# UNCONSTRAINED OPTIMIZATION
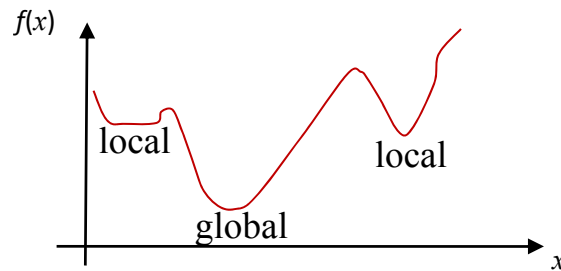
## (USING DERIVATIVES)

Consider the problems <u>Minimize $f(x)$, st $x \in \mathbb{R}^n$.</u>

<u>DEFINITION 1</u>: A point $x^* \in \mathbb{R}^n$ is said to be a *LOCAL MINIMUM POINT* of $f$ over

$\mathbb{R}^n$ if there exists $\varepsilon > 0$ such that $f(x) \geq f(x^*)$ for all $x$ such that $\|x - x^*\| < \varepsilon$ (i.e., within a

distance $\varepsilon$ of $x$). If $f(x) > f(x^*)$ then it is *STRICT* local minimum point.

<u>DEFINITION 2</u>: A point $x^* \in \mathbb{R}^n$ is said to be a *GLOBAL MINIMUM POINT* of $f$ over

$\mathbb{R}^n$ if $f(x) \geq f(x^*)$ for <u>all</u> $x \in \mathbb{R}^n$. If $f(x) > f(x^*)$ then it is a *STRICT* global minimum point.



Assuming that $f$ is differentiable, the gradient vector $\nabla f(x)$ is the vector

$$\nabla f(x) = \begin{bmatrix} \dfrac{\partial f}{\partial x_1}(x) \\ \dfrac{\partial f}{\partial x_2}(x) \\ \vdots \\ \dfrac{\partial f}{\partial x_n}(x) \end{bmatrix}$$

<u>THEOREM 1</u>: **Necessary Conditions**

If $f$ is smooth real-valued function then a <u>necessary</u> condition for the point $x^*$ to be a local minimizer (or maximizer) is that $\nabla f(x^*) = \mathbf{0}$, i.e.,

$$\begin{cases} \frac{\partial f}{\partial x_1}(x^*) = 0 \\ \frac{\partial f}{\partial x_2}(x^*) = 0 \\ \quad \vdots \\ \frac{\partial f}{\partial x_n}(x^*) = 0 \end{cases} \qquad \text{(in general, a nonlinear system of } n \text{ equations in } n \text{ unknowns...)}$$

Furthermore, the Hessian $\nabla^2 f(x^*)$ is positive (negative) semidefinite at a minimum (maximum).

<u>THEOREM 2</u>: **Sufficient Conditions**

Assuming that $f$ is twice continuously differentiable,

If        1) $\nabla f(x^*) = \mathbf{0}$                and,

          2) $y^T \nabla^2 f(x^*) y > 0$ for all $y \neq \mathbf{0}$

then $x^*$ is a local minimizer.

NOTE:

1.  Condition 2 states that the Hessian at $x^*$ should be positive definite

2.  If the gradient vanishes and the Hessian is <u>negative</u> definite at $x^*$, then $x^*$ is a local maximum

Note that there may be points other than $x^*$ that satisfy the <u>necessary</u> conditions; however, these will not satisfy the <u>sufficient</u> conditions unless they are local optima. (e.g. $f(x)=x^3$ at $x=0$).

EXAMPLE: Consider the function $f(x,y,z) = 2x^2 + xy + y^2 + yz + z^2 - 6x - 7y - 8z + 9$

*Points satisfying the **necessary** conditions ($[\nabla f(x)] = [0]$)*

$$\partial f/\partial x = 4x + y - 6 = 0$$

$$\partial f/\partial y = x + 2y + z - 7 = 0$$

$$\partial f/\partial z = y + 2z - 8 = 0$$

which yields         $[x=6/5; y=6/5; z=17/5]$

*Do these satisfy the **sufficient** conditions?*

$$H = \nabla^2 f(x^*) = \begin{bmatrix} 4 & 1 & 0 \\ 1 & 2 & 1 \\ 0 & 1 & 2 \end{bmatrix}$$

$$|4| = 4 > 0; \qquad \begin{vmatrix} 4 & 1 \\ 1 & 2 \end{vmatrix} = 7 > 0; \qquad \begin{vmatrix} 4 & 1 & 0 \\ 1 & 2 & 1 \\ 0 & 1 & 2 \end{vmatrix} = 10 > 0$$

Since all principal minors are strictly positive, it follows that $H$ is positive definite

over $\mathbb{R}^3$ and hence the sufficient conditions are satisfied and $x^*$ is a local minimum.

Is the point also a **global minimum**?

By the Taylor series expansion around $x^*$ we have

$$f(x^* + d) = f(x^*) + \nabla f(x^*)^T d + \tfrac{1}{2} d^T \nabla^2 f(x) d,$$

But $\nabla^2 f(x)$ is positive definite and hence, for all $d \neq 0$ we have $\tfrac{1}{2} d^T \nabla^2 f(x) d > 0$.

Further, $\nabla f(x^*) = 0 \Rightarrow f(x^* + d) - f(x^*) > 0$

$$\Rightarrow f(x^* + d) > f(x^*) \text{ for all } d \neq 0.$$

Therefore $x^*$ is also a **global** minimum.

Recall that a **convex** function has a positive semidefinite Hessian. Thus (along the same lines as the example above) the necessary conditions for a local minimum are also sufficient. Furthermore, the local minimum is also a **global** minimum. If the function is **strictly** convex then this global minimum is also unique. (WHY?). Similar remarks hold for **concave** functions and their **maxima**.

Most techniques are designed to find local minima (maxima); these of course, being adequate for convex (concave) functions.

# CONVERGENCE OF ALGORITHMS

Before getting into algorithms, we briefly examine the general conditions under which an algorithm will converge to an optimum solution.

<u>DEFINITION 1</u>: An algorithm $A$ is a point-to-point mapping defined on a space $X$ that assigns to every point $x \in X$, a point $y \in X$.

(An algorithm may also be viewed more generally as a point-to-**set** mapping; however we'll keep it simple here…). Thus an algorithm $A$ generates a sequence of points; i.e., given a starting point $x^0$, an iterative sequence is generated from

$$x^{k+1} = A(x^k)$$

# An Example

Minimize $x$-5, st $x \geq 0$

The optimum solution is at $x^* = 0$

One algorithm might be defined by the following point-to-point map:

$$x^{k+1} = A(x^k) = x^k/2$$

So for $x^0 = 5$ the sequence generated would be {5, 2.5, 1.25, 0.625, ...}

Clearly, this converges to the optimum

Another algorithm might use the point-to-set map:

$$x^{k+1} = A(x^k) = [x^k/4, x^k/2]$$

i.e., each new point is somewhere within the closed interval above. The exact

sequence cannot be predicted now, but starting at the same $x^0 = 5$ we might generate

the following sequence: {5, 2, 0.8, 0.36, 0.10,...}. Again, this is a sequence that

converges to $x^*$ (and actually, at least as fast the previous method).

Finally, consider the following algorithm:

$$x^{k+1} = A(x^k) = \begin{cases} (x^k/2) & if \ x^k < 1 \\ 1 + (x^k - 1)/2 & if \ x^k \geq 1 \end{cases}$$

For $x^0 < 1$ this converges to $x^*$ but for $x^0 \geq 1$ it converges to the sub-optimal point $x=1$

Associated with $A$ are (i) **a solution set**, and (ii) a **descent** (or **ascent**) **function**.

DEFINITION 2: A solution set $\Omega$ is defined as

$\quad\quad \Omega=\{x|\, x$ is a solution to our problem$\}$.

For instance consider the problem: Minimize $f(x)$, st $x\in\mathbb{R}^n$

Among other possibilities, we could define $\Omega$ as

1) $\Omega = \{x\in\mathbb{R}^n \,|\, x$ is a local minimum point of $f\}$,

2) $\Omega = \{x\in\mathbb{R}^n \,|\, \nabla f(x) = \mathbf{0}\}$

3) $\Omega = \{x\in\mathbb{R}^n \,|\, f(x) \le b\}$, where $b$ is some acceptable value

4) $\Omega = \{x\in\mathbb{R}^n \,|\, \nabla f(x) = \mathbf{0}$ and $\nabla^2 f(x)$ is positive definite$\}$.    etc. etc.


DEFINITION 3: Let $\Omega \subset X$ be a given solution set and let $A$ be an algorithm on $X$.

A continuous real valued function $z$ on $X$ is said to be a DESCENT FUNCTION for

$\Omega$ and $A$ if it satisfies

(i)    if $x\notin\Omega$ and $y= A(x)$, then $z(y)< z(x)$

(ii)    if $x\in\Omega$ and $y= A(x)$, then $z(y)\le z(x)$

In other words, a descent function decreases in value as each new point is generated

by $A$, and once an acceptable solution is found new points can never be any worse.

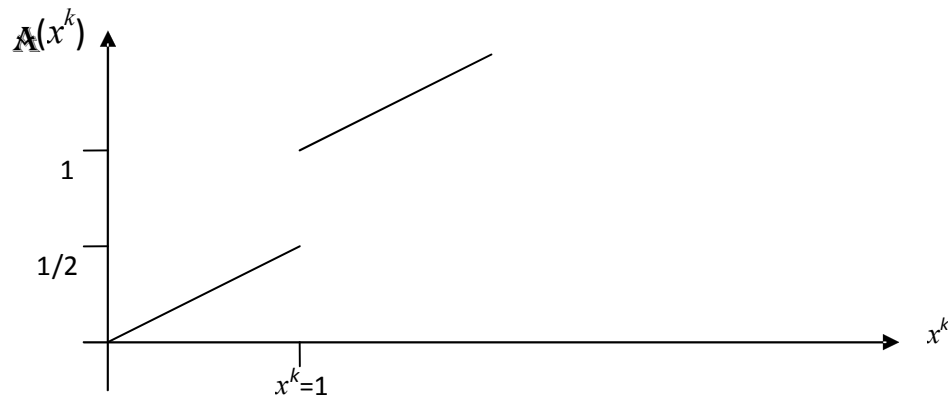For example, if $\Omega$ is defined as in 2), we could use $\|\nabla f(x)\|$ as a descent function.

For $\Omega$ as in 3) we could simply use $f(x)$ itself.

## Closed Maps

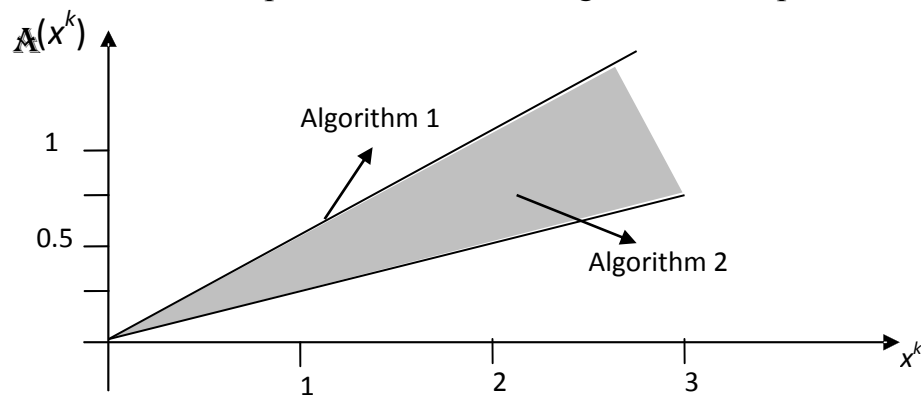Another important concept is continuity (or closedness) of algorithmic maps.

<u>DEFINITION 4</u>:  An algorithm $\mathbb{A}$ on $\mathbb{X}$ is closed (or continuous) at $\boldsymbol{x} \in \mathbb{X}$, if $\{\boldsymbol{x}^k\} \to \boldsymbol{x}$

implies that $\{\mathbb{A}(\boldsymbol{x}^k)\} \to \mathbb{A}(\boldsymbol{x})$.  The map (or algorithm) is closed on $\mathbb{Z} \subset \mathbb{X}$ if it is

closed at each point in $\mathbb{Z}$.

Consider the map for the third algorithm we saw earlier: As shown below, the map is

discontinuous at $x^k = 1$. This is why it doesn't converge to the optimum



For example, the sequence with $x^k = 1 - (1/k)$ converges to $x = 1$ but the sequence $\{\mathbb{A}(x^k)\}$

with $\mathbb{A}(x^k) = x^k/2$ converges to the value 1/2, which is not $\mathbb{A}(x) = \mathbb{A}(1) = 1$.

On the other hand, the maps for the first two algorithmic maps are both closed

Note that closedness of algorithmic maps is analogous to continuity in functions.

We may now state the main theorem of this section:

**CONVERGENCE THEOREM**: Let $A$ be an algorithm on $X$ and suppose that given $x^0$, the sequence $\{x^k\}_{k=0...\infty}$ is generated via $x^{k+1} = A(x^k)$. Let the optimal solution set $\Omega \subset X$ be given, and suppose that,

a) all points $x^k$ are contained in a compact set $S \subset X$

b) there exists a continuous descent function $z$ for $A$ and $\Omega$

c) the map of algorithm $A$ is closed at points outside $\Omega$.

Then the sequence $\{x^k\}_{k=0...\infty}$ converges to a point in $\Omega$, or the limit of any convergent subsequence of $\{x^k\}$ is a solution, i.e., the limit lies within $\Omega$. Also, $z(x^k)$ $\rightarrow z(x)$ for some $x \in \Omega$.

Note that the convergent subsequence may be the sequence itself. If not, at least some subset of the points generated converges to a solution.

NOTES ON THE CONVERGENCE THEOREM:

1.  It states "sufficient" conditions for convergence.  It might be possible to prove convergence under different (perhaps weaker) assumptions.  However, the notions of a continuous algorithmic map and of a descent function form the backbone of all proofs.

2.  By convergence it means that the algorithm will converge to a solution irrespective of the starting point $x^0$. (This is not to be confused with trying to find a global minimum).   Thus proof of convergence is an indication of robustness.

3.  In most engineering applications, the design variables are usually bounded as $x_{lower} \leq x \leq x_{upper}$.  This ensures that $\{x^k\}$ is generated in a compact set as per requirement (a).

4.  If $\Omega$ is the singleton $x^*$ then the whole sequence $\{x^k\}$ converges to $x^*$.

5.  The theorem is very useful for modifying existing algorithms so as to make them convergent.

# RATES OF CONVERGENCE

This refers to the "speed" with which convergence takes place.

