

Unconstrained Nonlinear Optimisation

Georgy Gimel'farb

COMPSCI 369 Computational Science

- ① Non-linear world
- ② Extremum points
- ③ Univariate search
- ④ Gradient methods
- ⑤ Direct search

RECOMMENDED READING:

- G. Strang, *Computational Science and Engineering*. Wellesley-Cambridge Press, 2007: Section 2.6
- W. H. Press e. a., *Numerical Recipes: The Art of Scientific Computing*. Cambridge Univ. Press, 2007: Section 15.5
- L. R. Foulds: *Optimization Techniques: An Introduction*. Springer-Verlag, 1981: Chapters 7, 8
- C. Woodford, C. Phillips: *Numerical Methods with Worked Examples*. Chapman & Hall, 1997: Chapter 7.
- M. T. Heath: *Scientific Computing: An Introductory Survey*. McGraw-Hill, 2002: Chapter 6

Learning Outcomes on Nonlinear Optimisation

Be familiar with unconstrained and constrained optimisation:

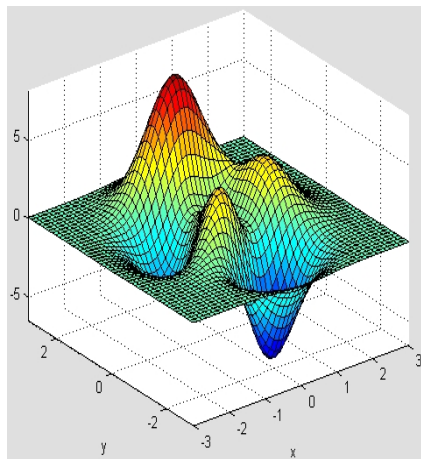
- Recognise discrete and continuous optimisation problems.
- Understand the method of Lagrange for optimising a function of many variables subject to a system of equality constraints.
- Be able to implement a simple search for the maximum (or minimum) of a function of one variable.
- Understand unconstrained gradient optimisation.
- Understand dynamic programming (DP) and be familiar with DP solution of the edit-distance problem.
- Understand that heuristic optimisation strategies must be used when no good exact algorithm is known.

Understand how a wide range of scientific questions can be answered by combining optimisation and other numerical methods.

Nature is nonlinear...



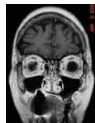
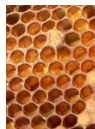
4896kj.com/journeying/wp-content/uploads/2007/06/



www.mathworks.com/matlabcentral/fx_files/

Why Optimisation?

- Many physical laws are formulated in terms of maximum or minimum of energy, entropy, or other object properties
- Optimisation is essential almost everywhere (in natural life, science, engineering, economics, business, etc).
 - Bees minimise the average amount of material per cell
 - Humans aim to the maximum effect under limited resources
- Many scientific and technological applications require large-scale optimisation with 10^4 .. 10^6 or more variables:
 - Medical imaging and computer assisted diagnostics
 - Shape design of mechanical objects
 - Control of industrial processes
 - Robotics and autonomous navigation, etc.



Extremum of a Function

One of most important applied problems: to find maximum or minimum value of a function $f(\mathbf{x})$ under constraints $\mathbf{x} \in \mathbf{X}$

- $f(\mathbf{x}) \equiv f(x_1, \dots, x_n)$ is a scalar function of n -dimensional vector argument
- \mathbf{X} is a certain subset of n -dimensional vector space \mathbb{R}_n

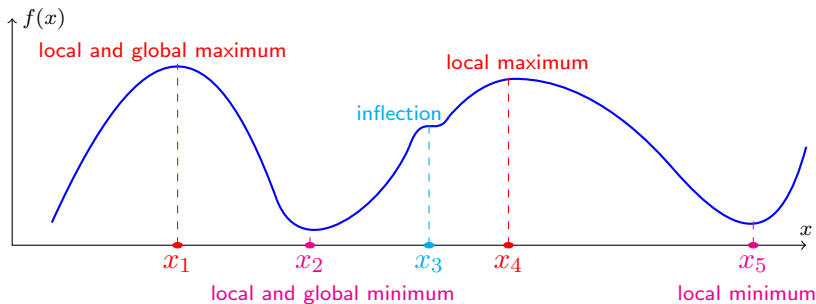
Unconstrained optimisation: if $\mathbf{X} = \mathbb{R}_n$

- Maximum / minimum function value: $f^* = \{\max_{\min}\} f(\mathbf{x})$
- Maximiser / minimiser: $\mathbf{x}^* = \arg \{\max_{\min}\} f(\mathbf{x})$

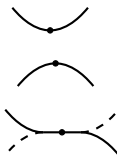
Constrained optimisation: if $\mathbf{X} \subset \mathbb{R}_n$

- Maximum / minimum function value: $f^* = \{\max_{\min}\}_{\mathbf{x} \in \mathbf{X}} f(\mathbf{x})$
- Maximiser / minimiser: $\mathbf{x}^* = \arg \{\max_{\min}\}_{\mathbf{x} \in \mathbf{X}} f(\mathbf{x})$

Univariate Function (Function of One Variable)



- Minimum: $\frac{df(x)}{dx} = 0$; $\frac{d^2f(x)}{dx^2} > 0$
- Maximum: $\frac{df(x)}{dx} = 0$; $\frac{d^2f(x)}{dx^2} < 0$
- Inflection point: $\frac{df(x)}{dx} = 0$; $\frac{d^2f(x)}{dx^2} = 0$



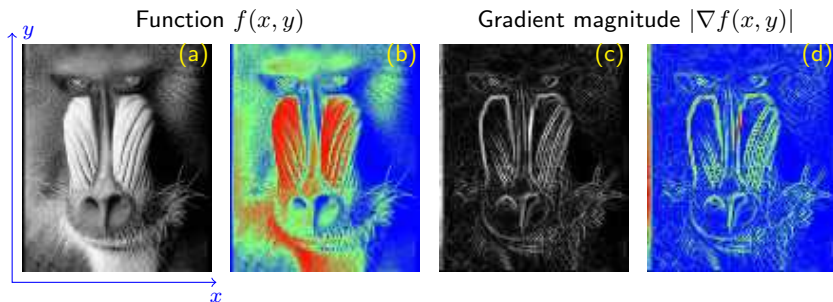
Functions of Many Variables $f(\mathbf{x})$

