



Casa abierta al tiempo

UNIVERSIDAD AUTÓNOMA METROPOLITANA
Unidad Iztapalapa

DIVISIÓN DE CIENCIAS BÁSICAS E INGENIERÍA

Funciones convexas no diferenciables

Tesis que presenta
Rafael Alejandro Nava Manzo
Para obtener el grado de
Maestro en Ciencias
(Matemáticas Aplicadas e Industriales)

Asesoras: Dra. Shirley Thelma Bromberg Silverstein

Dra. Patricia Saavedra Barrera

Jurado calificador:

Presidente: Dra. María Cristina Gigola Paglialunga

Secretario: Dra. Patricia Saavedra Barrera

Vocales: Dr. Joaquín Delgado Fernández

Dra. Shirley Thelma Bromberg Silverstein

Dr. Miguel Ángel Gutiérrez Andrade

México D. F., 13 de febrero de 2015



Casa abierta al tiempo

UNIVERSIDAD AUTÓNOMA METROPOLITANA
Unidad Iztapalapa

DIVISIÓN DE CIENCIAS BÁSICAS E INGENIERÍA

Funciones convexas no diferenciables

Tesis que presenta
Rafael Alejandro Nava Manzo
Para obtener el grado de
Maestro en Ciencias
(Matemáticas Aplicadas e Industriales)

Asesoras: Dra. Shirley Thelma Bromberg Silverstein

Dra. Patricia Saavedra Barrera

Jurado calificador:

Presidente: Dra. María Cristina Gigola Paglialunga

Secretario: Dra. Patricia Saavedra Barrera

Vocales: Dr. Joaquín Delgado Fernández

Dra. Shirley Thelma Bromberg Silverstein

Dr. Miguel Ángel Gutiérrez Andrade

México D. F., 13 de febrero de 2015

Rafael Alejandro Nava Manzo
Posgrado en Matemáticas
Universidad Autónoma Metropolitana

Funciones convexas no diferenciables

Quando el pasado te llame no lo atiendas ...
no tiene nada nuevo que contarte.

Agradecimientos

A mis queridos padres por ser un excelente ejemplo de sacrificio y esfuerzo para mí, ustedes me inspiran a ser mejor y me dan la fuerza necesaria para afrontar mis problemas con optimismo. A ustedes les dedico el presente trabajo.

A todos mis amigos que me acompañaron a lo largo de la maestría, por la compañía y el apoyo en todo momento.

A mis asesoras la Dra. Shirley Bromberg Silverstein y a la Dra. Patricia Barrera Barrera por todo su apoyo, tiempo y atención a lo largo de la maestría y en la realización de esta tesis.

Al Consejo Nacional de Ciencia y Tecnología (CONACYT) por su apoyo y patrocinio para mis estudios de maestría, los cuales concluyen con la realización de este proyecto de tesis.

Índice general

Agradecimientos	III
Introducción	VII
1. Análisis convexo	1
1.1. Preliminares	1
1.2. Conjuntos convexos e hiperplanos	3
1.3. Funciones convexas	5
1.3.1. Derivadas direccionales	7
1.3.2. Subgradientes	8
1.3.3. El ϵ -subdiferencial	12
1.3.4. Funciones convexas diferenciables	15
1.3.5. Mínimos de una función convexa	17
1.4. Problemas con restricciones y dualidad	19
1.5. Optimización no diferenciable	22
2. Métodos para optimización no diferenciable	27
2.1. Método de subgradiente	27
2.2. Método de planos de corte	31
2.3. Método <i>bundle</i>	32
2.3.1. Regularización de Moreau Yosida	32
2.3.2. Algoritmo del método <i>bundle</i>	36
2.3.3. Convergencia del método <i>bundle</i>	40
2.4. Métodos derivados del operador próximo	44
2.4.1. Operadores no expansivos	44
2.4.2. Sucesiones Féjer monótonas	48
2.4.3. Operador próximo	50
2.4.4. Puntos fijos	52
2.4.5. Descomposición de Moreau	53

2.4.6. Método de punto próximo	54
2.4.7. Convergencia del método de punto próximo	54
2.4.8. Método de gradiente próximo	56
2.4.9. Convergencia del método de gradiente próximo	57
3. Resultados numéricos	61
3.1. Formulación por penalización exacta	61
3.2. El problema de LASSO	64
3.3. Problema de LASSO modificado	65
3.4. El problema del portafolio	67
3.5. El problema de Minimax	69
3.6. ADMM	71
4. Conclusiones	75
Bibliografía	77

Introducción

La programación lineal fue planteada como un modelo matemático desarrollado durante la Segunda Guerra Mundial a fin de reducir los costos del ejército y aumentar las pérdidas del enemigo. Los fundadores de la técnica son George Dantzig, quien publicó el método simplex en 1947, John Von Neumann, que desarrolló la teoría de la dualidad en el mismo año y Leonid Kantoróvich, un matemático ruso que utilizó técnicas similares aplicadas a la economía antes de Dantzig y que ganó el premio Nobel en economía en 1975.

La teoría moderna de optimización comenzó esencialmente con el desarrollo del método simplex. Sin embargo, los enfoques modernos de la teoría de optimización deben sus orígenes al cálculo de variaciones que ha sido estudiado por más de tres siglos y que también fue crucial en el desarrollo del análisis funcional.

El estudio de la optimización convexa se inició con la publicación de la obra ya clásica *Convex Analysis* por R.T. Rockafellar [27]. No es posible minimizar el impacto que este libro tuvo en el desarrollo de la teoría y los métodos de optimización. A la optimización convexa se le adoptó el término de *análisis convexo*, el cual fue sugerido por Albert. W. Tucker (famoso por las condiciones de Kuhn-Tucker) [8].

Las funciones no diferenciables en todas partes, marcaron un paradigma en la teoría moderna de optimización. Se encontró que muchos problemas de optimización convexa no eran diferenciables en el punto mínimo. Así, un enfoque completamente diferente se desarrolló, en donde fue concebida la noción de subdiferencial. [8][9][12][27]. El subdiferencial de una función en un punto dado es un conjunto que se puede considerar como un sustituto de la noción de la derivada en puntos donde la función no es diferenciable. La importancia del subdiferencial fue que un conjunto de reglas de cálculo pudo ser desarrollado, y que se volvió muy útil para llevar a cabo el análisis de funciones no diferenciables. Fue entonces natural extender el dominio de optimización no diferenciable mas allá de la convexidad. La primera extensión natural era la clase de funciones localmente Lipschitz y esto fue iniciado por F.H. Clarke [16], que fue el segundo estudiante de R.T. Rockafellar. Él introdujo la noción de gradiente generalizado para funciones localmente Lipschitz. Clarke también

desarrolló un cálculo para el gradiente generalizado [16].

En lo que se refiere a la teoría, ahora es habitual considerar hechos básicos del análisis convexo para toda la clase de funciones convexas, incluyendo funciones no diferenciables; el uso de subdiferenciales se ha convertido en algo tan cotidiano como el uso de gradientes, sobre todo después de la obra de R.T. Rockafellar [27], B.N. Pshenichny [25] y otros.

Existen varias fuentes principales de problemas de optimización no diferenciable, algunas de ellas son:

1. Métodos de penalización no diferenciables (en particular la función de penalización exacta)[19]. Ciertos tipos de funciones de penalización no diferenciables son superiores a las funciones de penalización diferenciables de uso común, porque por lo general no requieren coeficientes de penalización que tienden al infinito, dando soluciones exactas para valores suficientemente grandes, pero finitos, de estos coeficientes [22].
2. Problemas Minimax, que son típicos de los modelos encontrados en la teoría de juegos, modelos multicriterio de la planificación óptima y la investigación de operaciones [33].
3. La función lagrangiana. Al tratar resolver un problema con restricciones, usualmente se recurre a resolver el problema dual por medio de la función lagrangiana, la cual a menudo es un problema más sencillo de resolver que el original. Esta función puede ser no diferenciable si las funciones que definen las restricciones son, aunque convexas, no diferenciables [8][14].

Los métodos computacionales para la optimización no diferenciable se desarrollaron en dos direcciones: (a) la investigación dirigida a resolver determinados tipos de problemas de minimización con funciones no diferenciables que tienen una estructura especial y que se define de forma explícita [22]; (b) la investigación sobre la elaboración de métodos para resolver clases más generales de problemas, que no suponen de antemano el conocimiento de la estructura específica de la función a minimizar pero requieren la evaluación de la función y sus gradientes (o sus análogos en el caso no diferenciable) en cualquier punto dado.

Para el primer grupo, se cuenta con numerosos trabajos sobre métodos para resolver problemas Minimax [20][33]. Por otro lado, una serie de obras se dedican a la minimización de funciones convexas lineales a trozos [22]. Para la solución de varios problemas de optimización no diferenciable en los que se dan explícitamente los puntos donde la función no es diferenciable (por ejemplo, las funciones con valores absolutos), se han desarrollado técnicas especiales de suavizado, como es el caso de la regularización de Moreau-Yosida [9].

En lo que se refiere a métodos generales de optimización no diferenciable, se pueden distinguir dos clases básicas en los cuales se requiere el cálculo del subdiferencial.

1. Métodos de gradiente generalizado. Para minimizar funciones diferenciables, se usan con frecuencia múltiples versiones del método de gradiente, que es natural, ya que la dirección negativa del gradiente en un punto dado es una dirección de descenso. La selección de un tamaño de paso en la mayoría de estos métodos tiene por objeto disminuir de manera significativa el valor de la función objetivo en cada iteración.

Un intento de extender estos métodos para funciones no diferenciables se encuentra con al menos dos dificultades. En primer lugar, es necesario definir un análogo de gradiente en los puntos en los que no hay diferenciabilidad (subdiferencial). En segundo lugar, se tiene que desarrollar nuevas formas de elegir las direcciones de búsqueda y tamaños de paso. En particular para estos casos se usa el método de subgradiente y sus variantes muy bien estudiadas por varios autores, entre ellos F.H. Clarke [16].

2. El método de planos de corte para resolver problemas convexos.
 - El método de planos de corte de Kelley [9], basado en aproximaciones lineales a la gráfica de una función convexa por medio de hiperplanos de soporte; en cada iteración el método resuelve un problema de programación lineal con un número creciente de restricciones en cada iteración.
 - Los métodos *bundle* [9][22]. Una variante de estos métodos se tiene cuando, aparte de tener aproximaciones lineales de la función, se aplica la regularización de Moreau-Yosida a estas aproximaciones lineales.

Una tercer clase de técnica es desarrollada con el objetivo de tratar de resolver problemas generales de optimización no diferenciable utilizando el operador próximo. El operador próximo de una función convexa es una extensión natural de la noción de proyección sobre un conjunto convexo [12][14]. Esta herramienta desempeña un papel central en el análisis y la solución numérica de problemas de optimización convexa como los siguientes:

1. El problema de LASSO. El LASSO es una herramienta popular para la regresión lineal dispersa [28], especialmente para problemas en los cuales el número de variables m exceden al número de observaciones n .
2. Problemas Minimax.

3. Solución de funciones de penalización exacta.
4. Problemas de proyección sobre un conjunto convexo, en particular cuando se trabaja con norma l_1 y l_∞ .

La variedad de problemas que pueden ser tratados se incrementa cuando se utiliza el operador próximo [12]. Los métodos que utilizan el operador próximo ocupan un nivel más alto de abstracción a comparación de métodos de optimización clásicos como el método de Newton. En este último, las operaciones básicas son de bajo nivel, que consisten en operaciones de álgebra lineal y el cálculo de gradientes y hessianos. En los métodos con el operador próximo, la operación básica es evaluar el operador próximo de una función, que consiste en resolver un pequeño problema de optimización convexa. Estos subproblemas pueden resolverse con métodos estándar (como es el caso del método de restricciones activas). Algunos de los métodos que usan el operador próximo son: el método de punto próximo y el método de gradiente próximo.

Objetivos de la tesis

Estudiar diversos métodos (subgradiente, bundle y métodos con operador próximo) utilizados para resolver problemas de optimización no diferenciable, aplicarlos a una serie de problemas y comparar los resultados con otros métodos de la literatura encontrados en [9][11][32][33].

Contenido de la tesis

En el primer capítulo de este trabajo se presentan algunas propiedades de las funciones convexas así como un teorema de separación que será importante a lo largo de este trabajo. También se presentan condiciones bajo las cuáles se garantiza la existencia de un mínimo de una función convexa.

En el segundo capítulo se presentan los métodos que se utilizaron para resolver nuestro problema de interés así como sus pruebas de convergencia. Se presentan el método de subgradiente, el método *bundle*, el método de punto próximo y el método de gradiente próximo.

En el último capítulo se presentan todos los resultados numéricos obtenidos al aplicar los métodos antes mencionados.

Análisis convexo

En este capítulo se definirán algunos conceptos y notaciones importantes que se utilizarán a lo largo de este trabajo.

1.1. Preliminares

La línea real extendida $[-\infty, +\infty] = \mathbb{R} \cup \{-\infty\} \cup \{+\infty\} = \overline{\mathbb{R}}$ se obtiene adjuntando los elementos $-\infty$ y $+\infty$ a la recta real y extendiendo el orden de la siguiente manera:

$$-\infty < \epsilon < +\infty, \quad \forall \epsilon \in \mathbb{R}.$$

Las reglas aritméticas son extendidas a los elementos $\{-\infty\}$ y $\{+\infty\}$ de la manera usual, dejando las expresiones como $+\infty + (-\infty)$, $0 * (\pm\infty)$, y $+\infty / \pm\infty$ como indefinidas.

El **límite inferior** de una sucesión $\{x_n\}_{n \in \mathbb{N}}$ en $[-\infty, +\infty]$ está definido como

$$\liminf x_n = \sup_{n \in \mathbb{N}} \inf_{\substack{m \in \mathbb{N} \\ m \geq n}} x_m$$

y el **límite superior** es

$$\limsup x_n = \inf_{n \in \mathbb{N}} \sup_{\substack{m \in \mathbb{N} \\ m \geq n}} x_m.$$

Es claro que $\liminf x_n \leq \limsup x_n$.

Se trabajará con funciones cuyo contradominio está contenido en la recta real extendida.

Definición 1.1.1. *Sea $C \neq \emptyset$ un subconjunto de \mathbb{R}^n y sea $f : C \rightarrow \overline{\mathbb{R}}$. Se definen los siguientes conjuntos:*

1. El dominio esencial de f

$$\text{dom } f = \{x \in C : f(x) < +\infty\}.$$

2. La gráfica de f

$$\text{gráf } f = \{(x, y) \in \mathbb{R}^{n+1} : f(x) = y\}.$$

3. La epigráfica de f

$$\text{epi } f = \{(x, y) \in \mathbb{R}^{n+1} : f(x) \leq y\}.$$

4. El conjunto de nivel de f a la altura $\alpha \in \mathbb{R}$

$$S_\alpha = \{x \in C : f(x) \leq \alpha\}.$$

Definición 1.1.2. Sea $C \neq \emptyset$ un subconjunto de \mathbb{R}^n . La función $f : C \rightarrow \bar{\mathbb{R}}$ se dice **propia** si $-\infty \notin f(C)$ y $\text{dom } f \neq \emptyset$.

Definición 1.1.3. Sea $f : \mathbb{R}^n \rightarrow \bar{\mathbb{R}}$ una función propia y $C \subseteq \mathbb{R}^n$. El **ínfimo** de f sobre C se denota por $\inf_C f$. Se dice que f **alcanza su mínimo** sobre C si existe $x^* \in C$ tal que $f(x^*) = \inf_C f$. En este caso, se escribe $f(x^*) = \min f(C)$ o $f(x^*) = \min_{x \in C} f(x)$. De manera similar, el **supremo** de f sobre C se denota por $\sup_C f$. Se dice que f **alcanza su supremo** sobre C si existe $x^* \in C$ tal que $f(x^*) = \sup_C f$. En este caso, se escribe $f(x^*) = \max f(C)$ o $f(x^*) = \max_{x \in C} f(x)$.

Definición 1.1.1. Una función $f : \mathbb{R}^n \rightarrow \bar{\mathbb{R}}$ se dice **inferiormente (resp. superiormente) semi continua (i.s.c.)** en $x \in \mathbb{R}^n$ si

$$f(x) \leq \liminf_{i \rightarrow \infty} f(x_i), \quad (\text{resp. } f(x) \geq \limsup_{i \rightarrow \infty} f(x_i))$$

para cualquier sucesión $\{x_i\}_{i \in \mathbb{N}} \subset \mathbb{R}^n$ tal que $x_i \rightarrow x$.

Cuando una función es inferiormente semi continua y superiormente semi continua en x entonces la función es **continua** en x . El siguiente teorema caracteriza geoméricamente a las funciones que son i.s.c.

Teorema 1.1.1. Sea $f : \mathbb{R}^n \rightarrow \bar{\mathbb{R}}$. Entonces las siguientes condiciones son equivalentes:

1. f es i.s.c. en todo \mathbb{R}^n .
2. S_α es cerrado para todo $\alpha \in \mathbb{R}$.
3. $\text{epi } f$ es un conjunto cerrado en \mathbb{R}^{n+1}

La demostración puede encontrarse en [27].

Por el teorema anterior a las funciones i.s.c. también se les llama **funciones cerradas**.

1.2. Conjuntos convexos e hiperplanos

Definición 1.2.1. Sea $S \subset \mathbb{R}^n$ un conjunto no vacío. Se dice que S es un **conjunto convexo** si para cualquier $x, y \in S$, $\lambda x + (1 - \lambda)y \in S$ para todo $\lambda \in (0, 1)$.

Definición 1.2.2. La suma vectorial

$$\lambda_1 x_1 + \dots + \lambda_m x_m$$

es llamada **combinación convexa** de x_1, \dots, x_m si los coeficientes λ_i son no negativos y $\sum \lambda_i = 1$.

Utilizando esta definición se pueden caracterizar los conjuntos convexos por medio del siguiente resultado.

Proposición 1.2.1. Un subconjunto de \mathbb{R}^n es convexo si y solo si contiene todas las combinaciones convexas de sus elementos.

Demostración. Se realiza fácilmente por inducción. □

Definición 1.2.3. Se dice que $f : \mathbb{R}^m \rightarrow \mathbb{R}^n$ es una **función afín** si existe una transformación lineal $A : \mathbb{R}^m \rightarrow \mathbb{R}^n$ y un vector $c \in \mathbb{R}^n$ tales que

$$f(x) = Ax + c.$$

Si la transformación lineal A es definida por $A(x) = \langle x, b \rangle$ donde $b \in \mathbb{R}^n$ y \langle, \rangle denota al producto interior en \mathbb{R}^n , se tiene que la función $A(x) + \beta = 0$ ($\beta \in \mathbb{R}$) es un conjunto de nivel de la función afín $f(x) = Ax - \beta$. Esta función da lugar a la definición del siguiente conjunto.

Definición 1.2.4. Un **hiperplano** $H \subset \mathbb{R}^n$ es el conjunto

$$H = \{x \in \mathbb{R}^n : \langle x, b \rangle = \beta\},$$

donde b es un vector en \mathbb{R}^n diferente de cero y $\beta \in \mathbb{R}$.

El hiperplano define dos semi espacios convexos $H^+ = \{x \in \mathbb{R}^n : \langle x, b \rangle \geq \beta\}$ y $H^- = \{x \in \mathbb{R}^n : \langle x, b \rangle \leq \beta\}$. Como observación, se tiene que un hiperplano es un conjunto convexo.

Definición 1.2.5. Sea S un subconjunto de \mathbb{R}^n y $x \in S$, un hiperplano en \mathbb{R}^n se dice **hiperplano de soporte de S en x** si se cumple que $x \in S \cap H$ y $S \subset H^+$ o $S \subset H^-$.

Ahora se presentarán dos teoremas de separación muy importantes.

Teorema 1.2.1 (Teorema de separación de Hahn-Banach). *Sea S un subconjunto convexo y cerrado de \mathbb{R}^n y $z \notin S$. Entonces existe un vector $\xi \in \mathbb{R}^n$ y un escalar α tal que $\langle z, \xi \rangle > \alpha$ y $\langle x, \xi \rangle \leq \alpha$ para cada $x \in S$, es decir, existe un hiperplano*

$$H = \{x \in \mathbb{R}^n : \langle x, \xi \rangle = \alpha\}$$

tal que $z \in H^+$ y $S \subset H^-$.

La demostración puede encontrarse en [8].

Teorema 1.2.2. *Sean S_1 y S_2 dos subconjuntos disjuntos, convexos y no vacíos de \mathbb{R}^n . Entonces existe un hiperplano que separa S_1 y S_2 ; esto es, existe un vector $\xi \in \mathbb{R}^n$ diferente de cero tal que*

$$\inf\{\langle \xi, x \rangle : x \in S_1\} \geq \sup\{\langle \xi, x \rangle : x \in S_2\}$$

La demostración puede encontrarse en [8].

Corolario 1.2.1. *Un subconjunto convexo, cerrado de \mathbb{R}^n es la intersección de todos los semi espacios que lo contienen.*

Demostración. Sea S un subconjunto convexo y cerrado de \mathbb{R}^n . Claramente S está contenido en la intersección de todos los semi espacios que lo contienen. Por otro lado, supongamos por contradicción que existe un elemento y en la intersección de todos los semi espacios que contienen a S pero $y \notin S$. Por el teorema de separación de Hahn-Banach existe un semi espacio que contiene a S pero no a y . Esta contradicción muestra el corolario. \square

Corolario 1.2.2. *Sea S un subconjunto cerrado y convexo de \mathbb{R}^n y sea z un elemento de la frontera de S . Entonces existe un hiperplano de soporte de S en z .*

Demostración. Como z está en la frontera de S existe una sucesión $\{z_n\}_{n \in \mathbb{N}}$ que cumple que $z_n \notin S$ para toda $n \in \mathbb{N}$ y $z_n \rightarrow z$. Por el teorema de separación de Hahn-Banach, para cada $n \in \mathbb{N}$ existe un vector ξ_n unitario tal que $\langle \xi_n, z_n \rangle > \langle \xi_n, x \rangle$ para todo $x \in S$. Dado que $\{\xi_n\}$ es acotada tiene una subsucesión convergente $\{\xi_{n_i}\}$ con límite el vector unitario ξ . Al considerar esta subsucesión, tomando $n_i \rightarrow \infty$ se obtiene

$$\langle \xi, z \rangle \geq \langle \xi, x \rangle$$

para todo $x \in S$. \square

1.3. Funciones convexas

Definición 1.3.1. Sea $f : S \rightarrow (-\infty, +\infty]$, donde S es un subconjunto convexo no vacío de \mathbb{R}^n . Se dice que f es una **función convexa** en S si

$$f(\lambda x + (1 - \lambda)y) \leq \lambda f(x) + (1 - \lambda)f(y) \quad (1.1)$$

para cada $x, y \in S$ y $\lambda \in (0, 1)$.

Si la desigualdad anterior es estricta entonces se dice que la función f es estrictamente convexa. La función $f : S \rightarrow (-\infty, +\infty]$ es **cóncava** si $-f$ es una función convexa.

En lo que sigue, el trabajo se concentrará en el estudio de las funciones convexas ya que los resultados para funciones cóncavas pueden ser deducidos fácilmente notando solamente que f es cóncava si y solo si $-f$ es convexa.

A continuación se dará una caracterización de las funciones convexas utilizando sus epigráficas.

Proposición 1.3.1. Sea S un conjunto convexo no vacío en \mathbb{R}^n y sea $f : S \rightarrow (-\infty, +\infty]$. Entonces f es convexa si y solo si $\text{epi } f$ es un conjunto convexo.

Demostración. Sea f una función convexa y sean $(x_1, y_1), (x_2, y_2) \in \text{epi } f$. Dado $\lambda \in (0, 1)$, entonces

$$\lambda y_1 + (1 - \lambda)y_2 \geq \lambda f(x_1) + (1 - \lambda)f(x_2) \geq f(\lambda x_1 + (1 - \lambda)x_2).$$

Como $\lambda x_1 + (1 - \lambda)x_2 \in S$, se sigue que $(\lambda x_1 + (1 - \lambda)x_2, \lambda y_1 + (1 - \lambda)y_2) \in \text{epi } f$. Recíprocamente, sean $(x_1, f(x_1)), (x_2, f(x_2)) \in \text{epi } f$, entonces por la convexidad de $\text{epi } f$, se tiene que

$$[\lambda x_1 + (1 - \lambda)x_2, \lambda f(x_1) + (1 - \lambda)f(x_2)] \in \text{epi } f,$$

es decir

$$f(\lambda x_1 + (1 - \lambda)x_2) \leq \lambda f(x_1) + (1 - \lambda)f(x_2)$$

para cada $\lambda \in (0, 1)$. Por lo tanto, f es una función convexa. \square

Una propiedad importante de las funciones convexas es que éstas son continuas en el interior de su dominio. Para probar esta aseveración se presenta el siguiente lema.

Lema 1.3.1. Sea $f : S \rightarrow (-\infty, +\infty]$, donde S es un subconjunto abierto y convexo de \mathbb{R}^n . Entonces f es localmente acotada.

