# Modern Optimization Techniques

## 2. Unconstrained Optimization / 2.1. Gradient Descent

Lars Schmidt-Thieme

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

# Syllabus

# Outline

1. Unconstrained Optimization

2. Iterative and Descent Methods

3. Gradient Descent

4. Line search

5. Convergence of Gradient Descent

# Outline

## 1. Unconstrained Optimization

## 2. Iterative and Descent Methods

## 3. Gradient Descent

## 4. Line search

## 5. Convergence of Gradient Descent

# Unconstrained Convex Optimization Problem

$$\underset{x \in X}{\arg\min} \ f(\mathbf{x})$$

where
- $f : X \to \mathbb{R}, X \subseteq \mathbb{R}^N$ is
  - convex

  - twice continuously differentiable

  - esp. dom $f = X = \mathbb{R}^N$ or convex and open.

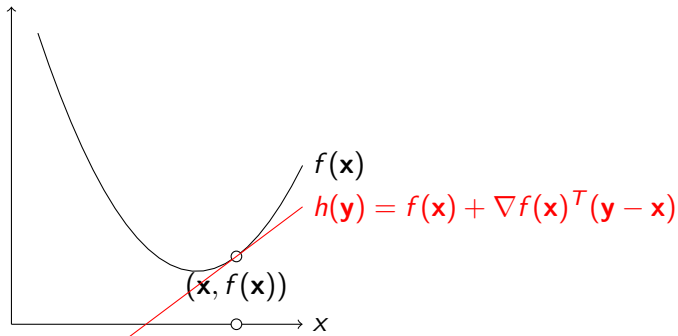- An optimal $\mathbf{x}^*$ exists and $p^* := f(\mathbf{x}^*)$ is finite

# Reminder: 1st-order condition

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

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

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

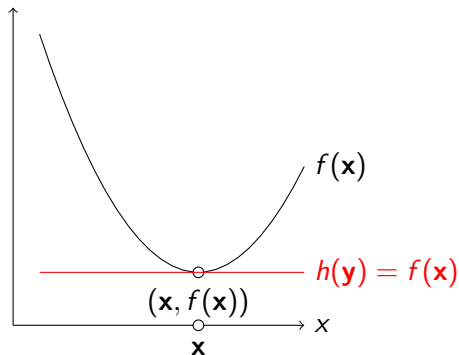(the function is above any of its tangents.)

# Minimality Condition

**x** is minimal iff

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

# Outline

# Iteative Methods

▶ Start with an initial (random) point: $\mathbf{x}^{(0)}$

▶ Generate a sequence of points: $\mathbf{x}^{(k)}$ with

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

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

# Iteative Methods

- ▶ Start with an initial (random) point: $\mathbf{x}^{(0)}$

- ▶ Generate a sequence of points: $\mathbf{x}^{(k)}$ with

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

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

# Convergence Criterion

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

▶ 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**:

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

  ▶ cheap to compute.
  ▶ can be used with any method.
  ▶ requires parameter $\epsilon \in \mathbb{R}^+$.
  ▶ may stop too early when the loss surface is too flat.

# 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** $\mu > 0$ such that

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

  ▶ always exists if the step size $\mu$ is sufficient small
    if the search direction $\Delta \mathbf{x}^{(k)}$ is a **descent direction**:

$$\nabla f(\mathbf{x}^{(k)})^T \Delta \mathbf{x}^{(k)} < 0$$

▶ search directions $\Delta \mathbf{x}^{(k)}$ can be computed different ways
  ▶ Gradient Descent
  ▶ Steepest Descent
  ▶ Newton's Method

# Descent Methods

```
1  min-descent(f, x^(0), K):
2    for k := 0 : K - 1:
3      Δx^(k) := search-direction(f, x^(k))
4      μ^(k) := step-size(f, x^(k), Δx^(k))
5      x^(k+1) := x^(k) + μ^(k) Δx^(k)
6      if converged(x^(k+1), x^(k), f):
7        return x^(k+1), f(x^(k+1))
8    raise exception "not converged in K iterations"
```

# Outline

# Gradient Descent

▶ The gradient of a function $f : X \to \mathbb{R}, X \subseteq \mathbb{R}^N$ 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:

