



# Deep Learning

## Lecture 5: Non-convex Optimization for Learning (Part 2)

University of Agder,  
Kristiansand, Norway

Prepared by: Hadi Ghauch<sup>\*</sup>, Hossein S. Ghadikolaei<sup>†</sup>

<sup>\*</sup> Telecom ParisTech

<sup>†</sup> Royal Institute of Technology, KTH

<https://sites.google.com/view/fundl/home>

# Outline

1. Recap
2. Optimization problems with structure
  - Coordinate Descent methods
  - Block Coordinate Descent Methods
  - Block Successive Upperbound Minimization
  - Applications
3. Supplements

# Outline

## 1. Recap

## 2. Optimization problems with structure

Coordinate Descent methods

Block Coordinate Descent Methods

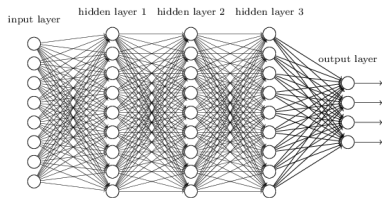
Block Successive Upperbound Minimization

Applications

## 3. Supplements

# Motivation

Resurgence of AI is due to **Deep Neural Networks (DNNs)**  
and countless variants (covered later)



DNNs is a composition of **non-linear layers**,  $\mathbf{W}_1, \dots, \mathbf{W}_J$ :

$$\mathbf{y}_i = \sigma_J(\mathbf{W}_J \cdots \sigma_2(\mathbf{W}_2 \sigma_1(\mathbf{W}_1 \mathbf{x}_i))), \quad i \in [N]$$

Let  $\mathbf{w} \in \mathbb{R}^d$  be the total number of weights in DNN (from all layers)  
 $d \geq 10^6$  in current deep learning application

The resulting optimization DNN training

$$\min_{\mathbf{w} \in \mathcal{D}} \frac{1}{N} \sum_{i=1}^N f_i((\mathbf{x}_i, \mathbf{y}_i); \mathbf{w}) + \lambda \|\mathbf{w}\|_2^2 := f(\mathbf{w})$$

**non-convex optimization** problem

# Recap: Non-convex Optimization

$$\min_{\mathbf{w} \in \mathcal{D}} f(\mathbf{w}) \quad (1)$$

$f(\mathbf{w}) : \mathcal{D} \subseteq \mathbb{R}^d \rightarrow \mathbb{R}$  is non-convex

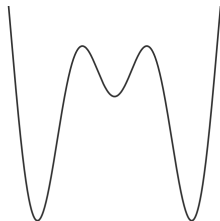
$\mathcal{D}$ : domain of  $f$  is convex

$f$  is Lipschitz continuous with constant

$L$ :  $\|f(\mathbf{w}_2) - f(\mathbf{w}_1)\|_2 \leq L\|\mathbf{w}_2 - \mathbf{w}_1\|_2$

$f$  is differentiable

less restrictive assumptions than convex



Local optimality **may not** necessarily imply global optimality

Proper initialization is very important in nonconvex optimization

First-order criteria ( $\nabla f(\mathbf{w}) = 0$ )

necessary and sufficient conditions for convex

only necessary condition for nonconvex

# Roadmap for Non-convex Optimization

**No generic method** for all non-convex problems ( $\neq$  convex case)  
solution approach depends on **structure of the optimization**

1. **No structure on  $f(w)$  nor constraint set (Recap from Lec 4):**  
first-order (GD and SCG) methods may be used.  
necessary conditions for GD/SGD to convergence to local min  
successive approximation: successively approx (1) with linear/convex bound
2.  **$f(w)$  is coordinate separable:** use **Coordinate Descent (CD)**  
if  $f(w) = f(w_1, \dots, w_d)$ , where  $(w_1, \dots, w_d)$  are **coordinates**  
 $f$  strongly convex in each coordinate,  $f$  non-convex jointly in all coordinates
3.  **$f(w)$  is block-coordinate separable:** use **Block-Coordinate Descent (BCD)**  
if  $f(w) = f(w_1, \dots, w_d)$ , where  $(w_1, \dots, w_d)$  **block of coordinates**  
 $f$  strongly convex in each block of coordinates,  
 $f$  non-convex jointly in all coordinates
4.  **$f(w)$  is block-coordinate separable**  
if  $f(w) = f(w_1, \dots, w_d)$ , where  $(w_1, \dots, w_d)$  **block of coordinates**  
 $f$  not necessarily convex in each block of coordinates,  
 $f$  non-convex jointly in all coordinates  
use **Block-Successive Upperbound minimization (BSUM)**

