

These are the slides of the lecture

Pattern Recognition
Winter term 2020/21
Friedrich-Alexander University of Erlangen-Nuremberg.

These slides are released under Creative Commons License Attribution CC BY 4.0.

Please feel free to reuse any of the figures and slides, as long as you keep a reference to the source of these slides at <https://lme.tf.fau.de/teaching/> acknowledging the authors Niemann, Hornegger, Hahn, Steidl, Nöth, Seitz, Rodriguez, Das and Maier.

Erlangen, January 8, 2021
Prof. Dr.-Ing. Andreas Maier

Pattern Recognition (PR)

Prof. Dr.-Ing. Andreas Maier

Pattern Recognition Lab (CS 5), Friedrich-Alexander-Universität Erlangen-Nürnberg

Winter Term 2020/21



Optimization



Motivation

- Optimization is crucial for many solutions in pattern recognition, pattern analysis, machine learning, artificial intelligence, etc.
- Optimization has many faces:
 - discrete optimization,
 - combinatorial optimization,
 - genetic algorithms,
 - gradient descent,
 - unconstrained and constrained optimization,
 - linear programming,
 - convex optimization, etc.
- There is no lecture on pattern recognition without a refresher course on optimization techniques.
- Each researcher has his own favorite optimization algorithm.

Convexity

Definition

A function $f : \mathbb{R}^d \rightarrow \mathbb{R}$ is *convex* if the domain $\text{dom}(f)$ of f is a convex set and if $\forall \mathbf{x}, \mathbf{y} \in \text{dom}(f)$, and θ with $0 \leq \theta \leq 1$, we have

$$f(\theta \mathbf{x} + (1 - \theta) \mathbf{y}) \leq \theta f(\mathbf{x}) + (1 - \theta) f(\mathbf{y})$$

Convexity

Definition

A function $f : \mathbb{R}^d \rightarrow \mathbb{R}$ is *convex* if the domain $\text{dom}(f)$ of f is a convex set and if $\forall \mathbf{x}, \mathbf{y} \in \text{dom}(f)$, and θ with $0 \leq \theta \leq 1$, we have

$$f(\theta \mathbf{x} + (1 - \theta) \mathbf{y}) \leq \theta f(\mathbf{x}) + (1 - \theta) f(\mathbf{y})$$

A function $f : \mathbb{R}^d \rightarrow \mathbb{R}$ is *concave* if $-f$ is convex.

Convexity

Definition

A function $f : \mathbb{R}^d \rightarrow \mathbb{R}$ is *convex* if the domain $\text{dom}(f)$ of f is a convex set and if $\forall \mathbf{x}, \mathbf{y} \in \text{dom}(f)$, and θ with $0 \leq \theta \leq 1$, we have

$$f(\theta \mathbf{x} + (1 - \theta) \mathbf{y}) \leq \theta f(\mathbf{x}) + (1 - \theta) f(\mathbf{y})$$

A function $f : \mathbb{R}^d \rightarrow \mathbb{R}$ is *concave* if $-f$ is convex.

Geometric interpretation:

The line segment between $(\mathbf{x}, f(\mathbf{x}))$ and $(\mathbf{y}, f(\mathbf{y}))$ lies above the graph of f .

Unconstrained Optimization

Let us assume in the following that we have to compute the minimum of a convex function

$$f : \mathbb{R}^d \rightarrow \mathbb{R}$$

that is twice differentiable.

Unconstrained Optimization

Let us assume in the following that we have to compute the minimum of a convex function

$$f : \mathbb{R}^d \rightarrow \mathbb{R}$$

that is twice differentiable.

The unconstrained optimization problem is just the solution of the minimization problem

$$\mathbf{x}^* = \underset{\mathbf{x}}{\operatorname{argmin}} f(\mathbf{x})$$

where \mathbf{x}^* denotes the optimal point.

Unconstrained Optimization (cont.)

For this particular family of functions, a necessary and sufficient condition for the minimum are the zero-crossings of the function's gradient:

$$\nabla f(\mathbf{x}^*) = 0.$$

Unconstrained Optimization (cont.)

Most methods follow an **iterative scheme**:

| | |
|----------------|--|
| initialization | $\mathbf{x}^{(0)}$ |
| iteration step | $\mathbf{x}^{(k+1)} = g(\mathbf{x}^{(k)})$ |

where $g : \mathbb{R}^d \rightarrow \mathbb{R}^d$ is the update function.

Unconstrained Optimization (cont.)

Most methods follow an **iterative scheme**:

$$\begin{array}{ll} \text{initialization} & \mathbf{x}^{(0)} \\ \text{iteration step} & \mathbf{x}^{(k+1)} = g(\mathbf{x}^{(k)}) \end{array}$$

where $g : \mathbb{R}^d \rightarrow \mathbb{R}^d$ is the update function.

The iterations **terminate**, if

$$\|\mathbf{x}^{(k+1)} - \mathbf{x}^{(k)}\| < \varepsilon,$$

i. e. no further significant change.

Descent Methods

We now consider iteration schemes that produce a sequence of estimates according to the update function

$$\mathbf{x}^{(k+1)} = g(\mathbf{x}^{(k)}) = \mathbf{x}^{(k)} + t^{(k)} \Delta \mathbf{x}^{(k)}.$$

where

$\Delta \mathbf{x}^{(k)} \in \mathbb{R}^d$: is the **search direction** in the k -th iteration
 $t^{(k)} \in \mathbb{R}$: denotes the **step length** in the k -th iteration

Descent Methods

We now consider iteration schemes that produce a sequence of estimates according to the update function

$$\mathbf{x}^{(k+1)} = g(\mathbf{x}^{(k)}) = \mathbf{x}^{(k)} + t^{(k)} \Delta \mathbf{x}^{(k)}.$$

where

$\Delta \mathbf{x}^{(k)} \in \mathbb{R}^d$: is the **search direction** in the k -th iteration
 $t^{(k)} \in \mathbb{R}$: denotes the **step length** in the k -th iteration

and where we expect

$$f(\mathbf{x}^{(k+1)}) < f(\mathbf{x}^{(k)}), \quad \text{i.e. } \nabla f(\mathbf{x}^{(k)})^T \Delta \mathbf{x}^{(k)} < 0$$

except $\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} = \mathbf{x}^*$.

Taylor Approximation

For many problems it is always good to know the **second order Taylor approximation**:

$$f(\mathbf{x} + t \cdot \Delta \mathbf{x}) \approx f(\mathbf{x}) + t \cdot \nabla f(\mathbf{x})^T \Delta \mathbf{x} + \frac{1}{2} t^2 \cdot \Delta \mathbf{x}^T \nabla^2 f(\mathbf{x}) \Delta \mathbf{x}$$

Descent Methods (cont.)

Input: function f , initial estimate $\mathbf{x}^{(0)}$

Initialize: $k := 0$

repeat

 Select (or compute) descent direction

 Line search (1-D optimization):

$$t^{(k)} = \underset{t \geq 0}{\operatorname{argmin}} f(\mathbf{x}^{(k)} + t \cdot \Delta \mathbf{x}^{(k)})$$

Update:

$$\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} + t^{(k)} \Delta \mathbf{x}^{(k)}.$$

$k := k + 1$

until $\|\mathbf{x}^{(k)} - \mathbf{x}^{(k-1)}\| < \varepsilon$

Output: $\mathbf{x}^{(k)}$

Line Search Methods

- Multivariate optimization in its described form requires a proper line search method.
- Exact line search along the straight line $\{\mathbf{x} + t\Delta\mathbf{x} \mid t \geq 0\}$ has to solve

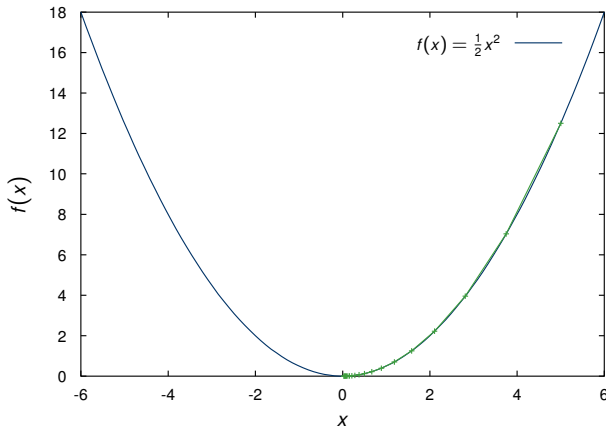
$$t^* = \operatorname{argmin}_{t \geq 0} f(\mathbf{x} + t\Delta\mathbf{x})$$

and is rarely used.

- An overview of methods can be found in numerical recipes.

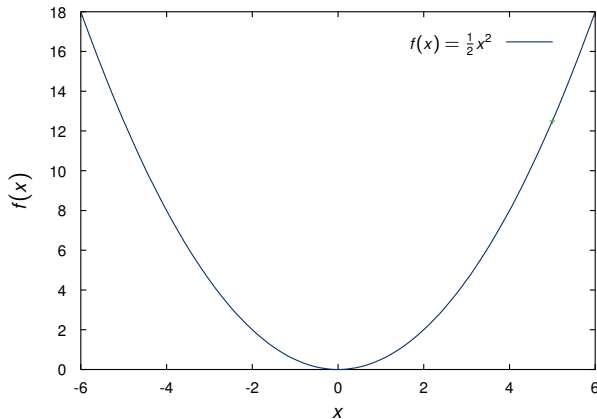
Line Search Methods (cont.)

