# Lecture Notes in Mathematics

## 1133

# Krzysztof C. Kiwiel

# Methods of Descent for Nondifferentiable Optimization

# Lecture Notes in Mathematics

## 1133

## Krzysztof C. Kiwiel

## Methods of Descent for Nondifferentiable Optimization

Springer-Verlag
Berlin Heidelberg New York Tokyo

**Author**

Krzysztof C. Kiwiel
Systems Research Institute, Polish Academy of Sciences
ul. Newelska 6, 01-447 Warsaw, Poland

# Lecture Notes in Mathematics

## PREFACE

This book is about numerical methods for problems of finding the largest or smallest values which can be attained by functions of several real variables subject to several inequality constraints. If such problems involve continuously differentiable functions, they can be solved by a variety of methods well documented in the literature. We are concerned with more general problems in which the functions are locally Lipschitz continuous, but not necessarily differentiable or convex. More succintly, this book is about numerical methods for nondifferentiable optimization.

Nondifferentiable optimization, also called nonsmooth optimization, has many actual and potential applications in industry and science. For this reason, a great deal of effort has been devoted to it during the last decade. Most research has gone into the theory of nonsmooth optimization, while surprisingly few algorithms have been proposed, these mainly by C.Lemaréchal, R.Mifflin and P.Wolfe. Frequently such algorithms are conceptual, since their storage and work per iteration grow infinitely in the course of calculations. Also their convergence properties are usually weaker than those of classical methods for smooth optimization problems.

This book gives a complete state-of-the-art in general-purpose methods of descent for nonsmooth minimization. The methods use piecewise linear approximations to the problem functions constructed from several subgradients evaluated at certain trial points. At each iteration, a search direction is found by solving a quadratic programming subproblem and then a line search produces both the next improved approximation to a solution and a new trial point so as to detect gradient discontinuities. The algorithms converge to points satisfying necessary optimality conditions. Also they are widely applicable, since they require only a weak semismoothness hypothesis on the problem functions which is likely to hold in most applications.

A unifying theme of this book is the use of subgradient selection and aggregation techniques in the construction of methods for nondifferentiable optimization. It is shown that these techniques give rise in a totally systematic manner to new implementable and globally convergent modifications and extensions of all the most promising algorithms which have been recently proposed. In effect, this book should give the reader a feeling for the way in which the subject has developed and is developing, even though it mainly reflects the author's research.

This book does not discuss methods without a monotonic descent (or ascent) property, which have been developed in the Soviet Union.

The reason is that the subject of their effective implementations is
still a mystery. Moreover, these subgradient methods are well descri-
bed in the monograph of Shor (1979). We refer the reader to Shor´s
excellent book (its English translation was published by Springer-
Verlag in 1985) for an extensive discussion of specific nondifferent-
iable optimization problems that arise in applications. Due to space
limitations, such applications will not be treated in this book.

In order to make the contents of this book accessible to as wide
a range of readers as possible, our analysis of algorithms will use
only a few results from nonsmooth optimization theory. These, as well
as certain other results that may help the reader in applications, are
briefly reviewed in the introductory chapter, which also contains a
review of representative existing algorithms. The reader who has basic
familiarity with nonsmooth functions may skip this chapter and start
reading from Chapter 2, where methods for unconstrained convex minimi-
zation are described in detail. The basic constructions of Chapter 2
are extended to the unconstrained nonconvex case in two fundamentally
different ways in Chapters 3 and 4, giving rise to competitive methods.
Algorithms for constrained convex problems are treated in Chapter 5,
and their extensions to the nonconvex case are described in Chapter 6.
Chapter 7 presents new versions of the bundle method of Lemaréchal and
its extensions to constrained and nonconvex problems. Chapter 8 con-
tains a few numerical results.

The book should enable research workers in various branches of
science and engineering to use methods for nondifferentiable optimizat-
ion more efficiently. Although no computer codes are given in the text,
the methods are described unambiguously, so computer programs may rea-
dily be written.

