# Modern Optimization Techniques

## 2. Unconstrained Optimization / 2.5. Subgradient Methods

Lars Schmidt-Thieme

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

# Syllabus

# Outline

1. Subgradients

2. Subgradient Calculus

3. The Subgradient Method

4. Convergence

# Outline

# Motivation

▶ If a function is once differentiable
we can optimize it using
  ▶ Gradient Descent,

  ▶ Stochastic Gradient Descent,

  ▶ Quasi-Newton Methods
(1st order information)

▶ If a function is twice differentiable
we can optimize it using
  ▶ Newton's method
(2nd order information)

▶ What if the objective function is not differentiable?

# 1st-Order Condition for Convexity (Review)

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

- ▶ dom $f$ is a convex set and

- ▶ for all $\mathbf{x}, \mathbf{y} \in$ dom $f$

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

  - ▶ i.e., the tangent (= first order Taylor approximation) of $f$ at $\mathbf{x}$ is a global underestimator

# Tangent as a global underestimator

# Tangent as a global underestimator



$f(\mathbf{x})$

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

$x$

$x$

# Tangent as a global underestimator



$f(\mathbf{x})$

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

$x$

$x$

**What happens if $f$ is not differentiable?**

# Subgradient

Given a function $f$ and a point $\mathbf{x} \in \text{dom}\, f$,
$\mathbf{g} \in \mathbb{R}^N$ is called a **subgradient** of $f$ at $\mathbf{x}$ if:
the hypersurface with slopes $\mathbf{g}$ through $(\mathbf{x}, f(\mathbf{x}))$
is a global underestimator of $f$, i.e.

$$f(\mathbf{y}) \geq f(\mathbf{x}) + \mathbf{g}^T(\mathbf{y} - \mathbf{x}), \quad \text{for all } \mathbf{y} \in \text{dom}\, f$$

# Subgradient

Given a function $f$ and a point $\mathbf{x} \in \text{dom}\, f$,
$\mathbf{g} \in \mathbb{R}^N$ is called a **subgradient** of $f$ at $\mathbf{x}$ if:
the hypersurface with slopes $\mathbf{g}$ through $(\mathbf{x}, f(\mathbf{x}))$
is a global underestimator of $f$, i.e.

$$f(\mathbf{y}) \geq f(\mathbf{x}) + \mathbf{g}^T(\mathbf{y} - \mathbf{x}), \quad \text{for all } \mathbf{y} \in \text{dom}\, f$$

# Subgradient

Given a function $f$ and a point $\mathbf{x} \in \operatorname{dom} f$,
$\mathbf{g} \in \mathbb{R}^N$ is called a **subgradient** of $f$ at $\mathbf{x}$ if:
the hypersurface with slopes $\mathbf{g}$ through $(\mathbf{x}, f(\mathbf{x}))$
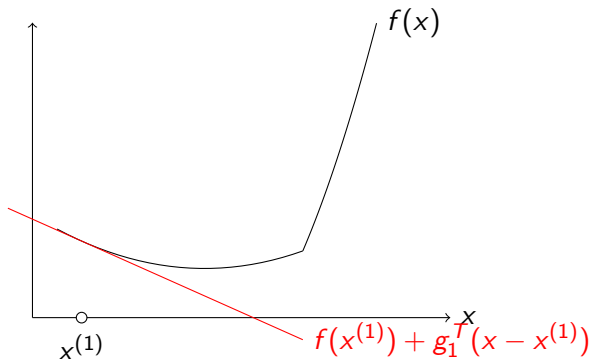is a global underestimator of $f$, i.e.

$$f(\mathbf{y}) \geq f(\mathbf{x}) + \mathbf{g}^T(\mathbf{y} - \mathbf{x}), \quad \text{for all } \mathbf{y} \in \operatorname{dom} f$$