Setting $t = 0.25$:



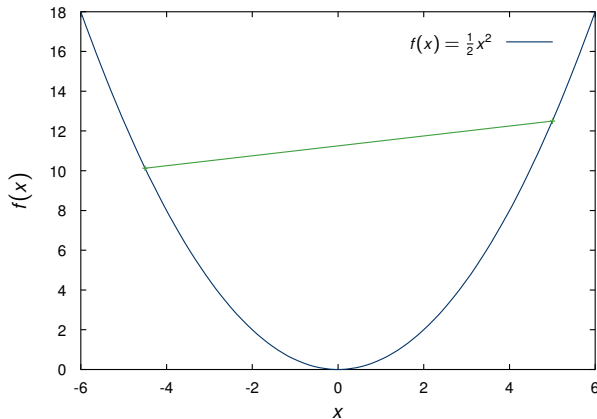
Line Search Methods (cont.)

Setting $t = 1.9$:



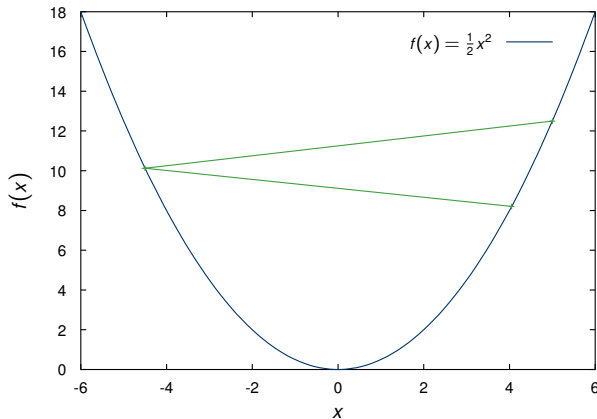
Line Search Methods (cont.)

Setting $t = 1.9$:



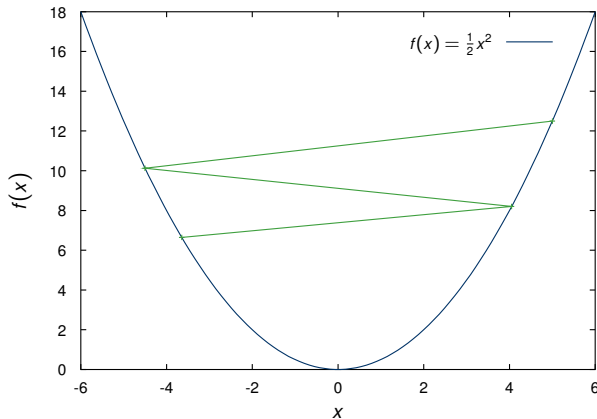
Line Search Methods (cont.)

Setting $t = 1.9$:



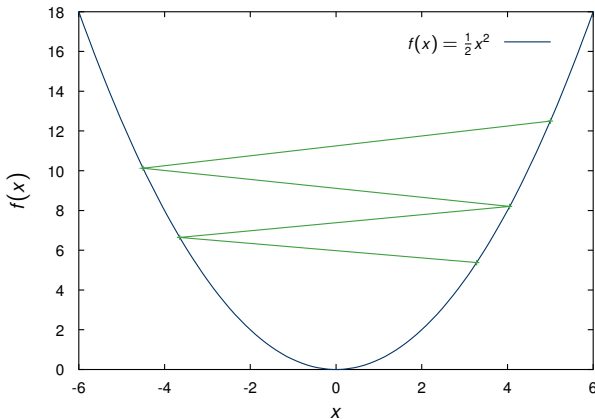
Line Search Methods (cont.)

Setting $t = 1.9$:



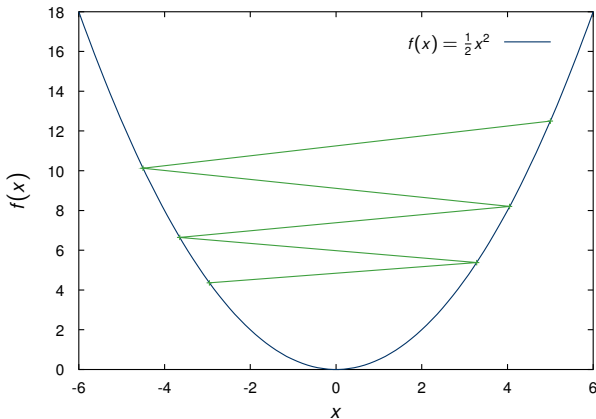
Line Search Methods (cont.)

Setting $t = 1.9$:



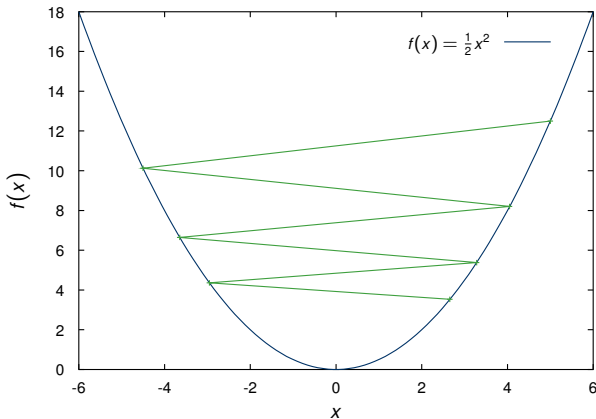
Line Search Methods (cont.)

Setting $t = 1.9$:



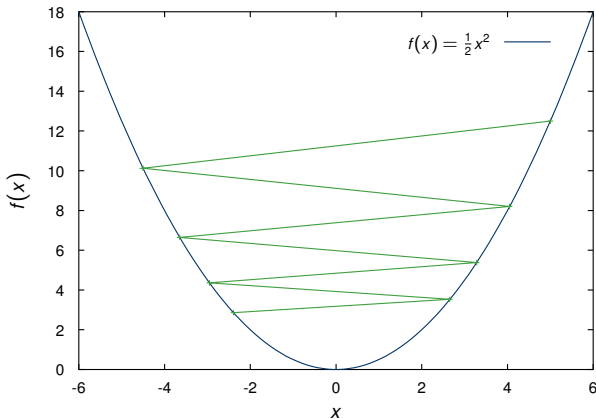
Line Search Methods (cont.)

Setting $t = 1.9$:



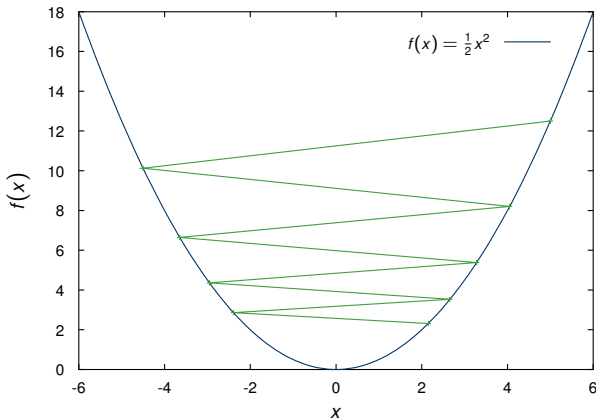
Line Search Methods (cont.)

Setting $t = 1.9$:



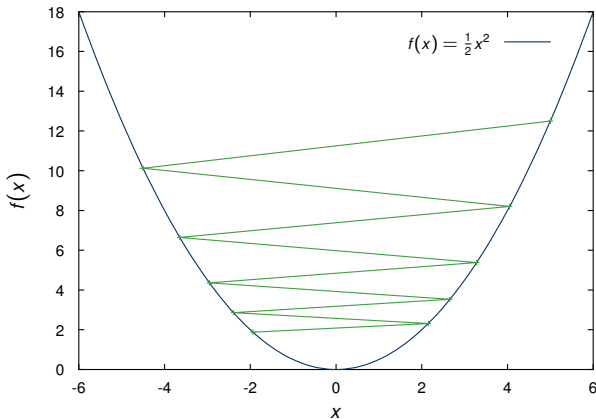
Line Search Methods (cont.)

Setting $t = 1.9$:



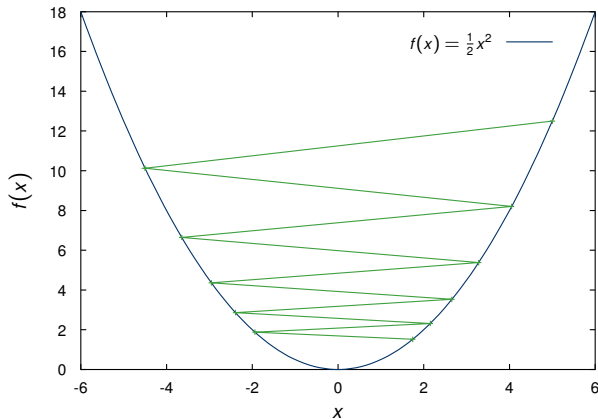
Line Search Methods (cont.)