TABLE OF CONTENTS

Chapter 6. Methods of Feasible Directions for Nonconvex
Constrained Problems

Chapter 7. Bundle Methods

Chapter 8. Numerical Examples

CHAPTER 1

Fundamentals


1. Introduction

The nonlinear programming problem, also known as the mathematical
programming problem, can be taken to have the form

P :  minimize  $f(x)$,  subject to  $F_i(x) \leq 0$  for  $i=1,\ldots,m$,

where the objective function  $f$  and the constraint functions  $F_i$  are
real-valued functions defined on the N-dimensional Euclidean space $R^N$.
The value of  $m \geq 0$  is finite; when  m=0  the problem is unconstrained.
Often the optimization problem  P  is smooth: the problem functions  $f$
and  $F_i$  are continuously differentiable, i.e. they have continuous gra-
dients  $\nabla f$  and  $\nabla F_i$, $i=1,\ldots,m$. But in many applications this is not
true. Nonsmooth problems are the subject of nonsmooth optimization, also
called nondifferentiable optimization.

Owing to actual and potential applications in industry and science,
recently much research has been conducted in the area of nonsmooth opti-
mization both in the East (see the excellent monographs by Gupal (1979),
Nurminski (1979) and Shor (1979)) and in the West (see the comprehensive
bibliographies of Gwinner (1981) and Nurminski (1982)).

Nonsmooth problems that arise in applications have certain common
features. They are more complex and have poorer analytical properties
than standard mathematical programming problems, cf. (Bazaraa and Shetty,
1979; Pshenichny  and Danilin, 1975). A single evaluation of the problem
functions usually requires solutions of auxiliary optimization subprob-
lems. In particular, it is very common to encounter a nondifferentiable
function which is the pointwise supremum of a collection of functions that
may themselves be differentiable - a max function.

Functions with discontinuous gradients, such as max functions, cannot
be minimized by classical nonlinear programming algorithms. This observa-
tion applies both to gradient-type algorithms (the method of steepest descent,
conjugate direction methods, quasi-Newton methods) and to direct search
methods which do not require calculation of derivatives (the method of
Nelder and Mead, the method of Powell, etc.), see (Lemarechal, 1978 and
1982; Wolfe, 1975).

This work is concerned with numerical methods for finding (approxi-
mate) solutions to problem P when the problem functions are locally Lip-
schitzian, i.e. Lipschitz continuous on each bounded subset of $R^N$, but not

necessarily differentiable.

The advent of F.H.Clarke's (1975) analysis of locally Lipschitzian functions provided a unified approach to both nondifferentiable and non-convex problems (Clarke, 1976). Clarke's subdifferential analysis, the pertinent part of which is briefly reviewed in the following section, suffices for establishing properties of a vast class of optimization problems that arise in applications (Pshenichny, 1980; Rockafellar, 1978).

## 2. Basic Results of Nondifferentiable Optimization Theory

In this section we describe general properties of nondifferentiable optimization problems that are the subject of this work. Basic familiarity is, however, assumed. Source material may be found in (Clarke, 1975; Clarke, 1976; Rockafellar, 1970; Rockafellar, 1978; Rockafellar, 1981).

The section is organized as follows. First, we review concepts of differentiability and elementary properties of the Clarke subdifferential. The proofs are omitted, because only simple results, such as Lemma 2.2, will be used in subsequent chapters. Other results, in particular the calculus of subgradients, should help the reader who is mainly interested in applications. Secondly, we study convex first order approximations to nondifferentiable functions. Such approximations are then used for deriving necessary conditions of optimality for nondifferentiable problems. Our approach is elementary and may appear artificial. However, it yields useful interpretations of the algorithms described in subsequent chapters.

The following notation is used. We denote by $<\cdot,\cdot>$ and $|\cdot|$, respectively, the usual inner product and norm in finite-dimensional, real Euclidean space. $R^N$ denotes Euclidean space of dimension $N<\infty$. We use $x_i$ to denote the i-th component of the vector x. Thus

$$<x,y> = \sum_{i=1}^{N} x_i y_i \quad \text{and} \quad |x| = <x,x>^{1/2} \quad \text{for} \quad x,y \in R^N.$$ Superscripts are used to denote different vectors, e.g. $x^1$ and $x^2$. All vectors are column vectors. However, for convenience a column vector in $R^{N+n}$ is sometimes denoted by $(x,y)$ even though x and y are column vectors in $R^N$ and $R^n$, respectively. $[x,y]$ denotes the line segment joining x and y in $R^N$, i.e. $[x,y] = \{z \in R^N : z = \lambda x + (1-\lambda)y \text{ for some } \lambda \text{ satisfying } 0 \le \lambda \le 1\}$.

