

# Modern Optimization Techniques

## 2. Unconstrained Optimization / 2.1. Gradient Descent

Lars Schmidt-Thieme

Information Systems and Machine Learning Lab (ISMLL)  
Institute of Computer Science  
University of Hildesheim, Germany

# Syllabus

Mon. 30.10.	(0)	0. Overview
		<b>1. Theory</b>
Mon. 6.11.	(1)	1. Convex Sets and Functions
		<b>2. Unconstrained Optimization</b>
Mon. 13.11.	(2)	2.1 Gradient Descent
Mon. 20.11.	(3)	2.2 Stochastic Gradient Descent
Mon. 27.11.	(4)	2.3 Newton's Method
Mon. 4.12.	(5)	2.4 Quasi-Newton Methods
Mon. 11.12.	(6)	2.5 Subgradient Methods
Mon. 18.12.	(7)	2.6 Coordinate Descent
	—	— <i>Christmas Break</i> —
		<b>3. Equality Constrained Optimization</b>
Mon. 8.1.	(8)	3.1 Duality
Mon. 15.1.	(9)	3.2 Methods
		<b>4. Inequality Constrained Optimization</b>
Mon. 22.1.	(10)	4.1 Primal Methods
Mon. 29.1.	(11)	4.2 Barrier and Penalty Methods
Mon. 5.2.	(12)	4.3 Cutting Plane Methods

# Outline

1. Unconstrained Optimization
2. Descent Methods
3. Gradient Descent
4. Line search
5. Convergence of Gradient Descent
6. Example: Linear Ridge Regression via Gradient Descent

# Outline

1. Unconstrained Optimization
2. Descent Methods
3. Gradient Descent
4. Line search
5. Convergence of Gradient Descent
6. Example: Linear Ridge Regression via Gradient Descent

# Unconstrained Convex Optimization Problem

$$\arg \min_{\mathbf{x} \in \mathbb{R}^N} f(\mathbf{x})$$

where

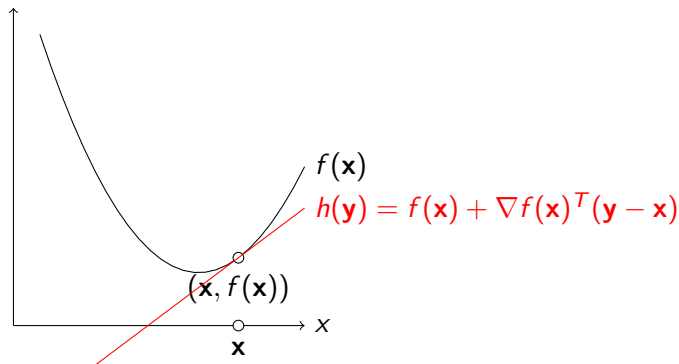
- ▶  $f : X \rightarrow \mathbb{R}, X \subseteq \mathbb{R}^N$  is
  - ▶ convex
  - ▶ twice continuously differentiable
  - ▶ esp.  $\text{dom } f = X = \mathbb{R}^N$  or open.
  
- ▶ An optimal  $\mathbf{x}^*$  exists and  $p^* := f(\mathbf{x}^*)$  is finite

# Reminder: 1st-order condition

**1st-order condition:** a differentiable function  $f$  is convex iff

- ▶  $\text{dom } f$  is a convex set
- ▶ for all  $\mathbf{x}, \mathbf{y} \in \text{dom } f$

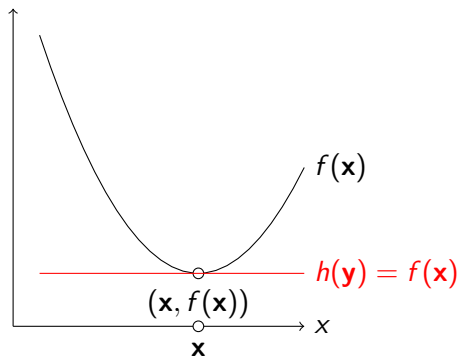
$$f(\mathbf{y}) \geq f(\mathbf{x}) + \nabla f(\mathbf{x})^T (\mathbf{y} - \mathbf{x})$$



# Optimality condition

$\mathbf{x}$  is optimal iff

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



# Methods for Unconstrained Optimization

- ▶ Start with an initial point:  $\mathbf{x}^{(0)}$
- ▶ Generate a sequence of points:  $\mathbf{x}^{(k)}$  with

$$f(\mathbf{x}^{(k)}) \rightarrow f(\mathbf{x}^*)$$

```
1 min-unconstrained( $f, \mathbf{x}^{(0)}$ ):  
2    $k := 0$   
3   repeat  
4      $\mathbf{x}^{(k+1)} := \mathbf{next-point}(f, \mathbf{x}^{(k)})$   
5      $k := k + 1$   
6   until converged( $\mathbf{x}^{(k)}, \mathbf{x}^{(k-1)}, f$ )  
7   return  $\mathbf{x}^{(k)}, f(\mathbf{x}^{(k)})$ 
```



# Methods for Unconstrained Optimization

- ▶ Start with an initial point:  $\mathbf{x}^{(0)}$
- ▶ Generate a sequence of points:  $\mathbf{x}^{(k)}$  with

$$f(\mathbf{x}^{(k)}) \rightarrow f(\mathbf{x}^*)$$

```
1 min-unconstrained( $f, \mathbf{x}^{(0)}, k^{\max}$ ):  
2   for  $k := 0 : k^{\max} - 1$ :  
3      $\mathbf{x}^{(k+1)} := \mathbf{next-point}(f, \mathbf{x}^{(k)})$   
4     if converged( $\mathbf{x}^{(k+1)}, \mathbf{x}^{(k)}, f$ ):  
5       return  $\mathbf{x}^{(k+1)}, f(\mathbf{x}^{(k+1)})$   
6   raise exception "not converged in  $k^{\max}$  iterations"
```