Suppose $\{x^k\} \to x^*$. We say that $\{x^k\}$ converges to $x^*$ with **order** (or **rate**) $r$ and a **rate constant** $C$, if $r$ is the supremum of the nonnegative values satisfying

$$0 \le \lim_{k \to \infty} \left( \frac{\left\| x^{k+1} - x^* \right\|}{\left\| x^k - x^* \right\|^r} \right) = C < \infty$$

e.g. if $(x^k - x^*)$ is the "error" $e^k$ at iteration $k$ (with $\lim_{k \to \infty} e^k = 0$), then this says

$$\lim_{k \to \infty} \left( \frac{\left\| e^{k+1} \right\|}{\left\| e^k \right\|^r} \right) = C.$$

- If $r=1$, $C$ must be $<1$ for convergence; if $0<C<1$ the sequence is then said to have a *linear* convergence rate.

- If $r=1$, but $C=0$ then the sequence is said to have *superlinear* convergence.

THEOREM: If $r>1$ then convergence is superlinear.

Proof: Suppose we have $\lim_{k \to \infty} \left( \frac{\left\| e^{k+1} \right\|}{\left\| e^k \right\|^r} \right) = C$ where $C < \infty$ and $r>1$.

Then $\lim_{k \to \infty} \left( \frac{\left\| e^{k+1} \right\|}{\left\| e^k \right\|} \right) = \lim_{k \to \infty} \left( \frac{\left\| e^{k+1} \right\|}{\left\| e^k \right\|^r} \left\| e^k \right\|^{r-1} \right) = C \cdot 0 = 0$

In particular, if $r=2$ and $C>0$ the sequence is said to have quadratic convergence; it will eventually always converge faster than a linearly convergent one…

EXAMPLE 1: Suppose $\left(\dfrac{\left\|e^{k+1}\right\|}{\left\|e^{k}\right\|^{r}}\right)=C$ for all $k$, i.e., $\left\|e^{k+1}\right\|=C\left\|e^{k}\right\|^{r}$. Note that with this

simplification, we don't need to worry about limits…

Case 1: $r=1$, $0<C<1$ (note that if $C>1$, then the sequence diverges!).

So if we start with $\left\|e^{0}\right\|=1$ and $C=10^{-1}$, then we have $\left\{\left\|e^{k}\right\|\right\}=\{1,\ 10^{-1},\ 10^{-2},\ 10^{-3},…\}$.

If $C=0.99$ then we have $\left\{\left\|e^{k}\right\|\right\}=\{1,\ .99,\ .9801,\ .9703,…\}$ – not so good!

Case 2: $r=2$

If we start with $\left\|e^{0}\right\|=1$ and $C=10^{-1}$, then $\left\{\left\|e^{k}\right\|\right\}=\{1,\ 10^{-1},\ 10^{-3},\ 10^{-7},…\}$.

If $C=0.99$ then we have $\left\{\left\|e^{k}\right\|\right\}=\{1,\ .99,\ .9703,\ 0.932,\ 0.86…\}$

Note that convergence is much faster: for $C=10^{-1}$ we can get the error to $10^{-7}$ in 4 steps now as opposed to eight with linear convergence!

With linear convergence rates the reduction $\left\|x^{k}-x^{*}\right\|$ at each step will depend significantly on the rate constant $C$. For example, when $C$ is small linear convergence can compare favorably with superlinear convergence.

EXAMPLE 2: Consider the sequence $\left\{ x^{\left(2^{-k}\right)} \right\}$ where $x \geq 0$ (each element is the square

root of the previous element). The limit of this sequence is 1.

If $r=1$ then $\lim\limits_{k \to \infty} \left( \dfrac{\left\| e^{k+1} \right\|}{\left\| e^{k} \right\|} \right) = \lim\limits_{k \to \infty} \left( \dfrac{x^{2^{-(k+1)}} - 1}{x^{2^{-k}} - 1} \right) = 0.5$. Thus the sequence has a

linear convergence rate with rate constant $C=0.5$

EXAMPLE 3: Consider the sequence $\{x^k\}$ with $x^k = a^k$ where $0<a<1$. The limit of

this sequence is 0.

If $r=1$ then $\lim\limits_{k \to \infty} \left( \dfrac{\left\| e^{k+1} \right\|}{\left\| e^{k} \right\|} \right) = \lim\limits_{k \to \infty} \left( \dfrac{a^{k+1} - 0}{a^{k} - 0} \right) = a$. Thus the sequence has a linear

convergence rate with rate constant $a$.

EXAMPLE 4: Consider the sequence $\left\{ x^{\left(2^{k}\right)} \right\}$ where $0 \leq x < 1$ (each element is the

square of the previous element). The limit of this sequence is 0.

If $r=2$ then $\lim\limits_{k \to \infty} \left( \dfrac{\left\| e^{k+1} \right\|}{\left\| e^{k} \right\|^2} \right) = \lim\limits_{k \to \infty} \left( \dfrac{x^{2^{(k+1)}}}{\left( x^{2^{k}} \right)^2} \right) = 1$. Thus the sequence has a

quadratic (and hence, superlinear) convergence rate with rate constant $C=1$

## Stopping Criteria

The algorithm $\mathbb{A}$ terminates when we reach a point in $\Omega$. Usually, convergence to a point in $\Omega$ occurs only in a limiting sense, and we need practical stopping criteria. Some common examples are,

1. Stop if the distance moved is less than $\varepsilon$ after $N$ successive iterations, that is,

   $$\|x^{k+N} - x^k\| < \varepsilon.$$

2. Stop if the change in $x$ from one iteration to the next is less than $\varepsilon\%$, i.e.,

   $$\frac{\|x^{k+1} - x^k\|}{\|x^k\|} < \varepsilon$$

3. Stop when the improvement in the descent function $z$ after $N$ iterations is less than $\varepsilon$, i.e., $z(x^k) - z(x^{k+N}) < \varepsilon$.

4. Stop if the improvement in the descent function from one iteration to the next is less than $\varepsilon\%$, i.e., $\dfrac{z(x^k) - z(x^{k+1})}{|z(x^k)|} < \varepsilon.$

5. If $z(x^*)$ is known in advance then stop when we are less than $\varepsilon$ units from this, i.e.,

   $$|z(x^k) - z(x^*)| < \varepsilon.$$

# ALGORITHMS FOR UNCONSTRAINED OPTIMIZATION

We have already seen simple search methods for unimodal functions. We now look at other methods.

## NEWTON'S METHOD

Looks for a point satisfying the necessary conditions, i.e., it tries to solve $\nabla f(x) = 0$; a system of nonlinear equations. To do this it uses the Newton-Raphson method for each equation $\partial f / \partial x_i = 0$ for each $i$,

Recall that to solve $g(x) = 0$ by the Newton-Raphson method we used the Taylor series expansion around $x^k$ at iteration $k$. Now consider a generalization to solve

$$g(x) = g(x_1, x_2, ..., x_n) = 0$$

The Taylor series expansion around $x^k$ is

$$g(x^k + \Delta x) = g(x^k) + [\nabla g(x^k)]^T [\Delta x] + \tfrac{1}{2} [\Delta x]^T \nabla^2 g(x^k) [\Delta x] + ...$$

$$\underbrace{\hspace{6cm}}$$

(negligible terms…)

To choose $\Delta x$ so that $g(x^k + \Delta x) \cong 0$ and then move from the current point $x^k$ to a better point $x^{k+1} = x^k + \Delta x$, we thus want $g(x^{k+1}) = g(x^k + \Delta x) = g(x^k) + [\nabla g(x^k)]^T [\Delta x] \cong 0$

Consider now when we have $n$ simultaneous nonlinear equations $g_1 = g_2 = g_n = 0$. With each $i$ we have the following system:

$$[\nabla g_1(x^k)]^T [\Delta x] = -g_1(x^k)$$

$$[\nabla g_2(x^k)]^T [\Delta x] = -g_2(x^k)$$

$$\vdots$$

$$[\nabla g_n(x^k)]^T [\Delta x] = -g_n(x^k)$$

where $\Delta x = \begin{bmatrix} \Delta x_1 \\ \Delta x_2 \\ \vdots \\ \Delta x_n \end{bmatrix}$. Thus we solve $J\Delta x = -g$ to obtain $\Delta x = -J^{-1}g$, i.e.,

$$\begin{bmatrix} \Delta x_1 \\ \Delta x_2 \\ \vdots \\ \Delta x_n \end{bmatrix} = - \begin{bmatrix} \partial g_1/\partial x_1 & \partial g_1/\partial x_2 & \dots & \partial g_1/\partial x_n \\ \partial g_2/\partial x_1 & \partial g_2/\partial x_2 & \dots & \partial g_2/\partial x_n \\ \vdots & \vdots & & \vdots \\ \partial g_n/\partial x_1 & \partial g_n/\partial x_2 & \dots & \partial g_n/\partial x_n \end{bmatrix}^{-1} \begin{bmatrix} g_1(x^k) \\ g_2(x^k) \\ \vdots \\ g_n(x^k) \end{bmatrix}$$

where each partial derivative above is evaluated at $x^k$,

Thus we have $\boxed{\Delta x = -[J(x^k)]^{-1} [g(x^k)]}$ where the matrix $J$ is called the JACOBIAN

matrix and $J(x^k)$ indicates that it is evaluated at the point $x^k$. The $x^{k+1} = x^k + \Delta x$

Let us apply this to solve $\nabla f(x) = 0$, i.e.,

$$g_1(x) \equiv \frac{\partial f}{\partial x_1}(x) = 0$$
$$g_2(x) \equiv \frac{\partial f}{\partial x_2}(x) = 0$$
$$\vdots$$
$$\vdots$$
$$g_n(x) \equiv \frac{\partial f}{\partial x_n}(x) = 0$$

Thus we have $g(x) \equiv \nabla f(x)$; [i.e., $g_i(x) = \frac{\partial f}{\partial x_i}(x)$]

$$\implies \quad \frac{\partial g_i}{\partial x_j}(x) \equiv \frac{\partial^2 f}{\partial x_j \partial x_i}(x)$$

So the Jacobian matrix for $g(x)$ is

$$
\begin{bmatrix}
\dfrac{\partial^2 f}{\partial x_1^2} & \dfrac{\partial^2 f}{\partial x_1 \partial x_2} & \cdots & \dfrac{\partial^2 f}{\partial x_1 \partial x_n} \\[2mm]
\dfrac{\partial^2 f}{\partial x_2 \partial x_1} & \dfrac{\partial^2 f}{\partial x_2^2} & \cdots & \dfrac{\partial^2 f}{\partial x_2 \partial x_n} \\[2mm]
\vdots & \vdots & \vdots & \vdots \\[1mm]
\vdots & \vdots & \vdots & \vdots \\[1mm]
\dfrac{\partial^2 f}{\partial x_n \partial x_1} & \dfrac{\partial^2 f}{\partial x_n \partial x_2} & \cdots & \dfrac{\partial^2 f}{\partial x_n^2}
\end{bmatrix} = \nabla^2 f
$$

Notice that this is nothing but the Hessian matrix ($H$) of the function $f$!
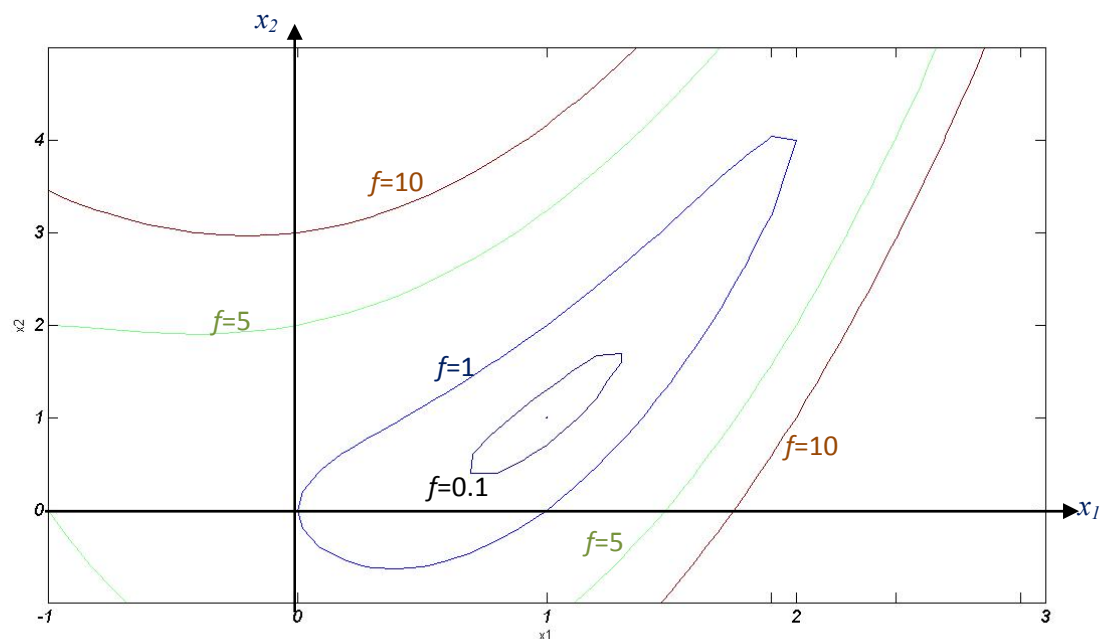
Thus if $x^i$ is the guess at iteration $i$

then $\Delta x = -[J(x^i)]^{-1} [g(x^i)]$ implies that

$$\Delta x = -[\nabla^2 f(x^i)]^{-1} [\nabla f(x^i)] = -H^{-1} \nabla f$$

and $x^{i+1} = x^i + \Delta x$ is the improved guess at iteration $i+1$.