# Subgradient

Given a function $f$ and a point $\mathbf{x} \in \operatorname{dom} f$,
$\mathbf{g} \in \mathbb{R}^N$ is called a **subgradient** of $f$ at $\mathbf{x}$ if:
the hypersurface with slopes $\mathbf{g}$ through $(\mathbf{x}, f(\mathbf{x}))$
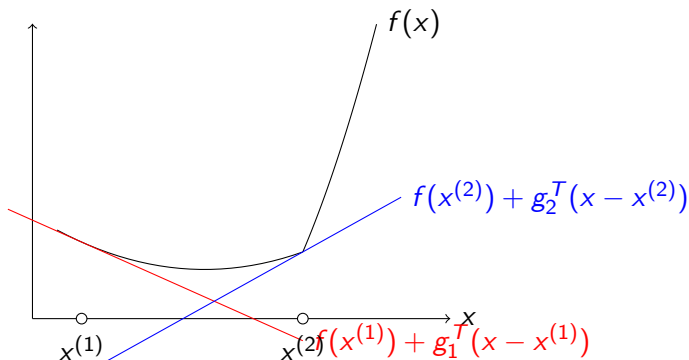is a global underestimator of $f$, i.e.

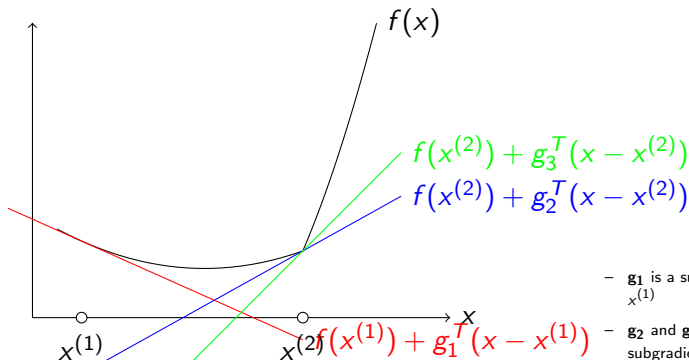$$f(\mathbf{y}) \geq f(\mathbf{x}) + \mathbf{g}^T(\mathbf{y} - \mathbf{x}), \quad \text{for all } \mathbf{y} \in \operatorname{dom} f$$



$f(x)$

$f(x^{(2)}) + g_3^T(x - x^{(2)})$

$f(x^{(2)}) + g_2^T(x - x^{(2)})$

$f(x^{(1)}) + g_1^T(x - x^{(1)})$

$x^{(1)}$

$x^{(2)}$

- $\mathbf{g_1}$ is a subgradient of $f$ at $x^{(1)}$

- $\mathbf{g_2}$ and $\mathbf{g_3}$ are subgradients of $f$ at $x^{(2)}$

# Example

For $f : \mathbb{R} \to \mathbb{R}$ and $f(x) = |x|$:

- ▶ For $x \neq 0$ there is one subgradient: $g = \nabla f(x) = \mathrm{sign}(x)$
- ▶ For $x = 0$ the subgradients are: $g \in [-1, 1]$

# Example

For $f : \mathbb{R} \to \mathbb{R}$ and $f(x) = |x|$:

▶ For $x \neq 0$ there is one subgradient: $g = \nabla f(x) = \text{sign}(x)$

▶ For $x = 0$ the subgradients are: $g \in [-1, 1]$

# Example

For $f : \mathbb{R} \to \mathbb{R}$ and $f(x) = |x|$:

▶ For $x \neq 0$ there is one subgradient: $g = \nabla f(x) = \text{sign}(x)$

▶ For $x = 0$ the subgradients are: $g \in [-1, 1]$

# Subdifferential

**Subdifferential** $\partial f(\mathbf{x})$: set of all subgradients of $f$ at $\mathbf{x}$

$$\partial f(\mathbf{x}) := \{\mathbf{g} \in \mathbb{R}^N \mid f(\mathbf{y}) \geq f(\mathbf{x}) + g^T(\mathbf{y} - \mathbf{x}) \; \forall \mathbf{y} \in \text{dom } f\}$$