# Convergence Criterion

$$\text{converged}(\mathbf{x}^{(k+1)}, \mathbf{x}^{(k)}, f)$$

- ▶ Different criteria in use
  - ▶ different optimization methods may use different criteria

- ▶ One would like to use the **optimality gap**:

$$\|\mathbf{x}^{(k+1)} - \mathbf{x}^{\star}\|_2^2 < \epsilon$$

- ▶ not possible as  $\mathbf{x}^{\star}$  is unknown

- ▶ **Minimum progress/change  $\epsilon$  in  $x$  in last iteration:**

$$\text{converged}(\mathbf{x}^{(k+1)}, \mathbf{x}^{(k)}, f) := \|\mathbf{x}^{(k+1)} - \mathbf{x}^{(k)}\|_2^2 < \epsilon$$

- ▶ cheap to compute
  - ▶ can be used with any method
  - ▶ requires parameter  $\epsilon \in \mathbb{R}^+$

# Outline

1. Unconstrained Optimization
2. Descent Methods
3. Gradient Descent
4. Line search
5. Convergence of Gradient Descent
6. Example: Linear Ridge Regression via Gradient Descent

# Descent Methods

- ▶ A class/template of methods

- ▶ The next point is generated as:

$$\mathbf{x}^{(k+1)} := \mathbf{x}^{(k)} + \mu \Delta \mathbf{x}^{(k)}$$

with

- ▶ a **search direction**  $\Delta \mathbf{x}^{(k)}$  and
- ▶ a **step size**  $\mu$  such that

$$f(\mathbf{x}^{(k)} + \mu \Delta \mathbf{x}^{(k)}) < f(\mathbf{x}^{(k)})$$

- ▶ Specific descent methods differ in how they compute the search direction  $\Delta \mathbf{x}^{(k)}$ 
  - ▶ Gradient Descent
  - ▶ Steepest Descent
  - ▶ Newton's Method

# Descent Methods

```
1 min-descent( $f, \mathbf{x}^{(0)}, k^{\max}$ ):  
2   for  $k := 0 : k^{\max} - 1$ :  
3      $\Delta \mathbf{x}^{(k)} := \mathbf{search-direction}(f, \mathbf{x}^{(k)})$   
4      $\mu^{(k)} := \mathbf{step-size}(f, \mathbf{x}^{(k)}, \Delta \mathbf{x}^{(k)})$   
5      $\mathbf{x}^{(k+1)} := \mathbf{x}^{(k)} + \mu^{(k)} \Delta \mathbf{x}^{(k)}$   
6     if converged( $\mathbf{x}^{(k+1)}, \mathbf{x}^{(k)}, f$ ):  
7       return  $\mathbf{x}^{(k+1)}, f(\mathbf{x}^{(k+1)})$   
8   raise exception "not converged in  $k^{\max}$  iterations"
```

# Computing the Step Size

The step size can be computed in various ways:

- ▶ constant value
- ▶ line search
- ▶ various heuristics depending on the specific algorithm

# Outline

1. Unconstrained Optimization
2. Descent Methods
3. Gradient Descent
4. Line search
5. Convergence of Gradient Descent
6. Example: Linear Ridge Regression via Gradient Descent

# Gradient Descent

- ▶ The gradient of a function  $f : X \rightarrow \mathbb{R}, X \subseteq \mathbb{R}^N$  at  $\mathbf{x}$  yields the direction in which the function is maximally growing locally.
- ▶ Gradient Descent is a descent method that searches in the opposite direction of the gradient:

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

- ▶ Gradient:

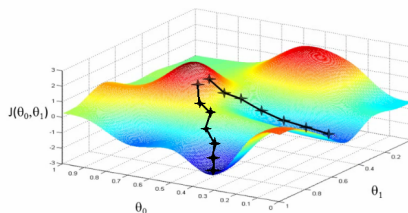
$$\nabla f(\mathbf{x}) := \nabla_{\mathbf{x}} f(\mathbf{x}) := \left( \frac{\partial f}{\partial x_n}(\mathbf{x}) \right)_{n=1:N}$$



# Gradient Descent

```

1 min-GD( $f, \mathbf{x}^{(0)}, k^{\max}$ ):
2   for  $k := 0 : k^{\max} - 1$ :
3      $\Delta \mathbf{x}^{(k)} := -\nabla f(\mathbf{x}^{(k)})$ 
4      $\mu^{(k)} := \text{step-size}(f, \mathbf{x}^{(k)}, \Delta \mathbf{x}^{(k)})$ 
5      $\mathbf{x}^{(k+1)} := \mathbf{x}^{(k)} + \mu^{(k)} \Delta \mathbf{x}^{(k)}$ 
6     if converged( $\mathbf{x}^{(k+1)}, \mathbf{x}^{(k)}, f$ ):
7       return  $\mathbf{x}^{(k+1)}, f(\mathbf{x}^{(k+1)})$ 
8   raise exception "not converged in  $k^{\max}$  iterations"
  
