

UN ESTUDIO SOBRE EL MÉTODO DE CAUCHY CON BÚSQUEDA LINEAL  
DE WOLFE SOBRE VARIEDADES DE RIEMANN

POR

MALON RAFAEL MENDOZA

UNIVERSIDAD CENTROCCIDENTAL “LISANDRO ALVARADO”

DECANATO DE CIENCIAS Y TECNOLOGIA

BARQUISIMETO

UN ESTUDIO SOBRE EL MÉTODO DE CAUCHY CON BÚSQUEDA LINEAL  
DE WOLFE SOBRE VARIEDADES DE RIEMANN

POR

MALON RAFAEL MENDOZA

TRABAJO DE ASCENSO PRESENTADO PARA OPTAR A LA  
CATEGORIA TITULAR EN EL ESCALAFON DEL PERSONAL  
DOCENTE Y DE INVESTIGACION.

UNIVERSIDAD CENTROCCIDENTAL "LISANDRO ALVARADO"

DECANATO DE CIENCIAS Y TECNOLOGIA

Barquisimeto, Noviembre 2013

## AGRADECIMIENTO

---

Mi obediencia y gratitud a la misma esencia inmanente en todo lo manifestado, el Ser Supremo. En consecuencia, a los que promovieron y ayudaron con este trabajo.... Om Tat Sat.



## Resumen

---

El objetivo planteado en este trabajo es, primero, estudiar y analizar el algoritmo de optimización de la búsqueda lineal con condición de Wolfe y mayor descenso en  $\mathbb{R}^n$ , desarrollado por J. Nocedal y S. Wright [2], y segundo, dar detalles de los pasos de los requerimientos que W. Ring y B. Wirth [2] realizan, para llevar esta búsqueda lineal con condición de Wolfe a variedades Riemannianas.



## Índice general

---

<b>AGRADECIMIENTO</b>	<b>III</b>
<b>Resumen</b>	<b>v</b>
<b>Introducción</b>	<b>1</b>
<b>1. Preliminares</b>	<b>5</b>
1.1. Nociones Básicas: En $\mathbb{R}^n$ . . . . .	5
1.2. Nociones Básicas en Optimización . . . . .	7
1.2.1. Tipos de problemas en Optimización . . . . .	7
1.2.2. Restricciones y Regiones Factibles . . . . .	8
1.2.3. Tipos de Solución Óptima . . . . .	9
1.2.4. Existencia de Solución para Problemas de Optimización . . . . .	10
1.2.5. Condiciones de Optimalidad para Problemas Sin Restricciones . . . . .	11
1.2.6. Condiciones de Optimalidad: Necesarias y suficientes . . . . .	12
1.2.7. Convexidad y Minimización . . . . .	16
<b>2. Método de Búsqueda Lineal con Condición de Wolfe en <math>\mathbb{R}^n</math></b>	<b>21</b>

2.1. Algoritmo de Optimización General . . . . .	21
2.2. La Condición de Wolfe . . . . .	26
<b>3. Variedades</b>	<b>33</b>
3.1. Aplicaciones diferenciables entre variedades . . . . .	34
3.2. Espacio Tangente . . . . .	35
3.3. Métricas Riemannianas . . . . .	37
3.4. Subvariedades Encajadas (Embedding) . . . . .	38
3.4.1. Teoría General . . . . .	39
3.5. Campos de vectores, conexiones afines y derivada covariante . . . . .	41
3.6. Curvatura de una variedad Riemanniana . . . . .	47
3.7. Preliminares de optimización sobre Variedades . . . . .	52
3.8. Funciones Convexas en Subvariedades de Riemann . . . . .	59
3.8.1. Gradiente y Hessiano en Subvariedades . . . . .	59
<b>4. Método De Descenso Clásico Sobre Variedades.</b>	<b>63</b>
4.1. Retracciones . . . . .	64
4.2. Búsqueda Lineal: minimización sobre Variedades . . . . .	72
4.3. Convergencia para el método de búsqueda lineal en variedades . . . . .	77
4.4. Aplicación . . . . .	84
<b>Bibliografía</b>	<b>95</b>



## Introducción

---

En palabras de Absil, Mahony y Sepulchre [16], la mayoría de las técnicas numéricas disponibles para la optimización y las ecuaciones no lineales asumen como base un espacio euclidiano. Sin embargo, muchos problemas de cálculo se plantean en espacios no euclidianos. Varios autores [3, 4, 5, 6, 7] han propuestos algoritmos abstractos que explotan la geometría subyacente (por ejemplo, la simétrica, la homogénea, la Riemanniana) de las variedades sobre las que se proyectan los problemas, pero la conversión de estos algoritmos geométricos abstractos en procedimientos numéricos en situaciones prácticas, es a menudo una tarea no trivial que depende fundamentalmente de una adecuada representación de la variedad.

El objetivo planteado en este trabajo es, primero, estudiar y analizar el algoritmo de optimización de la búsqueda lineal con condición de Wolfe y mayor descenso en  $\mathbb{R}^n$ , desarrollado por J. Nocedal y S. Wright [2], y segundo, dar detalles de los pasos de los requerimientos que W. Ring y B. Wirth [2] realizan, para llevar esta búsqueda lineal con condición de Wolfe a variedades Riemannianas.

Hay un número de problemas que puede ser expresado como una minimización de una función  $f : M \rightarrow \mathbb{R}$  sobre una variedad riemanniana  $M$ . Si la variedad  $M$  puede ser inmersa en un espacio de dimensiones superiores o si se define a través de una igualdad con limitaciones, entonces existe la posibilidad de emplear las herramientas muy avanzadas de optimización con restricciones (ver [2] para una introducción). A menudo, sin embargo, esta inmersión no está a la mano, y uno tiene que recurrir a métodos de optimización diseñados para variedades de Riemann. Aunque, si se obtuviera una inmersión, uno podría esperar, que un método de optimización Riemanniano se realizaría más eficientemente ya que aprovecharía la estructura geométrica de la variedad subyacente. Para ello, varios métodos se han ideado, desde el más simple, como el de descenso del gradiente en variedades [9], y hasta más sofisticado, como lo es el de *región de confianza* [8]. Previamente, se presentaron los primeros intentos de adaptar los métodos de optimización estándar para problemas en variedades por Gabay [3], quien presentó los algoritmos de los métodos de descenso más rápido, de Newton, y cuasi-Newton; indicando sus propiedades globales y locales de convergencia (sin embargo, sin dar detalles de los análisis para el caso cuasi-Newton). Udriste [7] también indica un algoritmo para el método de máximo descenso y para el método de Newton en variedades de Riemann; prueba convergencia (lineal) del primero bajo la suposición de la búsqueda lineal exacta. Muy recientemente, Yang tomó estos métodos y analizó la convergencia del método de máximo descenso y de Newton para el control del paso de Armijo [9].

En comparación con los métodos de búsqueda lineal estándar en los espacios vectoriales, en la variedad, como se ha dicho en el párrafo anterior, todo se aproxima a sustituir el paso lineal de la dirección de búsqueda por un paso a lo largo de una curva geodésica. Sin embargo, las geodésicas pueden ser difícil de obtener. En los enfoques alternativos, las geodésicas a menudo son reemplazadas por caminos más generales, sobre la base de las llamadas retracciones (una retracción  $R_x$  es una aplicación del espacio tangente,  $T_xM$ , a la variedad  $M$  en  $x \in M$ ; es decir,  $R_x : T_xM \rightarrow M$ , ver definición (4.1.1)). Usando esta retracción, por ejemplo, la función objetivo  $f : M \rightarrow \mathbb{R}$  es elevada a una función objetivo,  $\hat{f}_x = f \circ R_x$  sobre  $T_xM$ . Ya que  $T_xM$  es un espacio euclidiano, es posible definir un modelo cuadrático de  $\hat{f}_x$  y adaptar los métodos clásicos en  $\mathbb{R}^n$  para calcular (en general, aproximadamente) un minimizador del modelo dentro de una región factible alrededor de  $0_x \in T_xM$ . Este minimizador es entonces regresado de  $T_xM$  a  $M$  mediante la retracción  $R_x$ . Por lo tanto, tenemos un punto que será un candidato para la nueva iteración, el cual se acepta o se rechaza en función de la calidad de coincidencias entre el modelo cuadrático y la función  $f$  misma.

Nótese, como se verá en el capítulo 4, que teóricamente es posible escoger como retracción global en una variedad de Riemann, a la aplicación exponencial  $exp_x : T_xM \rightarrow M$ ,  $v \rightarrow exp_x v$ , definida como  $exp_x v = \gamma(1)$ , donde  $\gamma$  resuelve una ecuación diferencial ordinaria de segundo orden con condiciones iniciales  $\gamma(0) = x$  y  $\gamma'(0) = v$  (ver 3.7). Esto corresponde a una estrategia para utilizar métodos numéricos de optimización en varieda-

des de Riemann, cuando se calcula la exponencial sobre un vector director con el fin de obtener un nuevo iterado o dirección de movimiento.

Así que, el esquema de trabajo que seguiremos es el siguiente, en el capítulo 1, tendremos las nociones básicas requeridas; en el capítulo 2, veremos el métodos de optimización de búsqueda lineal con condición de Wolfe en  $\mathbb{R}^n$ ; en el capítulo 3, las variedades y elementos de optimización sobre las mismas y, en el capítulo 4, el métodos de optimización de búsqueda lineal con condición de Wolfe sobre variedades, siguiendo el procedimiento estándar de espacios vectoriales de dimensión finita [1]; además, se da el análisis del método de máximo descenso básico presentando convergencia. Y finalizaremos, con una aplicación del método del Armijo al cociente de Rayleigh sobre la esfera unitaria [16].

---

# Capítulo 1

## Preliminares

---

### 1.1. Nociones Básicas: En $\mathbb{R}^n$

**Funciones a valores reales:** Sea  $f : \mathbf{X} \rightarrow \mathbb{R}$ , donde  $\mathbf{X} \subset \mathbb{R}^n$  es abierto.

**Definición 1.1.1**  $f$  es diferenciable en  $\hat{x} \in \mathbf{X}$  si existe un vector  $\nabla f(\hat{x})$  (el gradiente de  $f$  por  $\hat{x}$ ) y una función  $\theta_{\hat{x}}(y) : \mathbf{X} \rightarrow \mathbb{R}$  con

$$\lim_{y \rightarrow 0} \theta_{\hat{x}}(y) = 0,$$

tal que para cada  $x \in \mathbf{X}$

$$f(x) = f(\hat{x}) + \nabla f(\hat{x})^T(x - \hat{x}) + \|x - \hat{x}\|\theta_{\hat{x}}(x - \hat{x}).$$

$f$  es diferenciable en  $\mathbf{X}$  si  $f$  es diferenciable para todo  $\hat{x} \in \mathbf{X}$ . El vector gradiente es el vector de las derivadas parciales

$$\nabla f(\hat{x}) = \left( \frac{\partial f(\hat{x})}{\partial x_1}, \dots, \frac{\partial f(\hat{x})}{\partial x_n} \right)^T.$$

La derivada direccional de  $f$  por  $\hat{x}$  en la dirección  $d$  es

$$\lim_{\lambda \rightarrow 0} \frac{f(\hat{x} + \lambda d) - f(\hat{x})}{\lambda} = \nabla f(\hat{x})^T d.$$

**Definición 1.1.2** La función  $f$  es dos veces diferenciable en  $\hat{x} \in \mathbf{X}$  si existe un vector  $\nabla f(\hat{x})$  y una matriz simétrica  $n \times n$ ,  $H(\hat{x})$  (el Hessiano de  $f$  por  $\hat{x}$ ) tal que para cada  $x \in \mathbf{X}$

$$f(x) = f(\hat{x}) + \nabla f(\hat{x})^T(x - \hat{x}) + \frac{1}{2}(x - \hat{x})^T H(\hat{x})(x - \hat{x}) + \|x - \hat{x}\|^2 \alpha(\hat{x}, x - \hat{x}).$$

y con

$$\lim_{y \rightarrow 0} \alpha(\hat{x}, y) = 0.$$

$f$  es dos veces diferenciable en  $\mathbf{X}$  si  $f$  es dos veces diferenciable para todo  $\hat{x} \in \mathbf{X}$ . El Hessiano, que denotamos por  $H(x)$  es la matriz de las segundas derivadas

$$[H(\hat{x})]_{ij} = \frac{\partial^2 f(\hat{x})}{\partial x_i \partial x_j},$$

y para funciones con segundas derivadas continuas, este será siempre simétrica

$$\frac{\partial^2 f(\hat{x})}{\partial x_i \partial x_j} = \frac{\partial^2 f(\hat{x})}{\partial x_j \partial x_i}$$

**Funciones a valores vectoriales:** Sea  $f : \mathbf{X} \rightarrow \mathbb{R}^m$ , donde  $\mathbf{X} \subset \mathbb{R}^n$  es abierto.

$$f(x) = f(x_1, \dots, x_n) = \begin{pmatrix} f_1(x_1, \dots, x_n) \\ f_2(x_1, \dots, x_n) \\ \vdots \\ f_m(x_1, \dots, x_n) \end{pmatrix}$$

donde cada una de las funciones  $f_i$  es una función a valores reales.

**Definición 1.1.3** El Jacobiano de la función  $f$  en  $\hat{x}$  es la matriz cuya  $j$ -ésima fila es el gradiente de  $f_j$  por  $\hat{x}$  traspuesta. Más específicamente, el Jacobiano de  $f$  en  $\hat{x}$  es definido como  $\nabla f(\hat{x})^T$ , donde  $\nabla f(\hat{x})$  es la matriz con entrada

$$[\nabla f(\hat{x})]_{ij} = \frac{\partial f_j(\hat{x})}{\partial x_i},$$

Nótese que la  $j$ -ésima columna de  $\nabla f(\hat{x})$  es el gradiente de  $f_j$  en  $\hat{x}$

## 1.2. Nociones Básicas en Optimización

### 1.2.1 Tipos de problemas en Optimización

Algunas terminologías(ver [10]): La función  $f(x)$  es la función objetivo. Las restricciones  $h_i(x) = 0$  se refieren a las *restricciones de igualdad*; mientras que  $g_i(x) \leq 0$  se refiere a las *restricciones de desigualdad*. Nótese que no se utiliza restricciones de la forma  $g_i(x) < 0$ .

#### Problema de Optimización sin restricciones

$$(P) \quad \min_x \quad f(x)$$

$$\text{sujeeto a } x \in \mathbf{X},$$

donde  $x = (x_1, \dots, x_n)^T \in \mathbb{R}^n$ ,  $f : \mathbb{R}^n \rightarrow \mathbb{R}$ , y  $\mathbf{X}$  es un abierto (usualmente  $\mathbf{X} = \mathbb{R}^n$ ).

### Problema de Optimización con restricciones

$$\begin{aligned}
 (P) \quad & \min_x \quad f(x) \\
 & \text{sujeto a } g_i(x) \leq 0 \quad i = 1, \dots, m \\
 & \quad \quad h_i(x) = 0 \quad i = 1, \dots, l \\
 & \quad \quad x \in \mathbf{X}
 \end{aligned}$$

Colocando  $g(x) = (g_1(x), \dots, g_m(x))^T : \mathbb{R}^n \rightarrow \mathbb{R}^m$ ,  $h(x) = (h_1(x), \dots, h_l(x))^T : \mathbb{R}^n \rightarrow \mathbb{R}^l$ . Entonces (P) se puede reescribir como

$$\begin{aligned}
 (P) \quad & \min_x \quad f(x) \\
 & \text{sujeto a } g(x) \leq 0 \\
 & \quad \quad h(x) = 0 \\
 & \quad \quad x \in \mathbf{X}
 \end{aligned}$$

#### 1.2.2 Restricciones y Regiones Factibles

Un punto  $x$  es *factible* para (P) si satisface todas las restricciones; en el caso de problemas sin restricciones,  $x \in \mathbf{X}$ . El conjunto de todos los puntos factibles forman *la región de puntos factibles*, o *conjunto factible* que denotaremos por  $\mathbf{S}$ .

Se dice que la restricción de desigualdad  $g_i(x) \leq 0$  es *activa* en el punto factible  $\hat{x}$ , si  $g_i(\hat{x}) = 0$ , y *no activa* si  $g_i(\hat{x}) < 0$  (todas las restricciones de desigualdad se consideraran activa en cualquier punto factible).



### 1.2.3 Tipos de Solución Óptima

Considere el problema de optimización general

$$(P) \quad \min_{x \in \mathbf{S}} \quad \text{o} \quad \max_{x \in \mathbf{S}} \quad f(x)$$

Recordemos que una  $\varepsilon$ -vecindad de  $\hat{x}$  o la bola centrada en  $\hat{x}$  de radio  $\varepsilon$  es:

$$B(\hat{x}, \varepsilon) = N_\varepsilon(\hat{x}) := \{x : \|x - \hat{x}\| \leq \varepsilon\}.$$

Tendremos la siguiente definición local/global, minimizador y maximizador estricto/no estricto.

**Definición 1.2.1** *En el problema de optimización (P),*

- $x \in \mathbf{S}$  es un minimizador global de (P) si  $f(x) \leq f(y)$  para todo  $y \in \mathbf{S}$ .
- $x \in \mathbf{S}$  es un minimizador global estricto de (P) si  $f(x) < f(y)$  para todo  $y \in \mathbf{S}$ , con  $y \neq x$ .
- $x \in \mathbf{S}$  es un minimizador local de (P) si existe  $\varepsilon > 0$  tal que  $f(x) \leq f(y)$  para todo  $y \in B(x, \varepsilon) \cap \mathbf{S}$ .
- $x \in \mathbf{S}$  es un minimizador local estricto de (P) si existe  $\varepsilon > 0$  tal que  $f(x) < f(y)$  para todo  $y \in B(x, \varepsilon) \cap \mathbf{S}$ ,  $y \neq x$ .
- $x \in \mathbf{S}$  es un maximizador global estricto de (P) si  $f(x) > f(y)$  para todo  $y \in \mathbf{S}$ , con  $y \neq x$ .

- $x \in \mathcal{S}$  es un maximizador global de  $(P)$  si  $f(x) \geq f(y)$  para todo  $y \in \mathcal{S}$ .
- $x \in \mathcal{S}$  es un maximizador local de  $(P)$  si existe  $\varepsilon > 0$  tal que  $f(x) \geq f(y)$  para todo  $y \in B(x, \varepsilon) \cap \mathcal{S}$ .
- $x \in \mathcal{S}$  es un maximizador local estricto de  $(P)$  si existe  $\varepsilon > 0$  tal que  $f(x) > f(y)$  para todo  $y \in B(x, \varepsilon) \cap \mathcal{S}$ , con  $y \neq x$ .

#### 1.2.4 Existencia de Solución para Problemas de Optimización

Nos interesamos en :

- la existencia de la solución óptima,
- caracterización de la solución óptima, y
- algoritmos para calcular la solución óptima.

Para el primer caso por ejemplo, considere el siguiente problema de optimización:

- $$(P) \quad \min_x \quad \frac{1+x}{2x}$$

$$\text{sujeto a } x \geq 1.$$

Aquí no existe solución óptima porque la región factible no es acotada.

- $$(P) \quad \min_x \quad \frac{1}{x}$$

$$\text{sujeto a } 1 \leq x < 2.$$

Aquí no existe solución óptima porque la región factible no es cerrada.

■

$$(P) \quad \min_x \quad f(x)$$

$$\text{sujeeto a } 1 \leq x < 2,$$

$$\text{donde } f(x) = \begin{cases} \frac{1}{x}, & x < 2 \\ 1, & x = 2 \end{cases}$$

Aquí no existe solución óptima porque la función  $f$  no es continua.

**Teorema 1.2.1 (Teorema de Weierstrass para sucesiones)** Sea  $\{x_k\}$  una sucesión infinita de puntos en el conjunto compacto  $\mathbf{S}$ . Entonces existe una subsucesión  $\{x_{k_j}\}$  de  $\{x_k\}$  que converge a un punto de  $\mathbf{S}$ .

**Teorema 1.2.2 (Teorema de Weierstrass para funciones)** Sea  $f$  una función continua sobre el conjunto compacto no vacío  $\mathbf{S} \subset \mathbb{R}^n$ . Entonces existe un punto en  $\mathbf{S}$  que minimiza (maximiza) a  $f$ ; es decir,  $f$  alcanza un mínimo absoluto (máximo absoluto) en  $\mathbf{S}$ .

### 1.2.5 Condiciones de Optimalidad para Problemas Sin Restricciones

Las condiciones para las soluciones local y global de problemas de optimización son intuitivas, y usualmente imposible de chequear directamente. En consecuencia, mostraremos más adelante como verificar fácilmente, las condiciones que son tanto necesarias para que un punto sea un minimizador local (lo que nos ayudaría a identificar los candidatos a ser minimizadores), como suficientes (que nos permitiría confirmar que el punto considerado es un minimizador local), o, algunas veces ambas.

$$(P) \quad \min \quad f(x)$$

$$\text{sujeto a } x \in \mathbf{X},$$

donde  $x = (x_1, \dots, x_n)^T \in \mathbb{R}^n$ ,  $f : \mathbb{R}^n \rightarrow \mathbb{R}$ , y  $\mathbf{X}$  es un abierto (usualmente  $\mathbf{X} = \mathbb{R}^n$ ).

### 1.2.6 Condiciones de Optimalidad: Necesarias y suficientes

Condición necesaria para optimalidad local: si “ $\hat{x}$  es un minimizador local de (P), then  $\hat{x}$  debe satisfacer...” Tal condición nos permite identificar todos los candidatos para la optimización local.