AN EXAMPLE:  Minimize $f(x_1, x_2) = (x_2 - x_1^2)^2 + (1-x_1)^2$

**Contours of $f(x)$ for various values**

For this problem $x^* = [1,1]$ and $f(x^*)=0$,

$$\nabla f(x) = \begin{bmatrix} \partial f/\partial x_1 \\ \partial f/\partial x_2 \end{bmatrix} = \begin{bmatrix} 4x_1^3 - 4x_1 x_2 + 2x_1 - 2 \\ 2(x_2 - x_1^2) \end{bmatrix}$$

and

$$\nabla^2 f(x) = \begin{bmatrix} \dfrac{\partial^2 f}{\partial x_1^2} & \dfrac{\partial^2 f}{\partial x_1 \, \partial x_2} \\ \dfrac{\partial^2 f}{\partial x_1 \, \partial x_2} & \dfrac{\partial^2 f}{\partial x_2^2} \end{bmatrix} = \begin{bmatrix} 12x_1^2 - 4x_2 + 2 & -4x_1 \\ -4x_1 & 2 \end{bmatrix}$$

<u>ITERATION 1</u>

Start with $x^0 = \begin{bmatrix} 0 \\ 0 \end{bmatrix}$ where $f(x^0) = 1$

$$\nabla f(x^0) = \begin{bmatrix} -2 \\ 0 \end{bmatrix}, \qquad H = \nabla^2 f(x^0) = \begin{bmatrix} 2 & 0 \\ 0 & 2 \end{bmatrix} \implies H^{-1} = \begin{bmatrix} 1/2 & 0 \\ 0 & 1/2 \end{bmatrix}$$

Therefore $\Delta x = -H^{-1}\nabla f = -\begin{bmatrix} 1/2 & 0 \\ 0 & 1/2 \end{bmatrix}\begin{bmatrix} -2 \\ 0 \end{bmatrix} = \begin{bmatrix} 1 \\ 0 \end{bmatrix}$

and $x^1 = x^0 + \Delta x = \begin{bmatrix} 0 \\ 0 \end{bmatrix} + \begin{bmatrix} 1 \\ 0 \end{bmatrix} = \begin{bmatrix} 1 \\ 0 \end{bmatrix}$

<u>ITERATION 2</u>         $[f(x^1) = 1]$

$$\nabla f(x^1) = \begin{bmatrix} 4 \\ -2 \end{bmatrix}, \qquad H = \nabla^2 f(x^1) = \begin{bmatrix} 14 & -4 \\ -4 & 2 \end{bmatrix} \implies H^{-1} = \begin{bmatrix} 1/6 & 1/3 \\ 1/3 & 7/6 \end{bmatrix}$$

$\Delta x = -H^{-1}\nabla f = -\begin{bmatrix} 1/6 & 1/3 \\ 1/3 & 7/6 \end{bmatrix}\begin{bmatrix} 4 \\ -2 \end{bmatrix} = \begin{bmatrix} 0 \\ 1 \end{bmatrix}$

and $x^2 = x^1 + \Delta x = \begin{bmatrix} 1 \\ 0 \end{bmatrix} + \begin{bmatrix} 0 \\ 1 \end{bmatrix} = \begin{bmatrix} 1 \\ 1 \end{bmatrix}$

<u>ITERATION 3</u>           $[f(x^2) = 0]$

$\nabla f(x^2) = \begin{bmatrix} 0 \\ 0 \end{bmatrix} \implies \begin{bmatrix} 1 \\ 1 \end{bmatrix}$ is a stationary point.           **STOP**

Suppose we had used a different starting point say $x^0 = \begin{bmatrix} 2 \\ 1 \end{bmatrix}$ where $f(x^0)=10$

ITERATION 1        $[f(x^0)=10]$

$$\nabla f(x^0) = \begin{bmatrix} 26 \\ -6 \end{bmatrix}, \qquad H = \begin{bmatrix} 46 & -8 \\ -8 & 2 \end{bmatrix} \Rightarrow H^{-1} = \begin{bmatrix} 0.0714 & 0.2857 \\ 0.2857 & 1.6428 \end{bmatrix}$$

$$\Delta x = -H^{-1}\nabla f = -\begin{bmatrix} 0.0714 & 0.2857 \\ 0.2857 & 1.6428 \end{bmatrix}\begin{bmatrix} 26 \\ -6 \end{bmatrix} = \begin{bmatrix} -0.1428 \\ 2.4286 \end{bmatrix}$$

and $x^2 = x^1 + \Delta x = \begin{bmatrix} 2 \\ 1 \end{bmatrix} + \begin{bmatrix} -0.1428 \\ 2.4286 \end{bmatrix} = \begin{bmatrix} 1.8572 \\ 3.4286 \end{bmatrix}$
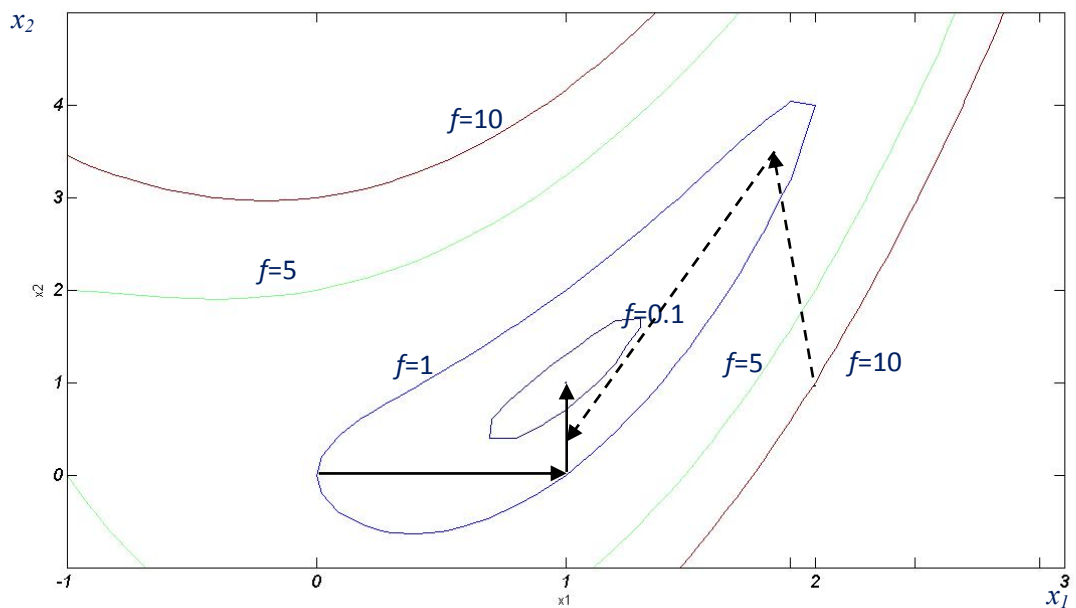
ITERATION 2        $[f(x^1)=0.735]$

$$\nabla f(x^1) = \begin{bmatrix} 1.8569 \\ -0.0408 \end{bmatrix} \quad H = \begin{bmatrix} 29.6735 & -7.4286 \\ -7.4286 & 2 \end{bmatrix} \Rightarrow H^{-1} = \begin{bmatrix} 0.4804 & 1.7843 \\ 1.7843 & 7.1275 \end{bmatrix}$$

$$\Delta x = -H^{-1}\nabla f = -\begin{bmatrix} 0.4804 & 1.7843 \\ 1.7843 & 7.1275 \end{bmatrix}\begin{bmatrix} 1.8569 \\ -0.0408 \end{bmatrix} = \begin{bmatrix} -0.8235 \\ -3.0384 \end{bmatrix}$$

and $x^2 = x^1 + \Delta x = \begin{bmatrix} 1.8572 \\ 3.4286 \end{bmatrix} + \begin{bmatrix} -0.8235 \\ -3.0384 \end{bmatrix} = \begin{bmatrix} 1.0336 \\ 0.3902 \end{bmatrix}$

Worked well with $\begin{bmatrix} 0 \\ 0 \end{bmatrix}$ as the starting point but with $x^0 = \begin{bmatrix} 2 \\ 1 \end{bmatrix}$ it gets messy!!
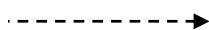
Solution Paths from two different starting points



$$x^0 = \begin{bmatrix} 0 \\ 0 \end{bmatrix}$$ ───────────►

$$x^0 = \begin{bmatrix} 2 \\ 1 \end{bmatrix}$$ ┄┄┄┄┄┄►

Newton's method is simple to use and understand but it has several drawbacks,

- Convergence depends strongly on starting point

- Requires matrix inversion

- Requires second partial derivatives to be computed

- Since it only searches for a stationary point, it may sometimes not converge to a minimum

## Newton's Method – An Alternative Viewpoint

Consider an approximation for $f$ using the first three terms of the Taylor series expansion for $f$ about the current iterate $x^k$:

$$f(x^k+d) \approx Q(d) = f(x^k) + d^T\nabla f(x^k) + \tfrac{1}{2} d^T\nabla^2 f(x^k)d$$

At every iteration Newton's method actually approximates $f(x)$ as above. Now suppose we try to minimize w.r.t. $d$ the above approximation to find an optimal $d$ by setting its gradient to zero (i.e. the necessary conditions $\nabla Q(d)=\mathbf{0}$). We get

$$\nabla f(x^k) + \nabla^2 f(x^k)d = \mathbf{0}, \text{ i.e., } \nabla^2 f(x^k)d = -\nabla f(x^k)$$

which yields the same value for $d$ that we saw earlier (the above system of linear equations is called the Newton Equations…).

This yields more insight into Newton's method:

- First, the approximation above may not be great as we move further away from $x^k$, or if the function is highly nonlinear, and thus if we start too far off from the optimum we may not converge at all.

- Second, even if we do converge, there's no guarantee we will converge to a minimum since we only look for a stationary point (i.e., one that satisfies the necessary conditions) – so we could wind up at a maximum or a saddle point.

- When the method DOES converge though it has a very good quadratic convergence rate; it works best when we're in the vicinity of the optimum.

# Unconstrained Optimization

There are two basic strategies: (i) Line Search, and (ii) Trust Region

**Line Search**

Here we choose a **search direction** $d^k$ along which the function $f$ will improve, and search along this direction from the current iterate $x^k$ for a better iterate. The distance to move along this direction (**step size**) is chosen by (approximately) solving the one-dimensional optimization problem
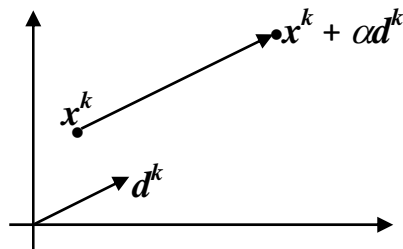
$$\text{Minimize } f(x^k + \alpha d^k), \text{ st } \alpha > 0$$

The value of $\alpha$ that solves the above (say $\alpha^*$) yields $x^{k+1=}\ x^k + \alpha^* d^k$ and the process continues from $x^{k+1}$

Also, a search direction $d^k$ often has the form

$$d^k = -B_k^{-1} \nabla f(x^k)$$

where $B_k$ is some symmetric and nonsingular matrix. In Newton's method $B_k$ is the exact Hessian $\nabla^2 f(x^k)$. For the steepest descent method it is simply the identity matrix $I_k$ and for quasi-Newton methods it is an approximation of the Hessian that is updated at each iteration.

**Trust Region**

In this strategy information about $f$ is used to construct a "model" $m_k$ that approximates the original function well in some limited area around the current iterate $x^k$ (called the "trust region") where we "trust" the approximation.   We then try to minimize the approximation in the trust region, i.e., we solve for $d$ via

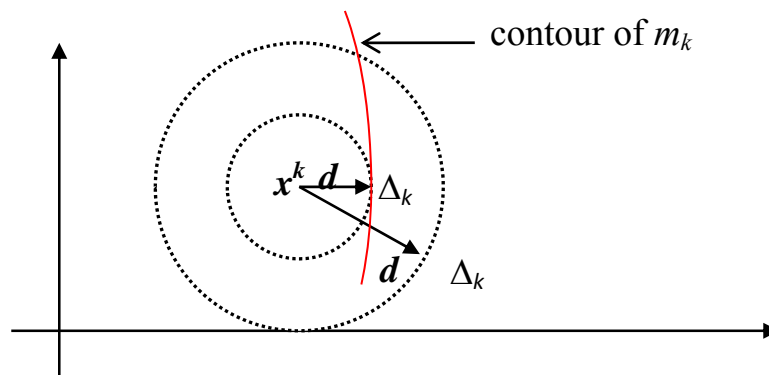"Minimize $m_k(x^k + d)$, st $(x^k + d)$ lies within the trust region"

If the solution is not satisfactory we shrink the trust region and try again.  Trust regions are typically balls, or boxes, or ellipsoids around $x^k$.

A typical model $m_k$ that approximates the function in the trust region around $x^k$ is

$$m_k(x^k+d) \approx f(x^k) + d^{\mathrm{T}} \nabla f(x^k) + \left(\tfrac{1}{2}\right)d^{\mathrm{T}}B(x^k)d$$

where the matrix $B$ is typically the Hessian (or some approximation of it) at $x^k$.

A typical constraint for the trust region might be $\|d\| \le \Delta_k$ so that we look inside a ball of radius $\Delta_k$.  If the optimum value $d^*$ is unsatisfactory we reduce the value of $\Delta_k$; otherwise we might choose to increase this value.  Note that when we shrink the trust region the optimum $d$ will in general, point in a different direction than the previous value of $d$.

In some sense the two approaches differ in how we choose the **direction** to move along and the **distance** we go along this direction:

- In line search we first fix a direction $d^k$ and then search along this direction for the distance $\alpha_k$ to move.

- In the trust region approach, we first choose a maximum distance $\Delta_k$ and then seek a direction and step that maximizes improvement, under the constraint that the step cannot be larger than the maximum distance $\Delta_k$. If the performance in moving to a new point is unsatisfactory, we reduce the value of $\Delta_k$ and try again in a smaller trust region. Conversely if the performance is good we expand the size of the trust region. "Performance" here refers to the quality of the approximation, usually measured by a comparison of the predicted improvement in the objective using the approximation with the actual improvement in the true objective value.

- The actual "amount" of work at each step of the trust region approach could be a little more since we have a constrained optimization and since line search methods seldom try and find the exact step length. However, they still perform very well in practice and both approaches seem to be equally popular.

# DESCENT METHODS

Suppose that $f$ is a continuously differentiable function defined on $\mathbb{R}^n$.

<u>DEFINITION</u>:  If $\nabla f^{\mathrm{T}}(\overline{x})\, d < 0$, then $d$ is called a *DIRECTION OF DESCENT* of the function $f$ at the point $\overline{x}$.

Suppose that $\delta$ is some (sufficiently small) positive number. Then for all $\theta > 0$

$$(1) \qquad f(\overline{x}+\theta d) = f(\overline{x}) + \nabla f^{\mathrm{T}}(\overline{x})\,(\theta d) + \theta\,\|d\|\,\beta(x^0,\,\theta d)$$

i.e., $\{f(\overline{x}+\theta d) - f(\overline{x})\}/\theta = \nabla f^{\mathrm{T}}(\overline{x})\,(d) + \|d\|\,\beta(x^0,\,\theta d)$

Now, by definition $\|d\|\,\beta(x^0,\,\theta d) \to 0$ as $\theta \to 0$, and since $\nabla f^{\mathrm{T}}(\overline{x})(d) < 0$, it follows that

$\{f(\overline{x}+\theta d) - f(\overline{x})\}/\theta < 0$, i.e., $f(\overline{x}+\theta d) < f(\overline{x})$ for <u>all</u> values of $\theta$ that are smaller than some (sufficiently small) positive value $\delta$, i.e. for all $\theta \in (0,\delta)$.

