

# Advanced computational statistics, lecture 1

Frank Miller, Department of Computer and Information Science,  
Linköping University

Department of Statistics; Stockholm University

March 16, 2023

# Course schedule

- Topic 1: **Gradient based optimisation**
  - Topic 2: **Stochastic gradient based optimisation**
  - Topic 3: **Gradient free optimisation**
  - Topic 4: **Optimisation with constraints**
  - Topic 5: **EM algorithm and bootstrap**
  - Topic 6: **Simulation of random variables**
  - Topic 7: **Importance sampling**
- Optimisation
- Simulation and Integration

Course homepage:

<http://www.adoptdesign.de/frankmillereu/adcompstat2023.html>

Includes schedule, reading material, lecture notes, assignments

# Optimisation in statistics

- Maximum Likelihood
  - Minimising risk in (Bayesian) decision theory
  - Minimising sum of squares (Least Squares Estimate)
  - Maximising information in experimental design
  - Machine learning
- Common problem in these examples:
    - $\mathbf{x}$   $p$ -dimensional vector,  $g: \mathbb{R}^p \rightarrow \mathbb{R}$  function
    - We search  $\mathbf{x}^*$  with  $g(\mathbf{x}^*) = \max g(\mathbf{x})$
- Typical:  $g = \sum_{i=1}^n g_i$  with a (large) sample size  $n$  where  $g_i: \mathbb{R}^p \rightarrow \mathbb{R}$
  - Minimisation problem turns into maximisation by considering  $-g$

# Least squares estimation (LSE)

- We search a Least Squares estimate  $\hat{\boldsymbol{\beta}}$  for  $\boldsymbol{\beta}$  minimising the distance  $g(\hat{\boldsymbol{\beta}}) = \|\hat{\mathbf{y}} - \mathbf{y}\|^2$  from  $\hat{\mathbf{y}} = \mathbf{X}\hat{\boldsymbol{\beta}}$  to  $\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}$
- $g(\hat{\boldsymbol{\beta}}) = \|\mathbf{X}\hat{\boldsymbol{\beta}} - \mathbf{y}\|^2 = (\mathbf{X}\hat{\boldsymbol{\beta}} - \mathbf{y})^T (\mathbf{X}\hat{\boldsymbol{\beta}} - \mathbf{y}) = \hat{\boldsymbol{\beta}}^T \mathbf{X}^T \mathbf{X} \hat{\boldsymbol{\beta}} - 2\hat{\boldsymbol{\beta}}^T \mathbf{X}^T \mathbf{y} + \mathbf{y}^T \mathbf{y}$
- Setting the derivative to 0 ( $\frac{\partial f}{\partial \hat{\boldsymbol{\beta}}} = 2\mathbf{X}^T \mathbf{X} \hat{\boldsymbol{\beta}} - 2\mathbf{X}^T \mathbf{y} = 0$ ), we get  $\hat{\boldsymbol{\beta}} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$
- Note that  $g(\hat{\boldsymbol{\beta}}) = \|\mathbf{X}\hat{\boldsymbol{\beta}} - \mathbf{y}\|^2 = \sum_{i=1}^n (\mathbf{x}_i^T \hat{\boldsymbol{\beta}} - y_i)^2 = \sum_{i=1}^n g_i(\hat{\boldsymbol{\beta}})$
- Optimisation problem:
  - $\hat{\boldsymbol{\beta}}$   $p$ -dimensional vector,  $g: \mathbb{R}^p \rightarrow \mathbb{R}$  function
  - We search  $\hat{\boldsymbol{\beta}}$  with  $g(\hat{\boldsymbol{\beta}}) = \min g(\mathbf{b}) = \min \sum_{i=1}^n g_i(\mathbf{b})$
- Here, we do not need to iteratively compute this minimum since we have an algebraic solution  $\hat{\boldsymbol{\beta}} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$

# Variations of least squares estimation

- Algebraic solution exists for the LSE, but not if we vary the problem
- Lasso estimate:  $g(\hat{\boldsymbol{\beta}}) = \|\mathbf{X}\hat{\boldsymbol{\beta}} - \mathbf{y}\|^2 + \lambda\|\hat{\boldsymbol{\beta}}\|_1 = \sum_{i=1}^n (\mathbf{x}_i\hat{\boldsymbol{\beta}} - y_i)^2 + \lambda\|\hat{\boldsymbol{\beta}}\|_1 = \sum_{i=1}^n g_i(\hat{\boldsymbol{\beta}})$
- $L_1$ -estimation:  $g(\hat{\boldsymbol{\beta}}) = \|\mathbf{X}\hat{\boldsymbol{\beta}} - \mathbf{y}\|_1 = \sum_{i=1}^n |\mathbf{x}_i\hat{\boldsymbol{\beta}} - y_i| = \sum_{i=1}^n g_i(\hat{\boldsymbol{\beta}})$
- Many further variations of estimates have been considered
- In all cases, we search  $\hat{\boldsymbol{\beta}}$  with  $g(\hat{\boldsymbol{\beta}}) = \min g(\mathbf{b}) = \min \sum_{i=1}^n g_i(\mathbf{b})$
- Recall: Norms for  $\mathbf{x} = (x_1, \dots, x_p)^T$ :  $\|\mathbf{x}\| = \|\mathbf{x}\|_2 = \sqrt{x_1^2 + \dots + x_p^2}$  (Euklid),  $\|\mathbf{x}\|_1 = |x_1| + \dots + |x_p|$ ,  $\|\mathbf{x}\|_\infty = \max\{|x_1|, \dots, |x_p|\}$  (max-norm)

# Maximizing information of experimental designs

- Regression model  $\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}$  (where  $\boldsymbol{\varepsilon}$  has iid components)
- $\mathbf{X}$  design matrix (depends on choice of observational points)
- Covariance matrix of Least Squares estimate  $\hat{\boldsymbol{\beta}}$  is
$$\text{Cov}(\hat{\boldsymbol{\beta}}) = (\mathbf{X}^T \mathbf{X})^{-1} \cdot \text{const}$$
- Choose design of an experiment such that  $\mathbf{X}^T \mathbf{X}$  “large”
- D-optimality:  $g(\text{"design"}) = \det(\mathbf{X}^T \mathbf{X})$
- We search  $\text{design}^*$  with  $g(\text{design}^*) = \max g(\text{design})$

# Maximizing information of experimental designs

- Regression model  $\mathbf{y} = \mathbf{X} \boldsymbol{\beta} + \boldsymbol{\varepsilon}$ ,  $\text{Cov}(\hat{\boldsymbol{\beta}}) = (\mathbf{X}^T \mathbf{X})^{-1} \cdot \text{const}$
- We search  $\text{design}^*$  with  $g(\text{design}^*) = \max g(\text{design})$
- Example: cubic regression,  $y = \beta_0 + \beta_1 x + \beta_2 x^2 + \beta_3 x^3 + \varepsilon$ ,  $n$  observations in each of following 4 points:  $-1, -a, a, 1$ . How should  $a \in (0,1)$  be chosen?

$$\mathbf{X} = \begin{pmatrix} 1 & -1 & 1 & -1 \\ 1 & -a & a^2 & -a^3 \\ 1 & a & a^2 & a^3 \\ 1 & 1 & 1 & 1 \end{pmatrix}$$

$$g(a) = \det(\mathbf{X}^T \mathbf{X}) = \det(\mathbf{X}(a)^T \mathbf{X}(a))$$

- We search  $a^*$  with  $g(a^*) = \max g(a)$

# Today's schedule

- Univariate Optimisation (bi-section, Newton, secant)
- Multivariate Optimisation
  - Analytical opt.
  - Newton
  - Steepest ascent
  - Accelerated steepest ascent
  - Quasi-Newton

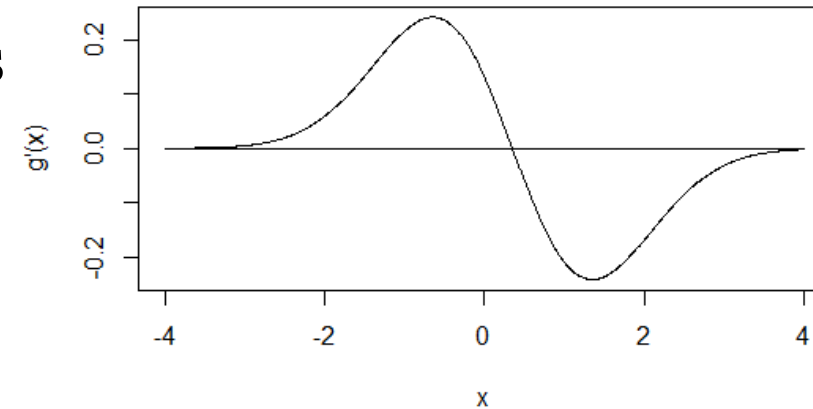
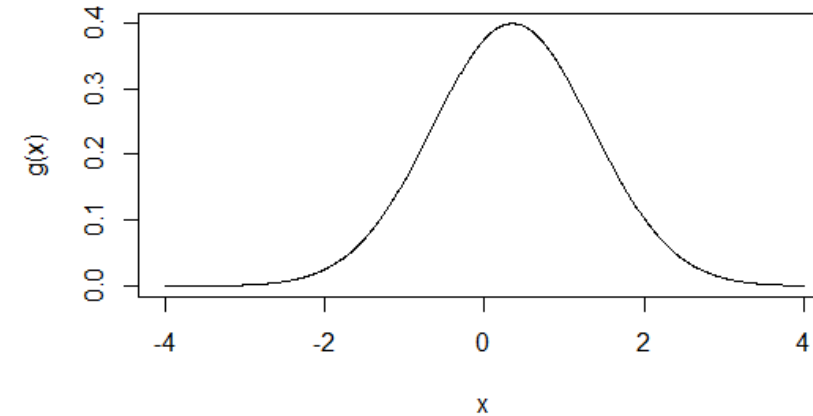


# Univariate optimisation

- $x$  real number,  $g: \mathbb{R} \rightarrow \mathbb{R}$  continuously differentiable function
- We search  $x^*$  with  $g(x^*) = \max g(x)$
- Compute  $g'(x)$  and search  $x^*$  with  $g'(x^*) = 0$
- One has then to check if the result is maximum, minimum, possibly local optimum...