A set $S \subset R^N$ is called <u>convex</u> if $[x,y] \subset S$ for all x and y belonging to S. A linear combination $\sum_{j=1}^{k} \lambda_j x^j$ is called a <u>convex combination</u> of points $x^1,...,x^k$ in $R^N$ if each $\lambda_j \ge 0$ and $\sum_{j=1}^{k} \lambda_j = 1$. The <u>convex hull</u> of a set $S \subset R^N$, denoted conv S, is the set of all convex combina-

tions of points in S. conv S is the smallest convex set containing S, and S is convex if and only if S=conv S. An important property of convex hulls is described in

Lemma 2.1 (Caratheodory's theorem; see Theorem 17.1 in (Rockafellar, 1970)).

If $S \subset R^N$ then $x \in$ conv S if and only if x is expressible as a convex combination of N+1 (not necessarily different) points of S.

Any nonzero vector $g \in R^N$ and number $\gamma$ define a hyperplane

$$H = \{x \in R^N : <g,x> = \gamma\},$$

which is a translation of the (N-1)-dimensional subspace $\{x \in R^N : <g,x>=0\}$ of $R^N$. H divides $R^N$ into two closed half-spaces $\{x \in R^N : <g,x> \leq \gamma\}$ and $\{x \in R^N : <g,x> \geq \gamma\}$, respectively. We say that H is a supporting hyperplane to a set $S \subset R^N$ at $\bar{x} \in S$ if $<g,\bar{x}> = \gamma$ and $<g,x> \leq \gamma$ for all $x \in S$. Any closed convex set S can be described as an intersection of all the closed half-spaces that contain S.

We use the set notation

$$S^1 + S^2 = \{z^1 + z^2 : z^1 \in S^1, z^2 \in S^2\},$$
$$\text{conv}\{S^i : i=1,2\} = \text{conv}\{z : z \in S^1 \cup S^2\}$$

for any subsets $S^1$ and $S^2$ of $R^N$.

A function $f:R^N \longrightarrow R$ is called convex if

$$f(\lambda x^1 + (1-\lambda)x^2) \leq \lambda f(x^1) + (1-\lambda)f(x^2) \quad \text{for all} \quad \lambda \in [0,1] \quad \text{and} \quad x^1, x^2 \in R.$$

This is equivalent to the epigraph of f

$$\text{epi } f = \{(x,\beta) \in R^{N+1} : \beta \geq f(x)\}$$

being a convex subset of $R^{N+1}$. A function $f:R^N \longrightarrow R^1$ is called concave if the function $(-f)(x)=-f(x)$ is convex. If $F_i:R^N \longrightarrow R$ is convex and $\lambda_i \geq 0$ for each $i=1,\ldots,k$, then the functions

$$\phi_1(x) = \sum_{i=1}^{k} \lambda_i f_i(x),$$

$$\phi_2(x) = \max\{f_i(x) : i=1,\ldots,k\}$$

$$(2.1)$$

are convex.

A function $f:R^N \longrightarrow R$ is strictly convex if $f(\lambda x^1 + (1-\lambda)x^2) < \lambda f(x^1) + (1-\lambda)f(x^2)$ for all $\lambda \in (0,1)$ and $x^1 \neq x^2$. For instance, the

function $|.|^2$ is strictly convex.

A function $f : R^N \longrightarrow R$ is said to be <u>locally Lipschitzian</u> if for each bounded subset $B$ of $R^N$ there exists a Lipschitz constant $L = L(B) < \infty$ such that

$$|f(x^1)-f(x^2)| \leq L| x^1-x^2 | \quad \text{for all} \quad x^1,x^2 \in B. \tag{2.2}$$

Then in particular $f$ is continuous. Examples of locally Lipschitzian functions include continuously differentiable functions, convex functions, concave functions and any linear combination or pointwise maximum of a finite collection of such functions, cf. (2.1).

Following (Rockafellar, 1978), we shall now describe differentiability properties of locally Lipschitzian functions. Henceforth let $f$ denote a function satisfying (2.2) and let $x$ be an interior point of $B$, i.e. $x \in \text{int } B$.