▶ the subdifferential $\partial f(\mathbf{x})$ is a convex set.
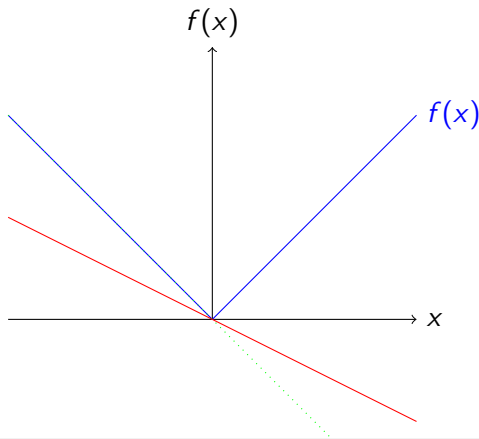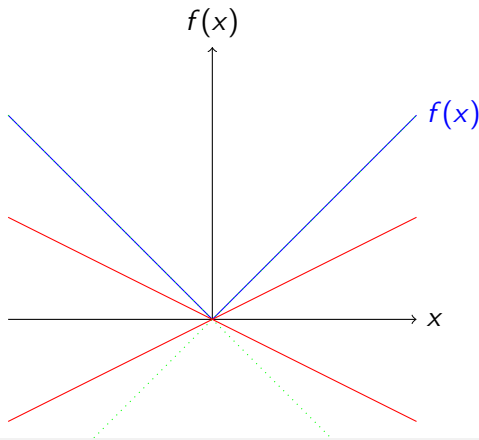
$$(\alpha\mathbf{g} + (1-\alpha)\mathbf{h})^T(\mathbf{y} - \mathbf{x}) = \alpha\mathbf{g}^T(\mathbf{y} - \mathbf{x}) + (1-\alpha)\mathbf{h}^T(\mathbf{y} - \mathbf{x})$$
$$\leq \alpha(f(\mathbf{y}) - f(\mathbf{x})) + (1-\alpha)(f(\mathbf{y}) - f(\mathbf{x}))$$
$$= f(\mathbf{y}) - f(\mathbf{x}) \quad \rightsquigarrow (\alpha\mathbf{g} + (1-\alpha)\mathbf{h}) \in \partial f(\mathbf{x})$$

▶ for a **convex** function $f$:
  ▶ subgradients always exist: $\partial f(\mathbf{x}) \neq \emptyset$

  ▶ $f$ is differentiable at $x$
     iff the subdifferential contains a single element (the gradient)

$$f \text{ differentiable at } x \iff \partial f(x) = \{\nabla f(x)\}$$

## Example

For $f(x) = |x|$:

▶ remember, $\partial f$ is set valued:

$$\partial f(x) = \{1\}, \quad \forall x > 0, \quad \partial f(x) = \{-1\}, \quad \forall x < 0, \quad \partial f(0) = [-1, +1]$$

# Subdifferential

For a **non-convex** function $f$:

- ▶ subgradients make less sense
    - ▶ see "generalized subgradients", defined on local information

# Outline

# Subgradient Calculus

▶ Assume $f$ convex and $\mathbf{x} \in \operatorname{dom} f$.

▶ Some algorithms require only **one** subgradient for optimizing nondifferentiable functions $f$.

▶ Other algorithms and optimality conditions require the **whole** subdifferential at $\mathbf{x}$.

▶ **Tools for finding subgradients:**
  ▶ **Weak subgradient calculus**: finding *one* subgradient $\mathbf{g} \in \partial f(\mathbf{x})$

  ▶ **Strong subgradient calculus**: finding the *whole* subdifferential $\partial f(\mathbf{x})$

# Subgradient Calculus

If $f$ is differentiable at $\mathbf{x}$: $\partial f(\mathbf{x}) = \{\nabla f(\mathbf{x})\}$.

