisface $\mathbf{p} = m\mathbf{q}$. Así pues, dependiendo del signo de m , tenemos que el primer parámetro η en satisfacer la restricción presupuestaria puede ser interpretado como el entero más pequeño o el entero más grande que satisface su respectiva desigualdad.

Lema 1.2.7. *Sea $\mathbf{p} \in \mathbb{R}^n$ un vector esencialmente entero y sea \mathbf{q} su múltiplo coprimo, de manera que $\mathbf{p} = m\mathbf{q}$ para algún escalar $m \neq 0$. Entonces la primera capa entera $H_{\mathbf{q},\eta\|\mathbf{q}\|^{-2}}$ que satisfacen la restricción (1.1b) está parametrizada por*

$$\eta := \begin{cases} \lceil u/m \rceil, & m < 0, \\ \lfloor u/m \rfloor, & m > 0. \end{cases} \quad (1.10)$$

Demostración. Se sigue inmediatamente de (1.9). □

Puesto que la gran mayoría de nuestros enunciados y algoritmos dependen de este primer parámetro η , tendremos que separarlos al menos en dos casos. Por ello, el autor creyó prudente considerar solamente el caso $m > 0$, aunque cabe mencionar que los enunciados y demostraciones para el caso $m < 0$

son completamente análogos, donde las diferencias recaen en que el orden de las desigualdades cambian o las funciones piso se reemplazan por funciones techo, por ejemplo.

En resumen, si $m > 0$, encontramos que las capas enteras que satisfacen la restricción (1.1b) están parametrizadas por $k \in \{\eta, \eta - 1, \dots\}$, donde η está definida en 1.2.7. Además, por el lema 1.2.6, todo $\mathbf{x} \in H_{\mathbf{q}, k\|\mathbf{q}\|}^{-2}$ satisface la ecuación $\mathbf{q}^T \mathbf{x} = k$.

Teorema 1.2.8. *Sea $\mathbf{p} \in \mathbb{R}^n$ un vector esencialmente entero y sea \mathbf{q} su múltiplo coprimo. Entonces el problema (1.1) es infactible si y solo si $\mathbf{q} \geq \mathbf{0}$ y el lado derecho u de (1.1b) es negativo.*

Demostración. Supongamos que $\mathbf{q} \geq \mathbf{0}$ y $u < 0$. Si $\mathbf{x} \in \mathbb{Z}_{\geq \mathbf{0}}^n$ entonces $\mathbf{q}^T \mathbf{x} \geq 0 > u$ y por lo tanto \mathbf{x} no es factible. Luego,

$$\mathbb{Z}_{\geq \mathbf{0}}^n \cap \{\mathbf{x} : \mathbf{q}^T \mathbf{x} \leq u\} = \emptyset,$$

y el problema no es factible.

Mostramos la otra implicación por contraposición. Si $u \geq 0$ observamos que $\mathbf{0} \in \mathbb{Z}^n$ es factible. Se debe cumplir $u < 0$. Similarmente, si $q_i < 0$ para algún $i \in \{1, \dots, n\}$, encontramos que $\lceil u/q_i \rceil \mathbf{e}_i \in \mathbb{Z}^n$ es factible:

$$\mathbf{q}^T \left\lceil \frac{u}{q_i} \right\rceil \mathbf{e}_i = q_i \left\lceil \frac{u}{q_i} \right\rceil \leq q_i \frac{u}{q_i} = u,$$

además, como $u < 0$, concluimos que $\lceil u/q_i \rceil \mathbf{e}_i$ es no negativo. \square

Debido al teorema anterior, somos capaces de determinar automáticamente si el problema (1.1) es infactible, por lo que supondremos de ahora en adelante que es factible. El siguiente teorema muestra que nuestro análisis para resolver este problema debe dividirse en dos casos.

Teorema 1.2.9. *Sea $\mathbf{p} \in \mathbb{R}^n$ un vector esencialmente entero y sea \mathbf{q} su múltiplo coprimo, de manera que $\mathbf{p} = m\mathbf{q}$ para alguna $m > 0$. Supongamos*

que el problema (1.1) es factible y tomemos η del lema 1.2.7. Entonces se satisface lo siguiente:

1. Si $q_i < 0$ para algún $i \in \{1, \dots, n\}$, entonces la η -ésima capa entera $H_{\mathbf{q}, \eta \|\mathbf{q}\|^{-2}}$ contiene un número infinito de puntos factibles.
2. Si $\mathbf{q} > \mathbf{0}$ entonces, para todo $k \in \{\eta, \eta - 1, \dots, 0\}$, la k -ésima capa entera $H_{\mathbf{q}, k \|\mathbf{q}\|^{-2}}$ contiene un número finito de puntos factibles.

Demostración.

1. En la subsección 1.2.2 mostraremos que, como \mathbf{q} es un vector cuyas entradas son coprimas, entonces existe un punto entero \mathbf{x} que satisface la ecuación lineal diofantina $\mathbf{q}^T \mathbf{x} = \eta$. Por el momento confiemos que esto es verdadero. Luego,

$$\mathbf{p}^T \mathbf{x} = m \mathbf{q}^T \mathbf{x} = m\eta = m \left\lfloor \frac{u}{m} \right\rfloor \leq m \frac{u}{m} = u,$$

y se satisface la restricción (1.1b).

Como no tenemos asegurada la no negatividad de \mathbf{x} , construiremos un vector entero \mathbf{x}^+ que sí satisface la restricción de no negatividad y también la restricción presupuestaria $\mathbf{q}^T \mathbf{x}^+ = \eta$, de manera que \mathbf{x}^+ sí será factible.

Definamos los siguientes conjuntos de índices:

$$I^+ := \{i : q_i > 0\}, \quad I^\circ := \{\ell : q_\ell = 0\}, \quad I^- := \{j : q_j < 0\}.$$

Podemos suponer sin pérdida de generalidad que I° es vacío. En efecto, si $x_k < 0$ para algún $k \in I^\circ$, esa entrada no sería factible, pero fácilmente podríamos definir $x_k^+ = 0$ para hacerla factible.

Por hipótesis, sabemos que \mathbf{q} tiene una entrada negativa y por lo tanto $I^- \neq \emptyset$. Además, por la definición 1.2.1, \mathbf{q} tiene una entrada positiva

y por lo tanto $I^+ \neq \emptyset$. Luego, ambos conjuntos I^+ e I^- forman una partición del conjunto $\{1, \dots, n\}$. Podemos escoger enteros positivos c_1, \dots, c_n que satisfagan simultáneamente

$$x_k + \sum_{i \in I^+} q_i c_i \geq 0, \quad \forall k \in I^-, \quad (1.11)$$

$$x_k - \sum_{j \in I^-} q_j c_k \geq 0, \quad \forall k \in I^+. \quad (1.12)$$

Definamos el vector $\mathbf{x}^+ \in \mathbb{Z}^n$ de manera que

$$x_k^+ := \begin{cases} x_k + \sum_{i \in I^+} q_i c_i, & k \in I^-, \\ x_k - \sum_{j \in I^-} q_j c_k, & k \in I^+. \end{cases}$$

Se verifica que \mathbf{x}^+ es no negativo y, además,

$$\begin{aligned} \mathbf{q}^T \mathbf{x}^+ &= \mathbf{q}^T \mathbf{x} + \sum_{k \in I^-} \sum_{i \in I^+} q_k q_i c_i - \sum_{k \in I^+} \sum_{j \in I^-} q_k q_j c_k \\ &= \eta + \sum_{j \in I^-} \sum_{i \in I^+} q_j q_i c_i - \sum_{i \in I^+} \sum_{j \in I^-} q_i q_j c_i \\ &= \eta. \end{aligned}$$

Así pues, tenemos existencia de un punto factible. Para concluir que hay un número infinito de puntos factibles, basta observar que si la elección de coeficientes c_1, \dots, c_n satisface ambas desigualdades (1.11) y (1.12), entonces cualquier múltiplo entero positivo de estos coeficientes también las satisface.

2. Se sigue que $u \geq 0$. Definamos

$$P_k := H_{\mathbf{q}, k \|\mathbf{q}\|^{-2}} \cap \mathbb{Z}_{\geq \mathbf{0}}^n = \{\mathbf{x} \in \mathbb{Z}^n : \mathbf{q}^T \mathbf{x} = k, \mathbf{x} \geq \mathbf{0}\}, \quad (1.13)$$

y observemos que $P_k = \emptyset$ para todo k negativo, pues $\mathbf{q} > \mathbf{0}$ y por

lo tanto $\mathbf{q}^T \mathbf{x} \geq 0$ para cualquier $\mathbf{x} \in \mathbb{Z}_{\geq 0}^n$. Esto implica que ningún punto sobre capas enteras con parámetros negativos es factible.

Sea $k \in \{\eta, \eta - 1, \dots, 0\}$. La capa entera $H_{\mathbf{q}, k\|\mathbf{q}\|^{-2}}$ interseca los ejes positivos en $\frac{k}{q_i} \mathbf{e}_i$. Definamos $\ell_i := \lceil k/q_i \rceil$. No es difícil ver que P_k está contenido en el prisma cuyas aristas son $[0, \ell_i]$ y, por lo tanto,

$$P_k \subseteq \prod_{i=1}^n [0, \ell_i] \cap \mathbb{Z}^n = \prod_{i=1}^n ([0, \ell_i] \cap \mathbb{Z}).$$

Pero $|[0, \ell_i] \cap \mathbb{Z}| = \ell_i + 1$. Así,

$$|P_k| \leq \prod_{i=1}^n (\ell_i + 1) < \infty.$$

Entonces la k -ésima capa entera contiene un número finito de puntos factibles.

□

Ciertamente el primer caso del teorema 1.2.9 es el menos interesante, pues conocemos inmediatamente el valor óptimo de estas instancias. No obstante, existen muchos elementos en común que comparten ambos casos. También es cierto que esta división dejará de existir una vez que introduzcamos múltiples restricciones en el capítulo 4.

Además, antes de analizar los dos casos que el teorema anterior impone, primero debemos mostrar que la ecuación lineal diofantina $\mathbf{q}^T \mathbf{x} = k$ admite soluciones enteras para toda $k \in \mathbb{Z}$ siempre que las entradas de \mathbf{q} sean coprimas. Habíamos supuesto esto en la demostración anterior.

Así también, la construcción de soluciones enteras de ecuaciones lineales diofantinas proveerá herramientas teóricas útiles para demostrar la gran mayoría de resultados que presentaremos. Cabe mencionar que la siguiente subsección se encarga de construir solamente soluciones enteras de estas

ecuaciones. Será cuestión de los capítulos 2 y 3 obtener soluciones que además sean no negativas.

1.2.2. Construcción de soluciones enteras

Debido al teorema 1.2.5, las soluciones del problema (1.1) se encuentran en una capa entera $H_{\mathbf{q}, k \|\mathbf{q}\|^{-2}}$. Luego, por el lema 1.2.6, los puntos $\mathbf{x} \in \mathbb{Z}^n$ que se encuentran sobre esa capa satisfacen la ecuación lineal diofantina

$$\mathbf{q}^T \mathbf{x} = q_1 x_1 + q_2 x_2 + \cdots + q_n x_n = k. \quad (1.14)$$

Como hemos mencionado previamente, podemos suponer sin pérdida de generalidad que ninguna entrada de \mathbf{q} es nula.

En la sección 1.1.1 mostramos condiciones de existencia de este tipo de ecuaciones, así como su construcción, cuando $n = 2$. Partimos de la observación que podemos resolver recursivamente esta ecuación. Definamos, por conveniencia, $g_1 := \text{mcd}\{q_1, \dots, q_n\}$ y también $\omega_1 := k$. Puesto que \mathbf{q} es un vector con entradas coprimas, sabemos que $g_1 = 1$. También definamos

$$\omega_2 := \frac{q_2}{g_2 \cdot g_1} x_2 + \cdots + \frac{q_n}{g_n \cdot g_1} x_n, \quad (1.15)$$

donde $g_2 := \text{mcd}\{q_2/g_1, \dots, q_n/g_1\}$. Como $q_n \neq 0$, tenemos que g_2 está bien definido y además es positivo. Así, la ecuación (1.14) es equivalente a

$$\frac{q_1}{g_1} x_1 + g_2 \omega_2 = \omega_1. \quad (1.16)$$

Observemos que

$$\begin{aligned} \text{mcd}\left\{\frac{q_1}{g_1}, g_2\right\} &= \text{mcd}\left\{\frac{q_1}{g_1}, \text{mcd}\left\{\frac{q_2}{g_1}, \dots, \frac{q_n}{g_1}\right\}\right\} \\ &= \text{mcd}\left\{\frac{q_1}{g_1}, \frac{q_2}{g_1}, \dots, \frac{q_n}{g_1}\right\} = 1. \end{aligned}$$

Por el teorema 1.1.8, existen soluciones enteras para todo $\omega_1 \in \mathbb{Z}$. Como

q_1/g_1 y g_2 son coprimos, encontramos que sus coeficientes de Bézout (ver definición 1.1.9) asociados x'_1, ω'_2 son soluciones particulares de la ecuación

$$\frac{q_1}{g_1}x_1 + g_2\omega_2 = 1.$$

Deducimos del teorema 1.1.10 que las soluciones de la ecuación (1.16) están dadas por

$$\begin{cases} x_1 = \omega_1 x'_1 + g_2 t_1, \\ \omega_2 = \omega_1 \omega'_2 - \frac{q_1}{g_1} t_1, \end{cases} \quad (1.17)$$

donde $t_1 \in \mathbb{Z}$ es una variable libre.

Observación. Los coeficientes de Bézout x'_1 y ω'_2 dependen exclusivamente de \mathbf{q} y no del punto \mathbf{x} . En efecto, x'_1 está asociado a q_1/g_1 y ω'_2 está asociado a g_2 . Pero ambos g_1 y g_2 son el máximo común divisor de q_1, \dots, q_n y de $q_1/g_1, \dots, q_n/g_1$, respectivamente.

Para el siguiente paso de la recursión, escogemos cualquier $t_1 \in \mathbb{Z}$ para fijar ω_2 . Tenemos de (1.15) que debemos resolver la ecuación

$$\frac{q_2}{g_2 \cdot g_1}x_2 + \frac{q_3}{g_2 \cdot g_1}x_3 + \dots + \frac{q_n}{g_2 \cdot g_1}x_n = \omega_2. \quad (1.18)$$

Como $g_2 = \text{mcd}\{q_2/g_1, \dots, q_n/g_1\}$, sabemos del lema 1.1.6 que

$$\text{mcd}\left\{\frac{q_2}{g_2 \cdot g_1}, \dots, \frac{q_n}{g_2 \cdot g_1}\right\} = 1.$$

En el mismo espíritu del primer paso de la recursión, definimos

$$\omega_3 := \frac{q_3}{g_3 \cdot g_2 \cdot g_1}x_3 + \dots + \frac{q_n}{g_3 \cdot g_2 \cdot g_1}x_n,$$

donde

$$g_3 := \text{mcd}\left\{\frac{q_3}{g_2 \cdot g_1}, \dots, \frac{q_n}{g_2 \cdot g_1}\right\}.$$

Nuevamente, como q_n es distinto de cero, g_3 está bien definido y además es

positivo. Así pues, la ecuación (1.18) es equivalente a

$$\frac{q_2}{g_2 \cdot g_1} x_2 + g_3 \omega_3 = \omega_2. \quad (1.19)$$

También se cumple que

$$\text{mcd} \left\{ \frac{q_2}{g_2 \cdot g_1}, g_3 \right\} = 1,$$

y entonces (1.19) tiene una infinidad de soluciones para todo $\omega_2 \in \mathbb{Z}$, las cuales están dadas por

$$\begin{cases} x_2 = \omega_2 x'_2 + g_3 t_2, \\ \omega_3 = \omega_2 \omega'_3 - \frac{q_2}{g_2 \cdot g_1} t_2, \end{cases}$$

donde $t_2 \in \mathbb{Z}$ es una variable libre, y x'_2, ω'_3 son los coeficientes de Bézout asociados a $\frac{q_2}{g_2 \cdot g_1}$ y g_3 , respectivamente.

De manera general, para $i \in \{1, \dots, n-2\}$, el i -ésimo paso de la recursión provee la ecuación

$$\frac{q_i}{\prod_{j=1}^i g_j} x_i + \frac{q_{i+1}}{\prod_{j=1}^i g_j} x_{i+1} + \dots + \frac{q_n}{\prod_{j=1}^i g_j} x_n = \omega_i, \quad (1.20)$$

donde

$$g_i := \text{mcd} \left\{ \frac{q_i}{\prod_{j=1}^{i-1} g_j}, \dots, \frac{q_n}{\prod_{j=1}^{i-1} g_j} \right\}, \quad (1.21)$$

por el lema 1.1.6 se sigue que

$$\text{mcd} \left\{ \frac{q_i}{\prod_{j=1}^i g_j}, \dots, \frac{q_n}{\prod_{j=1}^i g_j} \right\} = 1. \quad (1.22)$$

Ahora bien, definamos

$$g_{i+1} := \text{mcd} \left\{ \frac{q_{i+1}}{\prod_{j=1}^i g_j}, \dots, \frac{q_n}{\prod_{j=1}^i g_j} \right\}. \quad (1.23)$$

Como q_n es distinto de cero, se sigue que g_{i+1} está bien definido y es positivo. Definamos

$$\omega_{i+1} := \frac{q_{i+1}}{\prod_{j=1}^{i+1} g_j} x_{i+1} + \cdots + \frac{q_n}{\prod_{j=1}^{i+1} g_j} x_n,$$

de manera que la ecuación (1.20) es equivalente a

$$\frac{q_i}{\prod_{j=1}^i g_j} x_i + g_{i+1} \omega_{i+1} = \omega_i. \quad (1.24)$$

A partir de (1.22) y de (1.23), encontramos que

$$\text{mcd} \left\{ \frac{q_i}{\prod_{j=1}^i g_j}, g_{i+1} \right\} = \text{mcd} \left\{ \frac{q_i}{\prod_{j=1}^i g_j}, \frac{q_{i+1}}{\prod_{j=1}^i g_j}, \dots, \frac{q_n}{\prod_{j=1}^i g_j} \right\} = 1,$$

y del teorema 1.1.8 se sigue que la ecuación (1.24) tiene soluciones enteras para todo $\omega_i \in \mathbb{Z}$. Por el teorema 1.1.10, las soluciones enteras de (1.24) están dadas por

$$\begin{cases} x_i = \omega_i x'_i + g_{i+1} t_i, \\ \omega_{i+1} = \omega_i \omega'_{i+1} - \frac{q_i}{\prod_{j=1}^i g_j} t_i, \end{cases} \quad (1.25)$$

donde $t_i \in \mathbb{Z}$ es la i -ésima variable libre. Es valioso mencionar, otra vez, que los coeficientes de Bézout x'_i, ω'_{i+1} dependen exclusivamente de \mathbf{q} a través de sus entradas q_i y de los máximos común divisores entre ellas. En efecto, por el teorema 1.1.4, estos coeficientes son soluciones particulares de la ecuación

$$\frac{q_i}{\prod_{j=1}^i g_j} x'_i + g_{i+1} \omega'_{i+1} = 1. \quad (1.26)$$

Finalmente, en el último paso de la recursión obtenemos la ecuación lineal diofantina

$$\frac{q_{n-1}}{\prod_{j=1}^{n-1} g_j} x_{n-1} + \frac{q_n}{\prod_{j=1}^{n-1} g_j} x_n = \omega_{n-1}. \quad (1.27)$$

Por construcción, los coeficientes de x_{n-1} y x_n son coprimos. A causa del

teorema 1.1.10 las soluciones enteras están dadas por

$$\begin{cases} x_{n-1} = \omega_{n-1}x'_{n-1} + \frac{q_n}{\prod_{j=1}^{n-1} g_j} t_{n-1}, \\ x_n = \omega_{n-1}x'_n - \frac{q_{n-1}}{\prod_{j=1}^{n-1} g_j} t_{n-1}, \end{cases} \quad (1.28)$$

donde x'_{n-1}, x'_n son los coeficientes de Bézout asociados a $\frac{q_n}{\prod_{j=1}^{n-1} g_j}$ y $\frac{q_{n-1}}{\prod_{j=1}^{n-1} g_j}$, respectivamente, por lo que satisfacen

$$\frac{q_{n-1}}{\prod_{j=1}^{n-1} g_j} x'_{n-1} + \frac{q_n}{\prod_{j=1}^{n-1} g_j} x'_n = 1. \quad (1.29)$$

Hemos demostrado, que la ecuación lineal diofantina (1.14) tiene al menos una solución, siempre que \mathbf{q} sea un vector coprimo. Así pues, saldamos nuestra cuenta pendiente con respecto a una parte de la demostración 1.2.9. Además, construimos una infinidad de soluciones enteras, pues la elección de cada variable libre $t_i \in \mathbb{Z}$ provee una solución distinta. Aún más, por el teorema 1.1.10, sabemos que el conjunto de estas soluciones es exhaustiva.

Con respecto a la no negatividad de las soluciones mencionamos brevemente lo siguiente. Observamos de (1.25) que t_i debe satisfacer

$$t_i \geq \left\lceil -\frac{\omega_i x'_i}{g_{i+1}} \right\rceil, \quad (1.30)$$

para todo $i \in \{1, \dots, n-2\}$. Ahora bien, para asegurar la no negatividad de x_{n-1} y x_n , observamos de (1.27) que dependemos de los signos de q_{n-1} y de q_n . Debido al teorema 1.2.9, relegamos esta discusión para los siguientes dos capítulos.

1.2.3. Soluciones y variables libres

Hemos encontrado una relación entre un vector de variables libres $\mathbf{t} \in \mathbb{Z}^{n-1}$ y un vector solución $\mathbf{x} \in \mathbb{Z}^n$ de la ecuación (1.14). Sabemos de (1.25) que la relación está dada de manera recursiva. Puesto que deseamos establecer una

transformación lineal entre \mathbf{t} y \mathbf{x} , resulta sumamente conveniente determinar una forma cerrada de esta relación. Recordemos que habíamos definido, por construcción, $\omega_1 := k$. Combinando esto con la segunda igualdad de (1.25), obtenemos la relación de recurrencia

$$\begin{cases} \omega_1 = k, \\ \omega_{i+1} = \omega_i \cdot \omega'_{i+1} - \frac{q_i}{\prod_{\ell=1}^i g_\ell} \cdot t_i. \end{cases} \quad (1.31)$$

Lema 1.2.10. *La forma cerrada de la relación de recurrencia (1.31) está dada por*

$$\omega_i = k \cdot \prod_{j=2}^i \omega'_j - \sum_{j=1}^{i-1} t_j \cdot \frac{q_j}{\prod_{\ell=1}^j g_\ell} \cdot \prod_{\ell=j+2}^i \omega'_\ell. \quad (1.32)$$

donde, por conveniencia, le asignamos el valor de 0 a la suma vacía y el valor de 1 al producto vacío.