In other words, *by moving some finite distance $\theta$ along the direction of descent $d$ from the current point $\overline{x}$ to the new point $\overline{x}+\theta d$, we can reduce the value of $f(x)$ from $f(\overline{x})$ to $f(\overline{x}+\theta d)$.*

The question is how far to move, since if we move too far ($\theta > \delta$) we are no longer guaranteed a reduction in the value of $f$.

To formalize this:

Given the current approximation $x^i$ we have to answer the following:

1. Is $x^i$ "sufficiently" near optimal? If so, STOP.

2. In what direction $d^i$ should the next point $x^{i+1}$ be located?

3. How far along this direction $d^i$ should we move from $x^i$ to locate $x^{i+1}$?
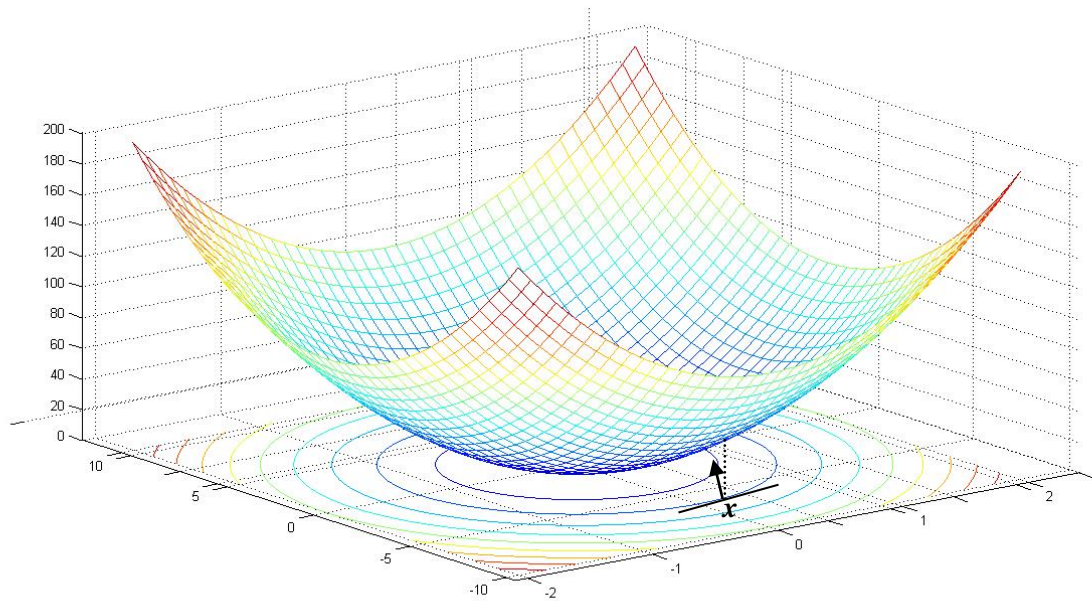
i.e., $x^{i+1} = x^i + \alpha_i d^i$ where

- $\alpha_i$ is a scalar quantity called the **step size**.

- $d^i$ is a vector quantity called the **search direction**.

In Newton's method $x^{i+1} = x^i + [-H^{-1} \nabla f]$, so that $\alpha_i$ is fixed at 1.0, and $d^i = -H^{-1}\nabla f = -[\nabla^2 f(x^i)]^{-1} \nabla f(x^i)$.
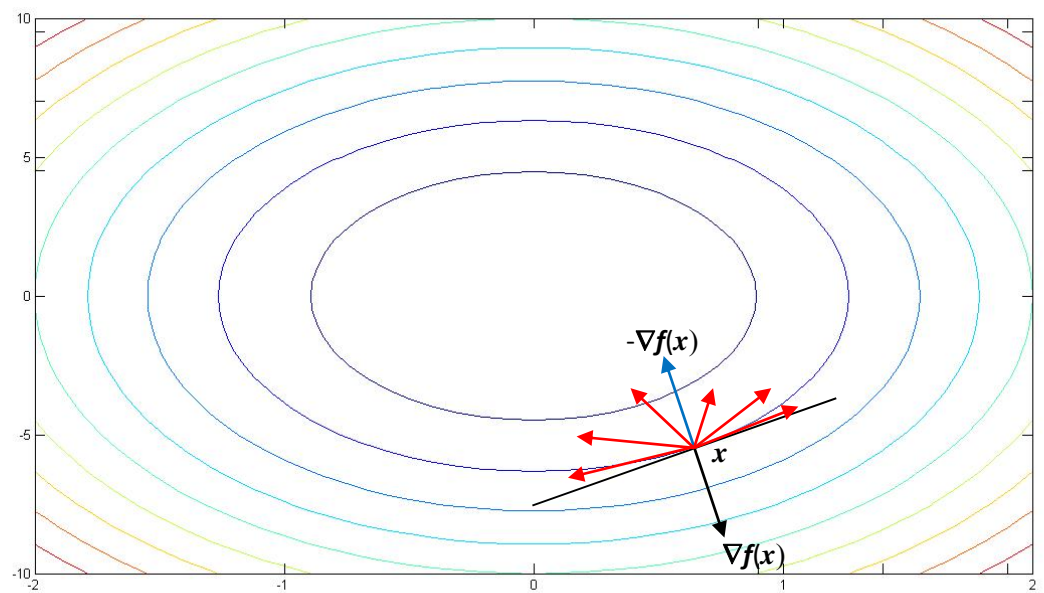
In the **STEEPEST DESCENT** method we choose $d^i$ differently. If $\Delta f = f(x^{i+1}) - f(x^i)$ and $[\Delta x] = x^{i+1} - x^i$, then from (1), $\Delta f \approx \nabla f^T(x^i) [\Delta x]$.

It is easy to see that $|\Delta f|$ can be maximized (i.e., the greatest reduction in $f$ obtained) by choosing $\Delta x = -\theta \nabla f^T(x^i)$. (The gradient vector points in the direction of steepest ascent).

Thus we select $d^i = -\nabla f(x^i)$.

*3-dimensional plot*



various directions of descent

direction of steepest descent

*Contour Plot*

$$f(x) = 25x_1^2 + x_2^2$$

The step size $\alpha_i$ is usually found by means of a one-dimensional search method (with $\alpha_i$ as the decision variable and $x^i + \alpha_i d^i$ as the function to be minimized for $x^i$ and $d^i$ given), such as the Golden-Section or Fibonacci or quadratic interpolation etc. (start with some initial interval for $\alpha_i$ and find the point inside that interval that minimizes $f(x^i + \alpha_i d^i)$.

The method continues until $\|\nabla f(x^k)\| \le \varepsilon$, where $\varepsilon$ is some user-specified tolerance.

AN EXAMPLE:  Minimize $f(x_1, x_2) = x_1^2 + 25x_2^2$

$$\nabla f(x) = \begin{bmatrix} 2x_1 \\ 50x_2 \end{bmatrix}$$

ITERATION 1:  Let $x^0 = \begin{bmatrix} 2 \\ 2 \end{bmatrix}$

(Note that Newton's method would find the minimum $f(0,0)=0$ in one step since $f$ is quadratic and $\nabla f=0$ is thus simply a set of linear equations!).

$$\nabla f(x^0) = \begin{bmatrix} 4 \\ 100 \end{bmatrix} \quad \|\nabla f(x^0)\| = 100.08 \quad d^0 = \begin{bmatrix} -4 \\ -100 \end{bmatrix}$$

$$x^1 = x^0 + \alpha d^0 = \begin{bmatrix} 2 \\ 2 \end{bmatrix} + \alpha \begin{bmatrix} -4 \\ -100 \end{bmatrix} = \begin{bmatrix} 2 - 4\alpha \\ 2 - 100\alpha \end{bmatrix}$$

$f(x^1) = (2\text{-}4\alpha)^2 + 25(2\text{-}100\alpha)^2 = 250016\alpha^2 - 10016\alpha + 104$

**Q.**  How large should $\alpha$ be to fix $x^1$?

**A.**  We try to choose $\alpha$ so that $f(x^1)$ is minimized.

The minimum of $f(2-4\alpha, 2-100\alpha)$ occurs at $\alpha_0 = 0.02003$.

$$\therefore \; x^1 = x^0 + \alpha_0 d^0 = \begin{bmatrix} 1.920 \\ -0.003 \end{bmatrix} \text{ and } f(x^1) = 3.19.$$

ITERATION 2: $\quad x^1 = \begin{bmatrix} 1.920 \\ -0.003 \end{bmatrix}$

$$\nabla f(x^1) = \begin{bmatrix} 3.84 \\ -0.15 \end{bmatrix} \quad \|\nabla f(x^1)\| = 3.84 \quad d^1 = \begin{bmatrix} -3.84 \\ 0.15 \end{bmatrix}$$

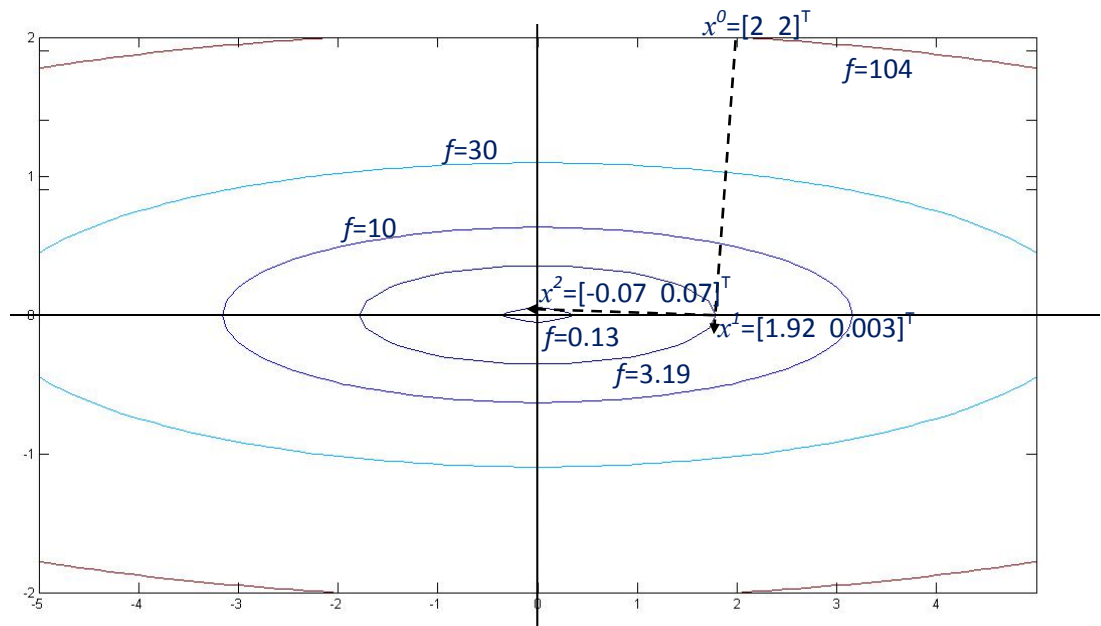$$x^2 = x^1 + \alpha d^1 = \begin{bmatrix} 1.92 - 3.84\alpha \\ -0.003 + 0.15\alpha \end{bmatrix}$$

$f(x^2)$ is minimized at $\alpha_1 \cong 0.48$.

$$\therefore \; x^2 = x^1 + \alpha_1 d^1 = \begin{bmatrix} -0.07 \\ 0.07 \end{bmatrix} \text{ and } f(x^2) = 0.13.$$

ITERATION 3: $\quad x^2 = \begin{bmatrix} -0.07 \\ 0.07 \end{bmatrix}$

$$\nabla f(x^2) = \begin{bmatrix} -0.14 \\ 3.50 \end{bmatrix} \quad \|\nabla f(x^2)\| = 3.50 \quad \text{etc. etc. etc.}$$

**THE QUADRATIC CASE**

Consider the convex quadratic function $f(x) = \frac{1}{2} x^T [Q]x + x^T b$, where $[Q]$ is an $(n \times n)$ symmetric positive definite matrix (the previous example falls into this category).

We have

$$\nabla f^T(x) = Qx + b$$

and the method of steepest descent is

$$x^{i+1} = x^i - \alpha \nabla f(x^i).$$

If we choose $\alpha$ so as to minimize the function w.r.t. $\alpha$

$$f(x^i - \alpha \nabla f(x^i)) = \frac{1}{2} [x^i - \alpha \nabla f(x^i)]^T Q [x^i - \alpha \nabla f(x^i)] - [x^i - \alpha \nabla f(x^i)]^T b$$

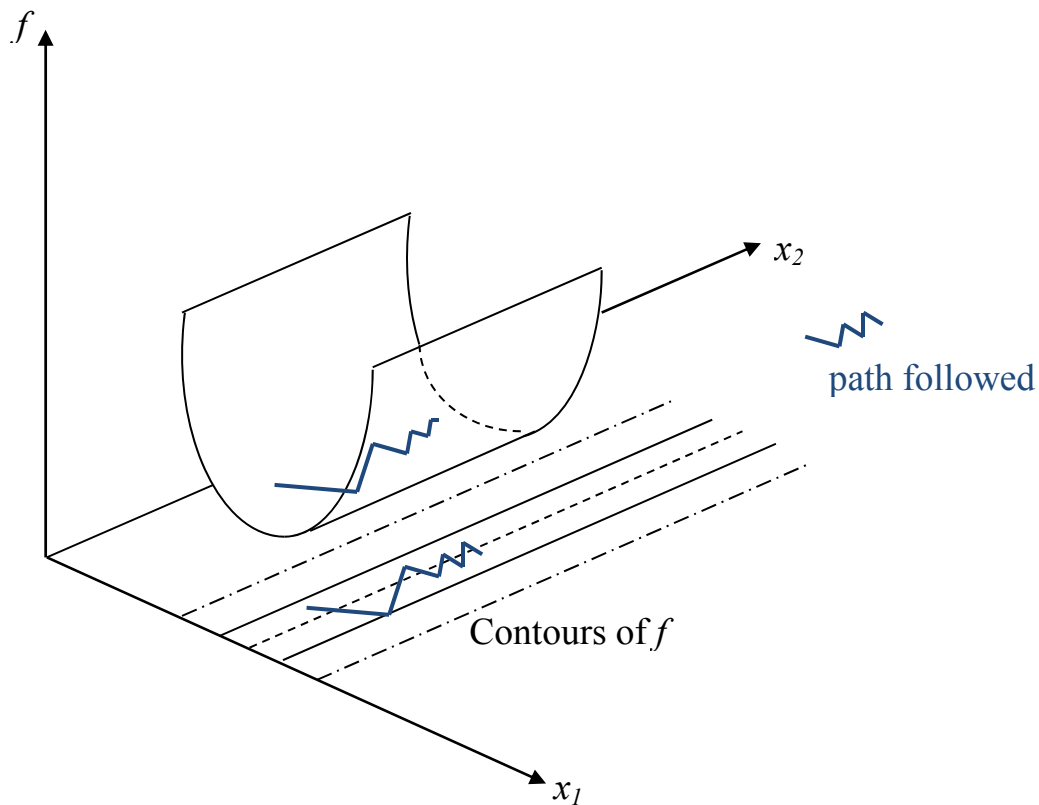For the sake of convenience, let us denote $\nabla f_i = \nabla f(x^i)$. It follows that

$$\alpha_i = \frac{\nabla f_i^T \nabla f_i}{\nabla f_i^T Q \nabla f_i}$$

Thus the new point is,

$$x^{i+1} = x^i - \left( \frac{\nabla f_i^T \nabla f_i}{\nabla f_i^T Q \nabla f_i} \right) \nabla f(x^i)$$

PROPOSITION: For the quadratic case successive search directions are mutually orthogonal. (PROVE and VERIFY using the previous example...)