# Univariate optimisation: bisection

- Search  $x^*$  with  $g'(x^*) = 0$ :
- See [video on course homepage](#)
- We iteratively improve approximations for  $x^*$ :  
 $x^{(0)} \rightarrow x^{(1)} \rightarrow x^{(2)} \rightarrow \dots$



# Optimisation: convergence criterion

- Compare  $x^{(t)}$  and  $x^{(t+1)}$  and stop if they are “close enough”
- Absolute convergence criterion:

$$|x^{(t+1)} - x^{(t)}| < \epsilon$$

- Relative convergence criterion:

$$\frac{|x^{(t+1)} - x^{(t)}|}{|x^{(t)}|} < \epsilon$$

# Univariate Newton(-Raphson)

- $x$  real number,  $g: \mathbb{R} \rightarrow \mathbb{R}$  twice differentiable function
- Search  $x^*$  with  $g(x^*) = \max g(x)$  by searching  $x^*$  with  $g'(x^*) = 0$

- Taylor expansion around  $x^*$  motivates:

$$0 = g'(x^*) \approx g'(x^{(t)}) + (x^* - x^{(t)})g''(x^{(t)})$$

$$-(x^* - x^{(t)})g''(x^{(t)}) \approx g'(x^{(t)})$$

$$x^* \approx x^{(t)} - g'(x^{(t)})/g''(x^{(t)})$$

- Therefore, the Newton-iteration works as:

$$x^{(t+1)} = x^{(t)} - g'(x^{(t)})/g''(x^{(t)})$$

# Univariate Newton(-Raphson)

- $x^{(t+1)} = x^{(t)} - g'(x^{(t)})/g''(x^{(t)})$
- Start with a  $x^{(0)}$
- Tangent in  $(x^{(0)}, g'(x^{(0)}))$  determines  $x^{(1)}$
- Tangent in  $(x^{(1)}, g'(x^{(1)}))$  determines  $x^{(2)}$
- ...
- until convergence criterion met

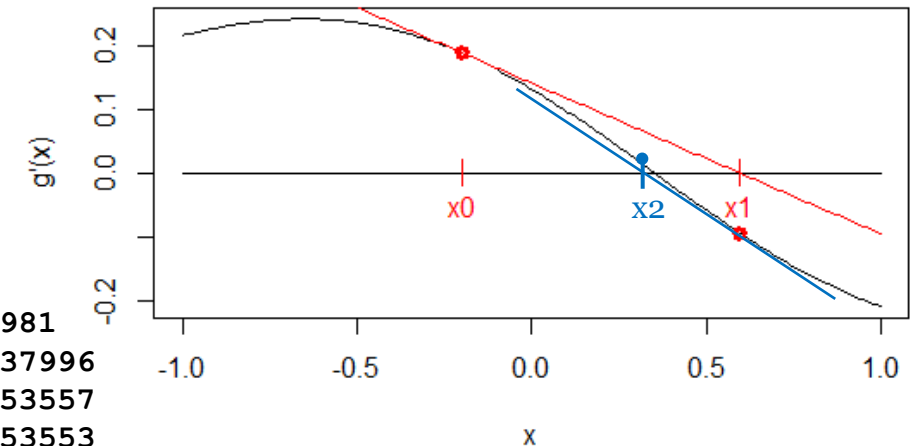
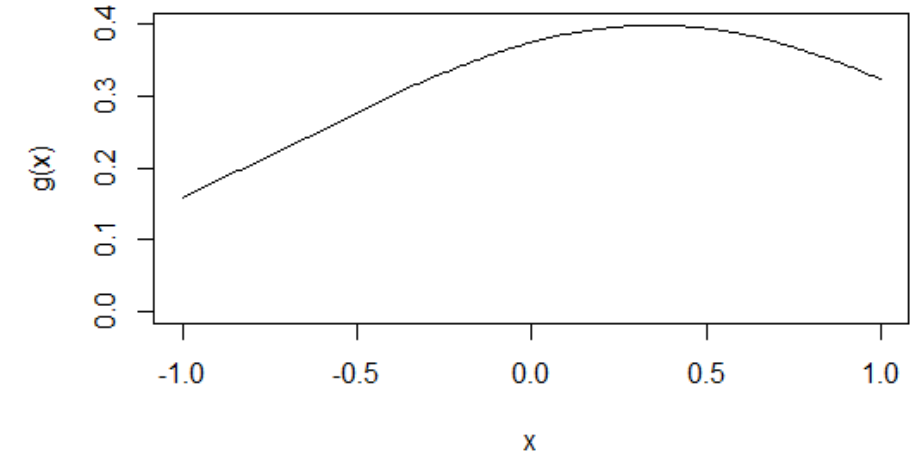
+Newton method is fast

- Requires existence and computation of  $g''$

```

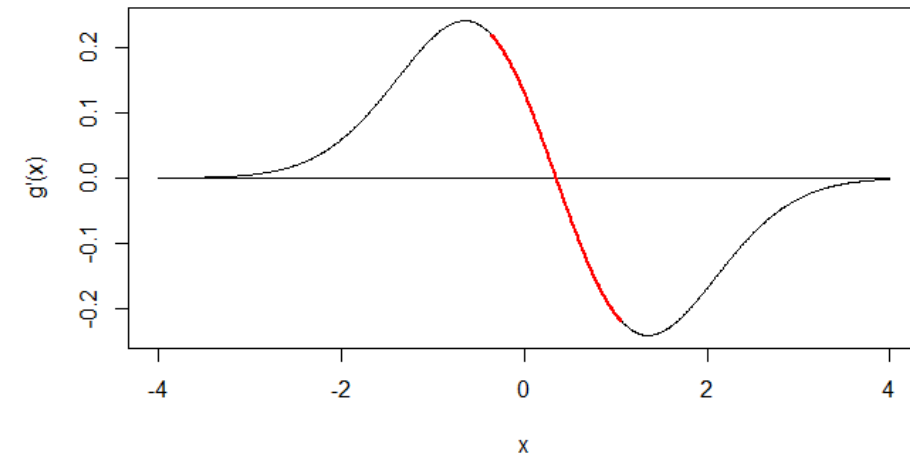
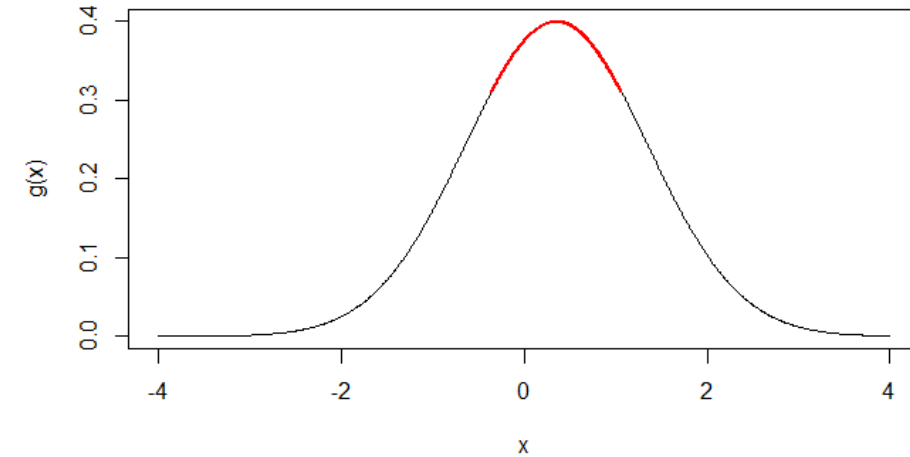
x0 -0.2
x1 0.5981
x2 0.337996
x3 0.353557
x4 0.353553
x5 0.353553
STOP

```



# Univariate Newton(-Raphson)

- $x^{(t+1)} = x^{(t)} - g'(x^{(t)})/g''(x^{(t)})$
- What about the starting value  $x^{(0)}$ ?



# Univariate secant method

- $x$  real number,  $g: \mathbb{R} \rightarrow \mathbb{R}$  once differentiable function
- Search  $x^*$  with  $g(x^*) = \max g(x)$  by searching  $x^*$  with  $g'(x^*) = 0$
- Recall: The Newton-iteration works as:
$$x^{(t+1)} = x^{(t)} - g'(x^{(t)})/g''(x^{(t)})$$
- Need to compute  $g''$  which might be difficult. Instead:
- Approximate  $g''(x^{(t)})$  by  $[g'(x^{(t)}) - g'(x^{(t-1)})]/(x^{(t)} - x^{(t-1)})$

# Univariate secant method

- $x^{(t+1)} = x^{(t)} - g'(x^{(t)}) \frac{x^{(t)} - x^{(t-1)}}{g'(x^{(t)}) - g'(x^{(t-1)})}$
- Start with  $x^{(0)}$  and  $x^{(-1)}$
- Secant through  $x^{(0)}$  and  $x^{(-1)}$  determines  $x^{(1)}$
- Secant through  $x^{(1)}$  and  $x^{(0)}$  determines  $x^{(2)}$
- ...
- until stopping crit. fulfilled

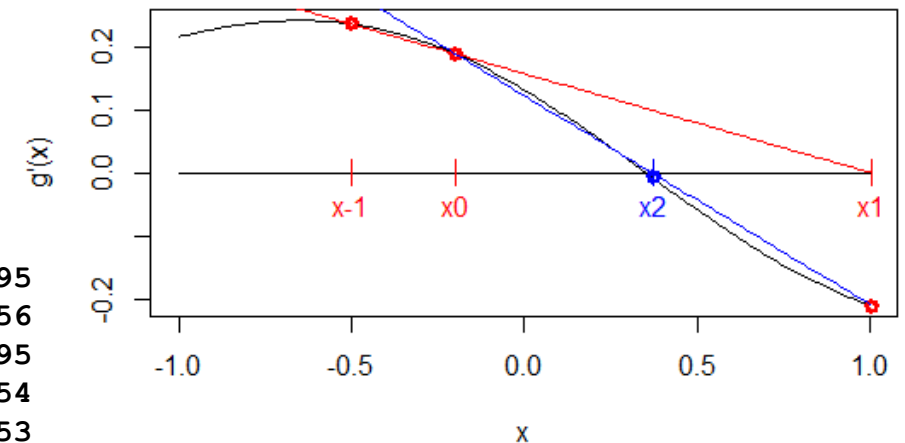
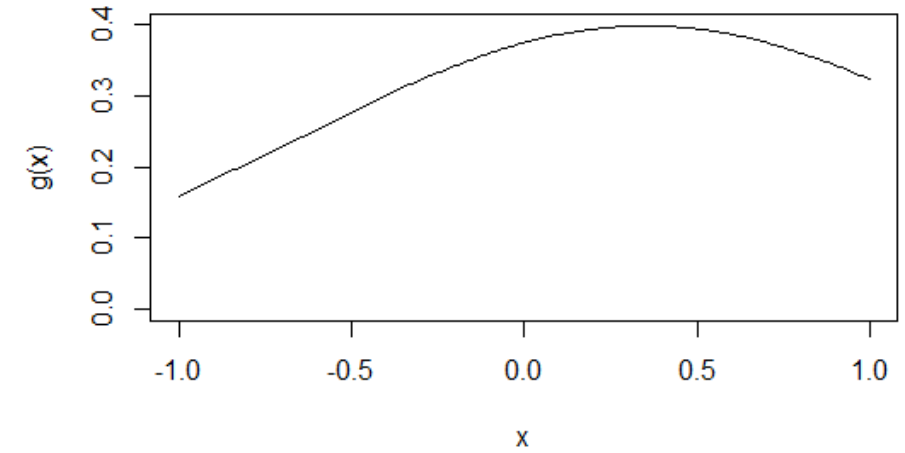
- Quite fast

- No 2<sup>nd</sup> derivative necessary

```

x0 -0.2
x1 1.006995
x2 0.371656
x3 0.349095
x4 0.353554
x5 0.353553
x6 0.353553
STOP

```





# Convergence speed of optimisation algorithms

- Convergence speed can be quantified by  $q$  and  $c$  as follows:
  - Let  $\varepsilon^{(t)} = x^{(t)} - x^*$ ,
  - Find  $q$  and  $c$  such that  $\lim_{t \rightarrow \infty} \varepsilon^{(t+1)} / (\varepsilon^{(t)})^q = c$
- $\varepsilon = 1, 0.5, 0.25, 0.125, 0.063, 0.031, \dots \rightarrow q=1, c=0.5,$
- $\varepsilon = 1, 0.1, 0.01, 0.001, 0.0001, \dots \rightarrow q=1, c=0.1,$
- If  $q=1$ , we say that convergence is "linear"
- $\varepsilon = 1, 0.5, 0.125, 0.008, 0.00003, \dots \rightarrow q=2, c=0.5.$
- If  $q=2$ , we say that convergence is "quadratic"

Convergence  
order

Convergence  
rate

Intuitively,  
 $\varepsilon^{(t+1)} \approx c \cdot (\varepsilon^{(t)})^q$

# Determine empirically convergence rate (and order) of optimisation algorithms

- You have a given optimisation algorithm and you have determined or know the maximiser  $\mathbf{x}^*$ . To check convergence speed in an optimisation-run, you can calculate

$$D^{(t)} = \frac{|x^{(t)} - x^*|}{|x^{(t-1)} - x^*|}$$

(see Givens and Hoeting, 2013, page 101/102, for an example)

- If  $D^{(t)} \rightarrow 1$ , there is not even linear convergence (bad, order  $q < 1$ ),  
If  $D^{(t)} \rightarrow c \in (0,1)$ , linear convergence (order  $q=1$ ) with rate  $c$ ,  
If  $D^{(t)} \rightarrow 0$ , better than linear convergence (order  $q > 1$ ).