The Clarke <u>generalized directional derivative</u> of $f$ at $x$ in a direction d

$$f^O(x;d) = \lim_{y \to x, t \downarrow 0} \sup \left[f(y+td)-f(y)\right]/t \tag{2.3}$$

is a finite, convex function of d and $f^O(x;d) \leq L|d|$. The Dini <u>upper directional derivative</u> of $f$ at $x$ in a direction d

$$f^D(x;d) = \lim_{t \downarrow 0} \sup \left[f(x+td)-f(x)\right]/t \tag{2.4}$$

exists for each $d \in R^N$ and satisfies

$$f(x+td) \leq f(x)+tf^D(x;d)+o(t), \tag{2.5}$$

where $o(t)/t \to 0$ as $t \downarrow 0$. The limit

$$f'(x;d) = \lim_{t \downarrow 0} \left[f(x+td)-f(x)\right]/t \tag{2.6}$$

is called the (one-sided) <u>directional derivative</u> of $f$ at $x$ with respect to d, if it exists. The two-sided derivative (the Gateaux derivative) corresponds to the case $f'(x;-d)=-f'(x;d)$. Clearly,

$$f^D(x;d) \leq f^O(x;d),$$
$$f'(x;d) \leq f^D(x;d), \tag{2.7}$$

whenever $f'(x;d)$ exists.

If $f'(x;d)$ is linear in d (Gateaux differentiable at x)

$$f'(x;d) = < g_f,d > \quad \text{for all} \quad d \in R^N, \tag{2.8}$$

then the vector $g_f$ is called the <u>gradient</u> of f at x and denoted by
$\nabla f(x)$. The components of $\nabla f(x)=(\frac{\partial f}{\partial x_1}(x),\dots,\frac{\partial f}{\partial x_N}(x))$ are the coordinate-
wise two-sided partial derivatives of f at x. The function f is (Frechet)
<u>differentiable</u> at x if

$$f(x+d)=f(x)+<\nabla f(x),d>+o(|d|) \quad \text{for all } d \in R^N , \qquad (2.9)$$

where $o(t)/t \to 0$ as $t \downarrow 0$. The above relation is equivalent to

$$\lim_{d' \to d, t \downarrow 0} \left[f(x+td')-f(x)\right]/t=<\nabla f(x),d> \quad \text{for all } d \in R^N. \qquad (2.10)$$

If

$$\lim_{y \to x, t \downarrow 0} \left[f(y+td)-f(y)\right]/t = <\nabla f(x), d> \quad \text{for all } d \text{ in } R^N , \qquad (2.11)$$

then f is called <u>strictly differentiable</u> at x. In this case f is diffe-
rentiable at x and the gradient $\nabla f:R^N \to R^N$ is continuous at x relative to
its domain

$$\text{dom } \nabla f = \{y \in R^N: f \text{ is differentiable at } y\}$$

It is known that a locally Lipschitzian function $f:R^N \to R$ is diffe-
rentiable at almost all points $x \in R^N$, and moreover that the gradient
mapping $\nabla f$ is locally bounded on its domain. Suppose that (2.2) holds
for some neighborhood B of a point $x \in R^N$. Then

$$<\nabla f(y),d> = f'(y;d) = \lim_{t \downarrow 0}\left[f(y+td)-f(y)\right]/t \le L|d|$$

for all $y \in B \cap \text{dom } \nabla f$ and $d \in R^N$, and this implies

$$|\nabla f(y)| \le L \text{ for all } y \in B \cap \text{dom } \nabla f. \qquad (2.12)$$

Since dom $\nabla f$ is dense in B, there exist sequences $\{y^j\}$ such that f is
differentiable at $y^j$ and $y^j \to x$. The corresponding sequence of gradients
$\{\nabla f(y^j)\}$ is bounded and has accumulation points (each being the limit
of some convergent subsequence). It follows that the set

$$M_f(x) = \{z \in R^N: \nabla f(y^j) \to z \text{ for some sequence } y^j \to z \text{ with f diffe-}$$
$$\text{rentiable at } y^j\} \qquad (2.13a)$$

is nonempty, bounded and closed. The set

$$\partial f(x) = \text{conv } M_f(x) \qquad (2.13b)$$

is called the _subdifferential_ of f at x (called the generalized gradient by Clarke (1975)). Each element $g_f \in \partial f(x)$ is called a _subgradient_ of f at x. Thus

$$\partial f(x) = \text{conv}\{\lim \nabla f(y^j) : y^j \to x, \quad f \text{ differentiable at } y^j\}. \tag{2.14}$$

In particular therefore, $\partial f(x) = -\partial(-f)(x)$. Three immediate consequences of the definition are listed in

_Lemma 2.2._ (i) $\partial f(x)$ is a nonempty convex compact set.
(ii)   The point-to-set mapping $\partial f(\cdot)$ is _locally bounded_ (bounded on bound-
        ed  subsets of $R^N$), i.e. if $B \subset R^N$ is bounded then the set
        $\{g_f \in \partial f(y) : y \in B\}$ is bounded.
(iii)  $\partial f(\cdot)$ is _upper semicontinuous_, i.e. if a sequence $\{y^j\}$ converges
        to x and $g_f^j \in \partial f(y^j)$ for each j then each accumulation point $g_f$ of
        $\{g_f^j\}$ satisfies $g_f \in \partial f(x)$.