Demostración. Lo demostramos inductivamente. Observemos que

$$\omega_1 = k \cdot \prod_{j=2}^1 \omega'_j - \sum_{j=1}^0 t_j \cdot \frac{q_j}{\prod_{\ell=1}^j g_\ell} \cdot \prod_{\ell=j+2}^1 \omega'_\ell = k,$$

debido a que definimos el producto vacío como 1 y la suma vacía como 0. Supongamos inductivamente que (1.32) se satisface para alguna $i \in \mathbb{N}$. Entonces, tenemos

$$\begin{aligned} & k \cdot \prod_{j=2}^{i+1} \omega'_j - \sum_{j=1}^i t_j \cdot \frac{q_j}{\prod_{\ell=1}^j g_\ell} \cdot \prod_{\ell=j+2}^{i+1} \omega'_\ell \\ &= k \cdot \prod_{j=2}^i \omega'_j \cdot \omega'_{i+1} - \sum_{j=1}^{i-1} t_j \cdot \frac{q_j}{\prod_{\ell=1}^j g_\ell} \cdot \prod_{\ell=j+2}^i \omega'_\ell \cdot \omega'_{i+1} - \frac{q_i}{\prod_{\ell=1}^i g_\ell} \cdot \prod_{\ell=i+2}^{i+1} \omega'_\ell \cdot t_i \\ &= \left(k \cdot \prod_{j=2}^i \omega'_j - \sum_{j=1}^{i-1} t_j \cdot \frac{q_j}{\prod_{\ell=1}^j g_\ell} \cdot \prod_{\ell=j+2}^i \omega'_\ell \right) \omega'_{i+1} - \frac{q_i}{\prod_{\ell=1}^i g_\ell} \cdot t_i \end{aligned}$$

$$\begin{aligned}
&= \omega_i \cdot \omega'_{i+1} - \frac{q_i}{\prod_{\ell=1}^i g_\ell} \cdot t_i \\
&= \omega_{i+1}.
\end{aligned}$$

Por el principio de inducción se sigue que (1.32) satisface (1.31) para todo $i \in \mathbb{N}$. Así, esta fórmula es la forma cerrada de la relación de recurrencia propuesta. \square

Ahora que encontramos una forma cerrada a la relación de recurrencia (1.31), somos capaces de establecer una transformación lineal entre el vector de variables libres $\mathbf{t} \in \mathbb{Z}^{n-1}$ y el vector solución $\mathbf{x} \in \mathbb{Z}^n$ de (1.14). Definamos, por conveniencia, los coeficientes $m_{ij} \in \mathbb{Z}$ con $i > j$ como

$$m_{ij} := \frac{q_j}{\prod_{\ell=1}^j g_\ell} \cdot \prod_{\ell=j+2}^i \omega'_\ell. \quad (1.33)$$

Sustituyendo en la forma cerrada (1.32), obtenemos la fórmula simplificada

$$\omega_i = k \cdot \prod_{j=2}^i \omega'_j - \sum_{j=1}^{i-1} m_{ij} t_j, \quad (1.34)$$

Así pues, juntando esto último con (1.25), obtenemos para $i \in \{1, \dots, n-2\}$:

$$\begin{aligned}
x_i &= \omega_i \cdot x'_i + g_{i+1} t_i \\
&= k \cdot \prod_{j=2}^i \omega'_j \cdot x'_i - \sum_{j=1}^{i-1} m_{ij} x'_i t_j + g_{i+1} t_i.
\end{aligned} \quad (1.35)$$

Similarmente, usando (1.34) y sustituyendo en (1.28), llegamos a

$$x_{n-1} = k \cdot \prod_{j=2}^{n-1} \omega'_j \cdot x'_{n-1} - \sum_{j=1}^{n-2} m_{n-1,j} x'_{n-1} t_j + \frac{q_n}{\prod_{j=1}^{n-1} g_j} t_{n-1}, \quad (1.36a)$$

$$x_n = k \cdot \prod_{j=2}^{n-1} \omega'_j \cdot x'_n - \sum_{j=1}^{n-2} m_{n-1,j} x'_n t_j - \frac{q_{n-1}}{\prod_{j=1}^{n-1} g_j} t_{n-1}. \quad (1.36b)$$

Así pues, definimos $\boldsymbol{\nu} \in \mathbb{Z}^n$ como

$$\nu_i := x'_i \cdot \prod_{j=2}^{\min\{i, n-1\}} \omega'_j. \quad (1.37)$$

También definimos la matriz $M \in \mathbb{Z}^{n \times (n-1)}$ a través de

$$M_{ij} := \begin{cases} -m_{ij}x'_i, & j < i, \\ g_{i+1}, & i = j < n-1, \\ \frac{q_n}{\prod_{k=1}^{n-1} g_k}, & i = j = n-1, \\ -\frac{q_{n-1}}{\prod_{k=1}^{n-1} g_k}, & i = n, j = n-1, \\ 0, & \text{e.o.c.} \end{cases} \quad (1.38)$$

Sustituyendo estas definiciones en (1.35) y (1.36) encontramos que

$$\mathbf{x} = k\boldsymbol{\nu} + M\mathbf{t}. \quad (1.39)$$

En la subsección 1.2.2 mencionamos a lo largo de la construcción de soluciones que los coeficientes de Bézout ω'_i, x'_i están asociados a términos exclusivamente dependientes de \mathbf{q} , por lo que no dependen de la elección $\mathbf{x} \in \mathbb{Z}$. Se sigue de (1.37) y de (1.38) que $\boldsymbol{\nu}$ y M dependen exclusivamente de \mathbf{q} y no de \mathbf{x} .

Lema 1.2.11. *Sea $\mathbf{q} \in \mathbb{Z}^n$ un vector coprimo. Entonces el vector $\boldsymbol{\nu} \in \mathbb{Z}^n$ definido en (1.37) satisface $\mathbf{q}^T \boldsymbol{\nu} = 1$.*

Demostración. Primero mostramos por inducción hacia atrás que se cumple

$$\sum_{j=i}^n q_j \nu_j = \prod_{j=2}^i \omega'_j \cdot \prod_{j=1}^i g_j, \quad (1.40)$$

para todo $i \in \{1, \dots, n-1\}$. Empezamos con el caso base $i = n-1$. De

(1.37), encontramos que

$$q_{n-1}\nu_{n-1} + q_n\nu_n = \prod_{j=2}^{n-1} \omega'_j \cdot (q_{n-1}x'_{n-1} + q_nx'_n). \quad (1.41)$$

Recordemos que x'_{n-1} y x'_n son coeficientes de Bézout asociados a los coeficientes del lado izquierdo de (1.27), los cuales son coprimos. Entonces se cumple, por el teorema 1.1.4, que

$$\frac{q_{n-1}}{\prod_{j=1}^{n-1} g_j} x'_{n-1} + \frac{q_n}{\prod_{j=1}^{n-1} g_j} x'_n = 1,$$

o, equivalentemente,

$$q_{n-1}x'_{n-1} + q_nx'_n = \prod_{j=1}^{n-1} g_j.$$

Sustituyendo en (1.41), obtenemos la base de la inducción:

$$q_{n-1}\nu_{n-1} + q_n\nu_n = \prod_{j=2}^{n-1} \omega'_j \cdot \prod_{j=1}^{n-1} g_j.$$

Supongamos que (1.40) se satisface para alguna $i \in \{2, \dots, n-1\}$. Reduciendo i , ocupando (1.37) y usando la hipótesis inductiva, obtenemos

$$\begin{aligned} \sum_{j=i-1}^n q_j\nu_j &= q_{i-1}\nu_{i-1} + \sum_{j=i}^n q_j\nu_j \\ &= \prod_{j=2}^{i-1} \omega'_j \cdot q_{i-1}x'_{i-1} + \prod_{j=2}^i \omega'_j \cdot \prod_{j=1}^i g_j \\ &= \prod_{j=2}^{i-1} \omega'_j \cdot \left(q_{i-1}x'_{i-1} + \omega'_i \prod_{j=1}^i g_j \right). \end{aligned}$$

Nuevamente, x'_{i-1} y ω'_i son coeficientes de Bézout asociados, respectivamente, a $\frac{q_{i-1}}{\prod_{j=1}^{i-1} g_j}$ y g_i , los cuales son coprimos. De esta manera satisfacen (1.26)

pero sustituyendo i por $i - 1$. Es decir, se satisface

$$\frac{q_{i-1}}{\prod_{j=1}^{i-1} g_j} x'_{i-1} + g_i \omega'_i = 1,$$

o, equivalentemente,

$$q_{i-1} x'_{i-1} + \omega'_i \prod_{j=1}^i g_j = \prod_{j=1}^{i-1} g_j.$$

Sustituyendo, obtenemos el resultado (1.40) para $i - 1$. Así, por inducción hacia atrás, (1.40) se cumple para todo $i \in \{1, \dots, n - 1\}$. Finalmente, para demostrar este lema, observamos que

$$\mathbf{q}^T \boldsymbol{\nu} = \sum_{j=1}^n q_j \nu_j = \prod_{j=2}^1 \omega'_j \cdot \prod_{j=1}^1 g_j = g_1 = 1.$$

El primer producto es uno por ser el producto vacío. Recordemos también que g_1 es el máximo común divisor de q_1, \dots, q_n , los cuales son coprimos por hipótesis, y entonces $g_1 = 1$. \square

Lema 1.2.12. *Si $q_n \neq 0$ entonces $\text{gen}\{\mathbf{q}\} = \ker\{M^T\}$, donde la matriz M está definida en (1.38).*

Demostración. La matriz M es triangular inferior y su diagonal principal es distinta de cero. En efecto, para todo $i \in \{1, \dots, n - 2\}$, tenemos

$$M_{ii} = g_{i+1} = \text{mcd} \left\{ \frac{q_i}{\prod_{j=1}^i g_j}, \dots, \frac{q_n}{\prod_{j=1}^i g_j} \right\}.$$

Pero el máximo común divisor entre cualesquiera enteros siempre es positivo. También tenemos por hipótesis que

$$M_{n-1, n-1} = \frac{q_n}{\prod_{j=1}^{n-1} g_j} \neq 0.$$

Se sigue que las columnas de M son linealmente independientes, y entonces su imagen tiene dimensión $n - 1$. Por lo tanto, M^T tiene $n - 1$ renglones linealmente independientes. Se sigue por el teorema de la Dimensión que $\dim \ker\{M^T\} = 1$, así que basta mostrar que $\mathbf{q} \in \ker\{M^T\}$.

Sea $\mathbf{x} \in \mathbb{Z}^n$. Por el teorema 1.2.5, existe una capa entera $H_{\mathbf{q}, k\|\mathbf{q}\|^{-2}}$ que contiene a \mathbf{x} . Así, por el lema 1.2.6, \mathbf{x} satisface la ecuación lineal diofantina $\mathbf{q}^T \mathbf{x} = k$. Por construcción en la subsección 1.2.2 y por la exhaustividad del teorema 1.1.10 sabemos que existe $\mathbf{t} \in \mathbb{Z}^{n-1}$ tal que $\mathbf{x} = k\boldsymbol{\nu} + M\mathbf{t}$. Luego, por el lema 1.2.11, tenemos

$$k = \mathbf{q}^T \mathbf{x} = k\mathbf{q}^T \boldsymbol{\nu} + \mathbf{q}^T M\mathbf{t} = k + (\mathbf{q}^T M)\mathbf{t}.$$

De donde obtenemos $(\mathbf{q}^T M)\mathbf{t} = 0$. Pero \mathbf{x} fue arbitrario, así que también lo fue \mathbf{t} . Entonces $\mathbf{q}^T M = \mathbf{0}^T$, lo que implica $\mathbf{q} \in \ker\{M^T\}$. \square

La gran mayoría de nuestra argumentación para demostrar los resultados ha sido fundamentada a través de las capas enteras $H_{\mathbf{q}, k\|\mathbf{q}\|^{-2}}$, al igual que por el teorema 1.2.5. Sin embargo, estas capas enteras contienen puntos que no son enteros y por lo tanto no son de nuestro interés. Nos gustaría concentrarnos exclusivamente en puntos enteros, al mismo tiempo que buscamos caracterizarlos por medio de \mathbf{q} .

Definición 1.2.13 ([Sch98]). Decimos que un subconjunto Λ de \mathbb{R}^n es un **grupo aditivo** si

1. $\mathbf{0} \in \Lambda$, y
2. si $\mathbf{x}, \mathbf{y} \in \Lambda$, entonces $\mathbf{x} + \mathbf{y} \in \Lambda$, y también $-\mathbf{x} \in \Lambda$.

Además, decimos que Λ es una **red** si existen vectores $\mathbf{v}_1, \dots, \mathbf{v}_n$ linealmente independientes tales que

$$\Lambda = \{\lambda_1 \mathbf{v}_1 + \dots + \lambda_n \mathbf{v}_n : \lambda_i \in \mathbb{Z}\}.$$

A los vectores $\mathbf{v}_1, \dots, \mathbf{v}_n$ los llamamos la **base de la red** Λ .

Ejemplo 1.2.14. No es difícil ver que \mathbb{Z}^n es un grupo aditivo. Si consideramos los vectores canónicos $\mathbf{e}_1, \dots, \mathbf{e}_n$, entonces encontramos que son linealmente independientes, pero también se cumple

$$\mathbb{Z}^n = \{\lambda_1 \mathbf{e}_1 + \dots + \lambda_n \mathbf{e}_n : \lambda_i \in \mathbb{Z}\}.$$

De esta, manera \mathbb{Z}^n es una red que tiene como base canónica a los vectores $\mathbf{e}_1, \dots, \mathbf{e}_n$.

Teorema 1.2.15. Sean $\mathbf{q} \in \mathbb{Z}^n$ un vector coprimo y supongamos que q_n es distinto de cero. Entonces $\boldsymbol{\nu}$ y las columnas de M (definidas en (1.37) y (1.38), respectivamente) forman una base de la red \mathbb{Z}^n .

Demostración. En el lema 1.2.12 mostramos que las columnas de M son linealmente independientes. Mostramos por contradicción que $\boldsymbol{\nu}$ es linealmente independiente de las columnas de M , así que supongamos que no lo es, por lo que existen escalares $\lambda_1, \dots, \lambda_{n-1}$ tales que

$$\boldsymbol{\nu} = \lambda_1 \mathbf{m}_1 + \dots + \lambda_{n-1} \mathbf{m}_{n-1},$$

donde $\mathbf{m}_1, \dots, \mathbf{m}_{n-1}$ son las columnas de M . De los lemas 1.2.11 y 1.2.12 obtenemos

$$1 = \mathbf{q}^T \boldsymbol{\nu} = \lambda_1 \mathbf{q}^T \mathbf{m}_1 + \dots + \lambda_{n-1} \mathbf{q}^T \mathbf{m}_{n-1} = 0,$$

lo cual es una contradicción. Se sigue que $\{\boldsymbol{\nu}, \mathbf{m}_1, \dots, \mathbf{m}_{n-1}\}$ es un conjunto de vectores linealmente independiente.

Ahora bien, sea $\mathbf{x} \in \mathbb{Z}^n$, por el teorema 1.2.5, sabemos que $\mathbf{x} \in H_{\mathbf{q}, k \|\mathbf{q}\|^{-2}}$ para alguna $k \in \mathbb{Z}$. Por el lema 1.2.6, \mathbf{x} satisface la ecuación lineal diofantina $\mathbf{q}^T \mathbf{x} = k$. Por construcción en la subsección 1.2.2 y por la exhaustividad del

teorema 1.1.10 sabemos que existe $\mathbf{t} \in \mathbb{Z}^{n-1}$ que satisface

$$\mathbf{x} = k\boldsymbol{\nu} + M\mathbf{t} = k\boldsymbol{\nu} + t_1\mathbf{m}_1 + \cdots + t_{n-1}\mathbf{m}_{n-1}.$$

Pero \mathbf{x} fue arbitrario, lo que implica que

$$\mathbb{Z}^n = \{k\boldsymbol{\nu} + t_1\mathbf{m}_1 + \cdots + t_{n-1}\mathbf{m}_{n-1} : k, t_1, \dots, t_{n-1} \in \mathbb{Z}\}$$

De esta manera, se cumple que $\{\boldsymbol{\nu}, \mathbf{m}_1, \dots, \mathbf{m}_{n-1}\}$ forma una base de la red \mathbb{Z}^n . \square

El siguiente corolario es presentado sin demostración, pero cabe mencionar que es una consecuencia directa del teorema anterior junto con las equivalencias encontradas en el teorema 4.3 de [Sch98].

Corolario 1.2.16. *Sea $\mathbf{q} \in \mathbb{Z}^n$ un vector coprimo y supongamos que q_n es distinto de cero. Consideremos $\boldsymbol{\nu} \in \mathbb{Z}^n$ y las columnas $\mathbf{m}_1, \dots, \mathbf{m}_{n-1} \in \mathbb{Z}^n$ de la matriz M (definidas en (1.37) y (1.38), respectivamente), entonces la matriz*

$$[\boldsymbol{\nu} \mid \mathbf{m}_1 \mid \cdots \mid \mathbf{m}_{n-1}] \in \mathbb{Z}^{n \times n}$$

es unimodular, es decir, su determinante es ± 1 .

Geométricamente, tenemos que si $\mathbf{q} \in \mathbb{Z}^n$ es un vector coprimo tal que $q_n \neq 0$, entonces \mathbf{q} induce una descomposición de \mathbb{Z}^n como la suma directa de las subredes Λ_p y Λ_h , donde

$$\Lambda_p := \{k\boldsymbol{\nu} : k \in \mathbb{Z}\}, \quad \Lambda_h := \{M\mathbf{t} : \mathbf{t} \in \mathbb{Z}^{n-1}\}. \quad (1.42)$$

Para todo vector $\mathbf{x} \in \Lambda_p$ existe una $k \in \mathbb{Z}$ tal que $\mathbf{x} = k\boldsymbol{\nu}$. Luego, por el lema 1.2.11, \mathbf{x} es una solución particular de la ecuación lineal diofantina $\mathbf{q}^T \mathbf{x} = k$. Similarmente, por el lema 1.2.12, todo vector $\mathbf{x} \in \Lambda_h$ es una solución de la ecuación lineal diofantina homogénea $\mathbf{q}^T \mathbf{x} = 0$.

Sea $\mathbf{x} \in \mathbb{Z}^n$. Por un lado, el teorema 1.2.5 indica que $\mathbf{x} \in H_{\mathbf{q}, k\|\mathbf{q}\|}^{-2}$ para alguna $k \in \mathbb{Z}$. Por el otro lado, la descomposición $\mathbb{Z}^n \cong \Lambda_p \oplus \Lambda_h$ es tal que \mathbf{x} puede ser escrito como $\mathbf{x}_p + \mathbf{x}_h$, donde $\mathbf{x}_p \in \Lambda_p$ y $\mathbf{x}_h \in \Lambda_h$, así que se satisface simultáneamente $\mathbf{q}^T \mathbf{x}_p = k$ y $\mathbf{q}^T \mathbf{x}_h = 0$.

De manera intuitiva, \mathbf{x}_p se encarga de que \mathbf{x} se encuentre sobre la k -ésima capa entera, mientras que \mathbf{x}_h se desliza sobre esta capa. Esta idea de descomponer el espacio a partir de soluciones particulares y homogéneas de una ecuación no es novedosa, pues es sumamente frecuente hacerlo para espacios vectoriales.

Observemos que si $q_n = 0$, es válido permutar las entradas de \mathbf{q} de manera que el vector permutado $\tilde{\mathbf{q}}$ cumpla el supuesto $\tilde{q}_n \neq 0$. Podemos preguntarnos cómo se relacionan las imágenes de las matrices M y \tilde{M} de estos dos vectores. Pero si $q_n = 0$, entonces puede ser que M no esté bien definida³. Requerimos de un supuesto más fuerte para responder la pregunta.

Corolario 1.2.17. *Sea \mathbf{q} un vector coprimo y sea $\tilde{\mathbf{q}}$ un vector con las entradas de \mathbf{q} permutadas. Supongamos que $q_n, \tilde{q}_n \neq 0$. Entonces sus respectivas matrices M y \tilde{M} definidas en (1.38) son isomorfas:*

$$\ker\{M^T\} \cong \ker\{\tilde{M}^T\}.$$

Demostración. Existe una matriz de permutación $P \in \mathbb{Z}^{n \times n}$ tal que $\tilde{\mathbf{q}} = P\mathbf{q}$. Puesto que P es invertible, tenemos $\text{gen}\{\mathbf{q}\} \cong \text{gen}\{\tilde{\mathbf{q}}\}$. Usando el lema 1.2.12 obtenemos

$$\ker\{M^T\} = \text{gen}\{\mathbf{q}\} \cong \text{gen}\{\tilde{\mathbf{q}}\} = \ker\{\tilde{M}^T\},$$

³Por ejemplo, si $q_n = q_{n-1} = 0$, encontramos que

$$g_{n-1} := \text{mcd}\left\{\frac{q_{n-1}}{\prod_{j=1}^{n-2} g_j}, \frac{q_n}{\prod_{j=1}^{n-2} g_j}\right\} = \text{mcd}\{0, 0\}.$$

Pero el máximo común divisor de dos números no está bien definido si ambos son cero. Esto implica que la entrada $M_{n-2, n-2} := g_{n-1}$ no está bien definida.

que es lo que queríamos demostrar. \square

Observación. Si $\tilde{\mathbf{q}} = P\mathbf{q}$ no es cierto que $\tilde{M} = PM$. Por ejemplo, consideremos el vector $\mathbf{q} := (1, 1, -2)^T$ y la matriz de permutación

$$P := \begin{pmatrix} 1 & 0 & 0 \\ 0 & 0 & 1 \\ 0 & 1 & 0 \end{pmatrix},$$

de donde obtenemos $\tilde{\mathbf{q}} = (1, -2, 1)^T$. Tenemos

$$M = \begin{pmatrix} 1 & 0 \\ 1 & -2 \\ 1 & -1 \end{pmatrix}, \quad \tilde{M} = \begin{pmatrix} 1 & 0 \\ 1 & 1 \\ 1 & 2 \end{pmatrix}, \quad PM = \begin{pmatrix} 1 & 0 \\ 1 & -1 \\ 1 & -2 \end{pmatrix}.$$

Se verifica inmediatamente que $\tilde{M} \neq PM$.

En esta última parte del capítulo, exploramos las consecuencias del corolario 1.2.17. Esto nos llevará de regreso a justificar de forma algebraica o geométrica por qué es posible ignorar las entradas nulas del vector \mathbf{q} , o tan siquiera por qué podemos concentrarnos en este vector coprimo en vez del vector esencialmente entero \mathbf{p} en el problema original (1.1).

Definición 1.2.18. Sea $\mathbf{q} \in \mathbb{Z}^n$ un vector coprimo con entradas distintas de cero, entonces definimos su **órbita** como

$$\text{orb}(\mathbf{q}) := \{P\mathbf{q} : P \in \mathbb{Z}^{n \times n} \text{ es matriz de permutación}\}.$$

Lema 1.2.19. Sean $\mathbf{q}, \tilde{\mathbf{q}} \in \mathbb{Z}^n$ vectores coprimos con entradas distintas de cero. Entonces $\tilde{\mathbf{q}} \in \text{orb}(\mathbf{q})$ si y solo si $\text{orb}(\tilde{\mathbf{q}}) = \text{orb}(\mathbf{q})$.

Demostración. Supongamos que $\tilde{\mathbf{q}} \in \text{orb}(\mathbf{q})$, luego existe una matriz de permutación $P \in \mathbb{Z}^{n \times n}$ tal que $\tilde{\mathbf{q}} = P\mathbf{q}$. Ahora bien, sea $\hat{\mathbf{q}} \in \text{orb}(\tilde{\mathbf{q}})$, por lo que existe una matriz de permutación $P' \in \mathbb{Z}^{n \times n}$ tal que $\hat{\mathbf{q}} = P'\tilde{\mathbf{q}} = (P'P)\mathbf{q}$.

Como el producto de matrices de permutación también es una matriz de permutación, se sigue que $\hat{\mathbf{q}} \in \text{orb}(\mathbf{q})$. Con esto mostramos la contención $\text{orb}(\tilde{\mathbf{q}}) \subseteq \text{orb}(\mathbf{q})$. La otra contención se muestra de manera análoga, pues $\mathbf{q} = P^{-1}\tilde{\mathbf{q}}$ y la inversa de una matriz de permutación también es una matriz de permutación.

Supongamos que $\text{orb}(\tilde{\mathbf{q}}) = \text{orb}(\mathbf{q})$. Entonces basta mostrar que $\tilde{\mathbf{q}} \in \text{orb}(\tilde{\mathbf{q}})$, pero esto es inmediatamente cierto puesto que $\tilde{\mathbf{q}} = I_n \tilde{\mathbf{q}}$ y la matriz identidad $I_n \in \mathbb{Z}^{n \times n}$ es una matriz de permutación. \square

Lema 1.2.20. *Sea $\mathbf{q} \in \mathbb{Z}^n$ un vector coprimo con entradas distintas de cero y sea $\tilde{\mathbf{q}} \in \text{orb}(\mathbf{q})$. Entonces las redes $\tilde{\Lambda}_h$ y Λ_h definidas en (1.42) son isomorfas. Similarmente, las redes $\tilde{\Lambda}_p$ y Λ_p son isomorfas.*

Demostración. Se sigue inmediatamente de la definición 1.2.18 junto con el corolario 1.2.17. \square

El lema anterior nos permite identificar la descomposición de \mathbb{Z}^n a partir de $\text{orb}(\mathbf{q})$ y no solamente de \mathbf{q} . Ahora bien, supongamos que $\mathbf{q} \in \mathbb{Z}^{n+\ell}$ es un vector coprimo con ℓ entradas nulas. Definamos el conjunto ordenado de índices no nulos de \mathbf{q} como $\sigma = (i : q_i \neq 0)$, y también definamos la proyección $\pi^{(n+\ell)} : \mathbb{Z}^{n+\ell} \rightarrow \mathbb{Z}^n$ a partir de $\pi^{(n+\ell)}(\mathbf{q})_i := q_{\sigma_i}$. Luego, $\pi^{(n+\ell)}(\mathbf{q})$ es un vector coprimo con entradas distintas de cero. Por el teorema 1.2.15 (y también por la discusión que le sucede), sabemos que $\pi^{(n+\ell)}(\mathbf{q})$ induce la descomposición

$$\mathbb{Z}^n \cong \Lambda_p \oplus \Lambda_h,$$

donde Λ_p y Λ_h están definidas en (1.42). Pero en el lema 1.2.20 vimos que todo vector $\tilde{\mathbf{q}} \in \text{orb}(\pi^{(n+\ell)}(\mathbf{q}))$ induce una descomposición isomorfa. Esto justifica el hecho de que podamos ignorar las entradas nulas del vector coprimo \mathbf{q} , y también motiva el siguiente resultado.

Teorema 1.2.21. Sean $\mathbf{q} \in \mathbb{Z}^n$ y $\tilde{\mathbf{q}} \in \mathbb{Z}^m$ vectores coprimos con $n < m$. Si $\text{orb}(\pi^{(n)}(\mathbf{q})) = \text{orb}(\pi^{(m)}(\tilde{\mathbf{q}}))$, entonces \mathbf{q} y $\tilde{\mathbf{q}}$ inducen descomposiciones isomorfas de una subred de \mathbb{Z}^n .

