

# Optimisation algorithms in Statistics I, lecture 1

Frank Miller, Department of Statistics; Stockholm University

October 2, 2020



# Course schedule

- Topic 1: **Gradient based algorithms**  
Lectures: October 2; Time 10-12, 13-15 (online, Zoom)
- Topic 2: **Stochastic gradient based algorithms**  
Lecture: October 13; Time: 9-12 (online, Zoom)
- Topic 3: **Gradient free algorithms**  
Lecture: October 23; Time 9-12 (online, Zoom)
- Topic 4: **Optimisation with restrictions**  
Lecture: November 6, Time 9-12 (online, Zoom)

Course homepage: <http://gauss.stat.su.se/phd/oasi/>

Includes 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  $g(\mathbf{x}^*) = \max g(\mathbf{x})$

# Optimisation in statistics

- Generic optimisation problem:
  - $x$   $p$ -dimensional vector,  $g: \mathbb{R}^p \rightarrow \mathbb{R}$  function
  - We search  $x^*$  with  $g(x^*) = \max g(x)$
- Typical situation in machine learning situations, for Maximum Likelihood, and Least Squares problems:  
 $g = \sum_{i=1}^n g_i$  with a (large) sample size  $n$  where  $g_i: \mathbb{R}^p \rightarrow \mathbb{R}$
- Note that a minimisation problem  $g(x^*) = \min g(x)$  can easily be solved by looking at  $-g$  if an algorithm for maximisation is available



# Least squares estimation (LSE)

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



# Variations of least squares estimation

- While we have a solution for the LSE, we have no algebraic solution 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  $g(\text{design}^*) = \max g(\text{design})$



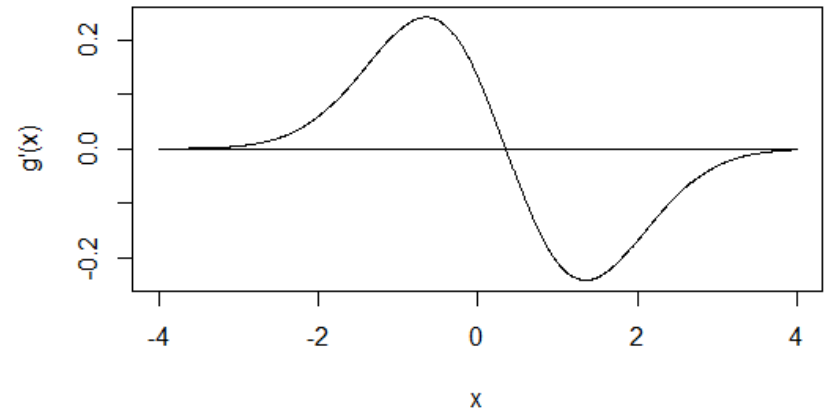
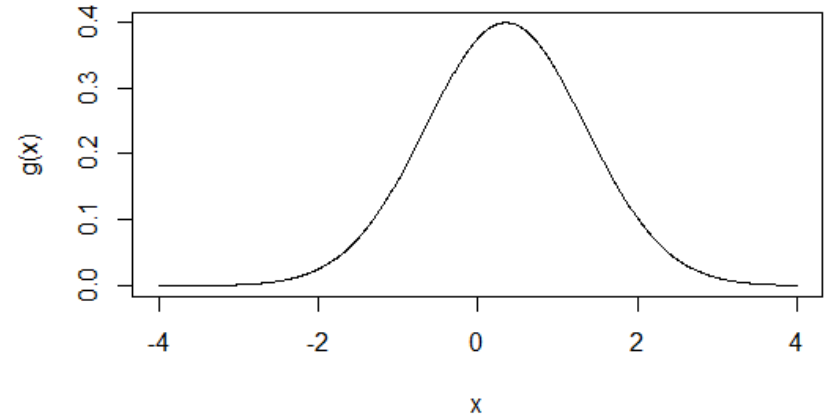


# Univariate optimisation

- $x$  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 method

- 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:

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

# Optimisation: Newton(-Raphson) method

- $x$  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:

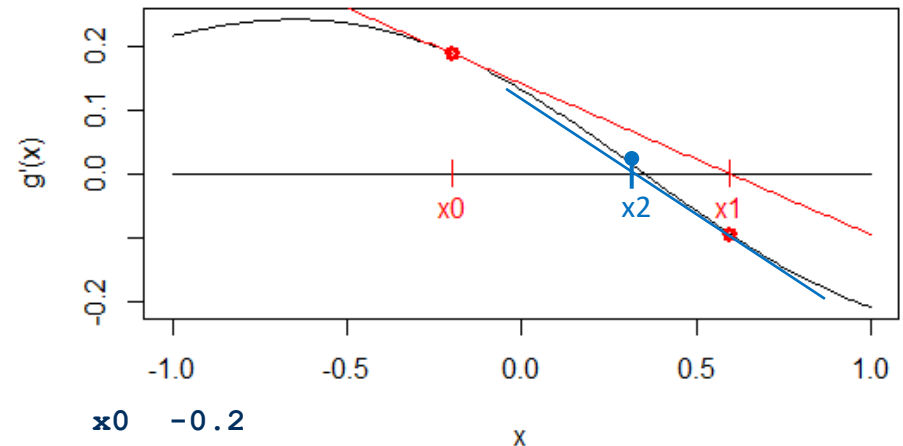
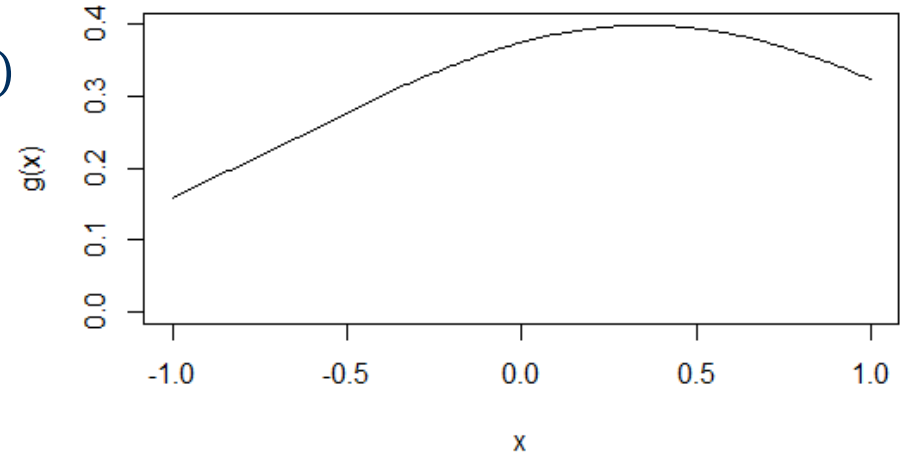
$$x^{(t+1)} = x^{(t)} - g'(x^{(t)})/g''(x^{(t)})$$