The main problem with steepest descent is the phenomenon of "zig-zagging" or "oscillation" when travelling along a narrow "valley", especially if $\alpha_k$ (the step size) is not exactly determined to minimize along the direction $\boldsymbol{d}^k$.

# VARIANTS OF THE STEEPEST DESCENT METHOD

*Alternate method for choosing step size $\alpha_k$*

Often it may be difficult to compute the step size based upon an exact line search, or the function may not be unimodal. Thus the step size may be chosen by some other logic:

STEP 0: Select an $\varepsilon$ such that $0<\varepsilon<1$ (usually $\varepsilon \leq 0.5$).

Also, choose any $\alpha>0$.

STEP 1: Compute $\boldsymbol{x} = \boldsymbol{x}^i - \alpha \nabla f(\boldsymbol{x}^i)$,

STEP 2: Compute $f(\boldsymbol{x})$ and check if,

$$|f(\boldsymbol{x}^i) - f(\boldsymbol{x})| \geq \varepsilon . \alpha . \left\| \nabla f(\boldsymbol{x}^i) \right\|^2$$

STEP 3: If the above inequality is satisfied, then set $\alpha_i = \alpha$ and end the line search.

Else, reduce $\alpha$ (perhaps by multiplying $\alpha$ by some $\lambda<1$) and return to Step 1.

*Note*: The initial choice of $\alpha$ is critical. Normally the initial values of $\alpha$ are relatively large during the initial iterations and then become successively smaller. This procedure ensures a "sufficient decrease" without exact computation of the minimum along the search direction.

# Guaranteeing Descent

Recall that a valid search direction $d$ is required to satisfy $\nabla f^T d < 0$.

With Newton's method $d = -(\nabla^2 f)^{-1}\nabla f = -H^{-1}\nabla f$

Thus we need to satisfy $\nabla f^T d = -\nabla f^T H^{-1}\nabla f < 0$, i.e. $\nabla f^T H^{-1}\nabla f > 0$

If $H$ is positive definite this condition is always satisfied. But what if $H$ is not

positive definite (e.g., it is indefinite)? Note that requiring that $H$ be positive definite

is a stronger requirement than just requiring that $\nabla f^T d < 0$.

STRATEGY: Replace $H$ with some "related" matrix ($A$) that is positive definite!

Justification:

- This guarantees descent

- If done properly this guarantees convergence with a line search

- Since $H$ is usually positive definite (and definitely, at least positive
  semidefinite) at the optimum $x^*$ we usually need to do this only at points far
  away from the optimum

- The related matrix $A$ is easily found with little extra effort

Consider the system $Hd = -\nabla f$ that is solved to find $d$.

If $H$ is positive definite it can always be factorized as $H=LDL^T$, where $D$ is a positive

diagonal matrix and $L$ is lower triangular so we just solve $LDL^T d = -\nabla f$.

This is easily done in two stages involving just back substitution (let $u=DL^T d$ and

find $u$ first, then solve $u=DL^T d$ for $d$).

If $H$ is not positive definite then some element $d_{ii}$ of $D$ is negative and we revise the factorization so that these values are replaced. It can be shown that this is equivalent to replacing $H$ with a matrix $A = H+E$ (where $E$ is diagonal) so that $A$ is positive definite. We then factor this matrix as usual.

Example:

Suppose $H = \begin{bmatrix} -1 & 2 & 4 \\ 2 & -3 & 6 \\ 4 & 6 & 22 \end{bmatrix}$. This is not positive definite. However if we define

$E = \begin{bmatrix} 5 & 0 & 0 \\ 0 & 12 & 0 \\ 0 & 0 & 0 \end{bmatrix}$, then $A=H+E = \begin{bmatrix} 4 & 2 & 4 \\ 2 & 9 & 6 \\ 4 & 6 & 22 \end{bmatrix}$ is positive definite and we can factorize

it as

$$A = \begin{bmatrix} 4 & 2 & 4 \\ 2 & 9 & 6 \\ 4 & 6 & 22 \end{bmatrix} = \begin{bmatrix} 1 & 0 & 0 \\ 1/2 & 1 & 0 \\ 1 & 1/2 & 1 \end{bmatrix} \begin{bmatrix} 4 & 0 & 0 \\ 0 & 8 & 0 \\ 0 & 0 & 16 \end{bmatrix} \begin{bmatrix} 1 & 1/2 & 1 \\ 0 & 1 & 1/2 \\ 0 & 0 & 1 \end{bmatrix}$$

And we would now use this to solve for $d$ via

$$\begin{bmatrix} 1 & 0 & 0 \\ 1/2 & 1 & 0 \\ 1 & 1/2 & 1 \end{bmatrix} \begin{bmatrix} 4 & 0 & 0 \\ 0 & 8 & 0 \\ 0 & 0 & 16 \end{bmatrix} \begin{bmatrix} 1 & 1/2 & 1 \\ 0 & 1 & 1/2 \\ 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} d_1 \\ d_2 \\ d_3 \end{bmatrix} = \begin{bmatrix} \nabla f_1 \\ \nabla f_2 \\ \nabla f_3 \end{bmatrix}$$

This is easily solved by using two back substitution stages…

# Globalization Strategies

Typically we compute a search direction based on a Taylor series approximation. "Globalization strategies" ensure that even when the approximation is not great (e.g. far from the solution) we can use the direction to pick a new point in a fashion that ensures global convergence to a stationary point.

Consider the typical line search formula for a new iterate: $x^{k+1} = x^k + \alpha d^k$

Suppose $d^k$ is a descent direction, i.e., $(d^k)^T \nabla f(x^k) < 0$. So ideally, we would like to exactly minimize $f(x^k + \alpha d^k)$, st $\alpha > 0$. But this could be too expensive, so in practice, we accept an *approximate* minimizer that reduces $f$ "sufficiently."

**Search Direction:** To guarantee convergence we need to make two assumptions

A. It produces "sufficient" descent

Assuming we can always choose $d^k$ such that $(d^k)^T \nabla f(x^k) < 0$, we want to be sure that the quantity $(d^k)^T \nabla f(x^k)$ is "sufficiently" negative, i.e., that it is not "too close" to zero and thus slows down progress substantially. Since we can write

$$(d^k)^T \nabla f(x^k) = \left\| d^k \right\| \left\| \nabla f(x^k) \right\| \cos\theta$$

this is equivalent to requiring that the angle $\theta$ be sufficiently obtuse, i.e. that

$$\cos\theta = \frac{(d^k)^T \nabla f(x^k)}{\left\| d^k \right\| \cdot \left\| \nabla f(x^k) \right\|} \leq -\varepsilon < 0, \text{ i.e.,} \qquad -\frac{(d^k)^T \nabla f(x^k)}{\left\| d^k \right\| \cdot \left\| \nabla f(x^k) \right\|} \geq \varepsilon > 0,$$

where $\varepsilon > 0$ is some specified tolerance. This is called the ***Angle Condition***.

B. It is "gradient related"

The search direction is said to be gradient-related if

$$\left\| d^k \right\| \geq m \left\| \nabla f(x^k) \right\|$$

for all $k$, where $m>0$ is some constant.  This ensures that the norm of the search

direction is not too much smaller than that of the gradient.  Note that for the

steepest descent method this is true with $m=1$.

**Step Size**:  Similarly, with the step size we need to make two assumptions as well

A. It produces a "sufficient" decrease in the objective $f$

Here "sufficiency" is related to the Taylor series: a linear (first order)

approximation would be $f(x^k+\alpha d^k) \approx f(x^k) + \alpha(d^k)^T \nabla f(x^k)$.  We will require that $\alpha$

produce a decrease that is *at least some fraction of the decrease* $-\alpha(d^k)^T \nabla f(x^k)$

*predicted by the above approximation*, i.e,

$$f(x^k) - f(x^k+\alpha d^k) \geq \varepsilon\{-\alpha(d^k)^T \nabla f(x^k)\},$$

i.e.,  $$\phi(\alpha) = f(x^k+\alpha d^k) \leq f(x^k) + \varepsilon\alpha(d^k)^T \nabla f(x^k)$$

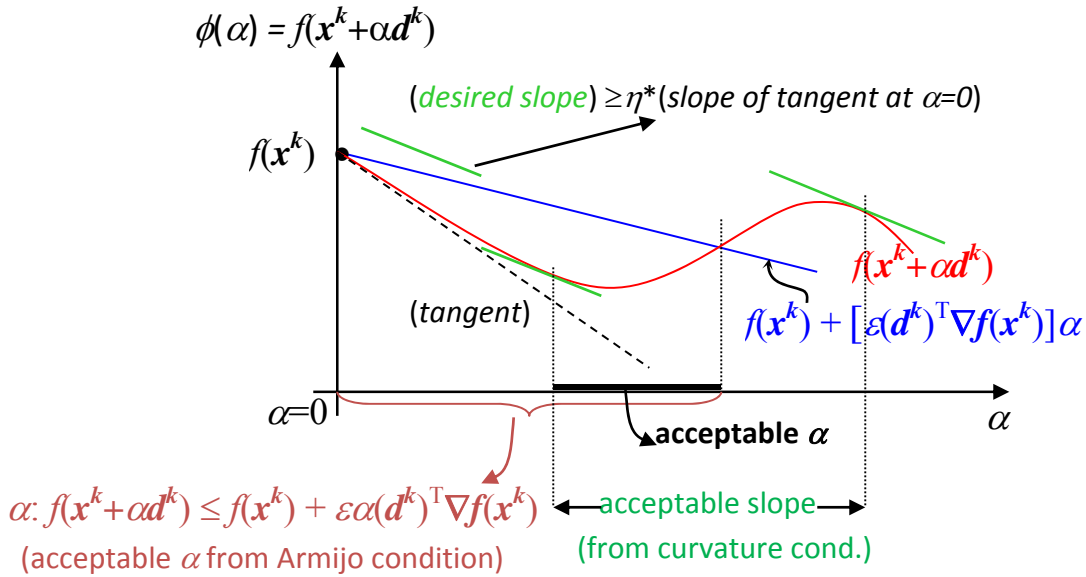where $0<\varepsilon<1$.  This is called the *Armijo* condition

B. It is not "too small"

In this condition we wish to ensure that $\alpha$ itself is not too small.  One simple

approach to ensure that we don't have unacceptably short steps is to start with

some value of $\alpha$ (e.g., $\alpha=1$ corresponding to the ideal Newton step) and if it does

not satisfy the Armijo condition, we reduce $\alpha$ by half and continue doing so until

we find an $\alpha$ for which the decrease in $f$ is sufficiently large.

There are other methods for setting $\alpha$ that are more efficient and do the same

job. For example, the **curvature condition** which states that $(d^k)^T \nabla f(x^k + \alpha d^k) \geq$

$\eta(d^k)^T \nabla f(x^k)$ where $0 < \varepsilon < \eta < 1$ and $\varepsilon$ is the value used for the Armijo condition.

Note that the LHS is the derivative of $\phi(\alpha)$, i.e., its slope and we just want this to

be greater than $\eta$ times the slope of $\phi(0)$.

The Armijo and curvature conditions together are called the **Wolfe Conditions**

$\phi(\alpha) = f(x^k + \alpha d^k)$

(*desired slope*) $\geq \eta *$(*slope of tangent at $\alpha = 0$*)

$f(x^k)$

$f(x^k + \alpha d^k)$

(*tangent*)

$f(x^k) + \left[\varepsilon(d^k)^T \nabla f(x^k)\right]\alpha$

$\alpha = 0$

$\alpha$

▲**acceptable** $\alpha$

$\alpha: f(x^k + \alpha d^k) \leq f(x^k) + \varepsilon\alpha(d^k)^T \nabla f(x^k)$

←acceptable slope→

(acceptable $\alpha$ from Armijo condition)

(from curvature cond.)

THEOREM: If the above four conditions (two on search direction and two on step

size) are satisfied for a general line-search type descent algorithm (along with some

additional assumptions such as compactness, etc.), then it can be proved that the

algorithm will converge to a stationary point having $\left\|\nabla f(x)\right\| = 0$

### Alternate Stopping Criteria

The termination criterion of $\|\nabla f(x)\| < \varepsilon$ is good only for small, well-scaled

problems. Most engineering design problems are however, not well-scaled. For

instance consider the problem

$$\text{Minimize } f(x) = 1 + \frac{10^{-9}}{2}(x\text{-}10^5)^2$$

At $x^0 = 0$, $f(x^0)=6$, $\|f'(x^0)\| = 10^{-4}$.

At the optimum point $x^* = 10^5$ we have $f(x^*)=1$,

Thus if the termination criterion had been $\|\nabla f(x)\| < \varepsilon$, where, say $\varepsilon=10^{-3}$, the

program would have stopped at $x^0$. However, by proceeding further a reduction of

$f(x^0)\text{-}f(x^*) = (6\text{-}1)$ units, i.e., $(6\text{-}1)/1 = 500\%$ is possible!


It should be emphasized that no one set of termination criteria is suitable for all cases.

Termination criteria depend on nonlinearity, scaling and magnitudes of functions and

variables. From an engineering standpoint $|f_i\text{-}f|$ is more important $\|x^i\text{-}x^*\|$.


Consequently we may use:

$$|f(x^k) - f(x^{k-1})| < \varepsilon \text{ for } k=i,i\text{-}1,...,i\text{-}m$$

i.e., change in objective is less than a specified tolerance for *m successive iterations*.

Gill, Murray & Wright suggest the following conditions be *simultaneously* satisfied:

(1) $f_{i-1} - f_i \leq \varepsilon(1+|f_i|)$

(2) $\| x^{i-1} - x^i \| \leq (\varepsilon)^{1/2} (1+ \| x^i \|)$

(3) $\| \nabla f(x^i) \| \leq (\varepsilon)^{1/3}(1+|f_i|)$

where $f_i \equiv f(x^i)$ and $\varepsilon$ is some prespecified tolerance.