```



# Gradient Descent / Implementations

- ▶ for analysis usually all updated variables are indexed

$$\mathbf{x}^{(k)}, \Delta \mathbf{x}^{(k)}, \mu^{(k)}$$

- ▶ in implementations, one usually does only need one copy
  - ▶ or two, to compare against the last one

```
1 min-GD( $f, \mathbf{x}, k^{\max}$ ):  
2   for  $k := 0 : k^{\max} - 1$ :  
3      $\Delta \mathbf{x} := -\nabla f(\mathbf{x})$   
4      $\mu := \mathbf{step-size}(f, \mathbf{x}, \Delta \mathbf{x})$   
5      $\mathbf{x}^{\text{old}} := \mathbf{x}$   
6      $\mathbf{x} := \mathbf{x}^{\text{old}} + \mu \Delta \mathbf{x}$   
7     if converged( $\mathbf{x}, \mathbf{x}^{\text{old}}, f$ ):  
8       return  $\mathbf{x}, f(\mathbf{x})$   
9   raise exception "not converged in  $k^{\max}$  iterations"
```

# Gradient Descent / Considerations

- ▶ Stopping criterion:  $\|\nabla f(\mathbf{x})\|_2 \leq \epsilon$

**converged**( $\mathbf{x}, \mathbf{x}^{\text{old}}, f$ ) :=

**converged**( $\nabla f(\mathbf{x})$ ) :=  $\|\nabla f(\mathbf{x})\|_2 \leq \epsilon$

- ▶ cheap to use as GD has to compute the gradient anyway
- ▶ GD is simple and straightforward
- ▶ GD has slow convergence
  - ▶ esp. compared to Newton's method
- ▶ Out-of-the-box, GD works only well for convex problems, otherwise will get stuck in local minima

# Gradient Descent Example

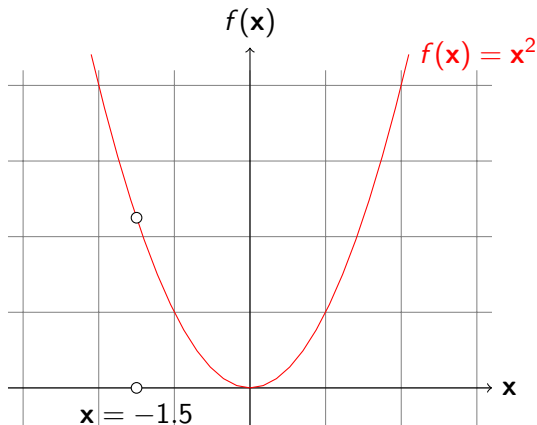
## Task:

minimize  $x^2$

►  $\mu = 0.3$

►  $-\nabla f(\mathbf{x}) = -2\mathbf{x}$

Initial point:  $x^0 = -1.5$



# Gradient Descent Example

## Task:

$$\text{minimize } \mathbf{x}^2$$

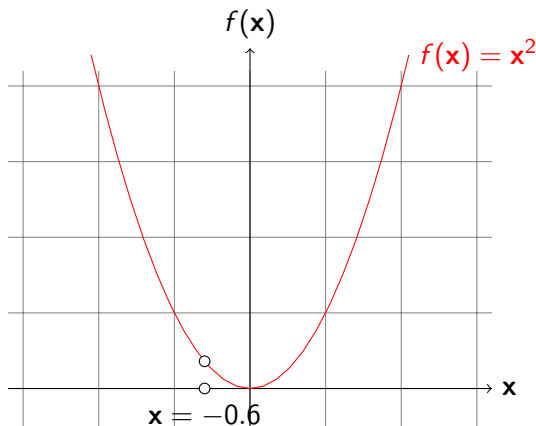
$$\blacktriangleright \mu = 0.3$$

$$\blacktriangleright -\nabla f(\mathbf{x}) = -2\mathbf{x}$$

$$\mathbf{x}^0 = -1.5$$

$$\mathbf{x} = -1.5 - 0.3 \cdot (2 \cdot -1.5)$$

$$\mathbf{x} = -0.6$$



# Gradient Descent Example

## Task:

$$\text{minimize } x^2$$

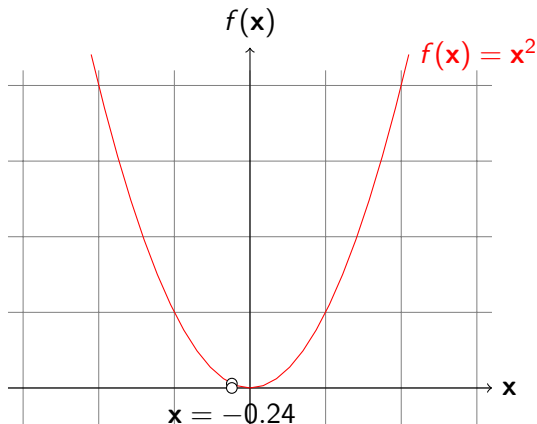
$$\triangleright \mu = 0.3$$

$$\triangleright -\nabla f(\mathbf{x}) = -2\mathbf{x}$$

$$\mathbf{x} = -0.6$$

$$\mathbf{x} = -0.6 - 0.3 \cdot (2 \cdot -0.6)$$

$$\mathbf{x} = -0.24$$



# Gradient Descent Example

## Task:

$$\text{minimize } x^2$$

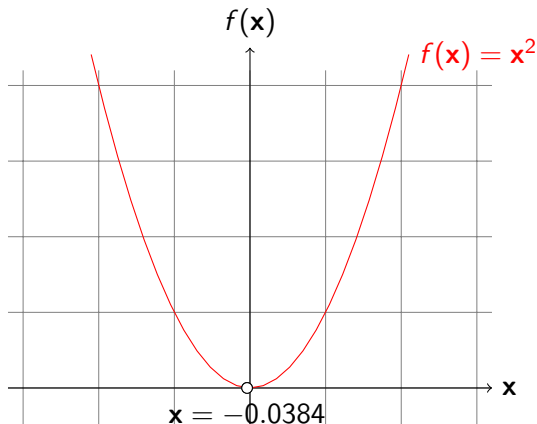
$$\blacktriangleright \mu = 0.3$$

$$\blacktriangleright -\nabla f(\mathbf{x}) = -2\mathbf{x}$$

$$x = -0.24$$

$$x = -0.24 - 0.3 \cdot (2 \cdot -0.24)$$

$$x = -0.0384$$



# Gradient Descent Example

## Task:

$$\text{minimize } x^2$$

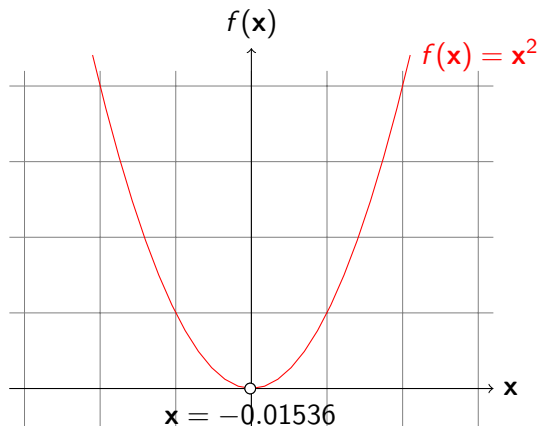
$$\blacktriangleright \mu = 0.3$$

$$\blacktriangleright -\nabla f(\mathbf{x}) = -2\mathbf{x}$$

$$x = -0.0384$$

$$x = -0.0384 - 0.3 \cdot (2 \cdot -0.0384)$$

$$x = -0.01536$$





# Considerations about the Step Size

- ▶ Crucial for the convergence of the algorithm
- ▶ Step size too small  $\rightsquigarrow$  slow convergence
- ▶ Step size too large  $\rightsquigarrow$  divergence!

# Gradient Descent Example - A perfect Step Size

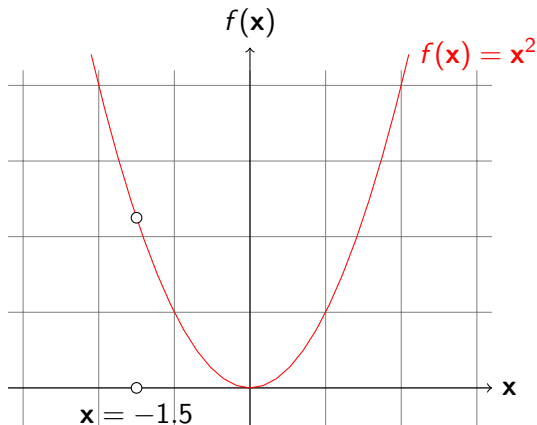
## Task:

minimize  $x^2$

►  $\mu = 0.5$

►  $-\nabla f(\mathbf{x}) = -2\mathbf{x}$

Initial point:  $x^0 = -1.5$



# Gradient Descent Example - A perfect Step Size

## Task:

$$\text{minimize } \mathbf{x}^2$$

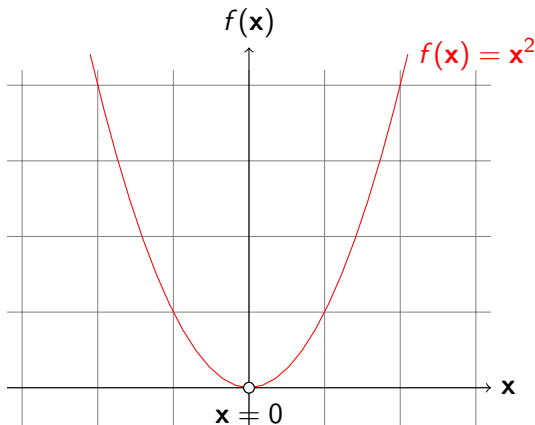
$$\blacktriangleright \mu = 0.5$$

$$\blacktriangleright -\nabla f(\mathbf{x}) = -2\mathbf{x}$$

$$\mathbf{x}^0 = -1.5$$

$$\mathbf{x} = -1.5 - 0.5 \cdot (2 \cdot -1.5)$$

$$\mathbf{x} = 0$$



# Gradient Descent Example - Too Large Step Size

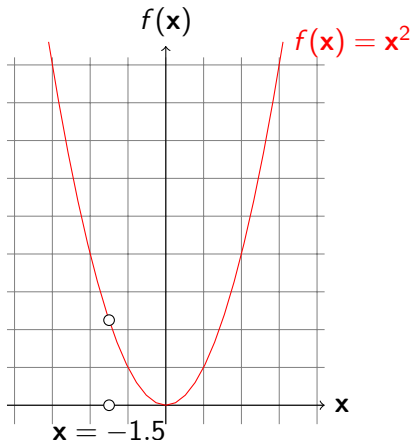
## Task:

minimize  $x^2$

►  $\mu = 1.5$

►  $-\nabla f(\mathbf{x}) = -2\mathbf{x}$

Initial point:  $\mathbf{x}^0 = -1.5$



# Gradient Descent Example - Too Large Step Size

## Task:

$$\text{minimize } \mathbf{x}^2$$

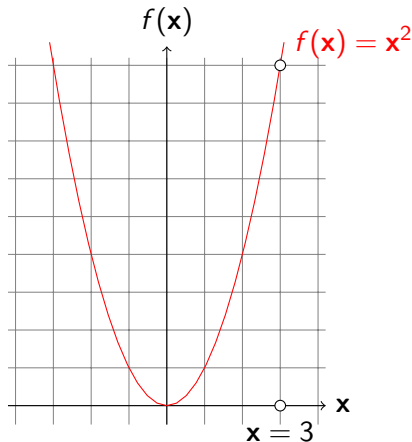
►  $\mu = 1.5$

►  $-\nabla f(\mathbf{x}) = -2\mathbf{x}$

$$\mathbf{x}^0 = -1.5$$

$$\mathbf{x} = -1.5 - 1.5 \cdot (2 \cdot -1.5)$$

$$\mathbf{x} = 3$$



# Gradient Descent Example - Too Large Step Size

## Task:

$$\text{minimize } \mathbf{x}^2$$

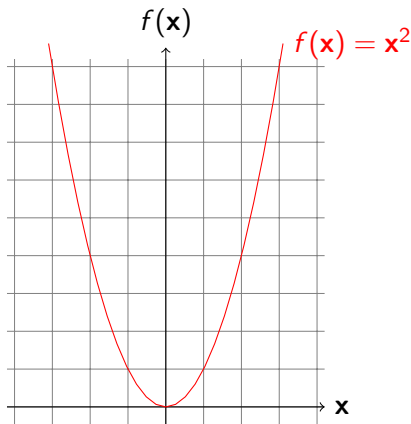
►  $\mu = 1.5$

►  $-\nabla f(\mathbf{x}) = -2\mathbf{x}$

$$\mathbf{x}^0 = 3$$

$$\mathbf{x} = 3 - 1.5 \cdot (2 \cdot 3)$$

$$\mathbf{x} = -6$$



# Outline

1. Unconstrained Optimization
2. Descent Methods
3. Gradient Descent
4. Line search
5. Convergence of Gradient Descent
6. Example: Linear Ridge Regression via Gradient Descent

# Line search

- ▶ **line search** is the task to compute the step length in a descent algorithm.
- ▶ a one-dimensional optimization problem in  $\mu$ :

$$\arg \min_{\mu \in \mathbb{R}^+} f(\mathbf{x} + \mu \Delta \mathbf{x})$$



# Line Search Methods

## ► exact line search

- Used if the problem can be solved analytically or with low cost
- e.g., for **unconstrained quadratic optimization**:

$$\arg \min_{x \in \mathbb{R}^N} f(x) := \frac{1}{2} x^T A x + b^T x, \quad A \in \mathbb{R}^{N \times N} \text{ pos. def.}, b \in \mathbb{R}^N$$

# Line Search Methods

## ► exact line search

- Used if the problem can be solved analytically or with low cost
- e.g., for **unconstrained quadratic optimization**:

$$\arg \min_{x \in \mathbb{R}^N} f(x) := \frac{1}{2} x^T A x + b^T x, \quad A \in \mathbb{R}^{N \times N} \text{ pos. def.}, b \in \mathbb{R}^N$$

## ► backtracking line search

- only approximative
- guarantees that the new function value is lower than a specific bound

# Backtracking Line Search

```
1 stepsize-backtracking( $f, \mathbf{x}, \Delta\mathbf{x}, \alpha \in (0, 0.5), \beta \in (0, 1)$ ):  
2    $\mu := 1$   
3   while  $f(\mathbf{x} + \mu\Delta\mathbf{x}) > f(\mathbf{x}) + \alpha\mu\nabla f(\mathbf{x})^T \Delta\mathbf{x}$ :  
4      $\mu := \beta\mu$   
5   return  $\mu$ 
```

# Backtracking Line Search

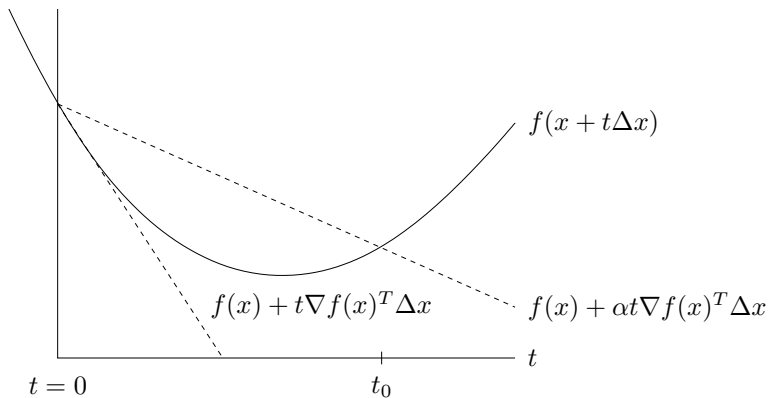
```
1 stepsize-backtracking( $f, \mathbf{x}, \Delta\mathbf{x}, \alpha \in (0, 0.5), \beta \in (0, 1)$ ):  
2    $\mu := 1$   
3   while  $f(\mathbf{x} + \mu\Delta\mathbf{x}) > f(\mathbf{x}) + \alpha\mu\nabla f(\mathbf{x})^T \Delta\mathbf{x}$ :  
4      $\mu := \beta\mu$   
5   return  $\mu$ 
```

Loop eventually terminates: for sufficient small  $\mu$ :

$$f(\mathbf{x} + \mu\Delta\mathbf{x}) \approx f(\mathbf{x}) + \mu\nabla f(\mathbf{x})^T \Delta\mathbf{x} < f(\mathbf{x}) + \alpha\mu\nabla f(\mathbf{x})^T \Delta\mathbf{x}$$

as for a descent direction:  $\nabla f(\mathbf{x})^T \Delta\mathbf{x} < 0$

# Backtracking Line Search



source: [Boyd and Vandenberghe, 2004, p. 465]

# Outline

1. Unconstrained Optimization
2. Descent Methods
3. Gradient Descent
4. Line search
5. Convergence of Gradient Descent
6. Example: Linear Ridge Regression via Gradient Descent

# Sublevel Sets

**sublevel set** of  $f : X \rightarrow \mathbb{R}, X \subseteq \mathbb{R}^N$  at level  $\alpha \in \mathbb{R}$ :

$$S_\alpha := \{x \in \text{dom } f \mid f(x) \leq \alpha\}$$

# Sublevel Sets

**sublevel set** of  $f : X \rightarrow \mathbb{R}, X \subseteq \mathbb{R}^N$  at level  $\alpha \in \mathbb{R}$ :

$$S_\alpha := \{x \in \text{dom } f \mid f(x) \leq \alpha\}$$

basic facts:

- ▶ if  $f$  is convex, then all its sublevel sets  $S_\alpha$  are convex sets.
  - ▶ useful to show that a set is convex
    - ▶ show that it can be represented as a sublevel set of a convex function.