Additional rules:

▶ **Scaling**: for $a > 0$: $\partial(a \cdot f) = a \cdot \partial f = \{a \cdot \mathbf{g} \mid \mathbf{g} \in \partial(f)\}$

▶ **Addition**: $\partial(f_1 + f_2) = \partial f_1 + \partial f_2$

▶ **Affine composition**: for $h(\mathbf{x}) = f(A\mathbf{x} + \mathbf{b})$ then

$$\partial h(\mathbf{x}) = A^T \partial f(A\mathbf{x} + \mathbf{b})$$

▶ **Finite pointwise maximum**: if $f(\mathbf{x}) = \max_{m=1\,...,M} f_m(\mathbf{x})$ then

$$\partial f(\mathbf{x}) = \text{conv}(\bigcup_{m:f_m(\mathbf{x})=f(\mathbf{x})} \partial f_m(\mathbf{x}))$$

The subdifferential is the convex hull of the union of subdifferentials of all **active functions** at $\mathbf{x}$.

# Subgradient Calculus / Pointwise Supremum

▶ **Pointwise Supremum**: if $f(\mathbf{x}) = \sup\limits_{a \in A} f_a(\mathbf{x})$ then

$$\partial f(\mathbf{x}) \supseteq \text{conv}(\bigcup_{a \in A: f_a(\mathbf{x}) = f(\mathbf{x})} \partial f_a(\mathbf{x}))$$

  ▶ "=" if $A$ is compact and $f$ continuous in $x$ and $a$.

# Subgradient Calculus / Function Composition

▶ **Function Composition**: if $f(\mathbf{x}) = h(g_1(\mathbf{x}), g_2(\mathbf{x}), \ldots, g_M(\mathbf{x}))$, then

$$\partial f(\mathbf{x}) \supseteq \text{conv}\{(b_1, b_2, \ldots, b_M)a \mid b_m \in \partial g_m(x), m = 1 : M,$$
$$a \in (\partial h)(g_1(x), g_2(x), \ldots, g_M(x))\}$$

   ▶ chain rule

   ▶ for differentiable $g_m$ and $h$:
      ▶ $Dg(x) = (b_1, b_2, \ldots, b_M)^T$ Jacobi matrix of $g := (g_1, g_2, \ldots, g_M)$

      ▶ $(\nabla h)(g(x)) = a$ gradient of $h$ at $g(x)$

# Subgradients / More Examples

$$f(x) := \|x\|_2$$

$$\partial f(x) =$$

# Subgradients / More Examples

$$f(x) := ||x||_2$$

$$\partial f(x) = \begin{cases} \{\frac{x}{||x||_2}\}, & \text{if } x \neq 0_N \\ \{g \in \mathbb{R}^N \mid ||g||_2 \leq 1\}. & \text{if } x = 0_N \end{cases}$$

proof:

use $||x||_2 = \max\limits_{z:||z||_2 \leq 1} z^T x$

$$\text{``} \leq \text{''} : z := \frac{x}{||x||_2}, \quad \text{``} \geq \text{''} : z^T x \leq ||z||_2 ||x||_2 \text{ Cauchy-Schwarz}$$

$$\partial(||x||_2) = \partial(\max\limits_{z:||z||_2 \leq 1} z^T x)$$

$$= \text{conv} \bigcup_{z:||z||_2 \leq 1, z^T x \text{ max.}} \{z\}, \quad \text{for } x = 0$$

$$= \text{conv} \bigcup_{z:||z||_2 \leq 1} \{z\} = \{z \in \mathbb{R}^N \mid ||z||_2 \leq 1\}$$

# Outline

# Descent Direction

- idea:
    - choose an arbitrary subgradient $g \in \partial f$
    - use its negative $-g$ as next direction

- negative subgradients are in general no descent directions
    - example:

$$f(x_1, x_2) := |x_1| + 3|x_2|$$

negative subgradients at $x := \begin{pmatrix} 1 \\ 0 \end{pmatrix}$:

# Descent Direction

- ▶ idea:
  - ▶ choose an arbitrary subgradient $g \in \partial f$
  - ▶ use its negative $-g$ as next direction