# Comparison of univariate optimisation methods

Bisection	Secant	Newton
$g'$ required	$g'$ required	$g''$ required
finds always an optimum between $a_0$ and $b_0$ (but could be local)	converges only when the two starting values "close" to optimum	converges only when starting value "close" to optimum
slow $q=1$	$q = \frac{1 + \sqrt{5}}{2} = 1.62$	fast $q=2$

# Today's schedule

- Univariate Optimisation (bi-section, Newton, secant)
- **Multivariate Optimisation**
  - Analytical opt.
  - Newton
  - Steepest ascent
  - Accelerated steepest ascent
  - Quasi-Newton

# Multivariate optimisation – gradient and Hessian

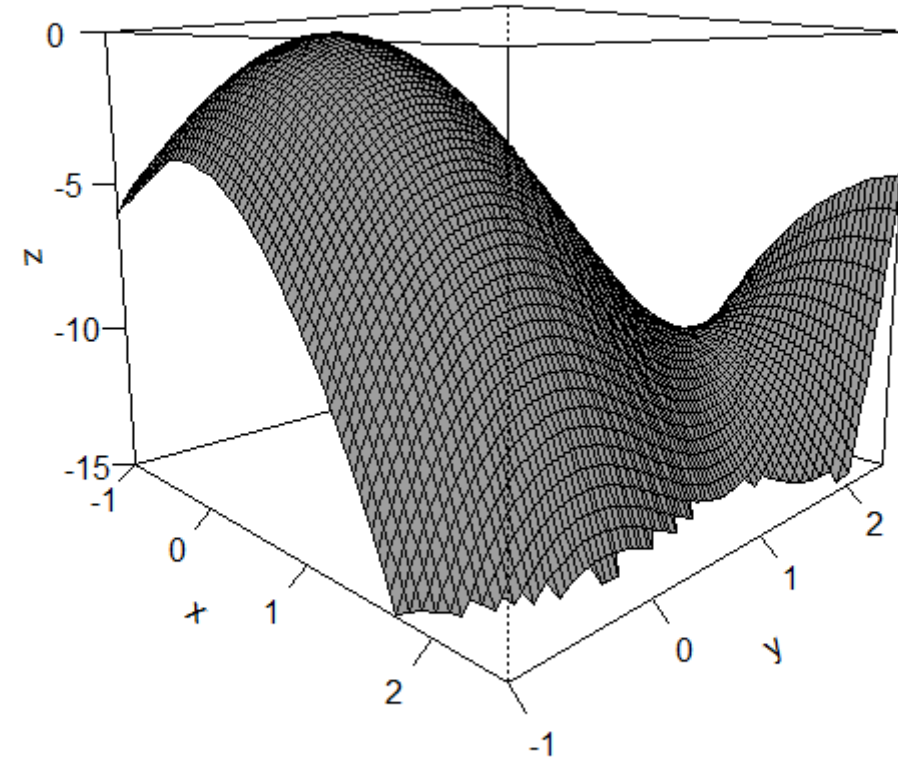
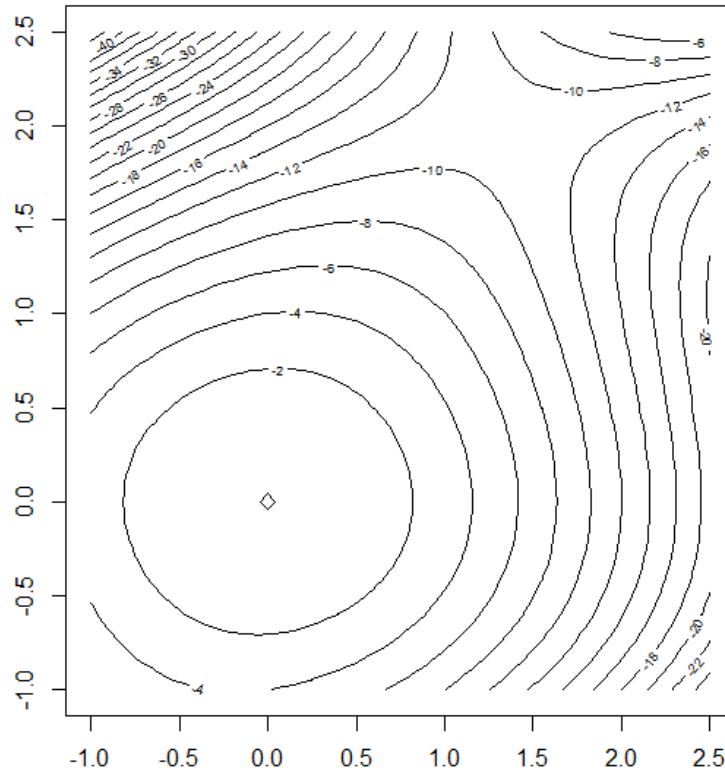
- $g \begin{pmatrix} x_1 \\ \vdots \\ x_p \end{pmatrix}$  is a real-valued function

- $g' \begin{pmatrix} x_1 \\ \vdots \\ x_p \end{pmatrix} = \begin{pmatrix} \frac{\partial g}{\partial x_1}(\mathbf{x}) \\ \vdots \\ \frac{\partial g}{\partial x_p}(\mathbf{x}) \end{pmatrix}$  is the gradient,  $\mathbf{x} = \begin{pmatrix} x_1 \\ \vdots \\ x_p \end{pmatrix}$

- $g'' \begin{pmatrix} x_1 \\ \vdots \\ x_p \end{pmatrix} = \begin{pmatrix} \frac{\partial^2 g}{\partial x_1 \partial x_1}(\mathbf{x}) & \dots & \frac{\partial^2 g}{\partial x_1 \partial x_p}(\mathbf{x}) \\ \vdots & & \vdots \\ \frac{\partial^2 g}{\partial x_1 \partial x_p}(\mathbf{x}) & \dots & \frac{\partial^2 g}{\partial x_p \partial x_p}(\mathbf{x}) \end{pmatrix}$  is the Hessian matrix

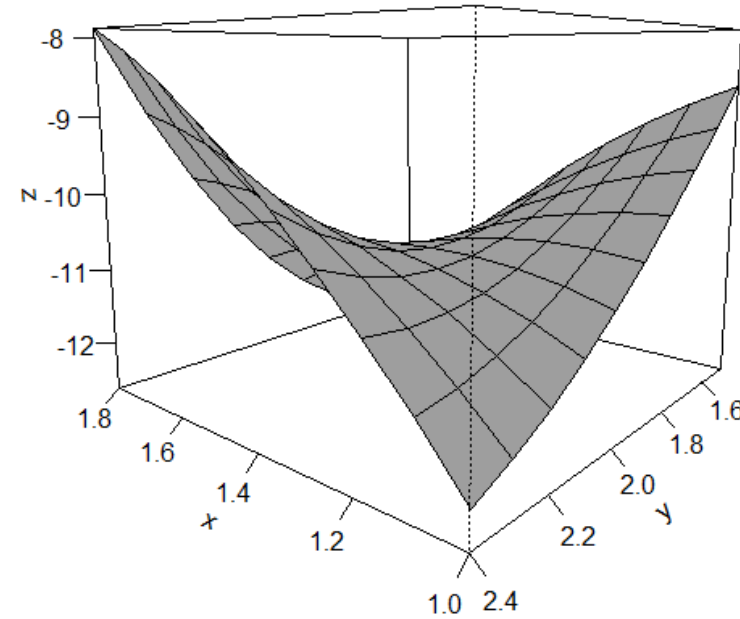
# Bivariate optimisation - visualisation

- $g \begin{pmatrix} x \\ y \end{pmatrix} = -3x^2 - 4y^2 + xy^3$



Figures can be drawn using R-core-functions `contour` and `persp`

# Multivariate optimisation - saddle points



# Multivariate optimisation - analytical optimisation

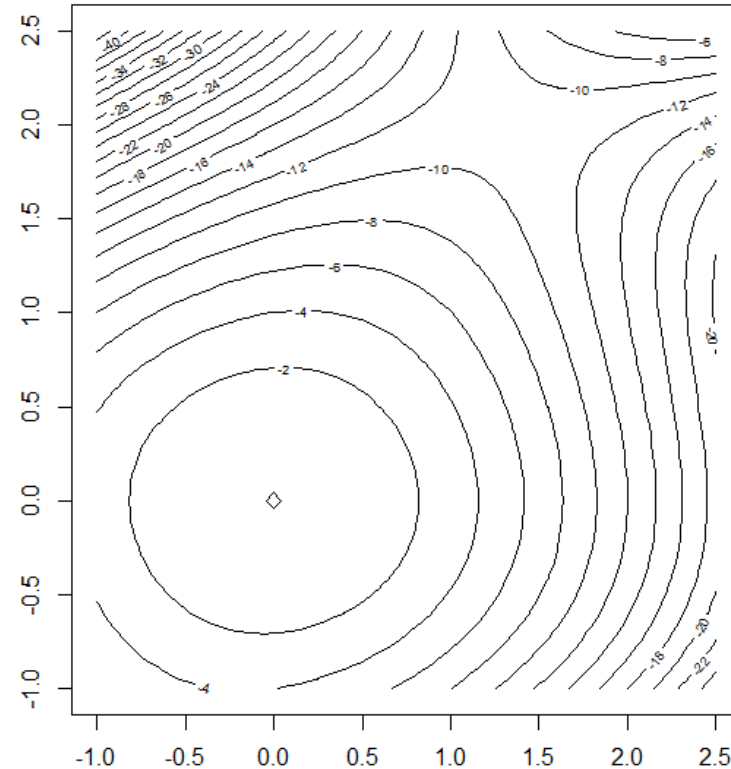
- $g \begin{pmatrix} x \\ y \end{pmatrix} = -3x^2 - 4y^2 + xy^3$

- $g' \begin{pmatrix} x \\ y \end{pmatrix} = \begin{pmatrix} -6x + y^3 \\ -8y + 3xy^2 \end{pmatrix}$

- $g'' \begin{pmatrix} x \\ y \end{pmatrix} = \begin{pmatrix} -6 & 3y^2 \\ 3y^2 & -8 + 6xy \end{pmatrix}$