Demostración [24]. Para $a \in S$ fijo, sea un cubo K en S , centrado en a , con vértices v_1, \dots, v_{2^n} . Claramente K es una vecindad de a . Cualquier $x \in K$ es una combinación convexa de vértices, entonces

$$f(x) = f\left(\sum_k \lambda_k v_k\right) \leq M = \sup_k f(v_k).$$

Luego, por la simetría de K para cada $x \in \text{int } K$ existe $y \in K$ tal que $a = (x + y)/2$, entonces

$$f(x) \geq 2f(a) - f(y) \geq 2f(a) - M.$$

Por lo tanto, f es acotada en K . □

Teorema 1.3.1. Sea $f : S \rightarrow (-\infty, +\infty]$ una función convexa, donde S es un subconjunto convexo y abierto de \mathbb{R}^n . Entonces f es localmente Lipschitz. En particular, f es continua.

Demostración [24]. De acuerdo con el lema anterior, se puede encontrar una bola $B_{2r} \subseteq S$ en la cual f es acotada por algún M . Ahora, sean $x, y \in B_r$ tales que $x \neq y$, y sea $z = y + (r/\alpha)(y - x)$, donde $\alpha = \|y - x\|$. Claramente, $z \in B_{2r}$. Como

$$y = \frac{r}{r + \alpha}x + \frac{\alpha}{r + \alpha}z,$$

de la convexidad de f se deduce que

$$f(y) \leq \frac{r}{r + \alpha}f(x) + \frac{\alpha}{r + \alpha}f(z).$$

Entonces

$$\begin{aligned} f(y) - f(x) &\leq \frac{\alpha}{r + \alpha}(f(z) - f(x)) \\ &\leq \frac{\alpha}{r}(f(z) - f(x)) \leq \frac{2M}{r}\|y - x\|, \end{aligned}$$

y la prueba termina intercambiando los roles de y y x . □

1.3.1. Derivadas direccionales

Definición 1.3.1. Sea S un subconjunto convexo no vacío de \mathbb{R}^n y sea $f : S \rightarrow \mathbb{R}$. Sea $\bar{x} \in S$ y d un vector diferente de cero tal que $\bar{x} + \lambda d \in S$ con $\lambda > 0$ suficientemente pequeño. La **derivada direccional** de f en \bar{x} en la dirección d , denotada por $f'(\bar{x}; d)$, está dada por el siguiente límite (cuando existe)

$$f'(\bar{x}; d) = \lim_{\lambda \rightarrow 0^+} \frac{f(\bar{x} + \lambda d) - f(\bar{x})}{\lambda}. \quad (1.2)$$

Este límite existe y es finito cuando f es una función convexa y $\bar{x} \in \text{int } S$ como se probará en el siguiente lema.

Lema 1.3.2. Sea $f : S \rightarrow \mathbb{R}$ una función convexa, donde S es un subconjunto convexo no vacío de \mathbb{R}^n . Si $\bar{x} \in \text{int } S$ entonces la derivada direccional $f'(\bar{x}; d)$ existe para cualquier dirección d diferente de cero.

Demostración [8]. Sean $\lambda_2 > \lambda_1 > 0$ tal que $\bar{x} + \lambda_2 d \in \text{int } S$. Considerando la convexidad de f tenemos

$$\begin{aligned} f(\bar{x} + \lambda_1 d) &= f \left[\frac{\lambda_1}{\lambda_2} (\bar{x} + \lambda_2 d) + \left(1 - \frac{\lambda_1}{\lambda_2} \right) \bar{x} \right] \\ &\leq \frac{\lambda_1}{\lambda_2} f(\bar{x} + \lambda_2 d) + \left(1 - \frac{\lambda_1}{\lambda_2} \right) f(\bar{x}). \end{aligned}$$

Esta desigualdad implica que

$$\frac{f(\bar{x} + \lambda_1 d) - f(\bar{x})}{\lambda_1} \leq \frac{f(\bar{x} + \lambda_2 d) - f(\bar{x})}{\lambda_2}.$$

Por lo tanto, el cociente $(f(\bar{x} + \lambda d) - f(\bar{x}))/\lambda$ es monótono decreciente cuando $\lambda \rightarrow 0^+$. Ahora, dado $\lambda \geq 0$ suficientemente pequeño, se tiene por la convexidad de f lo siguiente

$$\begin{aligned} f(\bar{x}) &= f \left[\frac{\lambda}{1 + \lambda} (\bar{x} - d) + \frac{1}{1 + \lambda} (\bar{x} + \lambda d) \right] \\ &\leq \frac{\lambda}{1 + \lambda} f(\bar{x} - d) + \frac{1}{1 + \lambda} f(\bar{x} + \lambda d). \end{aligned}$$

Así

$$\frac{f(\bar{x} + \lambda d) - f(\bar{x})}{\lambda} \geq f(\bar{x}) - f(\bar{x} - d).$$

Por lo tanto, dado que el cociente anterior está acotado inferiormente se tiene que

$$\lim_{\lambda \rightarrow 0^+} \frac{f(\bar{x} + \lambda d) - f(\bar{x})}{\lambda} = \inf_{\lambda > 0} \frac{f(\bar{x} + \lambda d) - f(\bar{x})}{\lambda}.$$

□

1.3.2. Subgradientes

La epigráfica de una función convexa es un conjunto convexo y cerrado en \mathbb{R}^{n+1} . Por lo tanto, debe tener hiperplanos de soporte en sus puntos frontera. Estos hiperplanos de soporte conducen a la noción de subgradiente, el cual es definido a continuación.

Definición 1.3.1. *Sea S un conjunto convexo de \mathbb{R}^n , y sea $f : S \rightarrow (-\infty, +\infty]$ una función convexa. Entonces ξ es llamado **subgradiente** de f en x si*

$$f(y) \geq f(x) + \langle \xi, y - x \rangle \quad (1.3)$$

para todo $y \in S$.

El subgradiente ξ define un hiperplano que soporta al conjunto epi f en el punto frontera $(x, f(x))$, es decir, si H es el hiperplano definido por la función afín $h(y) = f(x) + \langle \xi, y - x \rangle$ entonces epi $f \subseteq$ epi h .

El conjunto de todos los subgradientes en x es conocido como el **subdiferencial** de x , es decir, el subdiferencial de f en x está definido por el conjunto

$$\partial f(x) := \{ \xi : f(y) \geq f(x) + \langle \xi, y - x \rangle, \forall y \in S \}. \quad (1.4)$$

Es sencillo verificar que $\partial f(x)$ es cerrado y convexo. En efecto, dados $\xi_1, \xi_2 \in \partial f(x)$, $y \in S$ y $d_1, d_2 \in (0, 1)$ tales que $d_1 + d_2 = 1$, entonces

$$\begin{aligned} f(x) + \langle d_1 \xi_1 + d_2 \xi_2, y - x \rangle &= d_1 f(x) + d_1 \langle \xi_1, y - x \rangle + d_2 f(x) + d_2 \langle \xi_2, y - x \rangle \\ &\leq d_1 f(y) + d_2 f(y) = f(y). \end{aligned}$$

Por lo tanto, el conjunto $\partial f(x)$ es convexo. Para verificar que $\partial f(x)$ es cerrado se toma el hecho de que el producto punto es una función continua y de aquí se sigue inmediatamente de la definición de subgradiente que $\partial f(x)$ es cerrado. A continuación se demuestra que este conjunto también es acotado.

Lema 1.3.1. *Sea f una función convexa y propia. Sea $x \in \text{int}(\text{dom } f)$ y $r > 0$ tal que $B_r(x) \subseteq \text{int}(\text{dom } f)$. Si f es Lipschitz con constante L en $B_r(x)$ entonces $\partial f(x)$ es acotado.*

Demostración. Sea $\xi \in \partial f(x)$ con $\xi \neq 0$ y sea $\gamma > 0$ tal que $x + \gamma\xi \in B_r(x)$, entonces

$$f(x + \gamma\xi) \geq f(x) + \langle \xi, \gamma\xi \rangle = f(x) + \gamma\|\xi\|^2.$$

Tomando en cuenta la ecuación anterior y que f es Lipschitz se deduce que

$$\gamma L\|\xi\| \geq f(x + \gamma\xi) - f(x) \geq \gamma\|\xi\|^2.$$

Por lo tanto

$$\|\xi\| \leq L.$$

□

Como una consecuencia directa de este lema se deduce que $\partial f(x)$ es **compacto**.

El siguiente teorema muestra que cualquier función convexa tiene al menos un subgradiente en cada uno de los puntos del interior de su dominio.

Teorema 1.3.2. *Sea S un subconjunto convexo de \mathbb{R}^n , y sea $f : S \rightarrow \mathbb{R}$ convexa. Entonces para $\bar{x} \in \text{int } S$ se tiene que $\partial f(\bar{x}) \neq \emptyset$.*

Demostración [8]. Dado que f es convexa, entonces $\text{epi } f$ es un conjunto convexo. El punto $(\bar{x}, f(\bar{x}))$ está en la frontera de $\text{epi } f$. Por el teorema 1.2.1 existe un vector diferente de cero $(\xi_0, \mu) \in \mathbb{R}^n \times \mathbb{R}$ tal que

$$\langle \xi_0, x - \bar{x} \rangle + \mu(y - f(\bar{x})) \leq 0 \quad (1.5)$$

para todo $(x, y) \in \text{epi } f$. Notemos primero que μ no es positivo, porque de otra manera, la desigualdad anterior no sería cierta para y suficientemente grande.

Ahora se probará que $\mu < 0$. Supongamos que $\mu = 0$, entonces $\langle \xi_0, x - \bar{x} \rangle \leq 0$ para toda $x \in S$. Dado que $\bar{x} \in \text{int } S$, existe un $\lambda > 0$ tal que $\bar{x} + \lambda\xi_0 \in S$ y por lo tanto $\lambda\langle \xi_0, \xi_0 \rangle \leq 0$. Esto implica que $\xi_0 = 0$ y $(\xi_0, \mu) = (0, 0)$, contradiciendo el hecho de que (ξ_0, μ) es un vector diferente de cero. Así, $\mu < 0$. Denotando $\xi_0/|\mu|$ por ξ y dividiendo la desigualdad (1.5) por $|\mu|$ se obtiene

$$\langle \xi, x - \bar{x} \rangle - y + f(\bar{x}) \leq 0,$$

para todo $(x, y) \in \text{epi } f$. Haciendo $y = f(\bar{x})$ en la ecuación anterior se obtiene $f(x) \geq f(\bar{x}) + \langle \xi, x - \bar{x} \rangle$ para toda $x \in S$. □

Corolario 1.3.1. *Sea S un conjunto convexo no vacío en \mathbb{R}^n , y sea $f : S \rightarrow \mathbb{R}$ una función estrictamente convexa. Entonces, para $\bar{x} \in \text{int } S$ existe un vector ξ tal que*

$$f(x) > f(\bar{x}) + \langle \xi, x - \bar{x} \rangle \quad (1.6)$$

para todo $x \in S$ y $x \neq \bar{x}$.

Demostración [8]. Por el teorema anterior, existe un vector ξ tal que

$$f(x) \geq f(\bar{x}) + \langle \xi, x - \bar{x} \rangle, \quad (1.7)$$

para todo $x \in S$. Se hará esta prueba por contradicción, supongamos que existe un $x_0 \neq \bar{x}$ tal que $f(x_0) = f(\bar{x}) + \langle \xi, x_0 - \bar{x} \rangle$. Entonces, por la convexidad de f para $\lambda \in (0, 1)$

$$\begin{aligned} f(\lambda\bar{x} + (1 - \lambda)x_0) &< \lambda f(\bar{x}) + (1 - \lambda)f(x_0) \\ &= f(\bar{x}) + (1 - \lambda)\langle \xi, x_0 - \bar{x} \rangle. \end{aligned} \quad (1.8)$$

Haciendo $x = \lambda\bar{x} + (1 - \lambda)x_0$ en (1.7) se obtiene

$$f(\lambda\bar{x} + (1 - \lambda)x_0) \geq f(\bar{x}) + (1 - \lambda)\langle \xi, x_0 - \bar{x} \rangle$$

contradiciendo (1.8). □

El recíproco del corolario anterior no es cierto en general, es decir, si para cada $\bar{x} \in \text{int } S$ existe el subgradiente entonces no necesariamente la función es convexa en S . Sin embargo, aunque la función no necesariamente es convexa en S , el teorema siguiente garantiza que si todo punto en el interior de S tiene subgradiente, entonces la función es convexa en el interior de S .

Proposición 1.3.2. *Sea S un conjunto convexo de \mathbb{R}^n , y sea $f : S \rightarrow \mathbb{R}$. Supongamos que para cada punto $\bar{x} \in \text{int } S$ existe ξ tal que*

$$f(x) \geq f(\bar{x}) + \langle \xi, x - \bar{x} \rangle \quad (1.9)$$

para cada $x \in S$. Entonces, f es convexa en $\text{int } S$.

Demostración [8]. Sean $x_1, x_2 \in \text{int } S$ y $\lambda \in (0, 1)$. Se debe demostrar que $\lambda x_1 + (1 - \lambda)x_2 \in \text{int } S$. Por hipótesis, existe un subgradiente ξ de f en $\lambda x_1 + (1 - \lambda)x_2$. En particular, las siguientes desigualdades se cumplen

$$\begin{aligned} f(x_1) &\geq f(\lambda x_1 + (1 - \lambda)x_2) + (1 - \lambda)\langle \xi, x_1 - x_2 \rangle \\ f(x_2) &\geq f(\lambda x_1 + (1 - \lambda)x_2) + (1 - \lambda)\langle \xi, x_2 - x_1 \rangle. \end{aligned}$$

Multiplicando las dos desigualdades por λ y $(1 - \lambda)$ respectivamente y sumando se obtiene

$$\lambda f(x_1) + (1 - \lambda)f(x_2) \geq f(\lambda x_1 + (1 - \lambda)x_2).$$

□

Lema 1.3.3. *Sea $\{x_n\}_{n \in \mathbb{N}} \subset S$ una sucesión tal que $x_n \rightarrow x$ donde $x \in S$ y sea $\{\xi_n\}_{n \in \mathbb{N}}$ una sucesión tal que $\xi_n \in \partial f(x_n)$ para todo $n \in \mathbb{N}$. Entonces cualquier punto de acumulación de $\{\xi_n\}$ está en $\partial f(x)$.*

Demostración [19]. Para cualquier $y \in S$, se cumple por la definición de subgradiente que

$$f(y) \geq f(x_n) + \langle \xi_n, y - x_n \rangle.$$

Sea ξ un punto límite de la sucesión $\{\xi_n\}$, entonces existe una subsucesión $\{\xi_{n_i}\}$ tal que $\xi_{n_i} \rightarrow \xi$. Entonces, $x_{n_i} \rightarrow x$ y por la continuidad de f y del producto interior se deduce que

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

para toda $y \in S$. Por lo tanto, $\xi \in \partial f(x)$. \square

Se puede obtener una caracterización de la derivada direccional utilizando el subgradiente. Esta afirmación se concreta en el teorema siguiente.

Teorema 1.3.3. *Sea $f : \mathbb{R}^n \rightarrow (-\infty, +\infty]$ una función cerrada y convexa. Si $x \in \text{int}(\text{dom } f)$ entonces*

$$f'(x; d) = \text{máx}\{\langle \xi, d \rangle : \xi \in \partial f(x)\}.$$

Demostración [19]. Sea $\xi_n \in \partial f(x + \frac{1}{n}d)$, donde $d \in \mathbb{R}^n$ y $n \in \mathbb{N}$ tal que $x + \frac{1}{n}d \in \text{int}(\text{dom } f)$. Entonces

$$f(x) \geq f(x + \frac{1}{n}d) - \frac{1}{n}\langle \xi_n, d \rangle$$

y

$$f(x + \frac{1}{n}d) \geq f(x) + \frac{1}{n}\langle \xi, d \rangle, \quad \forall \xi \in \partial f(x).$$

Por lo tanto

$$\langle \xi_n, d \rangle \geq \frac{f(x + \frac{1}{n}d) - f(x)}{\frac{1}{n}} \geq \text{máx}_{\xi \in \partial f(x)} \langle \xi, d \rangle.$$

Dado que $\partial f(x + \frac{1}{n}d)$ es acotada en una vecindad de x , existe una subsucesión $\{\xi_{n_i}\}_{n_i \in \mathbb{N}}$ tal que $\xi_{n_i} \rightarrow \xi$, donde $\xi \in \partial f(x)$ por el lema 1.3.3. Tomando el límite cuando $n_i \rightarrow \infty$, se tiene

$$f'(x; d) = \text{máx}_{\xi \in \partial f(x)} \langle \xi, d \rangle.$$

\square

1.3.3. El ϵ -subdiferencial

En esta sección se construirá una aproximación del subdiferencial.

Definición 1.3.2. Para cualquier $\epsilon > 0$, el ϵ -subdiferencial de una función convexa $f : \mathbb{R}^n \rightarrow (-\infty, +\infty]$ en un punto x es el conjunto

$$\partial_\epsilon f(x) = \{\xi \in \mathbb{R}^n : f(y) \geq f(x) + \langle \xi, y - x \rangle - \epsilon, \forall y \in \mathbb{R}^n\}.$$

Algunas propiedades del ϵ -subdiferencial son:

1. $\partial f(x) \subseteq \partial_\epsilon f(x)$. Más aún, $\partial f(x) = \bigcap_{\epsilon > 0} \partial_\epsilon f(x)$.
2. $\partial_\epsilon f(x)$ es un conjunto cerrado y convexo.

Lema 1.3.4. Sean $x, x' \in \text{dom}(f)$ y $\xi' \in \partial f(x')$. Entonces,

$$s' \in \partial_\epsilon f(x) \iff f(x') \geq f(x) + \langle \xi', x' - x \rangle - \epsilon.$$

Demostración [12]. La necesidad se sigue usando $y = x'$ en la definición de $\partial_\epsilon f(x)$.

Para la suficiencia, dado que $\xi' \in \partial f(x')$

$$\begin{aligned} f(y) &\geq f(x') + \langle \xi', y - x' \rangle \\ &= f(x) + \langle \xi', y - x \rangle + [f(x') - f(x) + \langle \xi', x - x' \rangle] \\ &\geq f(x) + \langle \xi', y - x \rangle - \epsilon, \end{aligned}$$

donde la última desigualdad se sigue de la hipótesis. Por lo tanto, $x' \in \partial_\epsilon f(x)$. \square

Ahora se probará una forma de continuidad del ϵ -subdiferencial. Para hacer esto primero se darán algunos resultados previos.

Se denotará por $\mathcal{B}(\mathbb{R}^n)$ a la familia de conjuntos borelianos de \mathbb{R}^n .

Definición 1.3.3. Sea $M : \mathbb{R}^n \rightarrow \mathcal{B}(\mathbb{R}^n)$ una función.

1. Se dice que M es **superiormente semi continua** si para cualquier sucesión convergente $(x_k, s_k) \rightarrow (x, s)$ tal que $s_k \in M(x_k)$ para cualquier $k \in \mathbb{N}$, se tiene que $s \in M(x)$.
2. M se dice **inferiormente semi continua** si para cada $s \in M(x)$ y cualquier sucesión $x_k \rightarrow x$, existe una sucesión $\{s_k\}_{k \geq 1}$ tal que $s_k \in M(x_k)$ y $s_k \rightarrow s$.
3. Se dice que M es **continua** si es inferior y superiormente semi continua.

Lema 1.3.5. *Sea f una función convexa y propia. Sea $x \in \text{int}(\text{dom } f)$ y $\delta > 0$ tal que $B_\delta(x) \subseteq \text{int}(\text{dom } f)$. Si f es Lipschitz con constante L en $B_\delta(x)$, entonces para cualquier $\delta' < \delta$, $w \in B_{\delta'}(x)$ y $\xi \in \partial_\epsilon f(w)$ se tiene que*

$$\|\xi\| \leq L + \frac{\epsilon}{\delta - \delta'}.$$

Demostración [12]. Supongamos que $\xi \neq 0$ y $z = w + (\delta - \delta') \frac{\xi}{\|\xi\|}$. Como $\xi \in \partial_\epsilon f(w)$ se tiene que

$$f(z) \geq f(w) + \langle \xi, z - w \rangle - \epsilon.$$

Dado que $z, w \in B_\delta(x)$ y además f es Lipschitz con constante L en $B_\delta(x)$ se obtiene

$$L\|z - w\| \geq f(z) - f(w) \geq \langle \xi, z - w \rangle - \epsilon.$$

Ahora, usando que $\|z - w\| = \|(\delta - \delta') \frac{\xi}{\|\xi\|}\| = \delta - \delta'$ y la desigualdad anterior se tiene

$$(\delta - \delta')L \geq \langle \xi, (\delta - \delta') \frac{\xi}{\|\xi\|} \rangle - \epsilon = \|\xi\|(\delta - \delta') - \epsilon.$$

□

Como consecuencia de este lema, se tiene que el conjunto $\partial_\epsilon f(x)$ es acotado. Además, por la continuidad del producto punto, $\partial_\epsilon f(x)$ es cerrado. Por lo tanto, $\partial_\epsilon f(x)$ es un conjunto compacto de igual forma que $\partial f(x)$.

Teorema 1.3.4. *Sea $f : \mathbb{R}^n \rightarrow \mathbb{R}$ una función convexa y Lipschitz con constante L en \mathbb{R}^n . Entonces, existe una constante $K > 0$ tal que para toda $x, x' \in \mathbb{R}^n$, $\epsilon, \epsilon' > 0$, $\xi \in \partial_\epsilon f(x)$, existe $\xi' \in \partial_{\epsilon'} f(x')$ tal que*

$$\|\xi - \xi'\| \leq \frac{K}{\min\{\epsilon, \epsilon'\}} (\|x - x'\| + |\epsilon - \epsilon'|).$$

Demostración [12]. Dado que $\partial_\epsilon f(x)$ y $\partial_{\epsilon'} f(x')$ son conjuntos convexos (de otra manera, se podría separar $\xi \in \partial_\epsilon f(x)$ de $\partial_{\epsilon'} f(x')$ por medio de un hiperplano), será suficiente probar que

$$|\max\{\langle \xi, d \rangle : \xi \in \partial_\epsilon f(x)\} - \max\{\langle \xi', d \rangle : \xi' \in \partial_{\epsilon'} f(x')\}| \leq \frac{K}{\min\{\epsilon, \epsilon'\}} (\|x - x'\| + |\epsilon - \epsilon'|).$$

Sea d un vector unitario, se define la ϵ -derivada direccional de f en x en dirección d como

$$f'_\epsilon(x; d) = \inf_{t>0} \frac{f(x + td) - f(x) + \epsilon}{t}. \quad (1.10)$$

El cociente del lado derecho se denotará como

$$q_\epsilon(x, t, d) = \frac{f(x + td) - f(x) + \epsilon}{t}.$$

El teorema 1.3.3 se puede adaptar del tal manera que se obtenga lo siguiente

$$f'_\epsilon(x; d) = \text{máx}\{\langle s, d \rangle : s \in \partial_\epsilon f(x)\}.$$

Luego, para cualquier $\nu > 0$, existe $t_\nu > 0$ tal que

$$q_\epsilon(x, t_\nu, d) < f'_\epsilon(x; d) + \nu, \quad (1.11)$$

usando el lema anterior haciendo $\delta \rightarrow \infty$, tenemos que $f'_\epsilon(x; d) \leq L$, y así $q_\epsilon(x, t) \leq L + \nu$. Por otro lado

$$q_\epsilon(x, t_\nu, d) = \frac{f(x + t_\nu d) - f(x) + \epsilon}{t_\nu} \geq -L + \frac{\epsilon}{t_\nu},$$

entonces

$$\frac{1}{t_\nu} \leq \frac{2L + \nu}{\epsilon}. \quad (1.12)$$

Así, usando las ecuaciones 1.10 y 1.11 se obtiene que

$$\begin{aligned} f'_{\epsilon'}(x'; d) - f'_\epsilon(x; d) - \nu &\leq q_{\epsilon'}(x', t_\nu, d) - q_\epsilon(x; t_\nu, d) \\ &= \frac{f(x' + t_\nu d) - f(x + t_\nu d) + f(x) - f(x') + \epsilon' - \epsilon}{t_\nu} \\ &\leq \frac{2L\|x - x'\| + |\epsilon - \epsilon'|}{t_\nu} \\ &\leq \frac{2L + \nu}{\epsilon}(2L\|x - x'\| + |\epsilon - \epsilon'|), \end{aligned}$$

donde la última desigualdad es consecuencia de 1.12. Dado que ν es arbitrario y se pueden intercambiar x y x' entonces

$$|f'_{\epsilon'}(x'; d) - f'_\epsilon(x; d)| \leq \frac{2L}{\min\{\epsilon, \epsilon'\}}(2L\|x - x'\| + |\epsilon - \epsilon'|).$$

□

Ahora se puede probar el resultado deseado.

Teorema 1.3.5. *Sea $f : \mathbb{R}^n \rightarrow \mathbb{R}$ una función convexa e i.s.c., entonces la correspondencia $M : \mathbb{R}^n \times \mathbb{R}_+ \rightarrow \mathcal{B}(\mathbb{R}^n)$ definida por*

$$(x, \epsilon) \rightarrow \partial_\epsilon f(x)$$

es continua.