Setting $t = 1.9$:



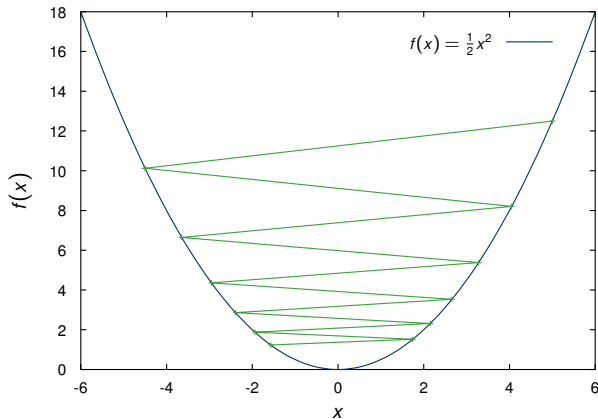
Line Search Methods (cont.)

Setting $t = 1.9$:



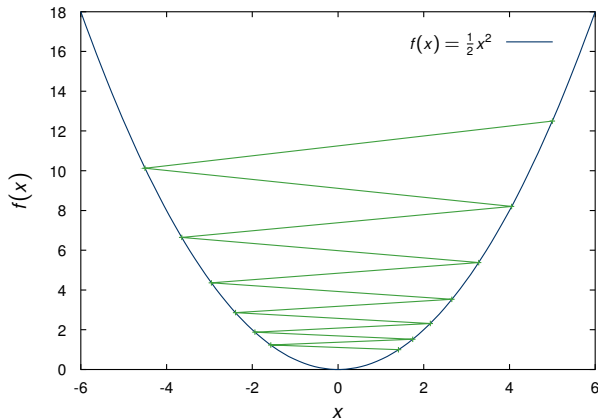
Line Search Methods (cont.)

Setting $t = 1.9$:



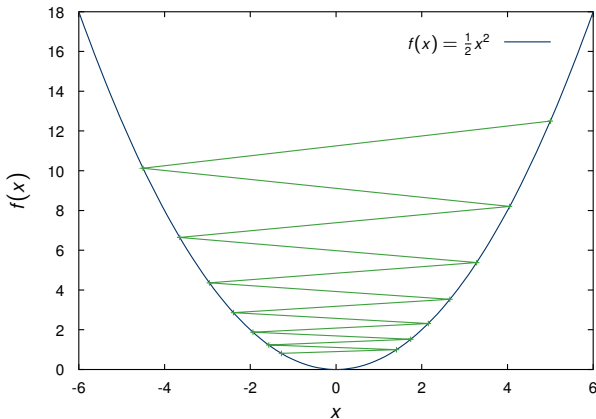
Line Search Methods (cont.)

Setting $t = 1.9$:



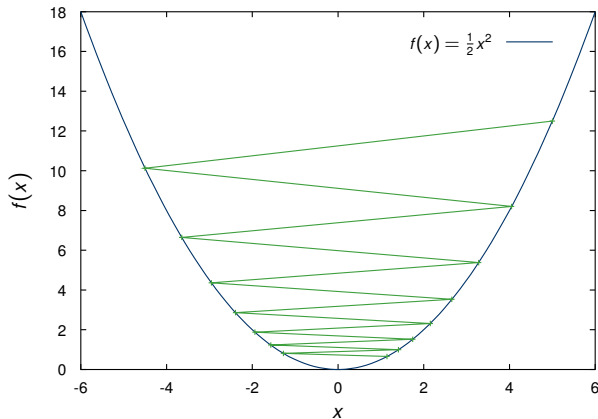
Line Search Methods (cont.)

Setting $t = 1.9$:



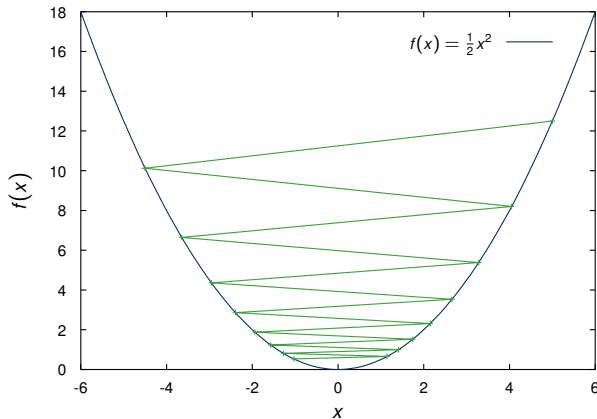
Line Search Methods (cont.)

Setting $t = 1.9$:



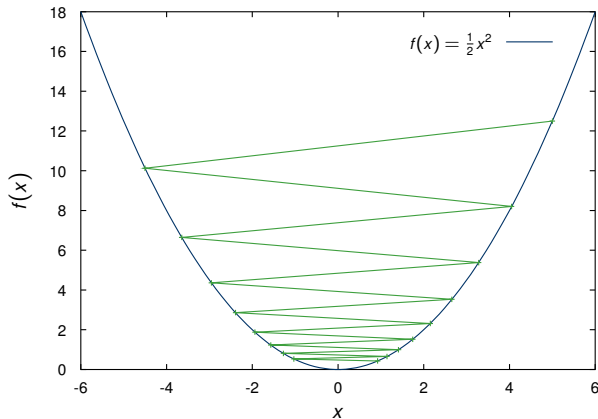
Line Search Methods (cont.)

Setting $t = 1.9$:



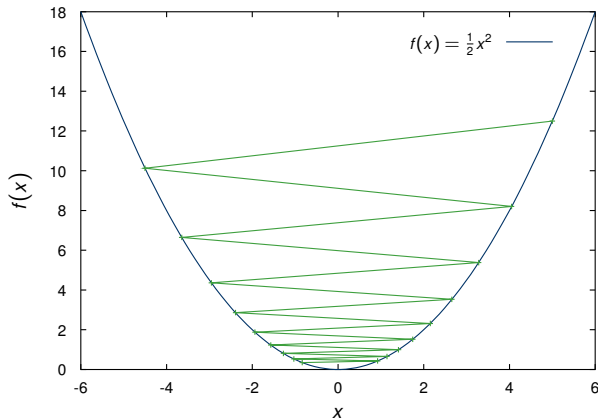
Line Search Methods (cont.)

Setting $t = 1.9$:



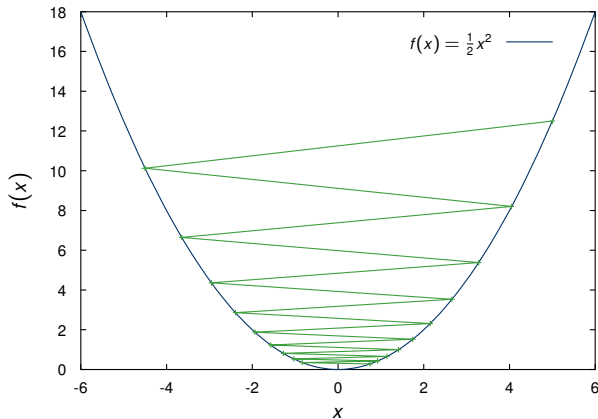
Line Search Methods (cont.)

Setting $t = 1.9$:



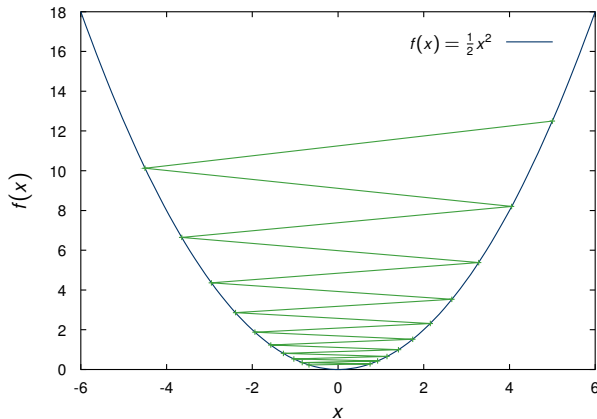
Line Search Methods (cont.)

Setting $t = 1.9$:



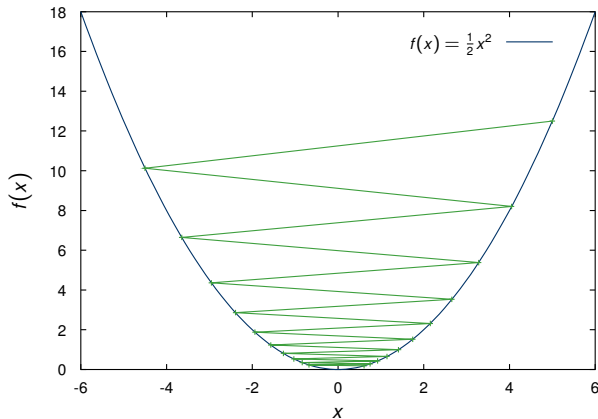
Line Search Methods (cont.)