- See calculation in following document:  
`AdvCompStat_AnalytOpt.pdf`

- ~~Maximum at (0,0), saddle point at (4/3,2)~~



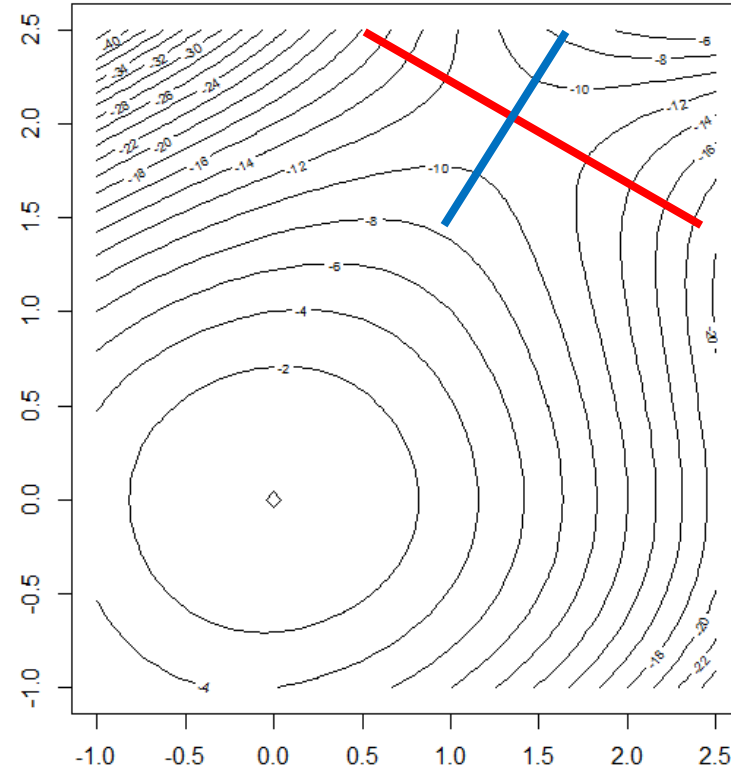


# Saddle point and eigenvectors of the Hessian

- $\mathbf{g} \begin{pmatrix} x \\ y \end{pmatrix} = -3x^2 - 4y^2 + xy^3$
- Saddle point at  $(4/3, 2)$

- $\mathbf{g}' \begin{pmatrix} 4/3 \\ 2 \end{pmatrix} = \begin{pmatrix} 0 \\ 0 \end{pmatrix}$

- $\mathbf{g}'' \begin{pmatrix} 4/3 \\ 2 \end{pmatrix} = \begin{pmatrix} -6 & 12 \\ 12 & 8 \end{pmatrix}$



- Eigenvalues 14.89, -12.89; eigenvectors  $\begin{pmatrix} 0.498 \\ 0.867 \end{pmatrix}$ ,  $\begin{pmatrix} -0.867 \\ 0.498 \end{pmatrix}$

# Multivariate Newton

- $\mathbf{x}$   $p$ -dimensional vector,  $g: \mathbb{R}^p \rightarrow \mathbb{R}$  function
- We search  $\mathbf{x}^*$  with  $g(\mathbf{x}^*) = \max g(\mathbf{x})$
- Now,  $\mathbf{g}'$  is  $p$ -dim. vector and  $\mathbf{g}''$  is  $p \times p$ -matrix (“Hessian”)
- The multivariate version of the Newton method is motivated by the multivariate Taylor expansion

$$0 = \mathbf{g}'(\mathbf{x}^*) \approx \mathbf{g}'(\mathbf{x}^{(t)}) + \mathbf{g}''(\mathbf{x}^{(t)})(\mathbf{x}^* - \mathbf{x}^{(t)})$$

- The Newton-iteration works as:

$$\mathbf{x}^{(t+1)} = \mathbf{x}^{(t)} - \left(\mathbf{g}''(\mathbf{x}^{(t)})\right)^{-1} \mathbf{g}'(\mathbf{x}^{(t)})$$

# Multivariate Newton

- $\mathbf{x}^{(t+1)} = \mathbf{x}^{(t)} - \left(\mathbf{g}''(\mathbf{x}^{(t)})\right)^{-1} \mathbf{g}'(\mathbf{x}^{(t)})$

- Example:

Let  $g_1$  be the density of  $N\left(\begin{pmatrix} 0 \\ 0 \end{pmatrix}, \begin{pmatrix} 0.6 & 0 \\ 0 & 0.6 \end{pmatrix}\right)$ ,  $g_2$  be density of  $N\left(\begin{pmatrix} 1.5 \\ 1.2 \end{pmatrix}, \begin{pmatrix} 0.5 & 0 \\ 0 & 0.5 \end{pmatrix}\right)$ ,

and  $g = \frac{g_1 + g_2}{2}$ , i.e.

$$g(x_1, x_2) = \frac{1}{4\pi} \left( \frac{1}{0.6} e^{-(x_1^2 + x_2^2)/1.2} + \frac{1}{0.5} e^{-((x_1 - 1.5)^2 + (x_2 - 1.2)^2)} \right)$$

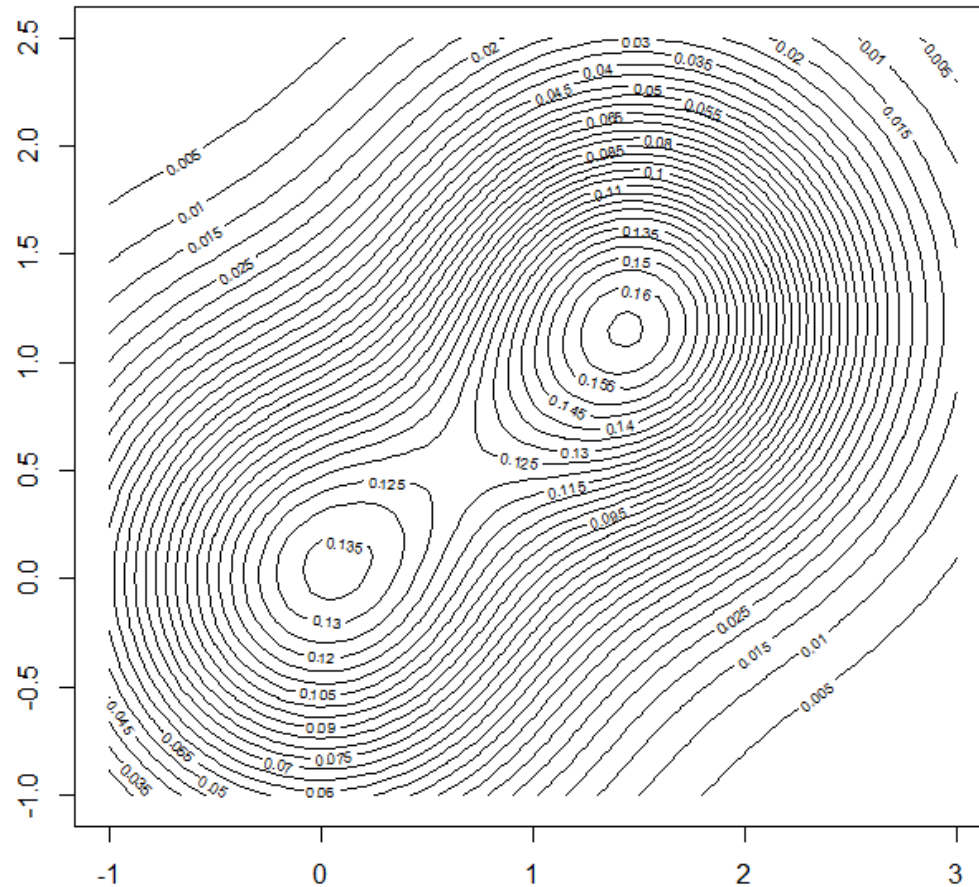
( $g$  is density of a normal mixture distribution).

- Compute point  $\mathbf{x}=(x_1, x_2)$  where density  $g(\mathbf{x})$  maximal.

- Do you have a guess?

# Multivariate Newton

- $g(x_1, x_2) = \frac{1}{4\pi} \left( \frac{1}{0.6} e^{-(x_1^2 + x_2^2)/1.2} + \frac{1}{0.5} e^{-((x_1 - 1.5)^2 + (x_2 - 1.2)^2)} \right)$



# Multivariate Newton

- $\mathbf{x}^{(t+1)} = \mathbf{x}^{(t)} - \left(\mathbf{g}''(\mathbf{x}^{(t)})\right)^{-1} \mathbf{g}'(\mathbf{x}^{(t)})$

- We need  $\mathbf{g}'$  and  $\mathbf{g}''$  of

$$g(x_1, x_2) = \frac{1}{4\pi} \left( \frac{1}{0.6} e^{-(x_1^2 + x_2^2)/(2 \cdot 0.6)} + \frac{1}{0.5} e^{-((x_1 - 1.5)^2 + (x_2 - 1.2)^2)} \right)$$

- $\frac{\partial g}{\partial x_1}(x_1, x_2) = \frac{1}{4\pi} \left( \frac{-2x_1}{1.2 \cdot 0.6} e^{-(x_1^2 + x_2^2)/1.2} + \frac{-2(x_1 - 1.5)}{0.5} e^{-((x_1 - 1.5)^2 + (x_2 - 1.2)^2)} \right)$

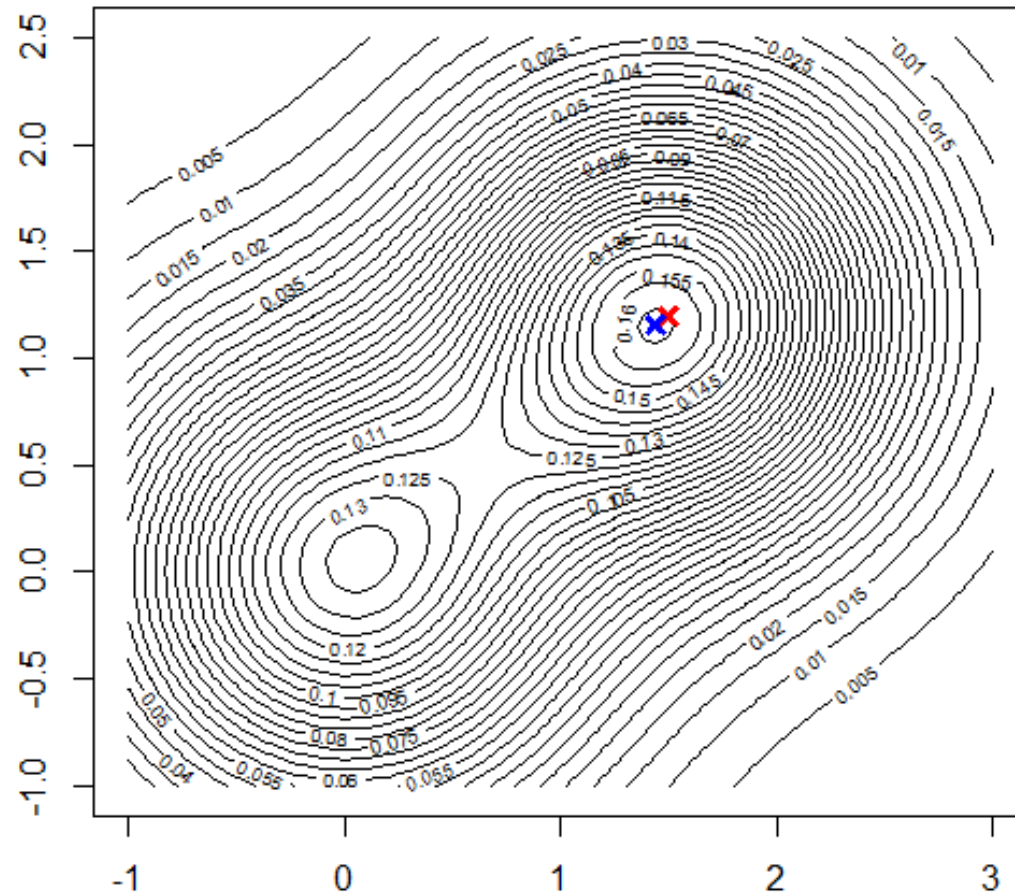
- $\frac{\partial g}{\partial x_2}(x_1, x_2) = \frac{1}{4\pi} \left( \frac{-2x_2}{1.2 \cdot 0.6} e^{-(x_1^2 + x_2^2)/1.2} + \frac{-2(x_2 - 1.2)}{0.5} e^{-((x_1 - 1.5)^2 + (x_2 - 1.2)^2)} \right)$

- $\mathbf{g}'(x_1, x_2) = \begin{pmatrix} \frac{\partial g}{\partial x_1}(x_1, x_2) \\ \frac{\partial g}{\partial x_2}(x_1, x_2) \end{pmatrix}$