# Outline

1. Recap
2. Optimization problems with structure
  - Coordinate Descent methods
  - Block Coordinate Descent Methods
  - Block Successive Upperbound Minimization
  - Applications
3. Supplements

# Coordinate Descent (CD) Methods

Assumes optimization problem is **coordinate-separable**:

$$\min_{\mathbf{w} \in \mathcal{D}} f(\mathbf{w}) = f(w_1, \dots, w_d),$$

$w_i$  is the  $i$ th coordinate

domain **convex and separable**  $\mathcal{D} = \prod_{i=1}^d D_i$

**Idea:** Minimize coordinate  $w_i$ , while fixing all other ones  
at iteration  $k$ , optimization for  $w_i$

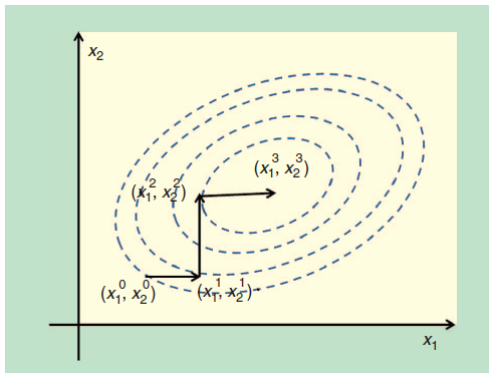
$$w_i^{k+1} := \arg \min_{w_i \in D_i} f(w_1^k, \dots, w_{i-1}^k, w_i, w_{i+1}^k, \dots, w_d^k) \quad (2)$$

assume that each **subproblem**, (2), **strongly convex problem**  
iteratively minimize  $d$  subproblems instead of joint problem

Update rule ? **cyclic update:**  $i = (k \bmod d), i \in \mathbb{N}_+$   
many other update rules: **random, greedy**



# Coordinate Descent Methods: Example



- descent direction along each coordinate (unlike GD)
- $f$  convex in  $x_1$  and  $x_2$  separately (not jointly)

# Outline

1. Recap
2. Optimization problems with structure
  - Coordinate Descent methods
  - Block Coordinate Descent Methods**
  - Block Successive Upperbound Minimization
  - Applications
3. Supplements

# Block Coordinate Descent Methods

**Block Coordinate Descent (BCD)** generalize CD:

- from minimizing coordinates to blocks of coordinates

$$\arg \min_{\mathbf{w}_1, \dots, \mathbf{w}_d} f(\mathbf{w}) = f(\mathbf{w}_1, \dots, \mathbf{w}_d)$$

$\mathbf{w}_i$  is the  $i$ th block of coordinates,

domain convex and separable:  $\mathcal{D} = \prod_{i=1}^d D_i$

**Low-rank MF:** Factorize large matrix  $\mathbf{X} \in \mathbb{R}^{M \times N}$  with low-rank matrices:  $\mathbf{P} \in \mathbb{R}^{M \times d}$ ,  $\mathbf{Q} \in \mathbb{R}^{N \times d}$ ,  $d \ll (M, N)$  such that  $\mathbf{X} \approx \mathbf{P}\mathbf{Q}^T$

$$\min_{\mathbf{P}, \mathbf{Q}} \|\mathbf{X} - \mathbf{P}\mathbf{Q}^T\|_F^2 \quad \text{s.t.} \quad \|\mathbf{P}\|_F^2 \leq \rho_1, \|\mathbf{Q}\|_F^2 \leq \rho_2$$

Not jointly convex in  $\mathbf{P}, \mathbf{Q}$ . But strongly convex in  $\mathbf{P}, \mathbf{Q}$  separately.  
constraints convex and separable

Similar optim problems for **training auto-encoders**, and **deep linear networks**

# Block Coordinate Descent Methods

**Idea:** Minimize block of coordinates,  $w_i$ , while fixing all other ones  
optimization for block  $w_i$ , at iteration  $k$

$$\begin{aligned} w_i^{k+1} &:= \arg \min_{w_i \in \mathcal{D}_i} f(w_1^k, \dots, w_{i-1}^k, w_i, w_{i+1}^k, \dots, w_d^k) \\ &= \arg \min_{w_i \in \mathcal{D}_i} f(w_i, w_{-i}^k) \end{aligned} \quad (3)$$

assume that (3) **strongly convex problem**

iteratively minimize  $d$  subproblems instead of joint problem

- $w_{-i}^k$  is the block of fixed variables at iteration  $k$ :  
 $w_{-i}^k := (w_1^k, \dots, w_{i-1}^k, w_{i+1}^k, \dots, w_d^k)$
- $f(w_i, w_{-i}^k)$  is the function when “looking” at block  $w_i$  only
- same type of update rules as CD