Setting $t = 1.9$:



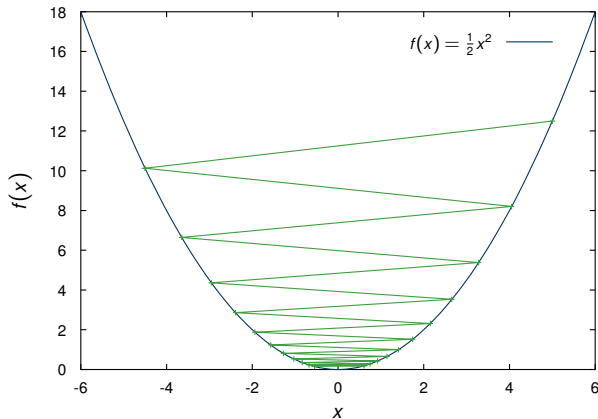
Line Search Methods (cont.)

Setting $t = 1.9$:



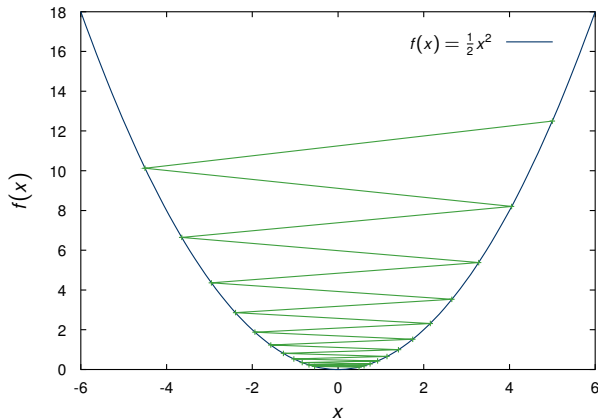
Line Search Methods (cont.)

Setting $t = 1.9$:



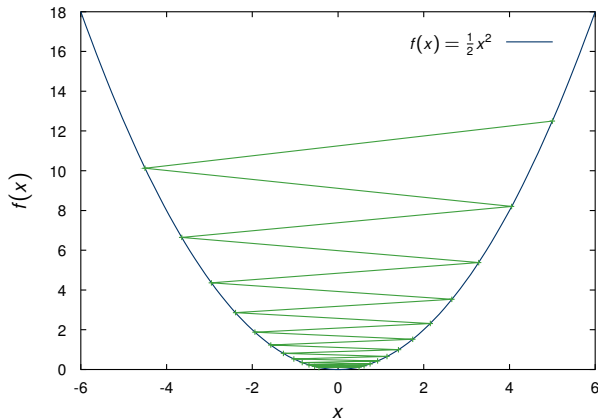
Line Search Methods (cont.)

Setting $t = 1.9$:



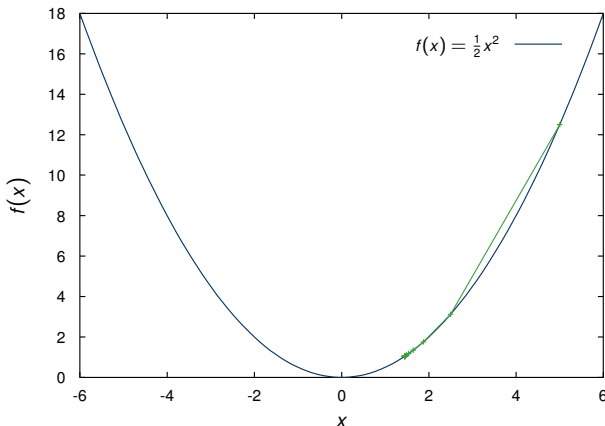
Line Search Methods (cont.)

Setting $t = 1.9$:

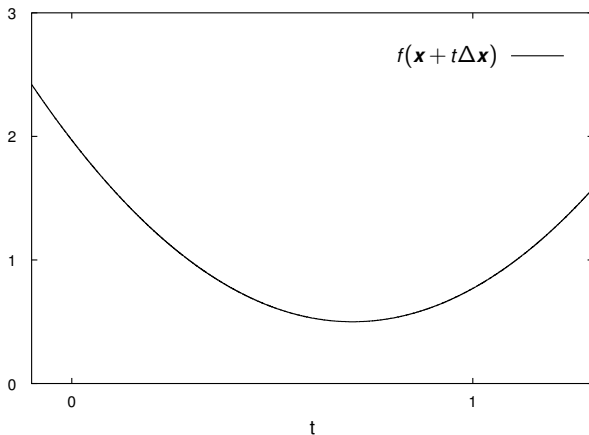


Line Search Methods (cont.)

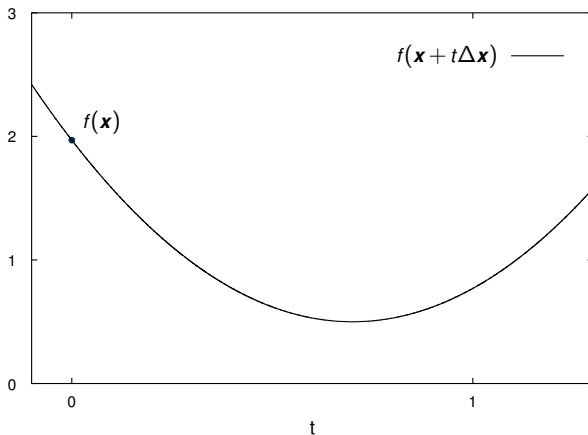
Setting $t^{(k+1)} = \frac{1}{2}t^{(k)}$ and starting with $t^{(0)} = 0.5$:



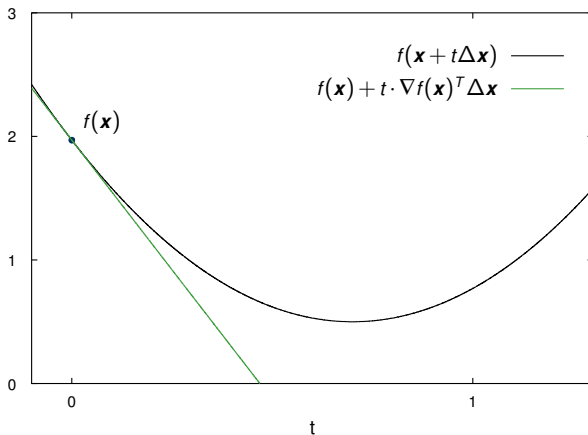
Backtracking Line Search



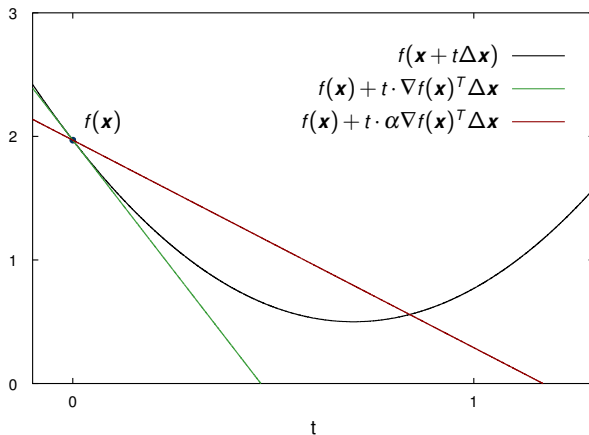
Backtracking Line Search



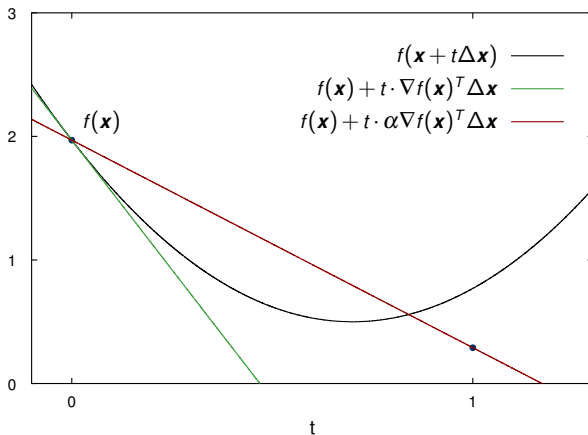
Backtracking Line Search



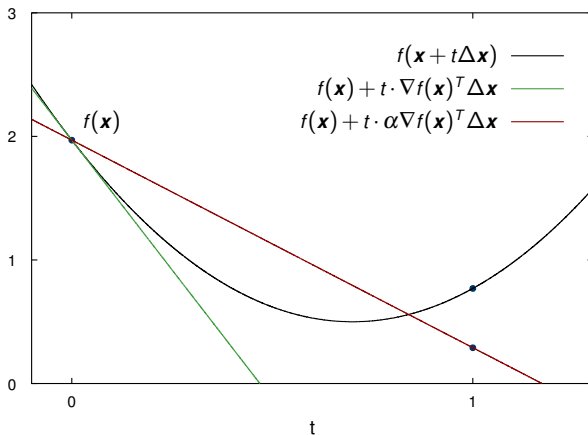
Backtracking Line Search



Backtracking Line Search



Backtracking Line Search



Backtracking Line Search (cont.)