# Optimisation: Newton(-Raphson) method

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

- + Newton method is fast
- Requires existence and computation of  $g''$



```

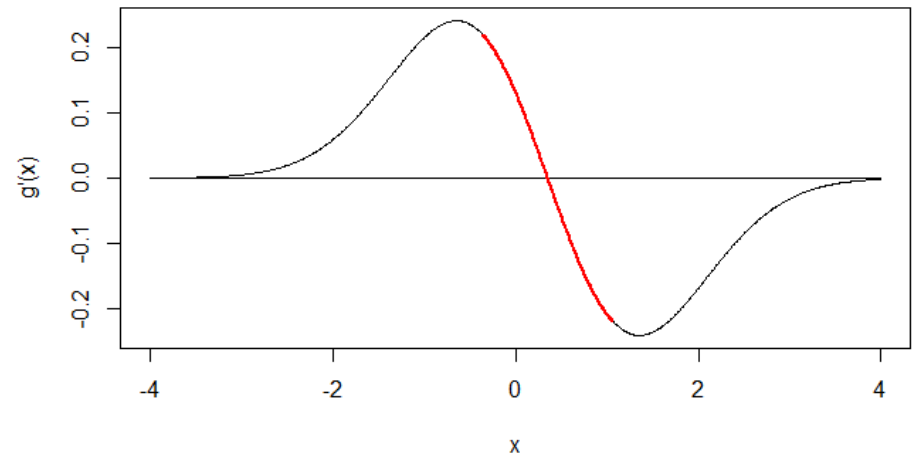
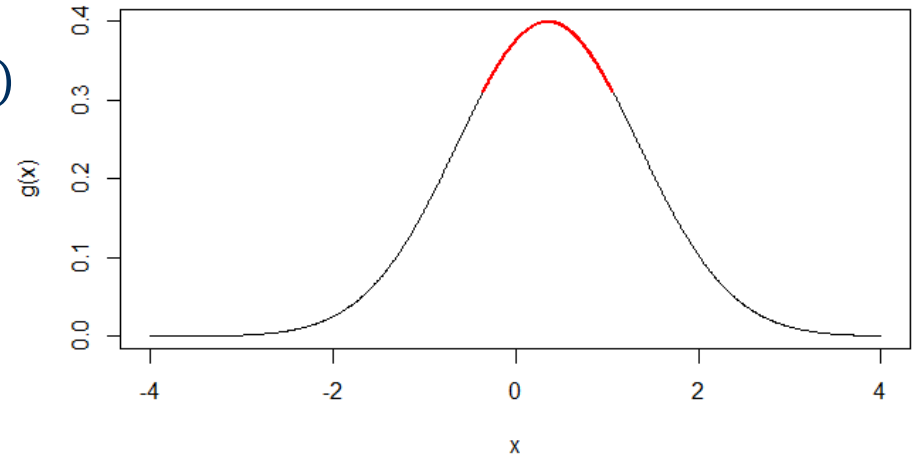
x0 -0.2
x1 0.5981
x2 0.337996
x3 0.353557
x4 0.353553
x5 0.353553

```

STOP

# Optimisation: Newton(-Raphson) method

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



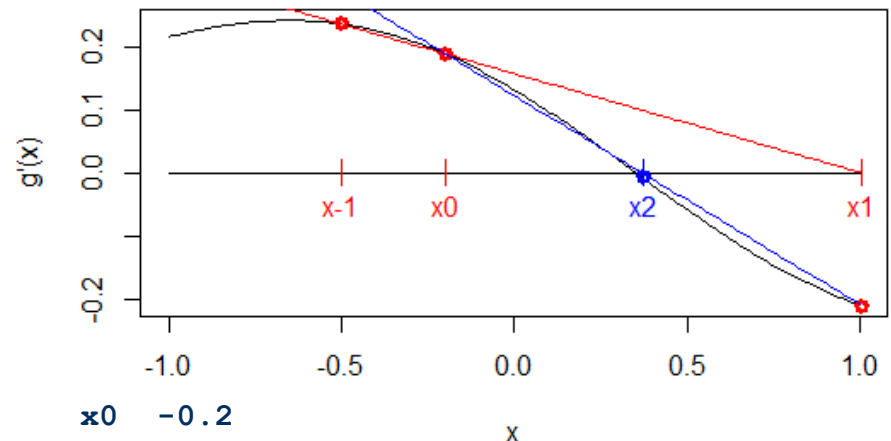
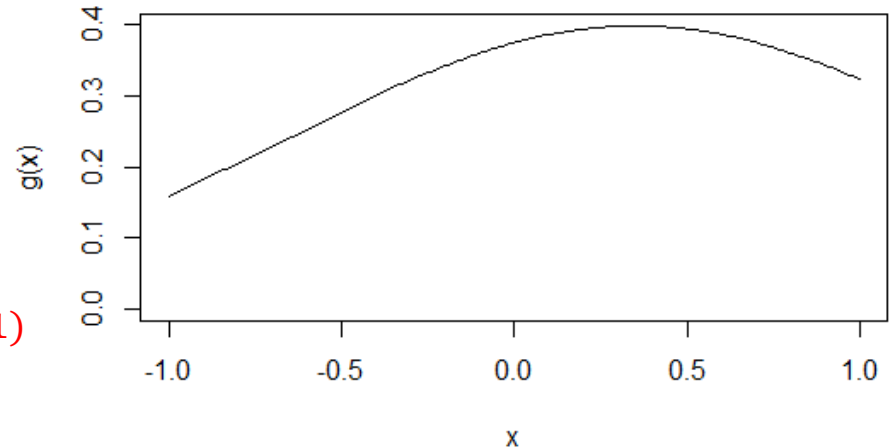
# Optimisation: Secant method

