



Optimisation algorithms in Statistics I, lecture 4

Frank Miller, Department of Statistics; Stockholm University

November 6, 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

Today's schedule

- Optimisation with constraints
 - Equality constraints
 - Inequality constraints

- Remarks
 - Simulated annealing

Optimisation with equality constraints

- Optimisation problem:
 - x p -dimensional vector, $g: \mathbb{R}^p \rightarrow \mathbb{R}$ function
 - We search x^* with $g(x^*) = \max g(x)$
 - Subject to $h_i(x^*) = 0, i = 1, \dots, m$ (equality constraints)



Optimisation with equality constraints

- Optimisation problem:
 - x p -dimensional vector, $g: \mathbb{R}^p \rightarrow \mathbb{R}$ function
 - We search x^* with $g(x^*) = \max g(x)$
 - Subject to $h_i(x^*) = 0, i = 1, \dots, m$ (equality constraints)
- Approaches:
 - Transformation to an unconstrained problem (problem specific approach)
 - Modification of iterative algorithm to reflect constraints (algorithm specific approach)
 - Lagrange multipliers (general approach)
- $S = \{x \in \mathbb{R}^p \mid h_i(x) = 0, i = 1, \dots, m\}$ called feasible points

Optimisation with equality constraints – transformation

- Example: Cubic regression model for fertilizer-yield-relationship with fertilizer $x \in [0,1.2]$. Experiment planned with
 - proportion w_1 of observations using $x_1 = 0$,
 - proportion w_2 using $x_2 = 0.4$,
 - proportion w_3 using $x_3 = 0.8$,
 - proportion w_4 using $x_4 = 1.2$.
- Note that $w_1 + w_2 + w_3 + w_4 = 1$.
- Information matrix M (proportional to inverse of covariance matrix for $(\hat{\beta}_0, \hat{\beta}_1, \hat{\beta}_2, \hat{\beta}_3)^T$): $M = X^T \text{diag}(w_1, \dots, w_4)X = \sum_{i=1}^4 w_i \mathbf{f}(x_i) \mathbf{f}(x_i)^T$ with $\mathbf{f}(x) = (1, x, x^2, x^3)^T$
- The D-optimal design maximises
$$g(\mathbf{w}) = \det\left(\sum_{i=1}^4 w_i \mathbf{f}(x_i) \mathbf{f}(x_i)^T\right)$$
 subject to
$$h_1(\mathbf{w}) = 1 - \sum_{i=1}^4 w_i = 0$$

Optimisation with equality constraints – transformation

- The D-optimal design maximises $g(\mathbf{w}) = \det(\sum_{i=1}^4 w_i \mathbf{f}(x_i) \mathbf{f}(x_i)^T)$ subject to $h_1(\mathbf{w}) = 1 - \sum_{i=1}^4 w_i = 0$
- Transformation: $1 - \sum_{i=1}^4 w_i = 0 \Rightarrow w_4 = 1 - w_1 - w_2 - w_3$
 $\tilde{g}(w_1, w_2, w_3)$
 $= \det(\sum_{i=1}^3 w_i \mathbf{f}(x_i) \mathbf{f}(x_i)^T + (1 - w_1 - w_2 - w_3) \mathbf{f}(x_4) \mathbf{f}(x_4)^T)$
- The constrained optimisation problem
max. $g(w_1, w_2, w_3, w_4)$ subj. to $h_1(w_1, w_2, w_3, w_4) = 1 - \sum_{i=1}^4 w_i = 0$
is equivalent to the unconstrained optimisation problem
maximise $\tilde{g}(w_1, w_2, w_3)$.
- Solution with **optim**: $(w_1, w_2, w_3) = (\frac{1}{4}, \frac{1}{4}, \frac{1}{4})$, $w_4 = 1 - \frac{3}{4} = \frac{1}{4}$