- $\frac{\partial^2 g}{\partial^2 x_1}(x_1, x_2) = \dots; \frac{\partial^2 g}{\partial x_1 \partial x_2}(x_1, x_2) = \dots; \frac{\partial^2 g}{\partial^2 x_2}(x_1, x_2) = \dots$  give  $\mathbf{g}''$

# Multivariate Newton

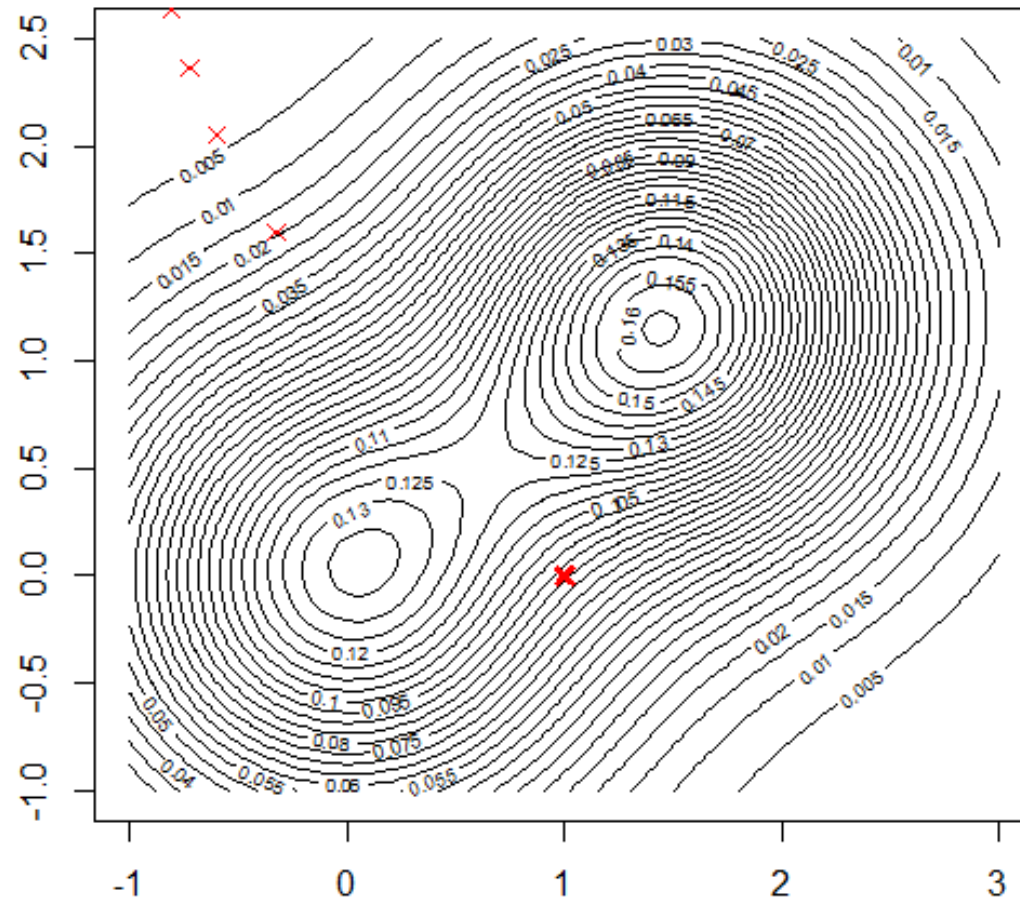
- $$g(x_1, x_2) = \frac{1}{4\pi} \left( \frac{1}{0.6} e^{-(x_1^2 + x_2^2)/1.2} + \frac{1}{0.5} e^{-((x_1 - 1.5)^2 + (x_2 - 1.2)^2)} \right)$$



- Start with  $\mathbf{x}^{(0)} = \begin{pmatrix} 1.5 \\ 1.2 \end{pmatrix}$
- $\mathbf{g}'(\mathbf{x}^{(0)}) = \begin{pmatrix} -0.0153 \\ -0.0123 \end{pmatrix}$
- $\mathbf{g}''(\mathbf{x}^{(0)}) = \begin{pmatrix} -0.2902 & 0.0306 \\ 0.0306 & -0.3040 \end{pmatrix}$
- $\left(\mathbf{g}''(\mathbf{x}^{(0)})\right)^{-1} \mathbf{g}'(\mathbf{x}^{(0)}) = \begin{pmatrix} 0.058 \\ 0.046 \end{pmatrix}$
- $\mathbf{x}^{(1)} = \begin{pmatrix} 1.5 \\ 1.2 \end{pmatrix} - \begin{pmatrix} 0.058 \\ 0.046 \end{pmatrix} = \begin{pmatrix} 1.442 \\ 1.154 \end{pmatrix}$
- $\mathbf{x}^{(2)} = \mathbf{x}^* = \begin{pmatrix} 1.441 \\ 1.153 \end{pmatrix}$

# Multivariate Newton

- $g(x_1, x_2) = \frac{1}{4\pi} \left( \frac{1}{0.6} e^{-(x_1^2+x_2^2)/1.2} + \frac{1}{0.5} e^{-((x_1-1.5)^2+(x_2-1.2)^2)} \right)$



- Start with  $\mathbf{x}^{(0)} = \begin{pmatrix} 1 \\ 0 \end{pmatrix}$

- $\mathbf{g}'(\mathbf{x}^{(0)}) = \begin{pmatrix} -0.0667 \\ +0.0705 \end{pmatrix}$

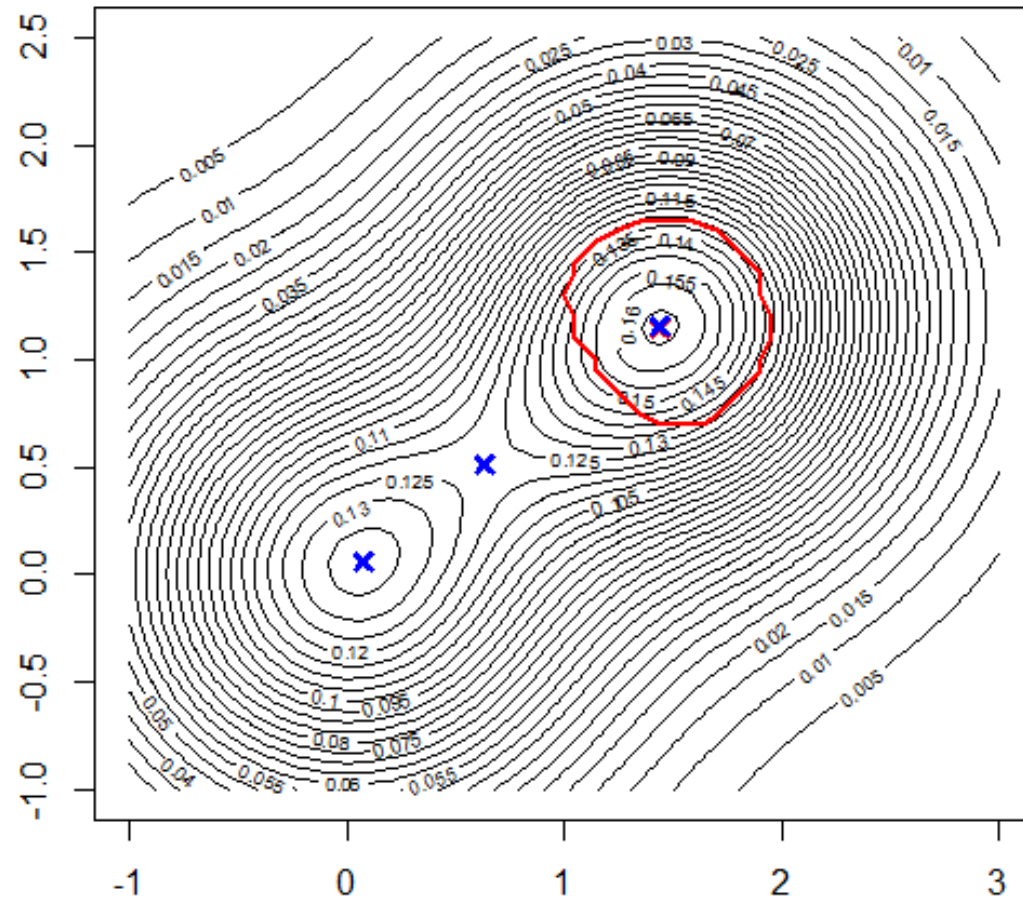
- $\mathbf{g}''(\mathbf{x}^{(0)}) = \begin{pmatrix} 0.0347 & 0.0705 \\ 0.0705 & 0.0144 \end{pmatrix}$

- $\left(\mathbf{g}''(\mathbf{x}^{(0)})\right)^{-1} \mathbf{g}'(\mathbf{x}^{(0)}) = \begin{pmatrix} 1.33 \\ -1.60 \end{pmatrix}$

- $\mathbf{x}^{(1)} = \begin{pmatrix} 1 \\ 0 \end{pmatrix} - \begin{pmatrix} 1.33 \\ -1.6 \end{pmatrix} = \begin{pmatrix} -0.33 \\ 1.6 \end{pmatrix}$

# Multivariate Newton

$$\bullet g(x_1, x_2) = \frac{1}{4\pi} \left( \frac{1}{0.6} e^{-(x_1^2 + x_2^2)/1.2} + \frac{1}{0.5} e^{-((x_1 - 1.5)^2 + (x_2 - 1.2)^2)} \right)$$



- Only starting values within the red-marked area converge to the right global maximum
- Convergence very quick
- Other starting values converge to the local maximum or saddle point (both blue-marked) or diverge while searching for a minimum



# Stopping criteria

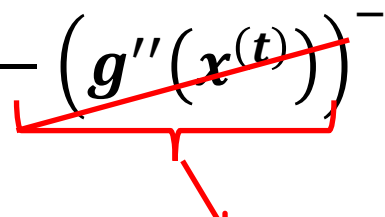
- Stopping criterion e.g.  $(\mathbf{x}^{(t+1)} - \mathbf{x}^{(t)})^T (\mathbf{x}^{(t+1)} - \mathbf{x}^{(t)}) < \epsilon$
- Other stopping criteria:
  - Absolut stopping criterion,  $\|\mathbf{x}^{(t+1)} - \mathbf{x}^{(t)}\| < \epsilon$ ,
  - Relative stopping criterion,  $\|\mathbf{x}^{(t+1)} - \mathbf{x}^{(t)}\| / \|\mathbf{x}^{(t+1)}\| < \epsilon$ ,
  - Modified rel. stopping crit.,  $\frac{\|\mathbf{x}^{(t+1)} - \mathbf{x}^{(t)}\|}{\|\mathbf{x}^{(t+1)}\| + \epsilon} < \epsilon$
  - Different norms  $\|\cdot\|$  can be used

# Today's schedule

- Univariate Optimisation (bi-section, Newton, secant)
- Multivariate Optimisation
  - Analytical opt.
  - Newton
  - **Steepest ascent**
  - **Accelerated steepest ascent**
  - Quasi-Newton

# Steepest ascent method

- When using Newton method, it is not guaranteed that  $g(x)$  increases in each step
- To compute the Hessian  $\mathbf{g}''$  can be difficult
- A method forcing improvements in each step is the steepest ascent method

$$\mathbf{x}^{(t+1)} = \mathbf{x}^{(t)} - \left( \mathbf{g}''(\mathbf{x}^{(t)}) \right)^{-1} \mathbf{g}'(\mathbf{x}^{(t)})$$


$$\mathbf{x}^{(t+1)} = \mathbf{x}^{(t)} + \alpha^{(t)} \mathbf{I} \mathbf{g}'(\mathbf{x}^{(t)})$$