- Unconditional local critical (or stationary) point: where the **gradient** of f is zero: $\nabla f(\mathbf{x}) = \left[\frac{\partial f(\mathbf{x})}{\partial x_1}, \dots, \frac{\partial f(\mathbf{x})}{\partial x_n} \right]^T = \mathbf{0}$
- Its properties depend on the **Hessian** matrix of the 2nd derivatives:

$$\mathbf{H}(\mathbf{x}) = \begin{bmatrix} \frac{\partial^2 f(\mathbf{x})}{\partial x_1^2} & \frac{\partial^2 f(\mathbf{x})}{\partial x_1 \partial x_2} & \cdots & \frac{\partial^2 f(\mathbf{x})}{\partial x_1 \partial x_n} \\ \frac{\partial^2 f(\mathbf{x})}{\partial x_2 \partial x_1} & \frac{\partial^2 f(\mathbf{x})}{\partial x_2^2} & \cdots & \frac{\partial^2 f(\mathbf{x})}{\partial x_2 \partial x_n} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial^2 f(\mathbf{x})}{\partial x_n \partial x_1} & \frac{\partial^2 f(\mathbf{x})}{\partial x_n \partial x_2} & \cdots & \frac{\partial^2 f(\mathbf{x})}{\partial x_n^2} \end{bmatrix}$$

- **Local minimum**: if the Hessian is positive definite (the quadratic form $\mathbf{e}^T \mathbf{H}(\mathbf{x}) \mathbf{e} > 0$ for any $\mathbf{e} \neq \mathbf{0}$)
- **Local maximum**: if the Hessian is negative definite (the quadratic form $\mathbf{e}^T \mathbf{H}(\mathbf{x}) \mathbf{e} < 0$ for any $\mathbf{e} \neq \mathbf{0}$)
- **Saddle point**: if the Hessian is indefinite

Facial Image as a Graph of Function of Two Variables



- (a,c) Linear grayscale coding:
 - (a) 0 (black) – 127 (grey) – 255 (white)
 - (c) min = 0 (black) – 55 (grey) – max = 110 (white)
- (b,d) Linear colour coding:
 - (b) 0 (blue) – 127 (green) – 255 (red)
 - (d) min = 0 (blue) – 55 (green) – max = 110 (red)

Quadratic Function $f(\mathbf{x}) = \mathbf{a}^\top \mathbf{x} + \frac{1}{2} \mathbf{x}^\top \mathbf{H} \mathbf{x}$ of n Variables

$$\begin{aligned}
 f(\mathbf{x}) &= \overbrace{a_1 x_1 + \dots + a_n x_n}^{\mathbf{a}^\top \mathbf{x}} + \overbrace{\frac{1}{2} (H_{11} x_1^2 + H_{12} x_1 x_2 + \dots + H_{1n} x_1 x_n \\
 &\quad + H_{21} x_2 x_1 + H_{22} x_2^2 + \dots + H_{2n} x_2 x_n \\
 &\quad \dots + H_{n1} x_n x_1 + H_{n2} x_n x_2 + \dots + H_{nn} x_n^2)}^{\frac{1}{2} \mathbf{x}^\top \mathbf{H} \mathbf{x}} \\
 &= \sum_{i=1}^n \left(a_i x_i + \frac{H_{ii}}{2} x_i^2 + \sum_{j=i+1}^n H_{ij} x_i x_j \right) \quad \text{where } H_{ij} = H_{ji}
 \end{aligned}$$

- Gradient $\nabla f(\mathbf{x}) = \mathbf{a} + \mathbf{H} \mathbf{x}$:

$$\begin{bmatrix} \frac{\partial f}{\partial x_1} \\ \vdots \\ \frac{\partial f}{\partial x_n} \end{bmatrix} = \begin{bmatrix} a_1 \\ \vdots \\ a_n \end{bmatrix} + \begin{bmatrix} H_{11} & H_{12} & \dots & H_{1n} \\ \vdots & \vdots & \ddots & \vdots \\ H_{n1} & H_{n2} & \dots & H_{nn} \end{bmatrix} \begin{bmatrix} x_1 \\ \vdots \\ x_n \end{bmatrix}$$

- Hessian $\frac{\partial \nabla f(\mathbf{x})}{\partial \mathbf{x}} \equiv \left[\frac{\partial^2 f}{\partial x_i \partial x_j} \right]_{i,j=1}^n \equiv [H_{ij}]_{i,j=1}^n \equiv \mathbf{H}$

Sylvester's Criterion: Sufficient Condition of an Extremum

Symmetric $n \times n$ matrix $\mathbf{A} = \begin{bmatrix} a_{11} & \dots & a_{1n} \\ \vdots & \ddots & \vdots \\ a_{n1} & \dots & a_{nn} \end{bmatrix}$ is:

- **Positive definite** if the determinants of all its upper-left submatrices are positive:

$$a_{11} > 0; \begin{vmatrix} a_{11} & a_{12} \\ a_{21} & a_{22} \end{vmatrix} > 0; \dots; \begin{vmatrix} a_{11} & \dots & a_{1n} \\ \vdots & \ddots & \vdots \\ a_{n1} & \dots & a_{nn} \end{vmatrix} > 0$$

- **Negative definite** if sign-alternate (starting from $a_{11} < 0$):

$$-a_{11} > 0; \begin{vmatrix} a_{11} & a_{12} \\ a_{21} & a_{22} \end{vmatrix} > 0; \dots; (-1)^n \begin{vmatrix} a_{11} & \dots & a_{1n} \\ \vdots & \ddots & \vdots \\ a_{n1} & \dots & a_{nn} \end{vmatrix} > 0$$

- **Indefinite** in all other cases

Quadratic Function: An Example with No Local Extremum

$$f(x_1, x_2, x_3) = 5x_1 - 3x_2 + x_3 + 4x_1^2 - 2x_1x_2 + 10x_1x_3 - x_2^2 + 3x_2x_3 + 9x_3^2$$

$$\frac{\partial f}{\partial x_1} = 5 + 8x_1 - 2x_2 + 10x_3 \quad \frac{\partial^2 f}{\partial x_1^2} = 8 \quad \frac{\partial^2 f}{\partial x_1 \partial x_2} = -2 \quad \frac{\partial^2 f}{\partial x_1 \partial x_3} = 10$$

$$\frac{\partial f}{\partial x_2} = -3 - 2x_1 - 2x_2 + 3x_3 \quad \frac{\partial^2 f}{\partial x_1 \partial x_2} = -2 \quad \frac{\partial^2 f}{\partial x_2^2} = -2 \quad \frac{\partial^2 f}{\partial x_2 \partial x_3} = 3$$

$$\frac{\partial f}{\partial x_3} = 1 + 10x_1 + 3x_2 + 18x_3 \quad \frac{\partial^2 f}{\partial x_1 \partial x_3} = 10 \quad \frac{\partial^2 f}{\partial x_2 \partial x_3} = 3 \quad \frac{\partial^2 f}{\partial x_3^2} = 18$$

Gradient:

$$\nabla f(x_1, x_2, x_3) = \begin{bmatrix} \frac{\partial f}{\partial x_1} \\ \frac{\partial f}{\partial x_2} \\ \frac{\partial f}{\partial x_3} \end{bmatrix} = \begin{bmatrix} 5 \\ -3 \\ 1 \end{bmatrix} + \underbrace{\begin{bmatrix} 8 & -2 & 10 \\ -2 & -2 & 3 \\ 10 & 3 & 18 \end{bmatrix}}_{\text{the indefinite Hessian } \mathbf{H}} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix}$$