The Armijo-Goldstein line search algorithm:

Input: function f , search direction $\Delta \mathbf{x}$

Initialize: $t := 1$

Select: $\alpha \in [0, 0.5]$ and $\beta \in [0, 1]$.

while $f(\mathbf{x} + t\Delta \mathbf{x}) > f(\mathbf{x}) + \alpha t \cdot \nabla f(\mathbf{x})^T \Delta \mathbf{x}$ **do**

$t := \beta t$

end while

Output: t

Gradient Descent Methods

A natural choice of the search direction is the **negative gradient**:

$$\Delta \mathbf{x}^{(k)} = -\nabla f(\mathbf{x}^{(k)})$$

Rule of thumb:

The negative gradient is the steepest descent direction.

Gradient Descent Methods (cont.)

Input: function f , initial estimate $\mathbf{x}^{(0)}$

initialize: $k := 0$

repeat

Set descent direction: $\Delta \mathbf{x}^{(k)} = -\nabla f(\mathbf{x}^{(k)})$

Line search (1-D optimization):

$$t^{(k)} = \underset{t \geq 0}{\operatorname{argmin}} f(\mathbf{x}^{(k)} + t \cdot \Delta \mathbf{x}^{(k)})$$

Update:

$$\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} + t^{(k)} \Delta \mathbf{x}^{(k)}.$$

$k := k + 1$

until $\|\mathbf{x}^{(k)} - \mathbf{x}^{(k-1)}\|_2 < \varepsilon$

Output: $\mathbf{x}^{(k)}$



**Pattern
Recognition
Lab**



**FRIEDRICH-ALEXANDER
UNIVERSITÄT
ERLANGEN-NÜRNBERG**

TECHNISCHE FAKULTÄT

Next Time in

Pattern Recognition



Steepest Descent Methods

(Normalized) steepest descent, what does it mean?

We search for the unit vector that shows the largest decrease in the linear approximation of f :

$$\Delta \mathbf{x} = \underset{\mathbf{u}}{\operatorname{argmin}} \{ \nabla f(\mathbf{x})^T \mathbf{u}; \|\mathbf{u}\|_p = 1 \}$$

Steepest Descent Methods

(Normalized) steepest descent, what does it mean?

We search for the unit vector that shows the largest decrease in the linear approximation of f :

$$\Delta \mathbf{x} = \underset{\mathbf{u}}{\operatorname{argmin}} \{ \nabla f(\mathbf{x})^T \mathbf{u}; \|\mathbf{u}\|_p = 1 \}$$

Conclusions:

- The steepest descent direction depends on the chosen norm.
- The negative gradient is not necessarily the best choice for the search direction.

Steepest Descent Methods (cont.)

We consider now the first order Taylor approximation of $f(\mathbf{x} + \mathbf{u})$ around the selected position \mathbf{x} :

$$f(\mathbf{x} + \mathbf{u}) \approx f(\mathbf{x}) + \nabla f(\mathbf{x})^T \mathbf{u}.$$

Steepest Descent Methods (cont.)

We consider now the first order Taylor approximation of $f(\mathbf{x} + \mathbf{u})$ around the selected position \mathbf{x} :

$$f(\mathbf{x} + \mathbf{u}) \approx f(\mathbf{x}) + \nabla f(\mathbf{x})^T \mathbf{u}.$$

- Here $\nabla f(\mathbf{x})^T \mathbf{u}$ is the directional derivative at \mathbf{x} in direction \mathbf{u} .
- The vector \mathbf{u} denotes a descent direction if the inner product with the gradient vector is negative, i. e.

$$\nabla f(\mathbf{x})^T \mathbf{u} < 0 .$$

Steepest Descent Methods (cont.)

Input: function f , initial estimate $\mathbf{x}^{(0)}$, norm $\|\cdot\|$

initialize: $k := 0$

repeat

 Compute highest descent direction:

$$\Delta \mathbf{x}^{(k)} = \underset{\mathbf{u}}{\operatorname{argmin}} \{ \nabla f(\mathbf{x}^{(k)})^T \mathbf{u}; \|\mathbf{u}\| = 1 \}$$

 Line search (1-D optimization):

$$t^{(k)} = \underset{t \geq 0}{\operatorname{argmin}} f(\mathbf{x}^{(k)} + t \cdot \Delta \mathbf{x}^{(k)})$$

 Update:

$$\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} + t^{(k)} \Delta \mathbf{x}^{(k)}.$$

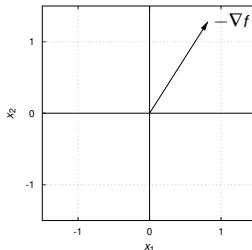
$k := k + 1$

until $\|\mathbf{x}^{(k)} - \mathbf{x}^{(k-1)}\| < \varepsilon$

Output: $\mathbf{x}^{(k)}$

L_2 -Norm

The unit ball for the L_2 -norm:

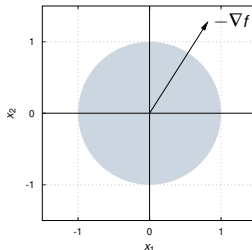


For the L_2 -norm the steepest descent direction is the negative gradient:

$$\Delta \mathbf{x} = \underset{\mathbf{u}}{\operatorname{argmin}} \{ \nabla f(\mathbf{x})^T \mathbf{u}; \|\mathbf{u}\|_2 = 1 \} = -\nabla f(\mathbf{x})$$

L_2 -Norm

The unit ball for the L_2 -norm:

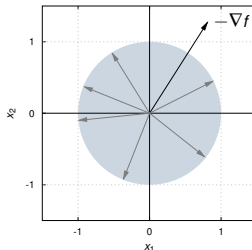


For the L_2 -norm the steepest descent direction is the negative gradient:

$$\Delta \mathbf{x} = \underset{\mathbf{u}}{\operatorname{argmin}} \{ \nabla f(\mathbf{x})^T \mathbf{u}; \|\mathbf{u}\|_2 = 1 \} = -\nabla f(\mathbf{x})$$

L_2 -Norm

The unit ball for the L_2 -norm:

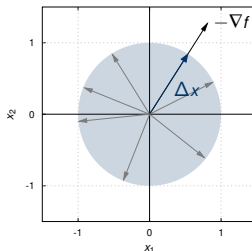


For the L_2 -norm the steepest descent direction is the negative gradient:

$$\Delta \mathbf{x} = \underset{\mathbf{u}}{\operatorname{argmin}} \{ \nabla f(\mathbf{x})^T \mathbf{u}; \|\mathbf{u}\|_2 = 1 \} = -\nabla f(\mathbf{x})$$

L_2 -Norm

The unit ball for the L_2 -norm:

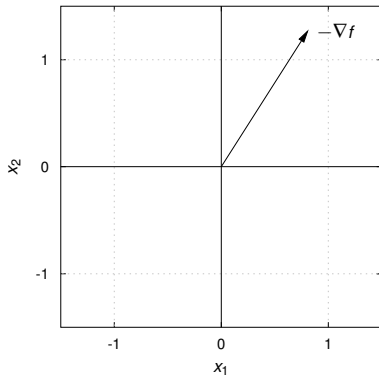


For the L_2 -norm the steepest descent direction is the negative gradient:

$$\Delta \mathbf{x} = \underset{\mathbf{u}}{\operatorname{argmin}} \{ \nabla f(\mathbf{x})^T \mathbf{u}; \|\mathbf{u}\|_2 = 1 \} = -\nabla f(\mathbf{x})$$

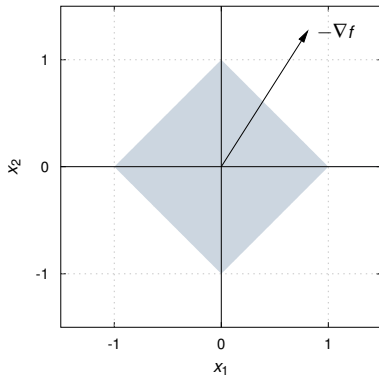
L_1 -Norm

The unit ball for the L_1 -norm:



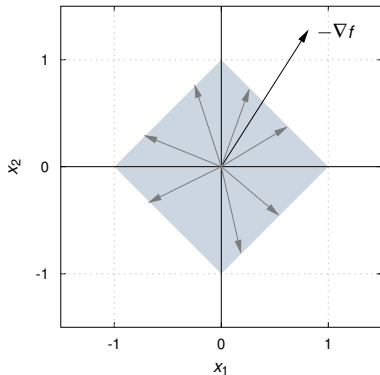
L_1 -Norm

The unit ball for the L_1 -norm:



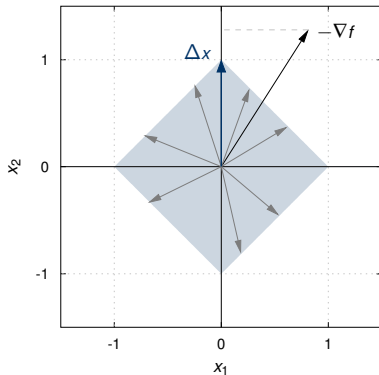
L_1 -Norm

The unit ball for the L_1 -norm:



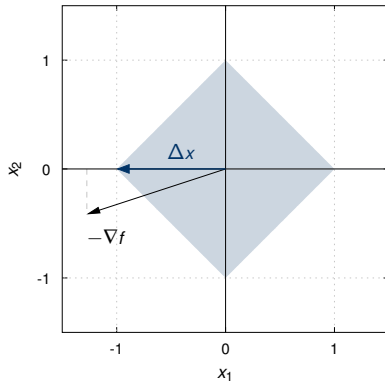
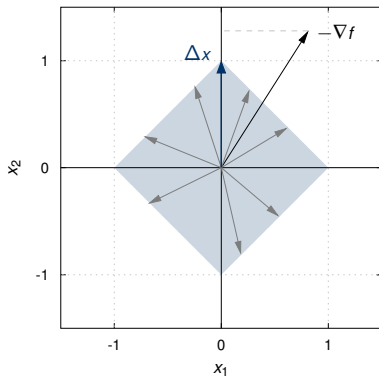
L_1 -Norm

The unit ball for the L_1 -norm:



L_1 -Norm

The unit ball for the L_1 -norm:



L_1 -Norm (cont.)