(1) and (2) check whether the sequence $\{x^i\}$ is converging while (3) is based on the

necessary condition $\nabla f^* = 0$. For well scaled problems satisfaction of (1) will imply

satisfaction of (2); for poorly scaled problems (2) will force the algorithm to try

harder for a better point. In engineering design we might even stop if the cost at the

current point is (say) 70% of the starting cost, i.e., $f_i/f_0 \leq 0.7$

## *SCALING*

Consider the problem:  Minimize $f(x)$ st $x \in \mathbb{R}^n$.

If we scale variables by the linear transformation using the ($n \times n$) matrix $\Lambda$:

$$\underset{n \times n}{\Lambda^{-1}} \cdot \underset{n \times 1}{x} = \underset{n \times 1}{y} \qquad \text{or equivalently} \qquad \underset{n \times 1}{x} = \underset{n \times n}{\Lambda} \cdot \underset{n \times 1}{y}$$

We now have the problem:  Minimize $g(y) \equiv f(\Lambda y)$.

So
$$\nabla_y g = \Lambda \cdot \nabla_y f \qquad \text{and} \qquad \nabla_y^2 g = \Lambda^T (\nabla_x^2 f) \Lambda$$

Scaling is recommended when the variables $x_i$ have widely differing magnitudes. It is

also useful if the matrix $\nabla_y^2 g$ has a better condition number (ratio of the largest to the

smallest eigenvalue) than the matrix $\nabla_x^2 f$.

Example:

Consider     Minimize $f(x) = 4x_1^2 + x_2^2$     (the optimum solution is at $x^* = [0 \ 0]^T$)

We have

$$\nabla_x f = [8x_1 \quad 2x_2]^T \qquad \text{and} \qquad \nabla_x^2 f = \begin{bmatrix} 8 & 0 \\ 0 & 2 \end{bmatrix}.$$

Thus $cond[\nabla_x^2 f] = \lambda_{max}/\lambda_{min} = 8/2 = 4$

So if $x^0 = [0.5 \quad 1]^T$ then the steepest descent method generates the sequence

$$x^1 = \begin{bmatrix} 0.5 \\ 1 \end{bmatrix} - \frac{5}{34}\begin{bmatrix} 4 \\ 2 \end{bmatrix} = \begin{bmatrix} -3/34 \\ 24/34 \end{bmatrix}, \qquad x^2 = \begin{bmatrix} -3/34 \\ 24/34 \end{bmatrix} - \alpha\begin{bmatrix} -12/34 \\ 48/34 \end{bmatrix} = \begin{bmatrix} ?? \\ ?? \end{bmatrix},$$

etc., etc.

Now consider the scaling transform $y_1 = 2x_1$, $y_2 = x_2$, so that $g(y) = y_1^2 + y_2^2$

$$\text{i.e.,} \qquad \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} = \begin{bmatrix} 0.5 & 0 \\ 0 & 1 \end{bmatrix}\begin{bmatrix} y_1 \\ y_2 \end{bmatrix}$$

Then     $\nabla_y^2 g = \begin{bmatrix} 2 & 0 \\ 0 & 2 \end{bmatrix}$ with $cond[\nabla_y^2 g] = 1$.

At $x^0 = [0.5 \quad 1]^T$, i.e., $y = [1 \quad 1]^T$ the steepest descent method yields

$$y^1 = \begin{bmatrix} 1 \\ 1 \end{bmatrix} - \frac{8}{16}\begin{bmatrix} 2 \\ 2 \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \end{bmatrix} = y^*$$

Single step convergence!

# A Trust-Region Algorithm

Suppose we start at the initial point $x^0$ with an initial trust region bound $\Delta_0$.

Specify a tolerance $\varepsilon$. Set $k=0$

<u>Step 1</u>: If $\left\|\nabla f(x^k)\right\| \leq \varepsilon$ is optimal, stop.

<u>Step 2</u>: Solve the following problem to obtain the trial step $d^k$:

$$\text{Minimize } m_k(d) = f(x^k) + d^T \nabla f(x^k) + (\tfrac{1}{2}) d^T \nabla^2 f(x^k) d$$

$$\text{st} \quad \|d\| \leq \Delta_k$$

<u>Step 3</u>: Compute $\rho_k = \dfrac{f(x^k) - f(x^k + d^k)}{f(x^k) - m_k(d^k)} = \dfrac{\text{actual reduction}}{\text{predicted reduction}}$

If $0 < \rho_k \leq \eta$ set $x^{k+1} = x^k + d$ (unsuccessful step) and $\Delta_{k+1} = 0.5\Delta_k$,

If $\rho_k \geq \mu$ set $x^{k+1} = x^k + d$ and $\Delta_{k+1} = 2\Delta_k$

If $\eta < \rho_k < \mu$ set $x^{k+1} = x^k + d$ and $\Delta_{k+1} = \Delta_k$

Set $k=k+1$ and return to Step 1.

*Typical values for $\eta$ and $\mu$ might be 0.25 and 0.75.

Note that $\rho_k$ indicates how well the model $m_k$ predicts the reduction in the function $f$.

If $\rho_k$ is small (say <0.25) then the actual reduction is much smaller than the prediction

and indicates the model cannot be "trusted" for a large $\Delta_k$. One approach might be to

altogether reject $d$ in this case (in the above algorithm we accept $d$ as long as there is

any reduction in $f$). Conversely, if $\rho_k$ is large (say > 0.75) then the model seems to be

predicting the reduction well – so we expand our trust region.

# A Trust-Region Algorithm (cont'd)

How do we solve the optimization problem in Step 2 to obtain the trial step $d^k$?

Minimize $m_k(d) = f(x^k) + d^T \nabla f(x^k) + (\frac{1}{2}) d^T \nabla^2 f(x^k) d$

st $\quad \|d\| \leq \Delta_k$

In general, this is not a trivial problem to solve. Note that a point satisfying the

necessary conditions for the unconstrained problem would be obtained from

$\nabla m_k(d) = 0$, i.e., $\nabla f(x^k) + \nabla^2 f(x^k) d = 0$, i.e., the Newton step $d = -[\nabla^2 f(x^k)]^{-1} \nabla f(x^k)$.

So if $\|d\| \leq \Delta_k$ then we just accept this as our trial step. If on the other hand $\|d\| > \Delta_k$,

it may be shown that we need to find a $\lambda$ such that the direction $d$ obtained from

$$d = -[\nabla^2 f(x^k) + \lambda I]^{-1} \nabla f(x^k) \qquad\qquad (\otimes)$$

has $\|d\| = \Delta_k$. This leads to a nonlinear equation in $\lambda$ which we solve for $\lambda$ and

substituting this into ($\otimes$) then yields the step $d$.

E.g., if $\nabla^2 f(x^k) = \begin{bmatrix} 120 & 0 \\ 0 & 36 \end{bmatrix}$, $\nabla f(x^k) = \begin{bmatrix} 152 \\ 28 \end{bmatrix}$ $\Delta_k = 1$, ($\|d\| = 1.49$ here...) then we solve

$$\left\| -\begin{bmatrix} 120+\lambda & 0 \\ 0 & 36+\lambda \end{bmatrix}^{-1} \begin{bmatrix} 152 \\ 28 \end{bmatrix} \right\| = 1, \text{ i.e., } \left(\frac{152}{(120+\lambda)}\right)^2 + \left(\frac{28}{(36+\lambda)}\right)^2 = 1 \text{ to get } \lambda = 42.655, \text{ so that}$$

$$d = \begin{bmatrix} -0.9345 \\ -0.356 \end{bmatrix} \text{ with } \|d\| = 1$$

In practice, most algorithms will find only approximate solutions to the subproblem.

# VARIABLE-METRIC (OR *QUASI-NEWTON*) METHODS

Recall that for Newton's method, at each iteration we have to

   (1) Compute the Hessian matrix $H$      *and*       (2) Invert the Hessian

In many situations, $H$ is either not available or very expensive (in terms of CPU time) to compute.  Also, inversion of a matrix involves a lot of computational effort (for an $(n \times n)$ matrix the CPU time required is of the order of $n^3$, which can become very large even for moderate values of $n$).

     The main idea behind quasi-Newton methods is to never *actually* compute the Hessian or its inverse, but rather to gradually construct an approximation to the latter using information gathered as the descent process progresses. That is, obtain second derivative information using only first derivative information. Then if $A_i$ is a matrix which is the current approximation (at iteration $i$) of the true Hessian inverse $H^{-1}$, we take steps in the direction $d^i = -A_i \cdot \nabla f(x^i)$.

## MODIFIED NEWTON'S METHOD

Let $x^{k+1} = x^k - \alpha_k A_k [\nabla f(x^i)]$, where $A_k$ is a symmetric $(n \times n)$ matrix and as usual we

Minimize $f\{ x^k - \alpha_k A_k [\nabla f(x^i)]\}$ to obtain $\alpha_k$.

Note that     if $A_k = H_k^{-1}$, $\alpha_k = 1$ this is Newton's method,

            if $A_k = I_n$ this is steepest descent method.

$A_k$ is required to be positive definite. Usually $A_k$ should be "close to" $H_k^{-1}$.  In the "classical" modified Newton's method $A_k = [H(x^o)]^{-1}$ , i.e., Hessian at $x^0$ is used throughout.

# Quasi-Newton Methods

Consider the Taylor series approximation for $f$ about the current iterate $x^k$:

$$f(x^k + d) \approx f_k + d^T \nabla f_k + (\tfrac{1}{2}) d^T H_k d$$

For ease of notation, let $f_k = f(x^k)$, $\nabla f_k = \nabla f(x^k)$ and $H_k = \nabla^2 f(x^k)$...

Recall that the Newton search direction is obtained by minimizing the above to obtain the Newton equations $H_k d = -\nabla f_k$ and solving for $d$.

In Quasi-Newton methods we either approximate $H_k$ via $D_k$ or, $H_k^{-1}$ via $A_k$.

To see how we get these approximations consider the following approximation for the second derivative of a univariate function:

$$f''(x_k) \approx \frac{f'(x_k) - f'(x_{k-1})}{x_k - x_{k-1}} \implies f''(x_k)(x_k - x_{k-1}) \approx f'(x_k) - f'(x_{k-1})$$

Extending this logic to $n$ dimensions we might rewrite this as follows:

$$H_k[x^k - x^{k-1}] \approx \nabla f_k - \nabla f_{k-1}$$

If we now approximate the matrix $H_k$ by the matrix $D_k$ we obtain the following

**Secant Condition** that is used to define the Quasi-Newton approximation:

$$\boxed{D_k[x^k - x^{k-1}] = \nabla f_k - \nabla f_{k-1}}$$

This is the basis for all Quasi-Newton methods. Note that in particular, if $f$ is

quadratic, so that we can write it as $f(x) = (\tfrac{1}{2}) x^T H_k x + c^T x$, we have

$$H_k[x^k - x^{k-1}] = \left(H_k x^k + c\right) - \left(H_k x^{k-1} + c\right) = \nabla f_k - \nabla f_{k-1}$$

i.e., $H_k$ satisfies the Secant Condition!

In Quasi-Newton methods we start with some matrix that satisfies the secant condition and update it at each iteration so that it continues to satisfy it.

If we define

$$p^k = x^{k+1} - x^k \text{ and } \qquad y^k = \nabla f_{k+1} - \nabla f_k$$

i.e., $p^k$ is the change in $x$ while $y^k$ is the change in the gradient when going from iteration $k$ to iteration $k+1$, then the Quasi-Newton conditions are $D_k p^{k-1} = y^{k-1}$, or more commonly

$$\boxed{D_{k+1} p^k = y^k}$$

An alternative way of saying the same thing is (recall that $A_k \approx H_k^{-1}$)

$$\boxed{p^k = A_{k+1} y^k}$$

Note: If we use a line search and get $x^{k+1} = x^k + \alpha d^k$ then $p^k = \alpha d^k$ !

A simple example of an update that satisfies the secant condition is the following (symmetric Rank-1 update):

$$D_{k+1} = D_k + \frac{(y^k - D_k p^k)(y^k - D_k p^k)^T}{(y^k - D_k p^k)^T p^k}$$

Note that

$$D_{k+1} p^k = D_k p^k + \frac{(y^k - D_k p^k)(y^k - D_k p^k)^T}{(y^k - D_k p^k)^T p^k} p^k$$

$$= D_k p^k + \frac{(y^k - D_k p^k)\left((y^k - D_k p^k)^T p^k\right)}{(y^k - D_k p^k)^T p^k} = D_k p^k + (y^k - D_k p^k) = y^k$$

This update maintains symmetry but not necessarily positive definiteness.

Better methods are based on the so-called Broyden class of updates that preserve

symmetry as well as positive definiteness as long as $(p^k)^T y^k > 0$. Positive definiteness

is important because presumably the Hessian is positive definite (and certainly,

positive semidefinite) at the optimum.

In this class of updates we can use the following (rather messy!) formulae for

updating either the approximation to the Hessian ($D_k$) or the approximation to its

inverse ($A_k$), depending on which one we need. These are defined as below:

For the approximation to $H_k$ we have

$$D_{k+1} = D_k + \frac{y^k (y^k)^T}{(y^k)^T p^k} - \frac{D_k p^k (D_k p^k)^T}{(p^k)^T D_k p^k} + \phi \left( (p^k)^T D_k p^k \right) u^k (u^k)^T, \qquad (\perp)$$

where $\phi$ is a scalar and $u^k = \dfrac{y^k}{(y^k)^T p^k} - \dfrac{D_k p^k}{(p^k)^T D_k p^k}$

NOTE: It is easily verified that $D_{k+1} p^k = y^k$

For the approximation to the $H_k^{-1}$ we have

$$A_{k+1} = A_k + \frac{p^k (p^k)^T}{(p^k)^T y^k} - \frac{A_k y^k (A_k y^k)^T}{(y^k)^T A_k y^k} + \varepsilon \left( (y^k)^T A_k y^k \right) v^k (v^k)^T \qquad (\perp\perp)$$

where $\varepsilon$ is a scalar and $v^k = \dfrac{p^k}{(p^k)^T y^k} - \dfrac{A_k y^k}{(y^k)^T A_k y^k}$

NOTE: It is easily verified that $A_{k+1} y^k = p^k$

Two of the best Quasi-Newton algorithms use the above schemes for updating the approximations to the Hessian (or its inverse) that is used in place of $H_k$ (or $H_k^{-1}$) when trying to solve the Newton equations $H_k d = -\nabla f_k$ for the search direction $d$, (or computing $d = -(H_k^{-1})\nabla f_k$).