Sylvester's criterion:

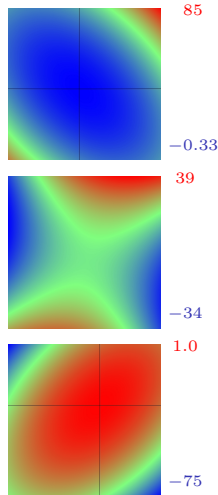
$$H_{11} = 8 > 0; \quad \begin{vmatrix} 8 & -2 \\ -2 & -2 \end{vmatrix} = -12 < 0; \quad \begin{vmatrix} 8 & -2 & 10 \\ -2 & -2 & 3 \\ 10 & 3 & 18 \end{vmatrix} = -120 < 0$$

Quadratic Functions of Two Variables

- $f(x, y) = x^2 - xy + y^2 + x - y$:
 the positive definite Hessian $\begin{bmatrix} 2 & -1 \\ -1 & 2 \end{bmatrix}$ as
 $2 > 0; \begin{vmatrix} 2 & 1 \\ 1 & 2 \end{vmatrix} = 2 \cdot 2 - 1 \cdot 1 = 4 - 1 = 3 > 0$
- $f(x, y) = -x^2 - xy + y^2 + x - y$:
 the indefinite Hessian $\begin{bmatrix} -2 & -1 \\ -1 & 2 \end{bmatrix}$ as
 $-2 < 0; \begin{vmatrix} -2 & 1 \\ 1 & 2 \end{vmatrix} = -2 \cdot 2 - 1 \cdot 1 = -4 - 1 = -5 < 0$
- $f(x, y) = -x^2 - xy - y^2 + x - y$:
 the negative definite Hessian $\begin{bmatrix} -2 & -1 \\ -1 & -2 \end{bmatrix}$ as
 $-2 < 0;$
 $\begin{vmatrix} -2 & -1 \\ -1 & -2 \end{vmatrix} = -2 \cdot (-2) - (-1) \cdot (-1) = 4 - 1 = 3 > 0$

Colour-coded range

$[f_{\min}, f_{\max}]$



Equivalent Definitions of a Positive Definite Matrix

Symmetric $n \times n$ matrix \mathbf{A} is positive definite if

- 1 All eigenvalues of \mathbf{A} are positive (> 0), or
- 2 Choleski decomposition $\mathbf{A} = \mathbf{L}\mathbf{L}^T$ exists where \mathbf{L} is a lower triangular matrix with $l_{ii} > 0$, or
- 3 Decomposition $\mathbf{A} = \mathbf{L}\mathbf{D}\mathbf{L}^T$ exists where \mathbf{L} is a lower triangular matrix with $l_{ii} = 1$ and \mathbf{D} is a diagonal matrix with $d_i > 0$, or
- 4 All pivots in Gaussian elimination without pivoting are positive (> 0).

For small matrices, the Sylvester's condition can be readily checked, but generally Conditions 2 or 3 are the most efficient and yield easy solutions to linear systems with coefficients \mathbf{A}

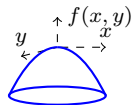
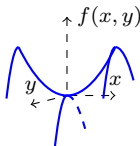
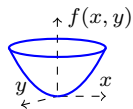
Extrema of $f(x, y) = ax^2 + by^2$; $a \neq 0$; $b \neq 0$

Gradient $\nabla f(x, y) = 0 \Rightarrow \frac{\partial f(x, y)}{\partial x} = 2ax = 0$; $\frac{\partial f(x, y)}{\partial y} = 2by = 0$

\Rightarrow Single extremum at the point $[0, 0]^T$

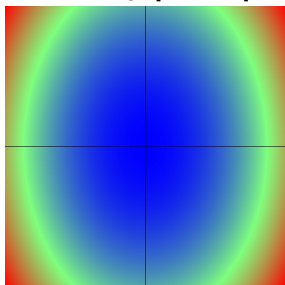
$$\text{Hessian } \mathbf{H} = \begin{bmatrix} 2a & 0 \\ 0 & 2b \end{bmatrix}$$

- Function $f(x, y)$ has the minimum in $[0, 0]^T$ if $a > 0$ and $b > 0$: **an elliptic paraboloid** (the Sylvester's criterion: $2a > 0$ and $4ab > 0$)
- If $a > 0$; $b < 0$ or $a < 0$; $b > 0$, there is no extremum: **a hyperbolic paraboloid** with a saddle point $[0, 0]^T$ (the Sylvester's criterion: $2a > 0$ and $4ab < 0$ or $2a < 0$ and $4ab < 0$)
- Function $f(x, y)$ has the maximum in $[0, 0]^T$ if $a < 0$ and $b < 0$: **an elliptic paraboloid** (the Sylvester's criterion: $2a < 0$ and $4ab > 0$)



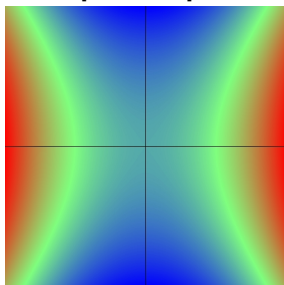
Colour-Coded Graphs of $f(x, y) = ax^2 + by^2$; $a \neq 0$; $b \neq 0$

Function range [0.0, 3.75]



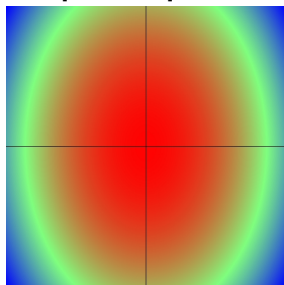
$$f(x, y) = 0.1x^2 + 0.05y^2$$

Function range [-1.25, 2.5]



$$f(x, y) = 0.1x^2 - 0.05y^2$$

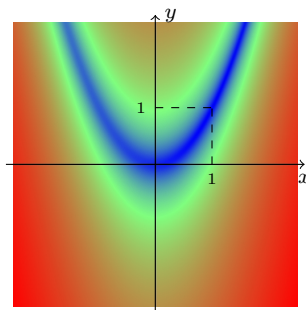
Function range [-3.75, 0.0]



$$f(x, y) = -0.1x^2 - 0.05y^2$$

- The single potentially extremum point $[x = 0, y = 0]$ such that $\nabla f(0, 0) = \mathbf{0}$
- (left) Minimum – the positive definite $\mathbf{H} = \begin{bmatrix} 0.2 & 0 \\ 0 & 0.1 \end{bmatrix}$
- (middle) Saddle – the indefinite $\mathbf{H} = \begin{bmatrix} 0.2 & 0 \\ 0 & -0.1 \end{bmatrix}$
- (right) Maximum – the negative definite $\mathbf{H} = \begin{bmatrix} -0.2 & 0 \\ 0 & -0.1 \end{bmatrix}$