**Teorema 1.2.3** *Suponga que  $f$  es diferenciable en  $\hat{x}$ . Si existe un vector  $d$  tal que  $\nabla f(\hat{x})^T d < 0$ , entonces para todo  $\lambda > 0$  suficientemente pequeño,  $f(\hat{x} + \lambda d) < f(\hat{x})$  ( $d$  se llama la dirección de descenso si satisface la última condición).*

**Prueba:** Tenemos que

$$f(\hat{x} + \lambda d) = f(\hat{x}) + \lambda \nabla f(\hat{x})^T d + \lambda \|d\| \alpha(\hat{x}; \lambda d),$$

donde  $\alpha(\hat{x}; \lambda d) \rightarrow 0$  cuando  $\lambda \rightarrow 0$ . Acomodando tenemos

$$\frac{f(\hat{x} + \lambda d) - f(\hat{x})}{\lambda} = \nabla f(\hat{x})^T d + \|d\| \alpha(\hat{x}; \lambda d).$$

Ya que  $\nabla f(\hat{x})^T d < 0$  y  $\alpha(\hat{x}; \lambda d) \rightarrow 0$  cuando  $\lambda \rightarrow 0$ ,  $f(\hat{x} + \lambda d) - f(\hat{x}) < 0$  para todo  $\lambda > 0$  suficientemente pequeño.  $\diamond$

**Corolario 1.2.1** *Suponga que  $f$  es diferenciable en  $\hat{x}$ . Si  $\hat{x}$  es un minimizador local, entonces  $\nabla f(\hat{x}) = 0$  (a tal punto se le llama estacionario).*

**Prueba:** Si  $\nabla f(\hat{x}) \neq 0$ , entonces  $d = -\nabla f(\hat{x})$  es una dirección de descenso, por lo cual  $\hat{x}$  no puede ser un minimizador local.  $\diamond$

El corolario arriba es una condición de optimalidad necesaria de primer orden para un problema de optimización sin restricciones. Sin embargo, un punto estacionario puede ser un minimizador local, un maximizador local o ninguno de los dos. El siguiente teorema probará una condición de optimalidad necesaria de segundo orden. Primero, una definición:

**Definición 1.2.2** Una matriz  $M$  de orden  $n \times n$  es llamada simétrica si  $M_{ij} = M_{ji}$ . Una matriz simétrica  $M$  de orden  $n \times n$  es llamada

- *definida positiva* si  $x^T M x > 0 \forall x \in \mathbb{R}^n, x \neq 0$
- *semidefinida positiva* si  $x^T M x \geq 0 \forall x \in \mathbb{R}^n$
- *definida negativa* si  $x^T M x < 0 \forall x \in \mathbb{R}^n, x \neq 0$
- *semidefinida negativa* si  $x^T M x \leq 0 \forall x \in \mathbb{R}^n$
- *indefinida* si  $\exists x, y \in \mathbb{R}^n : x^T M x > 0 \wedge y^T M y < 0$

Decimos que  $M$  es *SDP* si  $M$  es simétrica y definida positiva. Similarmente, decimos que  $M$  es *SSDP* si  $M$  es simétrica y semidefinida positiva.

### Ejemplo 1

$$M = \begin{pmatrix} 2 & 0 \\ 0 & 3 \end{pmatrix}$$

es definida positiva.

### Ejemplo 2

$$M = \begin{pmatrix} 8 & -1 \\ -1 & 1 \end{pmatrix}$$

es definida positiva. Para ver esto, nótese que para  $x \neq 0$ ,

$$x^T M x = 8x_1^2 - 2x_1x_2 + x_2^2 = 7x_1^2 + (x_1 - x_2)^2 > 0.$$

Ya que  $M$  es una matriz simétrica, todos sus autovalores son números reales. Se puede demostrar que  $M$  es *SSDP* si y sólo si todos sus autovalores son no negativos, definida positiva si todos sus autovalores son positivos, etc.

**Teorema 1.2.4** *Suponga que  $f$  es dos veces diferenciable en  $\hat{x}$ . Si  $\hat{x}$  es un minimizador local de entonces,  $\nabla f(\hat{x}) = 0$  y  $H(\hat{x})$  (el Hessiano en  $\hat{x}$ ) es semidefinido positivo.*

**Prueba:** De la condición necesaria de primer orden se tiene  $\nabla f(\hat{x}) = 0$ . Supongamos que  $H(\hat{x})$  no es semidefinido positivo. Entonces  $\exists d$  tal que  $\nabla f(\hat{x})d < 0$ . Tenemos que

$$\begin{aligned} f(\hat{x} + \lambda d) &= f(\hat{x}) + \lambda \nabla f(\hat{x})^T d + \frac{1}{2} \lambda^2 d^T H(\hat{x}) d + \lambda^2 \|d\|^2 \alpha(\hat{x}; \lambda d) \\ &= f(\hat{x}) + \frac{1}{2} \lambda^2 d^T H(\hat{x}) d + \lambda^2 \|d\|^2 \alpha(\hat{x}; \lambda d) \end{aligned}$$

donde  $\alpha(\hat{x}; \lambda d) \rightarrow 0$  cuando  $\lambda \rightarrow 0$ . Acomodando tenemos

$$\frac{f(\hat{x} + \lambda d) - f(\hat{x})}{\lambda^2} = \frac{1}{2} d^T H(\hat{x}) d + \|d\|^2 \alpha(\hat{x}; \lambda d).$$

Si  $d^T H(\hat{x}) d < 0$  y  $\alpha(\hat{x}; \lambda d) \rightarrow 0$  cuando  $\lambda \rightarrow 0$ ,  $f(\hat{x} + \lambda d) - f(\hat{x}) < 0$  para todo  $\lambda > 0$  suficientemente pequeño -contradicción.  $\diamond$

**Ejemplo 1** Sea

$$f(x) = \frac{1}{2}x_1^2 + x_1x_2 + 2x_2^2 - 4x_1 - 4x_2 - x_2^3$$

Entonces

$$\nabla f(x) = (x_1 + x_2 - 4, x_1 + 4x_2 - 4 - 3x_2^2)^T,$$

y

$$H(x) = \begin{pmatrix} 1 & 1 \\ 1 & 4 - 6x_2 \end{pmatrix}.$$

$\nabla f(x) = 0$  tiene exactamente dos soluciones:  $\hat{x} = (4, 0)$  y  $\hat{x} = (3, 1)$ . Pero

$$H(\hat{x}) = \begin{pmatrix} 1 & 1 \\ 1 & -2 \end{pmatrix}.$$

es indefinida, en consecuencia, el único candidato para un mínimo local es  $\hat{x} = (4, 0)$ .

La condición necesaria solamente nos permiten obtener una lista de puntos candidatos a ser minimizadores. Ahora, la condición suficiente para optimalidad: "si  $\hat{x}$  satisface..., entonces  $\hat{x}$  es un minimizador local de  $(P)$ ."

**Teorema 1.2.5** Suponga que  $f$  es dos veces diferenciable por  $\hat{x}$ . Si  $\nabla f(\hat{x}) = 0$  y  $H(\hat{x})$  es definido positivo, entonces  $\hat{x}$  es un minimizador local(estricto).

**Prueba:** Tenemos que

$$f(x) = f(\hat{x}) + \frac{1}{2}(x - \hat{x})^T H(\hat{x})(x - \hat{x}) + \|(x - \hat{x})\|^T \alpha(\hat{x}; x - \hat{x}).$$

Suponga que  $\hat{x}$  no es un minimizador local estricto. Entonces existe una sucesión  $x_k \rightarrow \hat{x}$  tal que  $x_k \neq \hat{x}$  y  $f(x_k) \leq f(\hat{x})$  para todo  $k$ . Defina  $d_k = \frac{x_k - \hat{x}}{\|x_k - \hat{x}\|}$ . Entonces

$$f(x_k) = f(\hat{x}) + \|x_k - \hat{x}\|^2 \left( \frac{1}{2} d_k^T H(\hat{x}) d_k + \alpha(\hat{x}; x_k - \hat{x}) \right),$$

así

$$\frac{1}{2} d_k^T H(\hat{x}) d_k + \alpha(\hat{x}; x_k - \hat{x}) = \frac{f(x_k) - f(\hat{x})}{\|x_k - \hat{x}\|^2} \leq 0.$$

Ya que  $\|d_k\| = 1$  para cualquier  $k$ , existe una subsucesión de  $d_k$  que converge hacia algún punto  $d$  tal que  $\|d\| = 1$  (por el teorema 1.2.1). Sin perder generalidad suponga que  $d_k \rightarrow d$ .

Entonces

$$\lim_{k \rightarrow \infty} \frac{1}{2} d_k^T H(\hat{x}) d_k + \alpha(\hat{x}; x_k - \hat{x}) = \frac{1}{2} d^T H(\hat{x}) d,$$

que contradice el hecho de que  $H(\hat{x})$  es definido positivo.  $\diamond$

Nota:

- Si  $\nabla f(x) = 0$  y  $H(x)$  definido negativo, entonces  $\hat{x}$  es un maximizador local.
- Si  $\nabla f(x) = 0$  y  $H(x)$  semidefinido positivo, no se puede asegurar que  $\hat{x}$  es un minimizador local.

### 1.2.7 Convexidad y Minimización

#### Definición 1.2.3

- Sean  $x, y \in \mathbb{R}^n$ , los puntos de la forma  $\lambda x + (1 - \lambda)y$  para  $\lambda \in [0, 1]$  son llamados combinación convexa de  $x$  con  $y$ . Más generalmente, el punto  $y$  es combinación convexa de los puntos  $x_1, \dots, x_k$  si  $y = \sum_{i=1}^k \alpha_i x_i$  donde  $\alpha_i \geq 0 \quad \forall i$  y  $\sum_{i=1}^k \alpha_i = 1$ .



- Un subconjunto  $\mathbf{S} \subset \mathbb{R}^n$  es llamado convexo si  $\forall x, y \in \mathbf{S}$  y  $\lambda \in [0; 1]$ ,  $\lambda x + (1 - \lambda)y \in \mathbf{S}$ .
- Una función  $f : \mathbf{S} \rightarrow \mathbb{R}^n$ , donde  $\mathbf{S}$  es un conjunto convexo no vacío es una función convexa si

$$f(\lambda x + (1 - \lambda)y) \leq \lambda f(x) + (1 - \lambda)f(y) \quad \forall x, y \in \mathbf{S}, \quad \forall \lambda \in [0; 1].$$

- Una función  $f$  como arriba es llamada estrictamente convexa si la desigualdad arriba es estricta para todo  $x \neq y$  y  $\lambda \in (0, 1)$ .
- Una función  $f : \mathbf{S} \rightarrow \mathbb{R}^n$  es llamada cóncava (estrictamente cóncava) si es convexa (estrictamente convexa).

Consideraremos el problema de optimización que escribiremos abajo y enunciaremos algunos teoremas cuyas demostraciones se pueden ver en [10]:

$$(CP) \quad \min_x f(x)$$

sujeto a  $x \in \mathbf{S}$ .

**Teorema 1.2.6** *Suponga que  $\mathbf{S}$  es un conjunto convexo no vacío,  $f : \mathbf{S} \rightarrow \mathbb{R}$  es una función convexa, y  $\hat{x}$  es un minimizador local de (CP). Entonces  $\hat{x}$  es un minimizador global de  $f$  sobre  $\mathbf{S}$ .*

**Prueba:**(ver [10])

**Nota:**

- Un problema de minimizar una función convexa sobre una región factible convexa (tal como se ha considerado en el teorema) es un problema de programación convexa.
- Si  $f$  es estrictamente convexa, un minimizador local es el único minimizador global.
- Si  $f$  es (estrictamente) cóncava, un maximizador local es un (único) maximizador global.

El siguiente resultado nos ayuda a determinar cuando una función es convexa.

**Teorema 1.2.7** *Suponga que  $X \subset \mathbb{R}^n$  es un conjunto convexo abierto no vacío. Y que  $f : X \rightarrow \mathbb{R}$  es diferenciable. Entonces  $f$  es convexa si y sólo si ésta satisface la desigualdad del gradiente:*

$$f(y) \geq f(x) + \nabla f(x)^T (y - x) \quad \forall x, y \in X$$

**Prueba:**(ver [10])

En una dimensión, la desigualdad de gradiente tiene la forma

$$f(y) \geq f(x) + f'(x)(y - x) \quad \forall x, y \in X$$

El siguiente teorema proporciona otra condición suficiente y necesaria, para el caso cuando  $f$  es dos veces continuamente diferenciable.

**Teorema 1.2.8** *Suponga que  $X \subset \mathbb{R}^n$  es un conjunto convexo abierto no vacío,  $f : X \rightarrow \mathbb{R}$  es dos veces diferenciable. Entonces  $f$  es convexa si el Hesiano de  $f$ ,  $H(x)$ , es semidefinido positivo  $\forall x \in X$ .*

**Prueba:**(ver [10])

En una dimensión, el Hesiano es la segunda derivada de la función, la condición de semi-definido positivo se puede declara como  $f'' \geq 0 \quad \forall x \in X$ .

Otro teorema relacionado al anterior (condición suficiente) es:

**Teorema 1.2.9** *Suponga que  $X \subset \mathbb{R}^n$  es un conjunto convexo abierto no vacío,  $f : X \rightarrow \mathbb{R}$  es dos veces diferenciable. Entonces  $f$  es convexa si el Hesiano de  $f$ ,  $H(x)$ , es definido positivo  $\forall x \in X$ .*

Para problemas de optimización (sin restricciones) convexos, las condiciones de optimalidad de lo dicho en la sección anterior puede ser simplificada significativamente, proporcionando una única condición necesaria y suficiente para optimalidad global:

**Teorema 1.2.10** *Suponga que  $f : X \rightarrow \mathbb{R}$  es convexa y diferenciable en  $X$ . Entonces  $\hat{x} \in X$  es un minimizador global si y sólo si  $\nabla f(\hat{x}) = 0$ .*

**Prueba:**

La necesidad de la condición  $\nabla f(\hat{x}) = 0$  fue establecida con independencia de la convexidad de la función.

Suponga que  $\nabla f(\hat{x}) = 0$ . Entonces, por la desigualdad del gradiente

$$f(y) \geq f(\hat{x}) + \nabla f(\hat{x})^T(y - \hat{x}) = f(\hat{x})$$

para todo  $y \in X$ , y así  $\hat{x}$  es un minimizador global.  $\diamond$

**Ejemplo 2** Sea

$$f(x) = -\ln(1 - x_1 - x_2) - \ln x_1 - \ln x_2.$$

Entonces

$$\nabla f(x) = \begin{pmatrix} \frac{1}{1-x_1-x_2} - \frac{1}{x_1} \\ \frac{1}{1-x_1-x_2} - \frac{1}{x_2} \end{pmatrix},$$

y

$$H(x) = \begin{pmatrix} \left(\frac{1}{1-x_1-x_2}\right)^2 + \left(\frac{1}{x_1}\right)^2 & \left(\frac{1}{1-x_1-x_2}\right)^2 \\ \left(\frac{1}{1-x_1-x_2}\right)^2 & \left(\frac{1}{1-x_1-x_2}\right)^2 + \left(\frac{1}{x_2}\right)^2 \end{pmatrix}$$

Es fácil probar que  $f(x)$  es una función estrictamente convexa, y por lo tanto  $H(x)$  es definido positivo en su dominio  $X = \{(x_1, x_2) : x_1 > 0, x_2 > 0, x_1 + x_2 < 1\}$ . En  $\hat{x} = (\frac{1}{3}, \frac{1}{3})$  tendremos que  $\nabla f(\hat{x}) = 0$ , y así  $\hat{x}$  es el único minimizador local de  $f$ .

En el capítulo 2, se dará a conocer en detalle, el método de búsqueda lineal con condición de Wolfe en  $\mathbb{R}^n$  y en el capítulo 4 en Variedades Riemannianas.

---

## Capítulo 2

### Método de Búsqueda Lineal con Condición de Wolfe en $\mathbb{R}^n$

---

#### 2.1. Algoritmo de Optimización General

Recordemos: se quiere resolver el problema

$$(P) \quad \min_x \quad f(x) \\ \text{sujeto a } x \in \mathbb{R}^n$$

Casi siempre, las soluciones de los problemas de optimización, son imposible de obtener directamente; salvo algunas excepciones. En consecuencia, en su mayor parte, se resolverán esos problemas con algoritmos iterativos. Estos algoritmos normalmente requieren que el usuario suministre un punto de partida  $x_0$ . Comenzando por  $x_0$ , un algoritmo iterativo generará una sucesión de puntos  $\{x_k\}_{k=0}^{\infty}$  llamados *iterados*. Ahora bien, en la decisión de como generar el próximo iterado,  $x_{k+1}$ , el algoritmo utiliza la información de como se comporta la función  $f$  en el iterado actual,  $x_k$ , y algunas veces en los iterados pasados  $x_0, \dots, x_{k-1}$ . En la práctica, más que construir una sucesión infinita de iterados, el algoritmo se detiene, cuando se satisface un criterio de finalización apropiado, indicando ya sea

que el problema ha sido resuelto dentro de una precisión deseada, o que no se puede hacer ningún progreso.

La discusión, en nuestro caso, caerá en la categoría de los *algoritmos de búsqueda lineal*:

### Algoritmo de Optimización de Búsqueda Lineal General

**Inicio** Especifique un valor inicial de la solución  $x_0$ .

**Iteración** Para  $k = 1, 2, \dots$

Si  $x_k$  es el óptimo, parar.

De lo contrario,

- Determine  $d_k$  – una *dirección de búsqueda*
- Determine  $\lambda_k > 0$  – *el tamaño del paso*
- Determine  $x_{k+1} = x_k + \lambda_k d_k$  – un nuevo estimado de la solución

**La elección de la dirección.** Típicamente, se requiere que  $d_k$  sea una dirección de descenso de  $f$  por  $x_k$ , es decir,

$$f(x_k + \lambda d_k) < f(x_k) \quad \forall \lambda \in (0, \varepsilon]$$

para algún  $\varepsilon > 0$ . Para el caso cuando  $f$  es diferenciable, se demostró en el Teorema 1.2 que si  $\nabla f(x_k) \neq 0$ , cualquier  $d_k$  que cumpla con  $\nabla f(x_k)^T d_k < 0$  es una dirección de descenso.

Frecuentemente, la dirección es de la forma

$$d_k = -D_k \nabla f(x_k), \quad (2.1)$$

donde  $D_k$  es una matriz simétrica no singular.

Los siguientes son dos métodos básicos para escoger la matriz  $D_k$  por cada iteración; ello da como resultado dos algoritmos clásicos de optimización sin restricciones:

- *Descenso de mayor pendiente*:  $D_k = I, k = 1, 2, \dots$
- *Método de Newton* :  $D_k = H(x_k)^{-1}$  (probado que  $H(x_k) = \nabla^2 f(x_k)$  es definido positivo. )

Ahora si tenemos que  $D_k$  es una aproximación del Hessiano, estamos hablando del *método Cuasi-Newton* que se va actualizando en cada iteración mediante una fórmula de bajo-rango. Cuando  $d_k$  es definido en (2.1) y  $D_k$  es definida positiva, tendremos

$$\nabla f(x_k)^T d_k = -\nabla f(x_k)^T D_k \nabla f(x_k) < 0,$$

y en consecuencia  $d_k$  es una dirección de descenso.

**La elección del Tamaño del Paso.** En el cálculo de la longitud del paso,  $\lambda_k$ , nos encontramos ante una disyuntiva. Nos gustaría elegir  $\lambda_k$  para dar una reducción sustancial de  $f$ , pero al mismo tiempo, no queremos gastar demasiado tiempo en hacer la elección. Esto es, después de que  $d_k$  sea fijado, lo ideal sería que  $\lambda_k$  resolviera el problema de optimización unidimensional

$$\min_{\lambda_k \geq 0} f(x_k + \lambda d_k)$$

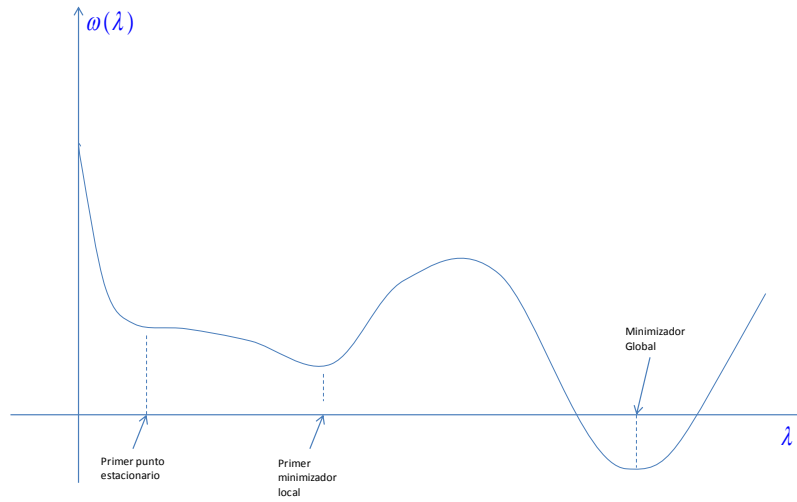


Figura 2.1: La longitud del paso ideal es el minimizador global.

Este problema de optimización suele ser también imposible de resolver exactamente. En su lugar,  $\lambda_k$  se calcula (a través de un procedimiento iterativo denominado búsqueda lineal), ya sea para resolver aproximadamente el anterior problema de optimización, o para asegurar una "suficiente" disminución en el valor de  $f$ . En detalle, la opción ideal sería el minimizador global de la función univariable  $\omega(\cdot)$  definida por

$$\omega(\lambda) = f(x_k + \lambda d_k), \quad \lambda > 0, \quad (2.2)$$

pero en general, cuesta mucho identificar este valor (véase la Figura 2.1). Para encontrar incluso un minimizador local del  $\omega$  con una precisión moderada, generalmente requiere demasiadas evaluaciones de la función objetivo  $f$  y, posiblemente, del gradiente  $\nabla f$ . Otras



estrategias prácticas realizan búsquedas lineales *inexactas* para identificar la longitud del paso que logran una adecuada reducción en  $f$  con un costo mínimo.

Típicamente el algoritmo de búsqueda lineal prueba una sucesión de valores candidatos para  $\lambda$  y, se detiene para aceptar uno de estos valores cuando se cumplen ciertas condiciones. La búsqueda lineal se realiza en dos fases: Una fase de horquillado o pruebas sucesivas para encontrar un intervalo que contiene longitudes de paso ideales, y una fase de bisección o interpolación que calcula una buena longitud de paso dentro de este intervalo.

Ahora discutiremos varias condiciones de finalización para el algoritmo de búsqueda lineal y mostrar que la longitud del paso efectiva no necesita yacer cerca del minimizador de la función real  $\omega$  definida en (2.2).

Una simple condición que podríamos imponer a  $\lambda_k$  es que produzca una reducción en  $f$ , es decir,  $f(x_k + \lambda_k d_k) < f(x_k)$ . Que este requisito no es suficiente para producir convergencia hacia  $\hat{x}$  se ilustra en la Figura (2.2), en la cual el valor mínimo de la función es  $\hat{f} = -1$ , pero una sucesión de *iterados*  $x_k$  para los cuales  $f(x_k) = \frac{5}{k}$ ,  $k = 1, 2, \dots$  admite un decrecimiento por cada iteración pero el valor del límite de función es cero. La dificultad es que no hay un *decrecimiento* suficiente de  $f$  en cada paso, pero este tema, lo discutiremos a continuación.

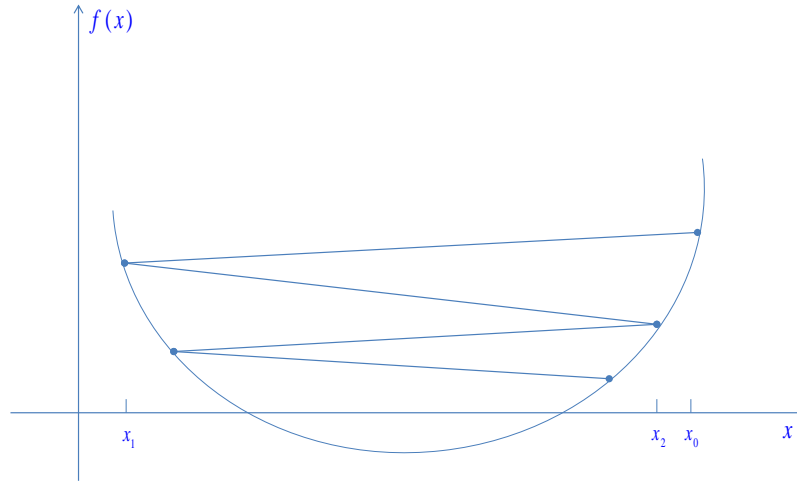


Figura 2.2: Insuficiente reducción en  $f$ .