- In the Broyden-Fletcher-Goldfarb-Shanno (BFGS) algorithm we approximate $H_k$

   via $D_k$ and <u>solve</u> the system $H_k d = -\nabla f_k \approx D_k d = -\nabla f_k$ for the vector $d$

   To update $D_k$ we use formula ($\perp$) with $\phi=0$. It may be shown that this is

   equivalent to setting $\varepsilon=1$ in ($\perp\perp$)…

- In the Davidon-Fletcher-Powell (DFP) algorithm we approximate $H_k^{-1}$ via $A_k$ and

   compute $d = -H_k^{-1}\nabla f_k \approx -A_k \nabla f_k.$

   To update $A_k$ we use formula ($\perp\perp$) with $\varepsilon=0$. It may be shown that this is

   equivalent to setting $\phi=1$ in ($\perp$)…


In fact, we may write

$$D_{k+1} = (1-\phi)D_{k+1}^{BFGS} + \phi D_{k+1}^{DFP}$$

or

$$A_{k+1} = \varepsilon A_{k+1}^{BFGS} + (1-\varepsilon)A_{k+1}^{DFP}$$


The above two algorithms from the Broyden class maintain both symmetry as well as positive semidefiniteness - BFGS and DFP are considered among the best methods for unconstrained optimization…

# BROYDEN-FLETCHER-GOLDFARB-SHANNO METHOD

At iteration $i$ the method is as follows:

<u>STEP 1</u>: Let $D_i \cong H_i$ (i.e., an approximation to the true Hessian at the current point $x^i$

(at $x^0$ choose $D_1 = I_n$ or any symmetric positive definite matrix).

<u>STEP 2</u>: If $x^i$ is optimal stop; otherwise obtain the search direction $d^i$ by solving

$$D_i d = -\nabla f_i$$

<u>STEP 3</u>: Minimize $f$ in the direction $d^i$, i.e.,

$$\text{Minimize}_\alpha \ \{f(x^i + \alpha d^i)\} \text{ to find } \alpha \text{ and hence } x^{i+1}.$$

<u>STEP 4</u>: Define

$$p^i = x^{i+1} - x^i = \Delta x^i \text{ (change in } x\text{)}$$

$$y^i = \nabla f_{i+1} - \nabla f_i = \Delta g^i \text{ (change in gradient)}$$

<u>STEP 5</u>: Compute

$$B_i = \frac{y^i (y^i)^T}{(y^i)^T p^i}$$

$$\frac{\blacksquare}{\blacksquare} = \frac{\blacksquare}{\bullet} = \frac{Matrix}{Scalar}$$

$$C_i = -\frac{D_i p^i (D_i p^i)^T}{(p^i)^T D_i p^i}$$

$$\frac{\blacksquare}{\blacksquare} = \frac{\blacksquare}{\bullet} = \frac{Matrix}{Scalar}$$

<u>STEP 6</u>: UPDATE the approximation to $H$ via

$$D_{i+1} = D_i + B_i + C_i$$

Set $i = i + 1$, and return to Step 2

# THE DAVIDON-FLETCHER-POWELL METHOD

STEP 1:  Let $A_i \cong H_i = [\nabla f(x^i)]^{-1}$ (an approximation to the true Hessian at the current point $x^i$ (at $x^0$ choose $A_0 = I_n$ or any symmetric positive definite matrix).

STEP 2: If $x^i$ is optimal stop; otherwise obtain the search direction $d^i$ by solving

$$d^i = -A_i \nabla f_i$$

STEP 3: Minimize $f$ in the direction $d^i$, i.e.,

$$\text{Minimize}_\alpha \ \{f(x^i + \alpha d^i)\} \ \text{to find } \alpha \text{ and hence } x^{i+1}.$$

STEP 4:  Define

$$p^i = x^{i+1} - x^i = \Delta x^i \ (\text{change in } x)$$

$$y^i = \nabla f_{i+1} - \nabla f_i = \Delta g^i \ (\text{change in gradient})$$

STEP 5: Compute

$$B_i = \frac{p^i (p^i)^\top}{(p^i)^\top y^i} \qquad \blacksquare = \frac{\blacksquare}{\centerdot} = \frac{Matrix}{Scalar}$$

$$C_i = -\frac{A_i y^i (A_i y^i)^\top}{(y^i)^\top A_i y^i} \qquad \blacksquare = \frac{\blacksquare}{\centerdot} = \frac{Matrix}{Scalar}$$

STEP 6: UPDATE the approximation to $H^{-1}$ via

$$A_{i+1} = A_i + B_i + C_i$$

Set $i=i+1$, and return to Step 2

If the gradient vector is readily available, this is generally regarded as one of the best algorithms available).

AN EXAMPLE:  Minimize $f(x) = (x_2 - x_1)^2 + (1-x_1)^2$

Thus  $\nabla f(x) = \begin{bmatrix} -2(x_2 - x_1) - 2(1 - x_1) \\ 2(x_2 - x_1) \end{bmatrix} = \begin{bmatrix} 4x_1 - 2x_2 - 2 \\ 2x_2 - 2x_1 \end{bmatrix}$

ITERATION 1

1)  Start with  $x^0 = \begin{bmatrix} 0 \\ 0 \end{bmatrix}$  and  $A_0 = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} \cong H^{-1}(x^0)$

2)  $d^0 = -A_0 \nabla f(x_0) = -\begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} \begin{bmatrix} -2 \\ 0 \end{bmatrix} = \begin{bmatrix} 2 \\ 0 \end{bmatrix}$

3)  $\text{Min}_\alpha f(x^0 + \alpha d^0) = \text{Min}_\alpha f\{(0+2\alpha),(0+0\alpha)\} = \text{Min}_\alpha \ 4\alpha^2 + (1-2\alpha)^2 = 8\alpha^2 + 1 - 4\alpha$

   This yields $\alpha_0 = 0.25$, so that $x^1 = \begin{bmatrix} 0 \\ 0 \end{bmatrix} + 0.25 \begin{bmatrix} 2 \\ 0 \end{bmatrix} = \begin{bmatrix} 0.5 \\ 0 \end{bmatrix}$, with $f(x^1) = 0.5$

4)  $p^0 = \Delta x^0 = \begin{bmatrix} 0.5 \\ 0 \end{bmatrix} - \begin{bmatrix} 0 \\ 0 \end{bmatrix} = \begin{bmatrix} 0.5 \\ 0 \end{bmatrix}$

   $y^0 = \Delta g^0 = \nabla f(x^1) - \nabla f(x^0) = \begin{bmatrix} 0 \\ -1 \end{bmatrix} - \begin{bmatrix} -2 \\ 0 \end{bmatrix} = \begin{bmatrix} 2 \\ -1 \end{bmatrix}$

5)  $B_0 = \dfrac{p^0 [p^0]^T}{[p^0]^T y^0} = \dfrac{\begin{bmatrix} 0.5 \\ 0 \end{bmatrix} [0.5 \ \ 0]}{[0.5 \ \ 0]\begin{bmatrix} 2 \\ -1 \end{bmatrix}} = \dfrac{\begin{bmatrix} 0.25 & 0 \\ 0 & 0 \end{bmatrix}}{1} = \begin{bmatrix} 0.25 & 0 \\ 0 & 0 \end{bmatrix}$

   $C_0 = -\dfrac{[A_0 y^0] \cdot [A_0 y^0]^T}{[y^0]^T A_0 y^0}; \ A_0 y^0 = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} \begin{bmatrix} 2 \\ -1 \end{bmatrix} = \begin{bmatrix} 2 \\ -1 \end{bmatrix}$

   $\therefore C_0 = -\dfrac{\begin{bmatrix} 2 \\ -1 \end{bmatrix} [2 \ -1]}{[2 \ -1]\begin{bmatrix} 2 \\ -1 \end{bmatrix}} = \begin{bmatrix} -0.8 & 0.4 \\ 0.4 & -0.2 \end{bmatrix}$

6)  $A_1 = A_0 + B_0 + C_0$

   $= \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} + \begin{bmatrix} 0.25 & 0 \\ 0 & 0 \end{bmatrix} + \begin{bmatrix} -0.8 & 0.4 \\ 0.4 & -0.2 \end{bmatrix} = \begin{bmatrix} 0.45 & 0.4 \\ 0.4 & 0.8 \end{bmatrix}$

ITERATION 2

2)  $d^1 = -A_1 \nabla f(x_0) = -\begin{bmatrix} 0.45 & 0.4 \\ 0.4 & 0.8 \end{bmatrix} \begin{bmatrix} 0 \\ -1 \end{bmatrix} = \begin{bmatrix} 0.4 \\ 0.8 \end{bmatrix}$

3)  $\text{Min}_\alpha f(x^1 + \alpha d^1) = \text{Min}_\alpha f\{(0.5 + 0.4\alpha),(0 + 0.8\alpha)\}$

   $= \text{Min}_\alpha (0.8\alpha - 0.5 - 0.4\alpha)^2 + (1 - 0.5 - 0.4\alpha)^2$

   This yields $\alpha_1 = 1.25$, and the optimal solution $(\nabla f(x^2) = 0)$ with

   $x^2 = \begin{bmatrix} 0.5 \\ 0 \end{bmatrix} + 1.25 \begin{bmatrix} 0.4 \\ 0.8 \end{bmatrix} = \begin{bmatrix} 1 \\ 1 \end{bmatrix}$, with $f(x^2) = 0$

# CONJUGATE DIRECTIONS

Better descent methods usually use *conjugate search directions.*

<u>DEFINITION</u>: Two directions $\hat{d}$ and $\bar{d}$ are conjugate with respect to the symmetric positive definite matrix $Q$ if,

$$(\hat{d})^T Q \bar{d} = 0$$

Conjugate directions are also linearly independent.

*NOTE*: If $Q=I$ then $(\hat{d})^T Q \bar{d} = (\hat{d})^T \bar{d} = 0$, i.e., the conditions for $\hat{d}$ and $\bar{d}$ to be conjugate is that they should be orthogonal;

## **Geometric Interpretation**

Suppose $f(x) = \frac{1}{2} x^T Q x + b^T x$

Thus $Q$ is the Hessian of $f(x)$ and of order $(n \times n)$

<u>Example</u>: $n=2$ and say $Q$ is positive definite (and of course, symmetric)

e.g., $\quad f(x) = ax_1^2 + bx_2^2 + cx_1x_2 = \frac{1}{2}[x_1 \quad x_2]\begin{bmatrix} 2a & c \\ c & 2b \end{bmatrix}\begin{bmatrix} x_1 \\ x_2 \end{bmatrix}.$
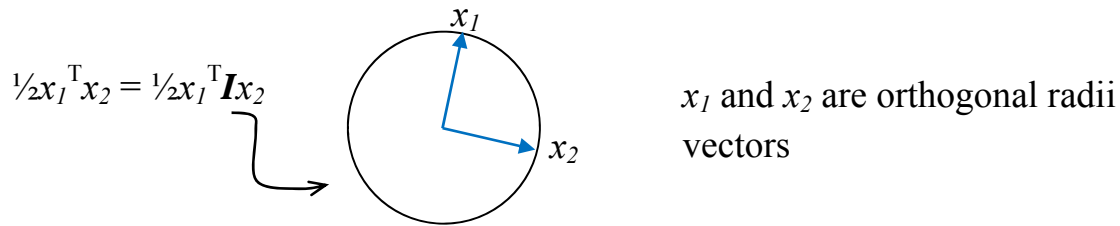
Here the contours of $f$ are ellipses centered at (0,0). In general, the contours of a quadratic function are $n$-dimensional ellipsoids with axes in the directions of the $n$ mutually orthogonal eigenvectors of the Hessian. The axis corresponding to the $i^{th}$ eigenvector is proportional to $1/\lambda_i$.
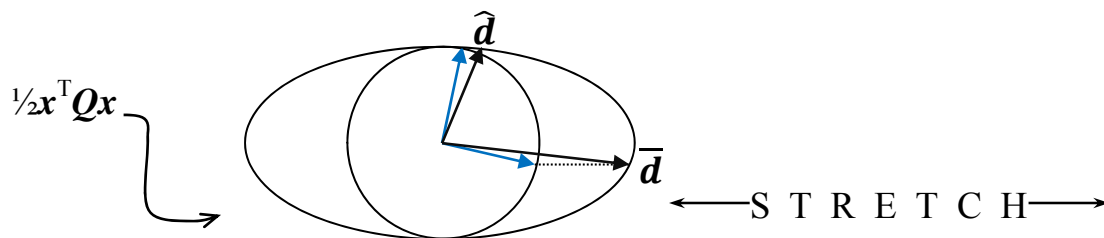
Verify that the eigenvalues are

$$\lambda_1 = (a+b) + \sqrt{(a-b)^2 + c^2} \quad \text{and} \quad \lambda_2 = (a+b) - \sqrt{(a-b)^2 + c^2}$$

E.g., if $a=b=0.5$, and $c=0$, then $\lambda_1=1$, $\lambda_2=1$, so we get a circle centered at $(0,0)$.

Now, if we start with

$$\tfrac{1}{2}x_1{}^\mathsf{T}x_2 = \tfrac{1}{2}x_1{}^\mathsf{T}\mathbf{I}x_2$$

$x_1$ and $x_2$ are orthogonal radii vectors

and we "stretch" this circle along one diameter as shown below

$$\tfrac{1}{2}\mathbf{x}^\mathsf{T}\mathbf{Q}\mathbf{x}$$
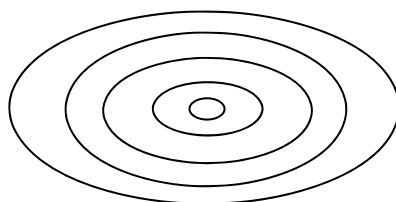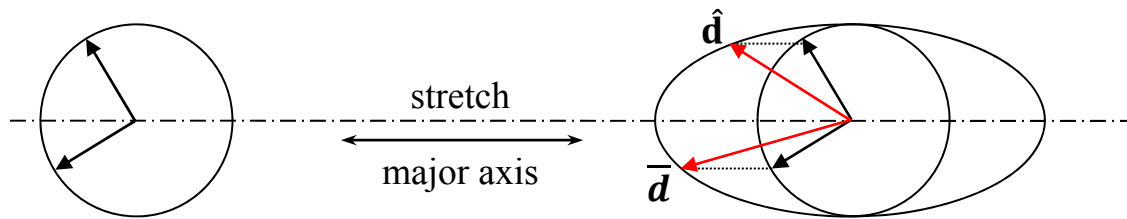
S T R E T C H

The two orthogonal radii vectors become conjugate with respect to some matrix $\mathbf{Q}$.