- The steepest descent for the L_1 -norm selects in each iteration the component of $\nabla f(\mathbf{x})$ with maximum absolute value and then decreases or increases dependent on the sign of the selected component.

L_1 -Norm (cont.)

- The steepest descent for the L_1 -norm selects in each iteration the component of $\nabla f(\mathbf{x})$ with maximum absolute value and then decreases or increases dependent on the sign of the selected component.
- Let i be the index of the gradient component with maximum absolute value, and let $\mathbf{e}_i \in \mathbb{R}^d$ denote the corresponding base vector. The steepest descent direction is given by:

$$\begin{aligned}\Delta \mathbf{x} &= \underset{\mathbf{u}}{\operatorname{argmin}} \{ \nabla f(\mathbf{x})^T \mathbf{u}; \|\mathbf{u}\|_1 = 1 \} \\ &= -\operatorname{sgn} \left(\frac{\partial}{\partial x_i} f(\mathbf{x}) \right) \mathbf{e}_i\end{aligned}$$

L_1 -Norm (cont.)

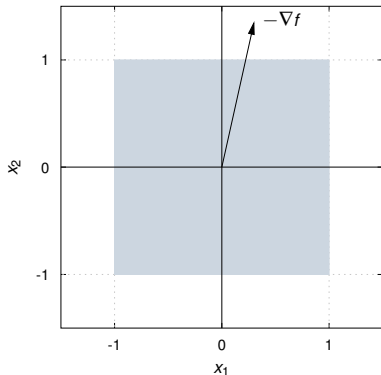
- The steepest descent for the L_1 -norm selects in each iteration the component of $\nabla f(\mathbf{x})$ with maximum absolute value and then decreases or increases dependent on the sign of the selected component.
- Let i be the index of the gradient component with maximum absolute value, and let $\mathbf{e}_i \in \mathbb{R}^d$ denote the corresponding base vector. The steepest descent direction is given by:

$$\begin{aligned}\Delta \mathbf{x} &= \underset{\mathbf{u}}{\operatorname{argmin}} \{ \nabla f(\mathbf{x})^T \mathbf{u}; \|\mathbf{u}\|_1 = 1 \} \\ &= -\operatorname{sgn} \left(\frac{\partial}{\partial x_i} f(\mathbf{x}) \right) \mathbf{e}_i\end{aligned}$$

- Note:** Steepest descent using the L_1 -norm results in the *coordinate descent algorithm*.

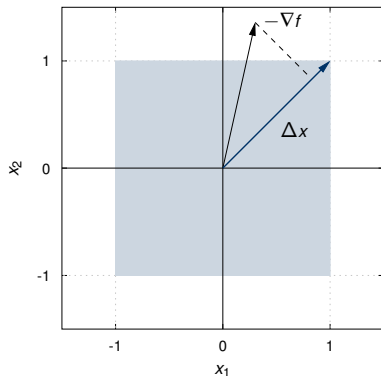
L_∞ -Norm

The unit ball for the L_∞ -norm:



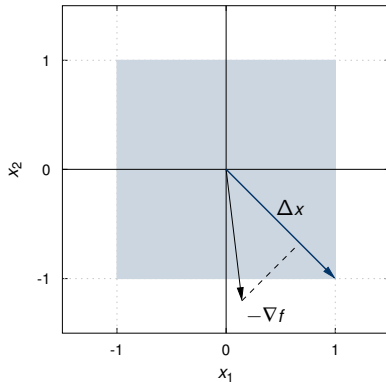
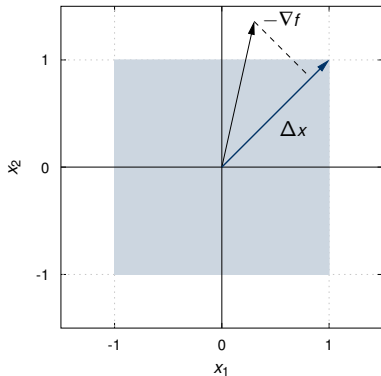
L_∞ -Norm

The unit ball for the L_∞ -norm:



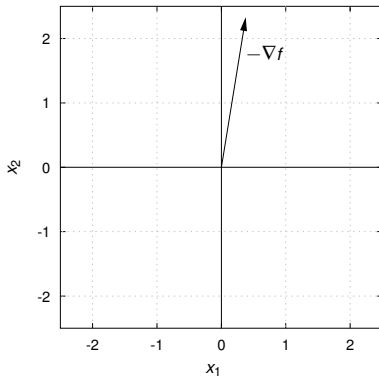
L_∞ -Norm

The unit ball for the L_∞ -norm:



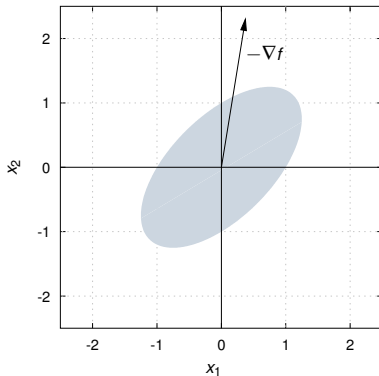
L_p -Norm

The unit ball for the L_p -norm:



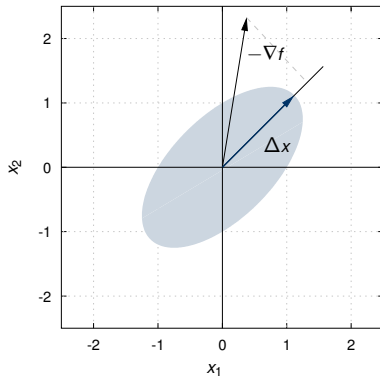
L_p -Norm

The unit ball for the L_p -norm:



L_p -Norm

The unit ball for the L_p -norm:



L_p -Norm (cont.)

The steepest descent for the L_p -norm is given by:

$$\Delta \mathbf{x} = \underset{\mathbf{u}}{\operatorname{argmin}} \{ \nabla f(\mathbf{x})^T \mathbf{u}; \|\mathbf{u}\|_p = 1 \}$$

L_P -Norm (cont.)

The steepest descent for the L_P -norm is given by:

$$\begin{aligned}\Delta \mathbf{x} &= \underset{\mathbf{u}}{\operatorname{argmin}} \{ \nabla f(\mathbf{x})^T \mathbf{u}; \|\mathbf{u}\|_P = 1 \} \\ &= \underset{\mathbf{u}}{\operatorname{argmin}} \{ \nabla f(\mathbf{x})^T \mathbf{u}; (\mathbf{u}^T \mathbf{P} \mathbf{u})^{\frac{1}{2}} = 1 \}\end{aligned}$$

L_P -Norm (cont.)

The steepest descent for the L_P -norm is given by:

$$\begin{aligned}\Delta \mathbf{x} &= \underset{\mathbf{u}}{\operatorname{argmin}} \{ \nabla f(\mathbf{x})^T \mathbf{u}; \|\mathbf{u}\|_P = 1 \} \\ &= \underset{\mathbf{u}}{\operatorname{argmin}} \{ \nabla f(\mathbf{x})^T \mathbf{u}; (\mathbf{u}^T \mathbf{P} \mathbf{u})^{\frac{1}{2}} = 1 \} \\ &= \underset{\mathbf{u}}{\operatorname{argmin}} \{ \nabla f(\mathbf{x})^T \mathbf{u}; \|\mathbf{P}^{\frac{1}{2}} \mathbf{u}\|_2 = 1 \}\end{aligned}$$

L_P -Norm (cont.)

The steepest descent for the L_P -norm is given by:

$$\begin{aligned}\Delta \mathbf{x} &= \underset{\mathbf{u}}{\operatorname{argmin}} \{ \nabla f(\mathbf{x})^T \mathbf{u}; \|\mathbf{u}\|_P = 1 \} \\ &= \underset{\mathbf{u}}{\operatorname{argmin}} \{ \nabla f(\mathbf{x})^T \mathbf{u}; (\mathbf{u}^T \mathbf{P} \mathbf{u})^{\frac{1}{2}} = 1 \} \\ &= \underset{\mathbf{u}}{\operatorname{argmin}} \{ \nabla f(\mathbf{x})^T \mathbf{u}; \|\mathbf{P}^{\frac{1}{2}} \mathbf{u}\|_2 = 1 \}\end{aligned}$$

As we did in the LDA-transform, we introduce a transform to get spherical data:

$$\mathbf{u}' = \mathbf{P}^{\frac{1}{2}} \mathbf{u}$$

L_P -Norm (cont.)