- $x$  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  $[g'(x^{(t)}) - g'(x^{(t-1)})]/(x^{(t)} - x^{(t-1)})$

# Optimisation: 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



# Optimisation: Convergence speed

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

# Optimisation: Comparison of 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 $\beta=1$	quite fast $\beta = \frac{1 + \sqrt{5}}{2}$	fast $\beta=2$

$$= 1.62$$

- $\lim_{t \rightarrow \infty} \frac{\varepsilon^{(t+1)}}{(\varepsilon^{(t)})^\beta} = c$

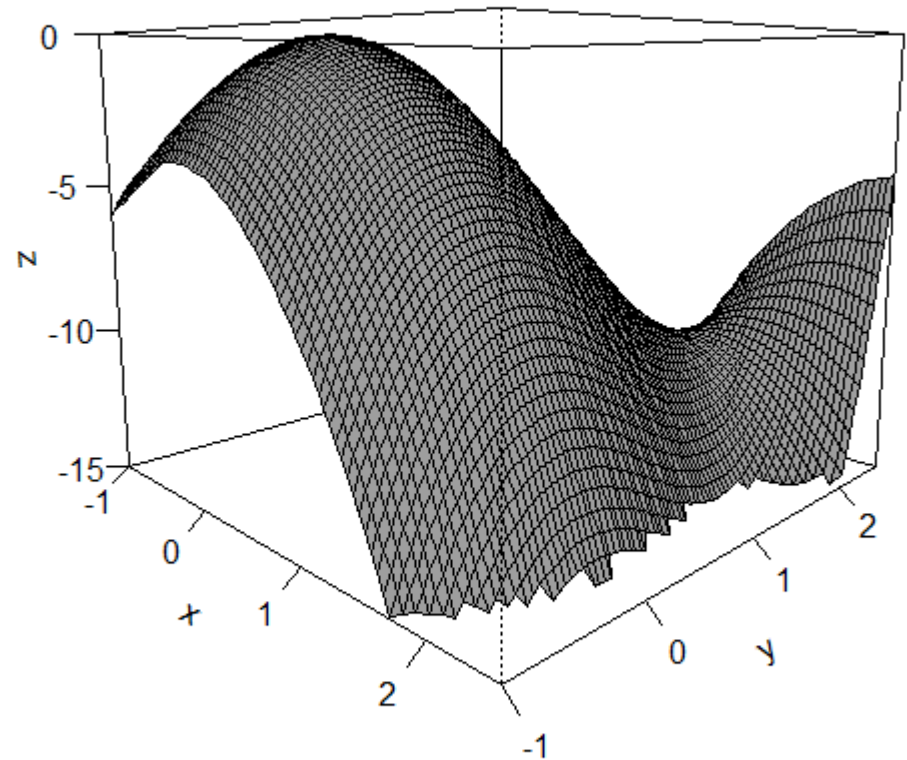
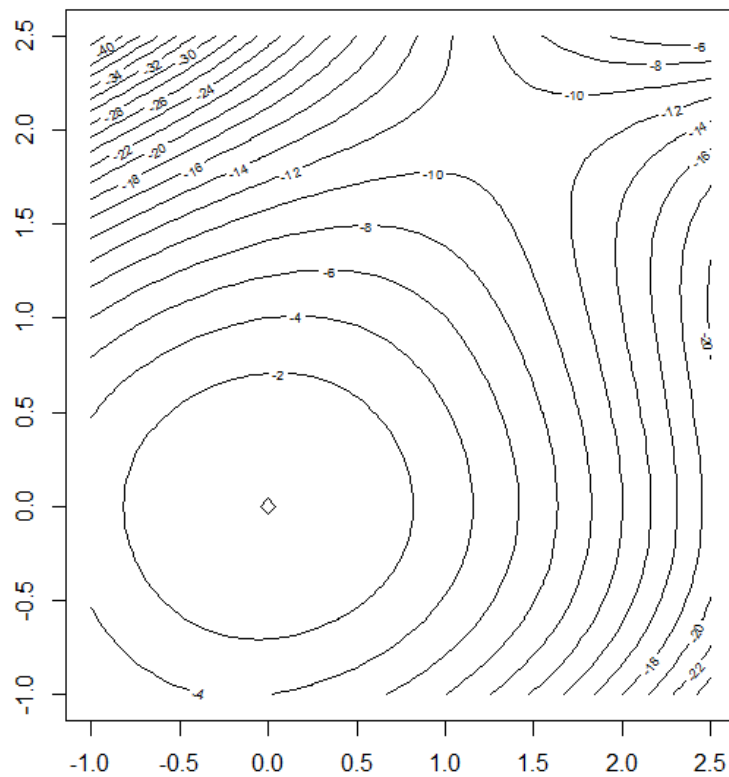