## 2.2. La Condición de Wolfe

Una condición común de búsqueda lineal inexacta estipula que,  $\lambda_k$  debe en primer lugar dar una *disminución suficiente* a la función objetivo  $f$ , tal como se mide por la siguiente desigualdad:

$$f(x_k + \lambda d_k) < f(x_k) + c_1 \lambda \nabla f(x_k)^T d_k, \quad (2.3)$$

para alguna constante  $c_1 \in (0, 1)$ . En otras palabras, la reducción en  $f$  debe ser proporcional tanto a la longitud del paso  $\lambda_k$  y la derivada direccional  $\nabla f(x_k)^T d_k$ . La desigualdad (2.3) se llama la *condición de Armijo*. La regla de Armijo es uno de varios métodos de

búsqueda de lineal inexactas que garantiza un grado de precisión suficiente para asegurar la convergencia del algoritmo. La condición de disminución suficiente se ilustra en la figura 2.3. Aquí,  $\omega(\lambda) = f(x_k + \lambda d_k)$  implica que  $\omega'(\lambda) = \nabla f(x_k + \lambda d_k)^T d_k$  y, como  $d_k$  es una dirección de descenso, se tiene que  $\omega'(0) < 0$ . Entonces, la aproximación de primer orden de  $\omega(\lambda)$  por  $\lambda = 0$  es dada por  $\omega(0) + \lambda\omega'(0)$ . En consecuencia, el lado derecho de (2.3), es una función lineal, que la denotamos por  $l(\lambda)$ . La función  $l(\cdot)$  tiene pendiente negativa  $c_1 \nabla f(x_k)^T d_k$ , pero a causa de que  $c_1 \in (0, 1)$ , su gráfica esta por arriba de la gráfica de  $\omega$  para pequeños valores positivos de  $\lambda$ . Entonces la condición de disminución suficiente declara que  $\lambda$  es aceptable solo sí  $\omega(\lambda) \leq l(\lambda)$ . Los intervalos sobre los cuales esta condición se satisface se muestran en la figura 2.3. En la práctica,  $c_1$  se escoge muy pequeño, digamos  $c_1 = 10^{-4}$ .

La condición de disminución suficiente no basta por sí misma para garantizar que el algoritmo haga un progreso razonable, ya que como vemos en la figura 2.3, ésta se satisface para valores suficientemente pequeños de  $\lambda$ . Para descartar avances cortos inaceptables introducimos un segundo requisito, que llamaremos *condición de curvatura*, la cual requiere que  $\lambda_k$  satisfaga la siguiente desigualdad

$$\nabla f(x_k + \lambda_k d_k)^T d_k \geq c_2 \nabla f(x_k)^T d_k, \quad (2.4)$$

para alguna constante  $c_2 \in (c_1, 1)$ , donde  $c_1$  es la constante de (2.3). Nótese que el lado izquierdo es simplemente la derivada  $\omega'(\lambda_k)$ , por lo que la condición de curvatura asegura que la pendiente de  $\omega$  en  $\lambda_k$  es mayor que  $c_2$  veces la pendiente inicial  $\omega'(0)$ . Esto tiene

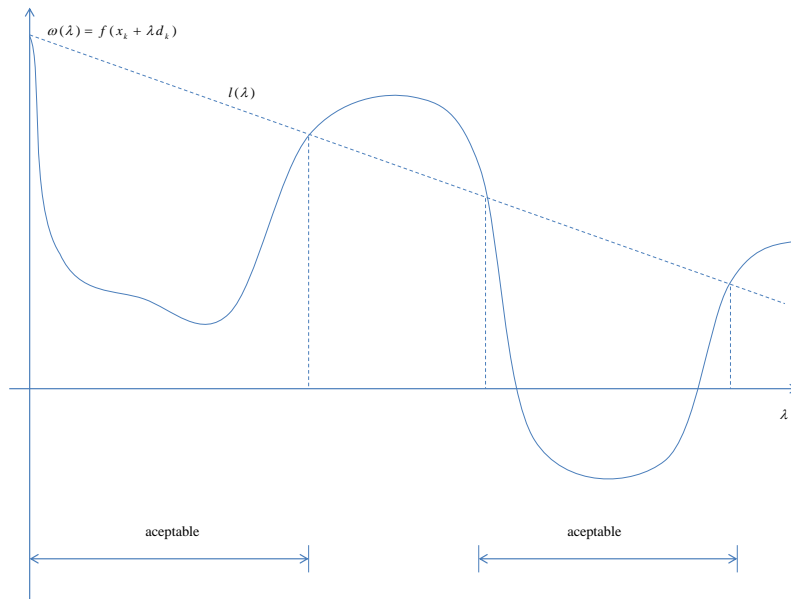


Figura 2.3: Condición de suficiente reducción.

sentido porque si la pendiente  $\omega'(\lambda_k)$  es fuertemente negativa, tenemos una indicación que podemos reducir significativamente a  $f$  moviéndonos más allá de la dirección elegida. Por otro lado, si  $\omega'(\lambda_k)$  es sólo ligeramente negativa, o incluso positiva, sería un signo que no podemos esperar más disminución de  $f$  en esta dirección, por lo que tendría sentido suspender la búsqueda lineal. La condición de curvatura se ilustra en la Figura 2.4. Los valores típicos de  $c_2$  son 0,9 cuando se elige la dirección de búsqueda  $d_k$  por un método de Newton o cuasi-Newton, y 0,1 cuando  $d_k$  se obtiene a partir de un método del gradiente conjugado no lineal.

Las condiciones de disminución suficiente y de curvatura se conocen colectivamente

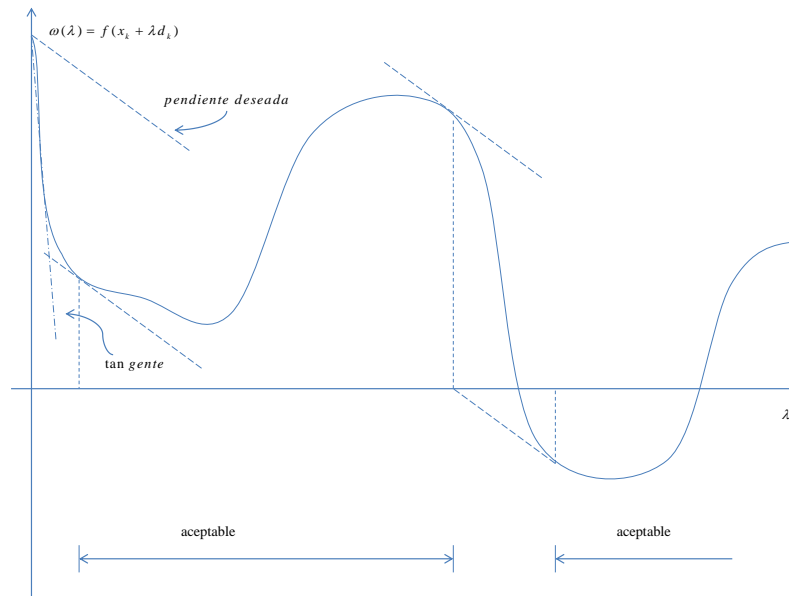


Figura 2.4: Condición de curvatura.

como *las condiciones de Wolfe*. Las ilustramos en la figura 3.5 y las reiteramos aquí para futuras referencias:

$$f(x_k + \lambda_k d_k) \leq f(x_k) + c_1 \lambda_k \nabla f(x_k)^T d_k, \quad (2.5a)$$

$$\nabla f(x_k + \lambda_k d_k)^T \lambda_k \geq c_2 \nabla f(x_k)^T d_k, \quad (2.5b)$$

con  $0 < c_1 < c_2 < 1$ . Una longitud de paso puede satisfacer las condiciones de Wolfe sin estar particularmente cerca de un minimizador de  $\omega$ , como se muestra en la Figura 2.5. Podemos, sin embargo, modificar la condición de curvatura y forzar a  $\lambda_k$  a caer en al menos en una vecindad amplia de un minimizador local o de un punto estacionario de  $\omega$ . También se tiene *las condiciones fuerte de Wolfe*, que requiere que  $\lambda_k$  satisfaga las

siguientes desigualdades:

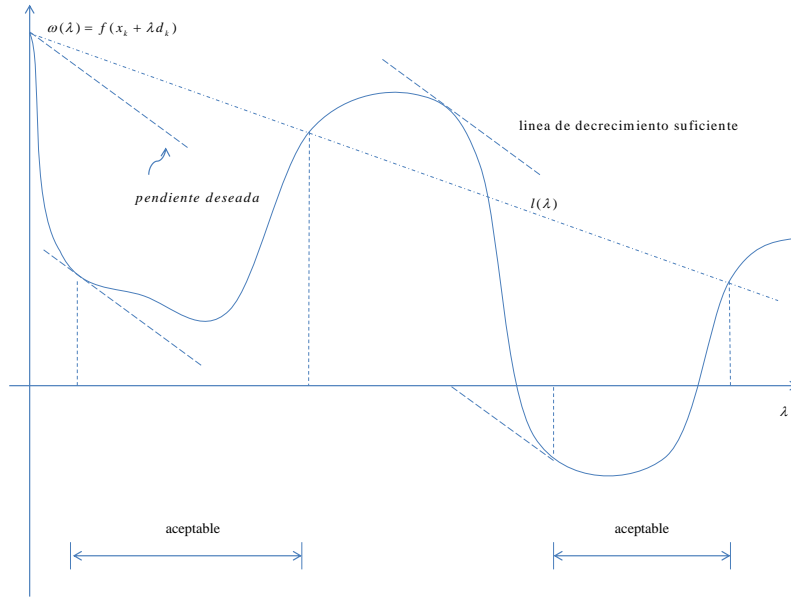


Figura 2.5: Longitudes de pasos que cumplen las condiciones de Wolfe .

$$f(x_k + \lambda_k d_k) \leq f(x_k) + c_1 \lambda_k \nabla f(x_k)^T d_k, \quad (2.6a)$$

$$| \nabla f(x_k + \lambda_k d_k)^T \lambda_k | \leq c_2 | \nabla f(x_k)^T d_k |, \quad (2.6b)$$

con  $0 < c_1 < c_2 < 1$ . La única diferencia con las condiciones de Wolfe es que ya no permitimos que la derivada  $\omega'(\lambda_k)$  sea demasiado positiva. Por lo tanto, se excluye puntos que están lejos de los puntos estacionarios de  $\omega$ .

No es difícil demostrar que existen longitudes de paso que satisfacen las condiciones

de Wolfe para cada función  $f$  que sea suave y acotada por debajo.

**Lema 2.2.1** *Suponga que  $f : \mathbb{R}^n \rightarrow \mathbb{R}$  es diferenciable continuamente. Sea  $d_k$  una dirección de descenso por  $x_k$ , con  $f$  acotada por abajo a lo largo del rayo  $\{x_k + \lambda d_k / \lambda > 0\}$ . Entonces si  $0 < c_1 < c_2 < 1$ , existen intervalos de longitud de paso que satisfacen las condiciones de Wolfe 2.5 y la condición fuerte de Wolfe 2.6*

**Demostración 2.2.1** *Ya que  $\omega(\lambda) = f(x_k + \lambda d_k)$  es acotada por abajo para toda  $\lambda > 0$  y tenemos que  $0 < c_1 < 1$ , entonces la recta  $l(\lambda) = f(x_k) + \lambda c_1 \nabla f(x_k)^T d_k$  debe interceptar el gráfico de  $\omega$  por lo menos una vez. Sea  $\lambda' > 0$  el más pequeño de estos valores, es decir,*

$$f(x_k + \lambda' d_k) = f(x_k) + \lambda' c_1 \nabla f(x_k)^T d_k. \quad (2.7)$$

*La condición de disminución suficiente (2.5a) claramente se cumple para toda longitud de paso menor que  $\lambda'$ .*

*Por el teorema del valor medio, existe  $\lambda'' \in (0, \lambda')$  tal que*

$$f(x_k + \lambda' d_k) - f(x_k) = \lambda' \nabla f(x_k + \lambda'' d_k)^T d_k. \quad (2.8)$$

*Por combinar (2.7) y (2.8), tenemos que*

$$\nabla f(x_k + \lambda'' d_k)^T d_k = c_1 \nabla f(x_k)^T d_k \geq c_2 \nabla f(x_k)^T d_k, \quad (2.9)$$

*ya que  $c_1 < c_2$  y  $\nabla f(x_k)^T d_k < 0$ . Por tanto,  $\lambda''$  satisface las condiciones de Wolfe 2.5, con desigualdad estricta en ambas ecuaciones (2.5a) y (2.5b). Así, ya que  $f$  es suave, existe un intervalo que contiene a  $\lambda''$  para el cual las condiciones de Wolfe se sostienen.*

Más aún, ya que el término del lado izquierdo de (2.9) es negativo, las condiciones fuertes de Wolfe 2.6 se cumplen en el mismo intervalo.◊



---

## Capítulo 3

### Variedades

---

En este capítulo daremos a conocer el ámbito en donde extenderemos el método de búsqueda lineal con condición de Wolfe, las variedades diferenciables:

**Definición 3.0.1** *Una variedad diferenciable de dimension  $n$ , es un conjunto  $M$  y una familia de aplicaciones inyectivas  $\mathcal{X}_\alpha : U_\alpha \rightarrow M$ ,  $\alpha \in I$ , definidas en abiertos  $U_\alpha$  de  $\mathbb{R}^n$  en  $M$  tales que cumplen las siguientes condiciones:*

- $M = \bigcup_{\alpha \in I} \mathcal{X}_\alpha(U_\alpha)$
- *Para todo par  $\mathcal{X}_\alpha$  y  $\mathcal{X}_\beta$  con  $\mathcal{X}_\alpha(U_\alpha) \cap \mathcal{X}_\beta(U_\beta) = W \neq \emptyset$ , los conjuntos  $\mathcal{X}_\alpha^{-1}(W)$  y  $\mathcal{X}_\beta^{-1}(W)$  son abiertos en  $\mathbb{R}^n$  y las aplicaciones  $\mathcal{X}_\beta^{-1} \circ \mathcal{X}_\alpha : \mathcal{X}_\alpha^{-1}(W) \rightarrow \mathcal{X}_\beta^{-1}(W)$  son diferenciables.*

*El par  $(U_\alpha, \mathcal{X}_\alpha)$  con  $p \in \mathcal{X}_\alpha(U_\alpha)$  es llamado una parametrización de  $M$  de  $p$ ,  $\mathcal{X}_\alpha(U_\alpha)$  es llamada vecindad coordinada de  $p$ . Una familia  $\{(U_\alpha, \mathcal{X}_\alpha)\}$  satisfaciendo los items 1 y 2 es llamada estructura diferenciable de  $M$ .*

**Proposición 3.0.1** *Toda superficie regular de  $\mathbb{R}^n$  de dimensión  $k$  es una variedad diferenciable de la misma dimensión.*

**Prueba:** ver [14]

**Proposición 3.0.2** *Si  $M_1$  y  $M_2$  son dos variedades diferenciables de dimensión  $m_1$  y  $m_2$  respectivamente, entonces el producto cartesiano  $M_1 \times M_2$  es una variedad de dimensión  $m_1 + m_2$ .*

**Prueba:** Ver [16]

### 3.1. Aplicaciones diferenciables entre variedades

**Definición 3.1.1** *Sea  $f : U \subset M \rightarrow \mathbb{R}$  una función definida en un subconjunto abierto  $U$  de una variedad diferenciable  $M$ . Diremos que  $f$  es diferenciable en  $p \in U$ , si para alguna parametrización  $\mathcal{X}_\alpha : U_\alpha \subset \mathbb{R}^n \rightarrow M$ , con  $p \in \mathcal{X}_\alpha(U_\alpha) \subset U$ , la composición  $f \circ \mathcal{X}_\alpha : U_\alpha \subset \mathbb{R}^n \rightarrow \mathbb{R}$  es diferenciable en  $\mathcal{X}_\alpha^{-1}(p)$ . Se dice que  $f$  es diferenciable en  $U$  si es diferenciable en todo punto de  $U$ .*

**Definición 3.1.2** *Una curva  $\gamma$  sobre una variedad diferenciable  $M$  es una aplicación  $\gamma : I \rightarrow M$  donde  $I = (-\varepsilon, \varepsilon)$ . Diremos que  $\gamma$  es diferenciable en  $t_0 \in I$  si para alguna parametrización  $\mathcal{X}_\alpha : U_\alpha \subset \mathbb{R}^n \rightarrow M$  con  $\gamma(t_0) \in \mathcal{X}_\alpha(U_\alpha)$ , la compuesta  $\mathcal{X}_\alpha^{-1} \circ \gamma : I \rightarrow U_\alpha$  es diferenciable en  $t_0$ , donde  $\gamma(I) \subset \mathcal{X}_\alpha(U_\alpha)$ . Se dice que  $\gamma$  es diferenciable en  $I$  si es di-*

ferenciabile en todo  $t \in I$ .

**Definición 3.1.3** Sean  $M_1$  y  $M_2$  variedades diferenciables de dimensión  $m$  y  $n$  respectivamente. Una aplicación  $\varphi : M_1 \rightarrow M_2$  es diferenciable en  $p \in V$ , si dadas  $\mathcal{X}_1 : U_1 \subset \mathbb{R}^n \rightarrow M_1$  parametrización de  $M_1$  en  $p$  y  $\mathcal{X}_2 : U_2 \subset \mathbb{R}^m \rightarrow M_2$  parametrización de  $M_2$  en  $\varphi(p)$ , con  $\varphi(\mathcal{X}_1(U_1)) \subset \mathcal{X}_2(U_2)$ , la aplicación  $\mathcal{X}_2^{-1} \circ \varphi \circ \mathcal{X}_1 : U_1 \subset \mathbb{R}^n \rightarrow \mathbb{R}^m$  es diferenciable en  $\mathcal{X}_1^{-1}(p)$

**Definición 3.1.4** Una aplicación  $\varphi : M_1 \rightarrow M_2$  es diferenciable entre dos variedades diferenciables. Decimos que  $\varphi$  es difeomorfismo si  $\varphi$  es biyectiva y  $\varphi^{-1}$  es diferenciable.

## 3.2. Espacio Tangente

**Definición 3.2.1** Sea  $M$  una variedad diferenciable. Consideremos una curva diferenciable  $\gamma : (-\varepsilon, \varepsilon) \rightarrow M$ , donde  $\gamma(0) = p$  y sea  $D_p$  el conjunto de las funciones de  $M$  en  $\mathbb{R}$  diferenciables en  $p$ ,  $D_p = \{f : M \rightarrow \mathbb{R} / f \text{ es diferenciable en } p\}$ . Se define el **Vector Tangente** a la curva  $\gamma$  en  $t = 0$  como la función  $\gamma'(0) : D_p \rightarrow \mathbb{R}$  dada por :

$$\gamma'(0)f = \left. \frac{d(f \circ \gamma)}{dt} \right|_{t=0}, f \in D_p. \quad (3.1)$$

**Definición 3.2.2** (Espacio tangente). El espacio tangente a una variedad  $M$  en un punto  $p$  representado por  $T_p M$ , es el conjunto de todos los vectores tangentes a  $M$  en  $p$ . Así,  $T_p M = \{v \in \mathbb{R}^m : v \text{ es un vector tangente en } p\}$ .

**Observación 3.2.1** Si escogemos una parametrización  $\mathcal{X}_\alpha : U \subset \mathbb{R}^n \rightarrow M$  con  $p = \mathcal{X}_\alpha(0)$  y  $q \in U$  podemos restringir la función  $f$  y la curva  $\gamma$  en esta parametrización por :

$$f \circ \mathcal{X}(q) = f(q) = f(q_1, \dots, q_n) \text{ y}$$

$$\mathcal{X}^{-1} \circ \gamma(t) = (q_1(t), \dots, q_n(t)),$$

restringiendo  $f$  a  $\gamma$ , tenemos

$$\begin{aligned} \gamma'(0)f &= \left. \frac{d}{dt}(f \circ \gamma) \right|_{t=0} = \left. \frac{d}{dt}f(q_1(t), \dots, q_n(t)) \right|_{t=0} \\ &= \sum_{i=1}^n q'_i(0) \left( \frac{\partial f}{\partial q_i} \right)_0 = \left( \sum_{i=1}^n q'_i(0) \left( \frac{\partial}{\partial q_i} \right)_0 \right) f. \end{aligned}$$

Así,

$$\gamma'(0) = \sum_{i=1}^n q'_i(0) \left( \frac{\partial}{\partial q_i} \right)_0 \tag{3.2}$$

es la expresión del vector tangente a  $f$  en  $p$  con relación a la parametrización  $\mathcal{X}$ .

**Observación 3.2.2** La expresión (3.2) muestra que el vector tangente a una curva  $\gamma$  en  $p$  depende de las derivadas,  $q'_i(0)$ , de  $\gamma$  en un sistema de coordenadas. De la elección de una parametrización obtenemos  $n$  vectores  $\left\{ \left( \frac{\partial}{\partial x_1} \right), \dots, \left( \frac{\partial}{\partial x_n} \right) \right\}$  en  $T_p M$  que generan, por (3.2), los vectores en  $T_p M$ .

**Proposición 3.2.1** El espacio tangente de una variedad diferenciable que es un subconjunto abierto de  $\mathbb{R}^n$  es el propio  $\mathbb{R}^n$ .

**Prueba:** Ver [14]

**Definición 3.2.3** Sean  $M_1$  y  $M_2$  dos variedades diferenciables de dimensión  $n$  y  $m$  respectivamente y sea  $\varphi : M_1 \rightarrow M_2$  una aplicación diferenciable. Para cada  $p \in M_1$  y cada  $v \in T_p M_1$ , escojamos una curva diferenciable  $\alpha : (-\varepsilon, \varepsilon) \rightarrow M$ , con  $\alpha(0) = p$  y  $\alpha'(0) = v$ . Definiendo  $\beta = \varphi \circ \alpha$  la aplicación:

$$d\varphi_p : T_p M_1 \rightarrow T_{\varphi(p)} M_2,$$

dada por  $d\varphi_p(v) = \beta'(0)$  es una aplicación lineal que no depende de la elección  $\alpha$ . Esta aplicación es llamada la diferencial de  $\varphi$  en  $p$ .