Demostración [12]. El teorema anterior prueba que la aplicación M es inferiormente semi continua para cualquier $\epsilon > 0$ y cualquier $x \in \mathbb{R}^n$.

Ahora, para demostrar que M es superiormente semi continua tomemos la sucesión $(x^k, s^k, \epsilon_k) \rightarrow (\bar{x}, \bar{s}, \bar{\epsilon})$, entonces para toda $y \in \mathbb{R}^n$

$$f(y) \geq f(x^k) + \langle s^k, y - x^k \rangle - \epsilon_k.$$

Dado que f es i.s.c., tomando el límite $k \rightarrow \infty$,

$$f(y) \geq f(\bar{x}) + \langle \bar{s}, y - \bar{x} \rangle - \bar{\epsilon}$$

para toda $y \in \mathbb{R}^n$. Así, $\bar{s} \in \partial_\epsilon f(\bar{x})$ y M es superiormente semi continua. \square

1.3.4. Funciones convexas diferenciables

Ahora se trabajará con funciones convexas diferenciables y lo que esto implica. Primero se considera la siguiente definición.

Definición 1.3.4. *Sea S un conjunto convexo diferente de vacío en \mathbb{R}^n y sea $f : S \rightarrow \mathbb{R}$. Se dice que f es diferenciable en $\bar{x} \in \text{int } S$ si existe $v \in \mathbb{R}^n$ tal que*

$$\lim_{x \rightarrow \bar{x}} \frac{f(x) - f(\bar{x}) - \langle v, x - \bar{x} \rangle}{\|x - \bar{x}\|} = 0.$$

En este caso v es único y se denota por $v = \nabla f(\bar{x})$.

El siguiente resultado muestra que el subdiferencial de una función es una generalización del gradiente.

Lema 1.3.6. *Sea S un conjunto convexo abierto diferente de vacío en \mathbb{R}^n , y sea $f : S \rightarrow \mathbb{R}$ convexa. Supongamos que f es diferenciable en $\bar{x} \in \text{int } S$. Entonces $\partial f(\bar{x}) = \{\nabla f(\bar{x})\}$.*

Demostración. Como f es convexa entonces $\partial f(\bar{x}) \neq \emptyset$. Sea $\xi \in \partial f(\bar{x})$, $d \in \mathbb{R}^n$ y $\lambda > 0$ suficientemente pequeño. Por definición de subgradiente

$$f(\bar{x} + \lambda d) > f(\bar{x}) + \lambda \langle \xi, d \rangle.$$

Así,

$$\frac{f(\bar{x} + \lambda d) - f(\bar{x})}{\lambda} \geq \langle \xi, d \rangle.$$

Luego, haciendo $\lambda \rightarrow 0$, el lado izquierdo converge dado que f es diferenciable. Por lo tanto

$$\langle \nabla f(\bar{x}), d \rangle \geq \langle s, d \rangle$$

para toda $d \in \mathbb{R}^n$. En consecuencia, $\nabla f(\bar{x}) = \xi$. □

En vista del lema anterior, se da la siguiente caracterización de convexidad para funciones que son diferenciables.

Teorema 1.3.6. *Sea S un conjunto abierto y convexo en \mathbb{R}^n , y sea $f : S \rightarrow \mathbb{R}$ diferenciable en S . Entonces f es convexa si y solo si para cualquier $\bar{x} \in S$ se satisface*

$$f(x) \geq f(\bar{x}) + \langle \nabla f(\bar{x}), x - \bar{x} \rangle \tag{1.13}$$

para cada $x \in S$.

Demostración. La demostración de suficiencia es directa del lema anterior. Para la necesidad se aplica la proposición 1.3.2. □

El siguiente teorema, igual que el anterior, proporciona una caracterización de las funciones convexas que son diferenciables.

Teorema 1.3.7. *Sea S un subconjunto abierto y convexo de \mathbb{R}^n y sea $f : S \rightarrow \mathbb{R}$ diferenciable en S . Entonces f es convexa si y solo si para cada $x_1, x_2 \in S$*

$$\langle \nabla f(x_2) - \nabla f(x_1), x_2 - x_1 \rangle > 0$$

Demostración. Supongamos que f es convexa y sean $x_1, x_2 \in S$, entonces

$$f(x_1) \geq f(x_2) + \langle \nabla f(x_2), x_1 - x_2 \rangle$$

$$f(x_2) \geq f(x_1) + \langle \nabla f(x_1), x_2 - x_1 \rangle.$$

Sumando las dos desigualdades se obtiene que

$$\langle \nabla f(x_2) - \nabla f(x_1), x_2 - x_1 \rangle \geq 0.$$

Para el recíproco, sean $x_1, x_2 \in S$. Por el teorema del valor medio, existe $\lambda \in (0, 1)$ tal que

$$f(x_2) - f(x_1) = \langle \nabla f(x), x_2 - x_1 \rangle, \quad (1.14)$$

donde $x = \lambda x_1 + (1 - \lambda)x_2$. Por hipótesis,

$$\langle \nabla f(x) - \nabla f(x_1), x - x_1 \rangle \geq 0,$$

esto implica que $\langle \nabla f(x), x_2 - x_1 \rangle \geq \langle \nabla f(x_1), x_2 - x_1 \rangle$. Por (1.14) se deduce que $f(x_2) \geq f(x_1) + \langle \nabla f(x_1), x_2 - x_1 \rangle$. Y por lo tanto f es una función convexa. \square

1.3.5. Mínimos de una función convexa

Sea $f : \mathbb{R}^n \rightarrow \mathbb{R}$ una función, entonces el problema P se define como

$$P = \min_{x \in S} f(x) \quad (1.15)$$

donde S es un subconjunto convexo de \mathbb{R}^n .

Definición 1.3.5. Sea $f : \mathbb{R}^n \rightarrow \mathbb{R}$ y consideremos el problema P . Un punto $x^* \in S$ es llamado **solución óptima**, **mínimo global** o **solución factible** si $f(x) \geq f(x^*)$ para todo $x \in S$.

Si $x^* \in S$ y existe una vecindad $V(x^*)$ de x^* tal que $f(x) \geq f(x^*)$ para cada $x \in S \cap V(x^*)$, entonces x^* es llamado **mínimo local (solución local)** del problema. Por otro lado, si $x^* \in S$ es el único mínimo local en $S \cap V(x^*)$ para alguna vecindad de x^* , entonces x^* es llamado **mínimo local aislado** del problema.

Hay dos observaciones que se deducen del teorema 1.3.6. Primero, si f es diferenciable, el mínimo de la función $g(y) = f(x) + \langle \nabla f(x), y - x \rangle$ sobre S produce una cota inferior del valor mínimo del problema P , la cual puede ser útil en una aproximación algorítmica.

Como segunda observación, la función afín $g(y)$ puede ser usada para producir aproximaciones poliédricas externas como resultado del corolario 1.2.1.

Teorema 1.3.8. Considerar el problema P . Sea S un subconjunto convexo no vacío de \mathbb{R}^n , y sea $f : S \rightarrow \mathbb{R}$ convexa en S . Supongamos que $x^* \in S$ es un mínimo local del problema. Entonces x^* es un mínimo global.

Demostración[8]. Dado que x^* es una solución local, existe una vecindad $V(x^*)$ de x^* tal que

$$f(x) \geq f(x^*)$$

para cada $x \in S \cap V(x^*)$. Supongamos, por el contrario que x^* no es un mínimo global del problema, es decir, existe $\hat{x} \in S$ tal que $f(\hat{x}) < f(x^*)$. Por la convexidad de f , para cada $\lambda \in [0, 1]$ se tiene que

$$f(\lambda\hat{x} + (1 - \lambda)x^*) \leq \lambda f(\hat{x}) + (1 - \lambda)f(x^*) < \lambda f(x^*) + (1 - \lambda)f(x^*) = f(x^*).$$

Así

$$f(\lambda\hat{x} + (1 - \lambda)x^*) < f(x^*)$$

y se llega a una contradicción cuando λ es suficientemente pequeño. Por lo tanto, x^* es un mínimo global. \square

Un punto muy importante que se prueba fácilmente es que si f es convexa, $x^* \in S$ es un mínimo global del problema P si y solo si $0 \in \partial f(x^*)$. Más aún, si S es abierto y f es convexa y diferenciable en S entonces $x^* \in S$ es un mínimo global del problema P si y solo si $\nabla f(x^*) = 0$.

También se puede caracterizar el mínimo del problema P por medio del siguiente teorema.

Teorema 1.3.9. *Sea $f : \mathbb{R}^n \rightarrow \mathbb{R}$ una función convexa, S un subconjunto convexo no vacío de \mathbb{R}^n , y se considera el problema P . El punto $x^* \in S$ es un mínimo global del problema si y solo si f tiene un subgradiente ξ en x^* tal que $\langle \xi, x - x^* \rangle \geq 0$ para toda $x \in S$.*

Demostración [8]. Primero supongamos que $\langle \xi, x - x^* \rangle \geq 0$ para todo $x \in S$, donde $\xi \in \partial f(x^*)$. Entonces, por hipótesis

$$f(x) \geq f(x^*) + \langle \xi, x - x^* \rangle \geq f(x^*)$$

para todo $x \in S$, y por lo tanto x^* es una solución óptima del problema dado. Para mostrar la necesidad, supongamos que x^* es una solución óptima del problema, se construyen dos conjuntos en \mathbb{R}^{n+1} de la siguiente manera

$$\Gamma_1 = \{(x - x^*, y) : x \in \mathbb{R}^n, y > f(x) - f(x^*)\}$$

$$\Gamma_2 = \{(x - x^*, y) : x \in S, y \leq 0\}.$$

Es fácil verificar que tanto Γ_1 y Γ_2 son conjuntos convexos. Además, $\Gamma_1 \cap \Gamma_2 = \emptyset$ porque de otra manera existiría un punto (x, y) tal que

$$x \in S, \quad 0 \geq y > f(x) - f(x^*),$$

contradiciendo la suposición de que x^* es un mínimo global del problema. Por lo tanto, por el teorema 1.2.2 existe un hiperplano que separa Γ_1 y Γ_2 ; esto es, existe un vector diferente de cero (ξ_0, μ) y un escalar α tal que

$$\langle \xi_0, x - x^* \rangle + \mu y \leq \alpha, \quad \forall x \in \mathbb{R}^n, \quad y > f(x) - f(x^*) \quad (1.16)$$

$$\langle \xi_0, x - x^* \rangle + \mu y \geq \alpha, \quad \forall x \in S, \quad y \leq 0 \quad (1.17)$$

Si se toma $x = x^*$ e $y = 0$ en la igualdad anterior se sigue que $\alpha \leq 0$. Ahora bien, al tomar $x = x^*$ e $y = \epsilon > 0$ en (1.16) se sigue que $\mu\epsilon \leq \alpha$. Dado que esto es verdad para cualquier $\epsilon > 0$, se tiene que $\mu \leq 0$ y $\alpha \geq 0$. Así, se ha demostrado que $\mu \leq 0$ y $\alpha = 0$. Si $\mu = 0$, de (2.19), $\langle \xi_0, x - x^* \rangle \leq 0$ para cada $x \in R^n$. Si se toma $x = x^* + \xi_0$, se sigue que

$$0 \geq \langle \xi_0, x - x^* \rangle = \|\xi_0\|^2$$

y por lo tanto $\xi_0 = 0$. Dado que $(\xi_0, \mu) \neq (0, 0)$, se debe tener que $\mu < 0$. Dividiendo (1.16) y (1.17) por $-\mu$ y denotando $-\xi_0/\mu$ por ξ , se tienen las siguientes desigualdades

$$y \geq \langle \xi, x - x^* \rangle, \quad \forall x \in \mathbb{R}^n, \quad y > f(x) - f(x^*) \quad (1.18)$$

$$\langle \xi, x - x^* \rangle - y \geq 0, \quad \forall x \in S, \quad y \leq 0. \quad (1.19)$$

Si $y = 0$ en la ecuación anterior, se obtiene que $\langle \xi, x - x^* \rangle \geq 0$ para todo $x \in S$. De (1.18) se sigue que

$$f(x) \geq f(x^*) + \langle \xi, x - x^* \rangle$$

para toda $x \in R^n$. Por lo tanto, ξ es un subgradiente de f en x^* con la propiedad de que $\langle \xi, x - x^* \rangle \geq 0$ para toda $x \in S$. \square

1.4. Problemas con restricciones y dualidad

En esta sección se considera el problema no lineal, el cual llamaremos el problema primario. Sea $f : \mathbb{R}^n \rightarrow \mathbb{R}$ una función convexa, sean $g : \mathbb{R}^n \rightarrow \mathbb{R}^m$ y $h : \mathbb{R}^n \rightarrow \mathbb{R}^l$ funciones, entonces el problema primario se define como

$$\begin{aligned} &\text{minimizar } f(x) \\ &\text{sujeto a } x \in C, \end{aligned}$$

donde $C = \{x \in \mathbb{R}^n : g(x) \leq 0, h(x) = 0\}$.

El problema dual se define como

$$\begin{aligned} &\text{maximizar } \theta(u, v) \\ &\text{sujeto a } u \geq 0 \end{aligned}$$

donde $\theta(u, v) = \inf\{f(x) + \sum_{i=1}^m u_i g_i(x) + \sum_{j=1}^l v_j h_j(x) : x \in \mathbb{R}^n\}$. En este caso, las restricciones han sido agregadas a la función f por medio de multiplicadores de Lagrange. Este proceso es llamado **dualización**.

A continuación se muestra la relación que tienen el problema primario y el dual.

Teorema 1.4.1 (Teorema débil de la dualidad). *Sea $x \in C$ y sea (u, v) tal que $u \geq 0$. Entonces $f(x) \geq \theta(u, v)$ y por lo tanto*

$$\inf\{f(x) : x \in C\} \geq \sup\{\theta(u, v) : u \geq 0\}.$$

Demostración. Es inmediata siguiendo la definición de θ . □

Corolario 1.4.1. *Si $f(x^*) = \theta(u^*, v^*)$, donde $u^* \geq 0$ y $x^* \in C$, entonces x^* y (u^*, v^*) resuelven respectivamente el problema primario y el problema dual.*

El teorema 1.4.1 y el corolario 1.4.1 nos garantizan que el valor óptimo del problema primario siempre es mayor o igual que el valor óptimo del problema dual.

Definición 1.4.1. *Se dice que existe una brecha en la dualidad si*

$$\inf\{f(x) : x \in C\} > \sup\{\theta(u, v) : u \geq 0\}.$$

Es de nuestro interés conocer cuándo los dos valores óptimos coinciden. El siguiente teorema nos dice bajo qué circunstancias el valor óptimo del problema primario y el del problema dual son iguales, es decir, bajo qué condiciones no existe una brecha en la dualidad.

Teorema 1.4.2 (Teorema fuerte de la dualidad). *Sean $f : \mathbb{R}^n \rightarrow \mathbb{R}$ y $g : \mathbb{R}^n \rightarrow \mathbb{R}^m$ funciones convexas, y sea $h : \mathbb{R}^n \rightarrow \mathbb{R}^l$ una función afín. Supongamos que existe $\hat{x} \in \mathbb{R}^n$ tal que $g(\hat{x}) < 0$, $h(\hat{x}) = 0$, y que $0 \in \text{int}(h(\mathbb{R}^n))$. Entonces*

$$\inf\{f(x) : x \in C\} = \sup\{\theta(u, v) : u \geq 0\}.$$

La demostración se puede encontrar en [8].

El teorema anterior muestra bajo qué condiciones el valor óptimo del problema primario y del problema dual coinciden. Más adelante se mostrará que una condición necesaria y suficiente para que ambos valores óptimos coincidan es que exista un punto silla.

Definición 1.4.2. *Dado el problema primario, se define la función lagrangiana como*

$$\mathcal{L}(x, u, v) = f(x) + \langle u, g(x) \rangle + \langle v, h(x) \rangle$$

donde $u \geq 0$.

Definición 1.4.3. *Un punto (x^*, u^*, v^*) es llamado punto silla de la función lagrangiana si $x^* \in \mathbb{R}^n$, $u^* \geq 0$, y*

$$\mathcal{L}(x^*, u, v) \leq \mathcal{L}(x^*, u^*, v^*) \leq \mathcal{L}(x, u^*, v^*)$$

para todo $x \in \mathbb{R}^n$ y (u, v) tal que $u \geq 0$.

Así, se tiene que x^* minimiza \mathcal{L} en \mathbb{R}^n cuando (u^*, v^*) es fijo y (u^*, v^*) maximiza \mathcal{L} en $\mathbb{R}^m \times \mathbb{R}^l$ cuando x^* es fija. El siguiente resultado caracteriza un punto silla y muestra que la existencia de este es una condición necesaria y suficiente para que no haya una brecha en la dualidad.

Teorema 1.4.3. *Sea el punto (x^*, u^*, v^*) tal que $x^* \in \mathbb{R}^n$ y $u^* \geq 0$. Entonces (x^*, u^*, v^*) es punto silla para la función Lagrangiana \mathcal{L} si y solo si*

1. $\mathcal{L}(x^*, u^*, v^*) = \min\{\mathcal{L}(x, u^*, v^*) : g(x) \leq 0, h(x) = 0, x \in \mathbb{R}^n\}$
2. $x^* \in C$.
3. $\langle u^*, g(x^*) \rangle = 0$.

Más aún, (x^*, u^*, v^*) es punto silla si y solo si x^* y (u^*, v^*) son, respectivamente, soluciones óptimas del problema primario y dual en donde no hay brecha en la dualidad, es decir, se cumple que $f(x^*) = \theta(u^*, v^*)$.

La demostración se puede encontrar en [8].

Corolario 1.4.2. *Supongamos que f , g son funciones convexas y h es una función afín. Supongamos además que $0 \in \text{int } h(\mathbb{R}^n)$ y que existe $\hat{x} \in \mathbb{R}^n$ tal que $g(\hat{x}) < 0$ y $h(\hat{x}) = 0$. Si x^* es una solución óptima del problema primario, entonces existe un vector (u^*, v^*) con $u^* \geq 0$ tal que (x^*, u^*, v^*) es un punto silla.*

Antes de terminar esta sección queda un resultado más por analizar: la consecuencia de que el problema primario y el problema dual no tengan una brecha en la dualidad. El valor óptimo del problema dual está dado por

$$\theta^* = \sup_{(u,v):u \geq 0} \inf_{x \in \mathbb{R}^n} \{\mathcal{L}(x, u, v)\}.$$

Luego, para cualquier función se cumple que

$$\theta^* \leq \inf_{x \in \mathbb{R}^n} \sup_{(u,v):u \geq 0} \{\mathcal{L}(x, u, v)\}.$$

El supremo de $\mathcal{L}(x, u, v)$ sobre (u, v) con $u \geq 0$ es infinito a menos que $g(x) \leq 0$ y $h(x) = 0$. Por lo tanto,

$$\begin{aligned} \theta^* &\leq \inf_{x \in \mathbb{R}^n} \sup_{(u,v):u \geq 0} \{\mathcal{L}(x, u, v)\} \\ &= \inf\{f(x) : x \in C\}, \end{aligned}$$

que es solución del problema primario. Por lo tanto, se observa que el valor óptimo del problema primario coincide con el del problema dual si y solo si se cumple $\sup \inf \mathcal{L} = \inf \sup \mathcal{L}$; por el teorema anterior, suponiendo que existe solución óptima, esto ocurre si y solo si existe un punto silla (x^*, u^*, v^*) para la función lagrangiana \mathcal{L} .

1.5. Optimización no diferenciable

La optimización no diferenciable trata problemas en los cuales la función objetivo no es diferenciable en todos lados. El gradiente no existe en todos lados, lo que implica que la función puede tener dobleces o puntos de esquina y por lo tanto no se puede aproximar localmente por un hiperplano tangente o por una aproximación cuadrática.

Definición 1.5.1. Sea $S \subset \mathbb{R}^n$, $x \in S$ y $d \in \mathbb{R}^n$ unitario. Se dice que d es una **dirección factible** con respecto a S en x si existe una sucesión $\{x_n\}_{n \in \mathbb{N}}$ contenida en S que cumple que $x_n \rightarrow x$ y que la sucesión

$$d_n = \frac{x_n - x}{\|x_n - x\|} \rightarrow d. \quad (1.20)$$

Ahora se dará un resultado análogo al teorema 1.3.3 utilizando direcciones factibles.

Lema 1.5.1. *Sea $f : S \rightarrow \mathbb{R}$ una función convexa, donde S es un subconjunto convexo y no vacío de \mathbb{R}^n y sea d una dirección factible con respecto a S en x , donde $\{x_n\}$ es una sucesión que cumple (1.20). Entonces*

$$\lim_{n \rightarrow \infty} \frac{f(x_n) - f(x)}{\|x_n - x\|} = \max_{\xi \in \partial f(x)} \langle \xi, d \rangle.$$

La demostración es similar a la del teorema 1.5.1.

El principal interés en introducir lo anterior, es estudiar las funciones construidas de la siguiente manera

$$\phi(x) = f(x) + h(c(x)) \quad (1.21)$$

donde $f : \mathbb{R}^n \rightarrow \mathbb{R}$ y $c : \mathbb{R}^n \rightarrow \mathbb{R}^m$ son funciones de clase C^1 y $h : \mathbb{R}^m \rightarrow \mathbb{R}$ es una función convexa no necesariamente diferenciable.

Si d es una dirección factible en \mathbb{R}^n , usando (1.20) y la serie de Taylor

$$f(x_n) = f(x) + \|x_n - x\| \langle \nabla f(x), d_n \rangle + o(\|x_n - x\|),$$

así

$$\frac{f(x_n) - f(x)}{\|x_n - x\|} \rightarrow \langle \nabla f(x), d \rangle.$$

De la misma manera, se obtiene para c que

$$\frac{c(x_n) - c(x)}{\|x_n - x\|} \rightarrow Dc(x)d.$$

Aplicando el teorema 1.5.1 a h se deduce

$$\lim_{\lambda \rightarrow 0} \frac{\phi(x_n) - \phi(x)}{\|x_n - x\|} = \max_{\xi \in \partial h(x)} \langle \nabla f(x) + D^*c(x)\xi, d \rangle. \quad (1.22)$$

Esto proporciona la derivada direccional de x en dirección factible d para la función $\phi(x)$. Si x^* es un mínimo local de ϕ , entonces $\phi(x) \geq \phi(x^*)$ para x suficientemente cerca de x^* , y por lo tanto

$$\max_{\xi \in \partial h(x^*)} \langle \nabla f(x^*) + D^*c(x^*)\xi, d \rangle \geq 0$$

para cualquier dirección factible d .

Esta es una condición de primer orden necesaria para un mínimo local, la cual puede ser interpretada como una derivada direccional positiva en todas las direcciones $d \in \mathbb{R}^n$. Más aún, como x^* es un mínimo local de ϕ entonces

$$0 \in \partial\phi(x^*) = \{\nabla f(x^*) + D^*c(x^*)\xi : \xi \in \partial h(x^*)\}. \quad (1.23)$$

La ecuación (1.22) y el conjunto (1.23) son resultado de una **generalización de la regla de la cadena** para funciones diferenciables [16]. El conjunto $\partial\phi(x^*)$, aunque convexo y compacto, no es la subdiferencial porque ϕ puede no ser convexa, pero por conveniencia se usa la misma notación.

Otra forma de establecer la condición (1.23) es por medio de la función lagrangiana

$$\mathcal{L}(x, \alpha) = f(x) + \langle \alpha, c(x) \rangle, \quad (1.24)$$

donde $\alpha \in \mathbb{R}^m$. Entonces un enunciado equivalente a (1.23) es el siguiente.

Teorema 1.5.1. *Si x^* minimiza $\phi(x)$ entonces existe un vector $\alpha^* \in \partial h(x^*)$ tal que*

$$\nabla \mathcal{L}(x^*, \alpha^*) = \nabla f(x^*) + D^*c(x^*)\alpha^* = 0. \quad (1.25)$$

Demostración. Es inmediata de (1.23). □

Ya que la función ϕ no es necesariamente convexa, es importante considerar condiciones de segundo orden para que x^* sea mínimo local de ϕ . Estas condiciones muestran la estrecha relación entre nuestro problema y el problema en el caso que ϕ sea diferenciable. A continuación se darán los antecedentes para deducir condiciones de segundo orden.

Sea α^* un vector que cumpla con el teorema 1.5.1. Se define el conjunto el conjunto X por

$$X = \{x : h(c(x)) = h(c(x^*)) + \langle c(x) - c(x^*), \alpha^* \rangle\}. \quad (1.26)$$

Se denotará como \mathcal{G}^* al conjunto de direcciones factibles con respecto a X en x^* . Es posible mostrar que estas direcciones están relacionadas con el conjunto G^* definido como

$$G^* = \{d : \max_{\xi \in \partial h(x^*)} \langle \nabla f(x^*) + D^*c(x^*)\xi, d \rangle = 0, \|d\| = 1\}. \quad (1.27)$$

Es decir, G^* es el conjunto de direcciones factibles que hacen $\phi'(x^*; d) = 0$.