# Block Coordinate Descent Methods

## Convergence of BCD and CD

- A1) **solution for each subproblem(block) is unique**  
 $\Leftrightarrow f$  **strongly convex** in each of its blocks.
- A2)  $f$  smooth with Lipschitz constant  $L$
- A3) The domain,  $\mathcal{D}_i$ , is closed and convex

**Theorem 3:** Suppose that A1)-A3) hold. The sequence of updates generated by the CD method, in (2), and BCD method, in (3), monotonically converges to a **stationary point** of  $f$ , as  $k \rightarrow \infty$

A1) is too strong. May not be true in many problems

# Outline

1. Recap
2. Optimization problems with structure
  - Coordinate Descent methods
  - Block Coordinate Descent Methods
  - Block Successive Upperbound Minimization
  - Applications
3. Supplements

# Block Successive Upperbound Minimization

What happens when  $f$  is not strongly convex in each block ?

- **Block Successive Upperbound Minimization (BSUM)** extends convergence of BCD

Include many known ML algo as special cases:

- Expectation Minimization
- Convex Concave Procedure
- Non-negative Matrix Factorization
- Majorisation Minimization
- Forward Backward Splitting algorithm

# Block Successive Upperbound Minimization

**Idea:**  $f$  **not convex** in each block.

**intuition:** do SLA/SCA on each block of coordinates

Minimize a **strongly convex upperbound** for each block

Minimize block of coordinates,  $w_i$ , while fixing all other ones  
optimization for block  $w_i$ , at iteration  $k$

$$\begin{aligned}(\mathcal{P}_i) \quad w_i^{k+1} &:= \arg \min_{w_i \in \mathcal{D}_i} u_i(w_1^k, \dots, w_{i-1}^k, w_i, w_{i+1}^k, \dots, w_d^k) \\ &= \arg \min_{w_i \in \mathcal{D}_i} u_i(w_i, w_{-i}^k)\end{aligned}$$

$u_i(w_i, w_{-i}^k)$  is upperbound for block  $w_i$   
minimize upperbound for each subproblem (not  $f$ )

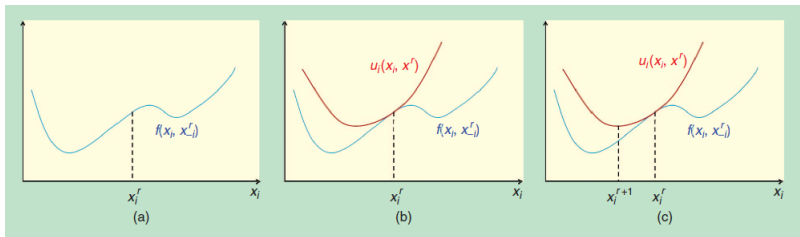
$w_{-i}^k$  is the block of fixed variables at iteration  $k$ :

$$w_{-i}^k := (w_1^k, \dots, w_{i-1}^k, w_{i+1}^k, \dots, w_d^k)$$

$u_i(w_i, w_{-i}^k)$  **strongly convex** in  $w_i$



# BSUM Example [Hong, 2016]



- a)  $f$  non-convex in coordinate  $x_i$
- b) construct upperbound fnc,  $u_i(\cdot)$ , at point  $x_i^r$
- c) find next iterate,  $x_i^{r+1}$ , by minimizing the upperbound  $u_i(\cdot)$  (not the function)

Hong et al “A Unified Algorithmic Framework for Block- Structured Optimization Involving Big Data”

# BSUM: Choosing the Upperbound

Several standard choices for upperbound,  $u_i(\mathbf{w}_i, \mathbf{w}_{-i}^k)$

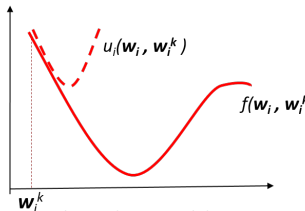
**first order proximal upperbound:**

$$u_i(\mathbf{w}_i, \mathbf{w}_{-i}^k) := f(\mathbf{w}_i^k) + \nabla f(\mathbf{w}_i^k)^T (\mathbf{w}_i - \mathbf{w}_i^k) + \gamma_i/2 \|\mathbf{w}_i - \mathbf{w}_i^k\|_2^2$$