Optimisation with equality constraints – modification of algorithms

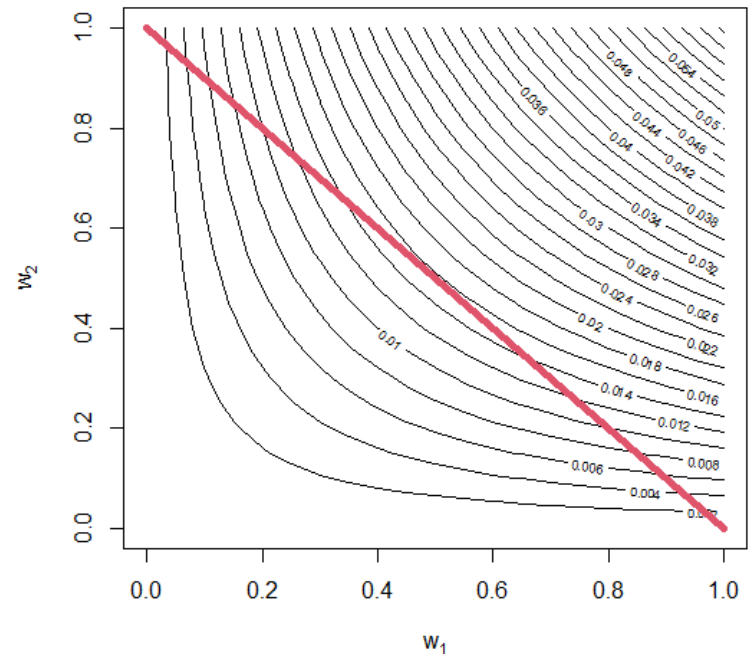
- Constrained optimisation problem:
 - x p -dimensional vector, $g: \mathbb{R}^p \rightarrow \mathbb{R}$ function
 - We search x^* with $g(x^*) = \max g(x)$
 - Subject to $Ax^* - b = \mathbf{0}$, $A \in \mathbb{R}^{m \times p}$, $b \in \mathbb{R}^m$ (**linear** equality constraints)
- Example: Particle Swarm Optimisation (see L3)
- Movement of particle i at iteration $t+1$:
 - $x_i^{(t+1)} = x_i^{(t)} + v_i^{(t+1)}$
 - $v_i^{(t+1)} = wv_i^{(t)} + c_1 R_1^{(t+1)} (p_{\text{best}, i}^{(t)} - x_i^{(t)}) + c_2 R_2^{(t+1)} (g_{\text{best}}^{(t)} - x_i^{(t)})$
- $R_1^{(t+1)}$ and $R_2^{(t+1)}$ are uniformly distributed, `runif()`
- Ensure that $Ax_i^{(0)} = b$ and $Av_i^{(0)} = \mathbf{0}$,
then $Ax_i^{(t)} = b$ for all i and t

Optimisation with equality constraints – Lagrange multipliers

- Example: D-optimal design for quadratic regression without intercept. Experiment planned on $x \in [0,1]$ with
 - prop. w_1 of observations using $x_1 = 0.5$,
 - prop. w_2 using $x_2 = 1$,
 - $w_1 + w_2 = 1$.

- $$g(\mathbf{w}) = \det\left(w_1 \begin{pmatrix} 1 & 1 \\ 4 & 8 \\ 1 & 1 \\ 8 & 16 \end{pmatrix} + w_2 \begin{pmatrix} 1 & 1 \\ 1 & 1 \end{pmatrix}\right)$$

- $$h(\mathbf{w}) = 1 - w_1 - w_2$$



Optimisation with equality constraints – Lagrange multipliers

- $g(\mathbf{w}) = \det\left(w_1 \begin{pmatrix} 1 & 1 \\ 4 & 8 \\ 1 & 1 \\ 8 & 16 \end{pmatrix} + w_2 \begin{pmatrix} 1 & 1 \\ 1 & 1 \end{pmatrix}\right)$