- ▶ negative subgradients are in general no descent directions
  - ▶ example:

$$f(x_1, x_2) := |x_1| + 3|x_2|$$

negative subgradients at $x := \begin{pmatrix} 1 \\ 0 \end{pmatrix}$:

# Descent Direction

- ▶ idea:
    - ▶ choose an arbitrary subgradient $g \in \partial f$
    - ▶ use its negative $-g$ as next direction

- ▶ negative subgradients are in general no descent directions
    - ▶ example:

$$f(x_1, x_2) := |x_1| + 3|x_2|$$

negative subgradients at $x := \begin{pmatrix} 1 \\ 0 \end{pmatrix}$:

# Descent Direction

- ▶ idea:
    - ▶ choose an arbitrary subgradient $g \in \partial f$
    - ▶ use its negative $-g$ as next direction

- ▶ negative subgradients are in general no descent directions
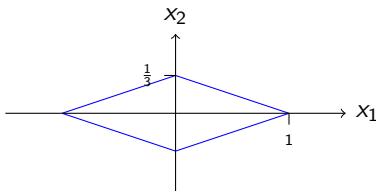    - ▶ example:

    $$f(x_1, x_2) := |x_1| + 3|x_2|$$

    negative subgradients at $x := \begin{pmatrix} 1 \\ 0 \end{pmatrix}$:
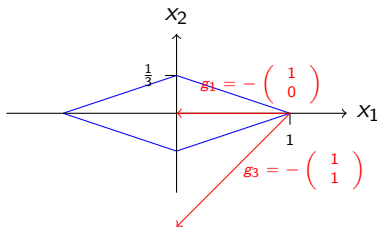
    $$-g_1 := - \begin{pmatrix} 1 \\ 0 \end{pmatrix} \quad \text{descent direction}$$

    $$-g_2 := - \begin{pmatrix} 1 \\ 3 \end{pmatrix} \quad \text{not a descent direction}$$

- ▶ thus cannot use stepsize controllers such as backtracking.

# Optimality Condition

For a convex $f : \mathbb{R}^N \to \mathbb{R}$:

$$\mathbf{x}^* \text{ is a global minimizer} \quad \Leftrightarrow \quad \mathbf{0} \text{ is a subgradient of } f \text{ at } \mathbf{x}^*$$
$$f(\mathbf{x}^*) = \min_{\mathbf{x} \in \text{dom } f} f(\mathbf{x}) \qquad\qquad \mathbf{0} \in \partial f(\mathbf{x}^*)$$

**Proof**:

If $\mathbf{0}$ is a subgradient of $f$ at $\mathbf{x}^*$, then for all $\mathbf{y} \in \mathbb{R}^N$:

$$f(\mathbf{y}) \geq f(\mathbf{x}^*) + \mathbf{0}^T(\mathbf{y} - \mathbf{x}^*)$$
$$f(\mathbf{y}) \geq f(\mathbf{x}^*)$$

# Gradient Descent (Review)

$$
\begin{aligned}
&1 \quad \textbf{min-gd}(f, \nabla f, x^{(0)}, \mu, \epsilon, K): \\
&2 \quad\quad \text{for } k := 1, \ldots, K: \\
&3 \quad\quad\quad \Delta x^{(k-1)} := -\nabla f(x^{(k-1)}) \\
&4 \quad\quad\quad \text{if } ||\nabla f(x^{(k-1)})||_2 < \epsilon: \\
&5 \quad\quad\quad\quad \text{return } x^{(k-1)} \\
&6 \quad\quad\quad \mu^{(k-1)} := \mu(f, x^{(k-1)}, \Delta x^{(k-1)}) \\
&7 \quad\quad\quad x^{(k)} := x^{(k-1)} + \mu^{(k-1)} \Delta x^{(k-1)} \\
&8 \quad\quad \text{return "not converged"}
\end{aligned}
$$

where
- ▶ $f$ objective function
- ▶ $\nabla f$ gradient of objective function $f$
- ▶ $x^{(0)}$ starting value
- ▶ $\mu$ step length controller
- ▶ $\epsilon$ convergence threshold for gradient norm
- ▶ $K$ maximal number of iterations