### 3.3. Métricas Riemannianas

**Definición 3.3.1** Sea  $M$  una variedad diferenciable. Una métrica Riemanniana es una aplicación que asocia a cada  $p \in M$  un producto interno (bilineal, simétrico y definido positivo)  $\langle \cdot, \cdot \rangle_p$ , dado por :

$$\langle \cdot, \cdot \rangle_p = T_p M \times T_p M \rightarrow \mathbb{R},$$

que varía diferenciablemente en el siguiente sentido: si  $\mathcal{X} : U \subset \mathbb{R}^n \rightarrow M$  es un sistema de coordenadas locales en torno de  $p$ , con  $\mathcal{X}(x_1, x_2, x_3, \dots, x_n) = q \in \mathcal{X}(U)$  y  $\frac{\partial}{\partial x_i}(q) = d\mathcal{X}_q(0, 0, \dots, 0, 1, 0, \dots, 0, 0)$ , entonces la función:  $g_{ij} : U \rightarrow \mathbb{R}$  definida por:

$$g_{ij}(x_1, x_2, \dots, x_n) = \left\langle \frac{\partial}{\partial x_i}(q), \frac{\partial}{\partial x_j}(q) \right\rangle_p,$$

es una función diferenciable en  $U$ .

Las funciones  $g_{ij}$  son llamadas expresiones de la métrica riemanniana en el sistema coordenado  $\mathcal{X}$  y la matriz  $G = (g_{ij})$  es la representación de la métrica riemanniana.

**Definición 3.3.2** (*Variedad riemanniana*). Una variedad diferenciable para la cual se define una métrica riemanniana se denomina una variedad riemanniana.

**Definición 3.3.3** El Producto interno de dos vectores  $u, v \in T_p M$  es definido por  $\langle u, v \rangle_p = g_x(u, v)$ , donde  $g_x$  es la métrica riemanniana evaluada en el punto  $x$ . La norma de un vector  $v \in T_p M$  es  $\|v\|_p = \sqrt{\langle v, v \rangle_p}$ .

**Definición 3.3.4** Un grupo de Lie es una variedad diferenciable  $M$  con una estructura de grupo  $*$ , tal que la aplicación  $M \times M \rightarrow M$  dada por  $(x, y) \rightarrow x * y^{-1}$  es diferenciable. Para un elemento  $y \in M$  la traslación a izquierda por  $x$  es el mapeo  $L_x : M \rightarrow M$  definido por  $L_x(y) = x * y$ , el cual es un difeomorfismo.

**Definición 3.3.5** Sea  $M$  un grupo de Lie. Una métrica riemanniana  $G$  en  $M$  es invariante por la izquierda si

$$\langle u, v \rangle_y = \langle d(L_x)u, d(L_x)v \rangle, \text{ para todo } x, y \in M \text{ y } u, v \in T_y M.$$

### 3.4. Subvariedades Encajadas (Embedding)

Un conjunto  $\mathcal{X}$  puede admitir varias estructuras múltiples. Sin embargo, si el conjunto  $\mathcal{X}$  es un subconjunto de una variedad  $(M, \mathcal{U}_\alpha)$ , entonces, esta admite a lo más una estructura de subvariedad.

### 3.4.1 Teoría General

Sean  $(M, \mathcal{U}_\alpha)$  y  $(N, \mathcal{V}_\alpha)$  variedades tal que  $N \subset M$ . La variedad  $(N, \mathcal{V}_\alpha)$  es llamada una *subvariedad inmersa* de  $(M, \mathcal{U}_\alpha)$  si la aplicación inclusión  $i : N \rightarrow M, i(x) = x$  es una inmersión.

Sea  $(N, \mathcal{V}_\alpha)$  una subvariedad de  $(M, \mathcal{V}_\alpha)$ . Ya que  $M$  y  $N$  son variedades, también son espacios topológicos con su variedad topológica. Si la variedad topológica de  $N$  coincide con su topología de subespacio inducida desde el espacio topológico  $M$ , entonces  $N$  se llama *subvariedad encajada o embedded*, o una *subvariedad regular*, o simplemente una subvariedad de la variedad  $M$ .

**Proposición 3.4.1** *Sea  $N$  un subconjunto de un variedad  $M$ . Entonces  $N$  admite a lo más una estructura diferenciable que hace que sea una subvariedad encajada o embedded de  $M$ .*

**Prueba:** Ver [16]

A  $M$  en la proposición anterior se le llama el *espacio encajo o embedding*. Así por ejemplo, cuando el espacio incrustado  $\mathbb{R}^{n+p}$  o un subconjunto abierto de  $\mathbb{R}^{n+p}$ , decimos que  $N$  es una *subvariedad matriz*

**Proposición 3.4.2** *Un subconjunto  $N$  de una variedad  $M$  es una subvariedad encajada  $d$ -dimensional de  $M$  si y sólo si, alrededor cada punto  $x \in N$ , existe una carta  $(\mathcal{U}, \varphi)$  de  $M$  tal que*

$$N \cap \mathcal{U} = \{x \in \mathcal{U} : \varphi(x) \in \mathbb{R}^d \times \{0\}\}.$$

En este caso, la carta  $(N \cap \mathcal{U}, \varphi)$ , donde  $\varphi$  es vista como una aplicación en  $\mathbb{R}^d$ , es una carta de la subvariedad encajada  $N$ .

**Prueba:** Ver [16]

Las siguientes proposiciones proporcionan condiciones suficientes para que los subconjuntos de variedades puedan ser subvariedades encajadas o embedded, simplemente una subvariedad de la variedad  $M$ . (Para las pruebas ver [16])

**Proposición 3.4.3** (Teorema de submersión) Sea Sea  $F : M_1 \rightarrow M_2$  una aplicación suave entre dos variedades de dimensión  $d_1$  y  $d_2$ ,  $d_1 > d_2$ , y sea  $c$  un punto de  $M_2$ . Si  $c$  es un valor regular (es decir, si el rango de  $F$  es igual a  $d_2$  para cada punto de  $F^{-1}(c)$ ) entonces  $F^{-1}(c)$  es una subvariedad encajada cerrada de  $M_1$ , y  $\dim(F^{-1}(c)) = d_1 - d_2$ .

**Proposición 3.4.4** (Teorema de subinmersión) Sea Sea  $F : M_1 \rightarrow M_2$  una aplicación suave entre dos variedades de dimensión  $d_1$  y  $d_2$ , y sea  $c$  un punto de  $F(M_1)$ . Si  $F$  tiene rango constante igual a  $k < d_1$  en una vecindad de  $(F^{-1}(c))$ , entonces  $(F^{-1}(c))$  es una subvariedad encajada cerrada de  $M_1$  de dimensión igual a  $d_1 - k$ .

Las funciones en subvariedades embebidas no plantean ninguna dificultad en particular. Sea  $N$  una subvariedad encajada de una variedad  $M$ . Si  $f$  es una función suave en  $M$ , entonces  $f|_N$ , la *restricción* de  $F$  sobre  $N$ , también es una función suave en  $N$ . Contrariamente, cualquier función suave en  $N$  se puede escribir localmente como una restricción de una función suave definida en un subconjunto abierto  $\mathcal{U} \subset M$ .



## 3.5. Campos de vectores, conexiones afines y derivada covariante

**Definición 3.5.1** (*Campo de vectores en una variedad diferenciable*). Un campo de vectores  $X$  en una variedad diferenciable  $M$  es una correspondencia que a cada punto  $p \in M$  asocia un vector  $X(p) \in T_pM$ .

Considerando una parametrización  $\mathcal{X} : U \subset \mathbb{R}^n \rightarrow M$ , es posible escribir:

$$X(p) = \sum_{i=1}^n a_i(p) \left( \frac{\partial}{\partial x_i} \right)_p,$$

donde cada  $a_i : M \rightarrow \mathbb{R}$  es una función en  $M$  y  $\left\{ \left( \frac{\partial}{\partial x_i} \right)_p \right\}$  es una base asociada a  $\mathcal{X}$ ,  $1 \leq i \leq n$ . Diremos que  $X$  es diferenciable si y solo si las funciones  $a_i$  son diferenciables.

**Definición 3.5.2** (*Campo de vectores a lo largo de curvas*). Un campo vectorial  $V$  a lo largo de una curva  $\alpha : I \rightarrow M$  es una aplicación que a cada  $\alpha(t) \in M$  asocia un vector tangente  $V(t) \in T_{\alpha(t)}M$ .

### Conexiones Afines.

Denotemos  $TM$  como el conjunto de espacios tangentes definidos en  $M$ . Sea  $\mathcal{H} = \mathcal{H}(M) = \{X : M \rightarrow TM : \text{para cada } p \in M, X(p) \in T_pM, \text{ y } X \in \mathcal{C}^\infty\}$  el conjunto de campo de vectores de clase  $\mathcal{C}^\infty$  y  $D = D(M) = \{f : M \rightarrow \mathbb{R} : f \in \mathcal{C}^\infty\}$  el conjunto de funciones reales de clase  $\mathcal{C}^\infty$ .

**Definición 3.5.3** Una conexión afín es una aplicación  $\nabla : \mathcal{H} \times \mathcal{H} \rightarrow \mathcal{H}$  donde a cada par de campos  $(X, Y)$  se asocia otro campo  $\nabla_X Y$  tal que para todo  $X, Y, Z \in \mathcal{H}$ , y  $f, g \in \mathcal{D}$  se cumple:

$$\text{i) } \nabla_{fX+gY} Z = f\nabla_X Z + g\nabla_Y Z;$$

$$\text{ii) } \nabla_X (Y + Z) = \nabla_X Y + \nabla_X Z;$$

$$\text{iii) } \nabla_X fY = f\nabla_X Y + X(f)Y, \text{ donde } X(f) = \sum_{i=1}^n a_i(\cdot) \frac{\partial f(\cdot)}{\partial x_i}$$

**Proposición 3.5.1** (Derivada Covariante) Sea  $M$  una variedad diferenciable con una conexión afín  $\nabla$ . Entonces existe una única aplicación que asocia a un campo vectorial  $V$  a lo largo de una curva diferenciable  $\alpha : I \rightarrow M$  otro campo vectorial  $\frac{DV}{dt}$  a lo largo de  $\alpha$ , denominado derivada covariante de  $V$  a lo largo de  $\alpha$ , tal que para todo  $V, W$  campo de vectores a lo largo de  $\alpha$  y  $f : I \rightarrow \mathbb{R}$  una función diferenciable en  $I$  se cumple:

$$\text{i) } \frac{D}{dt}(V + W) = \frac{DV}{dt} + \frac{DW}{dt}.$$

$$\text{ii) } \frac{D}{dt}(fV) = \frac{df}{dt}V + f\frac{DV}{dt}.$$

$$\text{iii) } \text{si } V(t) = Y(\alpha(t)), \text{ donde } Y \in \mathcal{H}, \text{ entonces } \frac{DV}{dt} = \nabla_{\frac{d\alpha}{dt}} Y.$$

**Prueba:** Ver [17]

**Expresión de la conexión afín en términos de las coordenadas locales.** Suponga que los campos de vectores  $X, Y \in \mathcal{H}$  son representados en una cierta vecindad local  $\mathcal{X} : U \subset \mathbb{R}^n \rightarrow M$  de un punto  $p$ , por:

$$X = \sum_{i=1}^m x_i X_i, Y = \sum_{i=1}^m y_i X_i,$$

donde  $X_i = \frac{\partial}{\partial x_i}$  representan los vectores de la base del sistema de coordenadas locales.

Usando las propiedades de la definición de la conexión afín:

$$\begin{aligned}\nabla_X Y &= \nabla_{\sum x_i X_i} \left[ \sum_j y_j X_j \right] = \sum_i x_i \left[ \nabla_{X_i} \left( \sum_j y_j X_j \right) \right] \\ &= \sum_i x_i \left[ \sum_j (y_j \nabla_{X_i} X_j) \right] + \sum_i x_i \left[ \sum_j \left( \frac{\partial y_j}{\partial x_i} X_j \right) \right],\end{aligned}$$

escribiendo  $\nabla_{X_i} X_j$  en función de la base local:

$$\nabla_{X_i} X_j = \sum_{k=1}^n \Gamma_{ij}^k X_k, \quad (3.3)$$

sustituyendo en la ecuación anterior obtenemos:

$$\nabla_X Y = \sum_{k=1}^n \left( \sum_{i,j=1}^n x_i y_j \Gamma_{ij}^k + \sum_i x_i \frac{\partial y_k}{\partial x_i} \right) X_k.$$

**Definición 3.5.4** *Los símbolos de Christoffel, o coeficientes de la conexión afín  $\nabla$  en  $U$ , son las funciones diferenciables  $\Gamma_{ij}^k : U \subset M \rightarrow R$  definidas por (3.3)*

**Expresión de la derivada covariante en términos de las coordenadas locales y de los símbolos de Christoffel.**

Sea  $\mathcal{X} : U \rightarrow M$  un sistema de coordenadas locales en torno de  $p \in M$ . Un resultado obtenido al demostrar la Proposición (3.5.1) es:

$$\frac{DV}{dt} = \sum_{j=1}^n \frac{dv^j}{dt} X_j + \sum_{i,j=1}^n v^j \frac{dx_i}{dt} \nabla_{X_i} X_j,$$

sustituyendo  $\nabla_{X_i} X_j = \sum_{k=1}^n \Gamma_{ij}^k X_k$  en la ecuación anterior tenemos

$$\begin{aligned} \frac{DV}{dt} &= \sum_{j=1}^n \frac{dv^j}{dt} X_j + \sum_{i,j=1}^n v^j \frac{dx_i}{dt} \left( \sum_{k=1}^n \Gamma_{ij}^k X_k \right), \\ &= \sum_{j=1}^n \frac{dv^j}{dt} X_j + \sum_{k=1}^n \sum_{i,j=1}^n v^j \frac{dx_i}{dt} \Gamma_{ij}^k X_k \end{aligned}$$

y así,

$$\frac{DV}{dt} = \sum_{k=1}^n \left( \frac{dv^k}{dt} + \sum_{i,j=1}^n v^j \frac{dx_i}{dt} \Gamma_{ij}^k \right) X_k, \quad (3.4)$$

es la expresión de la derivada covariante en términos de coordenadas locales y de los símbolos de Christoffel.

**Definición 3.5.5** (*Geodésicas*) Una curva parametrizada  $\alpha : I \rightarrow M$  es una geodésica si el campo tangente  $\frac{d\alpha}{dt}$  verifica:

$$\frac{D}{dt} \left( \frac{d\alpha}{dt} \right) = 0. \quad (3.5)$$

**Definición 3.5.6** (*Campos paralelos*) Dado  $M$  una variedad diferenciable, una conexión afín  $\nabla$  y un campo  $V$  a lo largo de una curva diferenciable  $\alpha : I \rightarrow M$ ,  $V$  es denominado campo paralelo si

$$\frac{DV}{dt} = 0, \quad \forall t \in I. \quad (3.6)$$

Así, si  $\alpha$  es una geodésica, entonces  $\frac{d\alpha}{dt}$  es paralelo.

**Proposición 3.5.2** Sea  $M$  una variedad diferenciable con una conexión afín  $\nabla$ . Sea  $\alpha : I \rightarrow M$  una curva diferenciable en  $M$  y  $V_0$  un vector tangente a  $M$  en  $\alpha(t_0)$ ,  $t_0 \in I$ .

Entonces existe un único campo de vectores paralelo  $V$  a lo largo de  $\alpha$ , tal que  $V(t_0) = V_0$  ( $V(t)$  es llamado transporte paralelo de  $V(t_0)$  a lo largo de  $\alpha$ ).

**Prueba:** Ver [17]

### Ecuaciones Geodésicas

De la expresión (3.4), un campo paralelo  $V$  es determinado por las ecuaciones

$$\sum_{k=1}^n \left( \frac{dv^k}{dt} + \sum_{i,j=1}^n v^j \frac{dx_i}{dt} \Gamma_{ij}^k \right) X_k = 0,$$

de forma equivalente,

$$\frac{dv^k}{dt} + \sum_{i,j=1}^n v^j \frac{d\alpha_i}{dt} \Gamma_{ij}^k = 0, \quad k = 1, 2, \dots, n.$$

Cuando se trata de una geodésica  $\alpha(t) = (\alpha_1(t), \dots, \alpha_n(t)(t))$ , se tiene  $v^i = \frac{d\alpha_i}{dt}$ , entonces esta última ecuación se transforma en:

$$\frac{d}{dt} \left( \frac{d\alpha_k}{dt} \right) + \sum_{i,j=1}^n \frac{d\alpha_j}{dt} \frac{d\alpha_i}{dt} \Gamma_{ij}^k = 0, \quad k = 1, 2, \dots, n,$$

así

$$\frac{d^2\alpha_k}{dt^2} + \sum_{i,j=1}^n \Gamma_{ij}^k \frac{d\alpha_j}{dt} \frac{d\alpha_i}{dt} = 0, \quad k = 1, 2, \dots, n, \quad (3.7)$$

el cual es un sistema de  $n$  ecuaciones diferenciales de 2do. orden, que posee solución única en algún intervalo  $I = [a, b]$ , verificando  $x(0) = \alpha(0) = p$  y  $\frac{dx}{dt}(0) = \alpha'(0) = v$ .

### Conexión afín en variedades riemannianas

**Definición 3.5.7** Sea  $M$  una variedad diferenciable con una conexión afín  $\nabla$  y una métrica riemanniana  $\langle, \rangle$ . Se dice que  $\nabla$  es compatible con la métrica  $\langle, \rangle$ , si para todo par de

campos de vectores  $V$  y  $W$  a lo largo de la curva diferenciable  $\alpha : I \rightarrow M$  se tiene:

$$\frac{d}{dt}\langle V, W \rangle = \left\langle \frac{DV}{dt}, W \right\rangle + \left\langle V, \frac{DW}{dt} \right\rangle \quad (3.8)$$

**Proposición 3.5.1** *Si la conexión afín  $\nabla$  es compatible con  $\langle, \rangle$  y  $V, W$  son campos paralelos a lo largo de una curva diferenciable  $\alpha : I \rightarrow M$  entonces,  $\langle V, W \rangle$  es constante. En particular si  $\alpha(t) = (\alpha_1(t), \dots, \alpha_n(t))$  es una geodésica,  $\langle \frac{d\alpha}{dt}, \frac{d\alpha}{dt} \rangle$  es constante.*

**Prueba:** Ver [17]

**Definición 3.5.8** *Una conexión afín  $\nabla$  en una variedad riemanniana  $M$  es compatible con  $\langle, \rangle$  si y solamente si:*

$$X\langle Y, Z \rangle = \langle \nabla_X Y, Z \rangle + \langle Y, \nabla_X Z \rangle, \quad \forall X, Y, Z \in \mathcal{H}.$$

**Definición 3.5.9** *Una conexión afín  $\nabla$  en una variedad riemanniana  $M$  es llamada simétrica si:*

$$\nabla_X Y - \nabla_Y X = [X, Y] \quad \forall X, Y \in \mathcal{H}, \text{ donde } [X, Y] = XY - YX.$$

**Definición 3.5.10** *(Levi-Civita). Dada una variedad riemanniana  $M$ , existe una única conexión afín  $\nabla$  en  $M$  satisfaciendo las condiciones:*

- a)  $\nabla$  es simétrica.
- b)  $\nabla$  es compatible con la métrica riemanniana.

*(Esta conexión es denominada conexión riemanniana).*

**Definición 3.5.11** Dado un sistema de coordenadas  $(U, \mathcal{X})$ , las funciones conocidas como *símbolos de Christoffel*  $\Gamma_{ij}^k : U \subset M \rightarrow \mathbb{R}$  que definen los coeficientes de conexión

$$\nabla_{X_i} X_j = \sum_{k=1}^n \Gamma_{ij}^k X_k,$$

están relacionados con la métrica de la siguiente manera :

$$\Gamma_{ij}^m = \frac{1}{2} \sum_k \left\{ \frac{\partial}{\partial x_i} g_{jk} + \frac{\partial}{\partial x_j} g_{ki} - \frac{\partial}{\partial x_k} g_{ij} \right\} g^{km}, \quad (3.9)$$

donde  $g_{ij} = \langle \frac{\partial}{\partial x_i}, \frac{\partial}{\partial x_j} \rangle$  son elementos de la matriz  $G(x)$  y  $g^{ij}$  los elementos de la inversa  $G(x)^{-1}$ .

## 3.6. Curvatura de una variedad Riemanniana

**Definición 3.6.1** Una curvatura  $K$  de una variedad Riemanniana  $M$  es una correspondencia que asocia a cada par  $X, Y \in \mathcal{H}$  una aplicación  $K(X, Y) : \mathcal{H} \rightarrow \mathcal{H}$  dada por :

$$K(X, Y)Z = \nabla_Y \nabla_X Z - \nabla_X \nabla_Y Z + \nabla_{[X, Y]} Z, \quad Z \in \mathcal{H},$$

donde  $\nabla$  es una conexión de Riemanniana de  $M$ .

**Observación 3.6.1** Si consideramos un sistema de coordenadas  $(U, \mathcal{X})$  en torno del punto  $p$  y  $\{X_i\}$ ,  $i = 1, 2, \dots, n$  es una base de  $T_p M$  obtenemos:

$$K(X_i, X_j)X_k = (\nabla_{X_j} \nabla_{X_i} - \nabla_{X_i} \nabla_{X_j})X_k$$

**Proposición 3.6.1** La curvatura  $K$  de una variedad Riemanniana cumple las siguientes propiedades:

i)  $K$  es bilineal en  $\mathcal{H} \times \mathcal{H}$ , esto es,

$$K(fX_1 + gX_2, Y_1) = fK(X_1, Y_1) + gK(X_2, Y_1),$$

$$K(X_1, fY_1 + gY_2) = fK(X_1, Y_1) + gK(X_1, Y_2),$$

donde  $f, g \in C^\infty(M)$ ,  $X_1, X_2, Y_1, Y_2 \in \mathcal{H}$ .

ii) Para todo  $X, Y \in \mathcal{H}$  el operador curvatura  $K(X, Y)$  es lineal, esto es,

$$K(X, Y)(Z + W) = K(X, Y)Z + K(X, Y)W,$$

$$K(X, Y)fZ = fK(X, Y)Z,$$

donde  $f \in C^\infty(M)$ ,  $Z, W \in \mathcal{H}$ .

iii) La curvatura es antisimétrica, esto es ,

$$K(X, Y) = -K(Y, X).$$

**Prueba:** Ver [17]

**Proposición 3.6.2** Sea  $(U, \mathcal{X})$  un sistema de coordenadas en torno de  $p \in M$  y  $\{X_i\}$  una base de  $T_pM$  en este sistema de coordenadas. Entonces:

$$K(X_i, X_j)X_k = \sum_{l=1}^n K_{ijk}^l X_l,$$

donde  $K_{ijk}^l$  son dadas por :

$$K_{ijk}^l = X_j \Gamma_{ik}^l - X_i \Gamma_{jk}^l + \sum_{s=1}^n \Gamma_{ik}^s \Gamma_{js}^l - \sum_{s=1}^n \Gamma_{jk}^s \Gamma_{is}^l.$$



**Prueba:** Ver [17]

**Observación 3.6.2** Si en las coordenadas  $(U, \mathcal{X})$  escribimos  $X = \sum_{i=1}^n u^i X_i$ ,  $Y = \sum_{j=1}^n v^j Y_j$ ,  $Z = \sum_{k=1}^n w^k Z_k$ , por la linealidad de  $K$  tenemos:

$$K(X, Y)Z = \sum_{i,j,k,l=1}^n K_{ijk}^l u^i v^j w^k X_l. \quad (3.10)$$

**Curvatura Seccional.** Íntimamente relacionado con el operador curvatura  $R$  está la curvatura seccional (o riemanniana) que definiremos a continuación.