Lema 1.5.2. $\mathcal{G}^* \subset G^*$.

Demostración [19]. Sea $d \in \mathcal{G}^*$, entonces d es una dirección factible con respecto a X en x^* y por lo tanto existe una sucesión $\{x_n\}_{n \in \mathbb{N}}$ contenida en X tal que $x_n \rightarrow x^*$. Utilizando la serie de Taylor, el lema (1.5.1) y la definición de ϕ

$$\begin{aligned} \max_{\xi \in \partial h(x^*)} \langle \nabla f(x^*) + Dc(x^*)\xi, d \rangle &= \lim_{n \rightarrow \infty} \frac{\phi(x_n) - \phi(x^*)}{\|x_n - x^*\|} \\ &= \lim_{n \rightarrow \infty} \frac{f(x_n) - f(x^*) + h(c(x_n)) - h(c(x^*))}{\|x_n - x^*\|} \\ &= \lim_{n \rightarrow \infty} \frac{f(x_n) - f(x^*) + \langle c(x_n) - c(x^*), \alpha^* \rangle}{\|x_n - x^*\|} \\ &= \langle \nabla f(x^*) + D^*c(x^*)\alpha^*, d \rangle. \end{aligned}$$

Luego, por el teorema 1.5.1 se obtiene lo enunciado. \square

En el caso en que la función f es diferenciable, para tener condiciones de segundo orden se necesita una condición de regularidad. En nuestro caso particular, se tiene una condición de regularidad como en el caso diferenciable para poder deducir condiciones de segundo orden.

La condición de regularidad para nuestro problema es:

$$\mathcal{G}^* = G^*. \quad (1.28)$$

Ahora es posible establecer condiciones de segundo orden. Para hacer esto supongamos que f y c son de clase \mathbb{C}^2 . Como es usual, las condiciones de regularidad solo son utilizadas para establecer condiciones necesarias.

Teorema 1.5.2. (Condiciones necesarias de segundo orden) Si x^* minimiza $\phi(x)$ y si la igualdad (1.28) se mantiene entonces

$$d^t \nabla^2 \mathcal{L}(x^*, \alpha^*) d \geq 0 \quad (1.29)$$

para todo $d \in G^*$, donde α^* cumple con el teorema 1.5.1.

Demostración [19]. Sea $d \in \mathcal{G}^*$, entonces d es una dirección factible con respecto a X en x^* , por lo tanto existe una sucesión $\{x_n\}_{n \in \mathbb{N}}$ contenida en X tal que $x_n \rightarrow x^*$, se denotara $\delta_n = \|x_n - x^*\|$. La expansión de Taylor de $\mathcal{L}(x, \alpha^*)$ alrededor de x^* es

$$\begin{aligned} \mathcal{L}(x_n, \alpha^*) &= f(x^*) + \langle c(x^*), \alpha^* \rangle + \langle \nabla \mathcal{L}(x^*, \alpha^*), x_n - x^* \rangle \\ &\quad + \frac{1}{2} e_n^t \nabla^2 \mathcal{L}(x^*, \alpha^*) e_n + o(\|x_n - x^*\|^2). \end{aligned}$$

donde $e_n = x_n - x^*$. Usando (1.25) y (1.20) se obtiene

$$\begin{aligned} \mathcal{L}(x_n, \alpha^*) &= f(x^*) + \langle c(x^*), \alpha^* \rangle + \frac{1}{2} \delta_n^2 d_n^t \nabla^2 \mathcal{L}(x^*, \alpha^*) d_n \\ &\quad + o(\|\delta_n\|^2). \end{aligned}$$

Como $\{x_n\}$ es una sucesión factible con respecto a X en x^* se sigue de la definición de ϕ que

$$\phi(x_n) = \phi(x^*) + \frac{1}{2} \delta_n^2 d_n^t \nabla^2 \mathcal{L}(x^*, \alpha^*) d_n + o(\|\delta_n\|^2).$$

Dado que x^* es un mínimo local, $\phi(x_n) \geq \phi(x^*)$ (para n suficientemente grande). Así, dividiendo entre $\frac{1}{2} \delta_n^2$ y tomando el límite cuando $n \rightarrow \infty$ se obtiene que

$$d^t \nabla^2 \mathcal{L}(x^*, \alpha^*) d \geq 0.$$

□

Ahora se consideran las condiciones bajo las cuáles $\mathcal{G}^* = G^*$, para esto se utilizará la siguiente definición y otras observaciones.

Definición 1.5.2. *Si existe una vecindad abierta $V(c(x^*))$ tal que*

$$h(c(x)) = h(c(x^*)) + \max_{\xi \in \partial h(x^*)} \langle c(x) - c(x^*), \xi \rangle \quad (1.30)$$

para todo $c(x) \in V(c(x^*))$, se dice que $h(c(x))$ es localmente lineal en $c(x^*)$.

Esta propiedad es cierta para cualquier función poliédrica convexa pero también permite manipular ciertos problemas formulados por medio funciones de penalización exacta con norma diferenciable [19].

Por otro lado, se denota a la dimensión de $\partial h(x^*)$ por l^* ($l^* \leq m$) y se define la matriz $D \in \mathbb{R}^{m \times l^*}$ cuyas columnas consisten en los elementos d_i^* , $i = 1, 2, \dots, l^*$ que son una base para $\partial h(c(x^*)) - \alpha^*$, que simplemente es una traslación de $\partial h(c(x^*))$. Bajo estas construcciones se enuncia el siguiente teorema que da condiciones suficientes para tener regularidad.

Lema 1.5.3. *(Condiciones suficientes de regularidad) Si x^* satisface las condiciones de primer orden, si $h(c)$ es localmente lineal en c^* , y si*

$$\text{rango}(\nabla f(x^*) D^*) = l^* \quad (1.31)$$

entonces $\mathcal{G}^* = G^*$

La demostración se puede encontrar en [19].

Estas no son las únicas condiciones suficientes para tener $\mathcal{G}^* = G^*$, por ejemplo, el mismo resultado se tiene si c es una función afín o si se puede probar que $G^* = \emptyset$ [19].

Métodos para optimización no diferenciable

En este capítulo se describirán los métodos estudiados para la minimización de funciones que son convexas pero no diferenciables. De igual forma se darán las pruebas de la convergencia de cada método.

2.1. Método de subgradiente

Este es el método más sencillo y es muy parecido al método de gradiente para funciones diferenciables pero cuenta con algunas modificaciones [12].

1. El método de subgradiente se aplica directamente a la función no diferenciable.
2. Los tamaños de paso no son elegidos por búsqueda lineal. En la mayoría de los casos, el tamaño de paso es fijo.
3. A diferencia del método de gradiente, el método de subgradiente no es un método de descenso; el valor de la función puede incrementar.

En este caso solo se tratará el problema de minimizar una función sin restricciones, es decir, minimizar $f : \mathbb{R}^n \rightarrow (-\infty, +\infty]$ donde f es convexa.

El método de subgradiente usa la iteración

$$x^{(k+1)} = x^k - \alpha_k g^k$$

donde x^k es el punto resultante de la k -ésima iteración, g^k es cualquier subgradiente de f en x^k y $\alpha_k > 0$ es el k -ésimo tamaño de paso. Así en cada iteración del método de subgradiente, se hace el paso en la dirección del subgradiente negativo.

Cuando la función es diferenciable, la única elección posible de g^k es $\nabla f(x^k)$, y el método del subgradiente se reduce al método de gradiente (excepto por la elección del tamaño de paso). Un problema en el método de subgradiente es que se toma cualquier subgradiente en cada iteración, por lo tanto no se toma en cuenta cual es la elección de subgradiente que hace decrecer el valor de la función en el nuevo punto.

Puede ocurrir que $-g^k$ no sea una dirección de descenso para f en x^k , es decir, que no se cumple necesariamente

$$\langle -g^k, \xi^k \rangle \leq 0.$$

para todo $\xi \in \partial f(x^*)$. Un ejemplo de esta aseveración se observa al tomar la función $f(x, y) = |x| + 2|y|$; $g = (1, 2) \in \partial f(1, 0)$, pero $-g$ no es una dirección de descenso en $(1, 0)$.

En el método de subgradiente se utilizan varios tamaños de paso los cuales influyen en la convergencia del método. Los tamaños de paso más usuales son los siguientes:

1. Tamaño de paso constante. Se toma $\alpha_k = \alpha$ una constante positiva que es independiente de k .
2. Longitud de paso constante. Se toma $\alpha_k = \gamma / \|g^k\|_2$. Esto implica que $\|x^{k+1} - x^k\|_2 = \gamma$.
3. Cuadrados sumables. En este caso se toma α_k de tal manera que

$$\alpha_k \geq 0, \quad \sum_{k=1}^{\infty} \alpha_k^2 < \infty, \quad \sum_{k=1}^{\infty} \alpha_k = \infty$$

4. Nonsummable diminishing. Se toma α_k tal que

$$\alpha_k \geq 0, \quad \lim_{k \rightarrow \infty} \alpha_k = 0, \quad \sum_{k=1}^{\infty} \alpha_k = \infty$$

Dado que el método de subgradiente no es un método de descenso se hará en cada iteración lo siguiente

$$f_{best}^k = \min\{f_{best}^{k-1}, f(x^k)\}$$

y se hace $i_{best}^k = k$ si $f(x^k) = f_{best}^k$, es decir, se guardará el mejor punto obtenido hasta la iteración k .

Ahora, para probar la convergencia del método se supondrá que existe un mínimo de la función f , digamos x^* . También se supondrá lo siguiente:

1. La norma de los subgradientes es acotada, es decir, existe $G \in \mathbb{R}$ tal que $\|g^k\|_2 \leq G$ para toda k .
2. Se conoce un número R tal que $R \geq \|x^1 - x^*\|_2$.

Teniendo en cuenta lo anterior y que x^* es cualquier punto óptimo entonces

$$\begin{aligned} \|x^{k+1} - x^*\|_2^2 &= \|x^k - \alpha_k g^k - x^*\|_2^2 \\ &= \|x^k - x^*\|_2^2 - 2\alpha_k g^{(k)t}(x^k - x^*) + \alpha_k^2 \|g^k\|_2^2 \\ &\leq \|x^k - x^*\|_2^2 - 2\alpha_k (f(x^k) - f(x^*)) + \alpha_k^2 \|g^k\|_2^2, \end{aligned} \quad (2.1)$$

donde la última desigualdad se obtuvo por la definición de subgradiente, es decir, g^k cumple que

$$f(x^*) \geq f(x^k) + \langle g^k, x^* - x^k \rangle.$$

Aplicando la desigualdad (2.1) recursivamente se obtiene que

$$\|x^{k+1} - x^*\|_2^2 \leq \|x^1 - x^*\|_2^2 - 2 \sum_{i=1}^k \alpha_i (f(x^i) - f(x^*)) + \sum_{i=1}^k \alpha_i^2 \|g^i\|_2^2.$$

Usando el hecho de que $\|x^{k+1} - x^*\|_2^2 \geq 0$ y $\|x^1 - x^*\|_2 \leq R$ se tiene que

$$2 \sum_{i=1}^k \alpha_i (f(x^i) - f(x^*)) \leq R^2 + \sum_{i=1}^k \alpha_i^2 \|g^i\|_2^2. \quad (2.2)$$

Por otro lado,

$$\begin{aligned} \sum_{i=1}^k \alpha_i (f(x^i) - f(x^*)) &\geq \left(\sum_{i=1}^k \alpha_i \right) \min_i (f(x^i) - f(x^*)) + \sum_{i=1}^k \|g^i\|_2^2 \\ &= \left(\sum_{i=1}^k \alpha_i \right) (f_{best}^k - f(x^*)). \end{aligned}$$

Combinando esta desigualdad con (2.2) se obtiene que

$$f_{best}^k - f(x^*) = \min_i (f(x^i) - f(x^*)) \leq \frac{R^2 + \sum_{i=1}^k \alpha_i^2 \|g^i\|_2^2}{2 \sum_{i=1}^k \alpha_i}. \quad (2.3)$$

Finalmente, usando la hipótesis de que $\|g^k\|_2 \leq G$, se consigue la siguiente desigualdad

$$f_{best}^k - f(x^*) \leq \frac{R^2 + G^2 \sum_{i=1}^k \alpha_i^2}{2 \sum_{i=1}^k \alpha_i}. \quad (2.4)$$

De esta desigualdad se desprenden varios resultados de convergencia, dependiendo el tamaño de paso que se esté usando se tendrá una mejor aproximación al punto óptimo.

1. Tamaño de paso constante. Cuando $\alpha_k = \alpha$, se tiene usando (2.4) que

$$f_{best}^k - f(x^*) \leq \frac{R^2 + G^2 \alpha^2 k}{2\alpha k}.$$

Cuando $k \rightarrow \infty$ entonces el lado derecho converge a $G^2 \alpha / 2$.

2. Longitud de paso constante. Usando $\alpha_k = \gamma / \|g^k\|_2$ y la desigualdad (2.3) se obtiene

$$f_{best}^k - f(x^*) \leq \frac{R^2 + \gamma^2 k}{2 \sum_{i=1}^k \alpha_i} \leq \frac{R^2 + \gamma^2 k}{2\gamma k / G},$$

usando $\alpha_i \geq \gamma / G$. Cuando $k \rightarrow \infty$ el lado derecho converge a $G\gamma/2$.

3. Cuadrados sumables. Ahora supongamos que

$$\|\alpha\|_2^2 = \sum_{i=1}^{\infty} \alpha_i^2 < \infty, \quad \sum_{i=1}^{\infty} \alpha_i = \infty.$$

Entonces se obtiene que

$$f_{best}^k - f(x^*) \leq \frac{R^2 + G^2 \|\alpha\|_2^2}{2 \sum_{i=1}^k \alpha_i}.$$

Cuando $k \rightarrow \infty$ el lado derecho converge a cero. Así el método de subgradiente converge.

4. Diminishing step. Si la sucesión α_k converge a cero y $\sum \alpha_i = \infty$, entonces el lado derecho de la desigualdad (2.4) converge a cero, lo cual implica que el método de subgradiente converge. Para probar esto, sea $\epsilon > 0$. Entonces existe un entero N_1 tal que $\alpha_i < \epsilon / G^2$ para $i > N_1$. También existe un natural N_2 tal que

$$\sum_{i=1}^{N_2} \alpha_i > \frac{1}{\epsilon} (R^2 + G^2 \sum_{i=1}^{N_2} \alpha_i^2).$$

Ahora, sea $N = \max\{N_1, N_2\}$, entonces para $k > N$ se tiene que

$$\begin{aligned} \frac{R^2 + G^2 \sum_{i=1}^k \alpha_i^2}{2 \sum_{i=1}^k \alpha_i} &\leq \frac{R^2 + G^2 \sum_{i=1}^{N_1} \alpha_i^2}{2 \sum_{i=1}^k \alpha_i} + \frac{G^2 \sum_{i=N_1+1}^k \alpha_i^2}{2 \sum_{i=1}^{N_1} \alpha_i + 2 \sum_{i=N_1+1}^k \alpha_i} \\ &\leq \frac{R^2 + G^2 \sum_{i=1}^{N_1} \alpha_i^2}{\frac{2}{\epsilon} (R^2 + G^2 \sum_{i=1}^{N_1} \alpha_i^2)} + \frac{G^2 \sum_{i=N_1+1}^k (\epsilon \alpha_i / G^2)}{2 \sum_{i=N_1+1}^k \alpha_i} \\ &= \frac{\epsilon}{2} + \frac{\epsilon}{2} = \epsilon. \end{aligned}$$

Así el método de subgradiente converge.

A partir de la ecuación (2.2) se puede obtener un criterio de paro para nuestro método. De nuevo, suponiendo que $R \geq \|x^1 - x^*\|_2$ y usando la desigualdad (2.2) se tiene que

$$f(x^*) \geq \frac{2 \sum_{i=1}^k \alpha_i f(x^i) - R^2 - \sum_{i=1}^k \alpha_i^2 \|g^k\|_2^2}{2 \sum_{i=1}^k \alpha_i} = l_k, \quad (2.5)$$

donde l_k puede ser calculado a partir de la k -ésima iteración. La iteración l_k no es necesariamente creciente [9] por lo que se puede guardar la mejor cota superior de $\{l_k\}$

$$l_{best}^k = \text{máx}\{l_1, \dots, l_k\}.$$

2.2. Método de planos de corte

Existen métodos en los cuales toda la información previa es usada para obtener una nueva iteración. Esta información puede ser utilizada para crear un modelo de la función f , de tal manera que se obtenga una aproximación lineal al problema original. Esto es usado en los métodos de planos de corte.

Al tomar en cuenta la información acumulada en las iteraciones anteriores $\{y^i, f(y^i), s^i \in \partial f(y^i)\}_{i=1}^l$, se puede construir la siguiente aproximación lineal a trozos de f

$$\hat{f}_l(y) = \text{máx}_{i=1, \dots, l} f(y^i) + \langle s^i, y - y^i \rangle.$$

Por construcción, $\hat{f}_l(y) \leq f(y)$ para toda $y \in \mathbb{R}^n$ y también se tiene que $\hat{f}_l(y) \leq \hat{f}_{l+1}(y)$ para toda $y \in \mathbb{R}^n$. Así, se modela la función a minimizar con aproximaciones lineales, convirtiendo el problema en uno de programación lineal. Este método está motivado por el lema 1.2.1 y el teorema 1.1.1.

Para que el método esté bien definido, se necesita especificar un conjunto compacto C y restringir nuestro modelo a C para asegurar que cada iteración esté bien definida.

El algoritmo del método de planos de corte es el siguiente:

1. Sea $\bar{\delta} > 0$ y C un conjunto compacto que contiene el punto mínimo de f . Sea $k = 1$, $y^1 \in C$ y $f_0 = -\infty$.
2. Calcular $f(y^k)$ y el subgradiente $s^k \in \partial f(y^k)$.
3. Definir $\delta_k = f(y^k) - \hat{f}_{k-1}(y^k) \geq 0$.
4. Si $\delta_k < \bar{\delta}$.

5. Actualizar el modelo $\hat{f}_k(y) = \max\{\hat{f}_{k-1}(y), f(x^k) + \langle s^k, y - y^k \rangle\}$.
6. Calcular $y^{k+1} \in \arg \min_{y \in C} \hat{f}_k(y)$
7. Hacer $k = k + 1$ e ir al paso 2.

A comparación del método de subgradiente, el método de planos de corte brinda un criterio de paro basado en δ_k el cual no existe en el método de subgradiente.

Ahora bien, por construcción el modelo satisface que

$$\hat{f}_k(y) \leq f(y)$$

y de aquí se sigue que

$$\min_y \hat{f}(y) \leq \min_y f(y).$$

Entonces, si el criterio de paro se satisface se tiene para alguna k que $f(y^k) - \hat{f}_{k-1}(y^k) < \bar{\delta}$, lo cual implica que

$$f(y^k) < \bar{\delta} + \hat{f}_{k-1}(y^k) = \bar{\delta} + \min_y \hat{f}_{k-1}(y) \leq \bar{\delta} + \min_y f(y).$$

2.3. Método *bundle*

Utilizando como motivación el algoritmo anterior, se construye un nuevo método utilizando la regularización de Moreau-Yosida, cuyas propiedades probadas a continuación son importantes para este trabajo.

2.3.1. Regularización de Moreau Yosida

Sea $f : \mathbb{R}^n \rightarrow (-\infty, +\infty]$ una función propia y convexa, además supongamos que f alcanza su mínimo. Se define la **regularización de Moreau-Yosida** de f para $\lambda > 0$ como sigue

$$F_\lambda(v) = \min_{x \in \mathbb{R}^n} \left\{ f(x) + \frac{1}{2\lambda} \|x - v\|^2 \right\}. \quad (2.6)$$

El punto $p(v) \in \mathbb{R}^n$ que cumple $f(p(v)) + \frac{1}{2\lambda} \|p(v) - v\|^2 = F_\lambda(v)$ está bien definido dado que la función que se está minimizando es estrictamente convexa y f alcanza su mínimo.

Lema 2.3.1. *Sea $f : \mathbb{R}^n \rightarrow (-\infty, +\infty]$ una función propia, cerrada y convexa, entonces la función $F_\lambda(v)$ es cerrada y estrictamente convexa.*

Demostración. Dado que f es cerrada y la norma es una función continua entonces se sigue inmediatamente que $F_\lambda(v)$ es cerrada. Ahora, por la convexidad de f y dado que la función norma es estrictamente convexa entonces F_λ es estrictamente convexa. \square

Proposición 2.3.1. *Si $f : \mathbb{R}^n \rightarrow (-\infty, +\infty]$ es una función propia, cerrada y convexa. Entonces F_λ es finita y diferenciable en cualquier parte con gradiente dado por*

$$\nabla F_\lambda(v) = \frac{1}{\lambda}(v - p(v)).$$

Más aún,

$$\|\nabla F_\lambda(v) - \nabla F_\lambda(v')\|^2 \leq \frac{1}{\lambda} \langle \nabla F_\lambda(v) - \nabla F_\lambda(v'), v - v' \rangle$$

para toda $v, v' \in \mathbb{R}^n$. Además

$$\|\nabla F_\lambda(v) - \nabla F_\lambda(v')\| \leq \frac{1}{\lambda} \|v - v'\|$$

para toda $v, v' \in \mathbb{R}^n$.

Demostración [12]. (Finita). Si f es una función propia, existe $x \in \mathbb{R}^n$ tal que $f(x) < +\infty$. Así F_λ es finita dado que

$$F_\lambda(v) \leq f(x) + \frac{1}{2\lambda} \|x - v\|^2 < \infty.$$

(Diferenciabilidad). Sea $d \in \mathbb{R}^n$ unitario, $v \in \mathbb{R}^n$ y $t > 0$. Por definición de F_λ

$$\begin{aligned} \frac{F_\lambda(v + td) - F_\lambda(v)}{t} &= \frac{\min_x [f(x) + \frac{1}{2\lambda} \|x - v - td\|^2] - \min_w [f(w) + \frac{1}{2\lambda} \|w - v\|^2]}{t} \\ &\geq \frac{[f(p(v + td)) + \frac{1}{2\lambda} \|p(v + td) - v - td\|^2]}{t} \\ &\quad - \frac{[f(p(v + td)) + \frac{1}{2\lambda} \|p(v + td) - v\|^2]}{t} \\ &= \frac{1}{2\lambda} \frac{\|p(v + td) - v - td\|^2 - \|p(v + td) - v\|^2}{t} \\ &= \frac{1}{2\lambda} \frac{\|p(v + td) - p(v) + p(v) - v - td\|^2}{t} \\ &\quad - \frac{\|p(v + td) - p(v) + p(v) - v\|^2}{t} \end{aligned}$$

$$= \frac{1}{2\lambda} \frac{\|p(v) - v - td\|^2 - \|p(v) - v\|^2}{t} - \frac{1}{\lambda} \langle p(v + td) - p(v), d \rangle,$$

donde se usó el hecho de que $F_\lambda(v) \leq f(p(v + td)) + \frac{1}{2\lambda} \|p(v + td) - v\|^2$. Ahora, al tomar el límite cuando $t \rightarrow 0$ se obtiene

$$\lim_{t \rightarrow 0} \frac{F_\lambda(v + td) - F_\lambda(v)}{t} \geq \frac{1}{\lambda} \langle v - p(v), d \rangle$$

Por otro lado, dado que $F_\lambda(v + td) \leq f(p(v)) + \frac{1}{2\lambda} \|p(v) - v - td\|^2$,

$$\begin{aligned} \frac{F_\lambda(v + td) - F_\lambda(v)}{t} &= \frac{\min_x [f(x) + \frac{1}{2\lambda} \|x - v - td\|^2] - \min_w [f(w) + \frac{1}{2\lambda} \|w - v\|^2]}{t} \\ &\leq \frac{[f(p(v)) + \frac{1}{2\lambda} \|p(v) - v - td\|^2] - [f(p(v)) + \frac{1}{2\lambda} \|p(v) - v\|^2]}{t} \\ &= \frac{1}{2\lambda} \frac{\|p(v) - v - td\|^2 - \|p(v) - v\|^2}{t}. \end{aligned}$$

De nuevo, al tomar $t \rightarrow 0$,

$$\lim_{t \rightarrow 0} \frac{F_\lambda(v + td) - F_\lambda(v)}{t} \leq \frac{1}{\lambda} \langle v - p(v), d \rangle.$$

Por lo tanto

$$\nabla F_\lambda(v) = \frac{1}{\lambda} (v - p(v)).$$

Para probar las aseveraciones restantes se observa lo siguiente

$$\begin{aligned} \|\nabla F_\lambda(v) - \nabla F_\lambda(v')\|^2 &= \frac{1}{\lambda} \langle \nabla F_\lambda(v) - \nabla F_\lambda(v'), v - p(v) - v' + p(v') \rangle \\ &= \frac{1}{\lambda} \langle \nabla F_\lambda(v) - \nabla F_\lambda(v'), v - v' \rangle + \frac{1}{\lambda} \langle \nabla F_\lambda(v) - \nabla F_\lambda(v'), p(v') - p(v) \rangle. \end{aligned}$$

Se mostrará que el último término es negativo. Para hacer esto primero se verificará lo siguiente: sea $\xi \in \partial f(v)$ y $\xi' \in \partial f(v')$, entonces $\langle \xi - \xi', v - v' \rangle \geq 0$. Para probar esto, se toman en cuenta las desigualdades siguientes, asociadas a ξ y ξ' , y posteriormente aplicadas a v' y v respectivamente

$$f(v') \geq f(v) + \langle \xi, v' - v \rangle$$

$$f(v) \geq f(v') + \langle \xi', v - v' \rangle,$$

sumando estas dos desigualdades se sigue que

$$\langle \xi - \xi', v - v' \rangle \geq 0. \tag{2.7}$$

Por optimalidad de $p(v)$ y de $p(v')$ se tiene que

$$0 \in \partial(f(p(v)) + \frac{1}{2\lambda}\|p(v) - v\|^2)$$

y

$$0 \in \partial(f(p(v')) + \frac{1}{2\lambda}\|p(v') - v'\|^2).$$

Por lo tanto, $\nabla F_\lambda(v) = \frac{1}{\lambda}(v - p(v)) \in \partial f(p(v))$ y $\nabla F_\lambda(v') = \frac{1}{\lambda}(v' - p(v')) \in \partial f(p(v'))$. Usando (2.7) se obtiene

$$(\nabla F_\lambda(v) - \nabla F_\lambda(v'), p(v') - p(v)) \leq 0. \quad (2.8)$$

La prueba de la última desigualdad es consecuencia de (2.8),

$$\begin{aligned} \|\nabla F_\lambda(v) - \nabla F_\lambda(v')\|^2 &\leq \frac{1}{\lambda} \langle \nabla F_\lambda(v) - \nabla F_\lambda(v'), v - v' \rangle \\ &\leq \frac{1}{\lambda} \|\nabla F_\lambda(v) - \nabla F_\lambda(v')\| \|v - v'\|. \end{aligned}$$

Por lo tanto, $\|\nabla F_\lambda(v) - \nabla F_\lambda(v')\| \leq \frac{1}{\lambda} \|v - v'\|$. \square

La regularización de Moreau-Yosida es importante en nuestro trabajo por la proposición anterior y porque tiene los mismos mínimos que la función f como se muestra en la siguiente proposición.