- Other choices instead  $\mathbf{I}$  in formula above possible
- We know that  $g$  will increase for small  $\alpha$

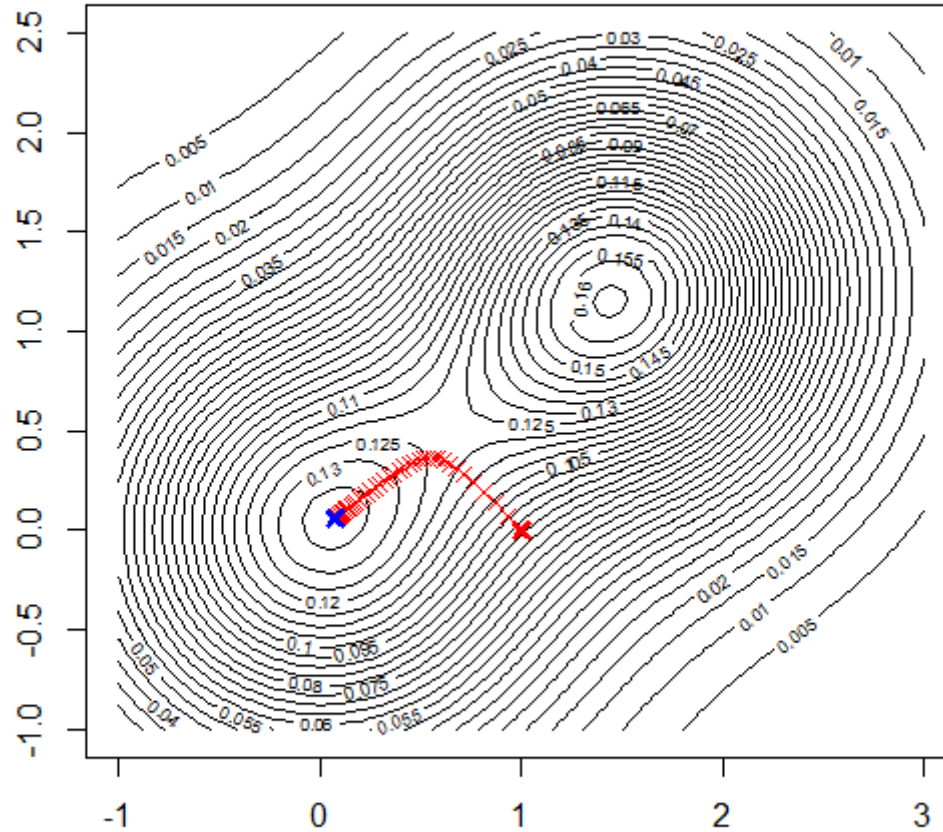
# Backtracking line search (for steepest ascent)

$$\mathbf{x}^{(t+1)} = \mathbf{x}^{(t)} + \alpha^{(t)} \mathbf{I} \mathbf{g}'(\mathbf{x}^{(t)})$$

- We know that  $g$  will increase for small  $\alpha$
- Try  $\alpha^{(t)} = 1$  first
- If  $g$  decreases, half  $\alpha^{(t)}$  until  $g(\mathbf{x}^{(t+1)})$  increases
- More sophisticated is to search  $\alpha$  such that  $g$  becomes maximal

# Steepest ascent

- $g(x_1, x_2) = \frac{1}{4\pi} \left( \frac{1}{0.6} e^{-(x_1^2+x_2^2)/1.2} + \frac{1}{0.5} e^{-((x_1-1.5)^2+(x_2-1.2)^2)} \right)$



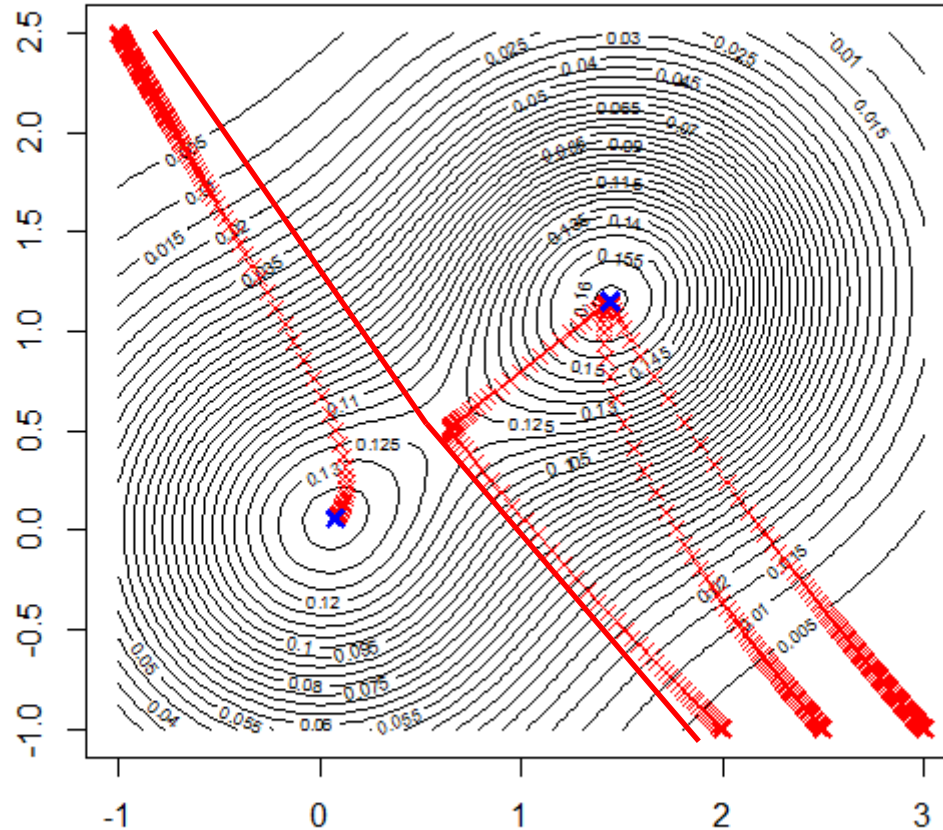
- Start with  $\mathbf{x}^{(0)} = \begin{pmatrix} 1 \\ 0 \end{pmatrix}$

- $\mathbf{g}'(\mathbf{x}^{(0)}) = \begin{pmatrix} -0.0667 \\ +0.0705 \end{pmatrix}$

- $\mathbf{x}^{(1)} = \begin{pmatrix} 1 \\ 0 \end{pmatrix} + \alpha^{(0)} \begin{pmatrix} -0.0667 \\ +0.0705 \end{pmatrix} = \begin{pmatrix} 0.9333 \\ 0.0705 \end{pmatrix}$

# Steepest ascent

$$g(x_1, x_2) = \frac{1}{4\pi} \left( \frac{1}{0.6} e^{-(x_1^2 + x_2^2)/1.2} + \frac{1}{0.5} e^{-((x_1 - 1.5)^2 + (x_2 - 1.2)^2)} \right)$$

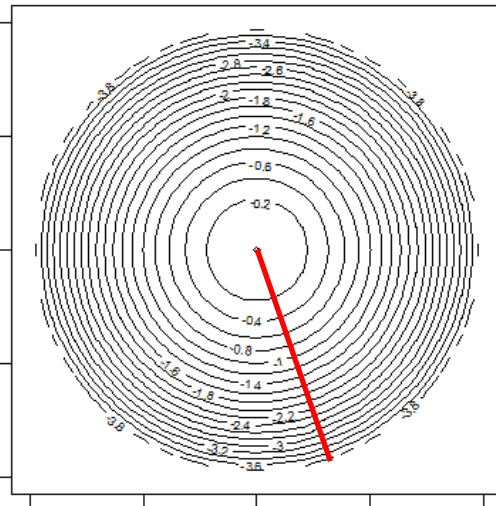


- Start with  $x^{(0)} = \begin{pmatrix} -1 \\ 2.5 \end{pmatrix}, \begin{pmatrix} 2 \\ -1 \end{pmatrix}, \begin{pmatrix} 2.5 \\ -1 \end{pmatrix}, \begin{pmatrix} 3 \\ -1 \end{pmatrix}$
- All these paths converge to either the global or local maximum
- Convergence is much slower than for Newton
- Depending on convergence criterion and alpha-rule, convergence not always guaranteed

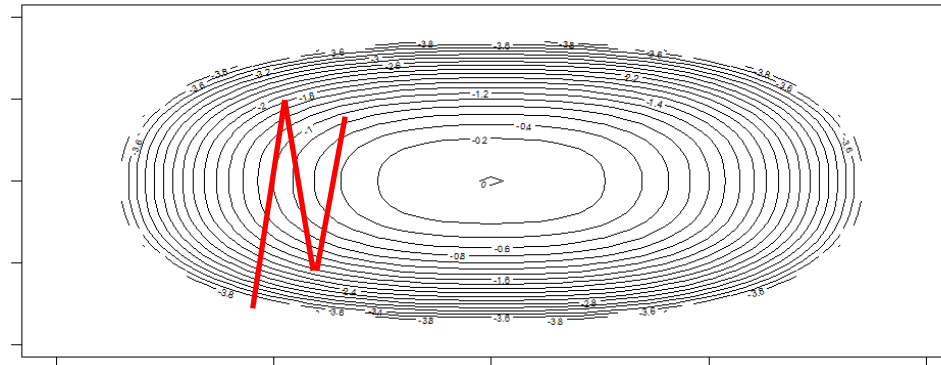
# Steepest ascent



Globen, Stockholm – by Arild Vågen,  
CC BY-SA 4.0,  
[https://commons.wikimedia.org/wiki/  
File:Globen\\_September\\_2014\\_02.jpg](https://commons.wikimedia.org/wiki/File:Globen_September_2014_02.jpg)

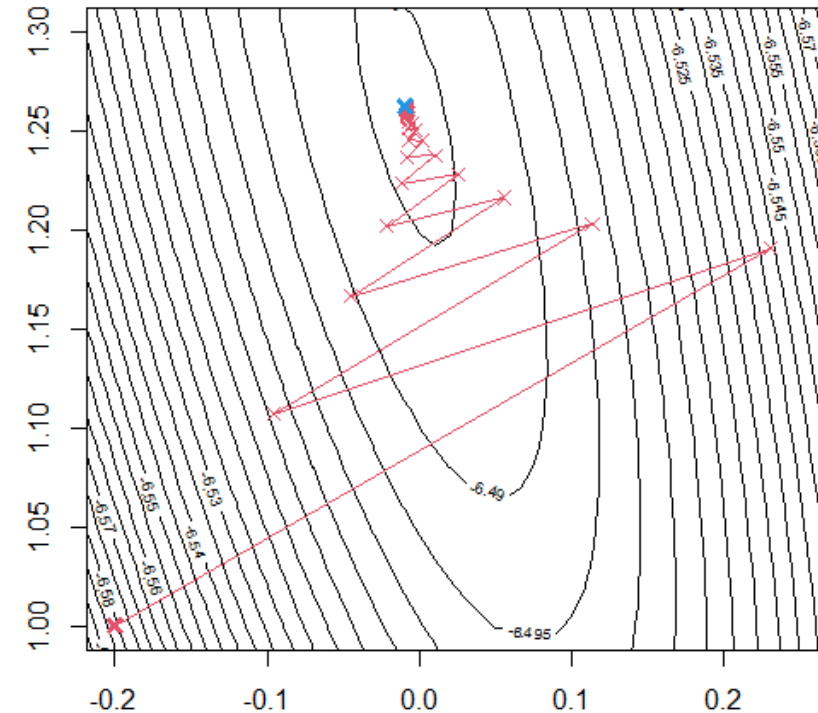


Uluru, Australia – by Stuart Edwards, CC BY-SA 3.0,  
<https://commons.wikimedia.org/w/index.php?curid=1650537>