# Subgradient Method

1   **min-subgrad**$(f, \partial f, x^{(0)}, \mu, K)$ :

2    $x_{\text{best}}^{(0)} := x^{(0)}$

3    for $k := 1, \ldots, K$:

4     if $0 \in \partial f(x^{(k-1)})$:

5      return $x_{\text{best}}^{(k-1)}$

6     choose $g \in \partial f(x^{(k-1)})$ arbitrarily

7     $\Delta x^{(k-1)} := -g$

8     $\mu^{(k-1)} := \mu_{k-1}$

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

10     $x_{\text{best}}^{(k)} := \begin{cases} x^{(k)}, & \text{if } f(x^{(k)}) < f(x_{\text{best}}^{(k-1)}) \\ x_{\text{best}}^{(k-1)}, & \text{else} \end{cases}$

11    return "not converged"

where
- $\mu \in \mathbb{R}^*$ step length schedule

# Outline

# Slowly Diminishing Stepsizes

Proof of convergence requires **slowly diminishing stepsizes**:

$$\lim_{k \to \infty} \mu^{(k)} = 0, \quad \sum_{k=0}^{\infty} \mu^{(k)} = \infty, \quad \sum_{k=0}^{\infty} (\mu^{(k)})^2 < \infty$$

Q: Which of the following stepsizes are slowly diminishing?
- constant $\mu^{(k)} := \mu_0$
- $\mu^{(k)} := \frac{1}{k+1}$
- $\mu^{(k)} := \frac{1}{(k+1)^2}$

# Slowly Diminishing Stepsizes

Proof of convergence requires **slowly diminishing stepsizes**:

$$\lim_{k \to \infty} \mu^{(k)} = 0, \quad \sum_{k=0}^{\infty} \mu^{(k)} = \infty, \quad \sum_{k=0}^{\infty} (\mu^{(k)})^2 < \infty$$

for example:

$$\mu^{(k)} := \frac{1}{k+1}$$

but not:
- ▶ constant stepsizes $\mu^{(k)} := \mu \in \mathbb{R}$
- ▶ too fast shrinking stepsizes, e.g., $\mu^{(k)} := \frac{1}{(k+1)^2}$
- ▶ adaptive stepsize chosen by a step length controller

Theorem (convergence of subgradient method)

*Under the assumptions*

I. $f : X \to \mathbb{R}$ is convex, $X \subseteq \mathbb{R}^N$ is open

II. $f$ is Lipschitz-continuous with constant $G > 0$, i.e.
$$|f(\mathbf{x}) - f(\mathbf{y})| \leq G||\mathbf{x} - \mathbf{y}||_2, \quad \forall \mathbf{x}, \mathbf{y} \in \mathbb{R}^N$$

   ▶ *Equivalently:* $||\mathbf{g}||_2 \leq G$ for any subgradient $\mathbf{g}$ of $f$ at any $\mathbf{x}$

III. *slowly diminishing stepsizes* $\mu^{(k)}$, i.e.,
$$\lim_{k \to \infty} \mu^{(k)} = 0, \quad \sum_{k=0}^{\infty} \mu^{(k)} = \infty, \quad \sum_{k=0}^{\infty} (\mu^{(k)})^2 < \infty$$

*the subgradient method converges and*

$$f(\mathbf{x}_{best}^{(k)}) - f(\mathbf{x}^*) \leq \frac{||\mathbf{x}^{(0)} - \mathbf{x}^*||^2 + G^2 \sum_{j=0}^{k} (\mu^{(j)})^2}{2 \sum_{j=0}^{k} \mu^{(j)}}$$