Demostración. Supongamos, sin pérdida de generalidad, que \mathbf{q} no tiene entradas nulas. Luego, $\pi^{(n)}(\mathbf{q}) = \mathbf{q}$. Entonces $\tilde{\mathbf{q}}$ tiene $m - n$ entradas no nulas, pues de otra forma las dimensiones de $\pi^{(n)}(\mathbf{q})$ y $\pi^{(m)}(\tilde{\mathbf{q}})$ serían distintas y por lo tanto sus órbitas no serían iguales. Tenemos por hipótesis y del lema 1.2.19 que $\pi^{(m)}(\tilde{\mathbf{q}}) \in \text{orb}(\mathbf{q})$. Finalmente, del lema 1.2.20 junto con el teorema 1.2.15 obtenemos lo que queríamos demostrar. \square

Definición 1.2.22. Sean $\mathbf{p} \in \mathbb{R}^n$ y $\tilde{\mathbf{p}} \in \mathbb{R}^m$ vectores esencialmente enteros (ver definición 1.2.1). Entonces decimos que \mathbf{p} y $\tilde{\mathbf{p}}$ son equivalentes si y solo si $\text{orb}(\pi^{(n)}(\mathbf{q})) = \text{orb}(\pi^{(m)}(\tilde{\mathbf{q}}))$, donde \mathbf{q} y $\tilde{\mathbf{q}}$ son sus respectivos múltiplos coprimos. En este caso escribimos $\mathbf{p} \sim \tilde{\mathbf{p}}$.

Puesto que el múltiplo coprimo de un vector esencialmente entero es único (ver la discusión al inicio de la subsección 1.2.1), esta relación está bien definida. Luego, por las propiedades de la igualdad, tenemos que esta relación es una de equivalencia sobre el conjunto de vectores esencialmente enteros, independientemente de su dimensión.

Con esta definición y el teorema 1.2.21 encontramos que si los vectores esencialmente enteros $\mathbf{p} \in \mathbb{R}^n$ y $\tilde{\mathbf{p}} \in \mathbb{R}^m$ son equivalentes, entonces el programa lineal (1.1) es esencialmente el mismo.

En conclusión, encontramos una suerte de clasificación de este tipo de programas lineales. Hemos visto que el único representante importante para resolver estos problemas es el múltiplo coprimo $\mathbf{q} \in \mathbb{Z}^n$ o, más bien, una proyección de este de manera que no tenga entradas nulas. Solamente cuando querramos relacionar los resultados con una instancia específica del problema (1.1) haremos mención del vector esencialmente entero $\mathbf{p} \in \mathbb{R}^n$.

Capítulo 2

El caso infinito

Sea $\mathbf{p} \in \mathbb{R}^n \setminus \{\mathbf{0}\}$ un vector esencialmente entero y recordemos de la definición 1.2.1 que tiene un único múltiplo coprimo $\mathbf{q} \in \mathbb{Z}^n$. Es decir, existe un único escalar $m \in \mathbb{R}$ que satisface tres cosas: $\mathbf{p} = m\mathbf{q}$, las entradas q_1, \dots, q_n son coprimas, y la primera entrada no nula q_i es positiva. Al igual que en el capítulo anterior, supondremos que m es positivo. Equivalentemente, supondremos que la primera entrada no nula p_i es positiva¹.

Retomemos el entero $\eta \in \mathbb{Z}$ del lema 1.2.7 que parametriza la primera capa entera que satisface el presupuesto (1.1b). A causa del teorema 1.2.9 sabemos que si $q_i \leq 0$ para alguna $i \in \{1, \dots, n\}$, entonces la η -ésima capa entera contiene un número infinito de puntos factibles. A partir de esto último, somos capaces de resolver automáticamente el problema de decisión de determinar si un escalar $u^* \in \mathbb{R}$ es el valor óptimo del programa (1.1).

Corolario 2.0.1. *Supongamos que $q_i \leq 0$ para algún $i \in \{1, \dots, n\}$. Entonces el valor óptimo del programa lineal entero (1.1) es $m\eta$. Además, si m es positivo, tenemos que η es el múltiplo de m más grande que satisface*

¹El autor hace recordar que esta es una cuestión puramente de comodidad y no hay pérdida de generalidad. Cuando m es negativo, los resultados se mantienen pero es necesario voltear las desigualdades y cambiar las funciones piso por las funciones techo, lo cual añadiría un número innecesario de casos a analizar.

$m\eta \leq u$, donde u es el lado derecho de la restricción presupuestaria (1.1b).

Demostración. Por el teorema 1.2.9 sabemos que existen una infinidad de soluciones en la η -ésima capa entera, así que sea \mathbf{x}^* una de ellas. Entonces $\mathbf{q}^T \mathbf{x}^* = \eta$, pero $\mathbf{p} = m\mathbf{q}$ por la definición 1.2.1, por lo que obtenemos $\mathbf{p}^T \mathbf{x}^* = m\mathbf{q}^T \mathbf{x}^* = m\eta$.

Ahora bien, si m es positivo, por el lema 1.2.7 tenemos que $\eta = \lfloor u/m \rfloor$. Supongamos que $\xi \in \mathbb{Z}$ satisface $m\xi \leq u$ y también $\lfloor u/m \rfloor < \xi$. Luego,

$$m \left\lfloor \frac{u}{m} \right\rfloor < m\xi \leq u \implies \left\lfloor \frac{u}{m} \right\rfloor < \xi \leq \frac{u}{m},$$

pero esto contradice las propiedades de la función piso. \square

Observación. Para ilustrar la conveniencia de restringir m a que sea positivo, consideremos el caso cuando $m < 0$. De una manera similar a la del lema 1.2.7, podemos demostrar que $\eta := \lceil u/m \rceil$ parametriza también la primera capa entera que satisface el presupuesto, pues ahora tenemos de la restricción (1.1b) que $\mathbf{p}^T \mathbf{x} \leq u$ si y solo si $\mathbf{q}^T \mathbf{x} \geq u/m$. Se sigue cumpliendo que el valor óptimo del problema (1.1) es $m\eta$. Sin embargo, η ahora es el múltiplo más chico de m que satisface $m\eta \geq u$.

Una vez resuelto el problema de decisión, podemos preguntarnos concretamente cómo obtener el punto óptimo. Por el teorema 1.2.9 sabemos que debemos resolver la ecuación lineal diofantina $\mathbf{q}^T \mathbf{x} = \eta$. Del teorema 1.2.15 sabemos que si $q_n \neq 0$, entonces existen $k \in \mathbb{Z}$ y $\mathbf{t} \in \mathbb{Z}^{n-1}$ tales que

$$\mathbf{x} = k\boldsymbol{\nu} + M\mathbf{t},$$

donde $\boldsymbol{\nu}$ y M están definidas por (1.37) y (1.38), respectivamente. De los lemas 1.2.11 y 1.2.12 sabemos que

$$\mathbf{q}^T (\eta\boldsymbol{\nu} + M\mathbf{t}) = \eta\mathbf{q}^T \boldsymbol{\nu} + \mathbf{q}^T M\mathbf{t} = \eta$$

para todo $\mathbf{t} \in \mathbb{Z}^{n-1}$. Así pues, debe ser el caso que $k = \eta$ y debemos encontrar condiciones suficientes en \mathbf{t} para asegurar la no-negatividad de \mathbf{x} . En primer lugar, sabemos que para todo $i \in \{1, \dots, n-2\}$, la entrada t_i debe satisfacer (1.30). En segundo lugar, recuperamos de (1.28) que las últimas dos soluciones de la ecuación $\mathbf{q}^T \mathbf{x} = \eta$ están dadas por

$$\begin{cases} x_{n-1} = \omega_{n-1} x'_{n-1} + \frac{q_n}{\prod_{j=1}^{n-1} g_j} t_{n-1}, \\ x_n = \omega_{n-1} x'_n - \frac{q_{n-1}}{\prod_{j=1}^{n-1} g_j} t_{n-1}, \end{cases}$$

donde los enteros g_i están definidos por (1.23) con $g_1 = 1$, ω_{n-1} está definida a través de la relación de recurrencia (1.31) con condición inicial $\omega_1 = \eta$ (o bien a partir del lema 1.32 con $k = \eta$), y x'_{n-1}, x'_n son coeficientes de Bézout que satisfacen (1.29).

En lo que se encuentra a continuación supondremos que ninguna entrada q_i es nula. Esto no constituye problema alguno debido a un razonamiento similar al de la demostración del teorema 1.2.9. Definimos

$$I^\circ := \{i : q_i = 0\},$$

y también definimos el vector $\tilde{\mathbf{q}}$ cuyas entradas son las entradas no nulas de \mathbf{q} . A partir de lo que sigue vamos a determinar un vector entero no nulo $\tilde{\mathbf{x}}$ que satisfaga $\tilde{\mathbf{q}}^T \tilde{\mathbf{x}} = \eta$. Luego, encontramos que el vector \mathbf{x} dado por

$$x_i := \begin{cases} \tilde{x}_i, & i \notin I^\circ, \\ 0, & i \in I^\circ, \end{cases}$$

es entero, no negativo, y también satisface $\mathbf{q}^T \mathbf{x} = \eta$. Así pues, la suposición de que $q_i \neq 0$ para todo $i \in \{1, \dots, n\}$ toma lugar sin pérdida de generalidad.

Para que se satisfagan las condiciones de no negatividad de x_{n-1} y de x_n , encontramos que $t_{n-1} \in \mathbb{Z}$ debe cumplir ciertas desigualdades según los

signos de q_{n-1} y de q_n . Definamos, por conveniencia,

$$b_1 := -\frac{\omega_{n-1}x'_{n-1}}{q_n} \cdot \prod_{j=1}^{n-1} g_j, \quad b_2 := \frac{\omega_{n-1}x'_n}{q_{n-1}} \cdot \prod_{j=1}^{n-1} g_j. \quad (2.1)$$

Entonces, para asegurar la no-negatividad de x_{n-1} y de x_n , debe ser el caso que

$$t_{n-1} \in \begin{cases} [\lceil \max\{b_1, b_2\} \rceil, \infty), & q_{n-1} < 0 < q_n, \\ (-\infty, \lfloor \min\{b_1, b_2\} \rfloor], & q_n < 0 < q_{n-1}, \\ [\lceil b_2 \rceil, \lfloor b_1 \rfloor], & q_{n-1}, q_n < 0, \\ [\lceil b_1 \rceil, \lfloor b_2 \rfloor], & 0 < q_{n-1}, q_n. \end{cases} \quad (2.2)$$

Podemos emplear la misma estrategia de permutar las entradas de q_i de manera que colapsemos estos cuatro casos distintos en uno solo. Como estamos en el caso infinito del teorema 1.2.9, naturalmente supondremos que $q_i < 0$ para alguna $i \in \{1, \dots, n\}$. Así pues, podemos permutar esta i -ésima entrada de \mathbf{q} con q_{n-1} , con lo que obtenemos $q_{n-1} < 0$. Luego, como \mathbf{q} es el múltiplo coprimo de \mathbf{p} y ninguna entrada de \mathbf{q} es nula, se sigue de la definición 1.2.1 que $q_1 > 0$. Así pues, podemos permutar la primera y última entrada de \mathbf{q} , de donde se sigue que $q_n > 0$. Juntándolo todo, obtenemos $q_{n-1} < 0 < q_n$. De esta manera, para asegurar la no negatividad de x_{n-1} y x_n , basta con que se satisfaga el primer caso:

$$t_{n-1} \geq \lceil \max\{b_1, b_2\} \rceil. \quad (2.3)$$

Lema 2.0.2. *Sea $\mathbf{p} \in \mathbb{R}$ un vector cuyas entradas son todas distintas de cero, y sea $\mathbf{q} \in \mathbb{Z}^n$ su múltiplo coprimo. Entonces existe un vector $\mathbf{t} \in \mathbb{Z}^{n-1}$ que satisface ambos (1.30) y (2.2).*

Demostración. Por la discusión anterior, podemos suponer sin pérdida de generalidad que $q_{n-1} < 0 < q_n$, así que basta mostrar la existencia de

$\mathbf{t} \in \mathbb{Z}^{n-1}$ que satisfaga (1.30) y (2.3). Si definimos

$$t_i := \begin{cases} \left\lceil -\frac{\omega_i x'_i}{g_{i+1}} \right\rceil, & i < n-1, \\ \lceil \max\{b_1, b_2\} \rceil, & i = n-1, \end{cases}$$

entonces se verifica automáticamente que estas condiciones se satisfacen. \square

En síntesis, por el lema 2.0.2 sabemos que existe un vector $\mathbf{t} \in \mathbb{Z}^{n-1}$ que satisface ambos (1.30) y (2.2). Al definir $\mathbf{x}^* := \eta \boldsymbol{\nu} + M \mathbf{t}$, encontramos que \mathbf{x}^* es entero y no negativo, y además por los lemas 1.2.11 y 1.2.12 encontramos que

$$\mathbf{q}^T \mathbf{x}^* = \eta \mathbf{q}^T \boldsymbol{\nu} + \mathbf{q}^T M \mathbf{t} = \eta.$$

Por el teorema 1.2.9, se sigue que \mathbf{x}^* es la solución al problema (1.1).

En la práctica es mejor usar la relación de recurrencia (1.25) y “construir” las entradas x_i al mismo tiempo que definimos t_i de manera que satisfaga (1.30) y (2.3). Si procedemos de esta forma no tenemos que encontrar primero $\boldsymbol{\nu}$ y M , determinar \mathbf{t} y luego recuperar \mathbf{x} . El algoritmo 1 muestra este procedimiento constructivo.

En el algoritmo 1 supusimos la existencia de una subrutina **Bezout** que, como su nombre lo indica, calcula los coeficientes de Bézout entre dos enteros. Es la creencia del autor que no es necesario escribir la subrutina en esta tesis, pero reitera, así como lo hizo en la Sección 1.1.1, que estos coeficientes se pueden calcular por medio del algoritmo Extendido de Euclides.

Lema 2.0.3. *El algoritmo 1 es correcto.*

Demostración. Basta observar que el algoritmo sigue la construcción recursiva de la Sección 1.2.2, donde escogemos las variables libres t_i como lo indica la demostración del lema 2.0.2 para asegurar que \mathbf{x} sea no negativo. El único punto de aclaración lo hacemos con respecto a las redefiniciones en la línea 11.

Sea \mathbf{q}' una copia del vector \mathbf{q} antes de realizar cualquier modificación. No es difícil ver, por medio de inducción y recordando $g_1 := \text{mcd}\{q'_1, \dots, q'_n\} = 1$, que

$$q_i = \frac{q'_i}{\prod_{j=1}^{\min\{i, n-1\}} g_j},$$

para todo $i \in \{1, \dots, n\}$. Luego, las definiciones en las líneas (5), (6) y (12) son consistentes con la construcción recursiva de la Sección 1.2.2. Juntando esto con el lema 2.0.2 encontramos que \mathbf{x} es no negativo y satisface la ecuación lineal diofantina $\mathbf{q}^T \mathbf{x} = \eta$.

□

El algoritmo 2 extiende el algoritmo 1. Solamente construimos un vector $\tilde{\mathbf{q}}$ a partir del vector coprimo \mathbf{q} de manera que se satisfagan las hipótesis del algoritmo 1. Esta construcción sigue la misma lógica con la que justificamos los supuestos $q_i \neq 0$ y $q_{n-1} < 0 < q_n$ antes de presentar el lema 2.0.2.

Al igual que en el algoritmo anterior, suponemos la existencia de las subrutinas **length** y **switch**, las cuales determinan la dimensión de un vector \mathbf{q} y permutan sus entradas, respectivamente. Ambas subrutinas son estándar en la literatura y por lo tanto diremos que son correctos sin proveer alguna demostración. Así también, la subrutina **NonNegativeIntSol** es el algoritmo 1, el cual es correcto a causa del lema 2.0.3.

Teorema 2.0.4. *El algoritmo 2 es correcto.*

Demostración. Primero mostramos que el vector $\tilde{\mathbf{q}}$ satisface las hipótesis del algoritmo 1. Por definición, en la línea 3, tenemos que ninguna entrada de $\tilde{\mathbf{q}}$ es nula.

Recordemos de la definición 1.2.1 que, como \mathbf{q} es el vector coprimo de un vector esencialmente entero \mathbf{p} , su primera entrada no nula es positiva. Así, es cierto que $\tilde{q}_1 > 0$. A partir de la permutación en la línea 5 encontramos que $\tilde{q}_m > 0$.

Del ciclo en la línea 6 recuperamos un índice j tal que $\tilde{q}_j < 0$ y lo permutamos con la $(m - 1)$ -ésima entrada de $\tilde{\mathbf{q}}$ en la línea (10), de manera que obtenemos $\tilde{q}_{m-1} < 0$.

Con los tres puntos anteriores, encontramos que el vector $\tilde{\mathbf{q}}$ satisface las hipótesis del algoritmo 1 y por lo tanto el vector $\tilde{\mathbf{x}}$ es no negativo y satisface la ecuación lineal diofantina $\tilde{\mathbf{q}}^T \tilde{\mathbf{x}} = \eta$, debido al lema 2.0.3.

Las siguientes dos líneas se encargan de invertir las permutaciones hechas previamente. Finalmente, en el ciclo (14) insertamos en \mathbf{x} las entradas i de $\tilde{\mathbf{x}}$ donde $q_{\sigma_i} \neq 0$. En otro caso tenemos $x_i = 0$. Así pues, el vector \mathbf{x} es no negativo y también tenemos

$$\mathbf{q}^T \mathbf{x} = \sum_{i=1}^n q_i x_i = \sum_{i=1}^m q_{\sigma_i} x_{\sigma_i} = \sum_{i=1}^m \tilde{q}_i \tilde{x}_i = \eta,$$

por lo que concluimos que el algoritmo 2 es correcto. \square

Teorema 2.0.5. *Sea $\mathbf{p} \in \mathbb{R}^n$ un vector esencialmente entero tal que su múltiplo coprimo $\mathbf{q} \in \mathbb{Z}^n$ tiene una entrada negativa. Entonces el problema (1.1) se puede resolver a través de encontrar la solución de una ecuación lineal diofantina en n incógnitas.*

Demostración. Como \mathbf{q} es el múltiplo coprimo de \mathbf{p} , existe un escalar $m \in \mathbb{R}$ tal que $\mathbf{p} = m\mathbf{q}$. Supongamos, sin pérdida de generalidad, que m es positivo. Recuperemos η del lema 1.2.7. Por hipótesis, una entrada de \mathbf{q} es negativa, y entonces este vector satisface las condiciones del algoritmo 2. Por el teorema 2.0.4 podemos encontrar, a partir de resolver solo una ecuación lineal diofantina, un vector entero no negativo \mathbf{x} que satisface $\mathbf{q}^T \mathbf{x} = \eta$. Observemos que

$$\mathbf{p}^T \mathbf{x} = m\mathbf{q}^T \mathbf{x} = m\eta.$$

Por el corolario 2.0.1 concluimos que \mathbf{x} no solo es factible para el problema (1.1), sino que también es un punto óptimo. \square

2.1. Experimentos numéricos

Algoritmo 1: NonNegativeIntSolInf

Datos:

$\mathbf{q} \in \mathbb{Z}^n$ coprimo tal que $q_i \neq 0$ para todo $i \in \{1, \dots, n\}$ y
 $q_{n-1} < 0 < q_n$.
 $\eta \in \mathbb{Z}_{\geq 0}$.

Resultado:

$\mathbf{x} \in \mathbb{Z}_{\geq 0}^n$ tal que $\mathbf{q}^T \mathbf{x} = \eta$.

inicio

| | |
|---|----|
| $\mathbf{x} \leftarrow \mathbf{0};$ | 1 |
| $\omega_1 \leftarrow \eta;$ | 2 |
| para $i \leftarrow 1$ a $n - 2$ hacer | 3 |
| $g_{i+1} \leftarrow \text{mcd}\{q_{i+1}, \dots, q_n\};$ | 4 |
| $x'_i, \omega'_{i+1} \leftarrow \text{Bezout}(q_i, g_{i+1});$ | 5 |
| $t_i \leftarrow \lceil -\omega_i x'_i / g_{i+1} \rceil;$ | 6 |
| $x_i \leftarrow \omega_i x'_i + g_{i+1} t_i;$ | 7 |
| $\omega_{i+1} \leftarrow \omega_i \omega'_{i+1} - q_i t_i;$ | 8 |
| para $j \leftarrow i$ a $n - 1$ hacer | 9 |
| $q_{j+1} \leftarrow q_{j+1} / g_{i+1};$ | 10 |
| $x'_{n-1}, x'_n \leftarrow \text{Bezout}(q_{n-1}, q_n);$ | 11 |
| $b_1 \leftarrow -\omega_{n-1} x'_{n-1} / q_n;$ | 12 |
| $b_2 \leftarrow \omega_{n-1} x'_n / q_{n-1};$ | 13 |
| $t_{n-1} \leftarrow \lceil \text{máx}\{b_1, b_2\} \rceil;$ | 14 |
| $x_{n-1} \leftarrow \omega_{n-1} x'_{n-1} + q_n t_{n-1};$ | 15 |
| $x_n \leftarrow \omega_{n-1} x'_n - q_{n-1} t_{n-1};$ | 16 |
| devolver $\mathbf{x};$ | 17 |
| | 18 |

Algoritmo 2: Dioph

Datos: $\mathbf{q} \in \mathbb{Z}^n$ coprimo tal que $q_i < 0$ para alguna $i \in \{1, \dots, n\}$. $\eta \in \mathbb{Z}_{\geq 0}$.**Resultado:** $\mathbf{x} \in \mathbb{Z}_{\geq 0}^n$ tal que $\mathbf{q}^T \mathbf{x} = \eta$.

| | |
|--|----|
| $\mathbf{x} \leftarrow \mathbf{0};$ | 1 |
| $\sigma \leftarrow (i: q_i \neq 0);$ | 2 |
| $\tilde{\mathbf{q}} \leftarrow (q_i: q_i \neq 0);$ | 3 |
| $m \leftarrow \text{length}(\tilde{\mathbf{q}});$ | 4 |
| $\text{switch}(\tilde{\mathbf{q}}, 1, m);$ | 5 |
| para $i \leftarrow 1$ a $m - 1$ hacer | 6 |
| si $\tilde{q}_i < 0$ entonces | 7 |
| $j \leftarrow i;$ | 8 |
| ir al paso 10; | 9 |
| $\text{switch}(\tilde{\mathbf{q}}, j, m - 1);$ | 10 |
| $\tilde{\mathbf{x}} \leftarrow \text{NonNegativeIntSolInf}(\tilde{\mathbf{q}}, \eta);$ | 11 |
| $\text{switch}(\tilde{\mathbf{x}}, j, m - 1);$ | 12 |
| $\text{switch}(\tilde{\mathbf{x}}, 1, m);$ | 13 |
| para $i \leftarrow 1$ a m hacer | 14 |
| $x_{\sigma_i} \leftarrow \tilde{x}_i;$ | 15 |
| devolver \mathbf{x} | 16 |

Capítulo 3

El caso finito

Nuevamente inspirados por el teorema 1.2.9, en este capítulo analizamos el caso en el que el vector coprimo \mathbf{q} tiene entradas estrictamente positivas. De esta manera, el problema (1.1) deviene una instancia particular del famoso Problema de la Mochila:

$$\max_{\mathbf{x} \in \mathbb{Z}^n} \mathbf{u}^T \mathbf{x}, \quad (3.1a)$$

$$\text{s.a. } \mathbf{w}^T \mathbf{x} \leq c, \quad (3.1b)$$

$$\mathbf{x} \geq \mathbf{0}, \quad (3.1c)$$

donde los vectores positivos $\mathbf{u}, \mathbf{w} \in \mathbb{Z}^n$ son conocidos como vector de útiles y vector de pesos, respectivamente. Puesto que no acotamos \mathbf{x} , el problema recibe el nombre de Problema de la Mochila no Acotado. Pero también como $\mathbf{u} = \mathbf{w}$, el problema puede ser considerado como un Problema de la Suma de Conjuntos no Acotado.

En la primera sección realizamos un análisis de capas enteras a fin de obtener un resultado análogo al teorema 2.0.5. En concreto, el teorema 3.1.19 enuncia que, para un presupuesto u suficientemente grande en el problema (1.1), la búsqueda de una solución se reduce a resolver solamente una ecuación lineal diofantina.

El resultado anterior, si bien interesante, es de existencia y no muestra

cómo obtener las soluciones enteras no negativas de ecuaciones lineales diofantinas. De manera similar a como lo hicimos en el capítulo anterior, la segunda sección se encarga de presentar tal construcción de soluciones a partir de los algoritmos 3 y 4.

Finalmente, en la tercera y última sección de este capítulo, realizamos algunos experimentos numéricos que comparan la eficacia de nuestros algoritmos recién desarrollados con la de Ramificación y Acotamiento, así como de una formulación alternativa de programación dinámica.