# Multivariate optimisation – gradient and Hessian

- $g\left(\begin{matrix} x_1 \\ \vdots \\ x_p \end{matrix}\right)$  is a real-valued function
- $g'\left(\begin{matrix} x_1 \\ \vdots \\ x_p \end{matrix}\right) = \begin{pmatrix} \frac{\partial g}{\partial x_1}(\mathbf{x}) \\ \vdots \\ \frac{\partial g}{\partial x_p}(\mathbf{x}) \end{pmatrix}$  is the gradient,  $\mathbf{x} = \begin{pmatrix} x_1 \\ \vdots \\ x_p \end{pmatrix}$
- $g''\left(\begin{matrix} x_1 \\ \vdots \\ x_p \end{matrix}\right) = \begin{pmatrix} \frac{\partial g}{\partial x_1 \partial x_1}(\mathbf{x}) & \dots & \frac{\partial g}{\partial x_1 \partial x_p}(\mathbf{x}) \\ \vdots & & \vdots \\ \frac{\partial g}{\partial x_1 \partial x_p}(\mathbf{x}) & \dots & \frac{\partial g}{\partial x_p \partial x_p}(\mathbf{x}) \end{pmatrix}$  is the Hessian matrix

# Bivariate optimisation – visualisation

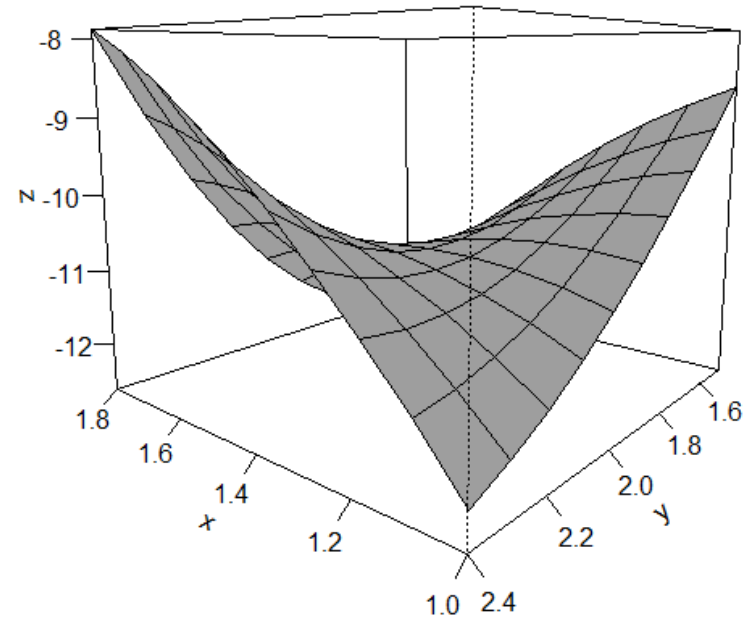
- $g\begin{pmatrix} x \\ y \end{pmatrix} = -3x^2 - 4y^2 + xy^3$



Figures can be drawn using R-core-

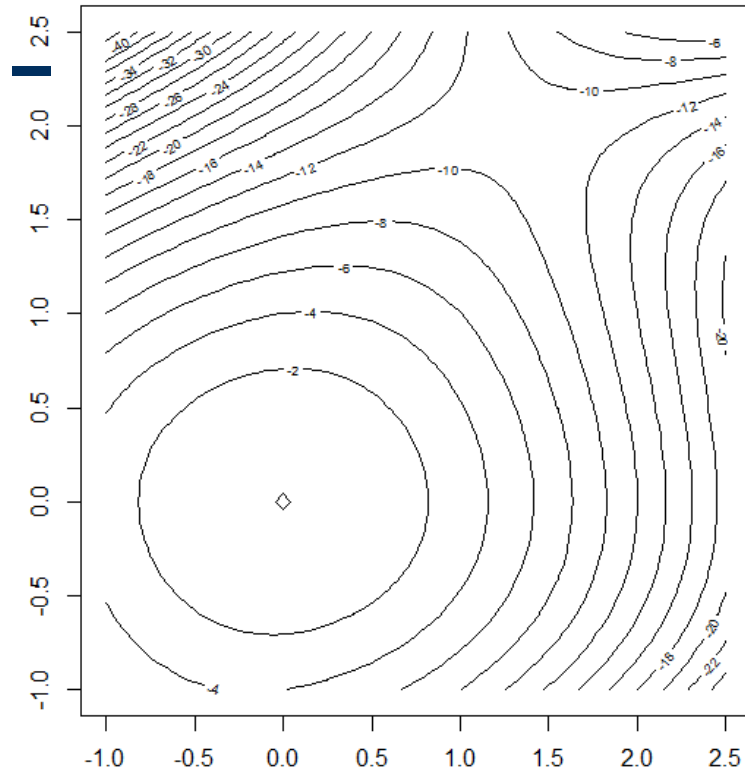
functions `contour` and `persp`

# 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:  
[AnalytOpt\\_OASI2020.pdf](#)
- Maximum at  $(0,0)$ , saddle point at  $(4/3,2)$



# Multivariate optimisation – Newton meth.

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

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

- The Newton-iteration works as:

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

# Multivariate optimisation – Newton meth.

- $\mathbf{x}^{(t+1)} = \mathbf{x}^{(t)} - \left(g''(\mathbf{x}^{(t)})\right)^{-1} g'(\mathbf{x}^{(t)})$
- 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  $g = (g_1 + g_2)/2$ , 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).