- $h(\mathbf{w}) = 1 - w_1 - w_2$

- $g'(\mathbf{w})$ direction of steepest ascent

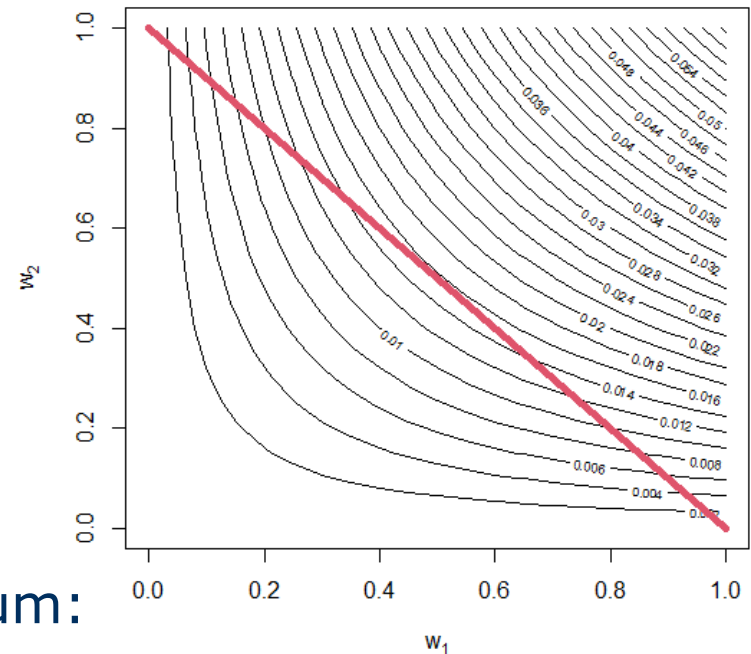
- $h'(\mathbf{w}) = (-1, -1)^\top$

- Condition for constrained maximum:

$$g'(\mathbf{w}) = \lambda h'(\mathbf{w})$$

- $g'(\mathbf{w}) - \lambda h'(\mathbf{w}) = 0$

- Define $\mathcal{L}(\mathbf{x}, \lambda) = g(\mathbf{w}) - \lambda h(\mathbf{w})$



Optimisation with equality constraints – Lagrange multipliers

- Constrained optimisation problem:
 - x p -dimensional vector, $g: \mathbb{R}^p \rightarrow \mathbb{R}$ function
 - We search x^* with $g(x^*) = \max g(x)$
 - Subject to $h_i(x^*) = 0$, $i = 1, \dots, m$ (equality constraints)
- Lagrange:

Let $\mathcal{L}(x, \lambda) = g(x) - \lambda^T \mathbf{h}(x)$, $\mathbf{h}(x) = (h_1(x), \dots, h_m(x))^T$, $\lambda \in \mathbb{R}^m$ and g, h_1, \dots, h_m are continuously differentiable. If g has a local maximum at some point x^* with $\mathbf{h}(x^*) = \mathbf{0}$ (i.e. in the constrained maximisation problem) and at which the gradients of h_1, \dots, h_m are linearly independent, then there exists a λ such that gradient $\mathcal{L}'(x^*, \lambda) = \mathbf{0}$ (i.e. stationary point in the unconstrained problem).



Optimisation with equality constraints – Lagrange multipliers

- Constrained optimisation problem:
 - x p -dimensional vector, $g: \mathbb{R}^p \rightarrow \mathbb{R}$ function
 - We search x^* with $g(x^*) = \max g(x)$
 - Subject to $h_i(x^*) = 0$, $i = 1, \dots, m$ (equality constraints)
- Unconstrained problem:
Search stationary point (x^*, λ) of $\mathcal{L}(x, \lambda) = g(x) - \lambda^T h(x)$.
- Note:
 - $\frac{\partial}{\partial \lambda_i} \mathcal{L}(x^*, \lambda) = 0$ ensures $h_i(x^*) = 0$
 - $\frac{\partial}{\partial x_i} \mathcal{L}(x^*, \lambda) = 0$ ensures that gradient $g'(x^*)$ is orthogonal to the set \mathcal{S} of feasible points at $x = x^*$



Optimisation with inequality constraints

- Constrained optimisation problem:
 - x p -dimensional vector, $g: \mathbb{R}^p \rightarrow \mathbb{R}$ function
 - We search x^* with $g(x^*) = \max g(x)$
 - Subject to $h_i(x^*) = 0, i = 1, \dots, m$
 - and $q_i(x^*) \leq 0, i = 1, \dots, n$ (inequality constraints)
- Set of feasible points
$$\mathcal{S} = \{x \in \mathbb{R}^p \mid h_i(x) = 0, i = 1, \dots, m; q_i(x) \leq 0, i = 1, \dots, n\}$$
- Approaches to handle inequality constraints:
 - Generalisation of Lagrange multipliers (Karush–Kuhn–Tucker approach)
 - penalty method
 - barrier method (also called: interior-point method)



