

Optimization →

given $f: \mathbb{R}^n \rightarrow \mathbb{R}$. find x s.t. $f(x) \leq f(\hat{x}) \forall \hat{x}$ near x .

review from last class

necessary conditions for x to be a minimizer

→ first order condition

$$\nabla f(x) = 0$$

→ second order condition

$\nabla^2 f(x)$ is positive semi-definite

↑
hessian

a sufficient condition (guarantee x is a local min)

$$\nabla f(x) = 0$$

$\nabla^2 f(x)$ positive-definite.

in fact, this $\Rightarrow f(x) < f(\hat{x}) \forall \hat{x}$ near x

Algorithms for finding minimizers:

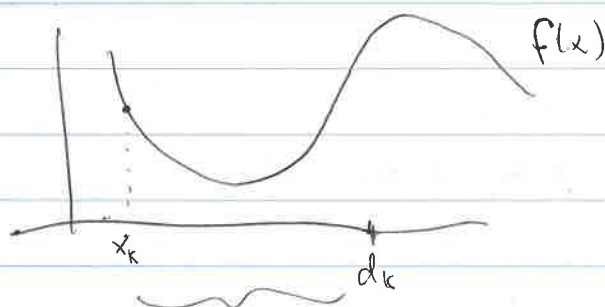
- generate x_1, x_2, \dots (starting w/ x_0)
- given x_k , find d_k (descent direction)
s.t. $(\nabla f)(x_k)^T d_k < 0$

- update $x_{k+1} = x_k + \alpha_k d_k$ for some $\alpha_k > 0$

so, 3 issues we need to work on:

1. choice of direction d_k
2. choice of scalar α_k .

let's consider α_k (interpreted as step length) for a moment: 1-d example.



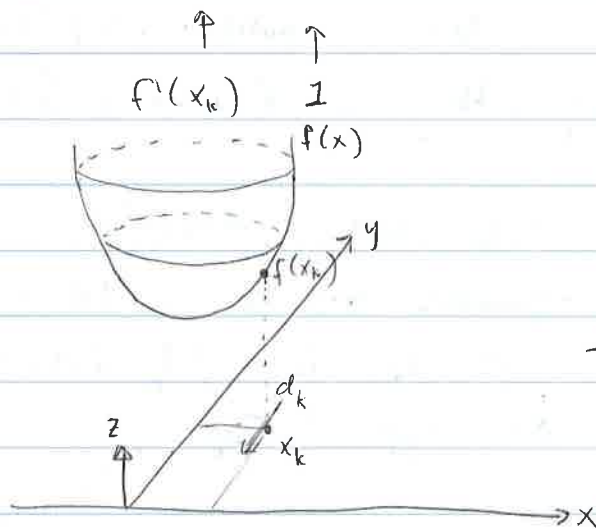
Want to find direction to move that will give us a smaller f^* value.

Suppose $|d_k| = 1$. also for x_k as shown,

$$\nabla f(x_k) = f'(x_k) < 0$$

$$\text{for } d_k = 1, f'(x_k)^T d_k < 0$$

2D example



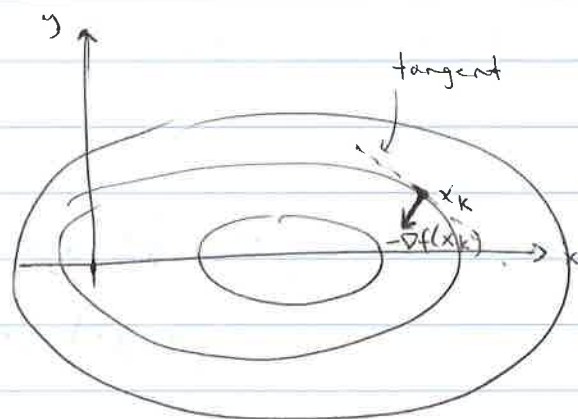
know

$$\nabla f(x_k)^T (-\nabla f(x_k)) < 0$$

$-\nabla f(x_k)$ is a

descent direction.

Contours

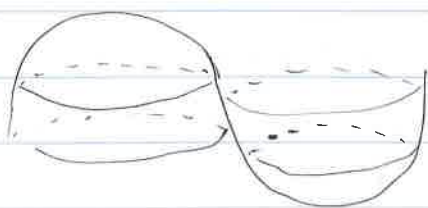


Contour Lines

$$\Rightarrow f(\vec{x}) = \text{constant}.$$

on the contour plot: $\nabla f(x_k) \perp$ tangent line
to the contour through x_k .

but, Suppose it later turned out the ~~bowl~~ bowl in the
graph turned over its surface, e.g.



then the step length
 d_k could make a
big difference.