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

# Gradient Descent

```
1  min-GD(f, x^(0), K):
2      for k := 0 : K − 1:
3          Δx^(k) := −∇f(x^(k))
4          μ^(k) := step-size(f, x^(k), Δx^(k))
5          x^(k+1) := x^(k) + μ^(k) Δx^(k)
6          if converged(x^(k+1), x^(k), f):
7              return x^(k+1), f(x^(k+1))
8      raise exception "not converged in K 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, x, K):
2    for k := 0 : K − 1:
3      Δx := −∇f(x)
4      μ := step-size(f, x, Δx)
5      x^old := x
6      x := x^old + μΔx
7      if converged(x, x^old, f):
8        return x, f(x)
9    raise exception "not converged in K iterations"
```

# Gradient Descent / Considerations

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

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

    ▶ 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 (see next chapter)

▶ Out-of-the-box, GD works only well for convex problems,
otherwise will get stuck in local minima.

# Gradient Descent Example

**Task:** minimize $f(x) := x^2$

- $\mu = 0.3$

- $-\nabla f(x) = -2x$

Initial point: $x^{(0)} = -1.5$



$f(x)$

$f(x) = x^2$

$x^{(0)} = -1.5$

$x$

# Gradient Descent Example

**Task:** minimize $f(x) := x^2$

▶ $\mu = 0.3$

▶ $-\nabla f(x) = -2x$

$x^{(0)} = -1.5$

$x^{(1)} = -1.5 - 0.3 \cdot (2 \cdot (-1.5))$

$\quad = -0.6$



$f(x)$

$f(x) = x^2$

$x^{(1)} = -0.6$

$x$

# Gradient Descent Example

**Task:** minimize $f(x) := x^2$

- $\mu = 0.3$

- $-\nabla f(x) = -2x$

$x^{(1)} = -0.6$

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

$\quad = -0.24$

# Gradient Descent Example

**Task:** minimize $f(x) := x^2$

- $\mu = 0.3$

- $-\nabla f(x) = -2x$

$x^{(2)} = -0.24$

$x^{(3)} = -0.24 - 0.3 \cdot (2 \cdot (-0.24))$

$\quad = -0.096$



$f(x)$

$f(x) = x^2$

$x^{(3)} = -0.096$

$x$

# Gradient Descent Example

**Task:** minimize $f(x) := x^2$

- $\mu = 0.3$

- $-\nabla f(x) = -2x$

$x^{(3)} = -0.096$

$x^{(4)} = -0.096 - 0.3 \cdot (2 \cdot (-0.096))$

$\quad = -0.0384$



$f(x)$

$f(x) = x^2$

$x^{(4)} = -0.0384$

$x$

# How About a Larger Step Size?



**Task:** minimize $f(x) := x^2$
- $\mu = 1.5$
- $-\nabla f(x) = -2x$

Initial point: $x^{(0)} = -1.5$

Figure labels:
- $f(x)$
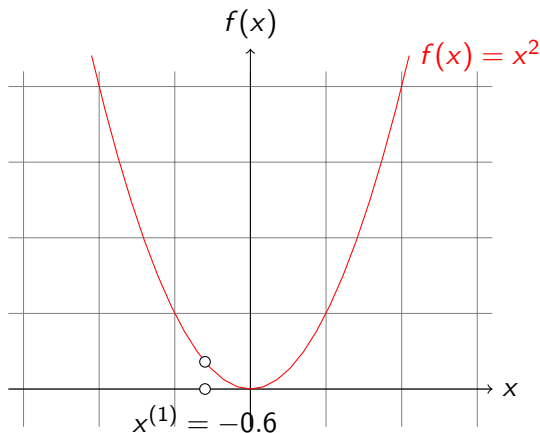- $f(x) = x^2$
- $x^{(0)} = -1.5$
- $x$

# How About a Larger Step Size?

**Task:** minimize $f(x) := x^2$

- $\mu = 1.5$
- $-\nabla f(x) = -2x$

$x^{(0)} = -1.5$

$x^{(1)} = -1.5 - 1.5 \cdot (2 \cdot (-1.5))$

$\quad\ = 3$



$f(x)$

$f(x) = x^2$

$x$

$x^{(1)} = 3$

# How About a Larger Step Size?

**Task:** minimize $f(x) := x^2$

- $\mu = 1.5$
- $-\nabla f(x) = -2x$

$$x^{(1)} = 3$$
$$x^{(2)} = 3 - 1.5 \cdot (2 \cdot 3)$$
$$= -6$$



$f(x)$

$f(x) = x^2$

$x$

# How About a Larger Step Size?

**Task:** minimize $f(x) := x^2$

- $\mu = 1.5$
- $-\nabla f(x) = -2x$

$x^{(1)} = 3$

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

$\quad\ = -6$

$\rightsquigarrow$ the algorithm diverges!