Optimisation with inequality constraints – lasso example

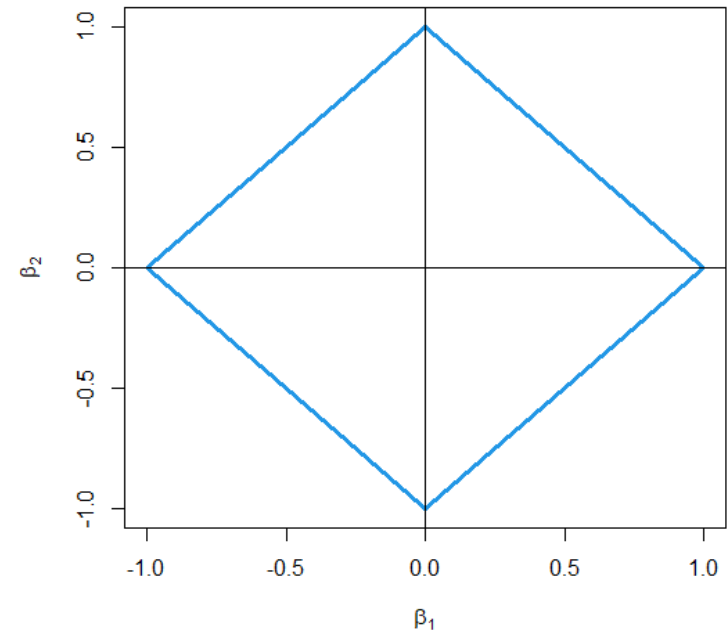
- Lasso's objective function to minimise:

$$g(\hat{\boldsymbol{\beta}}) = \|\mathbf{X}\hat{\boldsymbol{\beta}} - \mathbf{y}\|^2 + \lambda \sum_{i=1}^p |\hat{\beta}_i|$$

- Alternatively, one can solve the constrained problem:

$$\begin{aligned} &\text{minimise: } g(\hat{\boldsymbol{\beta}}) = \|\mathbf{X}\hat{\boldsymbol{\beta}} - \mathbf{y}\|^2 \\ &\text{subject to } \sum_{i=1}^p |\hat{\beta}_i| \leq t \end{aligned}$$

- For $p=2$ and $t=1$, the set of feasible points $\mathcal{S} = \{\hat{\boldsymbol{\beta}} \in \mathbb{R}^p \mid \sum_{i=1}^p |\hat{\beta}_i| \leq t\}$ is inside of the blue area



Optimisation with inequality constraints

– Karush–Kuhn–Tucker approach

- Constrained optimisation problem:
 - x p -dimensional vector, $g: \mathbb{R}^p \rightarrow \mathbb{R}$ function
 - We search x^* with $g(x^*) = \max g(x)$
 - Subject to $h_i(x^*) = 0, i = 1, \dots, m$
 - and $q_i(x^*) \leq 0, i = 1, \dots, n$ (inequality constraints)
- Karush–Kuhn–Tucker (KKT) approach uses generalised Lagrangian $\mathcal{L}(x, \lambda, \mu) = g(x) - \lambda^T \mathbf{h}(x) - \mu^T \mathbf{q}(x)$ with $\mathbf{h}(x) = (h_1(x), \dots, h_m(x))^T, \lambda \in \mathbb{R}^m, \mathbf{q}(x) = (q_1(x), \dots, q_n(x))^T, \mu \in \mathbb{R}^n$
- Instead of above constrained optimisation, search stationary point $(x^*, \lambda, \mu \geq \mathbf{0})$ of $\mathcal{L}(x, \lambda, \mu) = g(x) - \lambda^T \mathbf{h}(x) - \mu^T \mathbf{q}(x)$.
For x^* being a solution of the constrained problem, following condition required: “for all $i=1, \dots, n: q_i(x^*) = 0$ or $\mu_i = 0$ ”