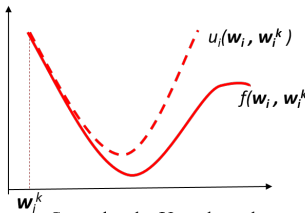
**second-order upperbound:**

$$u_i(\mathbf{w}_i, \mathbf{w}_{-i}^k) := f(\mathbf{w}_i^k) + \nabla f(\mathbf{w}_i^k)^T (\mathbf{w}_i - \mathbf{w}_i^k) + (1/2)(\mathbf{w}_i - \mathbf{w}_i^k)^T \nabla^2 f(\mathbf{w}_i^k) (\mathbf{w}_i - \mathbf{w}_i^k)$$

note similarity with SLA/SCA bounds



First-order Proximal  
Upperbound



Second-order Upperbound

# BSUM: Convergence

When and where does BSUM converge ?

Recall def of directional derivative and regular point.

- A1)  $u_i(-)$  is strongly convex, smooth, same directional derivatives as  $f$
- A2)  $f$  smooth with Lipschitz constant  $L$
- A3) The domain,  $\mathcal{D}_i$ , is closed and convex
- A4) The sequence of points,  $\{w_1^k, \dots, w_d^k\}_k$ , is regular

**Theorem 4:** Suppose that A1)-A4) hold. Assume that the solution of each subproblem,  $\mathcal{P}_i$ , is unique. Then, every limit point of  $\{w_1^k, \dots, w_d^k\}_k$  is a **stationary point** of  $f$ .

Convergence of BSUM is the most general among CD/BCD

## Variants of BSUM

**Parallel Successive Convex Approximation** [M. Razaviyayn, NIPS 2014]:

recent **variant of BSUM** to solve problems,

$$\arg \min_{\mathbf{w}_1, \dots, \mathbf{w}_d} f(\mathbf{w}_1, \dots, \mathbf{w}_d) + \sum_i g_i(\mathbf{w}_i)$$

$f$  is smooth and block-separable

$g$  non-smooth and strongly convex

several applications in **distributed learning**:

sparse dictionary learning, LASSO, K-SVD

current works apply this theoretical framework to show convergence

several DNN training algorithm: gradient methods, backprop, Newton, AdaGrad,

# Take-home Messages: BCD and BSUM

**CB, BCD and BSUM** ideal for optimization problems where

- $f$  is separable by (blocks of) coordinates
- $f$  not jointly convex in all blocks, but strongly convex in each block
- constraints are convex and separable
- strong convexity of each block **not a necessary cond** for convergence of BSUM

**Achilles' heel:** When does BCD and BSUM fail ?

- When the constraints are **not separable** (coupled constraints)
- Convergence to a stationary point **cannot be shown** for example:

$$\min_{x_1, x_2} x_1^2 + x_2^2 \quad \text{s. t.} \quad x_1 + x_2 = 1$$

# Outline

## 1. Recap

## 2. Optimization problems with structure

Coordinate Descent methods

Block Coordinate Descent Methods

Block Successive Upperbound Minimization

Applications

## 3. Supplements

# Application: Low rank Matrix Factorization

Factorize large matrix  $\mathbf{X} \in \mathbb{R}^{M \times N}$  with low-rank matrices:  
 $\mathbf{P} \in \mathbb{R}^{M \times d}$ ,  $\mathbf{Q} \in \mathbb{R}^{N \times d}$ ,  $d \ll (M, N)$  such that  $\mathbf{X} \approx \mathbf{P}\mathbf{Q}^T$

$$\min_{\mathbf{P}, \mathbf{Q}} \|\mathbf{X} - \mathbf{P}\mathbf{Q}^T\|_F^2 \quad \text{s. t.} \quad \|\mathbf{P}\|_F^2 \leq \rho_1, \quad \|\mathbf{Q}\|_F^2 \leq \rho_2$$

**not jointly convex in  $\mathbf{P}, \mathbf{Q}$ .** But strongly convex in  $\mathbf{P}, \mathbf{Q}$  separately.  
constraints convex and separable

Solved using BCD.

Given  $\mathbf{Q}_k$  optimize  $\mathbf{P}_{k+1} := \min_{\|\mathbf{P}\|_F^2 \leq \rho_1} \|\mathbf{X} - \mathbf{P}\mathbf{Q}_k^T\|_F^2$

Given  $\mathbf{P}_{k+1}$  update  $\mathbf{Q}_{k+1} := \min_{\|\mathbf{Q}\|_F^2 \leq \rho_2} \|\mathbf{X} - \mathbf{P}_{k+1}\mathbf{Q}^T\|_F^2$

subproblems are strongly convex: necessary cond for convergence of BCD hold  
a.k.a. Alternating Least Squares

Same BCD-based method used to train auto-encoder and deep linear nets.