The steepest descent for the L_P -norm is given by:

$$\begin{aligned}\Delta \mathbf{x} &= \underset{\mathbf{u}}{\operatorname{argmin}} \{ \nabla f(\mathbf{x})^T \mathbf{u}; \|\mathbf{u}\|_P = 1 \} \\ &= \underset{\mathbf{u}}{\operatorname{argmin}} \{ \nabla f(\mathbf{x})^T \mathbf{u}; (\mathbf{u}^T \mathbf{P} \mathbf{u})^{\frac{1}{2}} = 1 \} \\ &= \underset{\mathbf{u}}{\operatorname{argmin}} \{ \nabla f(\mathbf{x})^T \mathbf{u}; \|\mathbf{P}^{\frac{1}{2}} \mathbf{u}\|_2 = 1 \}\end{aligned}$$

As we did in the LDA-transform, we introduce a transform to get spherical data:

$$\mathbf{u}' = \mathbf{P}^{\frac{1}{2}} \mathbf{u}$$

and thus

$$f(\mathbf{u}) = f(\mathbf{P}^{-\frac{1}{2}} \mathbf{u}') = f'(\mathbf{u}')$$

L_P -Norm (cont.)

Instead of $f(\mathbf{x})$ we now minimize $f'(\mathbf{x}')$ using the L_2 -norm and back-transform the result:

$$\Delta \mathbf{x}' = \underset{\mathbf{u}}{\operatorname{argmin}} \{ \nabla f'(\mathbf{x}')^T \mathbf{u}'; \|\mathbf{u}'\|_2 = 1 \}$$

L_P -Norm (cont.)

Instead of $f(\mathbf{x})$ we now minimize $f'(\mathbf{x}')$ using the L_2 -norm and back-transform the result:

$$\begin{aligned}\Delta \mathbf{x}' &= \underset{\mathbf{u}}{\operatorname{argmin}} \{ \nabla f'(\mathbf{x}')^T \mathbf{u}' ; \|\mathbf{u}'\|_2 = 1 \} \\ &= -\nabla f'(\mathbf{x}')\end{aligned}$$

L_P -Norm (cont.)

Instead of $f(\mathbf{x})$ we now minimize $f'(\mathbf{x}')$ using the L_2 -norm and back-transform the result:

$$\begin{aligned}\Delta \mathbf{x}' &= \underset{\mathbf{u}}{\operatorname{argmin}} \{ \nabla f'(\mathbf{x}')^T \mathbf{u}' ; \|\mathbf{u}'\|_2 = 1 \} \\ &= -\nabla f'(\mathbf{x}') \\ &= -\mathbf{P}^{-\frac{1}{2}} \nabla f(\mathbf{P}^{-\frac{1}{2}} \mathbf{x}')\end{aligned}$$

L_P -Norm (cont.)

Instead of $f(\mathbf{x})$ we now minimize $f'(\mathbf{x}')$ using the L_2 -norm and back-transform the result:

$$\begin{aligned}
 \Delta \mathbf{x}' &= \underset{\mathbf{u}}{\operatorname{argmin}} \{ \nabla f'(\mathbf{x}')^T \mathbf{u}; \|\mathbf{u}'\|_2 = 1 \} \\
 &= -\nabla f'(\mathbf{x}') \\
 &= -\mathbf{P}^{-\frac{1}{2}} \nabla f(\mathbf{P}^{-\frac{1}{2}} \mathbf{x}') \\
 &= -\mathbf{P}^{-\frac{1}{2}} \nabla f(\mathbf{x})
 \end{aligned}$$

L_P -Norm (cont.)

Now we get for $\Delta \mathbf{x}$:

$$\Delta \mathbf{x} = \mathbf{P}^{-\frac{1}{2}} \Delta \mathbf{x}'$$

L_P -Norm (cont.)

Now we get for $\Delta \mathbf{x}$:

$$\begin{aligned}\Delta \mathbf{x} &= \mathbf{P}^{-\frac{1}{2}} \Delta \mathbf{x}' \\ &= \mathbf{P}^{-\frac{1}{2}} \left(-\mathbf{P}^{-\frac{1}{2}} \nabla f(\mathbf{x}) \right)\end{aligned}$$

L_P -Norm (cont.)

Now we get for $\Delta \mathbf{x}$:

$$\begin{aligned}\Delta \mathbf{x} &= \mathbf{P}^{-\frac{1}{2}} \Delta \mathbf{x}' \\ &= \mathbf{P}^{-\frac{1}{2}} \left(-\mathbf{P}^{-\frac{1}{2}} \nabla f(\mathbf{x}) \right) \\ &= -\mathbf{P}^{-1} \nabla f(\mathbf{x}).\end{aligned}$$

L_P -Norm (cont.)

Now we get for $\Delta \mathbf{x}$:

$$\begin{aligned}\Delta \mathbf{x} &= \mathbf{P}^{-\frac{1}{2}} \Delta \mathbf{x}' \\ &= \mathbf{P}^{-\frac{1}{2}} \left(-\mathbf{P}^{-\frac{1}{2}} \nabla f(\mathbf{x}) \right) \\ &= -\mathbf{P}^{-1} \nabla f(\mathbf{x}).\end{aligned}$$

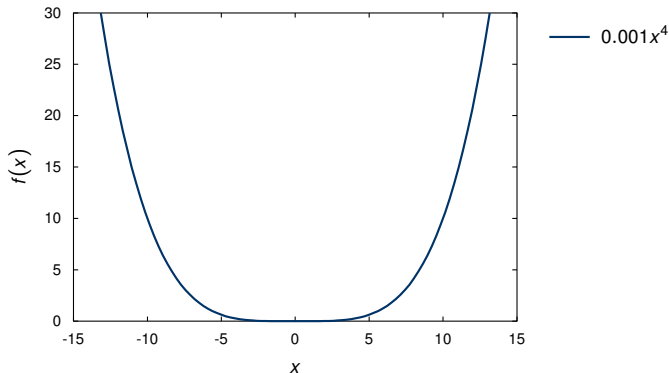
Conclusion: The steepest descent for the L_P -norm is given by

$$\Delta \mathbf{x} = -\mathbf{P}^{-1} \nabla f(\mathbf{x}) .$$

Newton's Method

The idea:

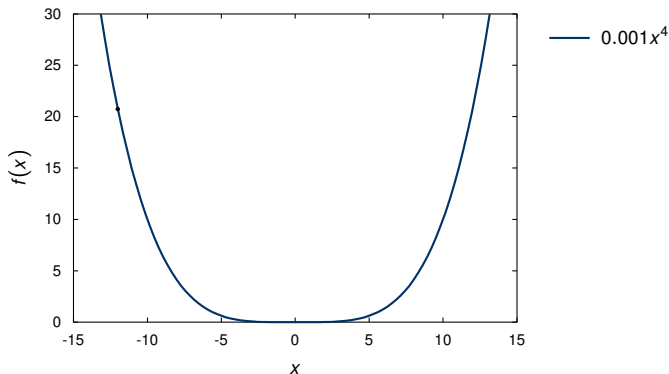
- Select a point.
- Compute the minimum of the second order Taylor approximation.



Newton's Method

The idea:

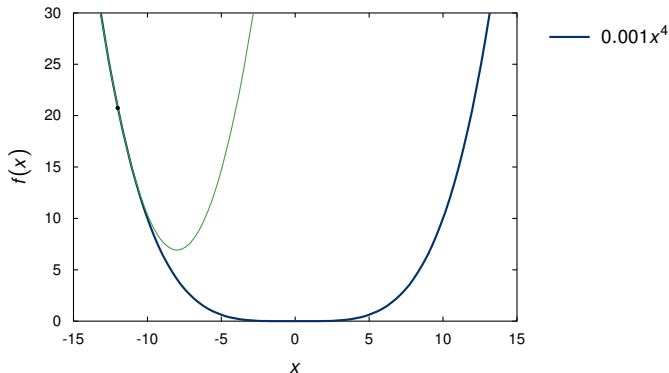
- Select a point.
- Compute the minimum of the second order Taylor approximation.



Newton's Method

The idea:

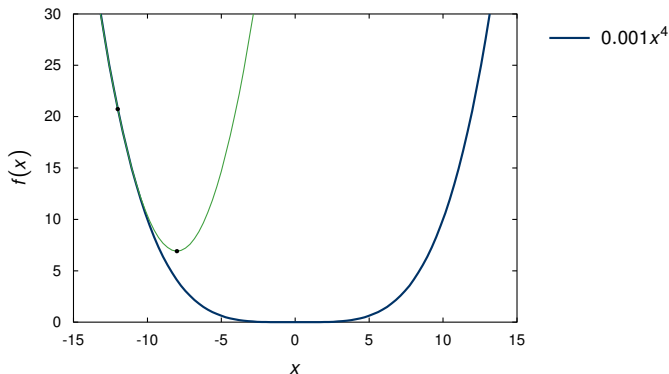
- Select a point.
- Compute the minimum of the second order Taylor approximation.