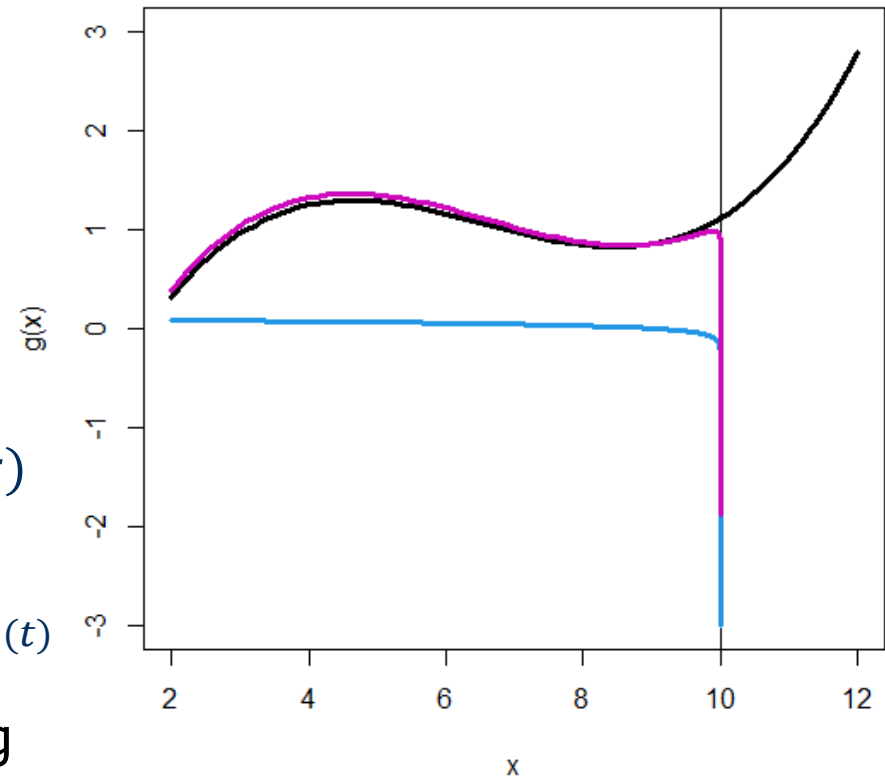
Optimisation with inequality constraints – penalty and barrier methods

- Constrained optimisation problem:
 - \mathbf{x} p -dimensional vector, $g: \mathbb{R}^p \rightarrow \mathbb{R}$ function
 - We search \mathbf{x}^* with $g(\mathbf{x}^*) = \max g(\mathbf{x})$
 - Subject to $q_i(\mathbf{x}^*) \geq 0, i = 1, \dots, n$ (inequality constraints)
- Idea: Modify g to \tilde{g} such that the algorithm finds only local maxima which fulfil $q_i(\mathbf{x}^*) \geq 0, i = 1, \dots, n$, even if optimisation done unconstrained
- Penalty methods: Set $\tilde{g} = g$ on $\mathcal{S} = \{\mathbf{x} | q_i(\mathbf{x}) \geq 0, i = 1, \dots, n\}$ and add a (negative) penalty if $q_i(\mathbf{x}) < 0$ for some i
- Barrier methods: Set $\tilde{g} = -\infty$ if $q_i(\mathbf{x}) < 0$ for some i and g is modified on $\mathcal{S} = \{\mathbf{x} | q_i(\mathbf{x}) \geq 0, i = 1, \dots, n\}$

Optimisation with inequality constraints

– Barrier method (interior-point method)