A More Complex Case: the Rosenbrock's "Banana" Test Function



Logarithmic colour scales

Left window

$$-2.5 \leq x, y \leq 2.5:$$

$$f(1, 1) = 0; f(0, 0) = 1;$$

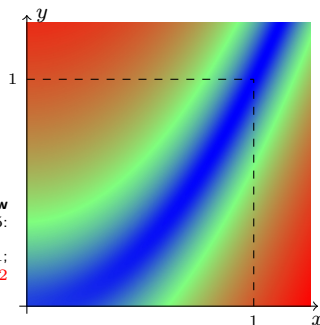
$$f(-2.5, -2.5) = 7668.5$$

Right window

$$0 \leq x, y \leq 1.25:$$

$$f(1, 1) = 0; f(0, 0) = 1;$$

$$f(1.25, 0) = 244.2$$



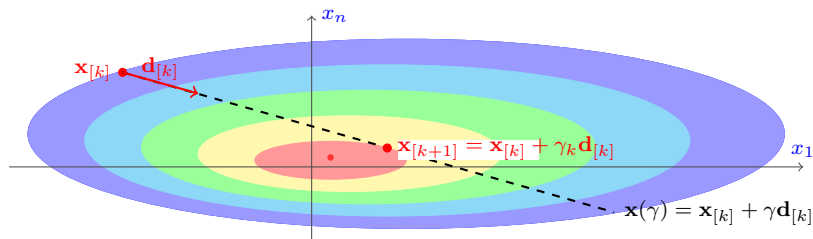
- Bivariate case: $f(x, y) = (1 - x)^2 + 100(y - x^2)^2$

$$\nabla f = \begin{bmatrix} -2(1 - x) - 400x(y - x^2) \\ 200(y - x^2) \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \end{bmatrix} \Rightarrow \text{Single minimum } f(1, 1) = 0$$
- Multivariate extension: $f(x_1, \dots, x_n) = \sum_{i=1}^{n-1} \left[(1 - x_i)^2 + 100(x_{i+1} - x_i^2)^2 \right]$
 - For $n = 3$: Single minimum at $[1, 1, 1]$
 - For $4 \leq n \leq 7$: Global minimum at all ones $(1, 1, \dots, 1)$ and local minimum near $(-1, 1, \dots, 1)$

Line Search for a Maximal Point

Find a maximiser of $f(\mathbf{x})$ along a direction $\mathbf{d}_{[k]}$ from a point $\mathbf{x}_{[k]}$, i.e. along the line $\mathbf{x}(\gamma) = \mathbf{x}_{[k]} + \gamma\mathbf{d}_{[k]}$:

$$\begin{cases} \gamma_{[k]} &= \arg \max_{\gamma \in \mathbb{R}} f(\mathbf{x}(\gamma)) \equiv \arg \max_{\gamma \in \mathbb{R}} f(\mathbf{x}_{[k]} + \gamma\mathbf{d}_{[k]}) \\ \mathbf{x}_{[k+1]} &= \mathbf{x}_{[k]} + \gamma_{[k]}\mathbf{d}_{[k]} \end{cases}$$



Line search is used repeatedly in many multivariate search methods

Line Search for a Maximal Point: An Example

Find a maximiser of the quadratic function

$$f(x, y) = -x^2 + xy - y^2 + 1$$

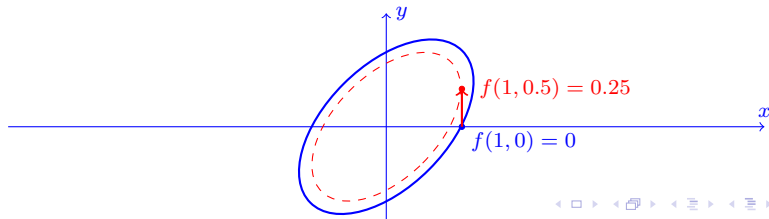
along a direction $[\delta_x = 0, \delta_y = 1]$ from the point $[x^\circ = 1, y^\circ = 0]$

Maximising along the line $[x^\circ + \gamma\delta_x, y^\circ + \gamma\delta_y]$, i.e. the line $[1 + 0 \cdot \gamma, 0 + 1 \cdot \gamma] = [1, \gamma]$ in this example:

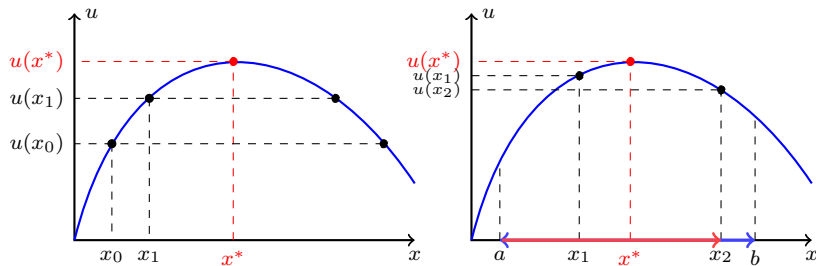
$$f(1, \gamma) = -1 + \gamma - \gamma^2 + 1 = \gamma - \gamma^2$$

$$\gamma^* = \arg \max_{\gamma \in \mathbb{R}} \{\gamma - \gamma^2\} = 0.5 \leftarrow \frac{d(\gamma - \gamma^2)}{d\gamma} = 1 - 2\gamma = 0$$

$$x^* = 1 \quad ; \quad y^* = 0.5 \Rightarrow f(1, 0.5) = 0.25 > f(1, 0) = 0$$



Maximising an Univariate Unimodal Function $u(x)$



Properties of the maximiser $x^* = \arg \max_x u(x)$

- If $x_0 < x_1 < x^*$ or $x_0 > x_1 > x^*$, then $u(x_0) < u(x_1) < u(x^*)$
- If $a \leq x^* \leq b$ and $a \leq x_1 < x_2 < b$ or $a < x_1 < x_2 \leq b$, then

$$\left. \begin{aligned} u(x_1) < u(x_2) &\Rightarrow x_1 < x^* \leq b \\ u(x_1) = u(x_2) &\Rightarrow x_1 < x^* < x_2 \\ u(x_1) > u(x_2) &\Rightarrow a \leq x^* < x_2 \end{aligned} \right\}$$

Therefore, the search interval $a \leq x^* \leq b$ is reduced

Golden Section Search

Initialisation: an interval $[a_0, b_0]$; $a_0 \leq x^* \leq b_0$, containing the maximiser x^* of a unimodal function $u(x)$

Iteration (step) $i = 0, 1, 2, \dots$, of reducing the search interval:

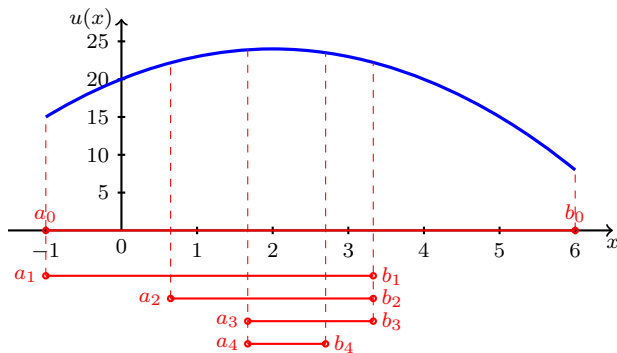
- Evaluate $u(x_i)$ and $u(\bar{x}_i)$ at symmetric internal points of the interval $[a_i, b_i]$ (below $\tau = \frac{1+\sqrt{5}}{2}$ is the Greek **golden section** ratio):

$$\begin{cases} x_i = a_i + (b_i - a_i)(2 - \tau) \approx a_i + 0.382(b_i - a_i) \\ \bar{x}_i = a_i + (b_i - a_i)(\tau - 1) \approx a_i + 0.618(b_i - a_i) \end{cases}$$

- If $u(x_i) > u(\bar{x}_i)$, then $a_{i+1} \leftarrow a_i$ and $b_{i+1} \leftarrow \bar{x}_i$
- If $u(x_i) < u(\bar{x}_i)$, then $a_{i+1} \leftarrow x_i$ and $b_{i+1} \leftarrow b_i$
- If $u(x_i) = u(\bar{x}_i)$, then $a_0 \leftarrow x_i$; $b_0 \leftarrow \bar{x}_i$, and initiate the search again from the new interval $[a_0, b_0]$ and $i = 0$

Stopping rule: Proceed until the interval $[a_0, b_0]$ is sufficiently small, or the next point is within the resolution distance of the last point

Golden Section Search: A Simple Example



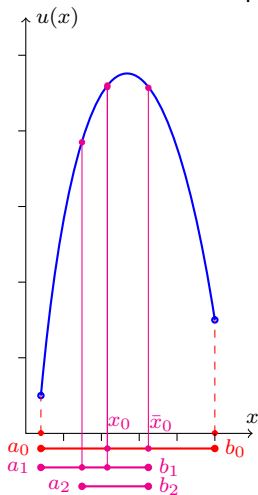
$$u(x) = 20 + 4x - x^2$$

$$x^* = 2; u(x^*) = 24$$

i	0	1	2	3	4	5	6	7	8	9	10
a_i	-1.0	-1.00	0.65	1.67	1.67	1.67	1.92	1.92	1.92	1.97	1.97
b_i	6.0	3.33	3.33	3.33	2.70	2.31	2.31	2.16	2.06	2.06	2.03
$u(a_i)$	15.	15.0	22.2	23.9	23.9	23.9	23.9 ₉	23.9 ₉	23.9 ₉	24.0	24.0
$u(b_i)$	8.	22.2	22.2	22.2	23.5	23.9	23.9 ₁	23.9 ₈	23.9 ₉	23.9 ₉	24.0

Golden Section vs. Fibonacci Search (optional)

Golden section search: an example



- Golden section search is less efficient than the **Fibonacci search**: for $i = 1, 2, \dots, n - 1$,

$$x_i = a_i + (b_i - a_i)F_{n-i}/F_{n+2-i}$$

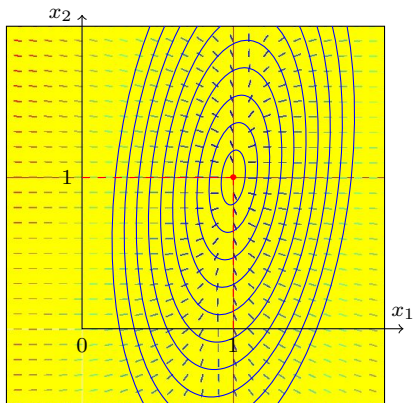
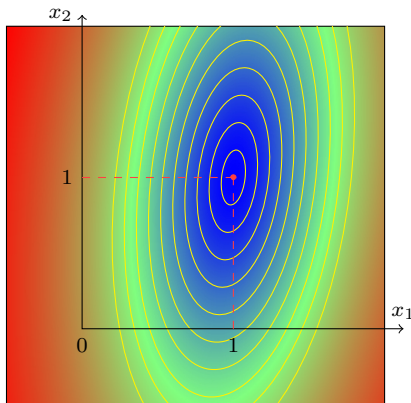
$$\bar{x}_i = a_i + (b_i - a_i)F_{n+1-i}/F_{n+2-i}$$

where F_k is the Fibonacci number: $F_0 = 0$;
 $F_1 = 1$; $F_k = F_{k-1} + F_{k-2}$, $k = 2, 3, \dots$

- Fibonacci search minimises the maximal interval of uncertainty about the maximiser x^* (in that sense it is optimal)
- But the number of points n to be evaluated in the Fibonacci search has to be prescribed
- Search for the root x^* of the first derivative, $\frac{du}{dx}(x^*) = 0$, be it available, is even more efficient

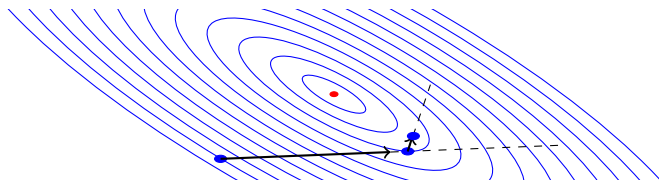
Gradient of a Function of Several Variables

Gradient vector $\nabla f(\mathbf{x}) = \left[\frac{\partial f}{\partial x_1}, \frac{\partial f}{\partial x_2}, \dots, \frac{\partial f}{\partial x_n} \right]^T$ is directed to the greatest slope of the function f at any point.



$$f(x_1, x_2) = 10x_1^2 - 2x_1x_2 + 2x_2^2 - 18x_1 - 2x_2; \quad \frac{\partial f}{\partial x_1} = 20x_1 - 2x_2 - 18; \quad \frac{\partial f}{\partial x_2} = -2x_1 + 4x_2 - 2$$

Gradient Search



Gradient methods for seeking a maximum (minimum) for f :

- 1 Evaluate the gradient at an initial point
- 2 Move along the gradient direction for a computable distance
- 3 Repeat this process until the maximum (minimum) is found

If exact partial derivatives are unknown, the gradients may be numerically approximated: for a small $\varepsilon > 0$

$$\frac{\partial f}{\partial x_i} \approx \frac{f(x_1, \dots, x_{i-1}, x_i + \varepsilon, x_{i+1}, \dots, x_n) - f(x_1, \dots, x_{i-1}, x_i - \varepsilon, x_{i+1}, \dots, x_n)}{2\varepsilon}$$