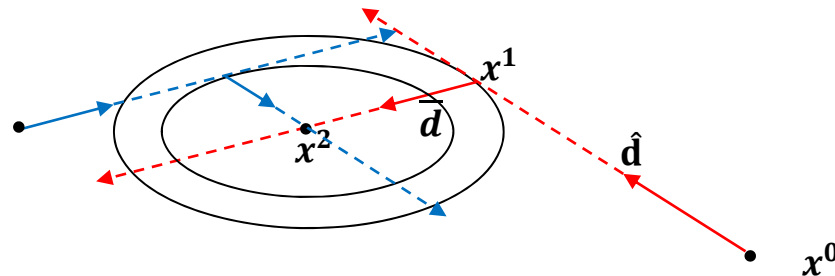**Question**. How do conjugate directions relate to minimization?

<u>THEOREM</u>: The minimum of a convex quadratic function of $n$ variables $f(x_1, x_2, \dots, x_n) = \tfrac{1}{2}\,\mathbf{x}^\mathsf{T}\mathbf{H}\mathbf{x} + \mathbf{c}^\mathsf{T}\mathbf{x} + k$, (where $\mathbf{H}$ is the positive definite Hessian) can be found in **at most** $n$ steps from **any** point via a series of minimizations along each of a set of $\mathbf{H}$-conjugate directions.

Example: $n=2$; contours of $f(\mathbf{x})$ are as shown below:

Directions $\hat{d}$ and $\bar{d}$ are **H**-conjugate



(1) Start at $x^0$ and search along $\hat{d}$ and minimize to obtain $x^1$

(2) Search along $\bar{d}$ and minimize to get $x^2$ (the optimum point)    (**2 steps!**)

(Same thing in second case (but first along $\bar{d}$ and then along $\hat{d}$ )

Normally

- A set of H-conjugate directions $d^1, d^2, ..., d^n$ is not usually known beforehand, but generated by some (often complicated) scheme which depends upon available information (Hessian, gradient vectors, function values at previous points in the search, etc.).

- Conjugate direction methods are "designed" for convex, quadratic functions

- Best for large problems: storage needs are low & "work" per iteration is low

- Unless $f$ is <u>actually</u> quadratic the Hessian matrix is not constant, but changes as the search proceeds -- thus the "conjugate" directions computed are seldom <u>actually</u> **H**-conjugate for any **H**. Despite this, even for non-quadratic functions, the so-called "conjugate" directions generated are good search directions.

# CONJUGATE GRADIENT METHOD

# (For solving $Ax=b$ when $A$ is symmetric & positive definite)

First note that solving $Ax=b$, is equivalent to minimizing $f(x) = \frac{1}{2}x^{\mathrm{T}}Ax - b^{\mathrm{T}}x$ (since the necessary conditions $-\nabla f(x^0) = Ax - b = 0$ are sufficient).

Now consider a set of vectors $p^0, p^1, \ldots, p^m$ that are conjugate w.r.t. $A$, i.e., $(p^i)^{\mathrm{T}}Ap^j = 0$ if $i \neq j$, and suppose we define $y = \sum_{i=0}^{m} \alpha_i p^i$.

Then
$$f(y) = \frac{1}{2}\left(\sum_{i=0}^{m} \alpha_i p^i\right)^{\mathrm{T}} A\left(\sum_{j=0}^{m} \alpha_j p^j\right) - b^{\mathrm{T}}\left(\sum_{i=0}^{m} \alpha_i p^i\right)$$

$$= \frac{1}{2}\left(\sum_{i=0}^{m}\sum_{j=0}^{m} \alpha_i \alpha_j (p^i)^{\mathrm{T}} Ap^j\right) - \left(\sum_{i=0}^{m} \alpha_i b^{\mathrm{T}} p^i\right)$$

$$= \frac{1}{2}\left(\sum_{i=0}^{m} \alpha_i^2 (p^i)^{\mathrm{T}} Ap^i\right) - \left(\sum_{i=0}^{m} \alpha_i b^{\mathrm{T}} p^i\right) \qquad \text{(from conjugacy)}$$

$$= \left(\sum_{i=0}^{m} \left(\frac{1}{2}\alpha_i^2 (p^i)^{\mathrm{T}} Ap^i - \alpha_i b^{\mathrm{T}} p^i\right)\right)$$

Minimizing this is easy since each term is separable and setting the derivative w.r.t. each $\alpha_i$ to 0 we get $\alpha_i (p^i)^{\mathrm{T}} Ap^i - b^{\mathrm{T}} p^i = 0$, i.e. $\alpha_i = \dfrac{b^{\mathrm{T}} p^i}{(p^i)^{\mathrm{T}} Ap^i}$

Since we don't know the conjugate vectors $p^i$ in advance the conjugate gradient method <u>determines these iteratively, along with their coefficients</u> $\alpha_i$.

Suppose we define $r^i = b - Ax^i$ as the residual vector at iteration $i$ and a scalar quantity $\beta_i$ that will be used to determine $p^i$.

Then the conjugate gradient algorithm may be defined as follows:

1. Start with $x^0=0$ so that $r^0=b$. Let $\beta_0=0$, $p^0=b$ and $\alpha_0 = \dfrac{b^{\mathrm{T}}p^0}{(p^0)^{\mathrm{T}}Ap^0}$

2. Find $x^1= x^0+\alpha_0 p^0$ and compute $r^1=b-Ax^1$. Then set iteration counter $i=1$

3. STOP if $\|r^i\| < \varepsilon$ (some suitable tolerance)

4. Set $\beta_i=\dfrac{(r^i)^{\mathrm{T}}r^i}{(r^{i-1})^{\mathrm{T}}r^{i-1}}$

5. Set $p^i = r^i + \beta_i p^{i-1}$

6. Set $\alpha_i = \dfrac{(r^i)^{\mathrm{T}}r^i}{(p^i)^{\mathrm{T}}Ap^i}$

7. Set $x^{i+1}= x^i +\alpha_i p^i$

8. Set $r^{i+1}= r^i-\alpha_i Ap^i$ and return to Step 3.

Notes:

1) If we wished to, we could also compute (using the original formulae...)

$$r^{i+1}=b-Ax^{i+1}, \text{ because } r^i-\alpha_i Ap^i = (b-Ax^i) -\alpha_i Ap^i = b-A(x^i+\alpha_i p^i)= b-Ax^{i+1}$$

$$\alpha_i = \frac{b^{\mathrm{T}}p^i}{(p^i)^{\mathrm{T}}Ap^i}, \text{ because } b^{\mathrm{T}}p^i = (r^i)^{\mathrm{T}}r^i$$

We use the formulae in (6) and (8) only for computational efficiency.

2) $x^{i+1}= x^i +\alpha_i p^i+\alpha_{i-1} p^{i-1}+\alpha_{i-2} p^{i-2}...+\alpha_0 p^0$

So these approximate solutions at each step have the same form as the vector $y$ on the previous page.

**THEOREM**: For $i>j$, the vectors $\{p^i\}$ and $\{r^i\}$ computed in the above algorithm satisfy

$(r^i)^{\mathrm{T}}r^j= 0$            (residuals are orthogonal to each other)

$(r^i)^{\mathrm{T}}p^j= 0,$           (residuals are orthogonal to the search directions)

$(p^i)^{\mathrm{T}}Ap^j= 0$          (search directions are conjugate w.r.t. $A$)

# THE FLETCHER REEVES ALGORITHM

### (also called the nonlinear CONJUGATE-GRADIENT method)

This is an algorithm that uses conjugate search directions:

Here $r^i \equiv b - Ax^i = -\nabla f(x^0)$ if the function is quadratic ($\frac{1}{2} x^T Ax - b^T x$)

Start with $d^0 = -\nabla f(x^0)$          (steepest descent dir.),

Minimize along $d^0$, i.e.,

Minimize $f(\alpha) = \{f(x^0 + \alpha d^0)\}$ to obtain $\alpha_0$ and hence $x^1 = x^0 + \alpha_0 d^0$

The next search direction is computed by

$$d^1 = -\nabla f(x^1) + \frac{\|\nabla f(x^1)\|^2}{\|\nabla f(x^0)\|^2} \cdot d^0$$

(steepest descent dir.) *plus* (a multiple of the previous search dir.)

Here $\|\nabla f(x)\|$ denotes the length of the vector $\nabla f(x)$,

i.e., $\|\nabla f(x)\| = \sqrt{\left(\frac{\partial f}{\partial x_1}\right)^2 + \left(\frac{\partial f}{\partial x_2}\right)^2 + \dots + \left(\frac{\partial f}{\partial x_n}\right)^2}$.

The procedure could be generalized (at iteration $i$),

$$d^{i+1} = -\nabla f(x^{i+1}) + \frac{\|\nabla f(x^{i+1})\|^2}{\|\nabla f(x^i)\|^2} \cdot d^i$$

How does this relate to conjugate directions?

Suppose that $f(x)$ is quadratic with Hessian $H$.

Then for the vectors $d^0, d^1, ..., d^{n-1}$

$$[d^i]^T H d^j = 0 \quad \text{for } i \neq j, \; i,j \in \{0,1,...,n\text{-}1\},$$

i.e., the vectors are all $H$-conjugate. Note that the directions are not specified

beforehand; rather they are determined sequentially at each step.

Some Features of the method

(1) $\nabla f(x^i)$ is always non-zero and linearly independent of previous directions, unless

the solution is obtained in less than $n$ steps (in which case of course, we are done!)

(2) The formula for finding $d^{i+1}$ is the simplest of all conjugate directions methods.

(3) The path taken usually shows *uniformly* good progress, and very often avoids

"zig-zagging."

(4) The minimization procedure along direction $d^k$ has a closed form solution with

$\alpha_k \equiv$ step size in direction $d^k$ at iteration $k$,

$$= \frac{-[\nabla f(x^k)]^T d^k}{[d^k]^T H d^k} \qquad \textit{as long as } f(x) \textit{ is quadratic}$$

(5) It is well-suited to large problems because there are no matrix operations and

storage requirements are not very high (we only need the last two gradients).

(6) If a "bad" direction is generated at some stage, there is a tendency for the next

direction to also be bad; so the algorithm will tend to crawl when this happened.

To avoid this, there is always some "restart" strategy used when this method is

implemented (e.g., restart with the steepest descent dir. if we seem to be stuck).

<u>AN EXAMPLE</u>    Consider the previous example solved by the steepest descent

method:              $\text{Min } f(x) = x_1{}^2 + 25x_2{}^2$

Recall that $\nabla f(x) = \begin{bmatrix} 2x_1 \\ 50x_2 \end{bmatrix}$ and $H = \begin{bmatrix} 2 & 0 \\ 0 & 50 \end{bmatrix}$

Start with $[2\ \ 2]^T$

<u>ITERATION 1</u>

$$x^0 = \begin{bmatrix} 2 \\ 2 \end{bmatrix} \qquad d^0 = -\nabla f(x) = \begin{bmatrix} -4 \\ -100 \end{bmatrix}; \quad x^1 = x^0 + \alpha_0 d^0 = \begin{bmatrix} 2 - 4\alpha_0 \\ 2 - 100\alpha_0 \end{bmatrix}$$

To obtain $\alpha_0$ we minimize $f(2-4\alpha_0, 2-100\alpha_0)$. Thus,

$$\alpha_0 = \frac{-[\nabla f(x^0)]^T]d^0}{[d^0]^T H d^0} = \frac{-[4\ \ 100]\begin{bmatrix} -4 \\ -100 \end{bmatrix}}{[-4 \ -100]\begin{bmatrix} 2 & 0 \\ 0 & 50 \end{bmatrix}\begin{bmatrix} -4 \\ -100 \end{bmatrix}} = 0.02003$$

$$\therefore \ x^1 = \begin{bmatrix} 2 \\ 2 \end{bmatrix} + 0.02003 \begin{bmatrix} -4 \\ -100 \end{bmatrix} = \begin{bmatrix} 1.91988 \\ -0.00307 \end{bmatrix}$$

To summarize: $x^0 = \begin{bmatrix} 2 \\ 2 \end{bmatrix}$; $d^0 = \begin{bmatrix} -4 \\ -100 \end{bmatrix}$; $\alpha_0 = 0.02003$

<u>ITERATION 2</u>

$$x^1 = \begin{bmatrix} 1.91988 \\ -0.00307 \end{bmatrix} \qquad d^1 = -\nabla f(x^1) + \frac{\|\nabla f(x^1)\|^2}{\|\nabla f(x^0)\|^2} \cdot d^0, \text{ i.e.,}$$

$$= -\begin{bmatrix} 2 * 1.91988 \\ 50 * (-0.00307) \end{bmatrix} + \left\{ \frac{\{2(1.91988)\}^2 + \{50(-0.00307)\}^2}{4^2 + 100^2} \right\} \cdot \begin{bmatrix} -4 \\ -100 \end{bmatrix}$$

$$= \begin{bmatrix} -3.84565 \\ 0.00615 \end{bmatrix}$$

$$\alpha_1 = \frac{-[\nabla f(x^1)]^T]d^1}{[d^1]^T H d^1} = 0.499$$

$$\therefore \ x^2 = \begin{bmatrix} 1.91988 \\ -0.00307 \end{bmatrix} + 0.499 \begin{bmatrix} -3.84565 \\ 0.00615 \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \end{bmatrix}$$

This must be the optimum solution since $\nabla f(x^2)=0$.

To summarize:

| Iteration ($i$) | $x^{i-1}$ | $d^{i-1}$ | $\alpha_{i-1}$ | $x^1$ |
|---|---|---|---|---|
| 1 | $\begin{bmatrix} 2 \\ 2 \end{bmatrix}$ | $\begin{bmatrix} -4 \\ -100 \end{bmatrix}$ | 0.02003 | $\begin{bmatrix} 1.91988 \\ -0.00307 \end{bmatrix}$ |
| 2 | $\begin{bmatrix} 1.91988 \\ -0.00307 \end{bmatrix}$ | $\begin{bmatrix} -3.84565 \\ 0.00615 \end{bmatrix}$ | 0.499 | $\begin{bmatrix} 0 \\ 0 \end{bmatrix}$ |

Note that with $n=2$, only 2 search directions were required.  (Try with a different $x^0$...)


## THE POLAK-RIBIERE METHOD

This is similar to the Fletcher-Reeves method except that the multiple of the old

search direction used is

$$\frac{\left[\nabla f(x^{i+1}) - \nabla f(x^i)\right]^{\mathrm{T}} \nabla f(x^{i+1})}{\|\nabla f(x^i)\|^2}$$

instead of $\dfrac{\|\nabla f(x^{i+1})\|^2}{\|\nabla f(x^i)\|^2}$


## THE HESTENES-STEIFEL METHOD

This is also similar but the multiple of the old search direction used is

$$\frac{\left[\nabla f(x^{i+1}) - \nabla f(x^i)\right]^{\mathrm{T}} \nabla f(x^{i+1})}{[\nabla f(x^{i+1}) - \nabla f(x^i)]^{\mathrm{T}} \nabla f(x^i)}$$

instead of $\dfrac{\|\nabla f(x^{i+1})\|^2}{\|\nabla f(x^i)\|^2}$