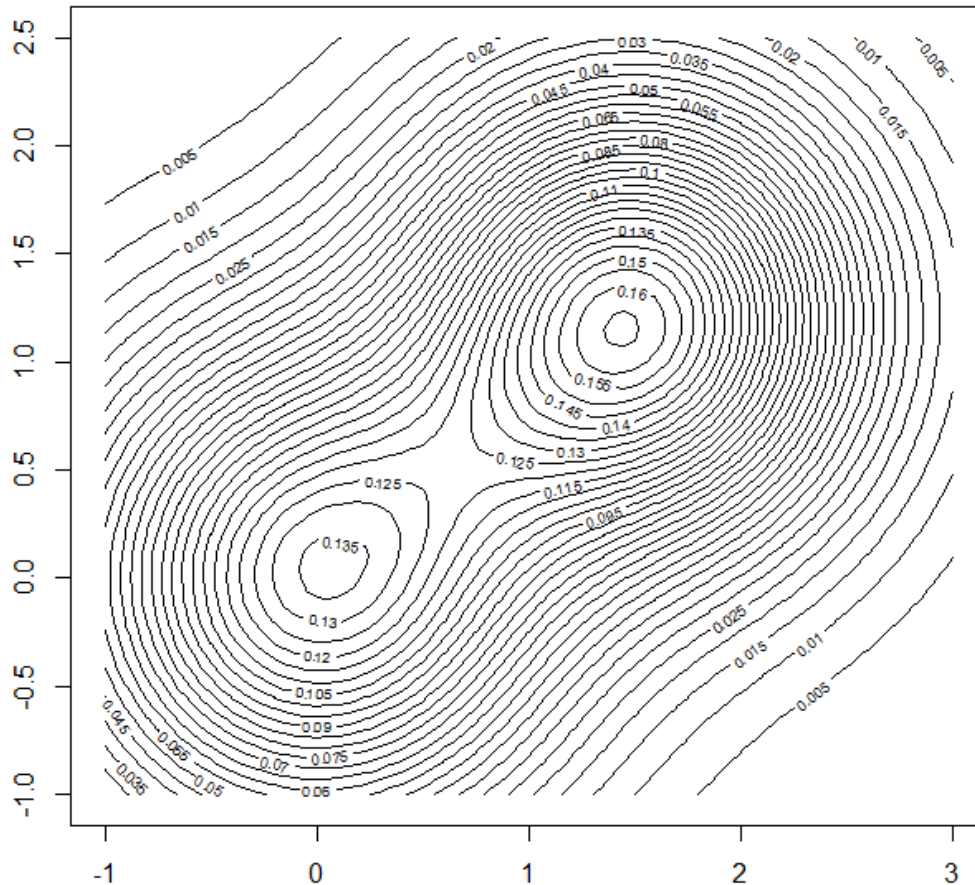
- Compute point  $\mathbf{x} = (x_1, x_2)$  where density  $g(\mathbf{x})$  maximal.
- Do you have a guess?





# Multivariate optimisation – Newton meth.

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



# Multivariate optimisation – Newton meth.

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

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

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

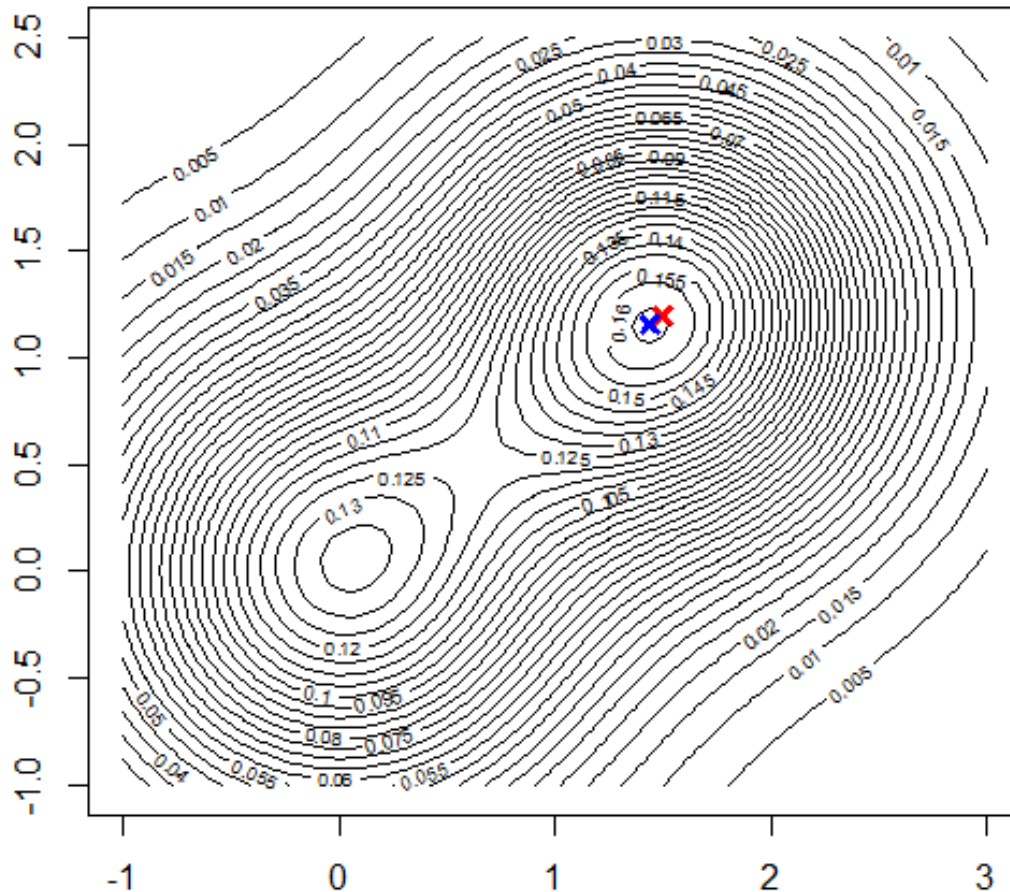
- $\mathbf{g}'(x_1, x_2) = \begin{pmatrix} \frac{\partial g}{\partial x_1}(x_1, x_2) \\ \frac{\partial g}{\partial x_2}(x_1, x_2) \end{pmatrix}$

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



# Multivariate optimisation – Newton meth.

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



- Start with  $x^{(0)} = \begin{pmatrix} 1.5 \\ 1.2 \end{pmatrix}$

- $g'(x^{(0)}) = \begin{pmatrix} -0.0153 \\ -0.0123 \end{pmatrix}$

- $g''(x^{(0)}) = \begin{pmatrix} -0.2902 & 0.0306 \\ 0.0306 & -0.3040 \end{pmatrix}$