3.1. Análisis de capas enteras

De acuerdo al segundo caso del teorema 1.2.9, el número de puntos enteros no negativos sobre la k -ésima capa entera es finito y, por lo tanto, puede ser cero. Sea $k \in \{\eta, \dots, 0\}$. Sabemos de la Sección 1.2.2 que deseamos resolver la ecuación lineal diofantina (1.14), por lo que implementamos la misma estrategia para plantear una formulación recursiva.

Debido al supuesto $\mathbf{q} > \mathbf{0}$, observemos de (1.20) que podemos agregar la condición $\omega_i \geq 0$. En efecto, buscamos que \mathbf{x} sea no negativo y recordemos que g_i es un máximo común divisor (ver (1.21)), por lo que es estrictamente positivo. Juntando esto con el supuesto $\mathbf{q} > \mathbf{0}$, encontramos que ω_i es no negativo para toda $i \in \{1, \dots, n-1\}$. Así pues, despejando t_i de (1.25) obtenemos los intervalos de factibilidad

$$\left\lceil -\frac{\omega_i x'_i}{g_{i+1}} \right\rceil \leq t_i \leq \left\lfloor \frac{\omega_i \omega'_{i+1}}{q_i} \prod_{j=1}^i g_j \right\rfloor,$$

para todo $i \in \{1, \dots, n-2\}$. Luego, como $0 < q_{n-1}, q_n$, se sigue de (1.28) que

$$\left\lceil -\frac{\omega_{n-1} x'_{n-1}}{q_n} \cdot \prod_{j=1}^{n-2} g_j \right\rceil \leq t_{n-1} \leq \left\lfloor \frac{\omega_{n-1} x'_n}{q_{n-1}} \cdot \prod_{j=1}^{n-2} g_j \right\rfloor.$$

Consecuentemente, el número de elecciones que podemos realizar para el vector de variables libres $\mathbf{t} \in \mathbb{Z}^{n-1}$ es, como lo confirma el teorema 1.2.9, finito. Si determinamos que no existe tal punto en la k -ésima capa entera, descendemos a la $(k - 1)$ -ésima capa entera y continuamos con nuestra búsqueda.

Ahora bien, en esta primera parte de la sección nos encargamos de calcular una cota superior para el número de capas enteras que debemos analizar de manera que garanticemos la existencia de un punto entero no negativo sobre una de estas capas enteras.

Lema 3.1.1. *Sea $\mathbf{p} \in \mathbb{R}^n$ un vector esencialmente entero y sea $\mathbf{q} \in \mathbb{Z}^n$ su múltiplo coprimo, por lo que existe $m \in \mathbb{R}$ tal que $\mathbf{p} = m\mathbf{q}$. Supongamos que $m > 0$ y que $\mathbf{q} > \mathbf{0}$. Sea $q^* := \max\{q_1, \dots, q_n\}$, y sea*

$$\tau := \left\lfloor \left\lfloor \frac{u}{q^*} \right\rfloor \frac{q^*}{m} \right\rfloor, \quad (3.2)$$

donde u es el lado derecho de (1.1b). Entonces la solución del problema (1.1), de ser factible, se encuentra en una capa entera parametrizada por $k \in \{\eta, \eta - 1, \dots, \tau\}$, donde recuperamos η del lema 1.2.7.

Demostración. Definamos $i^* := \arg \max\{q_1, \dots, q_n\}$ y consideremos el vector

$$\mathbf{v} := \left\lfloor \frac{u}{q^*} \right\rfloor \mathbf{e}_{i^*}.$$

Por hipótesis tenemos $q^* > 0$ y, además, como el problema (1.1) es factible, se sigue del teorema 1.2.9 que el presupuesto u es no negativo. De esto obtenemos que $\mathbf{v} \geq \mathbf{0}$. Así también,

$$\mathbf{q}^T \mathbf{v} = \left\lfloor \frac{u}{q^*} \right\rfloor q^* \leq \frac{u}{q^*} q^* = u,$$

y entonces \mathbf{v} es factible. De aquí se sigue que este vector provee una cota inferior para el problema (1.1). Así pues, todo vector \mathbf{x} candidato a ser el

óptimo del problema satisface

$$\mathbf{q}^T \mathbf{x} = \frac{\mathbf{p}^T \mathbf{x}}{m} \geq \frac{\mathbf{q}^T \mathbf{v}}{m} = \left\lfloor \frac{u}{q^*} \right\rfloor \frac{q^*}{m}.$$

Nos interesa determinar el entero τ más pequeño tal que todo punto sobre la capa entera $H_{\mathbf{q},k\|\mathbf{q}\|^{-2}}$ con $k \in \{\tau, \tau + 1, \dots\}$ satisfaga esta desigualdad. Del lema 1.2.4, encontramos que k debe satisfacer

$$\frac{k}{\|\mathbf{q}\|^2} = \frac{\mathbf{q}^T \mathbf{x}}{\|\mathbf{q}\|^2} \geq \left\lfloor \frac{u}{q^*} \right\rfloor \frac{q^*}{m \|\mathbf{q}\|^2},$$

equivalentemente,

$$k \geq \left\lfloor \frac{u}{q^*} \right\rfloor \frac{q^*}{m}.$$

Consecuentemente,

$$\tau = \left\lfloor \left\lfloor \frac{u}{q^*} \right\rfloor \frac{q^*}{m} \right\rfloor.$$

Finalmente, recordemos del lema 1.2.7 que η es la primera capa en satisfacer la restricción presupuestaria. Por lo tanto, el óptimo del problema (1.1) se encuentra en una capa entera parametrizada por $\tau \leq k \leq \eta$. \square

Observación. Siempre se cumple que $\tau \leq \eta$. En efecto,

$$\left\lfloor \frac{u}{q^*} \right\rfloor q^* \leq \frac{u}{q^*} q^* = u,$$

como $m > 0$, tenemos

$$\left\lfloor \frac{u}{q^*} \right\rfloor \frac{q^*}{m} \leq \frac{u}{m}.$$

Aplicando la función piso a ambos lados de la desigualdad y comparando con (3.2) y el lema 1.2.7 encontramos que $\tau \leq \eta$.

Observación. Nuevamente, la suposición de que el escalar m sea positivo ocurre sin pérdida de generalidad. Así como mencionamos en el capítulo anterior que si m es negativo entonces existe un parámetro η' análogo a η ,

también existe $\tau' \geq \eta'$ tal que la solución del problema (1.1) se encuentra en una capa parametrizada por $\eta' \leq k \leq \tau'$.

Lema 3.1.2. *Sean q y m enteros distintos de cero. Entonces la función $\Delta: \mathbb{R} \rightarrow \mathbb{R}$ dada por*

$$\Delta(x) := \left\lfloor \frac{x}{m} \right\rfloor - \left\lfloor \left\lfloor \frac{x}{q} \right\rfloor \frac{q}{m} \right\rfloor,$$

es periódica con periodo $\text{mcm}\{q, m\}$.

Demostración. Tenemos

$$\Delta(x + \text{mcm}\{q, m\}) = \left\lfloor \frac{x}{m} + \frac{\text{mcm}\{q, m\}}{m} \right\rfloor - \left\lfloor \left\lfloor \frac{x}{q} + \frac{\text{mcm}\{q, m\}}{q} \right\rfloor \frac{q}{m} \right\rfloor,$$

pero $q, m \mid \text{mcm}\{q, m\}$, por lo que $\text{mcm}\{q, m\}/m$ y $\text{mcm}\{q, m\}/q$ son enteros. Por las propiedades de la función piso obtenemos:

$$\begin{aligned} \Delta(x + \text{mcm}\{q, m\}) &= \left\lfloor \frac{x}{m} \right\rfloor + \frac{\text{mcm}\{q, m\}}{m} - \left\lfloor \left\lfloor \frac{x}{q} \right\rfloor \frac{q}{m} + \frac{\text{mcm}\{q, m\}}{q} \cdot \frac{q}{m} \right\rfloor \\ &= \left\lfloor \frac{x}{m} \right\rfloor + \frac{\text{mcm}\{q, m\}}{m} - \left\lfloor \left\lfloor \frac{x}{q} \right\rfloor \frac{q}{m} \right\rfloor - \frac{\text{mcm}\{q, m\}}{m} \\ &= \left\lfloor \frac{x}{m} \right\rfloor - \left\lfloor \left\lfloor \frac{x}{q} \right\rfloor \frac{q}{m} \right\rfloor \\ &= \Delta(x), \end{aligned}$$

que es lo que queríamos demostrar. \square

Definición 3.1.3. Sea $\mathbf{p} \in \mathbb{R}^n$ un vector esencialmente entero y sea $\mathbf{q} \in \mathbb{Z}^n$ su múltiplo coprimo. Consideremos los parámetros η y τ (c.f. lemas 1.2.7 y 3.1.1) como funciones del presupuesto u . Entonces decimos que la función $\Delta^*: \mathbb{R} \rightarrow \mathbb{R}$ dada por

$$\Delta^*(u) := \eta(u) - \tau(u) \tag{3.3}$$

denota el **número de capas enteras a revisar** dado el presupuesto u .

Si queremos aplicar el lema 3.1.2 a la función de la definición anterior, debemos reducir nuestra atención a vectores \mathbf{p} enteros. Esto se debe a que debemos asegurar que el múltiplo m sea entero¹. Independientemente del comportamiento periódico de Δ^* , tenemos que esta función varía significativamente ante cambios en m . Esto último implica que el número de capas enteras a revisar depende del número de cifras decimales usadas para especificar \mathbf{p} . Véase la Figura 3.1 o el Ejemplo 3.1.4.

Ejemplo 3.1.4. Si tenemos $\mathbf{p} := (9.6, 7.2, 5.6)^T$, entonces $m = 0.8$ y por lo tanto el número de capas a revisar dado $u := 119$ es $\Delta^*(u) = 14$. En cambio, si tenemos $\mathbf{p} := (9.60, 7.28, 5.68)^T$, obtenemos $m = 0.08$, por lo que el número de capas a revisar dado u es $\Delta^*(u) = 1499$. Es decir, si usamos una cifra decimal más en cada entrada, entonces $\Delta^*(u)$ se multiplica por 100, aproximadamente.

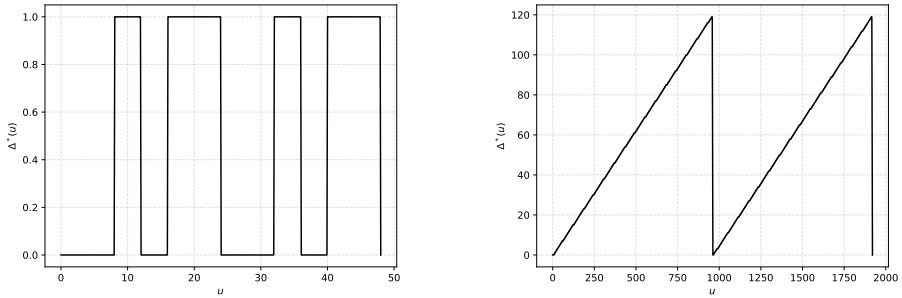


Figura 3.1.: Número de capas a revisar en función del presupuesto. *Izquierda:* Para los parámetros $m = 8$ y $q^* = 12$ encontramos que hay un máximo de una capa a revisar. *Derecha:* A medida que m se vuelve fraccionario, las capas a revisar aumentan. En este caso tenemos $m = 0.08$ y $q^* = 960$.

Observaremos en el análisis de resultados que el número de capas enteras

¹Es la creencia del autor que el lema 3.1.2 puede ser generalizado para múltiplos m racionales, mas esto no agrega demasiado valor en lo que sigue de la tesis.

que nuestro algoritmo revisa en realidad disminuye a medida que aumenta el presupuesto u .

En esta segunda parte de la sección, demostraremos que para un presupuesto u suficientemente grande, la solución del problema (1.1) se encuentra en la η -ésima capa entera donde, como siempre, η es recuperada del lema 1.2.7. Este resultado será análogo al teorema 2.0.5. No obstante, para lograr aquello, necesitamos de un par de definiciones y lemas preliminares.

Para motivar al lector, primero mostramos que existe una vecindad fija de todo punto en \mathbb{R}^n de manera que esa vecindad contiene al menos un punto entero. Esto es especificado en el teorema 3.1.6.

Luego, observamos que el “trozo” no negativo de una capa entera $H_{\mathbf{q},k\|\mathbf{q}\|^{-2}}$ crece a medida que k aumenta. Así pues, si k es lo suficientemente grande, habrá un punto sobre ese trozo no negativo cuya vecindad también se encuentra contenida en ese trozo y, por lo tanto, habrá un punto entero no negativo sobre ese trozo. Esto es especificado en el teorema 3.1.17.

Finalmente, relacionamos el punto entero que se encuentra sobre el pedazo no negativo de $H_{\mathbf{q},k\|\mathbf{q}\|^{-2}}$ con el problema (1.1). Así pues, concluimos esta sección con los teoremas 3.1.18 y 3.1.19.

Definición 3.1.5. Sea $\mathbf{q} \in \mathbb{Z}^n$ un vector coprimo y sea k un entero. Entonces definimos la **bola cerrada** sobre la k -ésima capa entera $H_{\mathbf{q},k\|\mathbf{q}\|^{-2}}$ con radio $r > 0$ y centro $\mathbf{x} \in H_{\mathbf{q},k\|\mathbf{q}\|^{-2}}$ como

$$B_r^{(k)}(\mathbf{x}) := \{\mathbf{y} \in \mathbb{R}^n : \|\mathbf{y} - \mathbf{x}\| \leq r\} \cap H_{\mathbf{q},k\|\mathbf{q}\|^{-2}}. \quad (3.4)$$

Teorema 3.1.6. Sea $\mathbf{q} \in \mathbb{Z}^n$ un vector coprimo y supongamos que $q_n \neq 0$. Sea k un entero. Entonces existe $r > 0$ tal que la familia de bolas

$$\left\{ B_r^{(k)}(\mathbf{x}) : \mathbf{x} \in H_{\mathbf{q},k\|\mathbf{q}\|^{-2}} \cap \mathbb{Z}^n \right\}$$

es una cubierta de la k -ésima capa entera $H_{\mathbf{q},k\|\mathbf{q}\|^{-2}}$.

Demostración. Como $q_n \neq 0$, recordemos del teorema (1.2.15) que

$$\mathbf{x} \in H_{\mathbf{q},k\|\mathbf{q}\|^{-2}} \cap \mathbb{Z}^n \iff \mathbf{x} = k\boldsymbol{\nu} + M\mathbf{t}$$

para algún vector $\mathbf{t} \in \mathbb{Z}^{n-1}$, donde recuperamos $\boldsymbol{\nu}$ y M de (1.37) y (1.38), respectivamente. Así, tenemos

$$\left\{ B_r^{(k)}(\mathbf{x}) : \mathbf{x} \in H_{\mathbf{q},k\|\mathbf{q}\|^{-2}} \cap \mathbb{Z}^n \right\} = \left\{ B_r^{(k)}(k\boldsymbol{\nu} + M\mathbf{t}) : \mathbf{t} \in \mathbb{Z}^{n-1} \right\}.$$

Por la definición 3.1.5 sabemos que $B_r^{(k)}(k\boldsymbol{\nu} + M\mathbf{t}) \subseteq H_{\mathbf{q},k\|\mathbf{q}\|^{-2}}$ para todo $r > 0$ y para todo $\mathbf{t} \in \mathbb{Z}^{n-1}$. Luego, para cualquier $r > 0$ tenemos

$$\bigcup_{\mathbf{t} \in \mathbb{Z}^{n-1}} B_r^{(k)}(k\boldsymbol{\nu} + M\mathbf{t}) \subseteq H_{\mathbf{q},k\|\mathbf{q}\|^{-2}}. \quad (3.5)$$

Ahora bien, sea \mathbf{y} un punto sobre la k -ésima capa entera. Por el lema 1.2.12 sabemos que las columnas de M son linealmente independientes, y entonces existe $\mathbf{t} \in \mathbb{R}^{n-1}$ tal que

$$\mathbf{y} = k\boldsymbol{\nu} + M\mathbf{t}.$$

Sea $\lfloor \mathbf{t} \rfloor \in \mathbb{Z}^{n-1}$ el vector resultante de redondear cada entrada de \mathbf{t} al entero más cercano. Luego, $\mathbf{t} = \lfloor \mathbf{t} \rfloor + \boldsymbol{\delta}$, para alguna $\boldsymbol{\delta} \in \mathbb{R}^{n-1}$ que satisface $\|\boldsymbol{\delta}\|_\infty \leq 1/2$. Definamos

$$\mathbf{x} := k\boldsymbol{\nu} + M\lfloor \mathbf{t} \rfloor \in \mathbb{Z}^{n-1},$$

de donde se sigue que

$$\begin{aligned} \|\mathbf{y} - \mathbf{x}\|^2 &= \|M\boldsymbol{\delta}\|^2 \\ &\leq \sum_{i=1}^{n-1} |\delta_i|^2 \|M\mathbf{e}_i\|^2 \\ &\leq \frac{1}{4} \sum_{i=1}^{n-1} \|M\mathbf{e}_i\|^2 \end{aligned} \quad (3.6)$$

$$= \frac{1}{4} \|M\|_F^2,$$

donde $\|M\|_F$ denota la norma Frobenius de M . Por lo tanto, si definimos

$$r := \frac{1}{2} \|M\|_F, \quad (3.7)$$

encontramos que $\mathbf{y} \in B_r^{(k)}(\mathbf{x})$. Luego, como $\mathbf{y} \in H_{\mathbf{q},k\|\mathbf{q}\|^2}$ fue genérico, tenemos que si r está definido por (3.7), entonces

$$H_{\mathbf{q},k\|\mathbf{q}\|^{-2}} \subseteq \bigcup_{\mathbf{t} \in \mathbb{Z}^{n-1}} B_r^{(k)}(k\boldsymbol{\nu} + M\mathbf{t}). \quad (3.8)$$

Juntando esto con (3.5) obtenemos lo que queríamos demostrar. \square

Lo que se encuentra a continuación se encarga de formalizar y caracterizar lo que nos referíamos anteriormente como “trozo no negativo” de la k -ésima capa entera. Todas las definiciones fueron tomadas de [BV04] a excepción del baricentro, mientras que las demostraciones de los lemas y teoremas fueron realizadas completamente por el autor.

Definición 3.1.7. Sean $\mathbf{v}_1, \dots, \mathbf{v}_m \in \mathbb{R}^n$ una colección de vectores, entonces definimos su **combinación afina** a partir de

$$\text{aff}\{\mathbf{v}_1, \dots, \mathbf{v}_m\} := \{\theta_1 \mathbf{v}_1 + \dots + \theta_k \mathbf{v}_k : \theta_1 + \dots + \theta_m = 1\}.$$

Lema 3.1.8. Sean $\mathbf{v}_1, \dots, \mathbf{v}_m \in \mathbb{R}^n$ una colección de vectores. Entonces

$$\text{aff}\{\mathbf{v}_1, \dots, \mathbf{v}_m\} = \mathbf{v}_j + \text{gen}\{\mathbf{v}_i - \mathbf{v}_j\}_{i=1}^n = \mathbf{v}_j + \text{gen}\{\mathbf{v}_i - \mathbf{v}_j\}_{i \neq j}.$$

Demostración. Puesto que $\mathbf{v}_j - \mathbf{v}_j = \mathbf{0}$, se sigue inmediatamente que

$$\mathbf{v}_j + \text{gen}\{\mathbf{v}_i - \mathbf{v}_j\}_{i=1}^n = \mathbf{v}_j + \text{gen}\{\mathbf{v}_i - \mathbf{v}_j\}_{i \neq j}.$$

Sean $\theta_1, \dots, \theta_m$ escalares tales que $\theta_1 + \dots + \theta_m = 1$. Por un lado, tenemos

$$\begin{aligned} \sum_{i=1}^m \theta_i \mathbf{v}_i &= \sum_{i=1}^m \theta_i \mathbf{v}_j + \sum_{i=1}^m \theta_i (\mathbf{v}_i - \mathbf{v}_j) \\ &= \mathbf{v}_j + \sum_{i \neq j} \theta_i (\mathbf{v}_i - \mathbf{v}_j). \end{aligned}$$

De donde se sigue que $\text{aff}\{\mathbf{v}_1, \dots, \mathbf{v}_m\} \subseteq \mathbf{v}_j + \text{gen}\{\mathbf{v}_i - \mathbf{v}_j\}_{i \neq j}$.

Ahora bien, sea $\{\lambda_i\}_{i \neq j}$ un conjunto de $m - 1$ escalares y definamos

$$\lambda_j = 1 - \sum_{i \neq j} \lambda_i.$$

Observemos que $\lambda_1 + \dots + \lambda_m = 1$ y, además,

$$\begin{aligned} \mathbf{v}_j + \sum_{i \neq j} \lambda_i (\mathbf{v}_i - \mathbf{v}_j) &= \left(1 - \sum_{i \neq j} \lambda_i\right) \mathbf{v}_j + \sum_{i \neq j} \lambda_i \mathbf{v}_i \\ &= \lambda_j \mathbf{v}_j + \sum_{i \neq j} \lambda_i \mathbf{v}_i \\ &= \sum_{i=1}^m \lambda_i \mathbf{v}_i. \end{aligned}$$

De donde se sigue que $\mathbf{v}_j + \text{gen}\{\mathbf{v}_i - \mathbf{v}_j\}_{i \neq j} \subseteq \text{aff}\{\mathbf{v}_1, \dots, \mathbf{v}_m\}$. Puesto que hemos mostrado ambas contenciones, obtenemos lo que queríamos demostrar. \square

Ejemplo 3.1.9. Si $\mathbf{q} > \mathbf{0}$ es un vector coprimo y k un entero positivo, entonces la k -ésima capa entera $H_{\mathbf{q}, k\|\mathbf{q}\|^{-2}}$ es la combinación afina de un conjunto de vectores. En efecto, recordemos de la definición 1.2.3 que esta capa entera es simplemente un hiperplano afino. Como \mathbf{q} es el vector normal a este hiperplano, se sigue que puede ser escrito como $\mathbf{v} + \ker\{\mathbf{q} \mapsto \mathbf{q}^T \mathbf{x}\}$ para alguna $\mathbf{v} \in H_{\mathbf{q}, k\|\mathbf{q}\|^{-2}}$.

Sean $\mathbf{u}_1, \dots, \mathbf{u}_n$ las intersecciones de la k -ésima capa entera con cada uno

de los ejes. Es decir, sean, para cada $i \in \{1, \dots, n\}$,

$$\mathbf{u}_i := \frac{k}{q_i} \mathbf{e}_i. \quad (3.9)$$

Como cada \mathbf{u}_i está en la k -ésima capa entera, se verifica que $\mathbf{q}^T \mathbf{u}_i = k$ y por lo tanto $\mathbf{u}_i - \mathbf{u}_j \in \ker\{\mathbf{q} \mapsto \mathbf{q}^T \mathbf{x}\}$. No es difícil ver entonces que el conjunto de vectores $\{\mathbf{u}_i - \mathbf{u}_j\}_{i \neq j}$ forma una base del espacio nulo de la transformación lineal $\mathbf{q} \mapsto \mathbf{q}^T \mathbf{x}$, por lo que obtenemos

$$H_{\mathbf{q}, k\|\mathbf{q}\|^{-2}} = \mathbf{u}_j + \text{gen}\{\mathbf{u}_i - \mathbf{u}_j\}_{i \neq j}.$$

Por el lema 3.1.8 concluimos que la k -ésima capa entera es la combinación afina de los vectores $\mathbf{u}_1, \dots, \mathbf{u}_n$.

Definición 3.1.10. Sean $\mathbf{v}_1, \dots, \mathbf{v}_m \in \mathbb{R}^n$ vectores linealmente independientes. Entonces definimos el **símplice** σ como la combinación convexa de estos vectores:

$$\sigma = \text{conv}\{\mathbf{v}_1, \dots, \mathbf{v}_m\} := \{\theta_1 \mathbf{v}_1 + \dots + \theta_m \mathbf{v}_m : \theta_1 + \dots + \theta_m = 1, \theta_i \geq 0\}.$$

Decimos entonces que σ es generado por $\mathbf{v}_1, \dots, \mathbf{v}_m$. También definimos la **j -ésima faceta** σ_j de σ como el símplice generado por los vectores $\{\mathbf{v}_i\}_{i \neq j}$.

Observación. Comparando con la definición 3.1.7, encontramos que todo símplice σ generado por $\mathbf{v}_1, \dots, \mathbf{v}_m$ está contenido en la combinación afina de estos vectores. Es decir,

$$\text{conv}\{\mathbf{v}_1, \dots, \mathbf{v}_m\} \subseteq \text{aff}\{\mathbf{v}_1, \dots, \mathbf{v}_m\}. \quad (3.10)$$

Observación. Si σ es un símplice generado por m vectores, entonces tiene $\binom{m}{m-1} = m$ facetas. Tomaremos por hecho, puesto que de otra manera arriesgamos desviarnos por una tangente, que estas facetas constituyen la frontera relativa del símplice. Es decir, tomaremos por hecho que las facetas

constituyen las “aristas” o “caras” de σ .