$f(x)$

$f(x) = x^2$

$x$

# Gradient Descent Example — Optimal Step Size

**Task:** minimize $f(x) := x^2$

- $\mu = 0.5$
- $-\nabla f(x) = -2x$

Initial point: $x^0 = -1.5$

# Gradient Descent Example — Optimal Step Size

**Task:** minimize $f(x) := x^2$

- $\mu = 0.5$
- $-\nabla f(x) = -2x$

$x^{(0)} = -1.5$

$x^{(1)} = -1.5 - 0.5 \cdot (2 \cdot (-1.5))$

$\quad = 0$

# Gradient Descent Example — Optimal Step Size

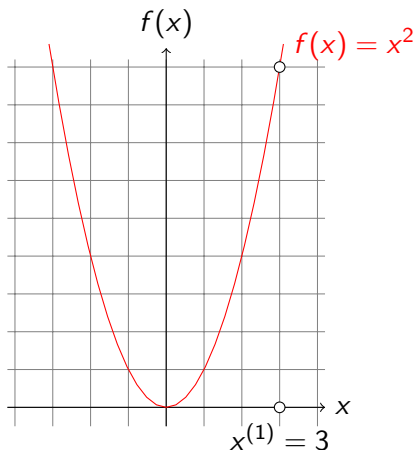**Task:** minimize $f(x) := x^2$

- $\mu = 0.5$
- $-\nabla f(x) = -2x$

$x^{(0)} = -1.5$

$x^{(1)} = -1.5 - 0.5 \cdot (2 \cdot (-1.5))$

$\quad = 0$

$\rightsquigarrow$ the algorithm converges in 1 step!



$f(x)$

$f(x) = x^2$

$x$

# How to Choose the Step Size $\mu$?

▶ Step size $\mu$ is crucial for the convergence of the algorithm.
  ▶ Step size too small. $\rightsquigarrow$ slow convergence.

  ▶ Step size too large. $\rightsquigarrow$ divergence!

▶ How to choose a good step size?
  $\rightsquigarrow$ **line search** (aka **step size control**).

# Outline

# Computing the Step Size

The step size can be computed in various ways:
- ▶ constant value
  - ▶ e.g., 1

- ▶ decreasing sequence, e.g., $\gamma^k$ for $\gamma \in (0, 1)$
  - ▶ e.g., for $\gamma = \frac{1}{2}$: $1, \frac{1}{2}, \frac{1}{4}, \frac{1}{8}, \ldots$

- ▶ line search

- ▶ various heuristics depending on the specific algorithm

# Line Search

▶ **line search** is the task to compute the step size in a descent algorithm.

▶ itself a one-dimensional optimization problem in $\mu$:

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

# Line Search Methods

- **exact line search**:
  - Used if the problem can be solved analytically or with low cost.

  - e.g., for **unconstrained quadratic optimization**:

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

# 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, x, Δx, α ∈ (0, 0.5), β ∈ (0, 1)):
2    μ := 1
3    while f(x + μΔx) > f(x) + αμ∇f(x)ᵀΔx:
4        μ := βμ
5    return μ
```

Q: Why does the backtracking condition guarantee $f(\mathbf{x}^{\text{next}}) < f(\mathbf{x})$ ?

# Backtracking Line Search