**Definición 3.6.2** Dado un espacio vectorial  $V$ , indicaremos por  $|x \wedge y|$  a la expresión

$$\sqrt{|x|^2|y|^2 - \langle x, y \rangle^2},$$

que representa el área del paralelogramo bidimensional definido por un par de vectores  $x, y \in V$

**Definición 3.6.3** Sea  $\sigma \subset T_p M$  un subespacio bidimensional del espacio vectorial  $T_p M$  y sean  $x, y \in \sigma$ , dos vectores linealmente independientes. Entonces,

$$K(x, y) = \frac{(K(x, y)x, y)}{|x \wedge y|},$$

no depende de la elección de los vectores  $x$  e  $y$ .

**Definición 3.6.4** (Curvatura Seccional). Dado un punto  $p \in M$  y un subespacio bidimensional  $\sigma \subset T_p M$ . El número  $K(x, y) = K(\sigma)$ , donde  $\{x, y\}$  es una base de  $\sigma$ , es llamado Curvatura Seccional de  $M$ .

Si  $K(x, y) \leq 0$  para todo  $x, y \in \sigma$  entonces, la curvatura seccional de la variedad riemanniana es no positiva.

Si  $K(x, y) \geq 0$  para todo  $x, y \in \sigma$  entonces, la curvatura seccional de la variedad riemanniana es no negativa.

### Variedades completas

**Definición 3.6.5** Una variedad Riemanniana  $M$  es llamada (geodésicamente) completa si para todo  $p \in M$ , la aplicación exponencial,  $\exp_p$ , están definidas para todos  $v \in T_p M$ ; es decir, las geodésicas que parten de  $p$  están definidas para todos los valores del parámetro  $t \in \mathbb{R}$ .

**Definición 3.6.6** Dados dos puntos  $p$  y  $q$  en  $M$ , la distancia Riemanniana de  $p$  a  $q$  en la variedad, denotada por  $d(p, q)$ , es definida por

$$dist(p, q) = \inf_{\substack{\gamma: [a, b] \rightarrow M \\ p = \gamma(a), q = \gamma(b)}} \int_a^b \|\gamma'(t)\| dt \quad (3.11)$$

**Proposición 3.6.1** Con la distancia geodésica (3.11)  $M$  es un espacio métrico.

**Teorema 3.6.1** (Hopf-Rinow) Sea  $M$  una variedad Riemanniana y sea  $p \in M$ . Las siguientes afirmaciones son equivalentes:

- a)  $\exp_p$  está definida en todo  $T_p M$  para algún  $p \in M$
- b) Los limitados y cerrados de  $M$  son compactos.

c)  $M$  es completo como espacios métrico.

d)  $M$  es geodésicamente completa.

e) Para todo  $q \in M$  existe una geodésica uniendo  $p$  y  $q$  con  $d(p, q) = \inf_{\gamma} \int_a^b \|\gamma'(t)\| dt$ .

**Prueba:** Ver [18].

**Teorema 3.6.2** *En una variedad Riemanniana completa  $M$  de dimensión finita con curvatura seccional no-negativa, tenemos:*

$$l_3^2 \leq l_1^2 + l_2^2 - 2l_1l_2 \cos \theta, \quad (3.12)$$

donde  $l_i$  denota la longitud de  $\alpha_i$ ,  $i = 1, 2$ ,  $l_3 = d(\alpha_1(l_1), \alpha_2(l_2))$  y  $\theta = \angle(\alpha_1'(0), \alpha_2'(0))$ .

**Prueba:** Ver [18].

**Definición 3.6.7** *Sea  $M$  una variedad Riemanniana y  $f : M \rightarrow \mathbb{R}$ , se define campo vectorial  $\mathbf{grad}(f)$  (gradiente de la función  $f$ ) como el único campo vectorial que cumple:*

$$df_p(X(p)) = \langle \mathbf{grad}f(p), X(p) \rangle_x, \forall x \in \mathcal{H}. \quad (3.13)$$

**Observación 3.6.3** *Sea  $M \subset \mathbb{R}^n$  una variedad Riemanniana con la métrica definida por  $\langle v, w \rangle_x = v^T G(x)w$  donde  $G(x)$  es una matriz simétrica definida positiva. Se puede caracterizar el campo gradiente como:*

$$\mathbf{grad}f(q) = G^{-1}(q)f'(q), \quad (3.14)$$

donde  $G^{-1} = g^{ij}(q)$  es la matriz inversa de  $G(q)$  y  $f'(x) = (\frac{\partial f}{\partial x_1}, \dots, \frac{\partial f}{\partial x_n})$  es el vector de derivadas parciales de la función  $f \circ X$ . En efecto,

$$df_p(X(p)) = f'(q)^T v = f'(q)^T (G^{-1}(q))^T G(q)v = (G^{-1}(q)f'(q))^T G(q)v = \langle G^{-1}(q)f'(q), v \rangle_q.$$

**Definición 3.6.8** Sea  $M$  una variedad Riemanniana y  $f : M \rightarrow \mathbb{R}$  una función de clase  $C^k, k \geq 2$ . La Hessiana de  $f$ , denotada por  $H^f$  es definida como la derivada covariante del campo vectorial gradiente, esto es,

$$H^f = \frac{D}{dt}(\text{grad}f).$$

La Hessiana en el punto  $p$ , en la dirección de  $v \in T_p M$  es:

$$H_p^f(v, w) = \left\langle \frac{D_v}{dt}(\text{grad}f)(p), w \right\rangle_p. \quad (3.15)$$

### 3.7. Preliminares de optimización sobre Variedades

**Definición 3.7.1** Sea  $M$  una variedad Riemanniana y  $f : M \rightarrow \mathbb{R}$ , definimos por **Punto Crítico**, los puntos donde :

$$\langle \text{grad}f(x), v \rangle_x = 0, x \in M \text{ y } \forall v \in T_x(M). \quad (3.16)$$

**Definición 3.7.2** Sea  $f : M \rightarrow \mathbb{R}, x^* \in M$  y una  $\mathbb{U}$  vecindad abierta de  $x^*$ , entonces,  $x^*$  es un **Punto Mínimo Local de  $f$** , si se cumple que  $f(x^*) \leq f(x) \forall x \in \mathbb{U} \cap A$ ; si la desigualdad se cumple para todo  $x \in M$ , entonces  $x^*$  es un **Punto Mínimo Global de  $f$  en  $A$** .

**Definición 3.7.3** *Se dice que una función  $f : M \rightarrow \mathbb{R}$  es semicontinua inferiormente en  $\hat{x} \in M$ , si para toda sucesión  $\{x^k\}$  de  $M$  convergente a  $\hat{x}$  se tiene que:*

$$\liminf_{k \rightarrow \infty} f(x^k) \geq f(\hat{x}).$$

*Si  $f$  es semicontinua inferiormente para todo  $x \in M$ , entonces decimos que  $f$  es semicontinua inferiormente en  $M$ .*

**Teorema 3.7.1** (Weierstrass) *si  $f : M \rightarrow \mathbb{R}$  es semicontinua inferiormente y  $M$  es compacto, entonces existe un punto mínimo global de  $f$ .*

**Prueba:** Ver [13].

**Teorema 3.7.2** (Condición necesaria de primer orden). *Sea  $f : M \rightarrow \mathbb{R}$  de clase  $C^1$ . Si  $x^*$  es un punto de mínimo local, entonces  $\text{grad}f(x^*) = 0$ .*

**Demostración 3.7.1** *Sean  $v \in T_{x^*}M$  y  $\gamma : \mathbb{R} \rightarrow M$  una curva geodésica tal que  $\gamma(0) = x^*$  y  $\gamma'(0) = v$ . Definamos la aplicación  $h : \mathbb{R} \rightarrow \mathbb{R}$  tal que  $h(t) = f(\gamma(t))$ . Como  $x^*$  es punto de mínimo local de  $f$ , existe  $\varepsilon > 0$  tal que  $h(0) = f(x^*) \leq f(\gamma(t)) = h(t)$ , para todo  $t \in (-\varepsilon, \varepsilon)$ , por lo que en  $t = 0$  tenemos un punto de mínimo local de  $h$ . Por la condición necesaria de primer orden en  $\mathbb{R}$  se tiene  $h'(0) = \langle \text{grad}f(x^*), v \rangle = 0$ . En particular para  $v = \text{grad}f(x^*)$  tenemos que  $\text{grad}f(x^*) = 0$ .*

**Teorema 3.7.3** (Condición necesaria de segundo orden). *Sea  $f : M \rightarrow \mathbb{R}$  de clase  $C^2$ . Si  $x^*$  es un punto de mínimo local, entonces  $\langle v, H_{x^*}^f v \rangle \geq 0, \forall v \in T_{x^*}M$ .*

**Demostración 3.7.2** Sean  $v \in T_{x^*}M$  y  $\gamma : \mathbb{R} \rightarrow M$  una curva geodésica tal que  $\gamma(0) = x^*$  y  $\gamma'(0) = v$ . Definamos la aplicación  $h : \mathbb{R} \rightarrow \mathbb{R}$  tal que  $h(t) = f(\gamma(t))$ . De la demostración del teorema anterior, sabemos  $t = 0$  hay un punto de mínimo local de  $h$ , entonces por la condición necesaria de segundo orden en  $\mathbb{R}$  se tiene,  $h'(0) = 0$ , luego  $h''(0) \geq 0$ . Ahora bien,

$$\begin{aligned} h'(t) &= \langle \text{grad}f(\gamma(t)), \gamma'(t) \rangle \\ h''(t) &= \left\langle \frac{D}{dt} \text{grad}f(\gamma(t)), \gamma'(t) \right\rangle + \left\langle \text{grad}f(\gamma(t)), \frac{D}{dt} \gamma'(t) \right\rangle \\ &= \langle H_{\gamma(t)}^f \gamma'(t), \gamma'(t) \rangle \\ &= \langle H_v^f v, v \rangle \geq 0. \end{aligned}$$

**Teorema 3.7.4** (Condición suficiente de segundo orden). Sea  $f : M \rightarrow \mathbb{R}$  de clase  $C^2$ . Si  $x^* \in M$  cumple:

- $\text{grad}f(x^*) = 0$ ,
- $H_{x^*}^f$  definida positiva,

entonces,  $x^*$  es un punto de mínimo local estricto de  $f$ .

**Prueba:** Ver [13].

**Definición 3.7.4** Sea  $A \subset M$  donde  $M$  es una variedad Riemanniana completa. Se dice que  $A$  es un **Conjunto Totalmente Convexo** si contiene toda geodésica  $\gamma$  de  $M$  que conecte a  $p$  y  $q$ , con  $p, q \in A$ .

**Definición 3.7.5** Sea  $A \subset M$  donde  $M$  es una variedad Riemanniana completa. Diremos que  $A$  es convexo si para todo par de puntos  $p$  y  $q$  de  $A$  existe una geodésica minimal que une  $p$  y  $q$  contenido en  $A$ .

**Definición 3.7.6** Sea  $f : M \rightarrow \mathbb{R}$ ; diremos que  $f$  es convexa si y solo si para toda geodésica  $\gamma : \mathbb{R} \rightarrow M$ , la compuesta  $f \circ \gamma : \mathbb{R} \rightarrow \mathbb{R}$  es convexa como una función real, es decir,  $\forall a, b \in \mathbb{R}, \forall \lambda \in [0, 1]$  se cumple:

$$f(\gamma(1 - \lambda)a + \lambda b) \leq (1 - \lambda)f(\gamma(a)) + \lambda f(\gamma(b)). \quad (3.17)$$

**Teorema 3.7.5** Sea  $f : M \rightarrow \mathbb{R}$  y  $M$  una variedad Riemanniana. Entonces  $f$  es convexa si y sólo si para todo  $p \in M$  y para toda geodésica  $\gamma : [0, +\infty) \rightarrow M$  tal que  $\gamma(0) = p$  se cumple:

$$f(\gamma(t)) - f(p) \geq t \langle \text{grad}f(p), \gamma'(0) \rangle \quad \text{con } t \in [0, +\infty). \quad (3.18)$$

**Demostración 3.7.3** Supongamos que  $f$  es convexa, entonces para toda geodésica  $\gamma : \mathbb{R} \rightarrow M$ ;  $f \circ \gamma$  es convexa, en particular para  $h : [0, +\infty) \rightarrow \mathbb{R}$  dada por  $h(t) = f(\gamma(t))$  con  $t \in [0, +\infty)$ .

Como  $h$  es convexa como función real, entonces se cumple

$$h(t) \geq h(0) + th'(0). \quad (3.19)$$

Por otro lado

$$h'(t) = (f \circ \gamma)' = \langle \text{grad}f(\gamma(t)), \gamma'(t) \rangle,$$

evaluando  $h'$  en 0

$$\begin{aligned} h'(0) &= \langle \text{grad}f(\gamma(0)), \gamma'(0) \rangle \\ &= \langle \text{grad}f(p), \gamma'(0) \rangle \end{aligned}$$

sustituyendo la igualdad anterior en 3.19 tenemos:

$$h(t) \geq h(0) + t\langle \text{grad}f(p), \gamma'(0) \rangle$$

lo que implica que:

$$f(\gamma(t)) - f(p) \geq t\langle \text{grad}f(p), \gamma'(0) \rangle.$$

Recíprocamente, supongamos que se cumple que:

$$f(\gamma(t)) - f(p) \geq t\langle \text{grad}f(p), \gamma'(0) \rangle \text{ con } t \in [0, +\infty),$$

con  $\gamma : \mathbb{R} \rightarrow M$  la geodésica tal que  $\gamma(0) = p$ .

Sea  $h : [0, +\infty) \rightarrow \mathbb{R}$  dada por  $h(t) = f(\gamma(t))$  con  $t \in [0, +\infty)$ .

como  $\gamma(0) = p$  y  $h'(0) = \langle \text{grad}f(p), \gamma'(0) \rangle$ , ocurre que  $h(t) - h(0) \geq th'(0)$ ,

de esto se tiene que  $h = f \circ \gamma$  es convexa y por lo tanto  $f$  es convexa.

**Corolario 3.7.1** Si  $f : M \rightarrow \mathbb{R}$  es convexa entonces todos los puntos críticos de  $f$  son mínimos globales de  $f$ .

**Demostración 3.7.4** Sea  $p \in M$  un punto crítico de  $f$ , entonces  $\text{grad}f(p) = 0$ . Dado  $q \in M$  con  $q \neq p$ , existe por teorema de Hopf-Rinow's una geodésica  $\gamma : [a, b] \rightarrow M$  que



conecta a  $p$  con  $q$  donde  $\gamma(a) = p$  y  $\gamma(b) = q$ . Como  $f$  es convexa por teorema de condición de primer orden para convexidad se cumple:  $f(\gamma(b)) - f(\gamma(a)) \geq \langle \text{grad}f(p), \gamma'(0) \rangle$  y dado que  $q$  es punto crítico  $f(q) \geq f(p) \forall q$ . Por lo tanto, todo punto crítico de una función convexa es un mínimo global.

**Teorema 3.7.6** Sea  $M$  una variedad Riemanniana, una función  $f : M \rightarrow \mathbb{R}$  es convexa si y solo si el Hessiano de  $f$  ( $H^f$ ) es semi-definido positivo.

**Demostración 3.7.5** Supongamos que  $f$  es convexa, sean  $p \in M$  y  $v \in T_pM$  y  $\gamma$  la única geodésica tal que  $\gamma(0) = p$  y  $\gamma'(0) = v$ . Definamos  $h : \mathbb{R} \rightarrow \mathbb{R}$  dada por  $f \circ \gamma = h$  por lo que  $h$  es convexa. Así, por la teoría de convexidad,  $h''(t) \geq 0 \forall t$ , ahora por definición de derivada

$$h'(t) = \langle \text{grad}f(\gamma(t)), v \rangle_{\gamma(t)},$$

luego

$$h''(0) = \langle H_p^f v, v \rangle \geq 0.$$

Recíprocamente, supongamos  $\langle H_p^f v, v \rangle \geq 0$ . Definiendo  $h : \mathbb{R} \rightarrow \mathbb{R}$  dada por  $f \circ \gamma = h$ , se tiene que  $h''(t) \geq 0$  y  $h = f \circ \gamma : \mathbb{R} \rightarrow \mathbb{R}$  es convexa, finalmente por la definición de convexidad  $f$  es convexa.

**Definición 3.7.7** Sea  $M$  una variedad Riemanniana completa y  $f : M \rightarrow \mathbb{R}$  una función real.  $f$  es llamada cuasi-convexa en  $M$  si para todo  $x, y \in M$ ,  $t \in [0, 1]$ , se cumple:

$$f(\gamma(t)) \leq \max\{f(x), f(y)\},$$

para toda curva geodésica  $\gamma : [0, 1] \rightarrow M$ , tal que  $\gamma(0) = x$  y  $\gamma(1) = y$ .

**Teorema 3.7.7** Sea  $f : M \rightarrow \mathbb{R}$  una función diferenciable y cuasi-convexa en una variedad Riemanniana completa  $M$  y sean  $x, y \in M$ . Si  $f(x) \leq f(y)$  entonces:

$$\langle \text{grad}f(y), \gamma'(0) \rangle \leq 0$$

donde  $\gamma$  es la curva geodésica tal que  $\gamma(0) = y$  y  $\gamma(1) = x$ .

**Prueba:** Ver [13].

**Definición 3.7.8** Sea  $\Gamma \geq 0$  una constante y  $f : M \rightarrow \mathbb{R}$ , si para cada  $p, q$  en  $M$  y el segmento geodésico  $\gamma : [0, a] \rightarrow M$  que conecta a  $p$  y  $q$  con  $\gamma(0) = p$  y  $\gamma'(a) = q$  dada por:

$l(\gamma_{0,t}) =$  longitud del segmento  $[\gamma(0), \gamma(t)] = \int_0^t \|\gamma'(t)\| dt$  cumple lo siguiente

$$\|\text{grad}f(\gamma(t)) - P_\gamma \text{grad}f(p)\| \leq \Gamma \|l(\gamma_{0,t})\|,$$

se dice que  $f$  tiene gradiente  $\Gamma$ -lipschitziana.

**Definición 3.7.9** Sea  $f : M \rightarrow \mathbb{R}$ , una función da la variedad  $M$  a los reales, se denomina conjunto de subnivel  $a$ :  $M^a = \{x \in M : f(x) \leq a\}$  con  $a \in \mathbb{R}$

## 3.8. Funciones Convexas en Subvariedades de Riemann

### 3.8.1 Gradiente y Hessiano en Subvariedades

Sea  $(M, g)$  una variedad Riemanniana de dimensión  $n + p$  y  $N$  una subvariedad de dimensión  $n$  cuya métrica inducida es denotada también por  $g$ . Denotamos por  $\xi_1, \dots, \xi_p$  el campo de vectores local sobre  $N$  que suponemos normal. Si  $\nabla'$  es la conexión Riemanniana sobre  $M$  y  $\nabla$  es la conexión Riemanniana sobre  $N$ , entonces la fórmula de Gauss se sostiene

$$\nabla'_X Y = \nabla_X Y + \alpha(X, Y),$$

donde  $X$  y  $Y$  son campos de vectores sobre  $N$ , y

$$\alpha(X, Y) = \sum_{s=1}^p \Omega^s(X, Y) \xi_s$$

es la *segunda forma fundamental* sobre  $N$ .

Sea  $f : M \rightarrow \mathbb{R}$  una función de clase  $\mathcal{C}^2$ . La restricción de  $f$  sobre  $N$  será denotada por  $f_N$ . Observe que

$$\text{grad}_M f = \text{grad}_N f_N + \sum_{s=1}^p a^s \xi_s, \quad (3.20)$$

donde

$$a^s = g(\text{grad}_M f, \xi_s) = df(\xi_s).$$

En consecuencia, si  $x_o \in M$  es un punto crítico de  $f$  y si  $x_o \in N$ , entonces  $x_o$  es un punto crítico de  $f_N$ . Más,  $x_o \in N$  es un punto crítico de  $f_N$  si y sólo si

$$\text{grad}_M f(x_o) = \sum_{s=1}^p a^s(x_o) \xi_s(x_o) \in T_{x_o} N^\perp.$$

**Teorema 3.8.1** *Si  $f : M \rightarrow \mathbb{R}$  es una función de clase  $\mathcal{C}^2$  y  $f_N$  es su restricción sobre  $N$ , entonces*

$$\text{Hess}_M f = \text{Hess}_N f_N - \sum_{s=1}^p a^s \Omega^s.$$

**Prueba:** Ver [7].

Ya que una subvariedad  $N$  puede ser accidentalmente un subconjunto convexo totalmente de  $M$ , aceptamos aquí que una función  $f$  es convexa sobre  $N$  si  $\text{Hess} f|_N \geq 0$ .

**Corolario 3.8.1** *1. Si  $f : M \rightarrow \mathbb{R}$  es convexa y  $\text{grad}_M f$  es tangente a  $N$ , entonces  $f_N$  es convexa.*

*2. Si  $f : M \rightarrow \mathbb{R}$  es convexa y  $N$  es una subvariedad geodésica totalmente, entonces  $f_N$  es convexa.*

*3. Si  $f : M \rightarrow \mathbb{R}$  es convexa y  $df(\alpha)$  es semidefinida positiva, entonces  $f_N$  es convexa.*

**Prueba:** Ver [7].

### Ejemplos.

1. La función  $f : \mathbb{R}^3 \rightarrow \mathbb{R}$ ,  $f(x, y, z) = e^z$  es convexa y tendremos que  $\text{grad} f = (0, 0, e^z)$ .

Sea  $N$  el cilindro circular  $x^2 + y^2 = 1$  en  $\mathbb{R}^3$ . El campo de vectores unitarios sobre  $N$  es  $\xi = (x, y, 0)$ . Entonces  $(\xi, \text{grad} f) = 0$  y en consecuencia  $f_N$  es convexa.