- Approximation errors can make the methods less attractive

Gradient Maximisation: The Steepest Ascent

Initialisation: Select an initial point \mathbf{x}_0

Iterative steps $i = 0, 1, \dots$:

- 1 Compute the gradient $\nabla f(\mathbf{x})$ at the point \mathbf{x}_i
- 2 Draw a line $\mathbf{x}_i + t\nabla f(\mathbf{x}_i)$ through \mathbf{x}_i in the gradient direction
- 3 Select the point \mathbf{x}_{i+1} on this line yielding the largest value for f of all points on the line:

$$f(\mathbf{x}_{i+1}) = \max_{t \in (0, \infty)} \{f(\mathbf{x}) : \mathbf{x} = \mathbf{x}_i + t\nabla f(\mathbf{x}_i)\}$$

Stopping rule: Terminate if the gradient is small ($|\nabla f(\mathbf{x}_i)| \approx \mathbf{0}$)
or successive points are close to each other ($|\mathbf{x}_{i+1} - \mathbf{x}_i| \approx \mathbf{0}$)

Gradient Minimisation: The Steepest Descent

Initialisation: Select an initial point \mathbf{x}_0

Iterative steps $i = 0, 1, \dots$:

- 1 Compute the gradient $\nabla f(\mathbf{x})$ at the point \mathbf{x}_i
- 2 Draw a line $\mathbf{x}_i - t\nabla f(\mathbf{x}_i)$ in the opposite gradient direction
- 3 Select the point \mathbf{x}_{i+1} on this line yielding the smallest value for f of all points on the line:

$$f(\mathbf{x}_{i+1}) = \min_{t \in (0, \infty)} \{f(\mathbf{x}) : \mathbf{x} = \mathbf{x}_i - t\nabla f(\mathbf{x}_i)\}$$

Stopping rule: Terminate if the gradient is small ($|\nabla f(\mathbf{x}_i)| \approx \mathbf{0}$) or successive points are close to each other ($|\mathbf{x}_{i+1} - \mathbf{x}_i| \approx \mathbf{0}$)

Finding the Best Point Along the Gradient Direction

Search for the extremum point along the line $f(\mathbf{x}) \pm t \cdot \nabla f(\mathbf{x})$

If the derivatives ∇f are computable and the function f is well-behaving **then**:

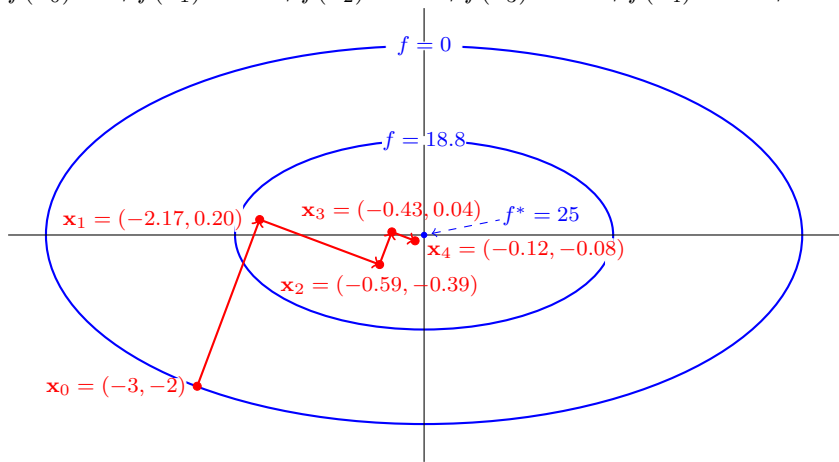
- 1 Substitute $\mathbf{x}_0 \pm t \cdot \nabla f(\mathbf{x}_0)$ into the equation for f ,
- 2 Differentiate with respect to t , and
- 3 Set the derivative equal to zero to find t

else if the function $u(t) \equiv f(\mathbf{x}_0 \pm t \nabla f(\mathbf{x}_0))$ is unimodal **then** use any one-dimensional **line search**, e.g.

- the golden section search or
- the Fibonacci section search

Steepest Ascent: An Example

$f(\mathbf{x}) \equiv f(x, y) = 25 - x^2 - 4y^2$; $\nabla f(\mathbf{x}) = (-2x, -8y)$; $\mathbf{x}_0 = (-3, -2)$:
 $f(\mathbf{x}_0) = 0$; $f(\mathbf{x}_1) = 20.1$; $f(\mathbf{x}_2) = 24.0$; $f(\mathbf{x}_3) = 24.8$; $f(\mathbf{x}_4) = 24.9$; ...



Steepest Descent: $f(x, y) = 10x^2 - 2xy + 2y^2 - 18x - 2y$

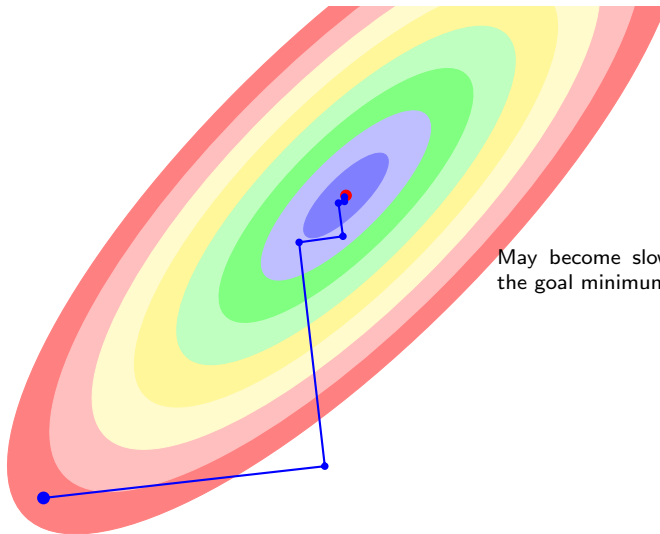
$$\begin{cases} \left. \frac{\partial f}{\partial x} \right|_{x_i, y_i} \equiv u_i = 20x - 2y - 18 \\ \left. \frac{\partial f}{\partial y} \right|_{x_i, y_i} \equiv v_i = -2x + 4y - 2 \end{cases} \Rightarrow \begin{cases} x_{i+1} = x_i - t_i u_i \\ y_{i+1} = y_i - t_i v_i \end{cases}$$

$$f(x_i + tu_i, y_i + tv_i) = f(x_i, y_i) - t(u_i^2 + v_i^2) + t^2 (au_i^2 + bu_i v_i + cv_i^2)$$

$$\Rightarrow t_i = \arg \min_t \{f(x_i + tu_i, y_i + tv_i)\} = \frac{1}{2} \frac{u_i^2 + v_i^2}{au_i^2 + bu_i v_i + cv_i^2}$$