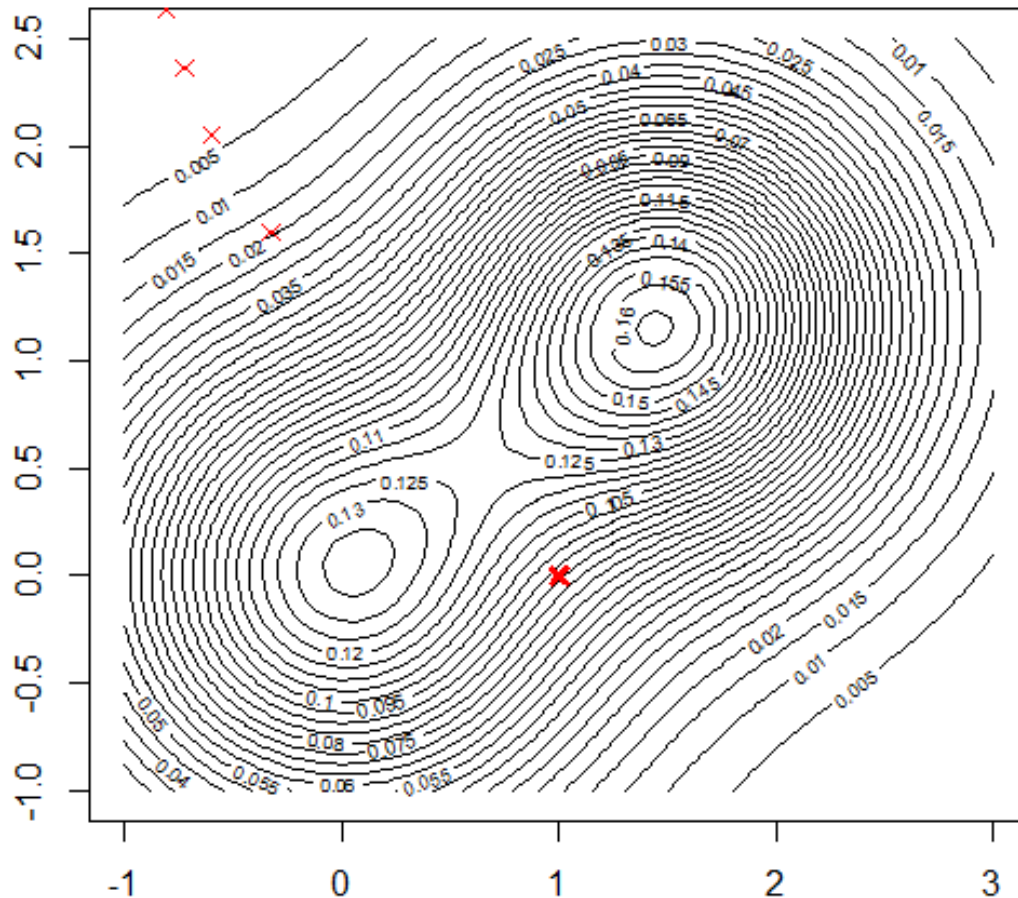
- $(g''(x^{(0)}))^{-1} g'(x^{(0)}) = \begin{pmatrix} 0.058 \\ 0.046 \end{pmatrix}$

- $x^{(1)} = \begin{pmatrix} 1.5 \\ 1.2 \end{pmatrix} - \begin{pmatrix} 0.058 \\ 0.046 \end{pmatrix} = \begin{pmatrix} 1.442 \\ 1.154 \end{pmatrix}$

- $x^{(2)} = x^* = \begin{pmatrix} 1.441 \\ 1.153 \end{pmatrix}$

# Multivariate optimisation – Newton meth.

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



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

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

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

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

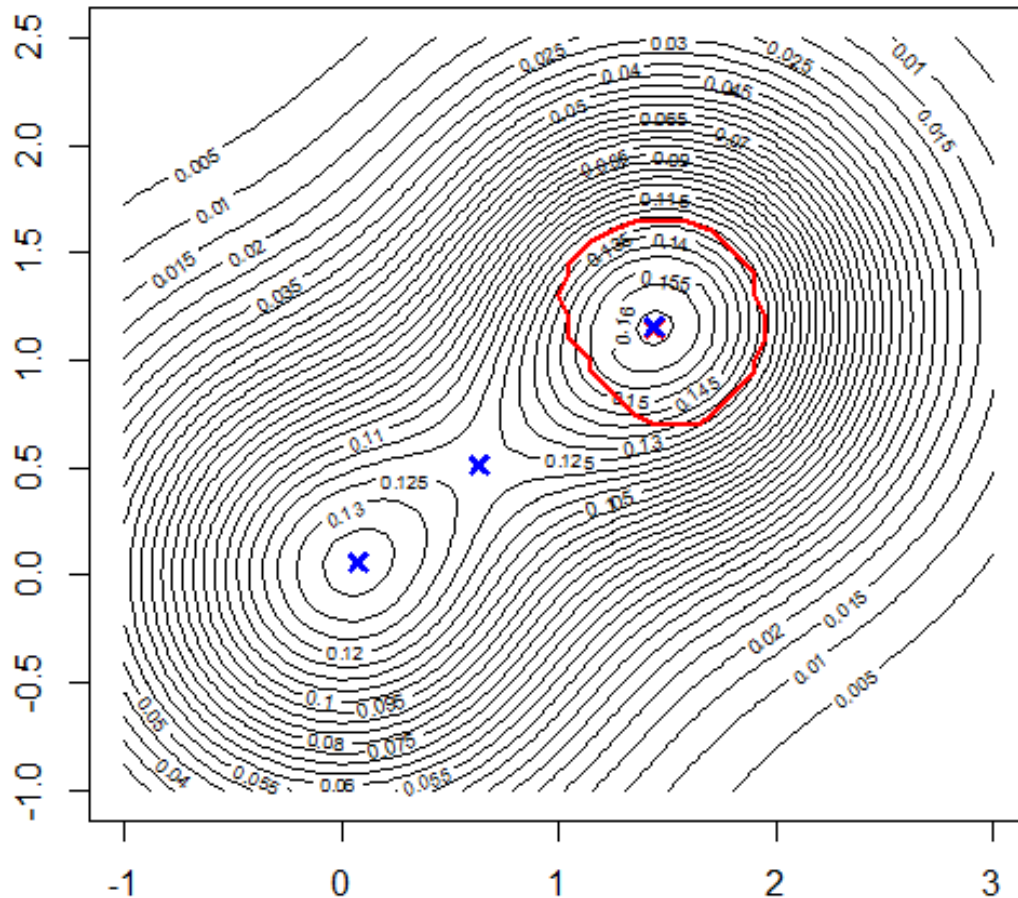
- $x^{(1)} = \begin{pmatrix} 1 \\ 0 \end{pmatrix} - \begin{pmatrix} 1.33 \\ -1.6 \end{pmatrix}$