2. Sea  $f : \mathbb{R}^3 \rightarrow \mathbb{R}$ ,  $f(x, y, z) = e^z$  y  $N : z = \frac{x^2 + y^2}{2}$ .

Como

$$df = e^z dz, \quad \xi = \frac{(-x, -y, 1)}{\sqrt{1 + x^2 + y^2}}, \quad \Omega = \frac{dx^2 + dy^2}{\sqrt{1 + x^2 + y^2}}.$$

Se sigue que

$$df(\xi) = \frac{e^z}{\sqrt{1 + x^2 + y^2}}$$

y

$$Hess_N f_N = Hess_{\mathbb{R}^3} f + df(\xi) \Omega = e^z \left( dz^2 + \frac{dx^2 + dy^2}{1 + x^2 + y^2} \right) \Big|_N.$$

En consecuencia  $f_N$  es convexa.

3. La función

$$f : \mathbb{R}^3 \rightarrow \mathbb{R}, \quad f(x, y, z) = \frac{x^2 + y^2}{2}$$

es convexa. Considerando el paraboloido de rotación

$$N : z = \frac{x^2 + y^2}{2},$$

orientado por

$$\xi = \frac{(-x, -y, 1)}{\sqrt{1 + x^2 + y^2}},$$

par el cual la segunda forma fundamental es

$$\Omega = \frac{dx^2 + dy^2}{\sqrt{1 + x^2 + y^2}}.$$

Se sigue que

$$Hess_N f_N = Hess_{\mathbb{R}^3} f + df(\xi) \Omega = \frac{dx^2 + dy^2}{1 + x^2 + y^2} \Big|_N,$$

y en consecuencia  $f_N$  es convexa.

**Generalización.** Si  $f : \mathbb{R}^2 \rightarrow \mathbb{R}$  es convexa, entonces  $\tilde{f} : \mathbb{R}^3 \rightarrow \mathbb{R}$ ,  $\tilde{f}(x, y, z) = f(x, y)$  es convexa. La restricción de  $\tilde{f}$  al gráfico de  $f$  es también convexa.

---

## Capítulo 4

### Método De Descenso Clásico Sobre Variedades.

---

En la optimización clásica en espacios vectoriales, los métodos de búsqueda lineal se utilizan ampliamente. Se basan en la actualización de la iteración por la elección de una dirección de búsqueda y luego se adiciona un múltiplo de ésta a la iteración anterior. Es decir, la búsqueda lineal en  $\mathbb{R}^n$  se basa en la fórmula iterativa

$$x_{k+1} = x_k + \lambda_k d_k, \quad (4.1)$$

donde  $d_k \in \mathbb{R}^n$  es la *dirección de búsqueda* y  $\lambda_k \in \mathbb{R}$  es la *longitud del paso*. La adición de un múltiplo a la dirección de búsqueda requiere, obviamente, la estructura de un espacio vectorial y, en general, esto no es posible en variedades. El objetivo de este capítulo es desarrollar una teoría similar para problemas de optimización planteados en variedades no lineales. La extensión natural a variedades es seguir la dirección de búsqueda a lo largo de una curva; más específicamente, la generalización de (4.1) a una variedad  $M$  consiste en seleccionar  $d_k$  como un vector tangente a  $M$  en  $x_k$ , y realizar la búsqueda a lo largo de una curva en  $M$  cuyo vector tangente en  $t = 0$  sea  $d_k$ . La elección de la curva se basa en el concepto de *retracción*, que veremos en la siguiente sección. Cabe destacar, que la

escogencia de una *retracción* eficiente computacionalmente es una decisión importante en el diseño de algoritmos numéricos de alto rendimiento en variedades no lineales.

## 4.1. Retracciones

Conceptualmente, el método más sencillo para la optimización de una función diferenciable es trasladar un punto de prueba  $x(t)$  continuamente en la dirección de máxima pendiente,  $-\text{grad } f(x)$ , sobre el conjunto de restricciones hasta que se llegue a un punto en que la pendiente se anule. Los puntos  $x$  donde  $\text{grad } f(x) = 0$  se llaman puntos estacionarios o puntos críticos de  $f$ . Una implementación numérica de la aproximación de descenso del gradiente de forma continua, requiere de la construcción de una curva  $\gamma$  de tal manera que  $\gamma'(t) = -\text{grad } f(\gamma(t))$  para todo  $t$ . Excepto en circunstancias muy especiales, la construcción de una curva utilizando métodos numéricos es poco práctico. La analogía numérica más cercana es la clase de métodos de optimización que utilizan procedimientos de *búsqueda lineal*, es decir, algoritmos iterativos que, dado un punto  $x$ , calcula una dirección de descenso  $\lambda = -\text{grad } f(x)$  (o alguna aproximación del gradiente) y se mueve en la dirección de  $\lambda$  hasta que se encuentre una disminución “razonable” de  $f$ . En  $\mathbb{R}^n$ , el concepto de movimiento en la dirección de un vector es sencillo. En una variedad, la noción de moverse en la dirección de un vector tangente, durante su estadía en la misma, es generalizada por la noción de la *función de retracción*.

Podemos ver a una retracción  $R$  por  $x$ , denotada por  $R_x$  (ver [8]), como una función de  $T_x M$  a  $M$  con una condición rígida que preserva el gradiente por  $x$ ; ver figura 4.1.



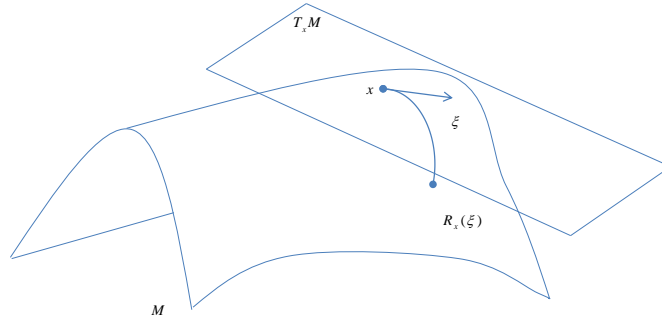


Figura 4.1: Retracción.

**Definición 4.1.1** Una retracción sobre una variedad  $M$  es una función  $R$  del fibrado tangente  $TM$  sobre  $M$  con las siguientes propiedades. Sea  $R_x$  la restricción de  $R$  a  $T_xM$ , así

1.  $R$  es continuamente diferenciable.
2.  $R_x(0_x) = x$ , donde  $0_x$  denota el elemento cero de  $T_xM$ .
3. Con la identificación canónica  $T_{0_x}T_xM \simeq T_xM$ ,  $R_x$  satisface

$$DR_x(0_x) = id_{T_xM}, \quad (4.2)$$

donde  $id_{T_xM}$  denota la función identidad sobre  $T_xM$ .

Por lo general, suponemos que el dominio de  $R$  es todo el fibrado tangente  $TM$ . En cuanto a la condición (4.2), nótese que, ya que  $R_x$  es una función de  $T_xM$  a  $M$  que envía a  $0_x$  a  $x$ , se sigue que  $DR_x(0_x)$  es una aplicación de  $T_{0_x}(T_xM)$  a  $T_xM$ . Ya que  $T_xM$  es un espacio vectorial, existe la identificación natural  $T_{0_x}(T_xM) \simeq T_xM$ . De aquí, por el teorema de la función inversa,  $R_x$  es un difeomorfismo local por  $0_x$ , que además, no solamente es  $C^1$  sino que también es biyectiva con inversa diferenciable en una vecindad  $V$  de  $0_x$  en  $T_xM$ . Aquí nos referiremos a la condición  $DR_x(0_x) = id_{T_xM}$  como la *condición de rigidez local*. Equivalentemente, para cada vector tangente  $\xi$  en  $T_xM$ , la curva  $\gamma_\xi : t \rightarrow R_x(t\xi)$  satisface  $\dot{\gamma}_\xi(0) = \xi$ . El moverse a lo largo de esta curva  $\gamma_\xi$  es pensado como ir en la dirección de  $\xi$  estando señado a la variedad  $M$ .

Además de convertir los elementos de  $T_xM$  en puntos de  $M$ , un segundo propósito importante de una retracción  $R_x$  es transformar la función costo definida en una vecindad de  $x \in M$  en una función costo sobre el espacio vectorial  $T_xM$ . Esto es, dada una función  $f$  a valores reales sobre una variedad  $M$  dotada con una retracción  $R$ , colocamos  $\hat{f} = f \circ R$  el cual denota el *pullback* de  $f$  a través de  $R$ . Así, para  $x \in M$  se tiene

$$\hat{f}_x = f \circ R_x, \quad (4.3)$$

que es la restricción de  $\hat{f}$  a  $T_xM$ . Nótese que  $\hat{f}_x$  es una función a valores reales sobre un espacio vectorial. También, por la condición de rigidez (4.2) tendremos (con la identificación canónica  $T_x\mathcal{E} \simeq \mathcal{E}$ ,  $\mathcal{E}$  espacio vectorial) que  $D\hat{f}_x(0_x) = Df(x)$ . Ahora, si  $M$  esta dotada con una métrica de Riemannian (y  $T_xM$  con un producto interno), tendremos que

$$grad \hat{f}_x(0_x) = grad f(x). \quad (4.4)$$

Posteriormente, retomaremos este segundo propósito de las retracciones, para conocer mejor su uso en combinación con el transporte paralelo de vectores.

Cada variedad que admite una métrica Riemannian también admite una retracción definida por la *aplicación exponencial Riemanniana* (lo que veremos más adelante). Más, cualquier otra retracción se puede pensar como una aproximación de la aplicación exponencial, aunque su dominio no necesariamente es todo  $TM$ ; en tal caso  $M$  tiene que ser *completa* (ver 3.6.5). Vamos a dar un ejemplo de retracción utilizando la aplicación exponencial. Antes, vamos exponer algunos detalles teóricos que hacen falta.

Recordemos del capítulo 3 que, una *geodésica* ( ver 3.5.5) sobre una variedad  $M$  dotada con una conexión afín  $\nabla$  (ver 3.5.3), es una curva con aceleración nula:

$$\frac{D^2}{dt^2}\gamma(t) = 0, \quad (4.5)$$

para todo  $t$  en el dominio de  $\gamma$ .

Cualquier conexión está ligada con la noción de transporte paralelo (ver 3.5.2). Dada una curva suave  $\gamma : [0, 1] \rightarrow M$ , el problema de valor inicial

$$\nabla_{\dot{\gamma}(t)}v = 0, v(0) = v_0, \quad (4.6)$$

define una manera de transportar el vector  $v_0 \in T_{\gamma(0)}M$  hacia el vector  $v(t) \in T_{\gamma(t)}M$ . Para la conexión Levi-Cevita considerada aquí, esto implica que el producto interior  $g_{\gamma(t)}(v(t), u(t))$  es constante para cualquier par de vectores  $v(t)$  y  $u(t)$ , transportados paralelamente a lo largo de  $\gamma$ . Para  $x = \gamma(0)$ ,  $y = \gamma(t)$ , denotaremos el transporte paralelo de  $v_0 \in T_xM$  a  $v(t) \in T_yM$  por  $v(t) = T_{x,y}^P v_0$  (si existieran más de un  $t$  con  $y = \gamma(t)$ , la

correcta interpretación llegará a ser clara en el contexto). Como se detalla más adelante  $T_{x,y}^{P_\gamma}$ , podrá ser interpretado como la derivada de la aplicación  $P_\gamma : T_x M \rightarrow M$ . Como consideraremos  $M$  completa, entonces por el teorema de Hopf-Rinow (ver 3.6.1 ) cualquier par de puntos  $x, y \in M$  se pueden conectar por una geodésica  $\gamma : [0, 1] \rightarrow M$  con  $x = \gamma(0)$ ,  $y = \gamma(1)$ , donde la longitud de esta curva es

$$L[\gamma] = \int_0^1 \sqrt{g_{\gamma(t)}(\dot{\gamma}(t), \dot{\gamma}(t))}.$$

Y la métrica inducida por esta longitud en  $M$  es

$$dist(x, y) = \inf_{\substack{\gamma: [0,1] \rightarrow M \\ x=\gamma(0), y=\gamma(1)}} L[\gamma]$$

Ya que la geodésica  $\gamma$  que conecta a  $x$  con  $y$  es única, entonces  $\dot{\gamma}(0) \in T_x M$  denota la aplicación logaritmo de  $y$  con respecto a  $x$ ,  $log_x y$ . Esta satisface  $dist(x, y) = \|log_x y\|_x$ . La geodésica puede calcularse vía la aceleración nula,  $\nabla_{\dot{\gamma}(t)} \dot{\gamma}(t) = 0$ ; esto es, resolviendo la ecuación (3.7). La aplicación exponencial  $exp_x : T_x M \rightarrow M$ ,  $v \rightarrow exp_x v$ , es definida como  $exp_x v = \gamma(1)$  donde  $\gamma$  es solución de la ecuación ordinaria dada (3.7) con condiciones iniciales  $\gamma_v(0) = x, \dot{\gamma}_v(0) = v$ . Además, por la propiedad de homogeneidad de  $\gamma$ , tenemos la siguiente igualdad  $\gamma_{kv}(t) = \gamma_v(kt)$ , de la que resulta

$$exp_x(tv) = \gamma_{kv}(1) = \gamma_v(t), \tag{4.7}$$

con lo que tenemos otra posibilidad de describir una geodésica. Obviamente se tiene

$$exp_x(0) = \gamma_v(0) = x$$

**Proposición 4.1.1** *Sea  $M$  una variedad completa con una conexión  $\nabla$  y  $p \in M$ . Entonces*

$$D \exp_x(0_x) = id_{T_x M}. \quad (4.8)$$

*Y en particular  $\exp_x$  es un difeomorfismo local en  $0$ .*

**Demostración:**

$$D \exp_x(0_x)(v) = \left. \frac{d}{dt} \right|_{t=0} \exp_x(tv) = \left. \frac{d}{dt} \right|_{t=0} \gamma_{tv}(1) = \left. \frac{d}{dt} \right|_{t=0} \gamma_v(t) = v, \quad (4.9)$$

para todo  $v \in T_x M$ , de donde se tiene la igualdad (4.8). Ahora, por el teorema de la función inversa,  $\exp_x$  es un difeomorfismo local en  $0$ .  $\diamond$

De la proposición (4.1.1), se tiene que la inversa de  $\exp_x$  sobre esta vecindad es  $\log_x$ ; con lo que la  $\exp_x$  es una parametrización (local) de  $M$  via  $T_x M$ . Por tanto, existen vecindades  $\widehat{U}$  y  $U$  de  $0_x \in T_x M$  y  $x \in M$ , respectivamente. Si  $\widehat{U}$  es *estrellada* (es decir,  $v \in \widehat{U}$  implica que  $tv \in \widehat{U}$  para  $0 \leq t \leq 1$ ), entonces  $U$  se llama *vecindad normal convexa de  $x$* . De esto se tiene el siguiente resultado.

**Proposición 4.1.2** *Sea  $M$  una variedad completa con conexión afín  $\nabla$ . La aplicación exponencial sobre  $M$  inducida por  $\nabla$  es una retracción, llamada la retracción exponencial.*

**Prueba:** Ver [14].

Bien, a continuación el ejemplo prometido de una retracción via la  $\exp_x$ : considere la geodésica  $\gamma : \mathbb{R} \rightarrow M$  con  $\gamma(0) = x$ , parametrizada por longitud de arco, defina la retracción  $P_\gamma : T_x M \rightarrow M$  como

$$P_\gamma(v) = \exp_{p_\gamma(v)} \left( T_{x, p_\gamma(v)}^{P_\gamma} [v - \pi_\gamma(v)] \right), \quad (4.10)$$

donde  $p_\gamma(v) = \exp_x(\pi_\gamma(v))$  y  $\pi_\gamma$  denota la proyección ortogonal sobre el *generador* de  $\dot{\gamma}(0)$ . Esta retracción, corresponde a moverse a lo largo de la línea geodésica  $\gamma$  de acuerdo con la componente de  $v$  paralela a  $\dot{\gamma}(0)$ ; luego, se sigue una nueva geodésica en la dirección del transporte paralelo de la componente de  $v$  ortogonal a  $\dot{\gamma}(0)$  (ver figura 4.2).

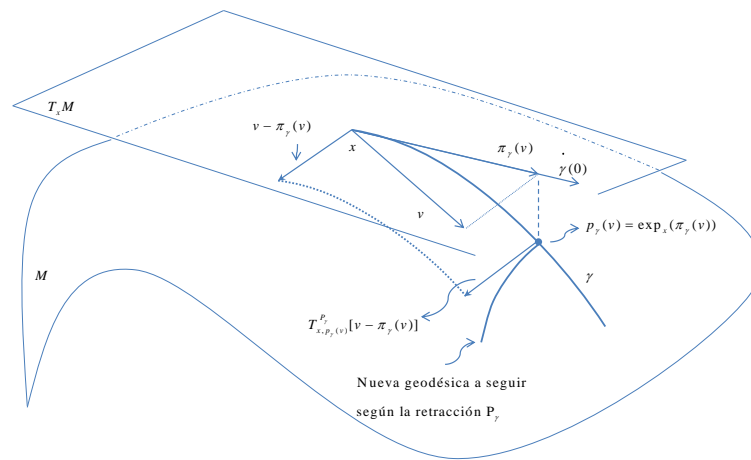


Figura 4.2: Retracción vía la exponencial.

Más adelante tendremos que considerar el transporte paralelo de un vector del espacio tangente  $T_x M$  a otro de  $T_y M$ ; es decir, consideraremos el isomorfismo  $T_{x,y} : T_x M \rightarrow T_y M$ . Pero estamos interesados en el operador  $T_{x,y}^{R_x}$  que representa la derivada  $D R_x(v)$  de  $R_x$  por  $v \in T_x M$  con  $R_x(v) = y$ . En ese sentido el transporte paralelo  $T_{x,y}^{P_\gamma}$  a lo largo de una geodésica  $\gamma$  que conecta a  $x$  con  $y$  pertenece a la retracción  $P_\gamma$ . Por ejemplo, de nuevo por la variación de la aplicación exponencial, podemos definir el transporte vectorial evaluado

por el representativo de  $y$  en  $T_x M$  como,

$$T_{x,y}^{exp_x} = D \exp_x(\log_x y), \quad (4.11)$$

que lleva a  $v_0 \in T_x M$  sobre el vector tangente  $\dot{\gamma}(0)$  de  $\gamma : t \rightarrow \exp_x(\log_x y + t v_0)$ . También se pueden considerar para el transporte vectorial  $T_{x,y}$ , los adjuntos  $(T_{y,x}^{P_\gamma})^*$ ,  $(T_{y,x}^{exp_y})^*$  (definidos por  $g_y(v, T_{y,x}^* w) = g_x(T_{y,x} v, w) \forall v \in T_y M, w \in T_x M$ ) o las inversas  $(T_{y,x}^{P_\gamma})^{-1}$ ,  $(T_{y,x}^{exp_y})^{-1}$ . Nótese que  $T_{x,y}^{P_\gamma}$  es una isometría con  $T_{x,y}^{P_\gamma} = (T_{y,x}^{P_\gamma})^* = (T_{y,x}^{P_\gamma})$ , donde  $\gamma$  es la geodésica que conecta  $x$  con  $y$  y  $\tilde{\gamma}(\cdot) = \gamma(-\cdot)$ . Además,  $\log_x y$  se puede transportar sobre el mismo vector  $T_{x,y} \log_x y = \dot{\gamma}(1)$  por  $T_{x,y}^{P_\gamma}$ , o por  $T_{x,y}^{exp_x}$ , o por sus adjuntos o por sus inversas.

Bien, volviendo al segundo propósito de las retracciones, en detalle sería: Dada una función suave  $f : M \rightarrow \mathbb{R}$ , definimos su derivada  $Df(x)$  en  $x \in M$  como un elemento del espacio dual de  $T_x M$  vía  $Df(x)v = v f$ . El teorema de representación de Riez (ver [19]), implica la existencia de un  $\nabla f(x) \in T_x M$  tal que  $Df(x)v = g_x(\nabla f(x), v)$  para todo  $v \in T_x M$ , que denota el *gradiente* de  $f$  en  $x$ . Luego, sobre  $T_x M$ , se define

$$f_{R_x} = f \circ R_x, \quad (4.12)$$

para la retracción  $R_x$ . Ya que  $DR_x(0) = id$  tendremos que  $Df_{R_x}(0) = Df_x$ . Además,

$f_{R_x} = f_{R_y} \circ R_y^{-1} \circ R_x$ , donde  $R_y^{-1}$  existe y así para  $y = R_x(v)$ , tenemos que

$$Df_{R_x}(v) = Df_{R_y}(0) DR_x(v) = Df(y) T_{x,y}^{R_x}, \quad \nabla f_{R_x}(v) = (T_{x,y}^{R_x})^* \nabla f(y).$$

## 4.2. Búsqueda Lineal: minimización sobre Variedades

Como se dijo al comienzo de este capítulo, los métodos de búsqueda lineal, están basados en la actualización de la iteración, por la elección de una dirección de búsqueda, hay que agregar un múltiplo (el tamaño o longitud del paso) a ésta dirección, lo que obviamente, requiere de la estructura de un espacio vectorial y, en general, esto no es posible en variedades. La extensión natural a variedades es seguir la dirección de búsqueda a lo largo de una curva. Entonces, vamos a considerar el algoritmo iterativo de la siguiente forma genérica.

### Algoritmo 1

**Entrar:**  $f : M \rightarrow \mathbb{R}$ , iterado inicial  $x_0 \in M$ ,  $k = 0$ .

**repetir**

- escoja una dirección de descenso  $d_k \in T_{x_k} M$
- escoja una retracción  $R_{x_k} \in T_{x_k} M \rightarrow M$
- escoja una longitud de paso  $\lambda_k \in \mathbb{R}$
- colocar  $x_{k+1} = R_{x_k}(\lambda_k d_k)$
- $k \leftarrow k + 1$

**hasta que**  $x_{k+1}$  minimice suficientemente a  $f$ .  $\diamond$

Como  $d_k$  es dirección de descenso entonces  $Df(x_k) d_k = \langle (\text{grad } f)(x_k), d_k \rangle < 0$ . Esta propiedad asegura que la función objetivo o costo,  $f$ , en verdad decrece en la dirección de



búsqueda.