Proposición 2.3.2. *Los siguientes enunciados son equivalentes*

1. $v \in \arg \min\{f(x) : x \in \mathbb{R}^n\}$
2. $v = p(v)$
3. $\nabla F_\lambda(v) = 0$
4. $v \in \arg \min\{F_\lambda(x) : x \in \mathbb{R}^n\}$
5. $f(v) = f(p(v))$
6. $f(v) = F_\lambda(v)$

Demostración [12]. (1) \Rightarrow (2) Dado que $v \in \arg \min\{f(x) : x \in \mathbb{R}^n\}$ se tiene que

$$f(v) = f(v) + \frac{1}{2\lambda}\|v - v\|^2 \leq f(x) + \frac{1}{2\lambda}\|x - v\|^2$$

para toda $x \in \mathbb{R}^n$. Entonces, $p(v) = v$.

(2) \Rightarrow (3) Supongamos que $p(v) = v$, entonces $\nabla F_\lambda(v) = \frac{1}{\lambda}(v - p(v)) = 0$.

(3) \Rightarrow (4) Dado que F_λ es convexa, entonces si $\nabla F_\lambda(v) = 0$ se debe tener que $v \in \arg \min\{F_\lambda(x) : x \in \mathbb{R}^n\}$.

(4) \Rightarrow (5) Usando el inciso 2, $v = p(v)$ y por lo tanto $f(v) = f(p(v))$.

(5) \Rightarrow (6) Por definición de F_λ e hipótesis se tiene la siguiente desigualdad

$$f(v) = f(p(v)) \leq f(p(v)) + \frac{1}{2\lambda}\|p(v) - v\|^2 = F_\lambda(v) \leq f(v) + \frac{1}{2\lambda}\|v - v\|^2.$$

Así, $f(v) = F_\lambda(v)$.

(6) \Rightarrow (1) Dado que $F_\lambda(v) = f(v)$, entonces $p(v) = v$. Así, por optimalidad de v

$$0 \in \partial f(v) + \frac{1}{\lambda}(v - p(v)) = \partial f(v).$$

□

2.3.2. Algoritmo del método *bundle*

Tomando como base el método de planos de corte se construye el método *bundle*. Se empezará agregando un punto x^k al conjunto de información que se acumula conforme se realizan iteraciones. Se seguirá usando el modelo lineal y se obtendrá un nuevo punto en la iteración usando la regularización de Moreau-Yosida.

El algoritmo para el método *bundle* (BA) es el siguiente:

1. Se fija $\hat{\delta} > 0$, $m \in (0, 1)$, x^0 (punto inicial), $y^0 = x^0$ y $k = 0$. Se calcula $f(x^0)$ y $s^0 \in \partial f(x^0)$ y se define $\hat{f}_0(y)$ como sigue

$$\hat{f}_0(y) = f(x^0) + \langle s^0, y - x^0 \rangle \quad (2.9)$$

es decir, \hat{f}_0 es un hiperplano de soporte de f en x^0 .

2. Se calcula la siguiente iteración

$$y^{k+1} \in \arg \min_{y \in \mathbb{R}^n} \hat{f}_k(y) + \frac{\mu_k}{2}\|y - x^k\|^2 \quad (2.10)$$

donde $\mu_k > 0$ es una sucesión creciente tal que $\mu_k \rightarrow 0$ y $\sum_{k=1}^{\infty} \mu_k = \infty$.

3. Se define δ_k como sigue

$$\delta_k = f(x^k) - [\hat{f}_k(y^{k+1}) + \frac{\mu_k}{2}\|y^{k+1} - x^k\|^2].$$

4. Si $\delta_k \leq \hat{\delta}$ entonces el algoritmo termina.
5. Se calcula $f(y^{k+1})$ y $s^{k+1} \in \partial f(y^{k+1})$.
6. Si $f(x^k) - f(y^{k+1}) > m\delta_k$ entonces $x^{k+1} = y^{k+1}$ (a esto se le llama paso serio), si no pasa entonces $x^{k+1} = x^k$ (paso nulo).
7. Se actualiza el modelo

$$\hat{f}_{k+1}(y) = \max_k \{ \hat{f}_k(y), f(y^{k+1}) + \langle s^{k+1}, y - y^{k+1} \rangle \}. \quad (2.11)$$

8. $k = k + 1$ y se regresa al paso 2.

Esta es una versión básica del método *bundle*, otras implementaciones se pueden tomar en cuenta tales como búsqueda lineal, actualización del parámetro μ_k , etc [9].

Ahora bien, antes de realizar las pruebas de convergencia será más conveniente en teoría y en práctica trabajar con el dual del problema (2.10) como se verá a continuación. Primero, se reescribirá el modelo (2.11) de manera más conveniente. Se define el error de linealización con centro en x^k como sigue

$$e_i = f(x^k) - [f(y^i) - \langle s^i, x^k - y^i \rangle]$$

para $i = 1, \dots, k + 1$. Utilizando el error de linealización, nuestro modelo (2.11) se vuelve de la siguiente forma

$$\begin{aligned} \hat{f}_k(y) &= \max_{i=1, \dots, k+1} \{ f(y^i) + \langle s^i, y - y^i \rangle \} \\ &= \max_{i=1, \dots, k+1} \{ f(x^k) - e_i - \langle s^i, x^k - y^i \rangle + \langle s^i, y - y^i \rangle \} \\ &= f(x^k) + \max_{i=1, \dots, k+1} \{ -e_i + \langle s^i, y - x^k \rangle \} \end{aligned} \quad (2.12)$$

en donde todo lo usado es conocido debido a la acumulación de información que se tiene.

Al considerar el problema (2.10), usando (2.11) se tiene que

$$\begin{aligned} \min_{y \in \mathbb{R}^n} \hat{f}_k(y) + \frac{\mu_k}{2} \|y - x^k\|^2 &= \min_{y, r} \{ r + \frac{\mu_k}{2} \|y - x^k\|^2 \} \\ r &\geq f(x^k) - e_i + \langle s^i, y - x^k \rangle \end{aligned} \quad (2.13)$$

para $i = 1, \dots, k+1$. Introduciendo multiplicadores de Lagrange $\alpha \in \mathbb{R}_+^{k+1}$, el lagrangiano se escribe como sigue

$$\begin{aligned} L(y, r, \alpha) &= r + \frac{\mu_k}{2} \|y - x^k\|^2 + \sum_{i=1}^{k+1} \alpha_i (f(x^k) - e_i + \langle s^i, y - x^k \rangle - r) \\ &= \left(1 - \sum_{i=1}^{k+1} \alpha_i\right) r + \frac{\mu_k}{2} \|y - x^k\|^2 + \sum_{i=1}^{k+1} \alpha_i (f(x^k) - e_i + \langle s^i, y - x^k \rangle) \end{aligned}$$

Así, por convexidad

$$\min_{y,r} \max_{\alpha} L(y, r, \alpha) = \max_{\alpha} \min_{y,r} L(y, r, \alpha). \quad (2.14)$$

Dado que ambos lados son finitos se debe tener que $1 - \sum_{i=1}^{k+1} \alpha_i = 0$. Así, por condiciones de optimalidad de primer orden,

$$\nabla L_y(r, \alpha) = \mu_k (y - x^k) + \sum_{i=1}^{k+1} \alpha_i s^i = 0.$$

Entonces, se tiene que

$$\mu_k (y - x^k) = - \left(\sum_{i=1}^{k+1} \alpha_i \right). \quad (2.15)$$

Al denotar $\Delta^{k+1} = \{\alpha \in \mathbb{R}_+^{k+1} : \sum_{i=1}^{k+1} \alpha_i = 1\}$, combinado la ecuación (2.14) y (2.15), se obtiene que resolver el problema (2.9) es equivalente a resolver

$$\begin{aligned} \max_{\alpha \in \Delta^{k+1}} & \frac{\mu_k}{2} \left\| \frac{-\sum_{i=1}^{k+1} \alpha_i s^i}{\mu_k} \right\|^2 + \sum_{i=1}^{k+1} \alpha_i (f(x^k) - e_i + \langle s^i, -\frac{\sum_{i=1}^{k+1} \alpha_i s^i}{\mu_k} \rangle) \\ &= f(x^k) + \max_{\alpha \in \Delta^{k+1}} \frac{1}{2\mu_k} \left\| \sum_{i=1}^{k+1} \alpha_i s^i \right\|^2 - \sum_{i=1}^{k+1} \alpha_i e_i. \end{aligned} \quad (2.16)$$

Definición 2.3.1. Sea $\alpha \in \Delta^{k+1}$ una solución óptima del problema (2.16) en la iteración k , se define el subgradiente agregado y el error de linealización agregado respectivamente como sigue

$$\hat{s}^k = \sum_{i=1}^{k+1} \alpha_i s^i \quad \text{y} \quad \hat{e}_k = \sum_{i=1}^{k+1} \alpha_i e_i.$$

Lema 2.3.2. Sea $\alpha \in \Delta^{k+1}$ una solución para el problema (2.16), entonces

1. $\hat{s}^k \in \partial \hat{f}_k(y^{k+1})$.

$$2. \hat{f}_k(y^{k+1}) = f(x^k) - \frac{1}{\mu_k} \|\hat{s}^k\|^2 - \hat{e}_k.$$

$$3. \delta_k = \frac{1}{2\mu_k} \|\hat{s}^k\|^2 + \hat{e}_k.$$

Demostración [9]. 1) De la ecuación (2.15), $-\mu_k(y^{k+1} - x^k) = \hat{s}^k$. Así, dado que y^{k+1} es óptimo para (2.9)

$$0 \in \partial \hat{f}_k(y^{k+1}) + \mu_k(y^{k+1} - x^k),$$

donde $\mu_k(y^{k+1} - x^k)$ es la derivada de la parte cuadrática de la regularización. Entonces

$$-\mu_k(y^{k+1} - x^k) \in \partial \hat{f}_k(y^{k+1}).$$

2) Por convexidad, no hay brecha de dualidad entre (2.13) y (2.16). Así,

$$\begin{aligned} \hat{f}_k(y^{k+1}) + \frac{\mu_k}{2} \|y^{k+1} - x^k\|^2 &= f(x^k) - \frac{1}{2\mu_k} \|\hat{s}^k\|^2 - \hat{e}_k \\ \hat{f}_k(y^{k+1}) &= f(x^k) - \frac{\mu_k}{2} \left\| \frac{-1}{\mu_k} \hat{s}^k \right\|^2 - \frac{1}{2\mu_k} \|\hat{s}^k\|^2 - \hat{e}_k \\ &= f(x^k) - \frac{1}{\mu_k} \|\hat{s}^k\|^2 - \hat{e}_k. \end{aligned}$$

3) Usando el inciso 2 y la definición de δ_k

$$\begin{aligned} \delta_k &= f(x^k) - \hat{f}_k(y^{k+1}) - \frac{\mu_k}{2} \|y^{k+1} - x^k\|^2 \\ &= f(x^k) - \frac{\mu_k}{2} \|y^{k+1} - x^k\|^2 - f(x^k) + \frac{1}{\mu_k} \|\hat{s}^k\|^2 + \hat{e}_k \\ &= \frac{1}{2\mu_k} \|\hat{s}^k\|^2 + \hat{e}_k. \end{aligned}$$

□

Lema 2.3.3. *Para el subgradiente agregado y el error de linealización agregado, se tiene que*

$$\hat{s}^k \in \partial_{\hat{e}_k} f(x^k)$$

Demostración [9]. Usando el inciso 1 del lema anterior, $\hat{s}^k \in \partial \hat{f}_k(y^{k+1})$, y por construcción $\hat{f}_k \leq f$. Entonces,

$$\begin{aligned} f(y) \geq \hat{f}_k(y) &\geq \hat{f}_k(y^{k+1}) + \langle \hat{s}^k, y - y^{k+1} \rangle \\ &= f(x^k) - \frac{1}{\mu_k} \|\hat{s}^k\|^2 - \hat{e}_k + \langle \hat{s}^k, y - x^k + x^k - y^{k+1} \rangle \\ &= f(x^k) + \langle \hat{s}^k, y - x^k \rangle - \hat{e}_k + \langle \hat{s}^k, x^k - y^{k+1} - \frac{1}{\mu_k} \|\hat{s}^k\|^2 \rangle \\ &= f(x^k) + \langle \hat{s}^k, y - x^k \rangle - \hat{e}_k. \end{aligned}$$

□

Los lemas anteriores son muy importantes para las pruebas de convergencia que se darán en la siguiente sección. Antes de eso, se define la pieza lineal agregada como sigue

$$f_{\hat{e}_k}(y) = f(x^k) + \langle \hat{s}^k, y - x^k \rangle - \hat{e}_k.$$

Lema 2.3.4. *Para $f_{\hat{e}_k}$ se tiene lo siguiente:*

1. $f_{\hat{e}_k}(y) = \hat{f}_k(y^{k+1}) + \langle \hat{s}^k, y - y^{k+1} \rangle$
2. $f_{\hat{e}_k}(y) \leq \hat{f}_k(y)$

La demostración se puede encontrar en [9].

2.3.3. Convergencia del método *bundle*

Se supondrá que f es una función cerrada y finita en \mathbb{R}^n . Además, sea

$$K_s = \{k \in \mathbb{N} : \text{la } k\text{-ésima iteración fue un paso serio}\}.$$

Este análisis se dividirá en dos partes, primero se supondrá que $|K_s| = \infty$ y después que $|K_s| < \infty$. Además, se definirá $\hat{\delta} = 0$, es decir, el algoritmo hará iteraciones una infinidad de veces.

Lema 2.3.5. *Se considera el método *bundle* y sea $\bar{f} = \min_x f(x) > -\infty$. Entonces,*

$$\sum_{k \in K_s} \delta_k \leq \frac{f(x^0) - \bar{f}}{m} < \infty$$

La demostración se puede encontrar en [9].

Lema 2.3.6. *Supongamos que $f^* = \lim_{k \in K_s} f(x^k) > -\infty$ y $|K_s| = \infty$.*

1. *si $\sum_{k \in K_s} \frac{1}{\mu_k} = \infty$, entonces el cero es un punto de acumulación de $\{\hat{s}^k\}_{k \in K_s}$, esto es, $\liminf \|\hat{s}^k\| = 0$.*
2. *Si $\mu_k \geq c > 0$ y $\text{argmin}_x f(x) \neq \emptyset$, entonces la sucesión $\{x^k\}_{k \in K_s}$ es acotada.*

Demostración [9]. Por el inciso 3 del lema 2.3.2, se tiene que $0 \leq \frac{1}{2\mu_k} \|\hat{s}^k\|^2 = \delta_k - \hat{e}_k \leq \delta_k$. Así, por el lema anterior, dado que $f(x^k)$ es una sucesión decreciente cuando $k \in K_s$

$$\sum_{k \in K_s} \frac{\|\hat{s}^k\|^2}{2\mu_k} \leq \sum_{k \in K_s} \delta_k \leq \frac{f(x^0) - f^*}{m}.$$

Así, $\hat{s}^k \rightarrow 0$ sobre $k \in K_s$.

Para probar 2, sea $x^* \in \arg \min_y f(y)$. Por definición $f(x^*) \leq f(y)$ para toda $y \in \mathbb{R}^n$. Ahora, para $k \in K_s$,

$$\begin{aligned}
\|x^* - x^{k+1}\|^2 &= \|x^* - x^k\|^2 + 2\langle x^* - x^k, x^k - x^{k+1} \rangle + \|x^k - x^{k+1}\|^2 \\
&= \|x^* - x^k\|^2 + \frac{2}{\mu_k} \langle x^* - x^k, \hat{s}^k \rangle + \frac{1}{\mu_k^2} \|\hat{s}^k\|^2 \\
&= \|x^* - x^k\|^2 + \frac{2}{\mu_k} (\langle x^* - x^k, \hat{s}^k \rangle + \frac{1}{2\mu_k} \|\hat{s}^k\|^2) \\
&\leq \|x^* - x^k\|^2 + \frac{2}{\mu_k} (\hat{f}_k(x^*) - f(x^k) + \hat{e}_k + \frac{1}{2\mu_k} \|\hat{s}^k\|^2) \\
&\leq \|x^* - x^k\|^2 + \frac{2}{\mu_k} (f(x^*) - f(x^k) + \delta_k) \\
&\leq \|x^* - x^k\|^2 + \frac{2}{\mu_k} \delta_k,
\end{aligned}$$

donde se usó el hecho de que $\hat{s}^k \in \partial_{\hat{e}_k} f(x^k)$ y el inciso 3 del lema 2.3.2. Haciendo lo anterior k veces

$$\|x^* - x^{k+1}\|^2 \leq \|x^* - x^0\|^2 + 2 \sum_{i=1}^{k+1} \frac{\delta_i}{\mu_i} \leq \|x^* - x^0\|^2 + \frac{2}{c} \sum_{k \in K_s} \delta_k,$$

usando el lema 2.3.5 se tiene que el lado derecho es acotado y por tanto, la sucesión $\{x^k\}_{k \in K_s}$ es acotada. \square

Teorema 2.3.1. *Supongamos que $f^* = \lim_{k \in K_s} f(x^k) > -\infty$ y que $|K_s| = \infty$. Si la sucesión $\{\mu_k\}$ es acotada y creciente, entonces $\{x^k\}_{k \in K_s}$ tiene al menos un punto de acumulación que es óptimo.*

Demostración [9]. Por el lema 2.3.5 y 2.3.2 se sigue que $0 \leq \hat{e}_k \leq \delta_k \rightarrow 0$ para $k \in K_s$. Luego, por el inciso 1 del lema 2.3.2

$$\hat{s}^k \in \partial_{\hat{e}_k} f(x^k)$$

para cualquier $k \in K_s$. Luego, por el inciso 2 del lema 2.3.2, existe una subsucesión $\{\hat{s}^{n_k}\}$ que converge a cero. Dado que $\{x^k\}_{k \in K_s}$ es acotada, entonces $\{x^{n_k}\}$ también es acotada. Tomando (si es necesario una subsucesión) $x^{n_k} \rightarrow x^*$, se tiene que

$$(x^{n_k}, \hat{s}^{n_k}, \hat{e}_{n_k}) \rightarrow (x^*, 0, 0).$$

El corolario 2.4.2 garantiza que la correspondencia $(x, \epsilon) \mapsto \partial_\epsilon f(x)$ es continua, así $0 \in \partial f(x^*)$. \square

En este momento solo se ha tomado en cuenta $|K_s| = \infty$, ahora se verá el caso cuando $|K_s| < \infty$. Se supondrá que k_0 es la última iteración donde se tiene un paso serio.

Lema 2.3.7. *Sea x^{k_0} el último paso serio y $\{y^{k+1}\}_{k \geq k_0}$ la sucesión de pasos nulos. Entonces, para toda $k > k_0$ y $y \in \mathbb{R}^n$*

$$f(x^{k_0}) - \delta_k + \frac{\mu_k}{2} \|y - y^{k+1}\|^2 = \hat{f}_k(y^{k+1}) + \langle \hat{s}^k, y - y^{k+1} \rangle + \frac{\mu_k}{2} \|y - x^{k_0}\|^2 \quad (2.17)$$

Demostración [9]. Se verifica lo siguiente

$$\begin{aligned} \|y - x^{k_0}\|^2 &= \|y - y^{k+1} + y^{k+1} - x^{k_0}\|^2 \\ &= \|y - y^{k+1}\|^2 + 2\langle y - y^{k+1}, y^{k+1} - x^{k_0} \rangle + \|y^{k+1} - x^{k_0}\|^2 \\ &= \|y - y^{k+1}\|^2 - \frac{2}{\mu_k} \langle y - y^{k+1}, \hat{s}^k \rangle + \|y^{k+1} - x^{k_0}\|^2. \end{aligned}$$

Usando la definición de δ_k se obtiene

$$\begin{aligned} f(x^{k_0}) - \delta_k + \frac{\mu_k}{2} \|y - y^{k+1}\|^2 &= \hat{f}_k(y^{k+1}) + \frac{\mu_k}{2} (\|y^{k+1} - x^{k_0}\|^2 + \|y - y^{k+1}\|^2) \\ &= \hat{f}_k(y^{k+1}) + \langle y - y^{k+1}, \hat{s}^k \rangle + \frac{\mu_k}{2} \|y - x^{k_0}\|^2. \end{aligned}$$

\square

Teorema 2.3.2. *Sea x^{k_0} el último paso serio, y sea $\{y^{k+1}\}_{k \in K_s}$ la sucesión de pasos nulos. Si $\{\mu_k\}_{k > k_0}$ es no decreciente, entonces $\delta_k \rightarrow 0$.*

Demostración [9]. Sea $y = y^{k+2}$ en el lema anterior, entonces

$$\begin{aligned} f(x^{k_0}) - \delta_k + \frac{\mu_k}{2} \|y^{k+2} - y^{k+1}\|^2 &= \hat{f}_k(y^{k+1}) + \langle \hat{s}^k, y^{k+2} - y^{k+1} \rangle + \frac{\mu_k}{2} \|y^{k+2} - x^{k_0}\|^2 \\ &= \hat{f}_{\hat{e}_k}(y^{k+1}) + \frac{\mu_k}{2} \|y^{k+2} - x^{k_0}\|^2 \\ &\leq \hat{f}_{k+1}(y^{k+2}) + \frac{\mu_k}{2} \|y^{k+2} - x^{k_0}\|^2 \\ &\leq \hat{f}_{k+1}(y^{k+2}) + \frac{\mu_{k+1}}{2} \|y^{k+2} - x^{k_0}\|^2 \\ &= f(x^{k_0}) - \delta_{k+1}, \end{aligned}$$

donde la última igualdad se deduce de que $\mu_k \leq \mu_{k+1}$. Por lo tanto

$$\frac{\mu_k}{2} \|y^{k+2} - y^{k+1}\|^2 + \delta_{k+1} \leq \delta_k. \quad (2.18)$$

Ahora probará que la sucesión $\{y^k\}$ es acotada. Usando de nuevo el lema anterior con $y = x^{k_0}$,

$$\begin{aligned} f(x^{k_0}) - \delta_k + \frac{\mu_k}{2} \|x^{k_0} - y^{k+1}\|^2 &= \hat{f}_k(y^{k+1}) + \langle \hat{s}^k, x^{k_0} - y^{k+1} \rangle \\ &= f_{\hat{e}_k}(x^{k_0}) \leq \hat{f}_k(x^{k_0}) \leq f(x^{k_0}). \end{aligned}$$

Así,

$$\|x^{k_0} - y^{k+1}\|^2 \leq \frac{2\delta_k}{\mu_k} \leq \frac{2\delta_{k_0}}{\mu_{k_0}}$$

Dado que $\{\delta_k\}$ es decreciente y $\{\mu_k\}$ es no decreciente, $\{y^k\}$ es una sucesión acotada.