# Steepest ascent: idea for acceleration

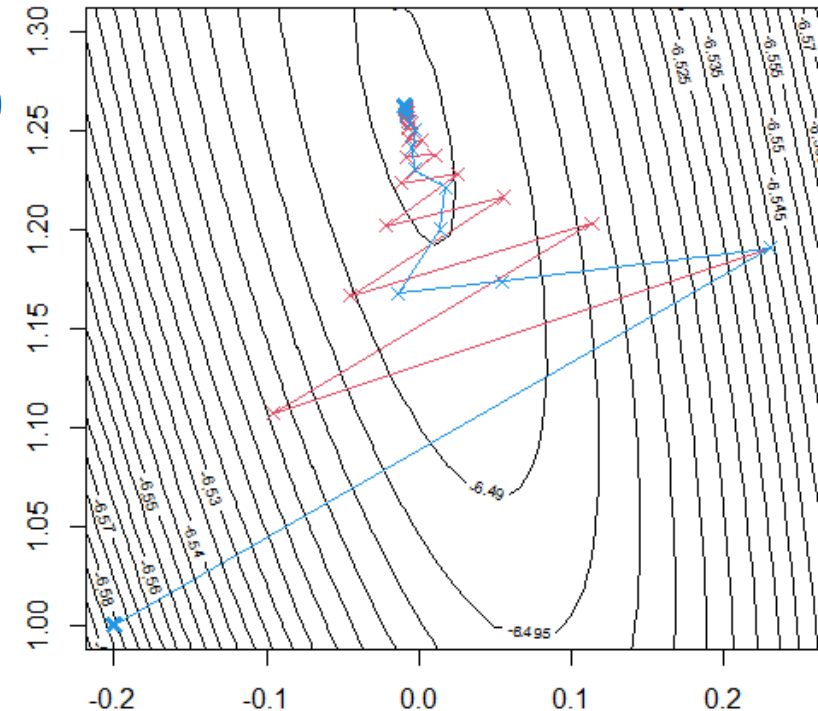
- Example: ML computation for a two-parameter model with steepest ascent, with fixed  $\alpha^{(t)} = 0.667$  (no backtracking)
- Zick-zack path is common and slows down convergence
- Idea to reduce/avoid this issue: use information from last iteration about "momentum" of search path
- Called: **Accelerated steepest ascent** (or steepest ascent with momentum)





# Accelerated steepest ascent: Polyak's momentum

- $\mathbf{x}^{(t+1)} = \mathbf{x}^{(t)} + \alpha^{(t)} \mathbf{g}'(\mathbf{x}^{(t)}) + \beta(\mathbf{x}^{(t)} - \mathbf{x}^{(t-1)})$
- Polyak="gradient+momentum"
- Steepest ascent ( $\alpha^{(t)} = 0.667$ )
- with momentum ( $\beta = 0.35$ )
- Called also *heavy-ball method*
- Adding momentum reduces number of iterations from **31** to **21** in this example
- Works well in many situations
- Examples exist where Polyak's method fails to converge



# Accelerated steepest ascent: Nesterov's momentum

- $\mathbf{x}^{(t+1)} = \mathbf{x}^{(t)} + \alpha^{(t)} \mathbf{g}'(\mathbf{x}^{(t)} + \beta(\mathbf{x}^{(t)} - \mathbf{x}^{(t-1)})) + \beta(\mathbf{x}^{(t)} - \mathbf{x}^{(t-1)})$
- Nesterov = “lookahead gradient + momentum”
- Ideally, this method has the capacity
  - to dampen oscillations and
  - to accelerate if the search path is in right direction
- Nesterov's accelerated ascent has better convergence rate as steepest ascent

# Parametrisation of accelerated methods

- Polyak's accelerated steepest ascent

$$\mathbf{x}^{(t+1)} = \mathbf{x}^{(t)} + \alpha \mathbf{g}'(\mathbf{x}^{(t)}) + \beta(\mathbf{x}^{(t)} - \mathbf{x}^{(t-1)})$$

can be written also as

$$\mathbf{x}^{(t+1)} = \mathbf{x}^{(t)} + \alpha \mathbf{v}^{(t+1)}$$

$$\mathbf{v}^{(t+1)} = \beta \mathbf{v}^{(t)} + \mathbf{g}'(\mathbf{x}^{(t)})$$

- Nesterov's accelerated steepest ascent

$$\mathbf{x}^{(t+1)} = \mathbf{x}^{(t)} + \alpha \mathbf{g}'(\mathbf{x}^{(t)} + \beta(\mathbf{x}^{(t)} - \mathbf{x}^{(t-1)})) + \beta(\mathbf{x}^{(t)} - \mathbf{x}^{(t-1)})$$

can be written also as

$$\mathbf{x}^{(t+1)} = \mathbf{x}^{(t)} + \alpha \mathbf{v}^{(t+1)}$$

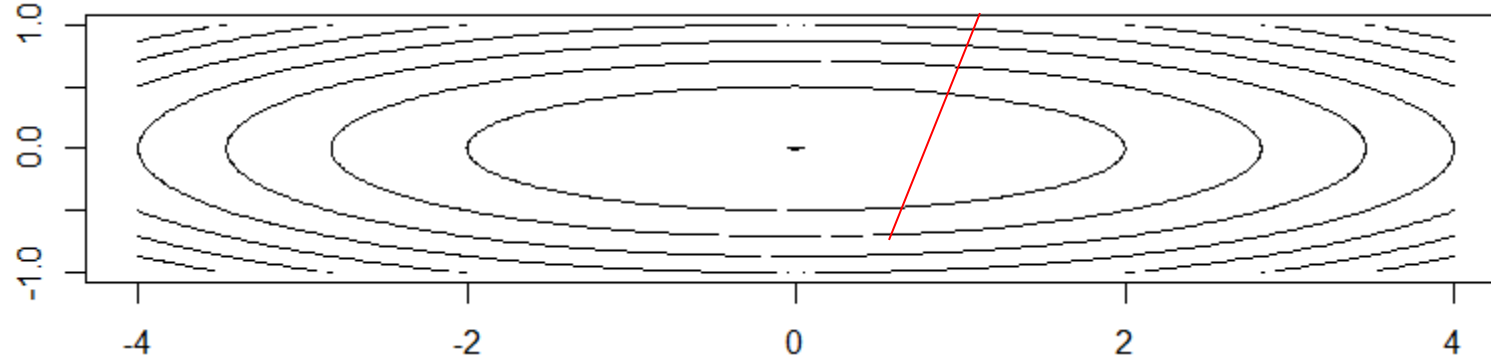
$$\mathbf{v}^{(t+1)} = \beta \mathbf{v}^{(t)} + \mathbf{g}'(\mathbf{x}^{(t)} + \alpha \beta \mathbf{v}^{(t)})$$

# Steepest ascent: optimal choice of step size

- $\mathbf{x}^{(t+1)} = \mathbf{x}^{(t)} + \alpha \mathbf{g}'(\mathbf{x}^{(t)})$
- Example:  
 $g(\mathbf{x}) = -\frac{1}{2} \mathbf{x}^T \mathbf{A} \mathbf{x} + \mathbf{b}^T \mathbf{x}$ ,  $\mathbf{A}$  symmetric  $p \times p$  and of full rank
- $\mathbf{g}'(\mathbf{x}) = \mathbf{b} - \mathbf{A} \mathbf{x}$
- To keep things simple (and to avoid a change of basis and some more linear algebra...), we use  $\mathbf{b} = \mathbf{0}$ ,  $\mathbf{A}$ =diagonal (i.e. eigenvalues in diagonal),  $p=2$
- $\mathbf{A} = \begin{pmatrix} \lambda_1 & 0 \\ 0 & \lambda_2 \end{pmatrix}$ ,  $\mathbf{g}'(\mathbf{x}) = \begin{pmatrix} -\lambda_1 x_1 \\ -\lambda_2 x_2 \end{pmatrix}$ ,  $\lambda_1, \lambda_2 > 0$
- Then, steepest ascent is:
- $x_i^{(t+1)} = (1 - \alpha \lambda_i) x_i^{(t)} = (1 - \alpha \lambda_i)^{t+1} x_i^{(0)}$

# Steepest ascent: optimal choice of step size

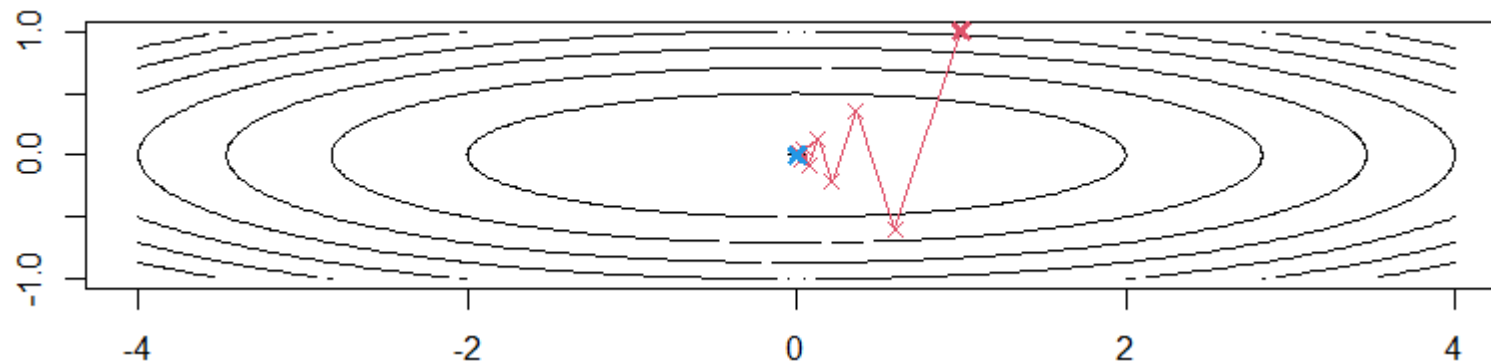
- $\mathbf{x}^{(t+1)} = \mathbf{x}^{(t)} + \alpha \mathbf{g}'(\mathbf{x}^{(t)})$
- Example:  $g(\mathbf{x}) = -\frac{1}{2} \mathbf{x}^T \begin{pmatrix} \lambda_1 & 0 \\ 0 & \lambda_2 \end{pmatrix} \mathbf{x}$ ,  $\mathbf{g}'(\mathbf{x}) = \begin{pmatrix} -\lambda_1 x_1 \\ -\lambda_2 x_2 \end{pmatrix}$ ,  $\mathbf{x}^{(0)} = \begin{pmatrix} 1 \\ 1 \end{pmatrix}$
- Steepest ascent:  $x_1^{(t+1)} = (1 - \alpha\lambda_1)^{t+1}$ ,  $x_2^{(t+1)} = (1 - \alpha\lambda_2)^{t+1}$
- For  $\lambda_1 = \frac{1}{2}$ ,  $\lambda_2 = 2$ :



- Fastest convergence attained if  $\alpha$  such that  $\rho = \max\{|1 - \alpha\lambda_1|, |1 - \alpha\lambda_2|\}$  is as small as possible