$$= \begin{pmatrix} -0.33 \\ 1.6 \end{pmatrix}$$



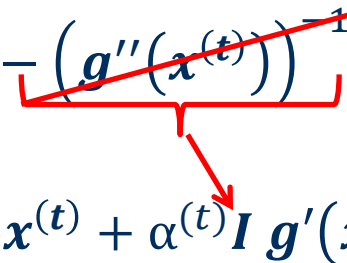
# Multivariate optimisation – Newton meth.

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



# Multivariate optimisation – 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  $\alpha$

# Multivariate optimisation – Steepest ascent method

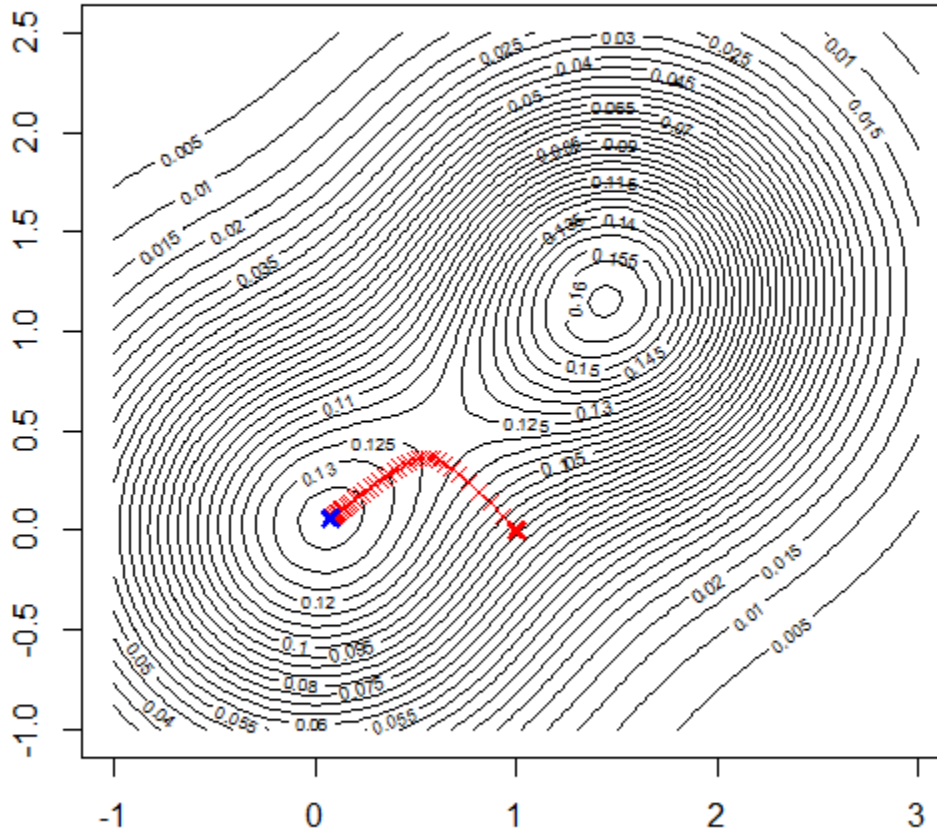
$$\mathbf{x}^{(t+1)} = \mathbf{x}^{(t)} + \alpha^{(t)} \mathbf{I} \mathbf{g}'(\mathbf{x}^{(t)})$$

- Other choices instead  $\mathbf{I}$  in formula above possible
- We know that  $g$  will increase for small  $\alpha$
- 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 (“line search method”)