Por otro lado, sea C constante de Lipschitz para f y \hat{f}_k en $B(x^{k_0}, \frac{\delta_{k_0}}{\mu_{k_0}})$. Combinando

$$-m\delta_k \leq f(y^{k+1}) - f(x^{k_0})$$

y

$$\delta_k \leq f(x^{k_0}) - \hat{f}_k(y^{k+1}),$$

y sustituyendo $y = y^k$ en (2.11) en la iteración $k - 1$ se tiene que $f(y^k) = \hat{f}_k(y^k)$, así

$$\begin{aligned} (1 - m)\delta_k &\leq f(y^{k+1}) - \hat{f}_k(y^{k+1}) \\ &= f(y^{k+1}) - f(y^k) + \hat{f}_k(y^k) - \hat{f}_k(y^{k+1}) \\ &\leq 2C \|y^{k+1} - y^k\|^2. \end{aligned}$$

Combinando esta desigualdad con (2.18)

$$\delta_k - \delta_{k+1} \geq \frac{\mu_k}{2} \|y^{k+2} - y^{k+1}\|^2 \geq \frac{(1 - m)^2}{8C^2} \mu_k \delta_k^2 \geq \frac{(1 - m)^2}{8C^2} \mu_{k_0} \delta_{k+1}^2.$$

Así, sumando sobre $k \geq k_0$,

$$\frac{(1 - m)^2}{8C^2} \mu_{k_0} \sum_{k \geq k_0} \delta_k^2 \leq \sum_{k \geq k_0} (\delta_k - \delta_{k+1}) \leq \delta_{k_0},$$

lo cual implica que $\delta_k \rightarrow 0$. □

Teorema 2.3.3. *Sea x^{k_0} la última iteración con paso serio del método bundle. Si $\{\mu_k\}_{k \geq k_0}$ es no decreciente, entonces x^{k_0} es una solución óptima.*

Demostración [9]. Dado que las suposiciones son las mismas que en el teorema anterior, entonces $\delta_k \rightarrow 0$ implica que $\hat{e}_k \rightarrow 0$ y $\|\hat{s}^k\| \rightarrow 0$ por el inciso 3 del lema 2.3.2. De nuevo, por el inciso 1 del lema 2.3.2

$$\hat{s}^k \in \partial_{\hat{e}_k} f(x^{k_0})$$

para todo $k > k_0$. Por el teorema 1.3.5 que asegura que la correspondencia $(x, \epsilon) \mapsto \partial_\epsilon f(x)$ es continua, se concluye que $0 \in \partial f(x^{k_0})$. \square

2.4. Métodos derivados del operador próximo

En este capítulo se construyen métodos basados en un operador llamado el operador próximo, el cual está íntimamente relacionado con la regularización de Moreau-Yosida. Antes de definirlo se darán algunos conceptos importantes que implicarán propiedades interesantes de dicho operador.

2.4.1. Operadores no expansivos

Los operadores no expansivos son operadores continuos Lipchitz con constante $L = 1$. Estos juegan un papel fundamental en matemáticas aplicadas porque muchos de los problemas de análisis no lineal se reducen a encontrar puntos fijos de operadores no expansivos [7]. En esta sección se discutirán estos operadores y algunas de sus propiedades.

Antes de dar la definición de operadores no expansivos se presenta un resultado que se utilizará en pruebas posteriores.

Proposición 2.4.1. *Sean $x, y \in \mathbb{R}^n$ y $\alpha \in \mathbb{R}$. Entonces*

$$\|\alpha x + (1 - \alpha)y\|^2 + \alpha(1 - \alpha)\|x - y\|^2 = \alpha\|x\|^2 + (1 - \alpha)\|y\|^2$$

Demostración. La demostración se deduce utilizando propiedades de la norma. \square

A partir de ahora, al conjunto de puntos fijos del operador T se denotará por **Fix** T .

Definición 2.4.1. *Sea D un subconjunto no vacío de \mathbb{R}^n y sea $T : D \rightarrow \mathbb{R}^n$. Entonces T se dice*

1. Firmemente no expansivo si

$$\|Tx - Ty\|^2 + \|(Id - T)x - (Id - T)y\|^2 \leq \|x - y\|^2 \quad (2.19)$$

$\forall x, y \in D$.

2. No expansivo si es Lipschitz con constante 1, es decir, $\forall x, y \in D$ se cumple

$$\|Tx - Ty\| \leq \|x - y\|. \quad (2.20)$$

3. Cuasi no expansivo si $\forall x \in D$ y $\forall y \in \text{Fix } T$ se tiene que

$$\|Tx - y\| \leq \|x - y\|. \quad (2.21)$$

Es claro que 2.19 implica 2.20 y 2.20 implica 2.21.

Proposición 2.4.2. Sea D un subconjunto no vacío de \mathbb{R}^n y sea $T : D \rightarrow \mathbb{R}^n$. Entonces los siguientes enunciados son equivalentes:

1. T es firmemente no expansivo.
2. $Id - T$ es firmemente no expansivo.
3. $2T - Id$ es no expansivo.
4. $\|Tx - Ty\|^2 \leq \langle x - y, Tx - Ty \rangle \quad \forall x, y \in D$.
5. $0 \geq \langle Tx - Ty, (Id - T)x - (Id - T)y \rangle \quad \forall x, y \in D$.
6. $\|Tx - Ty\| \leq \|\alpha(x - y) + (1 - \alpha)(Tx - Ty)\| \quad \forall x, y \in D$.

Demostración [7]. 1) \Leftrightarrow 2) Es inmediato de la definición de operador firmemente no expansivo.

2) \Leftrightarrow 3) Sean $x, y \in D$, se define

$$R = 2T - Id,$$

$$\mu = \|Tx - Ty\|^2 + \|(Id - T)x - (Id - T)y\|^2 - \|x - y\|^2$$

y

$$\nu = \|Rx - Ry\|^2 + \|x - y\|^2.$$

Luego, por la proposición 2.4.1 se tiene que

$$\begin{aligned}\|Rx - Ry\|^2 &= \|2(Tx - Ty) + (1 - 2)(x - y)\|^2 \\ &= 2\|Tx - Ty\|^2 - \|x - y\|^2 + 2\|(Id - T)x - (Id - T)y\|^2.\end{aligned}$$

Por lo tanto $\nu = 2\mu$. Así, R es no expansivo si y solo si

$$\begin{aligned}\nu \leq 0 &\Leftrightarrow \mu \leq 0 \\ &\Leftrightarrow T \text{ es no firmemente no expansivo.}\end{aligned}$$

1) \Leftrightarrow 4) Se escribe

$$\|(Id - T)x - (Id - T)y\|^2 = \|x - y\|^2 + \|Tx - Ty\|^2 - 2\langle x - y, Tx - Ty \rangle$$

y se sustituye en 2.19.

4) \Leftrightarrow 5) Es inmediato de 4).

5) \Leftrightarrow 6) De las propiedades de norma y producto interior se tiene el siguiente resultado:

$$\begin{aligned}\langle x, y \rangle \leq 0 &\Leftrightarrow \|x\| \leq \|x - \alpha y\| \quad \forall \alpha \in \mathbb{R}_+ \\ &\Leftrightarrow \|x\| \leq \|x - \alpha y\| \quad \forall \alpha \in [0, 1].\end{aligned}$$

Este resultado implica el inciso 6. □

Ahora se define el operador proyección que será de mucha utilidad en este trabajo.

Definición 2.4.2. Sea D un subconjunto no vacío de \mathbb{R}^n , sea $x \in \mathbb{R}^n$ y $p \in D$. Entonces p se dice **proyección** de x en D si $\|x - p\| = \inf\{\|y - x\| : y \in D\} = d_D(x)$.

Si cualquier punto en \mathbb{R}^n tiene exactamente una proyección sobre D , entonces se dice que D es un **conjunto de Chebyshev**. En este caso, el operador proyección sobre D es el operador definido por

$$\Pi_D(x) = \arg \min_{v \in D} \|x - v\|.$$

Proposición 2.4.3. Sea D un conjunto de Chebyshev de \mathbb{R}^n . Entonces Π_D es continua .

La demostración se puede encontrar en [7].

Teorema 2.4.1. *Sea D un subconjunto no vacío de \mathbb{R}^n . Entonces D es un conjunto de Chebyshev si para todo $x, p \in \mathbb{R}^n$*

$$p = \Pi_D(x) \Leftrightarrow p \in D \text{ y } \langle y - p, x - p \rangle \leq 0 \quad \forall y \in D.$$

La demostración se puede encontrar en [7].

Una propiedad muy interesante del operador proyección es que este es firmemente no expansivo como lo enuncia el siguiente resultado.

Proposición 2.4.4. *Sea D un subconjunto no vacío, cerrado y convexo de \mathbb{R}^n . Entonces la proyección Π_D es firmemente no expansiva.*

Demostración [7]. Sean $x, y \in \mathbb{R}^n$. Por el teorema anterior se tiene que

$$\langle \Pi_D y - \Pi_D x, x - \Pi_D x \rangle \leq 0$$

y además

$$\langle \Pi_D x - \Pi_D y, y - \Pi_D y \rangle \leq 0.$$

Sumando estas dos desigualdades se obtiene

$$\|\Pi_D x - \Pi_D y\|^2 \leq \langle x - y, \Pi_D x - \Pi_D y \rangle.$$

El resultado deseado se sigue de la proposición 2.4.2 inciso 4. \square

Ahora bien, si D es un conjunto convexo y cerrado, el operador proyección es no expansivo sobre el conjunto $\text{Fix } \Pi_D$ dado que

$$\text{Fix } \Pi_D = D.$$

Se discutirá a continuación bajo que circunstancias el conjunto de puntos fijos de un operador T es cerrado y convexo.

Proposición 2.4.5. *Sea D un subconjunto no vacío y convexo de \mathbb{R}^n , y sea $T : D \rightarrow \mathbb{R}^n$ un operador cuasi no expansivo. Entonces $\text{Fix } T$ es un conjunto convexo.*

Demostración [7]. Sean $x, y \in \text{Fix } T$, $\alpha \in (0, 1)$ y sea $z = \alpha x + (1 - \alpha)y$. Entonces $z \in D$ y por la proposición 2.4.1

$$\begin{aligned} \|Tz - z\|^2 &= \|\alpha(Tz - x) + (1 - \alpha)(Tz - y)\|^2 \\ &= \alpha\|Tz - x\|^2 + (1 - \alpha)\|Tz - y\|^2 - \alpha(1 - \alpha)\|x - y\|^2 \\ &\leq \alpha\|z - x\|^2 + (1 - \alpha)\|z - y\|^2 - \alpha(1 - \alpha)\|x - y\|^2 \\ &= \|\alpha(z - x) + (1 - \alpha)(z - y)\|^2 \\ &= 0. \end{aligned}$$

Por lo tanto, $z \in \text{Fix } T$. \square

Proposición 2.4.6. *Sea D un subconjunto no vacío y cerrado de \mathbb{R}^n y sea $T : D \rightarrow \mathbb{R}^n$ continua. Entonces $\text{Fix } T$ es cerrado.*

Demostración [7]. Sea $\{x_n\}$ una sucesión en $\text{Fix } T$ que converge a un punto en \mathbb{R}^n . Entonces $x \in D$ porque D es cerrado, mientras que $Tx_n \rightarrow Tx$ por la continuidad de T . Por otro lado, dado que $\{x_n\}$ está en $\text{Fix } T$, $Tx_n \rightarrow Tx$. Así, $Tx = x$. \square

Teorema 2.4.2 (Browder-Gohde-Kirk). *Sea D un subconjunto acotado, cerrado, convexo y diferente de vacío de \mathbb{R}^n y sea $T : D \rightarrow D$ un operador no expansivo. Entonces $\text{Fix } T \neq \emptyset$.*

La demostración se puede encontrar en [7].

2.4.2. Sucesiones Féjer monótonas

La siguiente definición es estudio central de varios métodos iterativos, en particular con la construcción de puntos fijos de operadores no expansivos.

Definición 2.4.3. *Sea D un subconjunto no vacío de \mathbb{R}^n y sea $\{x_n\}_{n \in \mathbb{N}}$ una sucesión en \mathbb{R}^n . Entonces $\{x_n\}_{n \in \mathbb{N}}$ es **Féjer monótona** con respecto a D si $\forall x \in D$ y $n \in \mathbb{N}$ se tiene*

$$\|x_{n+1} - x\| \leq \|x_n - x\|. \quad (2.22)$$

Si D es un subconjunto no vacío de \mathbb{R}^n , y $T : D \rightarrow \mathbb{R}^n$ un operador no expansivo tal que $\text{Fix } T \neq \emptyset$. Haciendo $x_{n+1} = Tx_n$ y $x_0 \in D$ se concluye que $\{x_n\}_{n \in \mathbb{N}}$ es Féjer monótona con respecto a $\text{Fix } T$.

Ahora se tratarán algunas propiedades básicas.

Proposición 2.4.7. *Sea $\{x_n\}_{n \in \mathbb{N}}$ una sucesión en \mathbb{R}^n y D un subconjunto no vacío de \mathbb{R}^n . Supongamos que $\{x_n\}_{n \in \mathbb{N}}$ es Féjer monótona con respecto a D . Entonces lo siguiente se cumple:*

1. $\{x_n\}_{n \in \mathbb{N}}$ es acotada.
2. Para cualquier $x \in D$, $\{\|x_n - x\|\}_{n \in \mathbb{N}}$ converge.
3. $\{d_D(x_n)\}_{n \in \mathbb{N}}$ es decreciente y converge.

Demostración. 1) Sea $x \in D$. Entonces (2.22) implica que $\{x_n\}_{n \in \mathbb{N}}$ está en $B_{\|x_0 - x\|}(x)$.
 2) Es inmediato de (2.22).
 3) Se toma el ínfimo en (2.22) sobre $x \in D$, entonces se tiene que $d_D(x_{n+1}) \leq d_D(x_n)$.
 \square

Dado un operador no expansivo T , es posible que la sucesión generada iterativamente por $x_{n+1} = Tx_n$ no produzca un punto fijo de T . Un simple ejemplo de esto es cuando $T = -Id$ y $x_0 \neq 0$. En este caso, es claro que la propiedad $Tx_n - x_n \rightarrow 0$ no se tiene.

Ahora se dará la definición de una clase especial de operadores no expansivos.

Definición 2.4.1. Sea D un subconjunto no vacío de \mathbb{R}^n , sea $T : \mathbb{R}^n \rightarrow \mathbb{R}^n$ un operador no expansivo y sea $\alpha \in (0, 1)$. Entonces T es α -promediado si existe un operador $R : D \rightarrow \mathbb{R}^n$ tal que $T = (1 - \alpha)Id + \alpha R$.

Como observaciones, si D es un subconjunto no vacío de \mathbb{R}^n y $T : D \rightarrow \mathbb{R}^n$ entonces:

1. Si T es α -promediado entonces es no expansivo.
2. Si T es no expansivo entonces no necesariamente es α -promediado. Como ejemplo se tiene el caso en que $T = -Id : \mathbb{R}^n \rightarrow \mathbb{R}^n$.
3. De la proposición 2.4.2 se verifica que T es firmemente no expansivo si y solo si T es 1/2-promediado.

Tomando en cuenta estos operadores se tiene el siguiente resultado.

Proposición 2.4.8. Sea $\alpha \in (0, 1)$, $T : \mathbb{R}^n \rightarrow \mathbb{R}^n$ un operador α -promediado tal que $\text{Fix } T \neq \emptyset$ y sea $\{\lambda_n\}_{n \in \mathbb{N}}$ una sucesión en $[0, 1/\alpha]$ tal que $\sum_{n \in \mathbb{N}} \lambda_n(1 - \alpha\lambda_n) = +\infty$. Partiendo de $x_0 \in \mathbb{R}^n$ se define la iteración

$$x_{n+1} = x_n + \lambda_n(Tx_n - x_n).$$

Entonces lo siguiente se cumple:

1. $\{x_n\}_{n \in \mathbb{N}}$ es Féjer monótona con respecto a $\text{Fix } T$.
2. $\{Tx_n - x_n\}_{n \in \mathbb{N}}$ converge a 0.
3. $\{x_n\}_{n \in \mathbb{N}}$ converge a un punto en $\text{Fix } T$.

La demostración se realizará más adelante para un caso particular de T .

Utilizando $\alpha = 1/2$ en el corolario anterior, si $T : \mathbb{R}^n \rightarrow \mathbb{R}^n$ es un operador firmemente no expansivo tal que $\text{Fix } T \neq \emptyset$, $x_0 \in \mathbb{R}^n$ y $x_{n+1} = Tx_n$. Entonces $\{x_n\}_{n \in \mathbb{N}}$ converge a un punto en $\text{Fix } T$.

2.4.3. Operador próximo

Sea $f : \mathbb{R}^n \rightarrow (-\infty, +\infty]$ una función propia, convexa y cerrada. Se define el operador próximo $\text{prox}_f : \mathbb{R}^n \rightarrow \mathbb{R}^n$ de f como

$$\text{prox}_f(v) = \arg \min_x \left\{ f(x) + \frac{1}{2} \|x - v\|^2 \right\} \quad (2.23)$$

donde $\|\cdot\|$ es la norma euclidiana. Es decir, el operador próximo es aquel que cumple

$$f(\text{prox}_f(v)) + \frac{1}{2} \|\text{prox}_f(v) - v\|^2 = F(v).$$

De aquí que el operador próximo y la regularización de Moreau-Yosida estén relacionados.

El operador próximo es estrictamente convexo dado que la norma euclidiana lo es y tiene valor finito en todas partes. Así, tiene un único minimizador para cualquier $v \in \mathbb{R}^n$ (incluso cuando $\text{dom } f \subsetneq \mathbb{R}^n$).

A menudo se encuentra el operador próximo de una función escalada λf ($\lambda > 0$) el cual puede ser expresado por

$$\text{prox}_{\lambda f}(v) = \arg \min_x \left\{ f(x) + \frac{1}{2\lambda} \|x - v\|^2 \right\} \quad (2.24)$$

En la figura (2.1) se representa lo que hace el operador próximo. Las líneas delgadas negras representan las curvas de nivel de la función convexa f y la línea gruesa representa la frontera del dominio esencial de f . Evaluar los puntos azules en el operador próximo tiene como resultado los puntos en rojo respectivamente. Los tres puntos en el dominio esencial de la función se mantienen en el dominio y se mueven hacia el mínimo de la función, mientras que los puntos fuera del dominio esencial se mueven hacia la frontera de este y hacia el mínimo de la función. El parámetro λ controla la proporción a la cual el operador próximo mapea puntos del dominio esencial de f hacia el punto mínimo.

Por la definición del operador próximo, $\text{prox}_f(v)$ es un punto que se encuentra entre $x^* = \min f$ y cerca de v . En el caso de $\text{prox}_{\lambda f}(v)$, el parámetro λ puede ser interpretado como un peso relativo entre estos dos términos.

Cuando f es la función indicadora de $C \subseteq \mathbb{R}^n$, es decir

$$\mathbb{I}_C(x) = \begin{cases} 0 & x \in C \\ +\infty & x \notin C \end{cases}$$

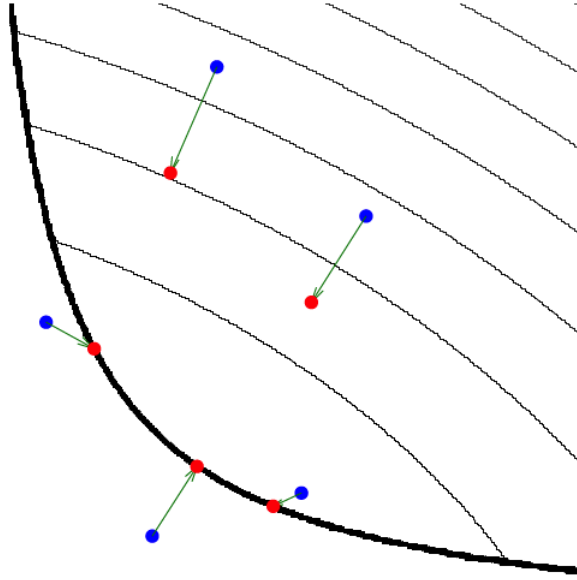


Figura 2.1: Operador próximo evaluado en varios puntos

entonces

$$\text{prox}_f(v) = \operatorname{argmin}_{x \in C} \{\|x - v\|\}. \quad (2.25)$$

Así, el operador próximo se puede ver como una proyección generalizada. Esta perspectiva sugiere varias propiedades que se espera que el operador próximo cumpla.

Si $f(x, y) = \phi(x) + \eta(y)$, es decir, f puede ser separada por dos funciones que dependen de cada una de las variables entonces

$$\text{prox}_f(v_1, v_2) = (\text{prox}_\phi(v_1), \text{prox}_\eta(v_2)).$$

Esto se observa directamente de la definición del operador próximo y del hecho que f es separable.

Una generalización de este resultado se tiene si $f(x) = \sum_{i=1}^n f_i(x_i)$, entonces, el operador próximo sería

$$(\text{prox}_f(v))_i = \text{prox}_{f_i}(v_i). \quad (2.26)$$

En otras palabras, este caso se reduce a evaluar el operador próximo de funciones escalares.

Ahora, si $f(x) = \alpha\phi(x) + b$, con $\alpha > 0$, entonces

$$\text{prox}_{\lambda f} = \text{prox}_{\alpha\lambda\phi}(v). \quad (2.27)$$

la prueba de esto es inmediata de la definición del operador próximo.

Proposición 2.4.9. *Sea f una función propia, convexa y cerrada, definida sobre \mathbb{R}^n . Entonces $p = \text{prox}_f(x)$ si y solo si $x - p \in \partial f(p)$.*

Demostración. Supongamos que $p = \text{prox}_f(x)$, esto se cumple si y solo si

$$0 \in \partial f(p) + (p - x)$$

$$\Leftrightarrow x - p \in \partial f(p).$$

□

2.4.4. Puntos fijos

Los puntos fijos del operador próximo juegan un papel fundamental cuando se trata de minimizar la función f como lo enuncia el siguiente resultado.

Teorema 2.4.3. *x^* minimiza a f si y solo si*

$$x^* = \text{prox}_f(x^*)$$

Demostración [12]. Supongamos que x^* minimiza a f , es decir, $f(x) \geq f(x^*)$ para cualquier x . Entonces

$$f(x) + \frac{1}{2}\|x - x^*\|^2 \geq f(x^*) = f(x^*) + \frac{1}{2}\|x^* - x^*\|^2$$

y por lo tanto x^* minimiza a $f(x) + \frac{1}{2}\|x - x^*\|^2$. Así, tenemos que $x^* = \text{prox}_f(x^*)$. Para mostrar el recíproco, como \hat{x} minimiza $f(x) + \frac{1}{2}\|x - v\|^2$ si y solo si

$$0 \in \partial f(\hat{x}) + (\hat{x} - v),$$

entonces si x^* es punto fijo de prox_f , $0 \in \partial f(x^*)$. □

Dado que los minimizadores de f son puntos fijos del operador prox_f , el problema de minimizar f se convierte en un problema de punto fijo. Si el operador prox_f fuera una contracción, la sucesión definida por $x_{n+1} = \text{prox}_f(x_n)$ con $x_0 \in S$ converge al único punto fijo del operador prox_f . El operador próximo no necesariamente es una contracción, pero es firmemente no expansivo, que es una propiedad suficiente para algoritmos de punto fijo como se observó anteriormente.

Proposición 2.4.10. *Sea f una función propia, convexa y cerrada. Entonces prox_f y $\text{Id} - \text{prox}_f$ son firmemente no expansivos.*

Demostración. Sean $x, y \in \mathbb{R}^n$ y sea $p = \text{prox}_f(x)$ y $q = \text{prox}_f(y)$. Entonces, por la proposición 2.4.9

$$x - p \in \partial f(p) \text{ y } y - q \in \partial f(q).$$

Por definición de subgradiente se tienen las siguientes desigualdades

$$\langle q - p, x - p \rangle + f(p) \leq f(q)$$

y

$$\langle p - q, y - q \rangle + f(q) \leq f(p).$$

Sumando estas dos desigualdades se obtiene que $0 \leq \langle p - q, (x - p) - (y - q) \rangle$ y (2.4.2) implica lo deseado. \square

2.4.5. Descomposición de Moreau

Definición 2.4.4. Sea f una función propia y convexa. La función conjugada f^* de f se define como

$$f^*(x^*) = \sup_x \{ \langle x, x^* \rangle - f(x) \}. \quad (2.28)$$

La función conjugada es el supremo (de forma puntual) de todas las funciones afines que están por debajo de f . La función conjugada tiene muchas propiedades útiles. En este trabajo solo se requiere la siguiente.

Lema 2.4.1. Sea $f : \mathbb{R}^n \rightarrow \mathbb{R}$ una función propia, convexa y cerrada; y sea $x \in \mathbb{R}^n$, entonces $x \in \partial f^*(x^*)$ si y solo si $x^* \in \partial f(x)$

La demostración se puede encontrar en [27].

Utilizando la función conjugada de f se tiene que la siguiente relación siempre se cumple

$$x = \text{prox}_f(x) + \text{prox}_{f^*}(x). \quad (2.29)$$

Esta propiedad es conocida como la descomposición de Moreau, y es una de las mayores relaciones entre el operador próximo y el dual de la función f .