Newton's Method

The idea:

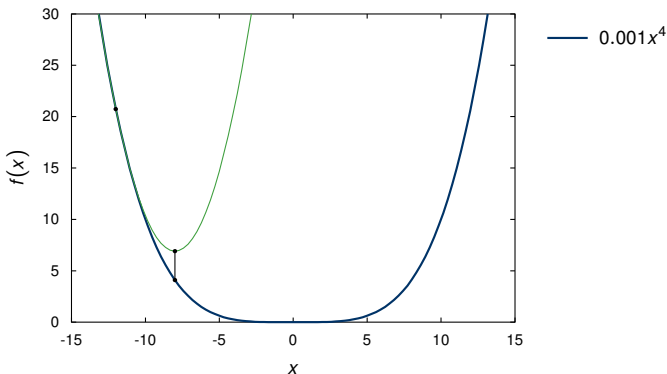
- Select a point.
- Compute the minimum of the second order Taylor approximation.



Newton's Method

The idea:

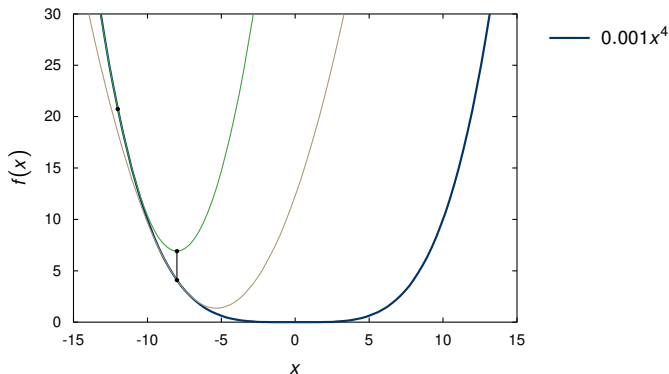
- Select a point.
- Compute the minimum of the second order Taylor approximation.



Newton's Method

The idea:

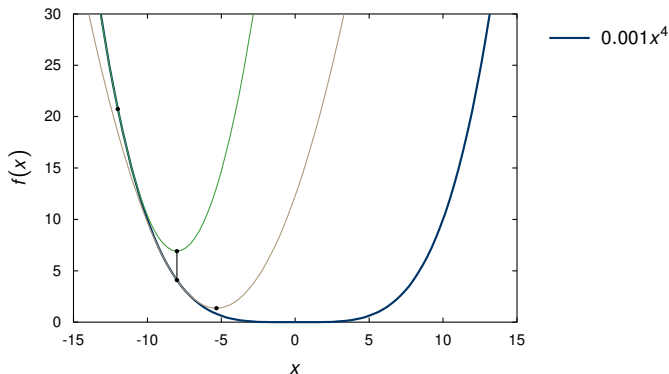
- Select a point.
- Compute the minimum of the second order Taylor approximation.



Newton's Method

The idea:

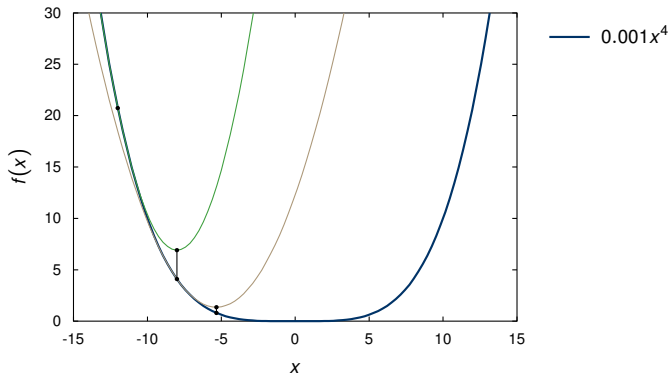
- Select a point.
- Compute the minimum of the second order Taylor approximation.



Newton's Method

The idea:

- Select a point.
- Compute the minimum of the second order Taylor approximation.



Newton's Method (cont.)

Second order Taylor approximation:

$$f(\mathbf{x} + \Delta\mathbf{x}) \approx f(\mathbf{x}) + \nabla f(\mathbf{x})^T \Delta\mathbf{x} + \frac{1}{2} \Delta\mathbf{x}^T (\nabla^2 f(\mathbf{x})) \Delta\mathbf{x}$$

Newton's Method (cont.)

Second order Taylor approximation:

$$f(\mathbf{x} + \Delta\mathbf{x}) \approx f(\mathbf{x}) + \nabla f(\mathbf{x})^T \Delta\mathbf{x} + \frac{1}{2} \Delta\mathbf{x}^T (\nabla^2 f(\mathbf{x})) \Delta\mathbf{x}$$

Now we select $\Delta\mathbf{x}$ such that

$$\nabla \{ f(\mathbf{x}) + \nabla f(\mathbf{x})^T \Delta\mathbf{x} + \frac{1}{2} \Delta\mathbf{x}^T (\nabla^2 f(\mathbf{x})) \Delta\mathbf{x} \} = 0$$

Newton's Method (cont.)

Second order Taylor approximation:

$$f(\mathbf{x} + \Delta\mathbf{x}) \approx f(\mathbf{x}) + \nabla f(\mathbf{x})^T \Delta\mathbf{x} + \frac{1}{2} \Delta\mathbf{x}^T (\nabla^2 f(\mathbf{x})) \Delta\mathbf{x}$$

Now we select $\Delta\mathbf{x}$ such that

$$\nabla \{ f(\mathbf{x}) + \nabla f(\mathbf{x})^T \Delta\mathbf{x} + \frac{1}{2} \Delta\mathbf{x}^T (\nabla^2 f(\mathbf{x})) \Delta\mathbf{x} \} = 0$$

Obviously the gradient is

$$\nabla f(\mathbf{x}) + \nabla^2 f(\mathbf{x}) \Delta\mathbf{x} = 0$$

Newton's Method (cont.)

Second order Taylor approximation:

$$f(\mathbf{x} + \Delta\mathbf{x}) \approx f(\mathbf{x}) + \nabla f(\mathbf{x})^T \Delta\mathbf{x} + \frac{1}{2} \Delta\mathbf{x}^T (\nabla^2 f(\mathbf{x})) \Delta\mathbf{x}$$

Now we select $\Delta\mathbf{x}$ such that

$$\nabla \{ f(\mathbf{x}) + \nabla f(\mathbf{x})^T \Delta\mathbf{x} + \frac{1}{2} \Delta\mathbf{x}^T (\nabla^2 f(\mathbf{x})) \Delta\mathbf{x} \} = 0$$

Obviously the gradient is

$$\nabla f(\mathbf{x}) + \nabla^2 f(\mathbf{x}) \Delta\mathbf{x} = 0$$

and thus

$$\Delta\mathbf{x} = -(\nabla^2 f(\mathbf{x}))^{-1} \nabla f(\mathbf{x})$$

Newton's Method (cont.)

Conclusion:

Newton's method is an \mathbf{x} -dependent steepest descent method regarding the $L_{\mathbf{P}}$ -norm, where $\mathbf{P} = \nabla^2 f(\mathbf{x})$ is the Hessian.

Damped Newton's Method

Input: function f , initial estimate $\mathbf{x}^{(0)}$

initialize: $k := 0$

repeat

 Compute Newton step:

$$\Delta \mathbf{x}^{(k)} = -\nabla^2 f(\mathbf{x}^{(k)})^{-1} \nabla f(\mathbf{x}^{(k)})$$

 Line search (1-D optimization):

$$t^{(k)} = \underset{t \geq 0}{\operatorname{argmin}} f(\mathbf{x}^{(k)} + t \cdot \Delta \mathbf{x}^{(k)})$$

 Update:

$$\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} + t^{(k)} \Delta \mathbf{x}^{(k)}.$$

$k := k + 1$

until $\|\mathbf{x}^{(k)} - \mathbf{x}^{(k-1)}\| < \varepsilon$

Output: $\mathbf{x}^{(k)}$

Lessons Learned

- Gradient descent is widely applied.
- Gradient descent and coordinate descent are special cases of steepest descent methods.
- Steepest descent method depends on the chosen norm.



**Pattern
Recognition
Lab**



**FRIEDRICH-ALEXANDER
UNIVERSITÄT
ERLANGEN-NÜRNBERG**

TECHNISCHE FAKULTÄT


Next Time in

Pattern Recognition



Further Readings

This chapter is basically copied from:

- S. Boyd, L. Vandenberghe:
[Convex Optimization](#),
Cambridge University Press, 2004.
 <http://www.stanford.edu/~boyd/cvxbook/>
- Jorge Nocedal, Stephen Wright:
[Numerical Optimization](#),
Springer, New York, 1999.

Comprehensive Questions

- What is the general formulation for an unconstrained optimization problem?
- Why do we need a line search in gradient descent approaches?
- What is the Armijo-Goldstein line search algorithm?
- What are the steepest descent directions if we apply the L_∞ , L_1 , L_2 , and L_p norm?