Lema 3.1.11. *Sea $\mathbf{q} > \mathbf{0}$ un vector coprimo y sea $H_{\mathbf{q},k\|\mathbf{q}\|^{-2}}$ la k -ésima capa entera, con parámetro k positivo. Consideremos el símplice σ generado por los vectores definidos en (3.9), entonces*

$$\sigma = H_{\mathbf{q},k\|\mathbf{q}\|^{-2}} \cap \mathbb{R}_{\geq \mathbf{0}}^n.$$

Demostración. En el Ejemplo 3.1.9 mostramos que

$$H_{\mathbf{q},k\|\mathbf{q}\|^{-2}} = \text{aff}\{\mathbf{u}_1, \dots, \mathbf{u}_n\}.$$

Sea $\mathbf{x} \in \sigma$, de la definición 3.1.10 y de (3.10) encontramos que \mathbf{x} se encuentra en la k -ésima capa entera. Además, existen escalares $\theta_1, \dots, \theta_n$ no negativos tales que

$$\mathbf{x} = \theta_1 \mathbf{u}_1 + \dots + \theta_n \mathbf{u}_n = k \begin{pmatrix} \theta_1/q_1 \\ \vdots \\ \theta_n/q_n \end{pmatrix}.$$

Como $\mathbf{q} > \mathbf{0}$ y $k > 0$ por hipótesis, tenemos que $\mathbf{x} \geq \mathbf{0}$, lo que implica que $\mathbf{x} \in H_{\mathbf{q},k\|\mathbf{q}\|^{-2}} \cap \mathbb{R}_{\geq \mathbf{0}}^n$.

El otro lado de la contención se muestra de manera completamente análoga. □

En el contexto del problema (1.1), sabemos que si σ es generado por los vectores en (3.9) entonces, por el lema anterior, todo punto entero sobre σ es un punto factible siempre que $0 < k \leq \eta$, donde recuperamos η del lema 1.2.7. Nos gustaría entonces garantizar la existencia de tal punto entero.

Adoptamos la siguiente estrategia: nos concentramos en un punto $\mathbf{x} \in \sigma$ y abrimos una bola (ver definición 3.1.5) con radio dado por (3.7). Si esa bola está contenida en el símplice σ , entonces el teorema 3.1.6 garantiza la existencia de un punto entero sobre σ . Por el lema anterior, garantizaríamos

la existencia de un punto entero no negativo sobre la k -ésima capa entera.

Lo que se encuentra a continuación es un análisis para determinar qué tan grande debe ser k para asegurar que la bola de radio (3.7) esté contenida en el símple σ , dado que la bola está centrada en un punto particular, a saber, en el baricentro del símple.

Definición 3.1.12. Sea σ un símple generado por $\mathbf{v}_1, \dots, \mathbf{v}_m \in \mathbb{R}^n$, definimos su **baricentro** $\hat{\sigma}$ como

$$\hat{\sigma} := \frac{1}{m} \sum_{i=1}^m \mathbf{v}_i.$$

Observación. El baricentro $\hat{\sigma}$ es un elemento de σ . Esto se debe a que $\hat{\sigma}$ es la combinación convexa de $\mathbf{v}_1, \dots, \mathbf{v}_m$, donde $\theta_1 = \dots = \theta_m = \frac{1}{m}$.

Definición 3.1.13. Sea σ un símple y sea $\hat{\sigma}$ su baricentro. Entonces definimos el **radio de la bola inscrita** en σ con centro $\hat{\sigma}$ como

$$r_\sigma := \max\{r > 0: B_r^{(k)}(\hat{\sigma}) \subseteq \sigma\}, \quad (3.11)$$

donde $B_r^{(k)}(\hat{\sigma})$ está dada en la definición 3.1.5.

Encontraremos que el radio de la bola inscrita está dado por el mínimo de las distancias entre el baricentro del símple con cada una de sus facetas. Puesto que $\hat{\sigma}_j \in \sigma_j$, sabemos bien por álgebra lineal, bien por optimización, que la distancia entre σ y su j -ésima faceta σ_j es

$$d(\hat{\sigma}, \sigma_j) = |\hat{\mu}_j^T(\hat{\sigma} - \hat{\sigma}_j)|, \quad (3.12)$$

donde $\hat{\mu}_j$ es un vector unitario y normal a la j -ésima faceta.

Lema 3.1.14. Sean $\mathbf{q} > \mathbf{0}$, $k > 0$ y retomemos el símple σ generado por los vectores $\{\mathbf{u}_i\}_{i=1}^n$ en (3.9). Definamos, para cada $j \in \{1, \dots, n\}$,

$$\mu_j := \mathbf{u}_j - \frac{\mathbf{q}^T \mathbf{u}_j}{\mathbf{q}^T \mathbf{q}} \mathbf{q} = \mathbf{u}_j - \frac{k}{\|\mathbf{q}\|^2} \mathbf{q}. \quad (3.13)$$

Entonces $\boldsymbol{\mu}_j$ es un vector normal a la j -ésima faceta σ_j del símplice σ .

Demostración. Debemos mostrar que si $\mathbf{x} \in \sigma_j$, entonces $\boldsymbol{\mu}_j^T \mathbf{y} = 0$ para todo $\mathbf{y} \in \sigma_j - \mathbf{x}$. Por la definición 3.1.10, tenemos que los vectores $\{\mathbf{u}_i\}_{i \neq j}$ generan la j -ésima faceta σ_j , y entonces basta mostrar que $\boldsymbol{\mu}_j^T \mathbf{y} = 0$ para todo $\mathbf{y} \in \sigma_j - \mathbf{u}_m$ con $m \neq j$.

Sea, pues, $m \in \{1, \dots, n\} \setminus \{j\}$. Tenemos de las definiciones 3.1.7 y 3.1.10, así como del lema 3.1.8 que

$$\sigma_j = \text{conv}\{\mathbf{u}_i\}_{i \neq j} \subseteq \text{aff}\{\mathbf{u}_i\}_{i \neq j} = \mathbf{u}_m + \text{gen}\{\mathbf{u}_i - \mathbf{u}_m\}_{i \neq j}.$$

De donde obtenemos

$$\sigma_j - \mathbf{u}_m \subseteq \text{gen}\{\mathbf{u}_i - \mathbf{u}_m\}_{i \neq j},$$

así que basta mostrar que $\boldsymbol{\mu}_j^T(\mathbf{u}_i - \mathbf{u}_m) = 0$ para todo $i \neq j$. Cabe mencionar que los vectores $\{\mathbf{u}_i\}_{i=1}^n$ son ortogonales entre sí (ver (3.9)). Sustituyendo con la definición de $\boldsymbol{\mu}_j$ en la hipótesis, obtenemos

$$\begin{aligned} \boldsymbol{\mu}_j^T(\mathbf{u}_i - \mathbf{u}_m) &= \mathbf{u}_j^T \mathbf{u}_i - \mathbf{u}_j^T \mathbf{u}_m - \frac{k}{\|\mathbf{q}\|^2}(\mathbf{q}^T \mathbf{u}_i - \mathbf{q}^T \mathbf{u}_m) \\ &= 0 - 0 - \frac{k}{\|\mathbf{q}\|^2}(k - k) \\ &= 0. \end{aligned}$$

De esta manera, concluimos que $\boldsymbol{\mu}_j$ es un vector normal a σ_j . □

Ahora que encontramos vectores normales $\boldsymbol{\mu}_j$ para cada faceta σ_j , podemos simplificar un poco más (3.12). Aprovechando el hecho de que $\{\mathbf{u}_i\}_{i=1}^n$ son todos ortogonales entre sí, obtenemos cálculos simples:

$$\boldsymbol{\mu}_j^T \hat{\boldsymbol{\sigma}} = \left(\mathbf{u}_j - \frac{k}{\|\mathbf{q}\|^2} \mathbf{q} \right)^T \frac{1}{n} \sum_{i=1}^n \mathbf{u}_i$$

$$\begin{aligned}
&= \frac{1}{n} \sum_{i=1}^n \mathbf{u}_j^T \mathbf{u}_i - \frac{k}{n \|\mathbf{q}\|^2} \sum_{i=1}^n \mathbf{q}^T \mathbf{u}_i \\
&= \frac{1}{n} \|\mathbf{u}_j\|^2 - \frac{k}{n \|\mathbf{q}\|^2} \sum_{i=1}^n k \\
&= \frac{k^2}{n q_j^2} - \frac{k^2}{\|\mathbf{q}\|^2}.
\end{aligned}$$

A través de un procedimiento similar, encontramos también que

$$\boldsymbol{\mu}_j^T \hat{\boldsymbol{\sigma}}_j = -\frac{k}{\|\mathbf{q}\|^2}, \quad (3.14)$$

y por lo tanto

$$\boldsymbol{\mu}_j^T (\hat{\boldsymbol{\sigma}} - \hat{\boldsymbol{\sigma}}_j) = \frac{k^2}{n q_j^2}. \quad (3.15)$$

Más adelante normalizaremos $\boldsymbol{\mu}$ de manera que este vector sea unitario. Cabe resaltar el hecho de que el lado derecho (3.15) es positivo. Geométricamente, lo anterior implica que los vectores normales $\boldsymbol{\mu}_j$ de cada faceta σ_j apuntan hacia el interior relativo del símple σ . Esto sugiere una caracterización alternativa de σ que nos permite interpretarlo como un poliedro.

Lema 3.1.15. *Sea $\mathbf{q} > \mathbf{0}$ un vector coprimo y sea σ el símple generado por los vectores definidos en (3.9), con $k > 0$. Entonces*

$$\sigma = \bigcap_{j=1}^n \{\mathbf{x} \in \mathbb{R}^n : \hat{\boldsymbol{\mu}}_j^T (\mathbf{x} - \hat{\boldsymbol{\sigma}}_j) \geq 0\} \cap H_{\mathbf{q}, k \|\mathbf{q}\|^{-2}}, \quad (3.16)$$

donde $\hat{\boldsymbol{\mu}}_j$ es el vector $\boldsymbol{\mu}_j$ definido en (3.13) normalizado.

Demostración. Denotemos por $\{\mathbf{u}_i\}_{i=1}^n$ los vectores ortogonales definidos en (3.9). Como $\hat{\boldsymbol{\mu}}_j$ es el vector $\boldsymbol{\mu}_j$ normalizado, se sigue que

$$\{\mathbf{x} \in \mathbb{R}^n : \hat{\boldsymbol{\mu}}_j^T (\mathbf{x} - \hat{\boldsymbol{\sigma}}_j) \geq 0\} = \{\mathbf{x} \in \mathbb{R}^n : \boldsymbol{\mu}_j^T (\mathbf{x} - \hat{\boldsymbol{\sigma}}_j) \geq 0\},$$

y entonces podemos trabajar con $\boldsymbol{\mu}_j$ sin normalizarlo.

Sea $\mathbf{x} \in \sigma$. Por el lema 3.1.11 sabemos que \mathbf{x} se encuentra en la k -ésima capa entera. También sabemos que existen escalares no negativos $\theta_1, \dots, \theta_n$ que suman 1 y que satisfacen $\mathbf{x} = \theta_1 \mathbf{u}_1 + \dots + \theta_n \mathbf{u}_n$. Tenemos entonces

$$\begin{aligned}
 \boldsymbol{\mu}_j^T \mathbf{x} &= \left(\mathbf{u}_j - \frac{k}{\|\mathbf{q}\|^2} \mathbf{q} \right)^T \left(\theta_j \mathbf{u}_j + \sum_{i \neq j} \theta_i \mathbf{u}_i \right) \\
 &= \theta_j \|\mathbf{u}_j\|^2 - \frac{k}{\|\mathbf{q}\|^2} \sum_{i \neq j} \theta_i \mathbf{q}^T \mathbf{u}_i \\
 &= \theta_j \frac{k^2}{q_j^2} - \frac{k}{\|\mathbf{q}\|^2} \sum_{i \neq j} k \theta_i \\
 &= \theta_j \frac{k^2}{q_j^2} - \frac{k^2}{\|\mathbf{q}\|^2} (1 - \theta_j) \\
 &= \theta_j \left(\frac{k^2}{q_j^2} + \frac{k^2}{\|\mathbf{q}\|^2} \right) - \frac{k^2}{\|\mathbf{q}\|^2}.
 \end{aligned}$$

Retomamos de (3.14) el valor de $\boldsymbol{\mu}_j^T \hat{\boldsymbol{\sigma}}_j$, así que obtenemos

$$\boldsymbol{\mu}_j^T (\mathbf{x} - \hat{\boldsymbol{\sigma}}_j) = \boldsymbol{\mu}_j^T \mathbf{x} - \boldsymbol{\mu}_j^T \hat{\boldsymbol{\sigma}}_j = \theta_j \left(\frac{k^2}{q_j^2} + \frac{k^2}{\|\mathbf{q}\|^2} \right),$$

lo cual es no negativo para todo $j \in \{1, \dots, n\}$.

Mostramos la otra contención por contrapositiva, así que supongamos que $\mathbf{x} \notin \sigma$. Por el lema 3.1.11 se sigue o bien que $\mathbf{x} \notin H_{\mathbf{q}, k\|\mathbf{q}\|^{-2}}$ o bien que $\mathbf{x} \notin \mathbb{R}_{\geq 0}^n$. En el primer caso obtenemos inmediatamente que \mathbf{x} no se encuentra en el lado derecho de (3.16).

Supongamos, pues, que \mathbf{x} está en la k -ésima capa entera pero que tiene al menos una entrada negativa con respecto a la base canónica. Como $\{\mathbf{u}_i\}_{i=1}^n$ es base de \mathbb{R}^n , existen escalares $\{\lambda_i\}_{i=1}^n$ tales que

$$\mathbf{x} = \sum_{i=1}^n \lambda_i \mathbf{u}_i.$$

Como las entradas de $\mathbf{u}_1, \dots, \mathbf{u}_n$ son todas no negativas y $x_j < 0$ para alguna $j \in \{1, \dots, n\}$, se sigue que $\lambda_j < 0$. Observemos que

$$\begin{aligned} \boldsymbol{\mu}_j^T \mathbf{x} &= \sum_{i=1}^n \lambda_i \boldsymbol{\mu}_j^T \mathbf{u}_i \\ &= \sum_{i=1}^n \lambda_i \left(\mathbf{u}_j - \frac{k}{\|\mathbf{q}\|^2} \mathbf{q} \right)^T \mathbf{u}_i \\ &= \sum_{i=1}^n \lambda_i \left(\mathbf{u}_j^T \mathbf{u}_i - \frac{k}{\|\mathbf{q}\|^2} \mathbf{q}^T \mathbf{u}_i \right) \\ &= \lambda_j \|\mathbf{u}_j\|^2 - \frac{k^2}{\|\mathbf{q}\|^2} \sum_{i=1}^n \lambda_i. \end{aligned}$$

Pero $\mathbf{x} \in H_{\mathbf{q}, k\|\mathbf{q}\|^{-2}} = \text{aff}\{\mathbf{u}_1, \dots, \mathbf{u}_n\}$ (ver Ejemplo 3.1.9) y entonces los escalares $\lambda_1, \dots, \lambda_n$ suman a 1. Sustituyendo,

$$\boldsymbol{\mu}_j^T \mathbf{x} = \lambda_j \frac{k^2}{q_j^2} - \frac{k^2}{\|\mathbf{q}\|^2},$$

retomando el valor de $\boldsymbol{\mu}_j^T \hat{\boldsymbol{\sigma}}_j$ en (3.14), encontramos que

$$\boldsymbol{\mu}_j^T (\mathbf{x} - \hat{\boldsymbol{\sigma}}_j) = \lambda_j \frac{k^2}{q_j^2} < 0$$

y entonces \mathbf{x} no es elemento del semi-espacio $\{\mathbf{x}: \boldsymbol{\mu}_j^T (\mathbf{x} - \hat{\boldsymbol{\sigma}}_j) \geq 0\}$, por lo que tampoco es elemento del lado derecho de (3.16). \square

Teorema 3.1.16. *Sea $\mathbf{q} > \mathbf{0}$ un vector coprimo y sea σ el s mplice generado por los vectores $\{\mathbf{u}_i\}_{i=1}^n$ definidos en (3.9). Entonces el radio r_σ de la bola inscrita (ver definici n 3.1.13) en σ con centro $\hat{\boldsymbol{\sigma}}$ est  dado por*

$$r_\sigma = \min_{1 \leq j \leq n} d(\hat{\boldsymbol{\sigma}}, \sigma_j) = \min_{1 \leq j \leq n} \hat{\boldsymbol{\mu}}_j^T (\hat{\boldsymbol{\sigma}} - \hat{\boldsymbol{\sigma}}_j),$$

donde $\hat{\boldsymbol{\mu}}_j$ es el vector $\boldsymbol{\mu}_j$ definido en (3.13) normalizado.

Demostración. Como $\hat{\sigma} \in \sigma$, tenemos del lema 3.1.15 que $\mu_j^T(\hat{\sigma} - \hat{\sigma}_j) \geq 0$ y, por lo tanto, deducimos de (3.12) que la distancia entre $\hat{\sigma}$ y la j -ésima faceta σ_j es

$$d(\hat{\sigma}, \sigma_j) = \hat{\mu}_j^T(\hat{\sigma} - \hat{\sigma}_j). \quad (3.17)$$

Supongamos que $r \leq d(\hat{\sigma}, \sigma_j)$ para todo $j \in \{1, \dots, n\}$ y sea $\mathbf{x} \in B_r^{(k)}(\hat{\sigma})$. Observemos que

$$\begin{aligned} \hat{\mu}_j^T(\mathbf{x} - \hat{\sigma}_j) &= \hat{\mu}_j^T(\mathbf{x} - \hat{\sigma}) + \hat{\mu}_j^T(\hat{\sigma} - \hat{\sigma}_j) \\ &= \hat{\mu}_j^T(\mathbf{x} - \hat{\sigma}) + d(\hat{\sigma}, \sigma_j). \end{aligned}$$

Por la desigualdad de Cauchy-Schwartz, tenemos

$$\hat{\mu}_j^T(\mathbf{x} - \hat{\sigma}) \geq -\|\hat{\mu}_j\| \|\mathbf{x} - \hat{\sigma}\| \geq -r,$$

pues $\hat{\mu}$ es unitario y $\mathbf{x} \in B_r^{(k)}(\hat{\sigma})$. Así pues, tenemos

$$\hat{\mu}_j^T(\mathbf{x} - \hat{\sigma}_j) \geq -r + d(\hat{\sigma}, \sigma_j) \geq 0,$$

pues supusimos que $r \leq d(\hat{\sigma}, \sigma_j)$ para todo $j \in \{1, \dots, n\}$. Además, como $\mathbf{x} \in B_r^{(k)}(\hat{\sigma})$, por la definición 3.1.5 tenemos que \mathbf{x} se encuentra en la k -ésima capa entera. Así pues,

$$\mathbf{x} \in \bigcap_{j=1}^n \{\mathbf{x} \in \mathbb{R}^n : \hat{\mu}_j^T(\mathbf{x} - \hat{\sigma}_j) \geq 0\} \cap H_{\mathbf{q}, k\|\mathbf{q}\|^{-2}} = \sigma,$$

donde la última igualdad se sigue del lema 3.1.15. Así pues, $B_r^{(k)}(\hat{\sigma}) \subseteq \sigma$ si $r \leq d(\hat{\sigma}, \sigma_j)$ para toda $j \in \{1, \dots, n\}$. De la definición 3.1.13 encontramos entonces que el radio r_σ de la bola inscrita satisface

$$r_\sigma \geq \min_{1 \leq j \leq n} d(\hat{\sigma}, \sigma_j). \quad (3.18)$$

Ahora bien, supongamos que $r > d(\hat{\sigma}, \sigma_j)$ para alguna $j \in \{1, \dots, n\}$.

Consideremos el punto $\mathbf{x} \in \sigma_j$ que satisface $d(\hat{\boldsymbol{\sigma}}, \sigma_j) = d(\hat{\boldsymbol{\sigma}}, \mathbf{x})$. Tal punto existe porque σ_j es cerrado. Luego, $\|\mathbf{x} - \hat{\boldsymbol{\sigma}}\| < r$. Entonces existe $\varepsilon > 0$ tal que

$$\|(\mathbf{x} - \varepsilon \hat{\boldsymbol{\mu}}_j) - \hat{\boldsymbol{\sigma}}\| \leq r,$$

lo que implica que $\mathbf{x} - \varepsilon \hat{\boldsymbol{\mu}}_j \in B_r^{(k)}(\hat{\boldsymbol{\sigma}})$. Observemos que

$$\hat{\boldsymbol{\mu}}_j^T((\mathbf{x} - \varepsilon \hat{\boldsymbol{\mu}}_j) - \hat{\boldsymbol{\sigma}}_j) = \hat{\boldsymbol{\mu}}_j^T(\mathbf{x} - \hat{\boldsymbol{\sigma}}_j) - \varepsilon.$$

Pero $\mathbf{x}, \hat{\boldsymbol{\sigma}}_j \in \sigma_j$, así que $\mathbf{x} - \hat{\boldsymbol{\sigma}}_j \in \sigma_j - \hat{\boldsymbol{\sigma}}_j$. Del lema 3.1.14 encontramos que

$$\hat{\boldsymbol{\mu}}_j^T(\mathbf{x} - \hat{\boldsymbol{\sigma}}_j) = 0,$$

de donde obtenemos

$$\hat{\boldsymbol{\mu}}_j^T((\mathbf{x} - \varepsilon \hat{\boldsymbol{\mu}}_j) - \hat{\boldsymbol{\sigma}}_j) = -\varepsilon < 0,$$

lo cual implica que $\mathbf{x} - \varepsilon \hat{\boldsymbol{\mu}}_j$ no se encuentra en el semi-espacio definido por $\{\mathbf{x}: \hat{\boldsymbol{\mu}}^T(\mathbf{x} - \hat{\boldsymbol{\sigma}}_j) \geq 0\}$. Así pues, por el lema 3.1.15, encontramos que $\mathbf{x} - \varepsilon \hat{\boldsymbol{\mu}}_j \notin \sigma$. Pero $\mathbf{x} - \varepsilon \hat{\boldsymbol{\mu}}_j \in B_r^{(k)}(\hat{\boldsymbol{\sigma}})$. De aquí se desprende que $B_r^{(k)}(\hat{\boldsymbol{\sigma}}) \not\subseteq \sigma$ si $r > d(\hat{\boldsymbol{\sigma}}, \sigma_j)$ para alguna $j \in \{1, \dots, n\}$. De la definición 3.1.13 obtenemos entonces

$$r_\sigma \leq \min_{1 \leq j \leq n} d(\hat{\boldsymbol{\sigma}}, \sigma_j). \quad (3.19)$$

De (3.18) y de (3.19) concluimos entonces con lo que queríamos demostrar. \square

De (3.12) tenemos

$$d(\hat{\boldsymbol{\sigma}}, \sigma_j) = \hat{\boldsymbol{\mu}}_j^T(\hat{\boldsymbol{\sigma}} - \hat{\boldsymbol{\sigma}}_j) = \frac{\boldsymbol{\mu}_j^T(\hat{\boldsymbol{\sigma}} - \hat{\boldsymbol{\sigma}}_j)}{\|\boldsymbol{\mu}_j\|}. \quad (3.20)$$

Recordemos de (3.15) que ya contamos con el numerador, así que ahora

debemos calcular la norma de $\boldsymbol{\mu}_j$. Tenemos

$$\begin{aligned}
\|\boldsymbol{\mu}_j\|^2 &= \boldsymbol{\mu}_j^T \boldsymbol{\mu}_j \\
&= \left(\mathbf{u}_j - \frac{k}{\|\mathbf{q}\|^2} \mathbf{q} \right)^T \left(\mathbf{u}_j - \frac{k}{\|\mathbf{q}\|^2} \mathbf{q} \right) \\
&= \|\mathbf{u}_j\|^2 - 2 \frac{k}{\|\mathbf{q}\|^2} \mathbf{q}^T \mathbf{u}_j + \frac{k^2}{\|\mathbf{q}\|^4} \mathbf{q}^T \mathbf{q} \\
&= \frac{k^2}{q_j^2} - 2 \frac{k^2}{\|\mathbf{q}\|^2} + \frac{k^2}{\|\mathbf{q}\|^2} \\
&= \frac{k^2}{q_j^2} - \frac{k^2}{\|\mathbf{q}\|^2}.
\end{aligned}$$

De donde obtenemos

$$\|\boldsymbol{\mu}_j\| = k \sqrt{\frac{1}{q_j^2} - \frac{1}{\|\mathbf{q}\|^2}}. \quad (3.21)$$

Usando (3.15) y (3.21) para sustituir en (3.20), obtenemos

$$d(\hat{\boldsymbol{\sigma}}, \sigma_j) = \frac{k}{n} \cdot \frac{1}{q_j^2 \sqrt{q_j^{-2} - \|\mathbf{q}\|^{-2}}} = \frac{k}{n} \cdot \frac{1}{Q_j},$$

donde definimos Q_j pertinentemente. Finalmente, del teorema 3.1.16 encontramos que el radio r_σ de la bola inscrita en el s mplice σ con centro $\hat{\boldsymbol{\sigma}}$ est  dado por

$$r_\sigma = \min_{1 \leq j \leq n} \{d(\hat{\boldsymbol{\sigma}}, \sigma_j)\} = \frac{k}{n} \cdot \frac{1}{\max_{1 \leq j \leq n} \{Q_j\}} \quad (3.22)$$

Teorema 3.1.17. *Sea $\mathbf{q} > \mathbf{0}$ un vector coprimo y sea k un entero positivo suficientemente grande. Entonces existe un punto entero sobre el s mplice σ generado por los vectores en (3.9).*

Demostraci n. Sea r el radio definido en (3.7) y sea r_σ el radio definido en (3.22). Por el teorema 3.1.6 sabemos que existe un punto entero \mathbf{x} en $B_r^{(k)}(\hat{\boldsymbol{\sigma}})$, y por el teorema 3.1.16 sabemos que la bola $B_{r_\sigma}^{(k)}(\hat{\boldsymbol{\sigma}})$ est  contenida

en σ . Entonces basta mostrar que existe k suficientemente grande tal que $r \leq r_\sigma$, pues esto implicaría la contención de σ en medio en la cadena

$$\mathbf{x} \in B_r^{(k)}(\hat{\boldsymbol{\sigma}}) \subseteq B_{r_\sigma}^{(k)}(\hat{\boldsymbol{\sigma}}) \subseteq \sigma.$$

De (3.7) y de (3.22) obtenemos que $r \leq r_\sigma$ si y solo si

$$k \geq \frac{n}{2} \|M\|_F \max_{1 \leq j \leq n} \{Q_j\}, \quad (3.23)$$

que es lo que queríamos demostrar. \square

De (3.23) parece que podemos concluir que hay una dependencia lineal entre la dimensión n y el parámetro de la capa entera k . No obstante, la norma $\|M\|_F$ depende implícitamente de n . Para ser más explícitos con respecto a esta dependencia, podemos rescatar de (3.6) la siguiente cota:

$$\frac{1}{4} \sum_{j=1}^{n-1} \|M\mathbf{e}_j\|^2 \leq \frac{n-1}{4} \max_{1 \leq j \leq n} \{\|M\mathbf{e}_j\|^2\},$$

de donde reemplazaríamos la cota (3.23) en el teorema 3.1.17 por

$$k \geq \frac{n\sqrt{n-1}}{2} \max_{1 \leq j \leq n} \{\|M\mathbf{e}_j\|\} \cdot \max_{1 \leq j \leq n} \{Q_j\}.$$

Esta cota, no obstante, es más grande que la propuesta inicialmente.