## More Applications

### Training a deep linear network

Deep Linear Network is a composition of **linear layers**,  $\mathbf{W}_1, \dots, \mathbf{W}_J$  (no activation function):  $\mathbf{y}_i = \mathbf{W}_J \cdots \mathbf{W}_2 \mathbf{W}_1 \mathbf{x}_i$ ,  $i \in [N]$

Training may be done using BCD:

$$\frac{1}{N} \sum_{i=1}^N \|\mathbf{y}_i - \mathbf{W}_J \cdots \mathbf{W}_2 \mathbf{W}_1 \mathbf{x}_i\|_2^2 := f(\mathbf{W}_1, \dots, \mathbf{W}_J)$$

notice similarity with MF optimization

**Auto-encoder Training in DL** : efficiently done using BCD

we will discuss these DL architectures in Lec 8.



# Conclusions

Solution method depends on structure of problem

## **Non-convex optimization w/out structure:**

- first-order methods: with tweaks to escape saddle pts
- Successive linear/convex approximation

## **Non-convex optimization with structure:**

- If variable and constraint is separable
- if function convex in each block: CD, BCD
- if function non-convex in each block: BSUM
- BSUM is most generic framework
- applications of BSUM in ML

## Useful references

- M. Hong, M. Razaviyayn, Z. Q. Luo and J. S. Pang, "A Unified Algorithmic Framework for Block-Structured Optimization Involving Big Data: With applications in machine learning and signal processing" in IEEE Signal Processing Magazine, vol. 33, no. 1, pp. 57-77, Jan. 2016.
- M. Razaviyayn, M. Hong, and Z.-Q. Luo, "A unified convergence analysis of block successive minimization methods for nonsmooth optimization", SIAM J. Optim., vol. 23, no. 2, pp. 1126-1153, 2013
- M. Razaviyayn, "Successive convex approximation: Analysis and applications", Ph.D. dissertation, Univ. of Minnesota, 2014;
- Stephen J. Wright, "Coordinate descent algorithm" , Math. Program., vol. 151, no. 1, pp 3-34, 2015
- M. Razaviyayn, M Hong, ZQ Luo, JS Pang, "Parallel successive convex approximation for nonsmooth nonconvex optimization", NIPS 2014

# Outline

1. Recap
2. Optimization problems with structure
  - Coordinate Descent methods
  - Block Coordinate Descent Methods
  - Block Successive Upperbound Minimization
  - Applications
3. Supplements

# Special Cases

Many known algorithm are special cases of BSUM

## Difference of Convex (DC) programming:

$\arg \min_{\mathbf{w}} f(\mathbf{w}) = g_1(\mathbf{w}) - g_2(\mathbf{w})$ , where  $g_1$  and  $g_2$  convex

$$\mathbf{w}^{k+1} = \arg \min_{\mathbf{w}} g_1(\mathbf{w}) - (\nabla g_2(\mathbf{w}^k)^T (\mathbf{w} - \mathbf{w}^k)) - g_2(\mathbf{w}^k)$$

- *Convex Concave Procedure*

## BCD: special case of BSUM

Select the upperbound in BSUM as the function itself:

$$u_i(\mathbf{w}_i, \mathbf{w}_{-i}^k) = f(\mathbf{w}_i, \mathbf{w}_{-i}^k)$$

we recover BCD

# More Applications

## Sparse Dictionary Learning

Given a data matrix  $\mathbf{D} \in \mathbb{R}^{N \times M}$ , find a dictionary  $\mathbf{X}\mathbf{Y}^T$ , that sparsely represents the data matrix,

$$\min_{\mathbf{X}, \mathbf{Y}} \|\mathbf{D} - \mathbf{X}\mathbf{Y}^T\|_F^2 + \lambda \|\mathbf{X}\|_1 \quad \text{s. t.} \quad \|\mathbf{Y}\|_F \leq \beta$$



# Deep Learning

## Lecture 5: Non-convex Optimization for Learning (Part 2)

University of Agder,  
Kristiansand, Norway

Prepared by: Hadi Ghauch<sup>\*</sup>, Hossein S. Ghadikolaei<sup>†</sup>

<sup>\*</sup> Telecom ParisTech

<sup>†</sup> Royal Institute of Technology, KTH

<https://sites.google.com/view/fundl/home>