        In general, $\partial f(x)$ does not reduce to $\nabla f(x)$ when the gradient $\nabla f$
is discontinuous at x.

_Lemma 2.3._ The following are equivalent:
(i)    $\partial f(x)$ consists of a single vector;
(ii)   $\nabla f(x)$ exists and $\nabla f$ is continuous at x relative to dom $\nabla f$;
(iii)  f is strictly differentiable at x.
Moreover, when these properties hold one has $\partial f(x) = \{\nabla f(x)\}$.

        Frequently $\partial f(x)$ is a singleton for almost every x. A locally Lip-
schitzian function $f : R^N \to R$ is _subdifferentially regular_ at $x \in R^N$ if for
every $d \in R^N$ the ordinary directional derivative (2.6) exists and coinci-
des with the generalized one in (2.3):

$$f'(x;d) = f^o(x;d) \quad \text{for all d.} \tag{2.15}$$

If (2.15) holds at each $x \in R^N$ then $\partial f(x)$ is actually single-valued at
almost every x. Below we give two important examples of subdifferential-
ly regular functions.

**Lemma 2.4.** If f is convex then f is subdifferentially regular and

$$f'(x;d) = \max\{<g_f,d> : g_f \in \partial f(x)\} \quad \text{for all } x,d. \tag{2.16}$$

**Lemma 2.5.** Suppose that

$$f(x) = \max\{f_u(x): u \in U\} \qquad \text{for all } x \in R^N, \tag{2.17}$$

where the index set U is a compact topological space (e.g. a finite set in the discrete topology), each $f_u$ is locally Lipschitzian, uniformly for u in U, and the mappings $f_u(x)$ and $\partial f_u(x)$ are upper semicontinuous in (x,u) (e.g. each $f_u$ is a differentiable function such that $f_u(x)$ and $\nabla f_u(x)$ depend continuously on (x,u)). Let

$$U(x) = \{u \in U: f_u(x) = f(x)\}. \tag{2.18}$$

Then f is locally Lipschitzian and

$$\partial f(x) \subset \text{conv} \{\partial f_u(x) : u \in U(x)\}. \tag{2.19}$$

If each $f_u$ is subdifferentially regular at x, then so is f, equality holds in (2.19), and

$$f'(x;d) = \max\{<g_u,d>: g_u \in \partial f_u(x), u \in U(x)\} \quad \text{for all } d. \tag{2.20}$$

**Corollary 2.6.** Suppose that

$$f(x) = \max\{f_i(x) : i \in I\} \qquad \text{for all } x \text{ in } R^N, \tag{2.21}$$

where the index set I is finite, and let $I(x)=\{i \in I: f_i(x)=f(x)\}$.
(i) If each $f_i$ is continuously differentiable then

$$f'(x;d) = \max\{<\nabla f_i(x),d> : i \in I(x)\} \qquad \text{for all } d,$$
$$\partial f(x) = \text{conv}\{\nabla f_i(x): i \in I(x)\}. \tag{2.22}$$

(ii) If each $f_i$ is convex then

$$f'(x;d) = \max\{<g_{f_i},d>: g_{f_i} \in \partial f_i(x), \quad i \in I(x)\} \quad \text{for all } d,$$
$$\partial f(x) = \text{conv}\{g_{f_i} \in \partial f_i(x): i \in I(x)\}. \tag{2.23}$$

When f is smooth, there exists an apparatus for computing $\nabla f$ in terms of the derivatives of other functions from which f is composed. The calculus of subgradients, which generalizes rules like $\nabla(f_1+f_2)(x) =$ $= \nabla f_1(x)+ \nabla f_2(x)$, is based on the following results.