Además, el resultado que obtuvimos es más fuerte de lo que aparenta. Hemos encontrado una cota inferior de manera que podamos asegurar la existencia de puntos enteros en una vecindad del baricentro $\hat{\boldsymbol{\sigma}}$. Este punto no es especial, pues en realidad podemos realizar el mismo procedimiento enfocándonos en otros puntos del símplex σ para asegurar soluciones en sus respectivas vecindades. Entonces, dependiendo del punto, podemos obtener mejores o peores cotas para k . El punto más interesante es aquel que provee

la cota inferior más pequeña².

De manera inmediata obtenemos también los siguientes teoremas. Cabe mencionar que estos resultados solamente muestran la existencia de una solución entera \mathbf{x} no negativa para la ecuación lineal diofantina $\mathbf{q}^T \mathbf{x} = k$. Será en la Sección 3.2 que discutiremos cómo encontrar esta solución.

Teorema 3.1.18. *Sea $\mathbf{q} > \mathbf{0}$ un vector coprimo. Entonces la ecuación lineal diofantina $\mathbf{q}^T \mathbf{x} = k$ tiene soluciones enteras no negativas para k suficientemente grande.*

Demostración. Consideremos el símlice σ generado por los vectores en (3.9) y supongamos que k satisface la cota en (3.23). Por el teorema 3.1.17 existe un punto entero no negativo $\mathbf{x} \in \sigma$, y esto implica que $x \in H_{\mathbf{q}, k \|\mathbf{q}\|^{-2}}$ por el lema 3.1.11. Luego, por el lema 1.2.6, \mathbf{x} satisface la ecuación lineal diofantina $\mathbf{q}^T \mathbf{x} = k$. \square

Teorema 3.1.19. *Sea $\mathbf{p} \in \mathbb{R}^n$ un vector esencialmente entero y supongamos que su múltiplo coprimo \mathbf{q} tiene entradas estrictamente positivas. Entonces el problema (1.1) se puede resolver a través de encontrar la solución de una sola ecuación lineal en n incógnitas para un presupuesto u suficientemente grande.*

Demostración. Por la definición 1.2.1 sabemos que existe un escalar m tal que $\mathbf{p} = m\mathbf{q}$. Supongamos, sin pérdida de generalidad, que m es positivo. Del lema 1.2.7 tenemos que el entero η parametriza la primera capa entera en satisfacer el presupuesto y que $\eta = \lfloor u/m \rfloor$. Por el teorema 3.1.18 sabemos que si η es suficientemente grande, entonces la ecuación lineal diofantina $\mathbf{q}^T \mathbf{x} = \eta$ tiene al menos una solución entera no negativa \mathbf{x} . Luego, \mathbf{x} es factible para el problema (1.1), pero por la maximalidad de η encontramos que

²Una hipótesis del autor es que el baricentro $\hat{\sigma}$ provee, en efecto, la mejor cota.

\mathbf{x} también es un punto óptimo. En conclusión, solo deviene necesario resolver una ecuación lineal diofantina para determinar la solución del problema (1.1). \square

El teorema 3.1.18 junto con la cota (3.23) provee, hasta donde llega el conocimiento del autor, nuevas cotas superiores para los números de Frobenius³. De manera resumida, dada una colección de enteros a_1, \dots, a_n coprimos, el número de Frobenius es el entero F más grande tal que F no pueda ser expresado como una combinación lineal entera no negativa de a_1, \dots, a_n . Un estudio sobre cómo se compara esta colección de cotas con respecto a la literatura existente, si bien interesante, queda fuera del propósito de esta tesis.

En último lugar, mencionamos que eventualmente es suficiente con revisar la primera capa entera. No hemos demostrado, empero, que el número de capas enteras a revisar eventualmente decrece en cuanto el presupuesto u aumenta. Observaremos en el análisis de resultados que hay un patrón periódico y decreciente en cuanto al número de capas enteras revisadas. Demostrar, en cambio, que este comportamiento siempre se cumple es mucho más difícil y queda fuera del propósito de esta tesis.

3.2. Construcción de soluciones

Sea $\mathbf{p} \in \mathbb{R}^n$ un vector esencialmente entero y supongamos que las entradas de su múltiplo coprimo \mathbf{q} son todas estrictamente positivas. Supongamos, sin pérdida de generalidad, que el escalar m que satisface $\mathbf{p} = m\mathbf{q}$ es también positivo. Bastante hemos discutido sobre cómo la solución del problema (1.1) se traduce a la búsqueda de una solución entera no negativa de la ecuación lineal diofantina $\mathbf{q}^T \mathbf{x} = k$ para alguna $k \leq \eta$, donde η es tomada del lema 1.2.7.

³Véase el Problema de la Moneda en https://en.wikipedia.org/wiki/Coin_problem.

En esta sección presentamos los algoritmos 3 y 4, los cuales se encargan de obtener estas soluciones enteras no negativas que tanto buscamos. Consecuentemente, estos algoritmos se encargan de resolver el problema (1.1).

Teorema 3.2.1. *El algoritmo 3 es correcto.*

Demostración. Hacemos la demostración por inducción en la dimensión n del vector \mathbf{q} . Supongamos, para el caso base, que $n = 2$. Luego, queremos encontrar soluciones enteras no negativas de la ecuación

$$q_1x_1 + q_2x_2 = k. \quad (3.24)$$

Por hipótesis sabemos que q_1 y q_2 son coprimos. Luego, del teorema 1.1.10 encontramos que las soluciones enteras de esta ecuación están dadas por

$$\begin{cases} x_1 = kx'_1 + q_2t, \\ x_2 = kx'_2 - q_1t, \end{cases} \quad (3.25)$$

donde $t \in \mathbb{Z}$ es una variable libre, y x'_1, x'_2 son los coeficientes de Bézout (c.f. definición 1.1.9) de q_1 y q_2 , respectivamente. Por claridad, escribimos x'_1 y x'_2 como x'_{n-1} y x'_n en la línea 4. Despejando de estas soluciones, encontramos que existen soluciones no negativas si y solo si existe $t \in \mathbb{Z}$ que satisfaga

$$\left\lceil -\frac{kx'_1}{q_2} \right\rceil \leq t \leq \left\lfloor \frac{kx'_2}{q_1} \right\rfloor.$$

Los enteros b_1 y b_2 en las líneas 5 y 6 representan el lado izquierdo y derecho de estas desigualdades, respectivamente. De esta manera, el algoritmo devuelve NIL si solo si este intervalo no está bien definido, es decir, si y solo si no existen soluciones enteras no negativas. Supongamos, pues, que este intervalo sí está bien definido. Entonces, podemos escoger que la variable libre t sea b_1 . Sustituyendo en (3.25) obtenemos una solución entera no negativa de la ecuación (3.24) (líneas 9 y 10) y entonces el algoritmo es

correcto para $n = 2$.

Haciendo uso de la hipótesis inductiva, queremos mostrar que el algoritmo también es correcto para $n \geq 3$ si lo es para $n - 1 \geq 2$. Entonces deseamos encontrar soluciones enteras no negativas de la ecuación (1.14) Haciendo la misma sustitución que en (1.16), recordando que q_1, \dots, q_n son coprimos por hipótesis, que definimos $\omega_1 := k$, y renombrando las variables (x en vez de x_1 , g en vez de g_2 y ω en vez de ω_2), obtenemos la ecuación

$$q_1x + g\omega = k. \quad (3.26)$$

Observemos que, como $g_1 = 1$, el entero $g = \text{mcd}\{q_2/g_1, \dots, q_n/g_1\}$, es equivalente a lo que se encuentra en la línea 12. Por definición de g , tenemos que q_1 y g son coprimos (ver el lema 1.1.6), así que por el teorema 1.1.10 tenemos que las soluciones enteras de esta ecuación están dadas por

$$\begin{cases} x = kx' + gt, \\ \omega = k\omega' - q_1t, \end{cases} \quad (3.27)$$

donde $t \in \mathbb{Z}$ es una variable libre, y x', ω' son los coeficientes de Bézout de x, ω . Recordemos de (1.16) que

$$\omega = \frac{q_2}{g}x_2 + \dots + \frac{q_n}{g}x_n. \quad (3.28)$$

Como $q > \mathbf{0}$ por hipótesis, $g > 0$ porque el máximo común divisor siempre es positivo, y exigimos que x_2, \dots, x_n sean no negativos, debe ser el caso que ω también sea no negativo. Luego, despejando de 3.27, existen soluciones no negativas de la ecuación (3.26) si y solo si existe $t \in \mathbb{Z}$ que satisfaga

$$\left\lceil -\frac{kx'}{g} \right\rceil \leq t \leq \left\lfloor \frac{k\omega'}{q_1} \right\rfloor. \quad (3.29)$$

Los enteros b_1 y b_2 en las líneas 14 y 15 representan el lado izquierdo y derecho de estas desigualdades, respectivamente. Si no existe tal variable

libre $t \in \mathbb{Z}$ es porque el intervalo $[b_1, b_2]$ no está bien definido y por lo tanto $b_2 < b_1$. El algoritmo entonces salta a la línea 27 y devuelve NIL.

Si el intervalo $[b_1, b_2]$ está bien definido, entonces podemos asegurar la no negatividad de x y de ω en (3.27) para cualquier elección de t en $[b_1, b_2]$ a causa de (3.29) y en la línea 21 nos encargamos entonces de encontrar soluciones enteras no negativas de la ecuación (3.28). Se verifica automáticamente que los coeficientes del lado derecho de esta ecuación son coprimos y constituyen justamente las entradas del vector \mathbf{q}^{tail} (c.f. línea 17). Como $g > 0$ se sigue que $\mathbf{q}^{\text{tail}} > \mathbf{0}$. Luego, \mathbf{q}^{tail} y ω satisfacen las hipótesis del algoritmo.

Por hipótesis inductiva, en la línea 21 tenemos o bien que \mathbf{x}^{tail} es entero no negativo y solución de (3.28), o bien es NIL. En el primer caso y definiendo \mathbf{x} como el vector de la línea 25 encontramos que

$$\mathbf{q}^T \mathbf{x} = q_1 x + g (\mathbf{q}^{\text{tail}})^T \mathbf{x}^{\text{tail}} = q_1 x + g \omega = k.$$

Pero $x \geq 0$ por construcción y $\mathbf{x}^{\text{tail}} \geq \mathbf{0}$ por este caso de la hipótesis inductiva. Así, \mathbf{x} también es no negativo.

Finalmente, en caso de que \mathbf{x}^{tail} sea NIL, iteramos sobre otra elección de la variable libre t y regresamos al caso pasado. En caso de que este vector sea NIL para todas las elecciones posibles de t en el intervalo de factibilidad $[b_1, b_2]$, se sigue por hipótesis inductiva que la ecuación (3.28) no tiene solución entera no negativa y por lo tanto tampoco la tiene la ecuación (1.14). Una vez agotadas estas elecciones finitas, devolvemos NIL en la línea 27.

En conclusión, si el algoritmo es correcto para vectores \mathbf{q} con dimensión $n-1 \geq 2$, entonces también es correcto para vectores \mathbf{q} con dimensión $n \geq 3$. Juntando esto con el caso base, se sigue por inducción que el algoritmo es correcto para toda $n \geq 2$, que es lo que queríamos demostrar. \square

Algoritmo 3: NonNegativeIntSolFin

Datos: $q \in \mathbb{Z}_{>0}^n$ coprimo tal que $\text{length}(q) \geq 2$. $k \geq 0$.**Resultado:** $x \in \mathbb{Z}_{\geq 0}^n$ tal que $q^T x = k$ o NIL.**inicio**

| | |
|--|----|
| $n \leftarrow \text{length}(q);$ | 1 |
| si $n = 2$ entonces | 2 |
| $x'_{n-1}, x'_n \leftarrow \text{Bezout}(q_1, q_2);$ | 3 |
| $b_1 \leftarrow \lceil -kx'_{n-1}/q_2 \rceil;$ | 4 |
| $b_2 \leftarrow \lfloor kx'_n/q_1 \rfloor;$ | 5 |
| si $b_2 < b_1$ entonces | 6 |
| devolver NIL; | 7 |
| $x_{n-1} \leftarrow kx'_{n-1} + b_1q_2;$ | 8 |
| $x_n \leftarrow kx'_n - b_1q_1;$ | 9 |
| devolver $(x_{n-1}, x_n);$ | 10 |
| $g \leftarrow \text{mcd}\{q_2, \dots, q_n\};$ | 11 |
| $x', \omega' \leftarrow \text{Bezout}(q_1, g);$ | 12 |
| $b_1 \leftarrow \lceil -kx'/g \rceil;$ | 13 |
| $b_2 \leftarrow \lfloor k\omega'/q_1 \rfloor;$ | 14 |
| $I \leftarrow \{b_1, b_1 + 1, \dots, b_2\};$ | 15 |
| $q^{\text{tail}} \leftarrow (q_{i+1}/g : 1 \leq i \leq n-1);$ | 16 |
| mientras $I \neq \emptyset$ hacer | 17 |
| elegir $t \in I;$ | 18 |
| $\omega \leftarrow k\omega' - tq_1;$ | 19 |
| $x^{\text{tail}} \leftarrow \text{NonNegativeIntSolFin}(q^{\text{tail}}, \omega);$ | 20 |
| si $x^{\text{tail}} \neq \text{NIL}$ entonces | 21 |
| $r \leftarrow \text{length}(x^{\text{tail}});$ | 22 |
| $x \leftarrow kx' + tg;$ | 23 |
| devolver $(x, x_1^{\text{tail}}, \dots, x_r^{\text{tail}});$ | 24 |
| $I \leftarrow I \setminus \{t\};$ | 25 |
| devolver NIL; | 26 |
| | 27 |

Observemos que la elección del parámetro libre t en el intervalo de factibilidad I definido en la línea 16 del algoritmo 3 es similar a la elección del subproblema S_i de optimización definido en la línea 6 del algoritmo 6. La diferencia radica en que, como todos los puntos enteros sobre la k -ésima capa entera tienen el mismo nivel de utilidad k , no es necesario desarrollar políticas de poda así como lo hicimos en el Ejemplo 1.1.18 en la Sección 1.1. De cierta manera, la única política de poda posible es la de infactibilidad por no respetar la no negatividad de un punto entero.

Teorema 3.2.2. *El algoritmo 4 es correcto.*

Demostración. A causa del teorema 3.2.1 basta verificar que el algoritmo termina y no devuelve NIL. Además, obtenemos la maximalidad de k debido a la manera en la que iniciamos el ciclo en la línea 2. Tenemos $0 \leq \eta$ por hipótesis y observemos que $\mathbf{0}$ es la única solución entera no negativa de la ecuación lineal diofantina $\mathbf{q}^T \mathbf{x} = 0$. De esta manera, si la ecuación $\mathbf{q}^T \mathbf{x} = k$ no tiene solución para $0 < k \leq \eta$, entonces el algoritmo devuelve $\mathbf{0}$ debido al teorema 3.2.1. \square

Sabemos, en realidad, por el lema 3.1.1 que el parámetro k definido en la línea 2 descenderá hasta 0 si y solo si el parámetro τ definido en (3.2) es nulo. No obstante, la demostración del teorema 3.2.2 deviene más simple cuando en el algoritmo 4 dejamos que k se encuentre en $[0, \eta]$ en vez de $[\tau, \eta]$. Esta modificación, sin embargo, no afecta en lo más mínimo la correctud o la complejidad del algoritmo.

Siguiendo la misma directriz, vale la pena mencionar lo siguiente con respecto al algoritmo 3. Varios lenguajes de programación, tales como Python, cuentan con un límite en las llamadas de recursión que el usuario puede realizar. Si bien este límite puede modificarse, aumenta la posibilidad de encontrarnos con un desbordamiento de pila, pues este algoritmo no está

expresado en forma de recursión terminal⁴.

Además, este algoritmo, por ejemplo, no minimiza el número de llamadas para calcular el máximo común divisor en la línea 12. En efecto, supongamos que un intervalo de factibilidad I definido en la línea 16 induce a que \mathbf{x}^{tail} sea NIL para todo $t \in I$. Entonces estaríamos haciendo $|I|$ llamadas recursivas a `NonNegativeIntSolFin` en la línea 21 con el mismo vector \mathbf{q}^{tail} y, por lo tanto, estaríamos calculando $|I|$ veces la misma g en la línea 12. Lo mismo ocurre con el cálculo de los coeficientes de Bézout x' y ω' en la línea 13.

A pesar de los puntos anteriores, el autor decidió escribir el algoritmo 3 de esa manera debido a que se simplificaba de manera significativa la demostración del teorema 3.2.1. Sin embargo, el autor realizó una implementación equivalente más eficiente a través de ciclos para obtener los resultados de la siguiente sección.

Algoritmo 4: Dioph

Datos:

$\mathbf{q} \in \mathbb{Z}_{>0}$ coprimo tal que $\text{length}(\mathbf{q}) \geq 2$.

$\eta \geq 0$.

Resultado:

$\mathbf{x} \in \mathbb{Z}_{\geq 0}^n$ tal que $\mathbf{q}^T \mathbf{x} = k$ con $0 \leq k \leq \eta$ maximal.

inicio

| | |
|---|---|
| para $k \leftarrow \eta$ a 0 hacer | 1 |
| $\mathbf{x} \leftarrow \text{NonNegativeIntSolFin}(\mathbf{q}, k);$ | 2 |
| si $\mathbf{x} \neq \text{NIL}$ entonces | 3 |
| devolver $\mathbf{x};$ | 4 |
| fin | 5 |

⁴Véase https://en.wikipedia.org/wiki/Tail_call.

3.3. Experimentos numéricos

Hemos mencionado que el algoritmo 3 fue escrito de esa manera para demostrar la correctud de nuestro método. El autor optó por realizar una implementación equivalente en ciclos debido al límite suave en las llamadas de recursión que permite Python. Este límite es de 3,000 llamadas recursivas en la máquina donde se realizaron los experimentos. Esto significa que, de manera predeterminada, solamente podríamos resolver problemas de dimensión $n \leq 3,000$. No obstante, en la subsección 3.3.1 realizamos experimentos con problemas de dimensión $n \leq 150,000$.

Vimos en la sección pasada que este algoritmo calcula repetidamente los mismos máximos común divisores g y los mismos coeficientes de Bézout x', ω' . Puesto que estos números dependen exclusivamente del vector \mathbf{q} , podemos calcularlos en una fase de preprocesamiento antes de llamar `NonNegativeIntSolFin`. Luego, esta rutina accede a ellos por medio de referencias.

La implementación por ciclos usa una pila de estados (i, b_1, b_2, k) , donde i indica el nivel o la variable t_i que debemos escoger; el resto de los parámetros están definidos en el algoritmo 3. La elección de t_i ciertamente es la más importante, pues determina el intervalo de factibilidad $I = [b_1, b_2]$ en el siguiente nivel $i + 1$. Existen, *a priori*, dos estrategias que podemos adoptar para este tipo de elecciones.

Por un lado, se encuentra la opción $t_i \leftarrow b_1$ o $t_i \leftarrow b_2$ para todo nivel i . Si, en tal nivel i , el intervalo de factibilidad es vacío, entonces retrocedemos en la pila y hacemos la sustitución $t_i \leftarrow t_i + 1$ o $t_i \leftarrow t_i - 1$ siempre que t_i se encuentre dentro del intervalo de factibilidad $[b_1, b_2]$. Así también, actualizamos nuestra variable x_i como lo hacemos en el algoritmo 3 y repetimos el proceso hasta obtener una solución \mathbf{x} entera no negativa.

Por el otro lado, podemos construir un árbol de la siguiente manera: en

el nivel i escogemos el punto medio $t_i \leftarrow \lfloor (b_1 + b_2)/2 \rfloor$, luego creamos los nodos adyacentes

$$\begin{cases} N_{i0} \leftarrow (i, t_i + 1, b_2, k) & \text{si } t_i + 1 \leq b_2, \\ N_{i1} \leftarrow (i, b_1, t_i - 1, k) & \text{si } t_i - 1 \geq b_1, \end{cases} \quad (3.30)$$

que representan subproblemas posiblemente a resolver en el mismo nivel. Ahora bien, dada esta elección del punto medio t_i , construimos el problema para el siguiente nivel $i+1$ con los nuevos parámetros calculados en las líneas 20, 14 y 15, en este orden. Repetimos este procedimiento hasta obtener una solución \mathbf{x} entera no negativa. Puesto que resolvemos el subproblema del siguiente nivel antes que los de los nodos adyacentes, realizamos una búsqueda *depth-first-search*.

El autor encontró en los experimentos preliminares que la segunda estrategia es mucho más eficaz. En los experimentos de la subsección 3.3.1, la primera estrategia tenía tiempos de terminación aproximadamente equivalentes a los de Ramificación y Acotamiento. Para los experimentos de la subsección 3.3.2, ambas estrategias tenían tiempos de terminación significativamente menores que Ramificación y Acotamiento. Por ello, el autor decidió realizar el análisis de resultados usando la estrategia del árbol.

Observemos de (3.30) que el orden en el que agregamos estos subproblemas determina si visitamos primero el lado izquierdo del árbol o el lado derecho. El autor realizó los experimentos con ambas posibilidades. Así pues, llamaremos `dioph_left` a la implementación que recorre primero el lado izquierdo y definimos análogamente `dioph_right`.

Al igual que en la sección 2.1, usamos el COIN-OR Branch-and-Cut (CBC) *solver* por medio de la interfaz de PuLP implementada en Python. También utilizamos dos configuraciones distintas de Ramificación y Acotamiento. La primera (R&A simple) es la implementación más cercana a una implementación “pura” de Ramificación y Acotamiento y prohíbe todo

tipo de cortes. La segunda (R&A config) permite agregar cortes de mochila, los cuales son específicamente diseñados para este tipo de problemas. Ambas configuraciones corren en un hilo y en un procesador para permitir comparaciones justas con nuestro método.

Con respecto a los resultados obtenidos en ambas subsecciones 3.3.1 y 3.3.2, mencionamos que cada observación de tiempo de cada método representa el promedio de 20 corridas. Para cada observación realizamos 2 corridas preliminares a fin de evitar sesgos en el tiempo por cuestión de temperatura en la computadora, o por cargar cosas a la memoria, o por compilación al momento de crear archivos `.pyc`, etcétera. En total, cada observación es el resultado de haber corrido el mismo experimento 22 veces.