Para probar esta igualdad, sea $u = \text{prox}_f(x)$, entonces por la proposición 2.4.9 y el lema 2.4.1 esto ocurre si y solo si

$$\begin{aligned} x - u \in \partial f(u) &\Leftrightarrow u \in \partial f^*(x - u) \\ &\Leftrightarrow x - u = \text{prox}_{f^*}(x). \end{aligned}$$

Así

$$x = u + (x - u) = \text{prox}_f(x) + \text{prox}_{f^*}(x).$$

La descomposición de Moreau nos brinda una manera simple de calcular el operador próximo de una función f en términos del operador próximo de f^* . Un ejemplo importante es cuando $f = \|\cdot\|$ es una norma en general, se sabe que $f^* = \mathbb{I}_{\mathcal{B}}$, donde

$$\mathcal{B} = \{x : \|x\|_* \leq 1\}$$

es la bola unitaria para la norma dual $\|\cdot\|_*$ definida por

$$\|x\|_* = \max\{\langle y, x \rangle : \|y\| \leq 1\}.$$

Luego, se tiene lo siguiente

$$\text{prox}_{f^*}(v) = \arg \min_{\|x\|_* \leq 1} \|x - v\| = \Pi_{\mathcal{B}}(v)$$

que es la proyección del punto v sobre el conjunto \mathcal{B} . Así

$$\text{prox}_f(v) = v - \Pi_{\mathcal{B}}(v). \quad (2.30)$$

En otras palabras, se puede evaluar fácilmente el operador próximo de f si se sabe como proyectar sobre \mathcal{B} .

2.4.6. Método de punto próximo

El método de minimización próxima o algoritmo de punto próximo es

$$x_{k+1} = \text{prox}_{\lambda f}(x_k) \quad (2.31)$$

donde $f : \mathbb{R}^n \rightarrow (-\infty, +\infty]$ es una función propia, convexa y cerrada, k es el contador de iteraciones y x_k denota la k -ésima iteración del algoritmo. Más adelante se mostrará que si f tiene un mínimo, entonces $\{x_k\}$ converge al conjunto de minimizadores de f y $\{f(x_k)\}$ converge al valor óptimo.

En cada iteración de este algoritmo se requiere minimizar la función f más un término cuadrático, por lo cual el algoritmo de punto próximo será útil cuando minimizar la función f sea difícil, pero fácil (o al menos más fácil) minimizar la función f más el término cuadrático.

2.4.7. Convergencia del método de punto próximo

Antes de dar la prueba de convergencia del método de punto próximo se definirán y probarán algunos resultados.

Definición 2.4.5. Sea $f : \mathbb{R}^n \rightarrow (-\infty, +\infty]$ y sea $\{x_n\}_{n \in \mathbb{N}}$ una sucesión en $\text{dom} f$. Se dice que $\{x_n\}_{n \in \mathbb{N}}$ es una sucesión minimizante de f si la sucesión $\{f(x_k)\}$ converge a $\inf f$.

Ahora se tiene el siguiente resultado de sucesiones en \mathbb{R}^n .

Lema 2.4.2. Sea $\{x_n\}_{n \in \mathbb{N}}$ una sucesión en \mathbb{R}^n y sea C un subconjunto no vacío de \mathbb{R}^n . Supongamos que para cada $x \in C$, la sucesión $\{\|x_n - x\|\}$ converge y que cualquier punto límite de $\{x_n\}_{n \in \mathbb{N}}$ está en C . Entonces $\{x_n\}$ converge en C .

Demostración. Como $\{x_n\}$ es acotada entonces posee al menos un punto límite, sea $x' \in C$ dicho punto límite. Ahora bien, al tomar una subsucesión $\{x_{n_k}\}$ que converge a x' se tiene que $\|x_{n_k} - x'\| \rightarrow 0$, pero por hipótesis la sucesión $\|x_n - x'\|$ es convergente, por lo tanto

$$\|x_n - x'\| \rightarrow 0.$$

Así, la sucesión $\{x_n\}$ converge en C . \square

Teorema 2.4.4 (Algoritmo de punto próximo). Sea f una función propia y cerrada en \mathbb{R}^n tal que $\arg \min f \neq \emptyset$, sea $\{\gamma_n\}_{n \in \mathbb{N}}$ una sucesión de números positivos tal que $\sum_{n \in \mathbb{N}} \gamma_n = +\infty$, y sea $x_0 \in \mathbb{R}^n$. Para todos los naturales se toma la iteración

$$x_{n+1} = \text{prox}_{\gamma_n f}(x_n) \quad (2.32)$$

Entonces se cumple lo siguiente:

1. $\{x_n\}_{n \in \mathbb{N}}$ es una sucesión minimizante de f ; de manera más precisa, $f(x_n) \downarrow \inf f$.
2. $\{x_n\}_{n \in \mathbb{N}}$ converge a un punto en $\arg \min f$.

Demostración [7]. 1. Sea $z \in \arg \min f$. De la proposición 2.4.9 y (2.31) se obtiene que

$$x_n - x_{n+1} \in \gamma_n \partial f(x_{n+1}).$$

Por la definición de subgradiente, para todo $n \in \mathbb{N}$

$$\langle z - x_{n+1}, x_n - x_{n+1} \rangle / \gamma_n \leq f(z) - f(x_{n+1}) \quad (2.33)$$

y

$$0 \leq \langle x_n - x_{n+1}, x_n - x_{n+1} \rangle / \gamma_n \leq f(x_n) - f(x_{n+1}). \quad (2.34)$$

Por lo tanto, (2.33) implica que

$$\begin{aligned}\|x_{n+1} - z\| &= \|x_n - z\|^2 + 2\langle x_n - z, x_{n+1} - x_n \rangle + \|x_{n+1} - x_n\|^2 \\ &= \|x_n - z\|^2 - \|x_{n+1} - z\|^2 + 2\langle x_{n+1} - z, x_{n+1} - x_n \rangle \\ &\leq \|x_n - z\|^2 - 2\gamma_n(f(x_{n+1}) - f(z)).\end{aligned}$$

Así se prueba que $\{x_n\}_{n \in \mathbb{N}}$ es Fejér monótona con respecto a $\arg \min f$. Al realizar iterativamente la desigualdad anterior, para todo $n \in \mathbb{N}$

$$\sum_{i=1}^{n+1} 2\gamma_i(f(x_{i+1}) - f(z)) \leq \|x_0 - z\|^2,$$

de aquí se tiene que

$$\sum_{n \in \mathbb{N}} \gamma_n(f(x_{n+1}) - f(z)) < +\infty.$$

Dado que $\sum_{n \in \mathbb{N}} \gamma_n = +\infty$, se tiene que $\liminf f(x_n) = 0$ y por la ecuación (2.34) $f(x_n) \downarrow f(z)$.

2. Dado que la sucesión $\{x_n\}_{n \in \mathbb{N}}$ es Fejér monótona con respecto a $\arg \min f$ entonces la sucesión $\|x_n - x\|$ es convergente para cualquier $x \in \arg \min f$ y además es una sucesión acotada, por lo tanto tiene al menos un punto límite, sea x' este punto límite y $\{x_{n_k}\}$ tal que $x_{n_k} \rightarrow x'$.

Luego, por el inciso 1 y el hecho de que f es cerrada se tiene que si $z \in \arg \min f$

$$f(z) \leq f(x') \leq \liminf f(x_{n_k}) = f(z).$$

Así, cualquier punto límite de $\{x_n\}$ está en $\arg \min f$. Por el lema 2.4.2 la sucesión $\{x_n\}$ converge a un elemento de $\arg \min f$. \square

2.4.8. Método de gradiente próximo

Considerar el siguiente problema

$$\min h(x) = f(x) + g(x),$$

donde $f : \mathbb{R}^n \rightarrow \mathbb{R}$ y $g : \mathbb{R}^n \rightarrow (-\infty, +\infty]$ son funciones propias, cerradas, convexas y f es diferenciable. En esta forma se puede partir el problema en dos, en donde una parte es diferenciable. La forma en la cual se parte el problema no es única y puede llevar a diferentes implementaciones del método de gradiente próximo.

El método de gradiente próximo es

$$x_{k+1} = \text{prox}_{\lambda_k g}(x_k - \lambda_k \nabla f(x_k))$$

donde $\lambda_k > 0$ es el tamaño de paso.

En la siguiente sección se mostrará que cuando ∇f es Lipschitz con constante $L \in (0, 1)$, el método converge con una rapidez de $o(1/k)$, cuando se elige un tamaño de paso fijo $\lambda_k = \lambda \in (0, 1/L]$.

2.4.9. Convergencia del método de gradiente próximo

Primero, supongamos que ∇f es Lipschitz con constante $L \in (0, 1)$, es decir

$$\| \nabla f(x) - \nabla f(y) \| \leq L \| x - y \|$$

para toda $x, y \in \mathbb{R}^n$. Además, supongamos que $x^* \in \arg \min h \neq \emptyset$. Se define

$$G_\lambda(x) = \frac{1}{\lambda}(x - \text{prox}_{\lambda g}(x - \lambda \nabla f(x))).$$

Esta expresión tiene la siguiente propiedad

$$G_\lambda(x) - \nabla f(x) \in \partial g(x - \lambda G_\lambda(x)),$$

En efecto, por la proposición 2.4.9

$$\begin{aligned} \text{prox}_{\lambda g}(x - \lambda \nabla f(x)) = p &\Leftrightarrow 0 \in \partial g(p) + \frac{1}{\lambda}(p - (x - \lambda \nabla f(x))) \\ &\Leftrightarrow \frac{(x - \lambda \nabla f(x)) - p}{\lambda} \in \partial g(p) \\ &\Leftrightarrow G_\lambda(x) - \nabla f(x) \in \partial g(x - \lambda G_\lambda(x)). \end{aligned}$$

Ahora se probará que $h(x^k) - h(x^*)$ decrece con una rapidez $1/k$ si un tamaño de paso $\lambda_k = 1/L$ es usado. Del hecho de que ∇f es Lipschitz se tiene que

$$\begin{aligned} f(y) &= f(x) + \langle \nabla f(x), v \rangle + \int_0^1 \langle \nabla f(x + \lambda v) - \nabla f(x), v \rangle d\lambda \\ &\leq f(x) + \langle \nabla f(x), v \rangle + \int_0^1 \| \nabla f(x + \lambda v) - \nabla f(x) \| \| v \| d\lambda \\ &\leq f(x) + \langle \nabla f(x), v \rangle + \int_0^1 L\lambda \| v \|^2 d\lambda \\ &= f(x) + \langle \nabla f(x), v \rangle + \frac{L}{2} \| v \|^2 \end{aligned}$$

donde $v = y - x$. Sustituyendo $y = x - \lambda G_\lambda(x)$ en la desigualdad anterior se obtiene

$$f(x - \lambda g_\lambda(x)) \leq f(x) - \lambda \langle \nabla f(x), G_\lambda(x) \rangle + \frac{\lambda^2 L}{2} \| G_\lambda(x) \|^2.$$

Dado que $0 < t \leq 1/L$ se tiene que

$$f(x - \lambda g_\lambda(x)) \leq f(x) - \lambda \langle \nabla f(x), G_\lambda(x) \rangle + \frac{\lambda}{2} \|G_\lambda(x)\|^2.$$

Luego, al tomar la función h y la desigualdad anterior

$$\begin{aligned} h(x - \lambda G_\lambda(x)) &\leq f(x) - \lambda \langle \nabla f(x), G_\lambda(x) \rangle + \frac{\lambda}{2} \|G_\lambda(x)\|^2 + g(x - \lambda G_\lambda(x)) \\ &\leq f(z) + \langle \nabla f(z), x - z \rangle - \lambda \langle \nabla f(x), G_\lambda(x) \rangle + \frac{\lambda}{2} \|G_\lambda(x)\|^2 \\ &\quad + g(z) + \langle G_\lambda(x) - \nabla f(x), x - z - \lambda G_\lambda(x) \rangle \\ &= f(z) + g(z) + \langle G_\lambda(x), x - z \rangle - \frac{\lambda}{2} \|G_\lambda(x)\|^2 \\ &= h(z) + \langle G_\lambda(x), x - z \rangle - \frac{\lambda}{2} \|G_\lambda(x)\|^2, \end{aligned} \quad (2.35)$$

donde la segunda desigualdad se tiene del hecho de que f y g son convexas y $G_\lambda(x) - \nabla f(x) \in \partial g(x - \lambda G_\lambda(x))$. Ahora, sea $x^+ = x - \lambda G_\lambda(x)$, la desigualdad anterior con $z = x$ implica que

$$h(x^+) \leq h(x) - \frac{\lambda}{2} \|G_\lambda(x)\|^2,$$

lo cual nos dice que se tiene un método de descenso.

Igualmente, al tomar la desigualdad (2.35) con $z = x^* \in \arg \min h$ se tiene que

$$\begin{aligned} 0 \leq h(x^+) - h(x^*) &\leq \langle G_\lambda(x), x - x^* \rangle - \frac{\lambda}{2} \|G_\lambda(x)\|^2 \\ &= \frac{1}{2\lambda} (\|x - x^*\|^2 - \|x - x^* - \lambda G_\lambda(x)\|^2) \\ &= \frac{1}{2\lambda} (\|x - x^*\|^2 - \|x^+ - x^*\|^2). \end{aligned}$$

Por lo tanto

$$\|x - x^*\| \leq \|x^+ - x^*\|,$$

es decir, la distancia al conjunto $\arg \min h$ decrece.

Al sumar las desigualdades anteriores con $x = x_{i-1}$, $x_+ = x_i$ y $t = 1/L$ se obtiene que

$$\begin{aligned} \sum_{i=1}^k (f(x_i) - f(x^*)) &\leq \frac{1}{2\lambda} \sum_{i=1}^k (\|x_{i-1} - x^*\|^2 - \|x_i - x^*\|^2) \\ &= \frac{1}{2\lambda} (\|x_0 - x^*\|^2 - \|x_k - x^*\|^2) \\ &\leq \frac{1}{2\lambda} \|x_0 - x^*\|^2. \end{aligned}$$

Dado que $f(x_i)$ es no decreciente,

$$f(x_k) - f(x^*) \leq \frac{1}{2\lambda} \sum_{i=1}^k (f(x_i) - f(x^*)) \leq \frac{1}{2k\lambda} \|x_0 - x^*\|^2.$$

Así se tiene el resultado deseado.

Resultados numéricos

En este capítulo se muestran los resultados numéricos obtenidos al aplicar los métodos antes vistos a una serie de problemas de optimización no diferenciable.

3.1. Formulación por penalización exacta

La estructura de los problemas de optimización con restricciones es esencialmente la siguiente:

$$\begin{aligned} &\text{minimizar } f(x) && (3.1) \\ &\text{sujeto a } c_i(x) = 0 \quad i \in E, \\ &\quad \quad \quad c_j(x) \leq 0 \quad j \in J. \end{aligned}$$

Cuando se intenta resolver un problema de programación no lineal con restricciones, una posibilidad es formular un problema equivalente sin restricciones cuyo mínimo se alcance en la región factible y sea el mismo al del problema original. Esto lleva a la idea de función de penalización la cual es una combinación de f y c_i que permite que f sea minimizada de tal manera que las restricciones no sean violadas.

Un enfoque atractivo para la programación no lineal es tratar de determinar una función de penalización exacta $\phi(x)$, es decir, una función definida en términos de f y c , y que es minimizada localmente por la solución x^* del problema inicial.

Para el problema de programación no lineal (3.1) la función de penalización L_1 asociada es

$$\phi(x) = \nu f(x) + \sum_{i \in E} |c_i(x)| + \sum_{j \in J} \text{máx}\{-c_j(x), 0\} \quad (3.2)$$

donde el parámetro ν ($\nu > 0$) proporciona un medio de ponderación para $f(x)$ [19].

Consideremos el problema

$$\begin{aligned} &\text{minimizar } -x_1 - x_2 \\ &\text{sujeto a } 1 - \|x\|_2^2 \geq 0, \end{aligned}$$

donde $x = (x_1, x_2)$. La solución a este problema es $(1/\sqrt{2}, 1/\sqrt{2})$. Aplicando penalización exacta, se tiene que la función objetivo es

$$\phi(x) = \nu(-x_1 - x_2) + \text{máx}\{\|x\|_2^2 - 1, 0\} \quad (3.3)$$

donde ν toma cualquier valor en $(0, \sqrt{2})$ [19]. Las curvas de nivel para $\nu = 1$ se muestran en la figura (3.1).

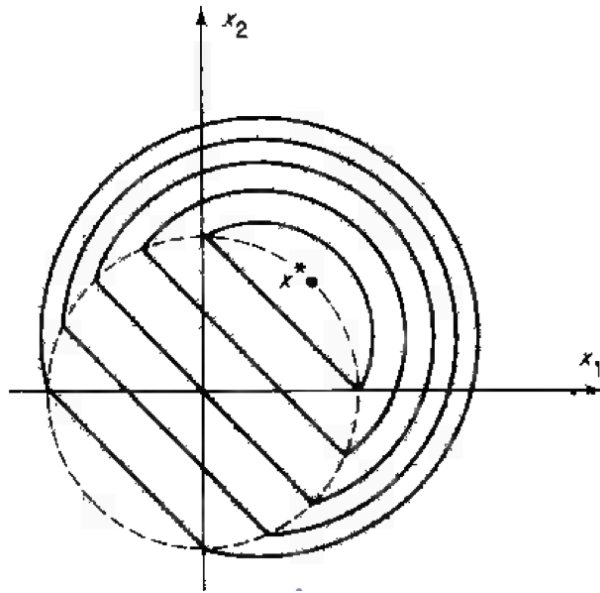


Figura 3.1: Curvas de nivel de la función de penalización exacta L_1

Esta función no es diferenciable en la frontera de la bola unitaria l_2 , pero es una función convexa lo que la hace una candidata ideal para aplicar los métodos anteriormente vistos.

Se observa que la función (3.3), puede ser separada de tal manera que se puede aplicar el método de gradiente próximo. El operador próximo de la función $g(x) = \text{máx}\{\|x\|_2^2 - 1, 0\}$ es por definición

$$\text{prox}_{\lambda g}(v) = \arg \min_x (\text{máx}\{\|x\|_2^2 - 1, 0\} + \frac{1}{2\lambda} \|x - v\|^2).$$

Para calcular el operador próximo de la función $g(x)$, primero notemos que la condición para que x^* sea mínimo del operador próximo en v es

$$0 \in \partial \text{máx}\{\|x^*\|_2^2 - 1, 0\} + \frac{1}{\lambda}(x^* - v). \quad (3.4)$$

Si $\|x\|_2 < 1$ entonces

$$\partial \text{máx}\{\|x\|_2^2 - 1, 0\} = 0, \quad (3.5)$$

y por lo tanto, si $\|v\| \leq 1$ entonces $x^* = v$.

Si $\|x\|_2 > 1$

$$\partial \text{máx}\{\|x\|_2^2 - 1, 0\} = 2x, \quad (3.6)$$

entonces por la ecuación 3.4 se tiene que cumplir

$$\begin{aligned} 2x^* + \frac{1}{\lambda}(x^* - v) &= 0 \\ \iff x^* &= \frac{1}{2\lambda + 1}v, \end{aligned}$$

para $v \in \mathbb{R}^n$ tal que $\|v\| \in [2\lambda + 1, +\infty)$.

Por último, si $\|x\|_2 = 1$

$$\partial \text{máx}\{\|x\|_2^2 - 1, 0\} = \{ax : a \in [0, 2]\}. \quad (3.7)$$

Así, se tiene que resolver la ecuación

$$ax^* + \frac{1}{\lambda}(x^* - v) = 0,$$

despejando x^* se obtiene

$$x^* = \frac{1}{a\lambda + 1}v.$$

Tomando el hecho de que $\|x^*\| = 1$, se verifica que $a\lambda + 1 = \|v\|$. Por lo tanto

$$x^* = \frac{v}{\|v\|}$$

para v tal que $\|v\| \in (1, 2\lambda + 1)$.

Así, usando lo anterior se llega a lo siguiente

$$\text{prox}_{\lambda \text{máx}\{\|x^*\|_2^2 - 1, 0\}}(v) = \begin{cases} v & \text{si } \|v\| \leq 1 \\ \frac{v}{\|v\|} & \text{si } \|v\| \in (1, 2\lambda + 1) \\ \frac{1}{2\lambda + 1}v & \text{si } \|v\| \geq 2\lambda + 1. \end{cases}$$

Por otro lado, para usar el método de subgradiente se necesita calcular en cada iteración el subgradiente de la función $g(x)$ en el punto x^k . Dependiendo del valor $\|x^k\|_2^2$, el subgradiente está dado por (3.5), (3.6) y (3.7). Aplicando los métodos mencionados se obtuvo la siguiente tabla

Método	Punto Inicial	Iteraciones	Tiempo	Error
Subgradiente	(3,3)	700	18.50 seg.	10^{-4}
<i>Bundle</i>	(3,3)	20	7.21 seg.	10^{-4}
Gradiente próximo	(3,3)	2	0.08 seg.	10^{-4}

Todos los métodos se programaron en Matlab utilizando como criterio de paro que la distancia del punto óptimo al punto obtenido x^k sea menor que 10^{-4} . Se tiene claramente que el método de gradiente próximo es superior a los otros dos métodos tanto por el número de iteraciones como por el tiempo de cálculo.

3.2. El problema de LASSO

El LASSO es una herramienta popular para la regresión lineal dispersa [28], especialmente para problemas en los cuales el número de variables m exceden al número de observaciones n . Cuando $m > n$, la función a minimizar en el problema de LASSO no es estrictamente convexa [32] y por lo tanto no es único el punto mínimo.

Se considera un problema con m variables y n observaciones. Sea $x \in \mathbb{R}^n$ el vector de resultados y $A \in \mathbb{R}^{m \times n}$ la matriz de predictores. El problema de optimización de LASSO es

$$\text{minimizar}_{x \in \mathbb{R}^n} (1/2) \|Ax - b\|^2 + \gamma \|x\|_1 \quad (3.8)$$

donde $\gamma > 0$ es un parámetro de ajuste. Es claro que la función objetivo de este problema es no diferenciable.

Este problema es un candidato ideal para aplicar el método de gradiente próximo, donde la parte no diferenciable del problema es la función $\gamma \|\cdot\|_1$. Primero, es necesario calcular el operador próximo de la función $\gamma \|\cdot\|_1$.

Si $f = \|\cdot\|$ es una norma en \mathbb{R}^n , entonces se tiene que $f^* = \mathbb{I}_{\mathcal{B}}$, donde \mathcal{B} es la bola unitaria para la norma dual. Luego, por la descomposición de Moreau, se sigue que

$$\begin{aligned} \text{prox}_f(v) &= v - \lambda \text{prox}_{1/\lambda f^*}(v/\lambda) \\ &= v - \lambda \Pi_{\mathcal{B}}(v/\lambda). \end{aligned} \quad (3.9)$$

Ahora, es bien conocido que la proyección sobre la bola unitaria de la norma l_∞ es

$$(\Pi_{\mathcal{B}}(v))_i = \begin{cases} 1 & \text{si } |v_i| > 1 \\ v_i & \text{si } |v_i| \leq 1 \\ -1 & \text{si } |v_i| < -1. \end{cases}$$

Dado que la norma l_∞ es la norma dual de la norma l_1 , entonces se tiene una manera de evaluar el operador próximo de $\|\cdot\|_1$ utilizando la proyección anterior. Así, el operador próximo de $\|\cdot\|_1$ es

$$(\text{prox}_{\gamma\|\cdot\|_1}(v))_i = \begin{cases} v_i - \gamma\lambda & \text{si } v_i > \gamma\lambda \\ 0 & \text{si } |v_i| < \gamma\lambda \\ v_i + \gamma\lambda & \text{si } \leq -\gamma\lambda. \end{cases}$$

El método de gradiente próximo se puede aplicar ahora para resolver el problema de LASSO.

El problema de LASSO que se resolvió fue el siguiente: se construyó una matriz $A \in \mathbb{R}^{m \times n}$ ($m = 100$, $n = 500$, $\gamma = 100$) donde las componentes de la matriz se eligieron de manera que $a_{ij} \sim \mathcal{N}(0, 1)$ y después se normalizaron las columnas tomando la norma 2. El vector b fue obtenido de la siguiente manera:

$$b = Ax + v$$

donde $x_i \sim \mathcal{N}(0, 1)$ y $v \sim \mathcal{N}(0, 10^{-3})$.

El programa cumplió el criterio de paro en 24 iteraciones, encontrando el valor mínimo de la función $f(x^*) = 3597.3$ en un tiempo de 0.0552 segundos. El criterio de paro que se utilizó fue el propuesto anteriormente en el método de gradiente próximo.

El programa fue comparado y validado utilizando los programas encontrados en [35]. En particular se comparó con el método de subgradiente para el cual se requirieron 153 iteraciones en un tiempo de 0.2893.

3.3. Problema de LASSO modificado

El problema de LASSO fue modificado para estudiar como se determina el operador próximo de la función $\|\cdot\|_\infty$. Se consideró el siguiente problema

$$\text{minimizar }_{x \in \mathbb{R}^n} (1/2)\|Ax - b\|^2 + \gamma\|x\|_\infty \quad (3.10)$$

donde $A \in \mathbb{R}^{m \times n}$ y $\gamma > 0$.