# Closed Functions

$f : X \rightarrow \mathbb{R}, X \subseteq \mathbb{R}^N$  **closed** :  $\Longleftrightarrow$  all its sublevel sets are closed.

# Closed Functions

$f : X \rightarrow \mathbb{R}, X \subseteq \mathbb{R}^N$  **closed** :  $\Longleftrightarrow$  all its sublevel sets are closed.

examples:

- ▶  $f(x) = x^2$  is closed.
- ▶  $f(x) = 1/x$  on  $\mathbb{R}^+$  is closed.
- ▶  $f(x) = x \log x$  on  $\mathbb{R}^+$  is not closed.
- ▶ but  $f$  on  $\mathbb{R}_0^+$  defined by

$$f(x) := \begin{cases} x \log x, & \text{if } x > 0 \\ 0, & \text{else} \end{cases}$$

is closed.

# Closed Functions

$f : X \rightarrow \mathbb{R}, X \subseteq \mathbb{R}^N$  **closed** :  $\Longleftrightarrow$  all its sublevel sets are closed.

examples:

- ▶  $f(x) = x^2$  is closed.
- ▶  $f(x) = 1/x$  on  $\mathbb{R}^+$  is closed.
- ▶  $f(x) = x \log x$  on  $\mathbb{R}^+$  is not closed.
- ▶ but  $f$  on  $\mathbb{R}_0^+$  defined by

$$f(x) := \begin{cases} x \log x, & \text{if } x > 0 \\ 0, & \text{else} \end{cases}$$

is closed.

Classes of closed functions:

- ▶ continuous functions on all of  $\mathbb{R}^N$
- ▶ continuous functions on an open set  
that go to infinity everywhere towards the border

# Semidefinite Matrices II

Let  $A, B \in \mathbb{R}^{N \times N}$  symmetric matrices:

$$A \succeq B : \Longleftrightarrow A - B \succeq 0$$

- ▶  $A \succeq mI, m \in \mathbb{R}^+$ :
  - ▶ all eigenvalues of  $A$  are  $\geq m$
- ▶  $A \preceq MI, M \in \mathbb{R}^+$ :
  - ▶ all eigenvalues of  $A$  are  $\leq M$

# Strongly Convex Functions

Let  $f : X \rightarrow \mathbb{R}$ ,  $X \subseteq \mathbb{R}^N$  be twice continuously differentiable.

$f$  is **strongly convex** :  $\Longleftrightarrow$

- ▶  $\text{dom } f = X$  is convex and
- ▶ the eigenvalues of the Hessian are uniformly bounded from below:

$$\nabla^2 f(x) \succeq ml, \quad \exists m \in \mathbb{R}^+ \quad \forall x \in \text{dom } f$$

# Strongly Convex Functions

Let  $f : X \rightarrow \mathbb{R}$ ,  $X \subseteq \mathbb{R}^N$  be twice continuously differentiable.

$f$  is **strongly convex** :  $\Longleftrightarrow$

- ▶  $\text{dom } f = X$  is convex and
- ▶ the eigenvalues of the Hessian are uniformly bounded from below:

$$\nabla^2 f(x) \succeq mI, \quad \exists m \in \mathbb{R}^+ \quad \forall x \in \text{dom } f$$

Every strongly convex function  $f$  is also strictly convex.

- ▶ but not the other way around
  - ▶  $f(x) = x^4$  on  $\mathbb{R}^+$  is strictly, but not strongly convex
- ▶ do not confuse strongly and strictly convex!

# Strongly Convex Functions / Basic Facts

(i)  $f$  is above a hyperbola:

$$f(y) \geq f(x) + \nabla f(x)^T (y - x) + \frac{m}{2} \|y - x\|_2^2$$

$$p^* \geq f(x) - \frac{1}{2m} \|\nabla f(x)\|_2^2$$

(ii) if  $f$  is closed and  $S$  one of its sublevel sets, then

a) the eigenvalues of the Hessian are also uniformly bounded from above on  $S$ :

$$\nabla^2 f(x) \preceq MI, \quad \exists M \in \mathbb{R}^+ \quad \forall x \in S$$

b)

$$f(y) \leq f(x) + \nabla f(x)^T (y - x) + \frac{M}{2} \|y - x\|_2^2, \quad x, y \in S$$

$$p^* \leq f(x) - \frac{1}{2M} \|\nabla f(x)\|_2^2$$

# Strongly Convex Functions / Basic Facts / Proofs

(i) for  $x, y \in \text{dom } f \exists z \in [x, y]$

(Taylor expansion with Lagrange mean value remainder):

$$f(y) = f(x) + \nabla f(x)^T (y - x) + \frac{1}{2} \underbrace{(y - x)^T \nabla^2 f(z) (y - x)}_{\geq m \|y - x\|_2^2}$$

$$\begin{aligned} f(y) &\geq f(x) + \nabla f(x)^T (y - x) + \frac{m}{2} \|y - x\|_2^2 \\ &\geq \min_y f(x) + \nabla f(x)^T (y - x) + \frac{m}{2} \|y - x\|_2^2 \end{aligned}$$

considered as function in  $y$  has

minimum at  $\tilde{y} := x - \frac{1}{m} \nabla f(x)$

$$= f(x) + \nabla f(x)^T (\tilde{y} - x) + \frac{m}{2} \|\tilde{y} - x\|_2^2$$

$$= f(x) - \frac{1}{2m} \|\nabla f(x)\|_2^2$$

$$\rightsquigarrow p^* = f(y = x^*) \geq f(x) - \frac{1}{2m} \|\nabla f(x)\|_2^2$$