```
1  stepsize-backtracking(f, x, Δx, α ∈ (0, 0.5), β ∈ (0, 1)):
2      μ := 1
3      while  f(x + μΔx) > f(x) + αμ∇f(x)^T Δx:
4          μ := βμ
5      return  μ
```

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

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

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

# Backtracking Line Search



$f(x + t\Delta x)$

$f(x) + t\nabla f(x)^T \Delta x$     $f(x) + \alpha t\nabla f(x)^T \Delta x$

$t = 0$     $t_0$     $t$

source: Boyd and Vandenberghe, 2004, p. 465

# Outline

# Sublevel Sets

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

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

# Sublevel Sets

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

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

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.

# Sublevel Sets / Examples

$$S_\alpha(x^2) =$$

$$S_\alpha(-\log x; \mathbb{R}^+) =$$

$$S_\alpha(\frac{1}{x}; \mathbb{R}^+) =$$

$$S_\alpha(x; \mathbb{R}^+) =$$

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

# Sublevel Sets / Examples

$$S_\alpha(x^2) = \begin{cases} [-\sqrt{\alpha}, \sqrt{\alpha}], & \alpha \geq 0 \\ \emptyset, & \text{else} \end{cases}$$

$$S_\alpha(-\log x; \mathbb{R}^+) = [e^{-\alpha}, \infty)$$

$$S_\alpha(\frac{1}{x}; \mathbb{R}^+) = \begin{cases} [\frac{1}{\alpha}, \infty), & \alpha \geq 0 \\ \emptyset, & \text{else} \end{cases}$$

$$S_\alpha(x; \mathbb{R}^+) = \begin{cases} (0, \alpha], & \alpha > 0 \\ \emptyset, & \text{else} \end{cases}$$

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

# Closed Functions

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

# Closed Functions

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

examples:

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

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

# Closed Functions

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

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

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

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 $f : X \to \mathbb{R}, X \subseteq \mathbb{R}^N$ 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:

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

## Strongly Convex Functions

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

$f$ is **strongly convex** $:\Longleftrightarrow$
- ▶ dom $f = X$ is convex and

- ▶ the eigenvalues of the Hessian are uniformly bounded from below:

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

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 / Examples



Q: Is $f$ convex, strictly or strongly convex?

(convex: $\forall x : \nabla^2 f(x) \succeq 0$, strictly convex: $\forall x : \nabla^2 f(x) \succ 0$, strongly convex: $\exists m > 0 \; \forall x : \nabla^2 f(x) \succeq mI$)

# Strongly Convex Functions / Basic Facts

(i) $f$ is above a parabola:

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

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

(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 $S$:

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

   b) $f$ is below a parabola ("sandwiched between two parabolas"):

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

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

# 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):

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

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

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

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)

# Theorem (Convergence of Gradient Descent — exact line search)

If (i) $f$ is strongly convex,

(ii) the initial sublevel set $S := \{x \in \operatorname{dom} f \mid f(x) \leq f(x^{(0)})\}$ is closed,

(iii) an exact line search is used,

then

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

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

$$\tilde{f}(t) := f(x - t\nabla f(x)), \quad t \in \{t \in \mathbb{R}_0^+ \mid x - t\nabla f(x) \in S\}$$

$$f(x^{\text{next}}) = \tilde{f}(t_{\text{exact}}) = \tilde{p}^*, \qquad\qquad \tilde{p}^* := \min_t \tilde{f}(t)$$

$$\leq \tilde{f}(0) - \frac{1}{2M}(\tilde{f}'(0))^2, \qquad \tilde{f} \text{ strongly convex (ii.b)}$$

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

$$\leq f(x) - \frac{m}{M}(f(x) - p^*)$$

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

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

# Convergence of Gradient Descent / in $x$

GD's convergence can also be described in $x$ (instead of in $f$):

$$
\begin{aligned}
||x^{(k)} - x^*||^2 &\underset{\text{s.c.(i)}}{\leq} \frac{2}{m}(f(x^{(k)}) - p^*) \\
&\underset{\text{conv}}{\leq} \frac{2}{m}(1 - \frac{m}{M})^k(f(x^{(0)}) - p^*) \\
&\underset{\text{s.c.(i)}}{\leq} (1 - \frac{m}{M})^k \frac{2}{m} \frac{1}{2m} ||(\nabla f(x))||^2 \\
&= (1 - \frac{m}{M})^k \frac{||(\nabla f(x^{(0)}))||^2}{m^2}
\end{aligned}
$$

# Theorem (Convergence of Gradient Descent — Backtracking)

If (i) $f$ is strongly convex,
   (ii) the initial sublevel set $S := \{x \in \operatorname{dom} f \mid f(x) \leq f(x^{(0)})\}$ is closed,
   (iii) a backtracking line search is used,

then

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

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)

# 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 size**: how far to go.

  - **convergence criterion**: when to stop.
    - small last step

    - small gradient

# Summary (2/2)

- ▶ step size (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

Boyd, Stephen and Lieven Vandenberghe (2004). *Convex Optimization*. Cambridge University Press.