Para la escogencia de la longitud del paso  $\lambda_k$ , existen varias posibilidades. En general, la escogencia de  $d_k$  tiene que satisfacer ciertos requerimiento de calidad, en nuestro caso  $d_k = -Df(x_k)$ . Para el objetivo de lo que estamos estudiando en este trabajo, nos concentraremos en la *condición de Wolfe*, es decir, para una dirección de descenso dada  $d \in T_x M$ , la escogencia de la longitud del paso  $\lambda$  tiene que satisfacer

$$f(R_x(\lambda d)) \leq f(x) + c_1 \lambda Df(x) d, \quad (4.13a)$$

$$Df(R_x(\lambda d)) T_{x, R_x(\lambda d)}^{R_x} d \geq c_2 Df(x) d, \quad (4.13b)$$

donde  $0 < c_1 < c_2 < 1$ . Nótese que ambas condiciones se pueden reescribir como

$$f_{R_x}(\lambda d) \leq f_{R_x}(0) + c_1 \lambda Df_{R_x}(0) d, \quad (4.14a)$$

$$Df_{R_x}(\lambda d) d \geq c_2 Df_{R_x}(0) d, \quad (4.14b)$$

que es la condición de Wolfe clásica para minimizar la función  $f_{R_x}$  del espacio vectorial  $T_x M$  a  $\mathbb{R}$ . Si la segunda condición (4.13b) es reemplazada por

$$|Df(R_x(\lambda d)) T_{x, R_x(\lambda d)}^{R_x} d| \leq -c_2 Df(x) d, \quad (4.15)$$

obtenemos la condición fuerte de Wolfe (generalmente se utiliza esta condición para un análisis último en un esquema cuasi-Newton). Ahora, en muchos algoritmos de optimización basta que la primera condición (4.13a) sea satisfecha, obviamente es la condición de Armijo; que en este caso es la generalizada y que trabajaremos mas adelante utilizando la

aplicación exponencial como retracción.

Similar al lema 2.2.1 para el caso  $\mathbb{R}^n$ , se tiene que, dada una dirección de descenso  $d$ , una longitud de paso factible  $\lambda > 0$  siempre se puede encontrar.

**Proposición 4.2.1** (*Longitud de paso factible*). *Sea  $x \in M$ ,  $d \in T_x M$  una dirección de descenso y  $f_{R_x} : \text{span}\{d\} \rightarrow \mathbb{R}$  diferenciable continuamente. Entonces existe  $\lambda > 0$  que satisface las condiciones (4.13) y (4.15).*

### Prueba

Al reescribir las condiciones de Wolfe como en (4.14) (4.15, respectivamente), el argumento clásico que se tiene para estas condiciones en espacios vectoriales se puede aplicar.

Bien, para ver los detalles, consideraremos los siguientes ingredientes:

Sean  $x_0 \in M$  (completa);  $\varepsilon > 0$ , los conjuntos  $B(x_0, \varepsilon) = \{x \in M : d(x, x_0) < \varepsilon\}$  y  $\bar{B}(x_0, \varepsilon) = \{x \in M : d(x, x_0) \leq \varepsilon\}$ , donde  $d$  es la distancia Riemanniana sobre  $M$ .  $\mathcal{D}_x$  es el conjunto de vectores  $v \in T_x M$  tal que la máxima geodésica  $\gamma_v(t)$  está definida en  $[0, 1]$ . El conjunto  $\mathfrak{X}(M)$  denotará los campos de vectores sobre  $M$  y  $\mathfrak{F}(M)$  el de funciones suaves sobre  $M$ .

Luego, para  $x \in \bar{B}(x_0, \varepsilon)$ , consideraremos,  $d_x \in \mathcal{D}_x$ ,  $\gamma_{d_x}(\lambda) = R_x(\lambda d_x)$  la geodésica que parte de  $x$  en dirección de  $d_x$  y por  $\xi_\lambda = \gamma'_{d_x}(\lambda)$  el campo vectorial de  $\mathfrak{X}(M)$ , tangente a lo largo de  $\gamma_{d_x}(\lambda)$  en donde  $\xi_0 = d_x$ . Así, colocando  $\omega(\lambda) = f(R_x(\lambda d_x)) - f(x)$  y como  $\frac{DR_x(\lambda d_x)}{\lambda} = \xi_\lambda$  resulta, por el desarrollo de Taylor de primer orden de  $\omega$  en 0 que

$$f(R_x(\lambda d_x)) - f(x) = \omega(\lambda) - \omega(0) = \lambda \omega'(0) + o(\lambda)$$

$$\begin{aligned}
&= \lambda \langle \nabla f(x), d_x \rangle + o(\lambda) = \lambda \langle \text{grad } f(x), d_x \rangle + o(\lambda) \\
&= \lambda c \langle \text{grad } f(x), d_x \rangle + \lambda [(1-c) \langle \text{grad } f(x), d_x \rangle + \frac{o(\lambda)}{\lambda}],
\end{aligned} \tag{4.16}$$

y ya que  $\lim_{x \rightarrow +\infty} \frac{o(\lambda)}{\lambda} = 0$ ,  $c < 1$  y  $d_k$  es dirección de descenso para la cual se cumple que  $Df(x_k) d_k = \langle (\text{grad } f)(x_k), d_k \rangle < 0$  en el **Algoritmo 1**, entonces, la última rama de (4.16) es negativa para  $\lambda$  suficientemente pequeño. En consecuencia, tenemos para  $d_x$  que  $f(R_x(\lambda d_x)) \leq f(x) + \lambda c \langle \text{grad } f(x), d_x \rangle$  con  $\lambda$  suficientemente pequeño. Es decir, la condición (4.13a) se cumple.

Por otra parte, para obtener la segunda condición (4.13b), observe que,  $f_{R_x}$  es acotada en  $\mathcal{D}_x$ , puesto que  $f$  lo es en el compacto  $\bar{B}(x_0, \varepsilon)$  y  $R_x$  en  $\mathcal{D}_x$ . Esto es, ya que  $d_x$  es dirección de descenso, se tiene que  $f$  es acota por debajo en  $\{x \in \bar{B}(x_0, \varepsilon) : f(R_x(\lambda d_x)) < f(x)\}$ , con lo que  $f_{R_x}$  es acotada por abajo en el conjunto

$$\{\xi_\lambda \in \mathcal{D}_x : f_{R_x}(\xi_\lambda) < f_{R_x}(0_x)\}. \tag{4.17}$$

Ahora, para cualquier  $\lambda$  no negativo definimos

$$\beta(\lambda) = f_{R_x}(\lambda d_x) + c_1 \lambda^2 \langle \text{grad } f(x), d_x \rangle \tag{4.18}$$

Entonces tenemos que

$$\lim_{\lambda \rightarrow 0^+} \frac{\beta(\lambda) - \beta(0)}{\lambda} = \lim_{\lambda \rightarrow 0^+} \frac{f_{R_x}(\lambda d_x) - f_{R_x}(0_x)}{\lambda} = \langle \text{grad } f(x), d_x \rangle < 0.$$

Luego, existe un  $\lambda' > 0$ , tal que para  $\lambda \in (0, \lambda']$  se tiene que

$$\frac{\beta(\lambda) - \beta(0)}{\lambda} \leq 0. \tag{4.19}$$

Por (4.17) tenemos que

$$\lim_{\lambda \rightarrow +\infty} \frac{\beta(\lambda) - \beta(0)}{\lambda} = +\infty,$$

puesto que este limite es

$$\lim_{\lambda \rightarrow +\infty} \left[ \frac{f_{R_x}(\lambda d_x) - f_{R_x}(0_x)}{\lambda} + c_1 \lambda \langle \text{grad } f(x), d_x \rangle \right],$$

donde el primer sumando es acotado y el segundo no.

Sea

$$\widehat{\lambda}' = \inf \left\{ \lambda > 0 : \frac{\beta(\lambda) - \beta(0)}{\lambda} = 0 \right\}.$$

Por el teorema del valor intermedio y (4.19), tenemos que  $\widehat{\lambda}'$  satisface

$$\frac{\beta(\widehat{\lambda}') - \beta(0)}{\widehat{\lambda}'} = 0. \quad (4.20)$$

Así, para cada  $\lambda \in (0, \widehat{\lambda}']$ , se tiene que

$$\frac{\beta(\lambda) - \beta(0)}{\lambda} \leq 0. \quad (4.21)$$

Aplicando el teorema del valor medio y (4.20) tenemos que existe  $\theta' \in [0, 1]$  tal que

$$\beta'(\theta' \widehat{\lambda}') = 0.$$

Por tanto,

$$Df_{R_x}(\theta' \widehat{\lambda}' d_x) d_x + 2 c_1 \theta' \widehat{\lambda}' \langle \text{grad } f(x), d_x \rangle = 0.$$

Es decir, si  $0 < c_1 < c_2 < 1$  tenemos que

$$Df_{R_x}(\theta' \widehat{\lambda}' d_x) d_x = -2 c_1 \theta' \widehat{\lambda}' \langle \text{grad } f(x), d_x \rangle \geq -2 c_2 \theta' \widehat{\lambda}' \langle \text{grad } f(x), d_x \rangle. \quad (4.22)$$

Al sustituir  $\lambda = \theta' \widehat{\lambda}' < \widehat{\lambda}'$  y cancelando términos semejantes en (4.22), tenemos la condición de Wolfe (4.13b); más, de (4.21) y (4.22) concluimos que  $\lambda = \theta' \widehat{\lambda}'$  es la longitud del paso que se requiere.  $\diamond$

### 4.3. Convergencia para el método de búsqueda lineal en variedades

El estudio de convergencia en variedades es una generalización directa del caso  $\mathbb{R}^n$ . Así que, una sucesión infinita  $\{x_k\}$  de puntos de una variedad  $M$  se dice que es *convergente* si existe una carta  $(\mathcal{U}, \varphi)$  de  $M$ , un punto  $x^* \in \mathcal{U}$ , y un conjunto  $K > 0$  tal que  $x_k \in \mathcal{U}$  para todo  $k \in K$  y la sucesión  $\{\varphi(x_k)\}$  converge a  $\varphi(x^*)$ . El punto  $\varphi^{-1}(\lim_{k \rightarrow \infty} \varphi(x_k))$  se llama el *limite* de la sucesión convergente  $\{x_k\}$ . Cada sucesión convergente de una variedad (Hausdorff) tiene uno y sólo un punto limite ( ver [16]).

La convergencia de los algoritmos de búsqueda lineal, naturalmente dependen, tanto de una buena escogencia de la longitud del paso como de una buena dirección de búsqueda. Con este objetivo en mente, vamos a utilizar el teorema de Zoutendijk dado en  $\mathbb{R}^n$  ( ver [2]) el cual asocia una serie numérica con direcciones de búsquedas que son válidas para el método del gradiente conjugado y en general para algoritmos de descenso, adaptado a una variedad Riemanniana  $M$  (ver [1]) . Igual que en  $\mathbb{R}^n$ , definimos en  $M$ , el ángulo  $\theta_k$  entre una dirección de búsqueda  $d_k$  y la dirección de descenso de mayor pendiente, el gradiente

negativo  $-Df(x)$ , como

$$\cos \theta_k = \frac{-Df(x_k) d_k}{\|Df(x_k)\|_k \|d_k\|_k} \quad (4.23)$$

**Teorema 4.3.1** (*Teorema de Zoutendijk*). Sea  $f : M \rightarrow \mathbb{R}$  acotada por abajo y de clase  $C^1$ . Suponga que  $\lambda_k$  en el **Algoritmo 1** satisfacen las ecuaciones (4.14). Si las funciones  $f_{R_{x_k}}$  son diferenciables continuamente sobre el  $\text{span } d_k$  y existe una constante Lipschitziana uniforme  $L > 0$  tal que

$$|Df_{R_{x_k}}(t d_k) d_k - Df_{R_{x_k}}(0_{x_k})(d_k)| \leq L t \quad d_k \in T_{x_k} M \text{ con } \|d_k\|_{x_k} = 1, x_k \in M, t \geq 0. \quad (4.24)$$

Entonces, la serie numérica

$$\sum_{k=0}^{\infty} \cos \theta_k \|Df(x_k)\|_{x_k}^2 < \infty, \quad (4.25)$$

converge.

**Prueba:**

De la condición de Wolfe (4.14b) tenemos que

$$(Df_{R_{x_k}}(\lambda_k d_k) d_k - Df_{R_{x_k}}(0) d_k) \geq (c_2 - 1) Df_{R_{x_k}}(0) d_k. \quad (4.26)$$

Luego, utilizando la condición de Lipschitz, resulta que

$$\lambda_k L \|d_k\|_{x_k}^2 \geq (Df_{R_{x_k}}(\lambda_k d_k) d_k - Df_{R_{x_k}}(0) d_k) \geq (c_2 - 1) Df(x_k) d_k,$$

de la cual obtenemos que

$$\lambda_k \geq \frac{(c_2 - 1) Df(x_k) d_k}{L \|d_k\|_{x_k}^2}.$$

Ahora, del **Algoritmo 1** y de la condición de Wolfe (4.14a) se tiene que

$$f(x_{k+1}) \leq f(x_k) - c_1 \frac{(1 - c_2)}{L} \cos^2 \theta_k \|Df(x_k)\|_{x_k}^2.$$

Colocando  $c = c_1 \frac{(1 - c_2)}{L} > 0$  y, sumando esta expresión sobre todos los índices menores que  $k$ , obtenemos

$$f(x_{k+1}) \leq f(x_0) - c \sum_{j=0}^k \cos^2 \theta_j \|Df(x_j)\|_{x_j}^2.$$

Ya que  $f$  es acotada por debajo, existe una constante  $\mathcal{C}$  tal que  $f(x) \geq \mathcal{C}$  para cualquier  $x \in M$ . Luego

$$\mathcal{C} \leq f(x_{k+1}) \leq f(x_0) - c \sum_{j=0}^k \cos^2 \theta_j \|Df(x_j)\|_{x_j}^2. \quad (4.27)$$

Así, tomando límite en (4.27) obtenemos

$$\sum_{k=0}^{\infty} \cos^2 \theta_k \|Df(x_k)\|_{x_k}^2 \leq \frac{f(x_0) - \mathcal{C}}{c} < \infty,$$

lo que concluye la prueba.  $\diamond$

Bien, la condición de Zoutendijk (4.25) implica que

$$\cos^2 \theta_k \|Df(x_k)\|_{x_k}^2 \rightarrow 0. \quad (4.28)$$

Si nuestro método para la elección de la dirección de búsqueda  $d_k$  en el **Algoritmo 1** garantiza que el ángulo  $\theta_k$  definido por (4.25) está acotado y alejado de  $90^\circ$ , existe una constante positiva  $\delta$  tal que

$$\cos \theta_k \geq \delta > 0 \quad \forall k.$$

Pero si esto es así, tenemos de (4.28), que

$$\lim_{k \rightarrow \infty} \|Df(x_k)\| = 0. \quad (4.29)$$

En otras palabras, podemos asegurar que la norma del gradiente  $\|Df(x_k)\|$  converge a cero, siempre que las direcciones de búsqueda no sean demasiado cercanas ortogonalmente al gradiente. En particular, el método de descenso de mayor pendiente (para el cual la dirección de búsqueda  $d_k$  es paralela al gradiente negativo) produce una sucesión de gradientes que converge a cero, si se utiliza una búsqueda lineal que satisface la condición de Wolfe.

**Corolario 4.3.1** (*Convergencia: Descenso de mayor pendiente generalizado*). *Suponga que la dirección de búsqueda en el **Algoritmo 1** es solución de  $\mathcal{B}_k(d_k, v) = -Df(x_k)v$  para todo  $v \in T_{x_k}M$ , donde las  $\mathcal{B}_k$  son formas bilineales acotadas y coercitivas uniformemente en  $T_{x_k}M$  (el caso  $\mathcal{B}_k(\cdot, \cdot) = g_{x_k}(\cdot, \cdot)$  admite la dirección de descenso de mayor pendiente). Entonces, bajo las mismas condiciones del teorema (4.3.1) se tiene que  $\|Df(x_k)\|_{x_k} \rightarrow 0$ .*

**Prueba:**

Siendo las  $\mathcal{B}_k$  acotadas y coercitivas uniformemente en  $T_{x_k}M$ , tenemos  $\mathcal{B}_k(d_k, v) \rightarrow +\infty$  cuando  $\|(d_k, v)\| \rightarrow +\infty$ ; con lo que  $\cos \theta_k = \frac{\mathcal{B}_k(d_k, d_k)}{\|\mathcal{B}_k(d_k, d_k)\| \|d_k\|_{x_k}}$  son uniformemente acotadas por encima de cero. Luego la convergencia se obtiene aplicando el teorema de Zoutendijk.  $\diamond$

Por la continuidad del gradiente  $Df$ , y el corolario anterior, se tiene que, cualquier punto



limite  $x^*$  de la sucesión  $\{x_k\}$  es un punto estacionario de la función  $f$ . Por lo tanto en variedades de dimensión finita, si  $\{x \in M : f(x) \leq f(x_0)\}$  es acotado,  $\{x_k\}$  se puede descomponer en subsucesiones que converjan a un punto estacionario. En variedades de dimensión infinita, donde sólo puede esperarse una convergencia débil de las subsucesiones, en general, esto no es cierto (lo cual no es sorprendente, dado que no se concede ni siquiera la existencia de puntos fijos, sin condiciones más fuertes sobre  $f$ , tales como secuencialmente débil semi-continua inferiormente). Más aún, el punto limite puede no ser único. A continuación, trataremos un caso especial de convergencia global, inspirada por la teoría clásica en  $\mathbb{R}$ , donde algunas restricciones se le harán a la dirección de descenso  $d_k$  y a la longitud del paso  $\lambda_k$  (ver [16]).

**Definición 4.3.1** (*Sucesión gradiente-relacionada*) Dada una función costo  $f$  sobre una variedad Riemanniana  $M$ , una sucesión  $\{d_k\}$ ,  $d_k \in T_{x_k}M$ , es gradiente-relacionada si, para cualquier subsucesión  $\{x_k\}_{k \in K}$  de  $\{x_k\}$  que converge a un punto no crítico de  $f$ , la correspondiente subsucesión  $\{d_k\}_{k \in K}$  es acotada y satisface que

$$\limsup_{k \rightarrow \infty, k \in K} \langle \text{grad } f(x_k), d_k \rangle < 0.$$

La siguiente definición, esta relacionada con el tamaño del paso  $\lambda_k$  en la condición de Armijo.

**Definición 4.3.2** (*Punto de Armijo*) Dada una función costo  $f$  sobre una variedad Riemanniana  $M$  con retracción  $R$ , un punto  $x \in M$ , un vector tangente  $d \in T_x M$ , y escalares  $\hat{\alpha} > 0$ ,  $\beta, c \in (0, 1)$ , el punto de Armijo es el número real  $\lambda^A$  tal que  $d^A =$

$\lambda^A d = \beta^m \hat{\alpha} d$ , donde  $m$  es el entero no negativo más pequeño que satisface la condición de Armijo

$$f(R_x(\beta^m \hat{\alpha} d)) - f(x) \leq c \langle \text{grad } f(x), \beta^m \hat{\alpha} d \rangle_x.$$

El número real  $\lambda^A = \beta^m \hat{\alpha}$  es el tamaño del paso de la condición de Armijo.

Nótese que, si escogemos  $x_{k+1} = R_{x_k}(\lambda_k^A d_k)$  en el paso 4 del **Algoritmo 1**, la condición de Armijo (4.13a) se cumple; pero no es obligatoria esta escogencia. Esta condición de Armijo da mucho margen de maniobra para sacar bastante información relacionada con un problema que puede conducir a un Algoritmo eficiente. Por ejemplo, se podría escoger  $x_{k+1} = R_{x_k}(\lambda_k^* d_k)$ , donde  $\lambda_k^* = \arg \min_{\lambda} f(R_{x_k}(\lambda d_k))$  satisface (4.13a); si la búsqueda lineal se lleva eficientemente.

**Teorema 4.3.2** (*Punto crítico de la función objetivo*) Sea  $\{x_k\}$  una sucesión infinita de iterados generados por el **Algoritmo 1**. Entonces cada punto de acumulación de  $\{x_k\}$  es un punto crítico de la función costo  $f$ .

**Prueba:**

Por reducción a lo absurdo, supongamos que existe una subsucesión  $\{x_k\}_{k \in K}$  convergiendo a un punto  $x^*$  con  $\text{grad } f(x^*) \neq 0$ . Ya que  $\{f(x_k)\}$  es no creciente, se sigue que toda ella converge a  $f(x^*)$ . En consecuencia  $f(x_{k+1}) - f(x_k)$  se va a cero. Ahora, de la construcción del algoritmo se tiene que,

$$f(x_{k+1}) - f(x_k) \leq c \sigma \lambda_k \langle \text{grad } f(x_k), d_k \rangle_x.$$

Ya que  $\{d_k\}$  es gradiente-relacionada, resulta que  $\{\lambda_k\}_{k \in K} \rightarrow 0$ . Como los  $\lambda_k$  son determinados por la condición de Armijo, entonces para todo  $k$  más grande que algún  $\widehat{k}$ , se tiene que  $\lambda_k = \beta^{m_k} \widehat{\alpha}$ , donde  $m_k$  es un entero más grande que cero. Pero esto significa que  $\frac{\lambda_k}{\beta} d_k$  no satisface la condición de Armijo. Por lo que

$$f(R_{x_k}\left(\frac{\lambda_k}{\beta} d_k\right)) - f(x_k) > \sigma \frac{\lambda_k}{\beta} \langle \text{grad } f(x_k), d_k \rangle_{x_k} \quad \forall k \in K, k \geq \widehat{k}.$$

Colocando

$$\widehat{d}_k = \frac{d_k}{\|d_k\|} \quad y \quad \widehat{\lambda}_k = \frac{\lambda_k \|d_k\|}{\beta},$$

en la desigualdad anterior tenemos

$$\frac{f_{R_{x_k}}(\widehat{\lambda}_k \widehat{d}_k) - f_{R_{x_k}}(0_{x_k})}{\widehat{\lambda}_k} > \sigma \langle \text{grad } f(x_k), \widehat{d}_k \rangle_{x_k}, \quad \forall k \in K, k \geq \widehat{k}$$

donde  $f_{R_{x_k}} = f \circ R_{x_k}$ . Por el teorema del valor medio, existe  $\tau \in [0, \widehat{\lambda}_k]$  tal que

$$Df_{R_{x_k}}(\tau \widehat{d}_k) \widehat{d}_k > \sigma \langle \text{grad } f(x_k), \widehat{d}_k \rangle_{x_k}, \quad \forall k \in K, k \geq \widehat{k}, \quad (4.30)$$

donde la diferencial se toma sobre el espacio vectorial  $T_{x_k} M$ . Ya que  $\{\lambda_k\}_{k \in K} \rightarrow 0$  y la sucesión  $\{d_k\}$  es gradiente-relacionada, por tanto acotada, resulta que  $\{\widehat{\lambda}_k\}_{k \in K} \rightarrow 0$ . Más aún, ya que  $\widehat{d}_k$  es de norma uno, éste pertenece a un conjunto compacto, con lo que existe un conjunto de índice  $\widehat{K} \subseteq K$  para el cual  $\{\widehat{d}_k\}_{k \in \widehat{K}} \rightarrow \widehat{d}_*$  para algún  $\widehat{d}_*$  con  $\|\widehat{d}_*\| = 1$ . Tomando limite en (4.30) sobre  $\widehat{K}$ , y ya que la métrica Riemanniana es continua,  $f \in \mathcal{C}^1$ , con  $Df_{R_{x_k}}(0_{x_k}) \widehat{d}_k = \langle \text{grad } f(x_k), \widehat{d}_k \rangle_{x_k}$ , tenemos que

$$\langle \text{grad } f(x_*), \widehat{d}_* \rangle_{x_*} \geq \sigma \langle \text{grad } f(x_*), \widehat{d}_* \rangle_{x_*}.$$

De donde se tiene que

$$\langle \text{grad } f(x_*), \widehat{d}_* \rangle_{x_*} (1 - \sigma) \geq 0.$$

Pero  $\sigma < 1$ , así que

$$\langle \text{grad } f(x_*), \widehat{d}_* \rangle_{x_*} \geq 0.$$

Por otra parte, la sucesión  $\{d_k\}$  es gradiente-relacionada, con lo que

$$\langle \text{grad } f(x_*), \widehat{d}_* \rangle_{x_*} < 0,$$

obteniéndose una contradicción.  $\diamond$

## 4.4. Aplicación

**Ejemplo1.** [Con Descenso (ver [7]). *Buscando el mínimo sobre la esfera unitaria  $S^2$*  ]

Sea  $S^2 : x^2 + y^2 + z^2 = 1$  la esfera en  $\mathbb{R}^3$ . Ya que  $S^2$  es un conjunto compacto, cualquier función que sea por lo menos continua sobre  $S^2$  tiene un mínimo y un máximo global.