- Example: maximise $g(x)$ on range $x \leq 10$
- Add barrier function $\mu^{(t)}b(x)$
- $\tilde{g}(x) = g(x) + \mu^{(t)}b(x)$ should be small close to 10, $x < 10$, and $-\infty$ for $x > 10$
- Log barrier: $b(x) = \log(10 - x)$
- Solve maximisation for $\tilde{g}(x)$
- Adapt barrier with smaller $\mu^{(t)}$
- If $\mu^{(t)} \rightarrow 0$, local maxima of g can be detected, both at the boundary and in the interior



$$\mu^{(t)} = 0.04$$



Optimisation with linear inequality constraints – R-function `constrOptim`

- Constrained optimisation problem:
 - x p -dimensional vector, $g: \mathbb{R}^p \rightarrow \mathbb{R}$ function
 - We search x^* with $g(x^*) = \max g(x)$
 - Subject to $Ux^* - c \geq \mathbf{0}$, $U \in \mathbb{R}^{n \times p}$, $c \in \mathbb{R}^n$ (**linear** inequality constraints; rows of U are u_i^T)
- The R-function `constrOptim` uses log barrier functions
- `constrOptim` calls repeatedly `optim` for function \tilde{g} with barrier; barrier adapted between iterations: $\mu^{(t)}$ decreases
- E.g: $\tilde{g}(x) = g(x) + \mu^{(t)} \sum_{i=1}^n \log(u_i^T x - c_i)$ (for maximisation; $g(x) - \mu^{(t)} \dots$ for minimisation)



Optimisation with linear inequality constraints – barrier method

- Example: Quadratic regression for fertilizer-yield-relationship with fertilizer $x \in [0,1.2]$. Experiment planned with
 - proportion w_i of observations using $x_i \in [0,1.2]$ (can be chosen by experimenter), $i=1,2,3$; $w_3 = 1 - w_1 - w_2$.
- Parameters to be optimised: $\mathbf{y} = (x_1, x_2, x_3, w_1, w_2)^T$
- D-optimal design maximises $g(\mathbf{y}) = \det(\sum_{i=1}^3 w_i \mathbf{f}(x_i) \mathbf{f}(x_i)^T)$ subject to
$$x_i \geq 0, 1.2 - x_i \geq 0, i=1,2,3, w_1 \geq 0, w_2 \geq 0, 1 - w_1 - w_2 \geq 0$$
- Construct \mathbf{U} and \mathbf{c} such that constraints can be written as
$$\mathbf{U}\mathbf{y} - \mathbf{c} \geq \mathbf{0}$$

Optimisation with linear inequality constraints – barrier method

- $\mathbf{y} = (x_1, x_2, x_3, w_1, w_2)^T$, $w_3 = 1 - w_1 - w_2$
- D-optimal design maximises $g(\mathbf{y}) = \det(\sum_{i=1}^3 w_i \mathbf{f}(x_i) \mathbf{f}(x_i)^T)$
subject to $x_i \geq 0$, $1.2 - x_i \geq 0$, $w_1 \geq 0$, $w_2 \geq 0$, $1 - w_1 - w_2 \geq 0$
- $U\mathbf{y} - \mathbf{c} \geq \mathbf{0}$ with

$$U = \begin{pmatrix} 1 & 0 & 0 & 0 & 0 \\ -1 & 0 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 & 0 \\ 0 & -1 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 & 0 \\ 0 & 0 & -1 & 0 & 0 \\ 0 & 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 0 & 1 \\ 0 & 0 & 0 & -1 & -1 \end{pmatrix}, \mathbf{c} = \begin{pmatrix} 0 \\ -1.2 \\ 0 \\ -1.2 \\ 0 \\ -1.2 \\ 0 \\ 0 \\ -1 \end{pmatrix}$$

Optimisation with linear inequality constraints – R-function `constrOptim`

- R-code:

```
U      <- matrix(0, nrow=9, ncol=5)
U[1,1] <- U[3,2] <- U[5,3] <- U[7,4] <- U[8,5] <- 1
U[2,1] <- U[4,2] <- U[6,3] <- U[9,4] <- U[9,5] <- -1
d      <- c(rep(c(0, -1.2), 3), 0, 0, -1)
startv <- c(0.2, 0.3, 0.4, 0.2, 0.2)
# Nelder-Mead as inner optimisation method:
res     <- constrOptim(startv, f=g, grad=NULL, ui=U, ci=d,
                       control=list(fnscale=-1))
round(res$par, 3)
```
- Result: 0.000 0.597 1.200 0.331 0.333
- Note: In this case, the solution is algebraically known based on optimal design theory