# Convergence / Proof (1/2)

$$\|\mathbf{x}^{(k+1)} - \mathbf{x}^*\|_2^2$$
$$= \|\mathbf{x}^{(k)} - \mu^{(k)}\mathbf{g}^{(k)} - \mathbf{x}^*\|_2^2$$
$$= \|\mathbf{x}^{(k)} - \mathbf{x}^*\|_2^2 - 2\mu^{(k)}(\mathbf{g}^{(k)})^T(\mathbf{x}^{(k)} - \mathbf{x}^*) + (\mu^{(k)})^2\|\mathbf{g}^{(k)}\|_2^2$$
$$\underset{\text{SG}}{\leq} \|\mathbf{x}^{(k)} - \mathbf{x}^*\|_2^2 - 2\mu^{(k)}(f(\mathbf{x}^{(k)}) - f(\mathbf{x}^*)) + (\mu^{(k)})^2\|\mathbf{g}^{(k)}\|_2^2$$
$$\underset{\text{rec}}{\leq} \|\mathbf{x}^{(0)} - \mathbf{x}^*\|_2^2 - 2\sum_{j=0}^{k}\mu^{(j)}(f(\mathbf{x}^{(j)}) - f(\mathbf{x}^*)) + \sum_{j=0}^{k}(\mu^{(j)})^2\|\mathbf{g}^{(j)}\|_2^2$$
$$\underset{\text{II}}{\leq} \|\mathbf{x}^{(0)} - \mathbf{x}^*\|_2^2 - 2\sum_{j=0}^{k}\mu^{(j)}(f(\mathbf{x}^{(j)}) - f(\mathbf{x}^*)) + G^2\sum_{j=0}^{k}(\mu^{(j)})^2 \quad (1)$$

# Convergence / Proof (2/2)

$$
\begin{aligned}
f(\mathbf{x}_{\text{best}}^{(k)}) - f(\mathbf{x}^*) &= \frac{\sum_{j=0}^{k}(f(\mathbf{x}_{\text{best}}^{(k)}) - f(\mathbf{x}^*))\mu^{(j)}}{\sum_{j=0}^{k}\mu^{(j)}} \\
&\leq \frac{\sum_{j=0}^{k}(f(\mathbf{x}^{(j)}) - f(\mathbf{x}^*))\mu^{(j)}}{\sum_{j=0}^{k}\mu^{(j)}} \\
&\leq \frac{2\sum_{j=0}^{k}(f(\mathbf{x}^{(j)}) - f(\mathbf{x}^*))\mu^{(j)} + ||\mathbf{x}^{(k+1)} - \mathbf{x}^*||_2^2}{2\sum_{j=0}^{k}\mu^{(j)}} \\
&\underset{(1)}{\leq} \frac{||\mathbf{x}^{(0)} - \mathbf{x}^*||_2^2 + G^2\sum_{j=0}^{k}(\mu^{(j)})^2}{2\sum_{j=0}^{k}\mu^{(j)}}
\end{aligned}
$$

$$
\lim_{k\to\infty} f(\mathbf{x}_{\text{best}}^{(k)}) - f(\mathbf{x}^*) \leq \lim_{k\to\infty} \frac{||\mathbf{x}^{(0)} - \mathbf{x}^*||_2^2 + G^2\sum_{j=0}^{k}(\mu^{(j)})^2}{2\sum_{j=0}^{k}\mu^{(j)}} \underset{\text{III}}{=} 0
$$

# Summary

- **Subgradients** generalize gradients (for convex functions):
    - any slope of a hypersurface that is global underestimator.
    - at a differentiable location: the gradient is the only subgradient.

- Example **absolute value**: $\partial(|x|)|(0) = [-1, +1]$

- **subgradient calculus**:
    - scalar multiplication, addition, affine composition, pointwise maximum

- The **subgradient method** generalizes gradient descent:
    - use an arbitrary subgradient
    - stop if 0 is among the subgradients
    - as subgradients generally are no descent direction, the best location so far has to be tracked.

- The subgradient method is converging.
    - for Lipschitz-continuous functions and slowly diminishing stepsizes.