i	x_i	y_i	$f(x_i, y_i)$	u_i	v_i	x_{i+1}	y_{i+1}
0	-1.0000	-1.0000	30.0000	-36.0000	-4.0000	0.8589	-0.7935
1	0.8589	-0.7935	-3.8741	0.7657	-6.8917	0.6937	0.6937
2	0.6937	0.6937	-9.0618	-5.5134	-0.6126	0.9784	0.7253
3	0.9784	0.7253	-9.8563	0.1173	-1.0555	0.9531	0.9531
4	0.9531	0.9531	-9.9780	-0.8444	-0.0938	0.9967	0.9579
5	0.9967	0.9579	-9.9966	0.0180	-0.1616	0.9928	0.9928
6	0.9928	0.9928	-9.9995	-0.1293	-0.0144	0.9995	0.9936
7	0.9995	0.9936	-9.9999	0.0028	-0.0248	0.9989	0.9989
8	0.9989	0.9989	-10.0000	-0.0198	-0.0022	0.9999	0.9990
9	0.9999	0.9990	-10.0000	0.0004	-0.0038	0.9998	0.9998

Steepest Descent: $f(x, y) = 10x^2 - 2xy + 2y^2 - 18x - 2y$

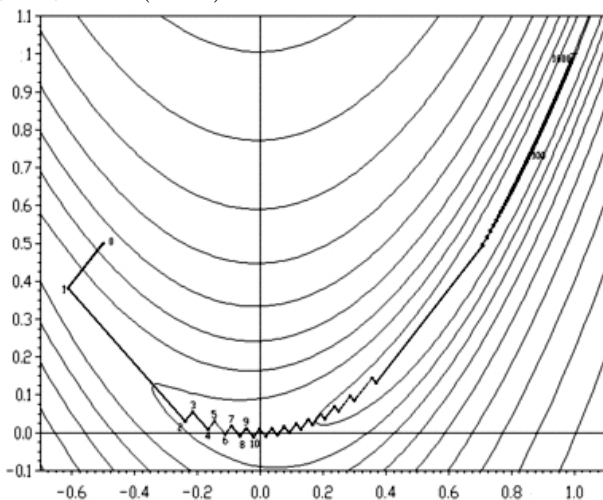


May become slow near the goal minimum...

Steepest Descent: Rosenbrock's "Banana"

$$f(x, y) = (1 - x)^2 + 100(y - x^2)^2$$

http://en.wikipedia.org/wiki/Gradient_descent



(c) 2006 P.A. Simianescu

Unconstrained Optimisation: An Informal Sketch

Minimum or maximum of a function $f(\mathbf{x})$ of many variables

- The goal: $\mathbf{x}^* = \arg \{ \underset{\min}{\max} \} f(\mathbf{x})$; $f^* \equiv f(\mathbf{x}^*) = \{ \underset{\min}{\max} \} f(\mathbf{x})$
- Mostly: a multi-modal $f(\mathbf{x})$
 - Modes correspond to the gradient's roots $\nabla f(\mathbf{x}^\circ) = \mathbf{0}$
 - Most optimistic case: an analytical solution of this system
 - The optimum may not be unique: $f^* = f(\mathbf{x}_1^*) = \dots = f(\mathbf{x}_k^*)$

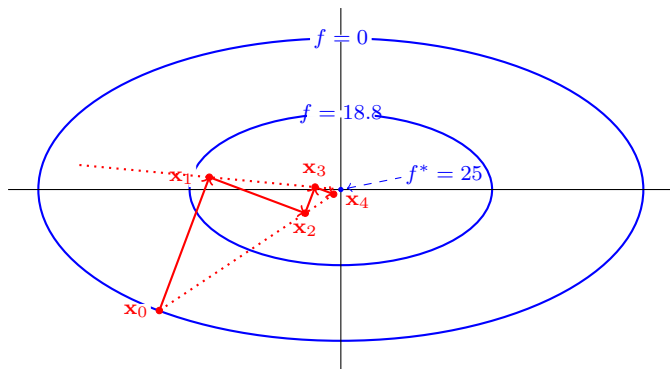
Global optimisation is typically a hard problem

- Feasible solutions - only for specific functions

Local optimisation: find the optimiser \mathbf{x}° close to a known \mathbf{x}_0

- Most optimistic case: $f(\mathbf{x})$; $\nabla f(\mathbf{x})$ – an analytical solution
- Typically – an iterative linear search starting from \mathbf{x}_0
 - Analytical or numerical selection of direction \mathbf{s}_i at each step i
 - Analytical or numerical step, e.g. $t^\circ = \arg \{ \underset{\min}{\max} \}_t f(\mathbf{x}_i + t\mathbf{s}_i)$
 - Accelerated linear search with interdependent \mathbf{s}_i and \mathbf{s}_j ; $j < i$

Accelerated Gradient Search



- Once $i > 2$, for i odd the point \mathbf{x}_i is found by gradient search from \mathbf{x}_{i-1} , and \mathbf{x}_{i+1} is found by an **accelerated step** by maximising over the line through \mathbf{x}_i and \mathbf{x}_{i-2}
- Global maximum of a negative definite quadratic function of n variables is provably found after $2n - 1$ steps of this procedure

Conjugate Directions (optional)

- Produce a sequence of points $\mathbf{x}_0, \mathbf{x}_1, \dots$, such that each point improves values in maximising a quadratic function $f(\mathbf{x}) = \mathbf{a}^T + \frac{1}{2}\mathbf{x}^T\mathbf{H}\mathbf{x}$
- All directions \mathbf{d}_i of search obey the relationship: $\mathbf{d}_i^T\mathbf{H}\mathbf{d}_j = 0$ for all $i, j, i \neq j$

General method of conjugate directions

- Choose \mathbf{x}_0 near an optimal point or randomly
- Carry out an one-dimensional search in the first conjugate direction \mathbf{d}_1 to find a new point \mathbf{x}_1
- For $i = 2, \dots, n$, search for a new point \mathbf{x}_i along the next conjugate direction \mathbf{d}_i such that $\mathbf{d}_j^T\mathbf{H}\mathbf{d}_k = 0; j, k \leq i, j < k$
- The maximum is located in at most n steps

Conjugate Gradients (optional)

Each new conjugate direction – from the gradient at the point concerned

Conjugate gradient method for maximising $f(\mathbf{x})$

- Choose a starting point \mathbf{x}_0
- Carry out an one-dimensional search in the gradient direction $\mathbf{d}_1 = \nabla f(\mathbf{x}_0)$ to find the maximum point \mathbf{x}_1
- For $i = 2, \dots, n$, form \mathbf{d}_i from $\nabla f(\mathbf{x}_i)$ to be conjugate to \mathbf{d}_{i-1} :
 $\mathbf{d}_i = \nabla f(\mathbf{x}_i) + \gamma_{i-1} \mathbf{d}_{i-1}$ and $\mathbf{d}_i^\top \mathbf{H} \mathbf{d}_{i-1} = 0$

$$\Rightarrow \mathbf{d}_i = \nabla f(\mathbf{x}_i) - \left(\frac{(\nabla f(\mathbf{x}_i))^\top \mathbf{H} \mathbf{d}_{i-1}}{\mathbf{d}_{i-1}^\top \mathbf{H} \mathbf{d}_{i-1}} \right) \mathbf{d}_{i-1}$$