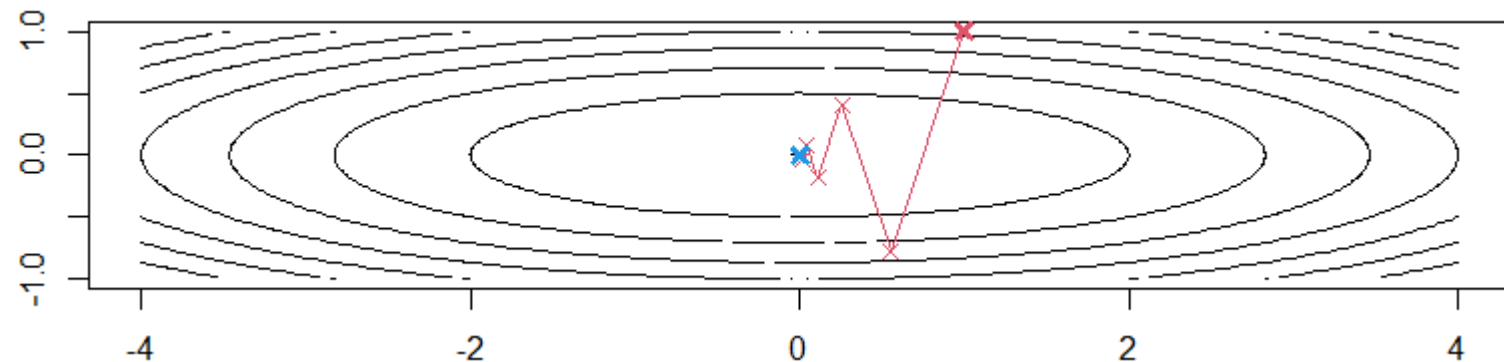
# Steepest ascent: optimal choice of step size

- Steepest ascent:  $x_1^{(t+1)} = (1 - \alpha\lambda_1)^{t+1}$ ,  $x_2^{(t+1)} = (1 - \alpha\lambda_2)^{t+1}$
- Fastest convergence attained if  $\alpha$  such that  $\rho = \max\{|1 - \alpha\lambda_1|, |1 - \alpha\lambda_2|\}$  is as small as possible
- Fulfilled for  $\alpha = \frac{2}{\lambda_1 + \lambda_2}$  and then  $\rho = \frac{\kappa - 1}{\kappa + 1}$  with  $\kappa = \lambda_2 / \lambda_1$
- $\rho$  is convergence rate;  $\kappa$  is condition number
- For example, with  $\lambda_1 = \frac{1}{2}$ ,  $\lambda_2 = 2$ :  $\rho = \frac{3}{5}$ ,  $\alpha = \frac{4}{5}$ .



# Accelerated steepest ascent: choice of hyperparameters

- Steepest ascent: convergence rate  $\rho = \frac{\kappa-1}{\kappa+1}$  with  $\kappa = \frac{\lambda_{max}}{\lambda_{min}}$
- Accelerated steepest ascent:
  - Best convergence rate:  $\rho = \left(\frac{\sqrt{\kappa}-1}{\sqrt{\kappa}+1}\right)$
  - Optimal step size:  $\alpha = \frac{(1+\rho)^2}{\lambda_{max}} = \frac{(1-\rho)^2}{\lambda_{min}}$
  - Optimal momentum:  $\beta = \rho^2$
- For example, with  $\lambda_1 = \frac{1}{2}, \lambda_2 = 2$ :  
 $\rho = \frac{1}{3}, \alpha = \frac{8}{9}, \beta = \frac{1}{9}$ .



# (Accelerated) steepest ascent: convergence

- Convergence rate for  $\kappa = \frac{\lambda_{max}}{\lambda_{min}}$ :
  - Steepest ascent:  $\rho = \frac{\kappa-1}{\kappa+1}$
  - Accelerated steepest ascent:  $\rho = \left(\frac{\sqrt{\kappa}-1}{\sqrt{\kappa}+1}\right)$
- $\lim_{t \rightarrow \infty} \frac{\|x^{(t+1)} - x^*\|}{\|x^{(t)} - x^*\|}^q = \rho$ 
  - convergence order; here  $q = 1$
  - convergence rate
- Example  $\kappa = 100$  (“ill-conditioned”):
  - $\frac{\kappa-1}{\kappa+1} = \frac{99}{101}$ ;  $\left(\frac{\kappa-1}{\kappa+1}\right)^t = 1, 0.98, \dots, 0.82, \dots, 0.14, \dots$ 
    - $t=10$
    - $t=100$
  - $\frac{\sqrt{\kappa}-1}{\sqrt{\kappa}+1} = \frac{9}{11}$ ;  $\left(\frac{\sqrt{\kappa}-1}{\sqrt{\kappa}+1}\right)^t = 1, 0.82, \dots, 0.13, \dots, 1.9 \cdot 10^{-9}, \dots$



# Today's schedule

- Univariate Optimisation (bi-section, Newton, secant)
- Multivariate Optimisation
  - Analytical opt.
  - Newton
  - Steepest ascent
  - Accelerated steepest ascent
  - **Quasi-Newton**

# Quasi-Newton

- Steepest ascent and Newton method have iteration

$$\mathbf{x}^{(t+1)} = \mathbf{x}^{(t)} - (\mathbf{M}^{(t)})^{-1} \mathbf{g}'(\mathbf{x}^{(t)})$$

with  $\mathbf{M}^{(t)} = \mathbf{g}''(\mathbf{x}^{(t)})$  for the Newton method and

with  $(\mathbf{M}^{(t)})^{-1} = -\alpha_t \mathbf{I}$  for the steepest ascent method

- A disadvantage of Newton is the need to calculate the Hessian  $\mathbf{g}''(\mathbf{x}^{(t)})$  in each iteration
- A disadvantage of steepest ascent is that no information about the curvature is used
- We can monitor the computed gradients  $\mathbf{g}'(\mathbf{x}^{(t)})$  and their change gives information about the curvature of  $g$

# Quasi-Newton

- Steepest ascent and Newton method have iteration

$$\mathbf{x}^{(t+1)} = \mathbf{x}^{(t)} - (\mathbf{M}^{(t)})^{-1} \mathbf{g}'(\mathbf{x}^{(t)})$$

- Newton ( $\mathbf{M}^{(t)} = \mathbf{g}''(\mathbf{x}^{(t)})$ ) was motivated with the multidimensional Taylor expansion

$$\mathbf{g}'(\mathbf{x}^*) \approx \mathbf{g}'(\mathbf{x}^{(t)}) + \mathbf{g}''(\mathbf{x}^{(t)})(\mathbf{x}^* - \mathbf{x}^{(t)})$$

or

$$\mathbf{g}'(\mathbf{x}^*) - \mathbf{g}'(\mathbf{x}^{(t)}) \approx \mathbf{g}''(\mathbf{x}^{(t)})(\mathbf{x}^* - \mathbf{x}^{(t)})$$

- We want to use approximations  $\mathbf{M}^{(t+1)}$  to  $\mathbf{g}''(\mathbf{x}^{(t)})$  which fulfil this relation when  $\mathbf{x}^*$  is replaced by  $\mathbf{x}^{(t+1)}$ :

$$\mathbf{g}'(\mathbf{x}^{(t+1)}) - \mathbf{g}'(\mathbf{x}^{(t)}) = \mathbf{M}^{(t+1)}(\mathbf{x}^{(t+1)} - \mathbf{x}^{(t)})$$

- This condition is called secant condition
- There are multiple solutions to the secant condition

# Quasi-Newton

- Steepest ascent and Newton method have iteration

$$\mathbf{x}^{(t+1)} = \mathbf{x}^{(t)} - (\mathbf{M}^{(t)})^{-1} \mathbf{g}'(\mathbf{x}^{(t)})$$

- Secant condition:

$$\mathbf{g}'(\mathbf{x}^{(t+1)}) - \mathbf{g}'(\mathbf{x}^{(t)}) = \mathbf{M}^{(t+1)} (\mathbf{x}^{(t+1)} - \mathbf{x}^{(t)})$$

- Or, with  $\mathbf{y}^{(t)} = \mathbf{g}'(\mathbf{x}^{(t+1)}) - \mathbf{g}'(\mathbf{x}^{(t)})$  and  $\mathbf{z}^{(t)} = \mathbf{x}^{(t+1)} - \mathbf{x}^{(t)}$ :

$$\mathbf{y}^{(t)} = \mathbf{M}^{(t+1)} \mathbf{z}^{(t)}$$

- Suggestion from Broyden, Fletcher, Goldfarb, and Shanno (BFGS; 4 publications in 1970) fulfilling secant condition:

$$\mathbf{M}^{(t+1)} = \mathbf{M}^{(t)} - \frac{\mathbf{M}^{(t)} \mathbf{z}^{(t)} (\mathbf{M}^{(t)} \mathbf{z}^{(t)})^T}{\mathbf{z}^{(t)T} \mathbf{M}^{(t)} \mathbf{z}^{(t)}} + \frac{\mathbf{y}^{(t)} \mathbf{y}^{(t)T}}{\mathbf{y}^{(t)T} \mathbf{z}^{(t)}}$$

# Quasi-Newton

- The BFGS (quasi-Newton) method has iteration

$$\mathbf{x}^{(t+1)} = \mathbf{x}^{(t)} - (\mathbf{M}^{(t)})^{-1} \mathbf{g}'(\mathbf{x}^{(t)})$$

and

$$\mathbf{M}^{(t+1)} = \mathbf{M}^{(t)} - \frac{\mathbf{M}^{(t)} \mathbf{z}^{(t)} (\mathbf{M}^{(t)} \mathbf{z}^{(t)})^T}{\mathbf{z}^{(t)T} \mathbf{M}^{(t)} \mathbf{z}^{(t)}} + \frac{\mathbf{y}^{(t)} \mathbf{y}^{(t)T}}{\mathbf{y}^{(t)T} \mathbf{z}^{(t)}}$$

- Ascent is not ensured but backtracking (stepsize-halving) can be used as for steepest ascent to ensure it:

$$\mathbf{x}^{(t+1)} = \mathbf{x}^{(t)} - \alpha_t (\mathbf{M}^{(t)})^{-1} \mathbf{g}'(\mathbf{x}^{(t)})$$

- The **R** function `optim` includes the quasi-Newton BFGS
- Convergence of quasi-Newton methods are faster than linear but slower than quadratic (some assumptions necessary; see e.g. Nocedal and Wright, 2006, Theorem 3.7)

# Convergence order for deterministic algorithms

- Recall: Convergence order and convergence rate

$$\frac{\{g(\mathbf{x}^{(t+1)}) - g(\mathbf{x}^*)\}}{\{g(\mathbf{x}^{(t)}) - g(\mathbf{x}^*)\}^q} \rightarrow c \quad (\text{for } t \rightarrow \infty)$$

- q is convergence order (q=1, 0<c<1 linear; q=2, 0<c<1 quadratic)
- c is convergence rate
- Under certain assumption, we have following orders:

<b>Uni-dimensional</b>	<b>Bisection</b> order = roughly 1*		<b>Secant</b> order = $(1 + \sqrt{5})/2$	<b>Newton</b> order = 2
<b>Multi-dimensional</b>		<b>Steepest ascent</b> order = 1	<b>Quasi-Newton</b> order > 1**	<b>Newton</b> order = 2

\*strictly, the above criterion cannot be proven for bisection

\*\*criterion above fulfilled for q=1 and c=0; “superlinear”

# Convergence speed for an example function

- The convergence of BFGS and Newton can be extremely fast in praxis compared to steepest ascent/descent
- Example from Nocedal and Wright (2006), chapter 6: Rosenbrock function  $g(\mathbf{x}) = 100(x_2 - x_1^2)^2 + (1 - x_1)^2$ , starting point  $(-1.2, 1)$ , optimum at  $(1,1)$ .

#iterations until error  $< 10^{-5}$ :

- Steepest descent            5264
- BFGS                            34
- Newton                         21

# Assignments

- Topic 1+2: March 17 until April 3
- Topic 3: April 4 until April 17
- Topic 4: April 18 until May 1
- Topic 5: May 2 until May 15
- Topic 6+7: May 17 until June 7
- Second chance for Topic 1-7: until **September 30 (no extension!)**