Así también, por cuestiones de tiempo, cada proceso tenía un tiempo límite de 300 segundos para terminar de correr. De esta manera, decimos que el método no encontró una solución en caso de que más de 10 observaciones sobrepasaran este límite. Además, comparamos los valores objetivo de nuestra implementación con los de Ramificación y Acotamiento y en ningún momento encontramos que estos fueran distintos. Además, el autor menciona que su implementación se encuentra libre en GitHub⁵ así como los tiempos de terminación de cada método junto con archivos de registro que certifican estos tiempos.

Finalmente, al igual que en la sección 2.1, los experimentos se realizaron en una computadora portátil Dell XPS 15 equipada con un procesador Intel Core i7-8750H (6 núcleos físicos y 12 hilos, frecuencia base de 2.20 GHz y frecuencia máxima de 4.10 GHz). El sistema cuenta con 12 CPU lógicos disponibles y una memoria RAM de 32 GB. Todos los cálculos fueron ejecutados bajo la arquitectura x86-64. El sistema operativo usado fue Fedora Linux 42 (Server Edition).

⁵TODO: LINK

3.3.1. Experimentos en la dimensión

Para obtener resultados informativos en cuanto varía la dimensión del problema (1.1), debemos ser cuidadosos para evitar tener soluciones triviales. A lo largo de este análisis suponemos sin pérdida de generalidad que el presupuesto u del lado derecho de (1.1b) es un entero.

En primer lugar, observemos de (1.14) que si alguna entrada del vector coprimo \mathbf{q} es tal que $q_j = 1$, entonces obtenemos la solución trivial

$$x_i^* := \begin{cases} u & i = j, \\ 0 & i \neq j. \end{cases}$$

Esto se vuelve aún más trivial para nuestro método cuando ordenamos las entradas de \mathbf{q} de manera ascendente. En términos del vector objetivo \mathbf{p} , esto se traduce a que no existe una entrada p_i tal que todas las entradas que le sigan sean múltiplo de p_i . De caso contrario, tendríamos $g_{i+1} = p_i$ y por lo tanto $q_i = 1$.

En segundo lugar, es posible mostrar que todo problema (1.1) tiene una reducción al problema binario

$$\max_{\mathbf{x} \in \{0,1\}^{n_b}} \{\hat{\mathbf{p}}^T \mathbf{x} : \hat{\mathbf{p}}^T \mathbf{x} \leq u\},$$

para alguna $\hat{\mathbf{p}} \in \mathbb{Z}^{n_b}$ que puede ser obtenida de \mathbf{p} (ver [MT90]). Observemos que si $u \geq \sum_{i=1}^{n_b} \hat{p}_i$, entonces obtenemos la solución trivial $\mathbf{x} = \mathbf{e} \in \mathbb{Z}^{n_b}$. Puesto que existen *solvers* que implícitamente reducen problemas como (1.1) a su forma binaria (el ejemplo más famoso es **KnapsackSolver** de Google OR-Tools), debemos ser cuidadosos con introducir este tipo de trivialidades. Afortunadamente, una forma de evitar esta situación es exigir que el lado derecho de (1.1b) satisfaga $u < \sum_{i=1}^n p_i$.

En tercer lugar, tenemos que si el vector \mathbf{p} contiene entradas repetidas, entonces podemos reducir trivialmente la dimensión del problema (1.1). En

efecto, si $p_j = p_\ell$, encontramos que

$$\sum_{i=1}^n p_i x_i = \sum_{\substack{i=1 \\ i \neq \ell}}^n p_i z_i,$$

donde $\mathbf{z} \in \mathbb{Z}^{n-1}$ está definida como

$$z_i := \begin{cases} x_j + x_\ell, & i \in \{j, \ell\}, \\ x_i, & \text{e.o.c.}, \end{cases}$$

y el problema (1.1) es equivalente a

$$\max_{\mathbf{z} \in \mathbb{Z}^{n-1}} \{\hat{\mathbf{p}}^T \mathbf{z} : \hat{\mathbf{p}}^T \mathbf{z} \leq u, \mathbf{z} \geq \mathbf{0}\},$$

donde $\hat{\mathbf{p}} \in \mathbb{R}^{n-1}$ resulta de remover del vector \mathbf{p} su ℓ -ésima entrada.

Así pues, sea n la dimensión del problema. Dejamos que $\mathbf{p} \in \mathbb{Z}^n$ sea un vector aleatorio tomado de una distribución uniforme discreta sobre $[10, 10n]^n$. Al calcular el vector coprimo \mathbf{q} checamos que $q_i \neq 1$ para todo $1 \leq i \leq n$. De caso contrario, calculamos otro vector aleatorio \mathbf{p} . Finalmente, escogemos el lado derecho de la restricción (1.1b) de manera que

$$u := \begin{cases} 0.5 \sum_{i=1}^n p_i & n \leq 20,000, \\ 0.1 \sum_{i=1}^n p_i & n > 20,000, \end{cases}$$

para evitar obtener un problema trivial de acuerdo a la reducción binaria.

La tabla 3.1 muestra los tiempos de terminación de los métodos utilizados para realizar cada experimento. Así también, la tabla 3.1 muestra los coeficientes de variación por experimento por método. Comparando rápidamente las cifras, observamos que ambas de nuestras implementaciones son significativamente más rápidas y estables que las de Ramificación y Acotamiento en cualesquiera de sus dos configuraciones.

Una de las hipótesis por las que el autor cree que `dioph_left` presenta

resultados aún más rápidos que su contraparte **dioph_right** es la que sigue. Puesto que estamos generando vectores aleatorios con dimensiones grandes, la probabilidad de que las últimas $n-i$ entradas tengan factores en común es pequeña. Por lo tanto, obtenemos $g_{i+1} = 1$. Recordando que los coeficientes de Bézout x'_i, ω'_{i+1} satisfacen (1.26), encontramos que $(x'_i, \omega'_{i+1}) = (0, 1)$ y también debe ser el caso que $\frac{q_i}{\prod_{j=1}^i g_j} = 1$. Sustituyendo en (1.25) tenemos

$$\begin{cases} x_i = t_i, \\ \omega_{i+1} = \omega_i - t_i. \end{cases}$$

De esta manera, nuestro método **dioph_left** recorre primero soluciones pequeñas en magnitud comparadas a las obtenidas por **dioph_right**. Es decir, el primero es menos voraz que el segundo. Ciertamente podemos obtener equivalencias entre estos dos métodos al imponer un orden ascendente o descendente sobre \mathbf{q} .

3.3.2. Experimentos en el presupuesto

| n | R&A simple | R&A configurado | dioph_right | dioph_left |
|---------|-----------------|-----------------|-----------------|-----------------|
| 50 | 0.020 (0.012) | 0.023 (0.017) | 0.000 (0.000) | 0.000 (0.000) |
| 100 | 0.013 (0.002) | 0.010 (0.001) | 0.001 (0.000) | 0.001 (0.000) |
| 200 | 0.056 (0.009) | 0.023 (0.004) | 0.001 (0.000) | 0.001 (0.000) |
| 500 | 0.127 (0.021) | 0.127 (0.019) | 0.005 (0.000) | 0.005 (0.000) |
| 1,000 | 0.049 (0.008) | 0.109 (0.016) | 0.013 (0.000) | 0.013 (0.000) |
| 2,000 | 0.142 (0.021) | 0.135 (0.022) | 0.043 (0.000) | 0.042 (0.000) |
| 5,000 | 1.616 (0.252) | 1.666 (0.196) | 0.265 (0.001) | 0.245 (0.001) |
| 10,000 | 2.137 (0.442) | 2.124 (0.333) | 1.164 (0.002) | 1.063 (0.003) |
| 20,000 | 7.817 (0.596) | 7.682 (0.339) | 4.991 (0.008) | 4.595 (0.007) |
| 30,000 | 19.786 (2.852) | 19.192 (1.176) | 11.553 (0.013) | 10.743 (0.016) |
| 40,000 | 24.317 (0.163) | 24.334 (0.193) | 20.994 (0.082) | 19.471 (0.023) |
| 50,000 | 38.907 (0.203) | 38.940 (0.179) | 33.451 (0.042) | 30.887 (0.026) |
| 60,000 | 51.915 (0.600) | 51.962 (0.526) | 48.545 (0.052) | 45.072 (0.048) |
| 70,000 | 80.847 (1.025) | 81.033 (1.259) | 66.345 (0.101) | 61.885 (0.036) |
| 80,000 | 116.751 (0.838) | 116.618 (0.932) | 87.726 (0.117) | 81.418 (0.057) |
| 90,000 | 141.288 (1.655) | 141.683 (2.102) | 111.553 (0.116) | 103.962 (0.035) |
| 100,000 | 170.992 (1.988) | 170.397 (1.750) | 139.652 (0.159) | 129.440 (0.047) |
| 150,000 | - | - | - | 299.154 (0.129) |

Tabla 3.1.: Tiempos en segundos de terminación promedio de cada método con una muestra de 20 observaciones por experimento por método. La desviación estándar se encuentra entre paréntesis.

| n | R&A simple | R&A configurado | dioph_right | dioph_left |
|---------|------------|-----------------|-------------|------------|
| 50 | 0.601 | 0.741 | 0.029 | 0.018 |
| 100 | 0.157 | 0.107 | 0.009 | 0.008 |
| 200 | 0.159 | 0.179 | 0.005 | 0.004 |
| 500 | 0.168 | 0.155 | 0.004 | 0.003 |
| 1,000 | 0.170 | 0.145 | 0.002 | 0.001 |
| 2,000 | 0.150 | 0.164 | 0.001 | 0.001 |
| 5,000 | 0.158 | 0.119 | 0.003 | 0.004 |
| 10,000 | 0.209 | 0.159 | 0.002 | 0.003 |
| 20,000 | 0.077 | 0.045 | 0.002 | 0.001 |
| 30,000 | 0.146 | 0.062 | 0.001 | 0.002 |
| 40,000 | 0.007 | 0.008 | 0.004 | 0.001 |
| 50,000 | 0.005 | 0.005 | 0.001 | 0.001 |
| 60,000 | 0.012 | 0.010 | 0.001 | 0.001 |
| 70,000 | 0.013 | 0.016 | 0.002 | 0.001 |
| 80,000 | 0.007 | 0.008 | 0.001 | 0.001 |
| 90,000 | 0.012 | 0.015 | 0.001 | 0.000 |
| 100,000 | 0.012 | 0.010 | 0.001 | 0.000 |
| 150,000 | - | - | - | 0.000 |

Tabla 3.2.: Coeficientes de variación en los tiempos de terminación de acuerdo a la tabla 3.1.

Capítulo 4

Múltiples restricciones

En este último capítulo construimos un método que permite resolver programas lineales enteros generales. Mostramos que la complejidad exponencial de este tipo de programas se reduce a resolver sistemas de ecuaciones lineales en los enteros. Además, encontramos una formulación alternativa a la manera tradicional de introducir programas lineales enteros que simplifica el árbol de subproblemas generado por el algoritmo de Ramificación y Acotamiento, lo cual podría resultar en mejores tiempos de terminación.

En la exposición de este capítulo dependemos en gran medida de la forma normal de Hermite y de Smith, las cuales son tratadas extensivamente en [Sch98] y [New72]. El estilo de nuestra discusión es menos formal aunque no por ello menos rigurosa.

Sea $\mathbf{p} \in \mathbb{R}^n$ esencialmente entero y sea $\mathbf{q} \in \mathbb{Z}^n$, de manera que $\mathbf{p} = m\mathbf{q}$ para algún escalar $m > 0$. Vimos en la subsección 1.2.2 que el problema (1.1) es equivalente a

$$\max_{\mathbf{x} \in \mathbb{Z}^n} \{\mathbf{q}^T \mathbf{x} : \mathbf{q}^T \mathbf{x} \in \{0, \dots, \eta\}, \mathbf{x} \geq \mathbf{0}\},$$

donde $\eta \in \mathbb{Z}$ está definida en el lema 1.2.7. Por ello, introducimos el problema con múltiples restricciones (4.1) en términos del vector \mathbf{q} y el parámetro η en vez del vector \mathbf{p} y el lado derecho u de (1.1b). Así pues, sea $A \in \mathbb{Q}^{m \times n}$

una matriz racional con renglones linealmente independientes y sea $\mathbf{b} \in \mathbb{Q}^m$ un vector. Consideremos el problema

$$\max_{\mathbf{x} \in \mathbb{Z}^n} \mathbf{q}^T \mathbf{x}, \quad (4.1a)$$

$$\text{s.a.} \quad \mathbf{q}^T \mathbf{x} \leq \eta, \quad (4.1b)$$

$$A\mathbf{x} = \mathbf{b}, \quad (4.1c)$$

$$\mathbf{x} \geq \mathbf{0}.$$

Observación. En contraste con el primer caso del teorema 1.2.9, no podemos asegurar que la solución se encuentre sobre la η -ésima capa entera $H_{\mathbf{q}, \eta \|\mathbf{q}\|^{-2}}$ aún cuando \mathbf{q} tenga una entrada negativa. En efecto, si $A = \mathbf{q}^T$ y $\mathbf{b} = \eta - 1$, la restricción (4.1b) se vuelve redundante y obtenemos un problema similar a (1.1). En caso de que \mathbf{q} tenga al menos una entrada negativa, el primer caso del teorema 1.2.9 nos indica que el problema es factible y que la solución se encuentra en la $(\eta - 1)$ -ésima capa entera.

Supongamos que el problema (4.1) es factible. Debido a la restricción presupuestaria (4.1b), sabemos que la solución se encuentra en alguna capa entera $H_{\mathbf{q}, k \|\mathbf{q}\|^{-2}}$ con parámetro entero $k \leq \eta$. Esto motiva el siguiente resultado.

Teorema 4.0.1. *Sea $M \in \mathbb{Z}^{n \times (n-1)}$ la matriz definida en (1.38) y $\boldsymbol{\nu} \in \mathbb{Z}^n$ el vector definido en (1.37). Entonces el problema (4.1) es equivalente al problema*

$$\max_{k \in \mathbb{Z}, \mathbf{t} \in \mathbb{Z}^{n-1}} k, \quad (4.2a)$$

$$\text{s.a.} \quad k \leq \eta, \quad (4.2b)$$

$$AM\mathbf{t} = \mathbf{b} - kA\boldsymbol{\nu}, \quad (4.2c)$$

$$M\mathbf{t} \geq -k\boldsymbol{\nu}. \quad (4.2d)$$

Demostración. Por el teorema 1.2.15 y la discusión que le sucede, sabemos que la transformación lineal

$$(k, \mathbf{t}) \mapsto \mathbf{x} := k\boldsymbol{\nu} + M\mathbf{t}$$

es un isomorfismo entre las redes $\Lambda_p \oplus \Lambda_h$ definidas en (1.42) y \mathbb{Z}^n . Así, tenemos

$$\begin{aligned} A\mathbf{x} = \mathbf{b} &\iff AM\mathbf{t} = \mathbf{b} - kA\boldsymbol{\nu}, \\ \mathbf{x} \geq \mathbf{0} &\iff M\mathbf{t} \geq -k\boldsymbol{\nu}, \end{aligned}$$

y por lo tanto basta mostrar que si un vector es factible para un problema, entonces satisface la correspondiente restricción presupuestaria (4.1b) o (4.2b) del otro problema.

Sea $\mathbf{x} \in \mathbb{Z}^n$ un vector factible de (4.1), entonces existe $(k, \mathbf{t}) \in \mathbb{Z}^n$ que satisface $\mathbf{x} = k\boldsymbol{\nu} + M\mathbf{t}$. Por los lemas 1.2.11 y 1.2.12 encontramos que

$$k = \mathbf{q}^T \mathbf{x} \leq \eta,$$

y entonces (k, \mathbf{t}) es factible. Como \mathbf{x} fue arbitrario, se sigue que la solución del problema (4.1) es una cota inferior del problema (4.2). La demostración de que la solución de (4.2) es una cota inferior de (4.1) es análoga usando el mismo isomorfismo.

Finalmente, supongamos que $(k, \mathbf{t}) \in \mathbb{Z}^n$ es solución de (4.2). Si existe $\tilde{\mathbf{x}}$ factible para (4.1) con utilidad $\mathbf{q}^T \tilde{\mathbf{x}} = \tilde{k}$ estrictamente mayor, entonces consideramos $(\tilde{k}, \tilde{\mathbf{t}})$ tal que $\tilde{\mathbf{x}} = \tilde{k}\boldsymbol{\nu} + M\tilde{\mathbf{t}}$. Este vector también es factible con utilidad $k < \tilde{k} \leq \eta$, y entonces (k, \mathbf{t}) no era la solución de (4.2). Obtenemos una contradicción. \square

Al inicio de esta tesis mencionamos que las políticas de poda de Ramificación y Acotamiento operan ineficientemente cuando el vector objetivo es

ortogonal a una de sus restricciones. Observemos en los problemas (1.1) y (4.1) que el vector objetivo \mathbf{q} es ortogonal a las restricciones presupuestarias (1.1b) y (4.1b). Esto también es cierto para el problema equivalente (4.2), pues el vector objetivo $k\mathbf{e}_1$ es ortogonal a la restricción (4.2b). A pesar de lo anterior, el problema equivalente induce a que las políticas de poda sean más eficientes.

Teorema 4.0.2. *Sea $(k_{PR}^*, \mathbf{t}_{PR}^*)$ el óptimo del problema relajado de (4.2) y supongamos que k_{PR}^* no es entero. Entonces el subproblema generado al añadir la restricción $k \geq \lceil k_{PR}^* \rceil$ es infactible.*

Demostración. Supongamos que el subproblema es factible. Puesto que k_{PR}^* no es entero, existe $\tau \in \mathbb{Z}$ tal que $\tau - 1 < k_{PR}^* < \tau$. Al añadir la restricción $k \geq \lceil k_{PR}^* \rceil = \tau$ al problema (4.2), encontramos que el valor óptimo de este subproblema es estrictamente mayor que k_{PR}^* . Pero esto es una contradicción ya que en problemas de maximización el valor óptimo de un problema es una cota superior del valor óptimo de cualesquiera de sus subproblemas. \square

Debido a este teorema, siempre es mejor priorizar ramificaciones en k_{PR}^* puesto que nos deshacemos de manera inmediata subproblemas infactibles.

A continuación desacoplamos el problema (4.2) en un subproblema de maximización y en otro de factibilidad. Supongamos, sin pérdida de generalidad, que las entradas de A y \mathbf{b} son enteras. En el capítulo 2 de [New72] es introducida la forma normal de Hermite de la matriz A , la cual afirma que existe una matriz unimodular $U \in \mathbb{Z}^{n \times n}$ que satisface $AU = [H, \mathbf{0}]$, donde $H \in \mathbb{Z}^{m \times m}$ es triangular inferior y no singular.

Con esto en mente, introducimos el subproblema de (4.2) como

$$\max_{k \in \mathbb{Z}, \mathbf{y} \in \mathbb{Z}^m} k, \quad (4.3a)$$

$$\text{s.a.} \quad k \leq \eta, \quad (4.3b)$$

$$A\tilde{\mathbf{y}} = \mathbf{b} - kA\boldsymbol{\nu}, \quad (4.3c)$$

donde

$$\tilde{\mathbf{y}} := U \begin{pmatrix} \tilde{\mathbf{y}}_m \\ \tilde{\mathbf{y}}_{n-m} \end{pmatrix} = U_m \tilde{\mathbf{y}}_m + U_{n-m} \tilde{\mathbf{y}}_{n-m} \in \mathbb{Z}^n, \quad (4.4)$$

con $\tilde{\mathbf{y}}_m \in \mathbb{Z}^m$ y $\tilde{\mathbf{y}}_{n-m} \in \mathbb{Z}^{n-m}$. Denotamos por U_m y U_{n-m} las primeras m columnas y últimas $n - m$ columnas de U , respectivamente. Observemos que para toda $k \in \mathbb{Z}$ se cumple

$$AU \begin{pmatrix} H^{-1}(\mathbf{b} - kA\boldsymbol{\nu}) \\ \tilde{\mathbf{y}}_{n-m} \end{pmatrix} = [H, 0] \begin{pmatrix} H^{-1}(\mathbf{b} - kA\boldsymbol{\nu}) \\ \tilde{\mathbf{y}}_{n-m} \end{pmatrix} = \mathbf{b} - kA\boldsymbol{\nu}, \quad (4.5)$$

lo cual sugiere definir

$$\tilde{\mathbf{y}}_m := H^{-1}(\mathbf{b} - kA\boldsymbol{\nu}). \quad (4.6)$$

No obstante, debemos asegurarnos que este vector sea entero. Observemos que $\tilde{\mathbf{y}}_{n-m}$ es un vector libre, así que en realidad este subproblema tiene dimensión $m + 1$. Definimos el conjunto de factibilidad

$$\mathcal{F} := \{k \in \mathbb{Z}: H^{-1}(\mathbf{b} - kA\boldsymbol{\nu}) \in \mathbb{Z}^m, k \leq \eta\} \quad (4.7)$$

Puesto que $H \in \mathbb{Z}^{m \times m}$ es no singular, para cada $k \in \mathbb{Z}$, existe una única solución $\tilde{\mathbf{y}}_m \in \mathbb{R}^m$ del sistema de ecuaciones $H\tilde{\mathbf{y}}_m = \mathbf{b} - kA\boldsymbol{\nu}$. Como, además, H es triangular inferior, podemos resolver rápidamente este sistema de ecuaciones y verificar si, para cada $k \in \mathbb{Z}$, la correspondiente solución $\tilde{\mathbf{y}}_m$ es entera o no.

Si \mathcal{F} es vacío, deducimos que el subproblema (4.3) es infactible y por lo tanto (4.2) también lo es. Supongamos, pues, que $\mathcal{F} \neq \emptyset$. No es difícil observar que \mathcal{F} tiene un elemento maximal k^* y que este elemento es la solución al subproblema (4.3). Luego, dada esta solución $k^* \in \mathbb{Z}$, buscamos

resolver el subproblema de (4.2)

$$M\mathbf{t} = \tilde{\mathbf{y}}, \quad (4.8a)$$

$$M\mathbf{t} \geq -k^*\boldsymbol{\nu}. \quad (4.8b)$$

Tenemos un sistema de n ecuaciones lineales con $2n - m - 1$ incógnitas, por lo que tendremos que lidiar con $n - m - 1$ variables libres. En efecto, sustituyendo (4.4) en (4.8a), obtenemos

$$\begin{aligned} M\mathbf{t} = \tilde{\mathbf{y}} &= U_m\tilde{\mathbf{y}}_m + U_{n-m}\tilde{\mathbf{y}}_{n-m} \\ \iff [M \mid -U_{n-m}] \begin{pmatrix} \mathbf{t} \\ \tilde{\mathbf{y}}_{n-m} \end{pmatrix} &= U_m\tilde{\mathbf{y}}_m. \end{aligned} \quad (4.9)$$

En el capítulo 2 de [New72] también se introduce la forma normal de Smith, de la cual obtenemos dos matrices unimodulares $S \in \mathbb{Z}^{n \times n}$ y $T \in \mathbb{Z}^{(2n-m-1) \times (2n-m-1)}$ que satisfacen

$$S[M, -U_{n-m}]T = D \in \mathbb{Z}^{n \times (2n-m-1)},$$

donde D es una matriz diagonal cuyas n primeras entradas son distintas de cero y las restantes $n - m - 1$ son cero. Si multiplicamos S por la izquierda en ambos lados de la ecuación (4.9), tenemos

$$DT^{-1} \begin{pmatrix} \mathbf{t} \\ \tilde{\mathbf{y}}_{n-m} \end{pmatrix} = SU_m\tilde{\mathbf{y}}_m. \quad (4.10)$$

Si d_i no divide a $(SU_m\tilde{\mathbf{y}}_m)_i$ para alguna $i \in \{1, \dots, n\}$, encontramos que la primera ecuación del subproblema (4.8) no tiene solución en los enteros, lo que implica que la elección de k^* fue la incorrecta para asegurar soluciones enteras a este subproblema. De ser este el caso, redefinimos nuestro conjunto de factibilidad \mathcal{F} (ver (4.7)) como $\mathcal{F} \leftarrow \mathcal{F} \setminus \{k^*\}$. Si \mathcal{F} ahora es vacío, entonces (4.2) es infactible, y en caso contrario escogemos el nuevo elemento

de maximal de \mathcal{F} y repetimos el proceso.

Supongamos, pues, que $d_i \mid (SU_m \tilde{\mathbf{y}}_m)_i$ para todo $i \in \{1, \dots, n\}$, por lo que obtenemos n soluciones enteras $\mathbf{r} \in \mathbb{Z}^n$ y $n - m - 1$ variables libres $\mathbf{s} \in \mathbb{Z}^{n-m-1}$:

$$T^{-1} \begin{pmatrix} \mathbf{t} \\ \tilde{\mathbf{y}}_{n-m} \end{pmatrix} = \begin{pmatrix} \mathbf{r} \\ \mathbf{s} \end{pmatrix}.$$

Por lo tanto, nuestro vector \mathbf{t} es una función lineal de \mathbf{s} , es decir, $\mathbf{t} = \mathbf{t}(\mathbf{s})$. En términos del problema original (4.1), hemos encontrado, hasta este punto, los vectores $\mathbf{x}(\mathbf{s}) := k^* \boldsymbol{\nu} + M\mathbf{t}(\mathbf{s})$ que maximizan la utilidad y que satisfacen todas las restricciones excepto, posiblemente, las de no negatividad.