- Can be proven by induction: all \mathbf{d}_i are mutually conjugate
- In actual implementation the directions \mathbf{d}_i can be computed by a simple recurrence relation, and only a few vectors and no matrices need be stored

Direct Search Methods

- If both the gradient and Hessian of $f(\mathbf{x})$ are too complicated to compute but f can be evaluated at any point $\mathbf{x} \in \mathbb{R}_n$

Pattern search of K. Hooke and T. A. Jeeves

- For $i = 1, \dots, n$ sequentially:
 - If $f(x_1, \dots, x_i + \varepsilon_i, \dots, x_n) > f(x_1, \dots, x_i, \dots, x_n)$, replace $x_i \leftarrow x_i + \varepsilon$
 - Else if $f(x_1, \dots, x_i - \varepsilon_i, \dots, x_n) > f(x_1, \dots, x_i, \dots, x_n)$, replace $x_i \leftarrow x_i - \varepsilon$
- Repeat this cycle of perturbations until no perturbations about \mathbf{x}_j bring about an improvement
- Halve the pre-defined perturbation sizes ε_i and repeat the process while the next point brings an improvement over \mathbf{x}_j

Sectioning and Accelerated Rosenbrock's Search

One-at-a-time search, or sectioning, from an initial point \mathbf{x}_0

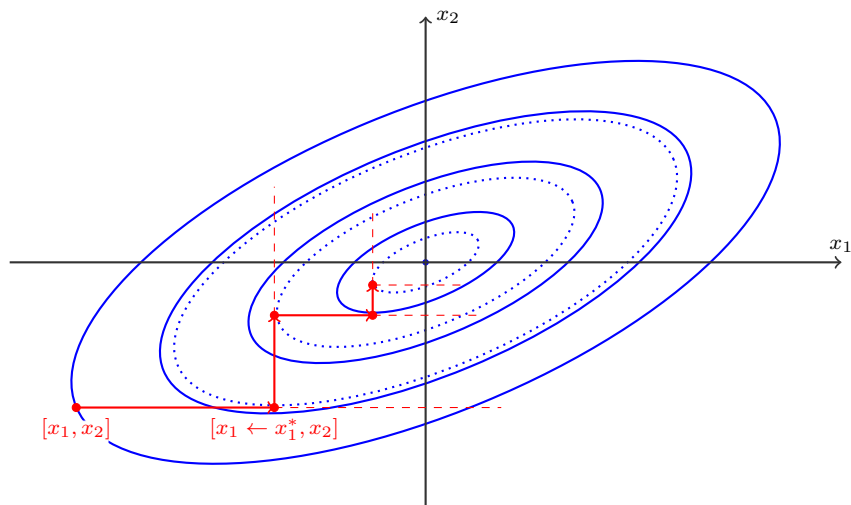
- For $i = 1, \dots, n$ sequentially, search for the maximum in the direction of the variable x_i by one of the one-dimensional search methods and replace x_i by the maximiser x_i^* : $x_i \leftarrow x_i^*$
- Repeat this cycle of one-dimensional searches until the steps $x_i - x_i^*$; $i = 1, \dots, n$ become less than a given threshold

Convergence rate is usually too slow and the search may halt far from the optimum

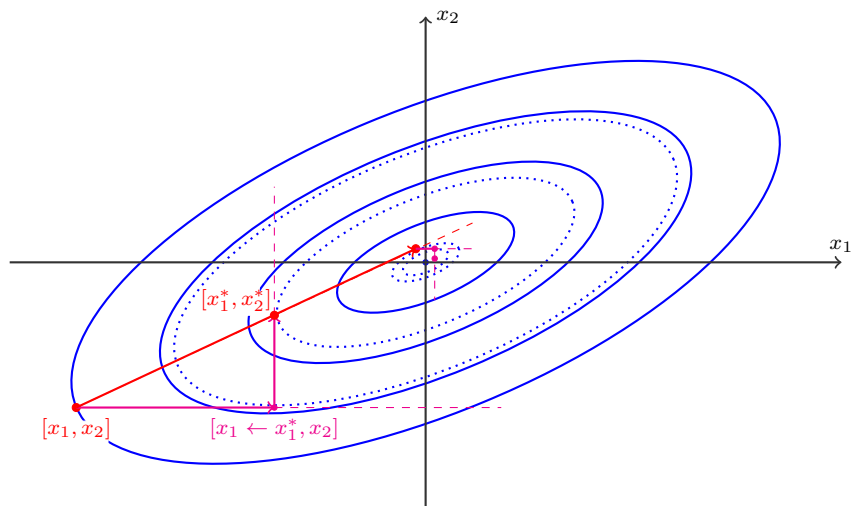
Accelerated search of H. H. Rosenbrock

- Use one-at-a-time search from \mathbf{x}_0 to find the next point \mathbf{x}_1^* and the direction $\boldsymbol{\delta}$ with components $\delta_i = x_{1:i}^* - x_{0:i}$
- Search for the maximum in the direction $\boldsymbol{\delta}$ and replace \mathbf{x}_0 by the maximiser \mathbf{x}_1 found
- Repeat this cycle until \mathbf{x}_t and \mathbf{x}_{t-1} are closer than a threshold

One-at-a-time Search: An Example



Rosenbrock's Search: An Example



Search Method of M. J. D. Powell (optional)

- Similar to the method of conjugate gradients, except that derivatives are not required
- Similar to the Rosenbrock's method, except that each search is carried out along a conjugate direction
 - Directions $\mathbf{d}_1, \dots, \mathbf{d}_n$ become conjugate w.r.t. an approximation of the Hessian matrix

If \mathbf{x}_0 is the initial estimate of the maximiser of $f(\mathbf{x})$ then

- 1 Set the search directions be equal to the coordinate directions
- 2 For $i = 1, \dots, n$ sequentially find the maximiser \mathbf{x}_i of f in the the direction \mathbf{d}_i from \mathbf{x}_{i-1}
- 3 Let $\mathbf{d}_i \leftarrow \mathbf{d}_{i+1}$ for $i = 1, \dots, n - 1$ and $\mathbf{d}_n = \mathbf{x}_n - \mathbf{x}_0$
- 4 Set \mathbf{x}_0 be equal to the maximiser of f in the \mathbf{d}_n direction from \mathbf{x}_n
- 5 Return to 2 unless some termination criterion is met