Igual que en el problema de LASSO original, se debe encontrar una expresión del operador próximo de la función $\|\cdot\|_\infty$, la cual no es tan inmediata como en el caso de la función $\|\cdot\|_1$.

Para encontrar el operador próximo de $\|\cdot\|_\infty$ se debe resolver el problema de como proyectar sobre la bola unitaria con norma l_1 . El problema se define como

sigue

$$\text{minimizar }_{\|x\|_1 \leq 1} \frac{1}{2} \|x - v\|^2. \quad (3.11)$$

Al introducir un multiplicador de Lagrange μ para la restricción $\|x\|_1 \leq 1$, se puede escribir la función Lagrangiana como sigue

$$\mathcal{L}(x, \mu) = \frac{1}{2} \|x - v\|^2 + \mu(\|x\|_1 - 1).$$

Sea x^* el punto óptimo del problema (3.10) y μ^* el punto óptimo del problema dual. Tanto x^* como μ^* deben satisfacer que $\|x^*\|_1 \leq 1$, $\mu^* \geq 0$ y además que

$$\mu^*(\|x^*\|_1 - 1) = 0. \quad (3.12)$$

Ahora, supongamos que μ^* es conocido, entonces x^* es la solución del problema

$$\text{mín } \mathcal{L}(x, \mu^*).$$

El problema anterior tiene solución única ya que $\mathcal{L}(\cdot, \cdot)$ es estrictamente convexa en su primer argumento. Dado que las variables en (3.11) se pueden descomponer, se tiene entonces que

$$x_i^* = \arg \min (1/2)(x_i - v_i)^2 + \mu^*(|x_i| - 1), \quad (3.13)$$

lo cual lleva a lo siguiente:

$$x_i^* = \text{signo}(v_i) \text{máx}\{|v_i| - \mu^*, 0\}. \quad (3.14)$$

Se consideran dos casos: $\|v\|_1 \leq 1$ y $\|v\|_1 > 1$. Se mostrará en el siguiente lema que para el primer caso $\mu^* = 0$.

Lema 3.3.1. *Si $\|v\|_1 \leq 1$, entonces el punto óptimo del problema dual μ^* es cero y el punto óptimo x^* está dado por $x^* = v$.*

Demostración. Supongamos que $\mu^* > 0$, entonces de (x_i^*) tenemos que

$$\|x^*\|_1 \leq 1,$$

por lo tanto $\mu^*(\|x^*\|_1 - 1) \neq 0$, lo que contradice (3.12). Así, $\mu^* = 0$ y de (3.14) se sigue que $x^* = v$. \square

Para el caso en que $\|v\|_1 > 1$, el siguiente teorema muestra que μ^* puede ser obtenida calculando la raíz de una función auxiliar.

Teorema 3.3.1. Si $\|v\|_1 > 1$, entonces μ^* es positivo y está dada por la única raíz de la función

$$f(\mu) = \sum_{i=1}^n \max(|v_i| - \mu, 0) - 1. \quad (3.15)$$

La demostración puede encontrarse en [31].

Se puede ahora calcular el operador próximo de la función $\|\cdot\|_\infty$ utilizando (3.14), (3.15) y la descomposición de Moreau. El operador próximo de $\|\cdot\|_\infty$ queda de la siguiente manera

$$(\text{prox}_{\gamma\|\cdot\|_\infty}(v))_i = \begin{cases} |v_i/\gamma\lambda| - \mu^* & \text{si } v_i > \mu^*\gamma\lambda \\ 0 & \text{si } |v_i| < \mu^*\gamma\lambda \\ |v_i/\gamma\lambda| + \mu^* & \text{si } v_i \leq -\mu^*\gamma\lambda. \end{cases}$$

donde μ^* es la raíz de la función

$$f(\mu) = \sum_{i=1}^n \max\{|v_i/\gamma\lambda| - \mu, 0\}.$$

El problema de LASSO que se resolvió fue el mismo que el problema de LASSO original. El programa realizado en Matlab cumplió el criterio de paro en 26 iteraciones, encontrando el valor mínimo de la función $f(x^*) = 615.371$ en un tiempo de 0.3876 segundos.

Este problema se validó utilizando la paquetería CVX [11] en Matlab. CVX realizó 14 iteraciones en un tiempo de 2.3342 en encontrando el valor $f(x^*) = 615.371$.

3.4. El problema del portafolio

El problema clásico de optimización de portafolio, en el que un agente escoge sus estrategias a modo de maximizar su utilidad esperada ha sido estudiado en profundidad por mucho tiempo. Sin embargo solo hace poco ha surgido el interés por plasmar el hecho de que la selección de un modelo particular al momento de hacer la toma de decisiones es en sí riesgosa.

El problema del portafolio se define de la siguiente manera

$$\begin{aligned} &\text{minimizar } \frac{1}{2}x^t\Sigma x + r^t x \\ &\text{sujeto a } \sum x_i = 1, \quad x_i \geq 0, \end{aligned}$$

donde $\Sigma \in \mathbb{R}^{n \times n}$ es la matriz de varianza covarianza y $r \in \mathbb{R}^n$. El problema del portafolio puede ser convertido al siguiente problema

$$\text{minimizar } \frac{1}{2}x^t \Sigma x + r^t x + \mathbb{I}_C,$$

donde \mathbb{I}_C es la función indicadora del conjunto $C = \{x \in \mathbb{R}^n : \sum x_i = 1, x_i \geq 0\}$. Se puede aplicar el método de gradiente próximo si se calcula el operador próximo de la función indicadora \mathbb{I}_C . Usando (2.25), se puede obtener éste operador conociendo como se proyecta sobre el conjunto C . La proyección sobre el conjunto C [15] está dada por

$$(\Pi_C(v))_i = \text{máx}\{v_i - \mu^*, 0\}$$

donde μ^* es la raíz de la función

$$f(\mu) = \sum_{i=1}^n \text{máx}\{v_i - \mu, 0\} - 1.$$

Utilizando (2.25), el operador próximo de \mathbb{I}_C está dada por

$$\text{prox}_{\mathbb{I}_C}(v) = \text{máx}\{v_i - \mu^*, 0\}.$$

Ejemplo. Sea P^1, \dots, P^5 los precios de 5 empresas de la rama de construcción: Se tomaron los datos del 2009-2014 de los precios diarios de las acciones. Los rendimientos diarios de la empresa i están dados por $r_j^i = \ln(\frac{P_j^{i+1}}{P_j^i})$. Σ es la matriz de varianza covarianza muestral

$$\Sigma = \begin{pmatrix} 0.0000778 & 0.00000796 & 0.000000645 & 0.0000541 & 0.00000346 \\ 0.00000796 & 0.000512 & -0.0000432 & 0.0000551 & 0.00000273 \\ 0.000000645 & -0.0000432 & 0.000315 & 0.000305 & 0.0000149 \\ 0.0000541 & 0.0000551 & 0.000305 & 0.0043 & 0.000116 \\ 0.00000346 & 0.00000273 & 0.0000149 & 0.000116 & 0.000208 \end{pmatrix}$$

y el vector r es la media muestral de los rendimientos de cada empresa

$$r = (-0.0004142, 0.0004127, 0.0018, -0.00411, 0.0008422).$$

Aplicando el método de gradiente próximo se llegó al punto

$$x^* = (0, 0, 0.9999, 0, 0)$$

en 10 iteraciones en un tiempo de 0.2140 segundos. Lo que significa que todo el dinero se debe invertir en la acción de la empresa 3.

Si se trata de resolver este problema por medio del método de restricciones activas, las condiciones Kuhn y Tucker (KT) correspondientes son: existen $\mu \in \mathbb{R}$ y $v \in \mathbb{R}^5$ tales que

$$\begin{aligned}\Sigma w + \mu \bar{1} - \sum_{i=1}^5 v_i \bar{e}_i &= 0, \\ \bar{1}^t w &= 1, \\ v_i w_i &= 0, \quad i = 1, \dots, 5\end{aligned}$$

con e_i , el i -ésimo vector de la base canónica de \mathbb{R}^5 . El método a seguir en este caso es clasificar las soluciones dependiendo de que las restricciones de desigualdad sean activas o pasivas. Para cada caso se tiene que comprobar que las condiciones KT se cumplen, es decir, se tiene que comprobar las condiciones KT para 25 casos posibles.

3.5. El problema de Minimax

El problema de minimización del Minimax (problema Minimax) es el siguiente

$$\text{minimizar } f(x),$$

donde

$$f(x) = \max_{i=1, \dots, n} f_i(x).$$

En este problema se supondrá que $f_1(x), \dots, f_m(x) : \mathbb{R}^n \rightarrow \mathbb{R}$ son funciones convexas y de clase C^2 .

Los problemas Minimax surgen en la ingeniería del diseño, diseño de circuitos y control óptimo [33]. La función objetivo tiene derivadas parciales discontinuas en los puntos donde dos o más funciones $f_i(x)$ son iguales a $f(x)$. En este problema, no es posible encontrar una expresión general del operador próximo ya que este cambia dependiendo de las funciones f_i . Sin embargo, el operador próximo se puede calcular de manera numérica replanteando el operador de la siguiente manera

$$\min_{r,x} r + \frac{1}{2\lambda} \|x - v\|^2 \tag{3.16}$$

$$\text{sujeto a } f_i(x) \leq r, \quad i = 1, \dots, m.$$

Dado que el problema (3.16) es un problema convexo, se puede resolver por medio del método de restricciones activas, para esto se aplicó la paquetería CVX

[11] en Matlab. Así, CVX se aplica a la función en cada iteración del algoritmo de punto próximo, es decir, en cada iteración del método de punto próximo se resuelve un problema cuya función objetivo es convexa y las funciones en las restricciones son de clase C^2 .

Nuestros resultados numéricos mejoran a los obtenidos en [33]. El método usado en [33], es un método de suavización para problemas Minimax. Está basado en la función de penalización exponencial de Kort y Bertsekas para optimización con restricciones. Utilizando el método de punto próximo con los puntos iniciales propuestos en [33] se obtiene lo siguiente:

Ejemplo 1. Sea el problema

$$\min_x \max\{f_1(x), f_2(x), f_3(x)\}$$

donde

$$\begin{aligned} f_1(x) &= x_1^2 + x_2^4 \\ f_2(x) &= (2 - x_1)^2 + (2 - x_2)^2, \\ f_3(x) &= 2 \exp^{-x_1+x_2}. \end{aligned}$$

El problema tiene como solución el punto (1.1390, 0.8986). El punto inicial es $x = (1, -0, 1)$. Después de 8 iteraciones y un tiempo de 7.2995 segundos se obtuvo la solución $x^* = (1.13907, 0.89953)$.

El método en [33] obtuvo el mismo punto óptimo después de 22 iteraciones.

Ejemplo 2. Sea el problema

$$\min_x \max\{f_1(x), f_2(x), f_3(x)\}$$

donde

$$\begin{aligned} f_1(x) &= x_1^4 + x_2^2 \\ f_2(x) &= (2 - x_1)^2 + (2 - x_2)^2, \\ f_3(x) &= 2 \exp^{-x_1+x_2}. \end{aligned}$$

El problema tiene como solución el punto (1,1). El punto inicial elegido es $x = (1, -0, 1)$. Después de 4 iteraciones y un tiempo de 4.1878 segundos se obtuvo la solución $x^* = (1, 0.99999)$.

El método en [33] obtuvo el punto óptimo $x^* = (1.0000, 1.0000)$ después de 25 iteraciones.

Ahora, se consideran los siguiente problemas: el problema de programación no lineal

$$\text{minimizar } F(x) \tag{3.17}$$

sujeto a $g_i(x) \geq 0 \quad i = 1, \dots, m,$

y el problema Minimax

$$\min \max_{1 \leq i \leq m} f_i(x) \quad (3.18)$$

donde

$$F(x) = f_1(x), \quad f_i(x) = F(x) - \alpha_i g_i(x), \quad 2 \leq i \leq m, \quad \alpha_i > 0.$$

Blander y Charalambous [2] probaron que para α_i suficientemente grande, el punto óptimo del problema de Minimax (3.18) coincide con el punto óptimo del problema (3.17).

Ejemplo 3. (El problema de Rosen-Suzuki). Sean

$$F(x) = x_1^2 + x_2^2 + 2x_3^2 + x_4^2 - 5x_1 - 5x_2 - 21x_3 + 7x_4$$

$$g_2(x) = -x_1^2 - x_2^2 - x_3^3 - x_4^2 - x_1 + x_2 - x_3 + x_4 + 8$$

$$g_3(x) = -x_1^2 - 2x_2^2 - x_3^2 - 2x_4^2 + x_1 + x_4 + 10$$

$$g_4(x) = -x_1^2 - x_2^2 - x_3^2 - 2x_1 + x_2 + x_4 + 5$$

y sea $\alpha_2 = \alpha_3 = \alpha_4 = 10$. La solución a este problema es el punto $x^* = (0, 1, 2, -1)$. Aplicando el mismo programa que los problemas anteriores, utilizando el punto inicial $x_0 = (1, 1, 1, 1)$, después de 2 iteraciones y 4.2527 segundos se llega a la solución

$$x^* = (0.0000, 1.0000, 2.0000, -0.9999).$$

El método en [33] obtuvo el punto

$$x^* = (0.0000, 1.0000, 2.0000, -1.0000)$$

después de 25 iteraciones.

3.6. ADMM

Se ha trabajado con funciones objetivo de la forma $f(x) + g(x)$ donde f y g son convexas pero solo f diferenciable. En el caso en que las dos funciones sean convexas pero no necesariamente diferenciables se tiene un método que permite resolver este tipo de problemas basado en el operador próximo. Se considera el problema

$$\min f(x) + g(x)$$

donde $f, g : \mathbb{R}^n \rightarrow \bar{\mathbb{R}}$ funciones propias, convexas y cerradas, no necesariamente diferenciables. Entonces el método ADMM [13](Alternating Directions Method of Multipliers) es

$$\begin{aligned}x_{k+1} &= \text{prox}_{\lambda f}(z_k - u_k) \\z_{k+1} &= \text{prox}_{\lambda g}(x_{k+1} + u_k) \\u_{k+1} &= u + x_{k+1} - z_{k+1}.\end{aligned}$$

Este método converge bajo las mismas condiciones generales que se han estado trabajando a lo largo de este trabajo [7]. Mientras x_k y z_k convergen hacia el mismo punto, x_k y z_k tienen distintas propiedades. Por ejemplo, $x_k \in \text{dom } f$ mientras que $z_k \in \text{dom } g$, si g es la función que abarca las restricciones entonces z_k satisface las restricciones para todo k , mientras que x_k satisface las restricciones solo en el límite.

La ventaja de trabajar con el ADMM es que los términos de la función objetivo son manejados por separado. El ADMM es útil cuando los operadores próximos de f y g son más fáciles de obtener que el operador próximo de la función $f + g$.

Ejemplo 1. Considerar el siguiente problema

$$\begin{aligned}\text{mín } & \|x\|_\infty \\ \text{sujeto a } & Ax \leq b\end{aligned}$$

donde

$$A = \begin{pmatrix} -1 & 1 \\ 1 & 1 \\ -1 & -1 \\ 1 & -1 \end{pmatrix}$$

y $b = (-1, 1, 1, 3)$.

El problema anterior puede plantearse de la siguiente manera

$$\text{mín } \|x\|_\infty + \mathbb{I}_C$$

donde $C = \{x : Ax \leq b\}$ y cuyo mínimo es $x^* = (0.5, -0.5)$. Las dos funciones son convexas pero no diferenciables. El operador próximo de la norma infinito fue deducido anteriormente y el de la función indicadora es la proyección sobre el conjunto C . De esta manera se aplica el método ADMM utilizando como puntos iniciales $u = (1, 1)$ y $z = (0, 0)$. Después de 9 iteraciones y 7.45 segundos se llega a la solución

$$x^* = (0.5001, -0.4999).$$

Ejemplo 2. Se considera el mismo problema anterior pero ahora sea

$$A = \begin{pmatrix} -1 & 1 & -1 & 1 \\ 1 & 1 & 1 & 1 \\ -1 & -1 & -1 & -1 \\ 1 & -1 & 1 & -1 \\ 1 & 1 & -1 & -1 \end{pmatrix}$$

y $b = (-1, 1, 1, 3, 4)$. El mínimo de este problema es el punto $x^* = (0.25, -0.25, 0.25, -0.25)$.

Se aplicó el método ADMM de la misma manera que el ejemplo 1 utilizando como puntos iniciales $u = (1, 1, 1, 1)$ y $z = (0, 0, 0, 0)$. Después de 22 iteraciones y 10.30 segundos se llega a la solución

$$x^* = (0.2499, -0.2500, 0.2500, -0.2499).$$

Los dos ejemplos fueron validados con la paquetería CVX [11] de Matlab.

Conclusiones

El método de subgradiente es fácil de implementar, sin embargo, determinar el subdiferencial de cualquier función no es simple, por lo que el método de subgradiente, aunque es fácil de implementar, presenta complicaciones al momento de encontrar el subgradiente de la función en cada iteración. Además, se encontró un ejemplo de una función que cumple que uno de sus subgradientes en dirección negativa no es una dirección de descenso, lo cual concreta lo argumentado en [12], de que el método de subgradiente no es un método de descenso. La versión del método de subgradiente que se utilizó en este trabajo no cuenta con un criterio de paro que no dependa del conocimiento del punto óptimo, lo cual representa otra dificultad.

El método *bundle* permite tener un método con muchas propiedades útiles, entre las cuáles está el tener un criterio de paro que no depende del conocimiento del punto óptimo. Además, aunque el método *bundle* requiere del cálculo de subgradientes en cada iteración, otorga una manera de elegir un subgradiente adecuado en cada iteración. El problema principal es que con cada iteración el subproblema a resolver crece, lo cual tiene repercusiones computacionales.

La regularización de Moreau-Yosida tiene la ventaja que puede obtenerse de manera explícita, lo cual es muy útil al momento de programar. Aún en el caso de que esto no sea posible, se puede manipular de manera que se tenga un problema convexo con restricciones convexas, que en general es un problema sencillo de resolver.

Íntimamente ligado a la regularización de Moreau-Yosida, el operador próximo hereda sus propiedades. Como se revisó en [7], este operador no necesariamente es una contracción; sin embargo, se probó que es firmemente no expansivo, la cual es una condición suficiente para tener convergencia en métodos de tipo punto fijo. Otra ventaja de este operador es que puede calcularse de manera explícita para varias funciones no diferenciables lo que simplifica su implementación computacional. El operador próximo, resulta ser una manera efectiva de obtener métodos que resuelven una amplia variedad de problemas no diferenciables. Sin embargo, no es la panacea

porque no se puede aplicar a todos los problemas y, en ocasiones, puede resultar demasiado complicado su implementación.

Los métodos de punto próximo y gradiente próximo son sencillos de implementar computacionalmente ya que solo necesitan el cálculo del operador próximo de la función a minimizar. La rapidez de convergencia del método de punto próximo y el método de gradiente próximo es $o(1/k)$, por lo que no son los métodos más rápidos en cuanto a convergencia. Sin embargo, es fácil su implementación computacional a comparación de otros métodos más generales que tienen una mayor rapidez de convergencia [12]. Los métodos que utilizan el operador próximo fueron superiores al método de subgradiente y al método *bundle*, pero también fueron superiores a otros métodos encontrados en la literatura como en el caso de [32][33].

Cabe destacar que la convergencia de todos los métodos estudiados en este trabajo no dependen del punto inicial elegido, lo cual es una ventaja muy grande al momento de aplicarlos a problemas no diferenciables (incluso problemas diferenciables).

Como trabajo futuro se pretende estudiar métodos que tengan un mayor rapidez de convergencia ($o(1/k^2)$). Más aún, que se puedan aplicar a una amplia variedad de problemas que surgen típicamente en las aplicaciones.

Bibliografía

- [1] Asaadi J., *A computational comparison of some non-linear programs*, Mathematical Programming, Vol. 4, Issue 1, pp 144-154, 1973.
- [2] Bandler J., Charalambous C., *Nonlinear programming using minimax techniques*, Journal of Optimization Theory Applications, Vol. 13, pp. 607-619, 1974.
- [3] Baptiste J., Lemarechal C., Urruty H., *Convex analysis and minimization algorithms I*, Springer, Second corrected printing, 1994.
- [4] Baptiste J., Lemarechal C., Urruty H., *Convex analysis and minimization algorithms II*, Springer, 1993.
- [5] Barbu V., Precupanu T., *Convexity and optimization in Banach spaces*, Springer monographs in mathematics, 4th edition, 2012.
- [6] Bauschke H., Burachik R., Combettes P., Elser V., Luke D., Wolkowicz H., *Fixed-Point Algorithms for Inverse Problems in Science and Engineering*, Springer Optimization and its Applications, Vol. 49, 2011.
- [7] Bauschke H., Combettes P., *Convex Analysis and Monotone Operator Theory in Hilbert Spaces*, Canadian Mathematical Society, Springer, 2011.
- [8] Bazaraa M., Sherali H., Shetty C., *Nonlinear Programming*, Wiley, 3rd Edition, 2006.
- [9] Belloni A., *Lecture notes for IAP 2005 course: Introduction to bundle methods*, Operation Research Center, M.I.T., Version of February 11, 2005.
- [10] Bertsekas D., *A new algorithm for solution of nonlinear resistive networks involving diodes*, IEEE Transactions on Circuit Theory, Vol. 23, Issue 10, pp. 599-608, 1976.

- [11] Boyd S., Grant M., *CVX: Matlab software for disciplined convex programming (web page and software)*, <http://cvxr.com/cvx/> , August 2008.
- [12] Boyd S., Mutapcic A., *Subgradient methods*, Notes for EE364b, Stanford University, Winter 2006-07.
- [13] Boyd S., Parikh N., *Proximal Algorithms*, Foundations and Trends in Optimization, Vol. 1, No. 3, pp 123-231, 2013.
- [14] Boyd S., Vandenberghe L., *Convex optimization*, Cambridge University press, 2004.
- [15] Carreira-Perpiñán M., Wang W., *Projection onto probability simplex: An efficient algorithm with a simple proof and an application*, <http://arxiv.org/abs/1309.1541v1>.
- [16] Clarke F., *Optimization and Nonsmooth Analysis*, CLASSICS in Applied Mathematics, SIAM, Wiley, 1990.
- [17] Combettes P., Reyes N., *Moreau's Decomposition in Banach Spaces*, <http://arxiv.org/abs/1103.3178v1>.
- [18] Combettes P., Wajs V., *Signal Recovery by Proximal Forward-Backward Splitting*, Multiscale Modeling and Simulation, 2005, Vol. 4, No. 4, pp. 1168-1200.
- [19] Fletcher R., *Practical Methods of Optimization*, Wiley, Second Edition, 2013.
- [20] Gigola C., Gómez S., A Regularization method for solving the finite convex Min-Max Problem, SIAM Journal of Numerical Analysis, Vol. 27, No. 6, USA: 1990.
- [21] Jahn J., *Duality in Vector Optimization*, Mathematical Programming, Vol. 25, Issue 3, pp 343-353, 1983.
- [22] Kiwiel K., *Methods on descent for nondifferentiable optimization*, Lecture Notes in Mathematics, Springer, 1985.
- [23] Luenberger D., *Optimization by Vector Space Methods*, Wiley Professional Paperback Series, 1968.
- [24] Niculescu C., Persson L., *Convex Functions and Their Applications*, Canadian Mathematical Society, Springer and Science Business Media, 2006.

- [25] Pshenichny B., *Necessary conditions for an extremum*, M. Dekker, Series: Pure and Applied Mathematics, Vol. 4, 1971.
- [26] Recht B., *Projected Gradient and Proximal Point Methods*, Department of Electrical Engineering and Computer Science, University of California, Berkeley, 2014.
- [27] Rockafellar R., *Convex Analysis*, Princeton University Press, 1970.
- [28] Rosenberg B., *Linear regression with randomly dispersed parameters*, Oxford Journals, Biometrika, Vol. 60, Issue 1, pp. 65-72, 1973.
- [29] Säljö R., *Implementation of bundle algorithm for convex optimization*, Göteborg University, Master's thesis, May 11, 2004.
- [30] Sion M., *On general minimax theorems*, Pacific Journal of Mathematics, Vol. 8, No. 1, pp. 171-176, March 1958.
- [31] Songsiri J., *Projection onto an l_1 -norm Ball with Application to Identification of Sparse Autoregressive Models*, Asian Symposium on Automatic Control, ASAC, 2011.
- [32] Tibshirani R., *The LASSO Problem and Uniqueness*, Electronic Journal of Statistics, Vol. 7, pp. 1456-1490, 2013.
- [33] Xu S., *Smoothing Methods for Minimax Problems*, Computational Optimization and Applications, Vol. 20, Issue 3, pp. 267-279, 2001.
- [34] Yu Y., *On Decomposing the Proximal Map*, Department of Computing Science, University of Alberta, Edmonton AB T6G 2E8, Canada, 2013.
- [35] <http://www.cs.ubc.ca/~schmidtm/Software/lasso.html>.