## Strongly Convex Functions / Basic Facts / Proofs (2/2)

- (ii.a)    ▶ due to (i) all sublevel sets are bounded
- ▶ the maximal eigenvalue of  $\nabla^2 f(x)$  is a continuous function on a closed bounded set and thus itself bounded,
- ▶ i.e., it exists  $M \in \mathbb{R}^+$ :  $\nabla^2 f(x) \preceq MI$
- (ii.b) as for (i), using (ii.a)

# Convergence of Gradient Descent / Exact Line Search

If

- ▶  $f$  is strongly convex,
- ▶ the initial sublevel set  $S := \{x \in \text{dom } f \mid f(x) \leq f(x^{(0)})\}$  is closed,
- ▶ an exact line search is used,

then

$$f(x^{(k)}) - p^* \leq \left(1 - \frac{m}{M}\right)^k (f(x^{(0)}) - p^*)$$

Equivalently, to guarantee  $f(x^{(k)}) - p^* \leq \epsilon$ , GD requires

$$k := \frac{\log \frac{f(x^{(0)}) - p^*}{\epsilon}}{\log \frac{1}{1 - \frac{m}{M}}} \quad \text{iterations.}$$

Especially,

- ▶ GD converges, i.e.,  $f(x^{(k)})$  approaches  $p^*$
- ▶ the convergence is exponential in  $k$  (with basis  $c := 1 - \frac{m}{M}$ )
  - ▶ called **linear convergence** in the optimization literature

# Convergence of Gradient Descent / Proof

$$\begin{aligned}\tilde{f}(t) &:= f(x - t\nabla f(x)), \quad t \in \{t \in \mathbb{R}_0^+ \mid x - t\nabla f(x) \in S\} \\ f(x^{\text{next}}) &= \tilde{f}(t_{\text{exact}}) \\ &\leq \tilde{f}(0) - \frac{1}{2M} \|\nabla \tilde{f}(0)\|_2^2, \quad \tilde{f} \text{ strongly convex (ii.b)}\end{aligned}$$

$$= f(x) - \frac{1}{2M} \underbrace{\|\nabla f(x)\|_2^2}_{\geq 2m(f(x) - p^*)}, \quad f \text{ strongly convex (i)}$$

$$f(x^{\text{next}}) - p^* \leq f(x) - p^* - \frac{1}{2M} 2m(f(x) - p^*) = \left(1 - \frac{m}{M}\right)(f(x) - p^*)$$

$$f(x^{(k)}) - p^* \leq \left(1 - \frac{m}{M}\right)^k (f(x^{(0)}) - p^*)$$

# Convergence of Gradient Descent / Backtracking

If

- ▶  $f$  is strongly convex,
- ▶ the initial sublevel set  $S := \{x \in \text{dom } f \mid f(x) \leq f(x^{(0)})\}$  is closed, and
- ▶ a **backtracking line search** is used,

then

$$f(x^{(k)}) - p^* \leq c^k (f(x^{(0)}) - p^*), \quad c := 1 - \min\{2\alpha m, 2\beta\alpha m/M\}$$

Equivalently, to guarantee  $f(x^{(k)}) - p^* \leq \epsilon$ , GD requires

$$k := \frac{\log \frac{f(x^{(0)}) - p^*}{\epsilon}}{\log \frac{1}{c}} \quad \text{iterations.}$$

Especially,

- ▶ GD converges, i.e.,  $f(x^{(k)})$  approaches  $p^*$
- ▶ the convergence is exponential in  $k$  (with basis  $c$ ; linear convergence)

# Outline

1. Unconstrained Optimization
2. Descent Methods
3. Gradient Descent
4. Line search
5. Convergence of Gradient Descent
6. Example: Linear Ridge Regression via Gradient Descent

# A More practical example

We do not want to always minimize parabolas so let us discuss a more practical example:

## Linear Regression!

- ▶ have  $m$  many data instances  $\mathbf{a} \in \mathbb{R}^n$  with  $n$  many features / predictors
- ▶ want to learn a linear model parametrized by a vector  $\beta \in \mathbb{R}^n$  to predict a real value  $y \in \mathbb{R}$

# Practical Example: Household Spending

If we have data about  $m$  households, we can represent it as:

$$A_{m,n} = \begin{pmatrix} 1 & a_{1,1} & a_{1,2} & a_{1,3} & a_{1,4} \\ 1 & a_{2,1} & a_{2,2} & a_{2,3} & a_{2,4} \\ \vdots & \vdots & \vdots & \vdots & \vdots \\ 1 & a_{m,1} & a_{m,2} & a_{m,3} & a_{m,4} \end{pmatrix} \mathbf{y} = \begin{pmatrix} y_1 \\ y_2 \\ \vdots \\ y_m \end{pmatrix}$$

We can model the household consumption is a linear combination of the household features with parameters  $\beta$ :

$$\hat{y}_i = \beta^T \mathbf{a}_i = \beta_0 1 + \beta_1 a_{i,1} + \beta_2 a_{i,2} + \beta_3 a_{i,3} + \beta_4 a_{i,4}$$