So far, our only tool (approach) is to choose the
negative gradient to go from there. but
it's not the best way.

ie, concerning item #1, one possible descent direction
is $d_k = -\nabla f(x_k)$.

For another approach, consider the Taylor series
for f .

$$f(x_k + d) = f(x_k) + \nabla f(x_k)^T d + \frac{1}{2} d^T \nabla^2 f(x_k) d + O(\|d\|^3) \text{ for small } d.$$

call this region of
interest $\Phi(d)$ s.t. $\Phi: \mathbb{R}^n \rightarrow \mathbb{R}$.

$$\nabla \Phi = (\nabla f)(x_k) + \underbrace{\left[\nabla^2 f(x_k) \right]}_{H_k} d \rightarrow \text{set} = 0$$

want to relate this back to the first order
condition that $\nabla f(x) = 0$.

$$\Rightarrow d = -H_k^{-1} (\nabla f) x_k$$

(but we never multiply by inverse - instead
we solve $H_k d = -\nabla f(x_k)$ for d)

$$\nabla_d^2 \Phi = H_k \quad (\text{b/c } \nabla \Phi \text{ ~~at~~ on is linear f' of d})$$

We know that d minimizes Φ if H_k is positive definite.

This strategy is known as Newton's method.

So - do we choose Newton's or steepest descent method?

Newton

expensive

Could lead you astray if

H_k is not

positive definite

much faster

steepest descent

cheap

slow

safe (guaranteed to work)

Strategy: start w/ something safe & then switch in a specific region of interest.

Convergence note: x_k converges to x with rate r

$$\text{if } \lim_{k \rightarrow \infty} \frac{\|x - x_{k+1}\|}{\|x - x_k\|^r} = c \quad \text{constant, independent of } k.$$

ie, concerning item #1, one possible descent direction
is $d_k = -\nabla f(x_k)$

For another approach, consider the Taylor series
for f .

$$f(x_k + d) = f(x_k) + \nabla f(x_k)^T d + \frac{1}{2} d^T \nabla^2 f(x_k) d + O(\|d\|^3) \text{ for small } d.$$

call this region of
interest $\Phi(d)$ s.t. $\Phi: \mathbb{R}^n \rightarrow \mathbb{R}$.

$$\nabla \Phi = (\nabla f)(x_k) + \underbrace{\left[\nabla^2 f(x_k) \right]}_{H_k} d \rightarrow \text{set} = 0$$

would to relate this back to the first order
condition that $\nabla f(x) = 0$.

$$\Rightarrow d = -H_k^{-1} (\nabla f) x_k$$

(but we never multiply by inverse - instead

we solve $H_k d = -\nabla f(x_k)$ for d)

$$\nabla_d^2 \Phi = H_k \quad (\text{b/c } \nabla \Phi \text{ ~~is~~ on is linear} \\ \text{f}^\circ \text{ of } d)$$

We know that d minimizes Φ if H_k is positive definite.

This strategy is known as Newton's method.

So - do we choose Newton's or steepest descent method?

Newton

expensive

could lead you astray if

H_k is not

positive definite

much faster

steepest descent

cheap

slow

safe (guaranteed to work)

Strategy: start w/ something safe & then switch in a specific region of interest.

Convergence note: x_k converges to x with rate r

$$\text{if } \lim_{k \rightarrow \infty} \frac{\|x - x_{k+1}\|}{\|x - x_k\|^r} = c \quad \text{constant, independent of } k.$$

3 types of rates we are concerned about:

$r=1$ linear convergence, in this case we need $c < 1$ (* true for steepest descent).

$r=2$ quadratic rate of convergence. don't need $c < 1$, but we need denominator to not affect c . (* this is true for Newton's method)

Suppose $c = 1/2$ for steepest descent.

Suppose $c = 2$ for Newton.

$c = 2$ / Newton ← Compare w/ steepest descent
much slower convergence

0	$1/4$
1	$2 \cdot 1/16 = 1/8$
2	$2 \cdot 1/64 = 1/32$
3	$2 \cdot 1/32^2 = 2 \cdot 1/1024 = 1/512$
4	$2 \cdot 1/2^{11} = 1/2^4$

One more thing: it's possible to have $1 < r < 2$ - this is superlinear convergence. obtained by combining Newton & other methods - "quasi-Newton" methods. saves some overhead.