**Lemma 2.7.** Let $g:R^{11}\to R$ and $h_i:R^N\to R$, $i=1,\ldots,n$, be locally Lipschitzian. Let $h(x)=(h_1(x),\ldots,h_n(x))$ and $(goh)(x)=g(h(x))$ for all $x\in E^N$. Then goh is locally Lipschitzian and

$$\partial(goh)(x) \subset conv \{ \sum_{i=1}^{n} u_i \partial h_i(x) : (u_1,\ldots,u_n)\in \partial g(h(x))\}. \qquad (2.24)$$

Moreover, equality holds in (2.24) if one of the following is satisfied:
(i) g is subdifferentially regular at $h(x)$, each $h_i$ is subdifferentially regular at x and $\partial g(h(x)) \subset R_+^n$ ($R_+^n =\{z\in R^n: z_i \geq 0$ for all i\});
(ii) g is subdifferentially regular at $h(x)$ and each $h_i$ is continuously differentiable at x;
(iii) Each $h_i$ is continuously differentiable at x and either g (or $-$ g) is subdifferentially regular at $h(x)$ or the Jacobian matrix of h at x is surjective;
(iv) n=1, g is continuously differentiable at $h(x)$ or g (or $-$ g) is subdifferentially regular at $h(x)$ and h is continuously differentiable at x. In cases (ii) $-$ (iv) the symbol "conv" is superfluous in (2.24). If (ii) holds then goh is subdifferentially regular at x.

**Corollary 2.8.** Suppose that $f_1$ and $f_2$ are locally Lipschitzian on $R^N$. For each $x\in R^N$ let $(f_1+f_2)(x)=f_1(x)+f_2(x)$, $(f_1f_2)(x)=f_1(x)f_2(x)$ and $(f_1/f_2)(x)=f_1(x)/f_2(x)$ if $f_2(x)\neq 0$. Then

$$\partial(f_1+f_2)(x) \subset \partial f_1(x)+\partial f_2(x), \qquad (2.25a)$$

$$\partial(f_1f_2)(x) \subset f_2(x) \, \partial f_1(x)+f_1(x)\partial f_2(x), \qquad (2.25b)$$

$$\partial(f_1/f_2)(x)\subset \frac{1}{(f_2(x))^2} [f_2(x)\partial f_1(x)-f_1(x)\partial f_2(x)]. \qquad (2.25c)$$

Equality holds in (2.25a) if each $f_i$ is subdifferentially regular at x, and in (2.25b) if in addition $f_i(x)\geq 0$.

Clarke (1975) established the following crucial relations between the subdifferential and the generalized directional derivatives of a lo-

cally Lipschitzian function f defined on $R^N$

$$f^O(x;d)= \max\{<g_f,d>:g_f \in \partial f(x)\} \text{ for all } x,d, \qquad (2.26)$$

$$\partial f(x)=\{g_f \in R^N: <g_f,d> \leq f^O(x;d) \text{ for all } d\} \text{ for all } x. \qquad (2.27)$$

We shall now interpret these relations in geometric terms. In what follows let x be a fixed point in $R^N$.

First, suppose that f is continuously differentiable at x. From Lemma 2.3, (2.26) and (2.8) we have

$$\partial f(x) = \{\nabla f(x)\}, \qquad (2.28a)$$

$$f^O(x;d) = f'(x;d) = <\nabla f(x),d> \quad \text{for all } d. \qquad (2.28b)$$

Suppose that $\nabla f(x) \neq 0$. Then $\nabla f(x)$ corresponds to the hyperplane

$$H_{\nabla f} = \{(z,\beta) \in R^{N+1}: \beta = f(x) + <\nabla f(x),z-x>\}$$

being tangent to the <u>graph</u> of f

$$\text{graph } f = \{(z,\beta) \in R^{N+1}:\beta = f(z)\}$$

at the point $(x,f(x))$. Here $\beta$ denotes the "vertical" coordinate of a point $(x,\beta) \in R^{N+1}$. Moreover, the hyperplane

$$H_C = \{z \in R^N: <\nabla f(x),z-x> = 0\}$$

is tangent at x to the <u>contour</u> of f at x

$$C = \{z \in R^N: f(z) = f(x)\}.$$

$\nabla f(x)$ is perpendicular to C at x and is the direction of steepest ascent for f at x. Define the following <u>linearization</u> of f at x

$$\bar{f}(z) = f(x) + <\nabla f(x),z-x> \quad \text{for all } z \text{ in } R^N \qquad (2.29)$$

and observe that $\nabla\bar{f}(z)=\nabla f(x)$ for all z (x is fixed). Therefore this linearization has the same differentiability properties as f at x in the sense that

$$\partial\bar{f}(x) = \partial f(x), \qquad (2.30a)$$

$$\bar{f}^O(x;d) = \bar{f}'(x;d) = f^O(x;d) \text{ for all } d, \qquad (2.30b)$$

cf. (2.28). In particular, by (2.28a), (2.9) and (2.30b), for any $d \in R^N$ we have