# Practical Example: Household Spending

We have:

$$\begin{pmatrix} 1 & a_{1,1} & a_{1,2} & a_{1,3} & a_{1,4} \\ 1 & a_{2,1} & a_{2,2} & a_{2,3} & a_{2,4} \\ \vdots & \vdots & \vdots & \vdots & \vdots \\ 1 & a_{m,1} & a_{m,2} & a_{m,3} & a_{m,4} \end{pmatrix} \cdot \begin{pmatrix} \beta_0 \\ \beta_1 \\ \beta_2 \\ \beta_3 \\ \beta_4 \end{pmatrix} \approx \begin{pmatrix} y_1 \\ y_2 \\ \vdots \\ y_m \end{pmatrix}$$

We want to find parameters  $\beta$  such that the measured error of the predictions is minimal:

$$\sum_{i=1}^m (\beta^T \mathbf{a}_i - y_i)^2 + \lambda \sum_{j=1}^n \beta_j^2 = \|A\beta - y\|_2^2 + \lambda \|\beta\|_2^2$$



# Linear Regression

Let us look at the function to optimize:

$$\begin{aligned}\mathcal{L}(\beta, A, y) + \lambda \text{Reg}(\beta) &= \sum_{i=1}^m (\beta^\top a_i - y_i)^2 + \lambda \|\beta\|_2^2 \\ &= \sum_{i=1}^m \left( \sum_{j=1}^n \beta_j a_{ij} - y_i \right)^2 + \lambda \sum_{j=1}^n \beta_j^2\end{aligned}$$

Then we can compute the gradient component wise:

$$\begin{aligned}\frac{\partial}{\partial \beta_k} \mathcal{L}(\beta, A, y) + \lambda \text{Reg}(\beta) &= \frac{\partial}{\partial \beta_k} \sum_{i=1}^m \left( \sum_{j=1}^n \beta_j a_{ij} - y_i \right)^2 + \lambda \sum_{j=1}^n \beta_j^2 \\ &= \sum_{i=1}^m 2 \cdot \left( \sum_{j=1}^n \beta_j a_{ij} - y_i \right) \cdot a_{ik} + 2\lambda \beta_k\end{aligned}$$

# Linear Regression

We obtain the update for every component of  $\beta$  as

$$\begin{aligned}\beta_k^{(k+1)} &= \beta_k^{(k)} - \mu \nabla_{\beta} (\mathcal{L}(\beta, A, y) + \lambda \text{Reg}(\beta)) \\ &= \beta_k^{(k)} - \mu \left( 2 \sum_{i=1}^m \cdot \left( \sum_{j=1}^n \beta_j a_{ij} - y_i \right) \cdot a_{ik} + 2\lambda \beta_k^{(k)} \right)\end{aligned}$$

- ▶ see that  $\left( \sum_{j=1}^n \beta_j a_{ij} - y_i \right)$  is actually the error of the model on the  $i$ -th instance
- ▶ error is the same for all  $k$ , can be precomputed

# Linear Regression

```
1: procedure LEARN LINEAR REGRESSION MODEL
   input: Data  $A$ , Labels  $y$ , initial parameters  $\beta^0$ , Step Size  $\mu$ ,
   Regularization constant  $\lambda$ , precision  $\epsilon$ 
2:   repeat
3:     Compute Error:  $e_i = \left( \sum_{j=1}^n \beta_j a_{ij} - y_i \right)$ 
4:     for  $k = 1, \dots, n$  do
5:        $\beta_k^{(k+1)} = \beta_k^{(k)} - \mu \left( \sum_{i=1}^m e_i a_{ik} + \lambda \beta_k^{(k)} \right)$ 
6:     end for
7:      $t = t + 1$ 
8:   until  $\|\nabla_{\beta} \mathcal{L}(\beta, A, y)\|_2^2 \leq \epsilon$ 
9:   return  $\beta, \mathcal{L}(\beta, A, y)$ 
10: end procedure
```

# Summary (1/2)

- ▶ **Unconstrained optimization** is the minimization of a function over all of  $\mathbb{R}^N$  or an open subset  $X \subseteq \mathbb{R}^N$ .
  - ▶ In **Unconstrained convex optimization**  $X$  also has to be convex (and  $f$ , too).
- ▶ **Descent methods** iteratively find a next iterate  $x^{(k+1)}$  with lower function value than the last iterate and require:
  - ▶ **search direction**: in which direction to search.
    - ▶ **Gradient Descent** (GD): negative gradient of the target function
  - ▶ **step length**: how far to go.
  - ▶ **convergence criterion**: when to stop.
    - ▶ small last step
    - ▶ small gradient

# Summary (2/2)

- ▶ step length (aka **line search**) in rare cases can be computed exactly.
  - ▶ one-dimensional optimization problem (**exact line search**)
- ▶ **backtracking line search**:
  - ▶ Choose the largest stepsize that guarantees a decrease in function value.
  - ▶ guaranteed to terminate
- ▶ GD has **linear convergence**
  - ▶ exponential in the number of steps
    - ▶ with basis  $1 - m/M$   
for smallest/largest eigenvalues  $m, M$  of the Hessian
  - ▶ if  $f$  is strongly convex, its initial sublevel set closed and exact line search is used.

# Further Readings

- ▶ Unconstrained minimization problems:
  - ▶ Boyd and Vandenberghe [2004], chapter 9.1
  
- ▶ Descent methods:
  - ▶ Boyd and Vandenberghe [2004], chapter 9.2
  
- ▶ Gradient descent:
  - ▶ Boyd and Vandenberghe [2004], chapter 9.3
  
- ▶ also accessible from here:
  - ▶ steepest descent — Boyd and Vandenberghe [2004], chapter 9.4

# References I

Stephen Boyd and Lieven Vandenberghe. *Convex Optimization*. Cambridge Univ Press, 2004.