# Multivariate optimisation – Steepest ascent method

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



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

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

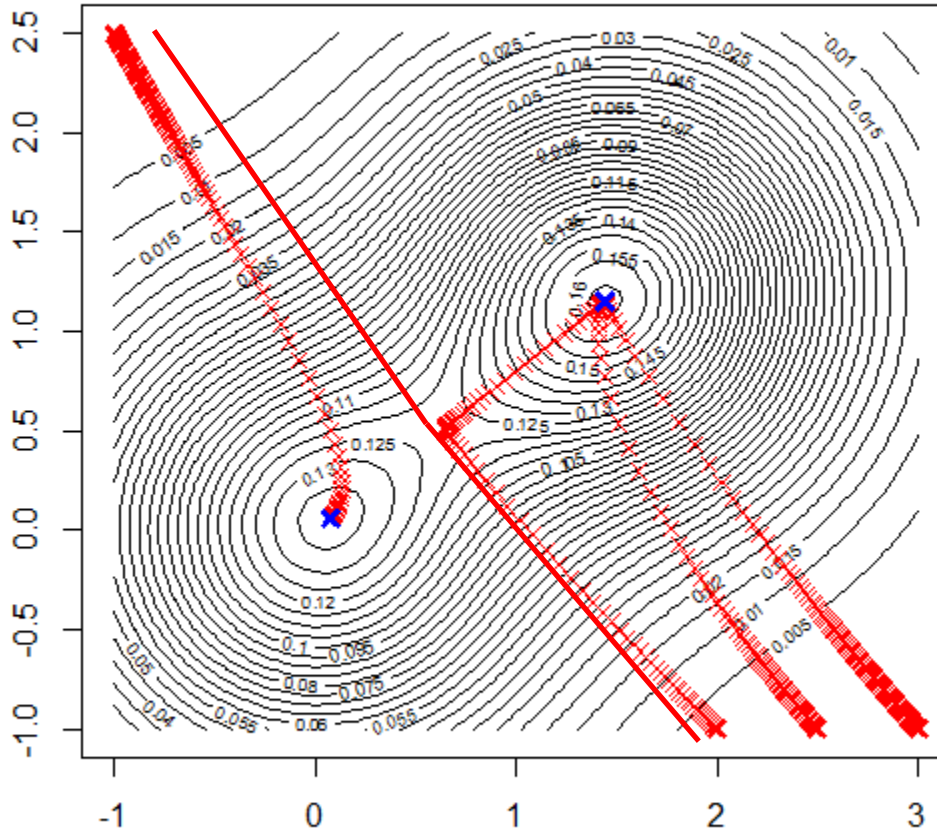
- $$\mathbf{x}^{(1)} = \begin{pmatrix} 1 \\ 0 \end{pmatrix} + \alpha^{(0)} \begin{pmatrix} -0.0667 \\ +0.0705 \end{pmatrix}$$

$$= \begin{pmatrix} 0.9333 \\ 0.0705 \end{pmatrix}$$



# Multivariate optimisation – Steepest ascent method

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



# 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
- 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)



# Convergence rate of optimisation algorithms

- Recall the definition for convergence speed/rate (e.g. bisection is approximately linear, Newton is quadratic)
- To check convergence speed in an optimisation-run, you can calculate

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

(after the algorithm has converged for the result  $\mathbf{x}^*$ ), see GH, page 101/102

- If  $D^{(t)} \rightarrow 1$ , there is not even linear convergence (bad),
- If  $D^{(t)} \rightarrow c \in (0,1)$ , linear convergence,
- If  $D^{(t)} \rightarrow 0$ , better than linear convergence.



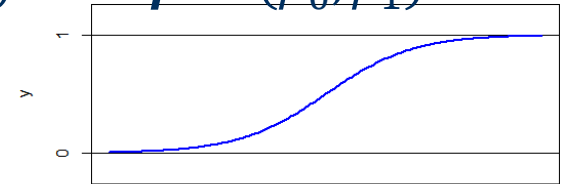
# Maximum likelihood estimator (MLE)

- The MLE is solution of  $g(\hat{\beta}) = \max g(\mathbf{b})$  with  
 $g(\hat{\beta}) = \log\text{-likelihood}(\hat{\beta}, \mathbf{X}, \mathbf{y}) = \sum_{i=1}^n \log\text{-likelihood}(\hat{\beta}, \mathbf{x}_i, y_i)$   
(the latter equation requires independence of observations)
- In the simple case of normally distributed observations,  
MLE=LSE and we have an algebraic solution
- Otherwise, we need usually iterative methods to compute  
the MLE

# MLE for logistic regression

- The MLE is solution of  $g(\hat{\boldsymbol{\beta}}) = \max g(\mathbf{b})$  with  
 $g(\hat{\boldsymbol{\beta}}) = \log\text{-likelihood}(\hat{\boldsymbol{\beta}}, \mathbf{X}, \mathbf{y}) = \sum_{i=1}^n \log\text{-likelihood}(\hat{\boldsymbol{\beta}}, \mathbf{x}_i, y_i)$
- Logistic regression model:  

$$p(\mathbf{x}_i) = P(Y = 1|\mathbf{x}_i) = \frac{1}{1 + \exp(-\mathbf{x}_i^T \boldsymbol{\beta})} = \frac{\exp(\mathbf{x}_i^T \boldsymbol{\beta})}{1 + \exp(\mathbf{x}_i^T \boldsymbol{\beta})}$$
- For simple logistic regression:  $\mathbf{x}_i^T = (1, x_{i1})$  and  $\boldsymbol{\beta} = (\beta_0, \beta_1)^T$
- For multiple logistic regression ( $p$  var.):  
 $\mathbf{x}_i^T = (1, x_{i1}, \dots, x_{ip})$  and  $\boldsymbol{\beta} = (\beta_0, \beta_1, \dots, \beta_p)^T$
- Log likelihood:  $g(\mathbf{b}) = \sum_{i=1}^n \log(\hat{y}_i)y_i + \log(1 - \hat{y}_i)(1 - y_i)$  where  
 $\hat{y}_i = 1/(1 + \exp(-\mathbf{x}_i^T \mathbf{b}))$
- We can write each term in the sum alternatively as follows:  
 $\log(1 - \hat{y}_i) + (\log(\hat{y}_i) - \log(1 - \hat{y}_i))y_i = -\log(1 + \exp(\mathbf{x}_i^T \mathbf{b})) + \mathbf{x}_i^T \mathbf{b} y_i$
- Derivatives with respect to  $b_j, j = 0, \dots, p$ , necessary



# Assignments

- Topic 1: October 2 until October 12
- Topic 2: October 13 until October 22
- Topic 3: October 23 until November 5
- Topic 4: November 6 until November 27
- Final assignment and second chance for Topic 1-4: November 13 until **January 25 (no extension!)**
  
- Some Problems (e.g. 1.3) will be discussed in upcoming lectures and developed further. If you do not solve them until the first deadline, I will provide a different problem as replacement to be done until January 25