Consideremos el conjunto de vectores $\mathbf{s} \in \mathbb{Z}^{n-m-1}$ que inducen a que $\mathbf{t}(\mathbf{s})$ satisfaga (4.8b):

$$\mathcal{S} := \{\mathbf{s} \in \mathbb{Z}^{n-m-1} : M\mathbf{t}(\mathbf{s}) \geq -k^* \boldsymbol{\nu}\}$$

Por un lado, es sabido que los programas enteros tales como (4.1) o (4.2) son problemas difíciles de resolver, a excepción de cuando la matriz de restricciones $A \in \mathbb{Z}^{m \times n}$ es totalmente unimodular. De manera superficial, decimos que un problema es difícil de resolver si no es conocida la existencia de un algoritmo con complejidad polinomial que lo pueda resolver.

Por el otro lado, a lo largo de este capítulo hemos resuelto todos los problemas en tiempo polinomial. En efecto, obtener M y $\boldsymbol{\nu}$ de (1.38) y (1.37) se reduce a multiplicar números y calcular coeficientes de Bézout, al igual que máximos común divisores. En [Sch98] y [New72] se muestra que realizar este tipo de cálculos, así como de obtener las formas normales de Hermite y de Smith, son problemas acotados en tiempo polinomial.

Entonces, la única deducción posible es que el problema de determinar si el conjunto \mathcal{S} es vacío, o cuántos elementos tiene, o cuáles son los elementos que contiene, son todos problemas difíciles de resolver. Esta complejidad se reduce drásticamente en dos casos especiales.

En primer lugar, si $m = n - 1$, entonces no hay parámetros libres. De manera gráfica, el politopo factible resultante es un semirrayo o un segmento de línea. Al momento de escoger la k^* -ésima capa entera, estamos agregando la ecuación $k = k^*$, con lo que obtenemos un sistema lineal entero de n ecuaciones con n incógnitas, y entonces la solución es única. Resta verificar que esta solución es entera y satisface (4.8b). Este caso se ilustra en el ejemplo 4.0.3.

En segundo lugar, si $m = n - 2$, obtenemos un solo parámetro libre $s \in \mathbb{Z}$, con lo que podemos determinar rápidamente la existencia o inexistencia de un conjunto de factibilidad en s que induce a que $\mathbf{t}(s)$ satisfaga (4.8b). Este caso se ilustra en el ejemplo 4.0.4.

A modo de resumen, mostramos en el pseudocódigo 5 la forma de resolver problemas del tipo (4.2). Por el teorema 4.0.1, este método también resuelve problemas del tipo (4.1). Después de presentar los ejemplos 4.0.3 y 4.0.4, mostramos una manera con la cual podemos deshacernos del ciclo infinito en la línea 7.

A fin de obtener las formas normales de Hermite y de Smith de la matriz de restricciones $A \in \mathbb{Z}^{m \times n}$ de los siguiente ejemplos, el autor utilizó la librería `hsnf` de Python¹.

Ejemplo 4.0.3. Consideremos el problema con $n = 2$ variables y $m = 1$ restricciones de igualdad

$$\begin{aligned} & \text{máx } x - y, \\ \text{s.a. } & x - y \leq 12, \\ & 3x + 5y = 25, \\ & x, y \geq 0. \end{aligned}$$

En este caso tenemos $A = (3, 5)$, $\mathbf{b} = 25$, y también $\mathbf{q} = (1, -1)^T$, al igual

¹Véase <https://hsnf.readthedocs.io/en/latest/index.html>.

que $\eta = 12$. De (1.37) y (1.38) calculamos

$$\boldsymbol{\nu} = \begin{pmatrix} 1 \\ 0 \end{pmatrix}, \quad M = \begin{pmatrix} -1 \\ -1 \end{pmatrix}.$$

De la forma normal de Hermite de A tenemos

$$H = 1, \quad U = \begin{pmatrix} 2 & -5 \\ -1 & 3 \end{pmatrix},$$

y de la forma normal de Smith de $[M \mid -U_{n-m}]$,

$$S = \begin{pmatrix} -1 & 0 \\ 1 & -1 \end{pmatrix}, \quad D = \begin{pmatrix} 1 & 0 \\ 0 & 8 \end{pmatrix}, \quad T = \begin{pmatrix} 1 & 5 \\ 0 & 1 \end{pmatrix}.$$

Como $H = 1$, se sigue que $H^{-1}(\mathbf{b} - kA\boldsymbol{\nu}) = 25 - 3k$ es entero para todo $k \in \mathbb{Z}$. Así, el conjunto factible \mathcal{F} definido en 4.7 está dado por

$$\mathcal{F} = \{k \in \mathbb{Z} : k \leq \eta = 12\}.$$

Entonces escogemos $k^* = 12$ por ser el elemento maximal de \mathcal{F} . Luego,

$$\mathbf{z} := SU_m \tilde{\mathbf{y}}_m = SU_m (H^{-1}(\mathbf{b} - k^* A\boldsymbol{\nu})) = \begin{pmatrix} 22 \\ -33 \end{pmatrix}.$$

Observemos que $D_{22} \nmid z_2$, y entonces el subproblema (4.8) no es factible para la elección de $k^* = 12$. Escogemos el segundo elemento de \mathcal{F} más grande, con lo que tenemos $k^* = 11$. Siguiendo con el mismo procedimiento, encontramos ahora que $\mathbf{z} = (16, -24)$. En este caso, la diagonal de D sí divide, elemento a elemento, las entradas de \mathbf{z} , y entonces $\mathbf{r} = (16, -3)$. Puesto que $n - m - 1 = 0$, no hay variables libres. Tenemos

$$\begin{pmatrix} \mathbf{t} \\ \tilde{\mathbf{y}}_{n-m} \end{pmatrix} = T \begin{pmatrix} 16 \\ -3 \end{pmatrix} = \begin{pmatrix} 1 \\ -3 \end{pmatrix},$$

y verificamos que se satisfaga (4.8b):

$$M\mathbf{t} + k^*\boldsymbol{\nu} = 1 \begin{pmatrix} -1 \\ -1 \end{pmatrix} + 11 \begin{pmatrix} 1 \\ 0 \end{pmatrix} = \begin{pmatrix} 10 \\ -1 \end{pmatrix} \not\geq \mathbf{0}.$$

Ahora la elección de $k^* = 11$ dio un punto entero pero con una entrada negativa. Es decir, el subproblema (4.8) es infactible dada esta elección.

Repetimos este procedimiento hasta llegar a $k^* = 3$. En este caso encontramos que $(\mathbf{t}, \tilde{\mathbf{y}}_{n-m}) = (-2, 6)^T$. Por lo tanto,

$$M\mathbf{t} + k^*\boldsymbol{\nu} = -2 \begin{pmatrix} -1 \\ -1 \end{pmatrix} + 3 \begin{pmatrix} 1 \\ 0 \end{pmatrix} = \begin{pmatrix} 5 \\ 2 \end{pmatrix} \geq \mathbf{0}.$$

Luego, $(k^*, \mathbf{t}) := (3, -2)$ es el óptimo del programa (4.2). Por el teorema 4.0.1, concluimos que $(x^*, y^*) = (5, 2)$ es el óptimo de (4.1).

Ejemplo 4.0.4. Ahora consideremos el problema con $n = 3$ variables y $m = 1$ restricciones de igualdad

$$\begin{aligned} & \text{máx } x - y + 2z, \\ & \text{s.a. } x - y + 2z \leq 10 \\ & \quad 3x + 4y - z = 15 \\ & \quad x, y, z \geq 0. \end{aligned}$$

En este caso tenemos $A = (3, 4, -1)$, $\mathbf{b} = 15$, y también $\mathbf{q} = (1, -1, 2)^T$, al igual que $\eta = 10$. De (1.37) y (1.38) calculamos

$$\boldsymbol{\nu} = \begin{pmatrix} 1 \\ 0 \\ 0 \end{pmatrix}, \quad M = \begin{pmatrix} 1 & 0 \\ -1 & 2 \\ -1 & 1 \end{pmatrix}.$$

De la forma normal de Hermite de A tenemos

$$H = 1, \quad U = \begin{pmatrix} 0 & 0 & 1 \\ 0 & 1 & 0 \\ -1 & 4 & 3 \end{pmatrix},$$

y de la forma normal de Smith de $[M \mid -U_{n-m}]$,

$$S = \begin{pmatrix} 1 & 0 & 0 \\ -1 & -1 & 0 \\ 3 & 4 & -1 \end{pmatrix}, \quad D = \begin{pmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 7 & 0 \end{pmatrix}, \quad T = \begin{pmatrix} 1 & 0 & 0 & 1 \\ 0 & 0 & 1 & 0 \\ 0 & 1 & 2 & -1 \\ 0 & 0 & 0 & 1 \end{pmatrix}.$$

Puesto que $H = 1$, tenemos de (4.7) que el conjunto factible es

$$\mathcal{F} = \{k \in \mathbb{Z} : k \leq \eta = 10\}.$$

Ahora bien, seguimos exactamente el mismo procedimiento que en el ejemplo 4.0.3 hasta llegar a $k^* = 5$. Llegando a la línea 18 del pseudocódigo 5, encontramos que

$$T^{-1} \begin{pmatrix} \mathbf{t} \\ \tilde{\mathbf{y}}_{n-m} \end{pmatrix} = \begin{pmatrix} 0 \\ 0 \\ 0 \\ s \end{pmatrix} \implies \begin{pmatrix} \mathbf{t} \\ \tilde{\mathbf{y}}_{n-m} \end{pmatrix} = s \begin{pmatrix} 1 \\ 0 \\ -1 \\ 1 \end{pmatrix},$$

donde $s \in \mathbb{Z}$ es la única variable libre. En este caso podemos determinar rápidamente un intervalo de existencia: tenemos $M\mathbf{t}(s) \geq -k^*\boldsymbol{\nu}$ si y solo si

$$s \begin{pmatrix} 1 \\ 0 \\ -1 \end{pmatrix} \geq \begin{pmatrix} -5 \\ 0 \\ 0 \end{pmatrix},$$

de donde se sigue inmediatamente que $s \in \{-5, -4, \dots, 0\}$. Sustituyendo

cada posible valor de s en $\mathbf{t}(s)$ y transformando a $\mathbf{x}^*(s) = k^*\boldsymbol{\nu} + M\mathbf{t}(s)$, encontramos que

$$\left\{ \begin{pmatrix} 0 \\ 5 \\ 5 \end{pmatrix}, \begin{pmatrix} 1 \\ 4 \\ 4 \end{pmatrix}, \begin{pmatrix} 2 \\ 3 \\ 3 \end{pmatrix}, \begin{pmatrix} 3 \\ 2 \\ 2 \end{pmatrix}, \begin{pmatrix} 4 \\ 1 \\ 1 \end{pmatrix}, \begin{pmatrix} 5 \\ 0 \\ 0 \end{pmatrix} \right\}$$

son las seis soluciones del problema (4.1). Ciertamente, todas alcanzan un nivel de utilidad $k^* = 5$.

Si el problema (4.1) es factible, por la equivalencia del teorema 4.0.1, existe $k^* \in \mathbb{Z}$ que es el valor óptimo del problema (4.2), así que eventualmente saldremos del ciclo infinito de la línea 7. En caso de que el problema (4.1) sea infactible, nada asegura, por el momento, que salgamos de este ciclo infinito. A continuación veremos cómo arreglar este problema, y en el proceso seremos capaces de eliminar la restricción presupuestaria (4.1b). Por lo tanto, en esta última parte, podremos encontrar soluciones a programas lineales enteros generales.

Sea $A \in \mathbb{Z}^{m \times n}$ una matriz con renglones linealmente independientes y sea $\mathbf{b} \in \mathbb{Z}^m$ un vector. Definamos el poliedro

$$P := \{\mathbf{x} \in \mathbb{R}^n : A\mathbf{x} = \mathbf{b}, \mathbf{x} \geq \mathbf{0}\}.$$

Sea $\mathbf{q} \in \mathbb{Z}^n$ un vector coprimo y consideremos ambos problemas de maximización y minimización sobre este poliedro

$$\ell^* := \min_{\mathbf{x} \in P} \{\mathbf{q}^T \mathbf{x}\}, \quad u^* := \max_{\mathbf{x} \in P} \{\mathbf{q}^T \mathbf{x}\}, \quad (4.11)$$

y definamos

$$\tau := \lceil \ell^* \rceil, \quad \eta := \lfloor u^* \rfloor. \quad (4.12)$$

Observemos de (4.11) que siempre se cumple que $\tau \leq \eta$. Ciertamente, la restricción $\tau \leq \mathbf{q}^T \mathbf{x} \leq \eta$ es válida para el programa lineal entero

$\max_{P \cap \mathbb{Z}^n} \{\mathbf{q}^T \mathbf{x}\}$ y, por lo tanto, este problema es equivalente a

$$\max_{\mathbf{x} \in \mathbb{Z}^n} \quad \mathbf{q}^T \mathbf{x}, \quad (4.13a)$$

$$\text{s.a.} \quad \tau \leq \mathbf{q}^T \mathbf{x} \leq \eta, \quad (4.13b)$$

$$A\mathbf{x} = \mathbf{b}, \quad (4.13c)$$

$$\mathbf{x} \geq \mathbf{0},$$

En primer lugar, si $\eta = \infty$, entonces el problema relajado es no acotado y no tiene sentido usar el pseudocódigo (5) para buscar una solución de este problema. Esta estrategia es de igual manera usada por todo *solver* de código libre o comercial antes de tan siquiera buscar una solución.

En segundo lugar, si $\tau = -\infty$ y $\eta < \infty$, entonces el problema (4.13) es factible y representa exactamente el mismo problema que (4.1). En este caso, como lo hemos discutido, siempre saldremos del ciclo infinito en la línea 7, por lo que el método delineado por el pseudocódigo (5) eventualmente terminará con una solución de este problema.

Finalmente, si $-\infty < \tau \leq \eta < \infty$ nos encontramos en la situación ideal. Esto se debe a que podemos reemplazar el ciclo en la línea 7 por algo del estilo “**para** $k \leftarrow \eta$ **a** τ **hacer**...”. Es decir, sabemos exactamente cuántas capas enteras debemos recorrer para que el método delineado por el pseudocódigo (5) termine. Observemos que, en este caso, existe la posibilidad de que $P \neq \emptyset$ pero $P \cap \mathbb{Z}^n = \emptyset$. Sin realizar modificaciones grandes al pseudocódigo (5), encontramos que, o bien termina con una solución \mathbf{x}^* del problema (4.13), o bien muestra en $\eta - \tau + 1$ pasos que este problema es infactible.

Cabe mencionar que en la subsección (1.1.2) indicamos que existen diversos algoritmos capaces de resolver rápidamente problemas lineales del estilo (4.11). Así pues, la parte de calcular los valores τ y η puede ser considerada como una parte de preprocesamiento. Recordemos que el método de Rami-

ficación y Acotamiento, en el peor de los casos, necesita resolver un número exponencial de problemas relajados de (4.13). Nuestro método, en cambio, solo necesita resolver, en el peor de los casos, dos problemas relajados.

A modo de conclusión, al autor le gustaría mencionar que futuras líneas de investigación podrían estar concentradas en resolver el problema de la línea 19. Esto se reduce a investigar sistemas de desigualdades lineales en los enteros. Existen tres posibilidades para estas investigaciones con respecto al vector de variables libres $\mathbf{s} \in \mathbb{Z}^{n-m-1}$:

1. Decidir la existencia de este vector: si bien no podríamos obtener la solución entera \mathbf{x}^* , sí podríamos concluir que k^* es el valor óptimo de (4.2) y, por el lema 1.2.11 así como del teorema 4.0.1, también es el valor óptimo de (4.1).
2. En caso de tener existencia, determinar el número de estos vectores: además de saber que k^* es el valor óptimo de (4.1), también conoceríamos el número de soluciones que tiene este problema.
3. En caso de tener existencia, calcular todos estos vectores: además de saber que k^* es el óptimo de (4.1) y de conocer cuántas soluciones tiene este problema, conoceríamos también cuáles son esas soluciones.

Otra posible futura línea de investigación, más aplicada pero no por ello menos interesante, es desarrollar las consecuencias del teorema 4.0.2. Es una creencia del autor que los tiempos de terminación de Ramificación y Acotamiento usando la formulación equivalente (4.2) serán menores que usando la formulación tradicional (4.1). Para lograr esto, necesitaremos calcular rápida y eficientemente la matriz $M \in \mathbb{Z}^{n \times (n-1)}$ y el vector $\boldsymbol{\nu} \in \mathbb{Z}^n$ definidos en (1.38) y (1.37), respectivamente.

Pseudocódigo 5:

Datos: Vector coprimo $\mathbf{q} \in \mathbb{Z}^n$, $\eta \in \mathbb{Z}$, $A \in \mathbb{Z}^{m \times n}$ y $\mathbf{b} \in \mathbb{Z}^m$.**Resultado:** Solución óptima \mathbf{x}^* de (4.1).

| | |
|--|----|
| inicio | 1 |
| Calcular M y $\boldsymbol{\nu}$ de (1.38) y (1.37); | 2 |
| Obtener U y H de la forma normal de Hermite de A ; | 3 |
| Particionar U en U_m y U_{n-m} tal que $[U_m \mid U_{n-m}] = U$; | 4 |
| Obtener S y T de la forma normal de Smith de $[M \mid -U_{n-m}]$; | 5 |
| $k \leftarrow \eta$; | 6 |
| mientras $1 + 1 = 2$ hacer | 7 |
| Obtener $\tilde{\mathbf{y}}_m$ de $H\tilde{\mathbf{y}}_m = \mathbf{b} - kA\boldsymbol{\nu}$; | 8 |
| si $\tilde{\mathbf{y}}_m \in \mathbb{Z}^m$ entonces | 9 |
| └ ir al paso 12; | 10 |
| $k \leftarrow k - 1$; | 11 |
| $\mathbf{z} \leftarrow SU_m\tilde{\mathbf{y}}_m$; | 12 |
| $\mathbf{r} \leftarrow \mathbf{0}_n$; | 13 |
| para $i \leftarrow 1$ a n hacer | 14 |
| si $D_{ii} \nmid z_i$ entonces | 15 |
| └ ir al paso 11; | 16 |
| $r_i \leftarrow z_i/D_{ii}$; | 17 |
| $(\mathbf{t}(\mathbf{s}), \tilde{\mathbf{y}}_{n-m}(\mathbf{s})) \leftarrow T(\mathbf{r}, \mathbf{s})^T$; | 18 |
| si existe \mathbf{s} tal que $M\mathbf{t}(\mathbf{s}) \geq -k\boldsymbol{\nu}$ entonces | 19 |
| $\mathbf{x}^* \leftarrow k\boldsymbol{\nu} + M\mathbf{t}(\mathbf{s})$; | 20 |
| devolver \mathbf{x}^* ; | 21 |
| ir al paso 11; | 22 |

Capítulo A

Algoritmo de Ramificación y Acotamiento

El Algoritmo 6 presenta una versión rudimentaria del algoritmo de ramificación y acotamiento. El rendimiento de este método depende en gran parte de la elección del subproblema (6) pues partir de su solución podemos obtener cotas que nos permitan podar subárboles lo más pronto posible. En la práctica, también debemos tomar en cuenta estrategias de selección que permitan paralelizar la solución de los problemas relajados, o que minimicen la sobrecarga computacional de “saltar” de un subproblema a otro.

Además del problema de selección de los nodos, también se encuentra el de creación de estos nodos. En la línea (15) ramificamos S_i usando una de las técnicas más básicas: elegir x_j^i fraccionario y generar S_{i0} , S_{i1} a partir de los cortes válidos $x_j \leq \lfloor x_j^i \rfloor$ y $x_j \geq \lceil x_j^i \rceil$. En realidad, existen muchas otras estrategias de corte, tales como los cortes de Gomory, cortes SOS1, cortes de pseudo costos, cortes fuertes, cortes de mochila, etcétera.

Implementaciones comerciales y de código abierto extienden el algoritmo de ramificación y acotamiento a partir de otros esquemas. Es común que estas cuenten con métodos de presolución para disminuir el tamaño del problema original o con heurísticas para generar nuevos tipos de cortes. Normalmente, en las implementaciones comerciales, las heurísticas no son

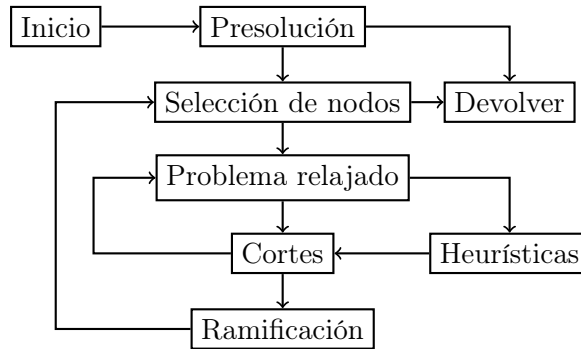


Figura A.1.: Flujo típico de algoritmos que resuelven problemas lineales mixtos. Adaptado de [Oli17].

de dominio público. Referirse a [AA95] para conocer algunas técnicas de presolución.

La Figura A.1 muestra el flujo típico de algoritmos que resuelven problemas lineales mixtos. Implementaciones comunes de código abierto son COIN-OR CBC, HiGHS y SCIP, mientras que algunas implementaciones comerciales son Gurobi Optimizer, IBM ILOG CPLEX Optimizer, y Fico Xpress Solver. Referirse a las documentaciones respectivas para obtener más información sobre el contexto en el que entra el algoritmo de ramificación y acotamiento en la resolución de problemas lineales.

Algoritmo 6: Algoritmo de Ramificación y Acotamiento. Adaptado de [Oli17].

Datos: Problema de maximización lineal S_0 .

Resultado: Solución óptima entera \mathbf{x}^* y valor óptimo z_{PE}^* .

```

inicio                                                                 1
     $\mathcal{L} \leftarrow \{S_0\};$                                          2
     $\mathbf{x}^* \leftarrow -\infty;$                                          3
     $z_{PE}^* \leftarrow -\infty;$                                          4
    mientras  $\mathcal{L} \neq \emptyset$  hacer                                   5
        elegir de  $\mathcal{L}$  subproblema  $S_i$ ;                               6
        obtener de  $S_i$  valor óptimo  $z_i^*$  y solución óptima  $\mathbf{x}^i$ ;      7
         $\mathcal{L} \leftarrow \mathcal{L} \setminus \{S_i\};$                          8
        si  $S_i = \emptyset$  o  $z_i^* \leq z_{PE}^*$  entonces                 9
            ir al paso (6);                                           10
        si  $\mathbf{x}^i \in \mathbb{Z}^n$  entonces                                   11
             $\mathbf{x}^* \leftarrow \mathbf{x}^i;$                                    12
             $z_{PE}^* \leftarrow z_i^*;$                                    13
            ir al paso (6);                                           14
        elegir  $x_j^i \notin \mathbb{Z}$  y generar subproblemas  $S_{i0}$  y  $S_{i1}$  con regiones 15
            factibles  $S_i \cup \{x_j \leq \lfloor x_j^i \rfloor\}$  y  $S_i \cup \{x_j \geq \lceil x_j^i \rceil\}$ ,
            respectivamente;
             $\mathcal{L} \leftarrow \mathcal{L} \cup \{S_{i1}, S_{i2}\}.$                  16
    devolver  $(\mathbf{x}^*, z_{PE}^*)$                                          17

```

Bibliografía

- [AA95] Erling Andersen and Knud Andersen. Presolving in linear programming. *Math. Program.*, 71:221–245, 12 1995.
- [BH09] Robert F. Bodi and Katrin Herr. Symmetries in integer programs. *arXiv: Combinatorics*, 2009.
- [BV04] Stephen Boyd and Lieven Vandenberghe. *Convex Optimization*. Cambridge University Press, 2004.
- [Lav14] Carmen Gómez Laveaga. *Álgebra Superior: Curso completo*. Programa Universitario del Libro de Texto. Facultad de Ciencias, Universidad Nacional Autónoma de México, Ciudad de México, México, primera edición edition, 2014. Primera reimpresión: julio de 2015.
- [MT90] Silvano Martello and Paolo Toth. *Knapsack problems: algorithms and computer implementations*. John Wiley & Sons, Inc., USA, 1990.
- [New72] Morris Newman. *Integral Matrices*, volume 45 of *Pure and Applied Mathematics*. Academic Press, New York, 1972.
- [NW06] Jorge Nocedal and Stephen J. Wright. *Numerical Optimization*. Springer Series in Operations Research and Financial Engineering. Springer, New York, 2 edition, 2006.

- [Oli17] Fabricio Oliveira. Linear optimisation notes. <https://github.com/gamma-opt/linopt-notes>, 2017. Accessed: 2025-07-14.
- [Sch98] Alexander Schrijver. *Theory of Linear and Integer Programming*. John Wiley & Sons, Chichester, UK, 1998.