# Further Readings

- ▶ Subgradient methods are not covered by Boyd and Vandenberghe, 2004

- ▶ Subgradients:
  - ▶ Bertsekas, 1999, ch. B.5 and 6.1

- ▶ Subgradient methods:
  - ▶ Bertsekas, 1999, ch. 6.3.1

# References

Bertsekas, Dimitri P. (1999). *Nonlinear Programming*. Springer.

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

# Example: Text Classification

Features **A**: normalized word frequecies in text documents

Category **y**: topic of the text documents

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

$$\hat{y}_i = \sigma(\mathbf{x}^T \mathbf{a_i})$$

# Text Classification: L1-Regularized Logistic Regression

For $\mathbf{x} \in \mathbb{R}^N$, $\mathbf{y} \in \mathbb{R}^m$ and $\mathbf{A} \in \mathbb{R}^{m \times n}$ we have the following problem

$$\text{minimize} \quad -\sum_{i=1}^{m} y_i \log \sigma(\mathbf{x}^T \mathbf{a_i}) + (1 - y_i) \log(1 - \sigma(\mathbf{x}^T \mathbf{a_i})) + \lambda ||\mathbf{x}||_1$$

Which can be rewritten as:

$$\text{minimize} \quad -\sum_{i=1}^{m} y_i \log \sigma(\mathbf{x}^T \mathbf{a_i}) + (1 - y_i) \log(1 - \sigma(\mathbf{x}^T \mathbf{a_i})) + \lambda \sum_{k=1}^{N} |x_k|$$

$f$ is convex and non-smooth

# Example: L1-Regularized Logistic Regression

The subgradients of
$f(\mathbf{x}) = -\sum_{i=1}^{m} y_i \log \sigma(\mathbf{x}^T \mathbf{a_i}) + (1 - y_i) \log(1 - \sigma(\mathbf{x}^T \mathbf{a_i})) + \lambda \|\mathbf{x}\|_1$ are:

$$\mathbf{g} = -\mathbf{A}^T(\mathbf{y} - \hat{\mathbf{y}}) + \lambda \mathbf{s}$$

where $\mathbf{s} \in \partial \|\mathbf{x}\|_1$, i.e.:

▶ $s_k = \text{sign}(\mathbf{x}_k)$ if $\mathbf{x}_k \neq 0$

▶ $s_k \in [-1, 1]$ if $\mathbf{x}_k = 0$

# Example - The algorithm

For $\mathbf{x} \in \mathbb{R}^N$, $\mathbf{y} \in \mathbb{R}^m$ and $\mathbf{A} \in \mathbb{R}^{m \times n}$ we have the following the problem

$$\text{minimize} \quad -\sum_{i=1}^{m} y_i \log \sigma(\mathbf{x}^T \mathbf{a_i}) + (1 - y_i) \log(1 - \sigma(\mathbf{x}^T \mathbf{a_i})) + \lambda \sum_{k=1}^{N} |x_k|$$

1. Start with an initial solution $\mathbf{x}^{(0)}$
2. $t \leftarrow 0$
3. $f_{\text{best}} \leftarrow f(\mathbf{x}^{(0)})$
4. Repeat until convergence

    4.1 $\mathbf{x}^{(k+1)} \leftarrow \mathbf{x}^{(k)} - \mu^{(k)}(-\mathbf{A}^T(\mathbf{y} - \hat{\mathbf{y}}) + \lambda \mathbf{s})$
    4.2 $t \leftarrow t + 1$
    4.3 $f_{\text{best}} \leftarrow \min(f(\mathbf{x}^{(k)}), f_{\text{best}})$

5. Return $f_{\text{best}}$

where $\mathbf{s} \in \partial ||\mathbf{x}||_1$, i.e.:

- $s_k = \text{sign}(\mathbf{x}_k)$ if $\mathbf{x}_k \neq 0$

- $s_k \in [-1, 1]$ if $\mathbf{x}_k = 0$