Sea  $f : \mathbb{R}^3 \rightarrow \mathbb{R}$ ,  $f(x, y, z) = e^z$ . Si queremos determinar los extremos de  $f$  sobre la esfera  $S^2 : x^2 + y^2 + z^2 = 1$ , podemos utilizar los multiplicadores de Lagrange y encontrar que,  $p = (0, 0, 1)$  y  $q = (0, 0, -1)$  serían los puntos críticos. En el punto  $p$  tendríamos el máximo de  $f$ ,  $\max f = f(p) = e$ , en el punto  $q$  el mínimo,  $\min f = f(q) = 1/e$ .

Ahora de acuerdo a lo estudiado anteriormente es este capítulo, el problema se puede resolver como sigue. Sea  $f : S^2 \rightarrow \mathbb{R}$ ,  $f(x, y, z) = e^z$ . Encontramos, aplicando la teoría de

subvariedades que, (ver [3.20])

$$\text{gra } f = (0, 0, e^z), \quad \xi = (x, y, z) \quad a = g(\text{grad}_{\mathbb{R}^3} f, \xi) = df(\xi) = (0, 0, e^z) \cdot (x, y, z) = z e^z.$$

Luego

$$\text{grad}_{S^2} f = \text{grad}_{\mathbb{R}^3} f - a \xi = (0, 0, e^z) - z e^z (x, y, z),$$

de donde tenemos que

$$\text{grad}_{S^2} f = e^z (-xz, -yz, 1 - z^2).$$

De aquí que los puntos críticos son  $p = (0, 0, 1)$  y  $q = (0, 0, -1)$  como se dijo antes.

Con respecto a la geodésica en la esfera  $S^2$ , tenemos las siguientes ecuaciones paramétricas

$$\begin{cases} x = a \cos t + d \sin t \\ y = b \cos t + e \sin t \\ z = c \cos t + f \sin t, \quad t \in \mathbb{R}, \end{cases}$$

donde  $(a, b, c)$  y  $(d, e, f)$  son versores ortogonales, es decir, de módulo uno. Colocando el punto inicial como

$$p_1 = \left( \frac{1}{\sqrt{2}}, \frac{1}{\sqrt{2}}, 0 \right).$$

Ya que  $\text{grad}_{S^2} f(p) = \vec{k} = (0, 0, 1)$ , la dirección y el sentido de descenso de  $f$  por  $p_1$  son indicado por  $-\vec{k}$ . La geodésica que inicia en  $p_1$  en dirección de  $-\vec{k}$  es

$$x(0) = \frac{1}{\sqrt{2}}, \quad y(0) = \frac{1}{\sqrt{2}}, \quad z(0) = 0, \quad x'(0) = 0, \quad y'(0) = 0, \quad z'(0) = -1$$

o equivalente a

$$a = \frac{1}{\sqrt{2}}, \quad b = \frac{1}{\sqrt{2}}, \quad c = 0, \quad d = 0, \quad e = 0, \quad f = -1;$$

con lo que

$$x = \frac{1}{\sqrt{2}} \cos t, \quad y = \frac{1}{\sqrt{2}} \cos t, \quad z = -\sin t, \quad t \in [0, \infty).$$

Sea

$$\omega(t) = f(x(t), y(t), z(t)) = e^{-\sin t}, \quad t \in [0, \infty).$$

Así que, en  $t = \pi/2$ ,  $\omega(\pi/2) = f(0, 0, -1) = f(q) = e^{-1}$ , obtiene su mínimo por  $q$ . Ahora, como  $\text{grad}_{S^2} f(q) = 0$ , se tiene que  $q$  es un punto crítico de  $f$ ; más, es un mínimo global, pues el Hessiano en este punto es definido positivo. Este extremo se obtuvo después de realizar un único paso. Bien para chequear este resultado, fijamos un punto de inicio  $p_2 = \left(\frac{1}{\sqrt{2}}, 0, \frac{1}{\sqrt{2}}\right)$ . Encontramos que

$$\text{grad}_{S^2} f(p_2) = e^{1/\sqrt{2}} \left( -\frac{1}{2} \vec{i} + \frac{1}{2} \vec{k} \right).$$

De esta forma el versor  $\vec{v} = \frac{\vec{i}}{\sqrt{2}} - \frac{\vec{k}}{\sqrt{2}}$ , indica la dirección y el sentido de  $f$  por  $p_2$ . Aquí, la geodésica que comienza por  $p_2$  en dirección de  $\vec{v}$  es

$$x = \frac{1}{\sqrt{2}} (\cos t + \sin t), \quad y = 0, \quad z = \frac{1}{\sqrt{2}} (\cos t - \sin t) \quad t \in [0, \infty).$$

Ya que la función

$$\phi(t) = f(x(t), y(t), z(t)) = e^{\sin(\pi/4-t)}$$

obtiene su mínimo  $e^{-1}$  por  $t = 3\pi/4$ , es decir, por  $x = 0$ ,  $y = 0$ ,  $z = -1$ ; de nuevo llegamos al punto  $q$  (ver figura 4.3).

**Ejemplo 2.**[Con Armijo (ver [16]). *Minimizando el Cociente de Rayleigh sobre  $S^{n-1}$* ]

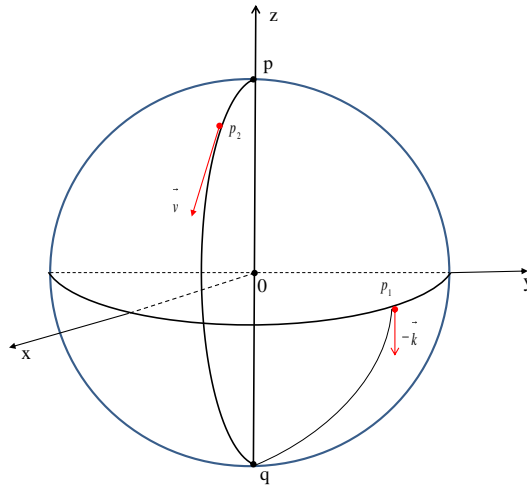


Figura 4.3: Método de descenso en  $S^2$ .

**Definición 4.4.1** (*Cociente de Rayleigh*) Si  $x$  es autovector de la matriz  $A$ , entonces su correspondiente autovalor es dado por

$$\mu = \frac{x^T A x}{x^T x}.$$

Este cociente se llama el **Cociente de Rayleigh**.

Se sabe que este cociente de Rayleigh puede acelerar convergencia en los métodos iterativos.

Con el objetivos de presentar una aplicación del **Algoritmo 1**, vamos a utilizarlo aquí.

Se aplicara el **Algoritmo 1** al problema de encontrar un minimizador de

$$f : S^{n-1} \rightarrow \mathbb{R} \quad \text{donde} \quad f(x) = x^T A x, \quad (4.31)$$

que denota el cociente de Rayleigh sobre la esfera unitaria. La matriz  $A$  es simétrica ( $A = A^T$ ) y no necesariamente definida positiva (ver 1.2.2). Sea  $\mu_1$  el autovalor más

pequeño de  $A$  cuyo autovector asociado de norma unitaria es  $v_1$ .

Considere la función costo

$$\widehat{f}: \mathcal{S}^{n-1} \rightarrow \mathbb{R} \quad \text{donde} \quad \widehat{f}(x) = x^T A x, \quad (4.32)$$

cuya restricción a la esfera unitaria  $\mathcal{S}^{n-1}$  admite (4.31). Aquí, estamos viendo a  $\mathcal{S}^{n-1}$  como una subvariedad Riemanniana del espacio euclídeo  $\mathbb{R}^n$ , dotada con la métrica Riemanniana

$$\widehat{g}(u, v) = u^T v.$$

Dado  $x \in \mathcal{S}^{n-1}$ , tendremos que

$$D\widehat{f}v = v^T A x + x^T A v = 2v^T A x,$$

para todo  $v \in T_x \mathbb{R}^n \simeq \mathbb{R}^n$ . De la definición de gradiente se tiene que

$$\text{grad } \widehat{f}(x) = 2 A x.$$

El espacio tangente a  $\mathcal{S}^{n-1}$ , visto como subvariedad de  $T_x \mathbb{R}^n \simeq \mathbb{R}^n$  es

$$T_x \mathcal{S}^{n-1} = \{v \in \mathbb{R}^n : x^T v = 0\}.$$

El espacio normal es

$$(T_x \mathcal{S}^{n-1})^\perp = \{x \rho : \rho \in \mathbb{R}\}.$$

Las proyecciones ortogonales al espacio tangente y al normal son

$$P_x v = v - x x^T v, \quad P_x^\perp v = x x^T v.$$



Recordemos que cualquier elemento  $v \in T_x \mathbb{R}^n$  se puede escribir como la suma directa

$$v = P_x v + P_x^\perp v.$$

Ahora, ya que en general, si  $\hat{f}$  es la función costo sobre una variedad  $\overline{M}$  y  $f$  su restricción a la subvariedad  $M$ , el gradiente de  $f$  es el gradiente de  $\hat{f}$  proyectado sobre  $T_x M$ , es decir,

$$\text{grad } f(x) = P_x \text{grad } \hat{f}(x).$$

Así que en nuestro caso

$$\text{grad } f(x) = 2P_x(Ax) = 2(Ax - xx^T Ax). \quad (4.33)$$

La retracción  $R_x$  sobre  $T_x \mathbb{R}^n \simeq \mathbb{R}^n$  es definida como  $R_x v = x + v$  así que en  $T_x \mathcal{S}^{n-1}$ , la podemos definir como

$$R_x v = \frac{x + v}{\|x + v\|}.$$

Pero ya que  $\mathcal{S}^{n-1}$  pertenece al grupo de matrices ortogonales, una alternativa para la retracción también es

$$R_x v = qf(x + v),$$

donde  $qf(A)$  denota el  $Q$  factor de la descomposición de  $A \in \mathbb{R}_*^{n \times p}$  como  $A = QR$  donde  $Q$  pertenece a las  $St(p, n) = \{X \in \mathbb{R}^{n \times p} : X^T X = I_p\}$  y  $R$  es una matriz triangular superior  $n \times p$  con los elementos de la diagonal estrictamente positivos.

Todas las formulas dadas arriba están resumidas en el cuadro (4.1).

	Variedad ( $\mathcal{S}^{n-1}$ )	Espacio de embedding ( $\mathbb{R}^n$ )
Función Costo	$f(x) = x^T A x, x \in \mathcal{S}^{n-1}$	$\hat{f}(x) = x^T A x, x \in \mathbb{R}^n$
Métrica	Métrica inducida	$\hat{g}(u, v) = u^T v$
Espacio Tangente	$v \in \mathbb{R}^n : x^T v = 0$	$\mathbb{R}^n$
Espacio Normal	$v \in \mathbb{R}^n : v = \alpha x$	$\emptyset$
Proyección sobre el Espacio Tangente	$P_x v = (I - x x^T) v$	Identidad
Gradiente	$\text{grad } f(x) = P_x \text{grad } \hat{f}(x)$	$\text{grad } \hat{f}(x) = 2 A x$
Retracción	$R_x v = q f(x + v)$	$R_x v = x + v$

Cuadro 4.1: Cociente de Rayleigh sobre la esfera unitaria

Las siguientes proposiciones nos hablan sobre los puntos críticos, mínimos, máximos y puntos sillas del cociente de Rayleigh. Las pruebas respectivas se pueden ver en [5].

**Proposición 4.4.1** *Sea  $A = A^T$  una matriz simétrica de orden  $n \times n$ . Un vector de norma unitaria  $x \in \mathbb{R}^n$  es un autovector de  $A$  si y sólo si este es un punto crítico del cociente de Rayleigh.*

**Proposición 4.4.2** *Sea  $A = A^T$  una matriz simétrica de orden  $n \times n$  con autovalores  $\mu_1 \leq \dots \leq \mu_n$  y autovectores ortonormales asociados  $v_1, \dots, v_n$ . Entonces*

- $\pm v_1$  son minimizadores locales y globales del Cociente de Rayleigh (4.31) si el autovalor  $\mu_1$  es simple (multiplicidad 1), con lo que serían los únicos minimizadores.

- $\pm v_n$  son maximizadores locales y globales del Cociente de Rayleigh (4.31) si el autovalor  $\mu_n$  es simple, con lo que serían los únicos minimizadores.
- $\pm v_t$  son los autovectores correspondientes a los autovalores interiores ( que están estrictamente entre  $\mu_1$  y  $\mu_n$  ), son los puntos sillas del Cociente de Rayleigh (4.31) .

Se desprende de la Proposición (4.4.1) y el análisis de convergencia global del métodos de búsqueda lineal (Proposición 4.3.2), que todos los métodos dentro de la clase del **Algoritmo 1**, producen iteraciones que convergen en el conjunto de vectores propios de  $A$ . Y como estamos considerando el método de descenso, si el autovalor  $\mu_1$  es simple, se tendrá convergencia estable para  $\pm v_1$  e inestable para todos los demás autovectores.

Luego, en el caso que nos interesa para el **Algoritmo 1**, vamos a colocar la dirección (gradiente-relacionada) como

$$d_k = -\text{grad } f(x_k) = -2 (A x_k - x_k x_k^T A x_k).$$

En lo siguiente, para no complicar el asunto referente a la retracción ya comentada, tomemos

$$R_x v = \frac{x + v}{\|x + v\|}, \quad (4.34)$$

donde  $\| \cdot \|$  representa la norma euclidiana de  $\mathbb{R}^n$  dada por,  $\| z \| = \sqrt{z^T z}$ . Otra posibilidad es

$$R_x v = x \cos \| v \| + \frac{v}{\|v\|} \sin \| v \|,$$

para la cual la curva  $t \rightarrow R_x(t v)$  en un círculo grande de la esfera.

Bien, ya con todos los ingredientes necesario planteamos el **Algoritmo 1** en el siguiente cuadro: Este algoritmo presenta resultados numéricos ya probados (ver [5]), por ejemplo

---

**Algoritmo 1\*** Búsqueda lineal de Armijo para el Cociente de Rayleigh sobre  $\mathcal{S}^{n-1}$

---

**Requiere:** La matriz simétrica  $A$ , escalares  $\hat{\lambda} > 0$ ,  $\beta, \sigma \in (0, 1)$ .

**Entrada:** Iterado inicial  $x_0, \|x_0\|$ .

**Salida:** Sucesión de iterados  $\{x_k\}$ .

1: **Para**  $k = 0, 1, 2, \dots$  **hacer**

2:   Calcular  $d_k = -2 (A x_k - x_k x_k^T A x_k)$ .

3:   Encontrar el entero más pequeño  $m \geq 0$  tal que

$$f(R_{x_k}(\hat{\lambda} \beta^m d_k)) \leq f(x_k) - \sigma \hat{\lambda} \beta^m d_k^T d_k,$$

con  $f$  definida en (4.31) y  $R$  en (4.34).

4: Colocar

$$x_{k+1} = R_{x_k}(\hat{\lambda} \beta^m d_k).$$

5: **Finalizar**

---

se puede tomar la matriz  $A$  como,  $A = \text{diag}(1, 2, \dots, 100)$ ,  $\sigma = 0,5$ ,  $\lambda = 1$  y  $\beta = 0,5$ ; el punto de inicio  $x_0$  se escoge normalizando un vector cuya entradas son seleccionadas de una distribución normal. Pero el asunto aquí, es que éste mismo algoritmo, por la teoría expuesta, también presenta convergencia al agregar la segunda condición de Wolfe (4.13b);

inclusive la condición fuerte. El mismo quedaría de la forma siguiente:

---

**Algoritmo 1\*** Búsqueda lineal de Armijo para el Cociente de Rayleigh sobre  $\mathcal{S}^{n-1}$

---

**Requiere:** La matriz simétrica  $A$ , escalares  $\widehat{\lambda} > 0$ ,  $\beta$ ,  $\sigma \in (0, 1)$ .

**Entrada:** Iterado inicial  $x_0$ ,  $\|x_0\|$ .

**Salida:** Sucesión de iterados  $\{x_k\}$ .

1: **Para**  $k = 0, 1, 2, \dots$  **hacer**

2: Calcular  $d_k = -2 (A x_k - x_k x_k^T A x_k)$ .

3: Encontrar el entero más pequeño  $m \geq 0$  tal que  $0 < c_1 = \sigma \beta^m < \frac{c_2}{m} < 1$

$$f(R_{x_k}(\widehat{\lambda} \beta^m d_k)) \leq f(x_k) - \sigma \widehat{\lambda} \beta^m d_k^T d_k,$$

y

$$D f(R_{x_k}(\widehat{\lambda} \beta^m d_k)) T_{x_k, R_{x_k}(\widehat{\lambda} \beta^m d_k)}^{R_{x_k}} d_k \leq c_2 d_k^T d_k,$$

con  $f$  definida en (4.31),  $R$  en (4.34) y  $T^{R_{x_k}}$  es el

transporte vectorial escalar que representa la derivada  $D R_{x_k}$  de la

retracción  $R_{x_k}$  sobre  $M$ .

4: Colocar

$$x_{k+1} = R_{x_k}(\widehat{\lambda} \beta^m d_k).$$

5: **Finalizar**

---

[Se hace la siguiente observación, el valor de  $m$  dado en el punto de Armijo (4.3.2), no es fácil de conseguir. En general se toma de tal forma que  $0 < c_1 = \sigma \beta^m < \frac{c_2}{m} < 1$  .]



## Conclusión

---

En cuanto al objetivo planteado en este trabajo, hemos presentado la búsqueda lineal con condición de Wolfe dada en  $\mathbb{R}^n$ , sobre variedades Riemannianas. Valiéndonos para ello, de instrumentos de la geometría diferencial como son las retracciones junto con el transporte vectorial paralelo, que nos permitieron caminar sobre la estructura no euclidiana e implementar las estrategias necesarias para conseguir convergencia en estos tipos de métodos. Por su puesto, se buscó destacar los detalles de las demostraciones que, en general, cuando se trata de variedades se opta por decir, que si se cumple en el  $\mathbb{R}^n$  correspondiente, de igual forma debe ocurrir en la misma.





## Bibliografía

---

- [1] W. RING AND B. WIRTH, *Optimization Methods Riemannian Manifold and their application to shape space*, SIAMJ. Optim, 22(2):596-627, 2012.
- [2] J. NOCEDAL AND S. WRIGHT, *Numerical Optimization*, Springer Series in Operations Research and Financial Engineering Springer, New York, Second Edition, 2006.
- [3] D. GABAY, *Minimizing a Differentiable Function over a Differential Manifold*, J. Optim. Theory Appl., 37(2):177-219, 1982.
- [4] R. E. MAHONY, *The Constrained Newton Method on a Lie Grupo and the Symmetric Eigenvalue Problem*, Linear Algebra Appl. 23(3):309-327, 2002.
- [5] R. MAHONY AND J. MANTON, *The Geometry of the Newton on non-compact Lie Grupos*, J Global Optim. 248:367-89, 1996.
- [6] IN: A. BLOCH(ED), *Halmition and Gradient Flows*, Algorithms and Control, Fields Institute Communication, Vol 3, Amer. Math. Soc, pp:113-136, 1994.

- [7] C. UDRISTE, *Convex Functions and Optimization Methods on Riemannian Manifold*, Kluwer Acad, Publ., Pordre 1994) pp.1-8 Rio de Janeiro: Sociedade Brasileira de Matemática, tercera edición (2005).
- [8] C. BAKER, P. ABSIL AND K. GALLIVAN *An implicit trust-region method on Riemannian Manifold*, IMA-J. Number Anal, 28(4):665-207 (1988).
- [9] Y. YANG, *Globally Convergent Optimization Algorithms on Riemannian Manifolds: Uniform framework for unconstrained and constrained optimization*, J. Optim. Theory Appl.,132(2): 245-265, 2007.
- [10] M. EPELMAN, *Continuos Optimization Methods*, Section 1.IOE 511 . Math 652. Fall,2000.
- [11] L. GRIPO,F. LAMPARIELLO AND S. LUCIDI, *A Nonmonotone Line Search Technique For Newton's Method*, Istituto di analisi dei Sistemi ed Informatica del CNR (1985).
- [12] F. CAMARGO, *Estudo Comparativo de Passos Espectrais e Buscas Lineares não monótonas* , Instituto De matemática E Estatística Da Universidade De São Paulo (2008).
- [13] E. QUIROZ, E. QUISPE AND P. OLIVEIRA, *steepest Descent Method for Quasiconvex Minimization on Riemannian Manifold*, Federal University of Rio de Janeiro, PESC-COPPE (2006).

- [14] S. DO CARMO, MANFREDO PERDIGAO, *Geometria Riemanniana*, Rio de Janeiro: IMPA-Projeto Euclides, segunda edición, (1988).
- [15] ERIK QUIROZ, *Un método no Euclidiano para Problemas de Optimización en Espacios Euclidianos*, Universidad Nacional Del Callao (2007).
- [16] A. ABSIL, R. MAHONY AND R. SEPULCHRE, *Riemanniana Geometry of Grassmann Manifold With a View on Algorithmic Computation*, Florida State University Tallahassee,FL 32306-4120,Usa (2000).
- [17] BOOTHBY WILLIAM, *An Introduction to Differentiable Manifold and Riemannian Geometry*, Florida USA: Edit. Adademic Press, First Edition, 1986.
- [18] LAWRENCE CONLON, *Differentiable Manifold a First Course*, Birkhäuser Advanced Texts.Boston.Basel.Berlin: Edited by Herbert Amann,Zürich (1993).
- [19] B. GONZÁLEZ,D. BENAVIDES, *Nociones de Análisis Funcional*, UNIVERSIDAD DE SEVILLA.Departamento de Análisis Matemático.Sevilla,España (2010).