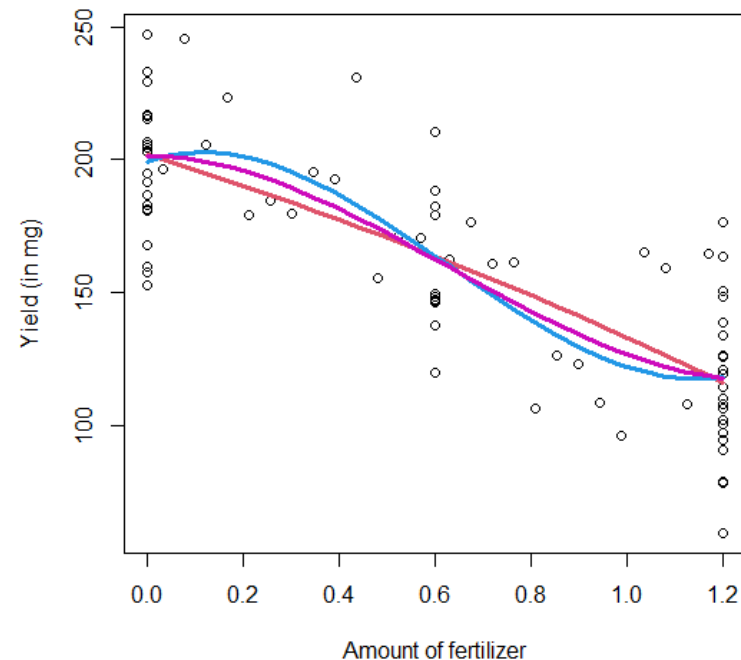
Optimisation with linear inequality constraints – barrier method

- Limitations of barrier method (Lange, 2010, page 301):
 - Iterations within iterations necessary
 - No obvious choice how fast $\mu^{(t)}$ should go to 0
 - A too small value $\mu^{(t)}$ can lead to numerical instability

Other topics

About the cress experiment – Problem 3.3

- Design chosen has some optimality property assuming a quadratic regression with some robustness if other model valid
- Regressions: **quadratic**; **cubic**; **cubic without linear term**



Gradient free optimisation – Simulated annealing

- Start value $x^{(0)}$; Stage $j=0,1,2,\dots$ has m_j iterations; set $j=0$
- Given iteration $x^{(t)}$, generate $x^{(t+1)}$ as follows:
 1. Sample a candidate x^* from a proposal distribution $p(\cdot|x^{(t)})$
 2. Compute $h(x^{(t)}, x^*) = \exp\left(\frac{g(x^*) - g(x^{(t)})}{\tau_j}\right)$ ← $g(x^{(t)}) - g(x^*)$ for minimisation
 3. Define next iteration $x^{(t+1)}$ according to
$$x^{(t+1)} = \begin{cases} x^*, & \text{with probability } \min\{h(x^{(t)}, x^*), 1\} \\ x^{(t)}, & \text{otherwise} \end{cases}$$
 4. Set $t \leftarrow t+1$ and repeat 1.-3. m_j times
 5. Update $\tau_j = \alpha(\tau_{j-1})$ and $m_j = \beta(m_{j-1})$; set $j \leftarrow j+1$; go to 1

τ_j is temperature; function α should slowly decrease it

Function β should be increasing

Markov Chain Monte Carlo (MCMC) – Metropolis-Hastings algorithm

- Metropolis-Hastings (MH) algorithm:
- A starting value $x^{(0)}$ is generated from some starting distribution
- Given observation $x^{(t)}$, generate $x^{(t+1)}$ as follows:
 1. Sample a candidate x^* from a proposal distribution $g(\cdot | x^{(t)})$
 2. Compute the MH ratio $R(x^{(t)}, x^*) = \frac{f(x^*) g(x^{(t)} | x^*)}{f(x^{(t)}) g(x^* | x^{(t)})}$
 3. Sample $x^{(t+1)}$ according to
$$x^{(t+1)} = \begin{cases} x^*, & \text{with probability } \min\{R(x^{(t)}, x^*), 1\} \\ x^{(t)}, & \text{otherwise} \end{cases}$$
 4. If more observations needed, set $t \leftarrow t+1$; go to 1



Remarks