

# The Proximal Alternating Linearized Minimization (PALM) algorithm

[Bolte et al., 2014]

Minimization of a sum of finite functions

$$\text{minimize}_{x,y} \Psi(x, y) := f(x) + g(y) + H(x, y)$$

- $f(x)$ ,  $g(y)$  are proper, nonconvex, lower semicontinuous functions
- $H(x,y)$  smooth,  $C^1$ , function

In a convex setting, a Gauss Seidel-type or coordinate descent optimization

$$x^{k+1} \in \operatorname{argmin}_x \Psi(x, y^k)$$

$$y^{k+1} \in \operatorname{argmin}_y \Psi(x^{k+1}, y)$$

converges when the minimum in each step is uniquely attained

- no convergence in the nonconvex setting, difficult subproblems

In the non-convex setting, the following proximal regularization of the Gauss Seidel scheme provides a non-increasing sequence of points

$$x^{k+1} \in \operatorname{argmin}_x \left\{ \Psi(x, y^k) + \frac{c_k}{2} \|x - x^k\|^2 \right\}$$
$$y^{k+1} \in \operatorname{argmin}_y \left\{ \Psi(x^{k+1}, y) + \frac{d_k}{2} \|y - y^k\|^2 \right\}$$

- More well-posed subproblems

Two drawbacks here

- solving for the minimum per iteration is a difficult problem, due to nonconvexity and nonsmoothness
- accumulation of computation errors

Assume a cost function with smooth  $h(x)$  and nonsmooth  $\sigma(x)$ .

$$\Psi(x) = \sigma(x) + h(x)$$

A proximal forward-backward scheme is expressed as

$$x^{k+1} \in \text{prox}_t^\sigma \left( x^k - \frac{1}{t} \nabla h(x^k) \right)$$

$$x^{k+1} \in \underset{x \in \mathbb{R}^d}{\text{argmin}} \left\{ \langle x - x^k, \nabla h(x^k) \rangle + \frac{t}{2} \|x - x^k\|^2 + \sigma(x) \right\}, \quad (t > 0)$$

PALM's main idea: replace  $\Psi$  with its linearization

$$\widehat{\Psi}(x, y^k) = \langle x - x^k, \nabla_x H(x^k, y^k) \rangle + \frac{c_k}{2} \|x - x^k\|^2 + f(x), \quad (c_k > 0)$$

$$\widehat{\widehat{\Psi}}(x^{k+1}, y) = \langle y - y^k, \nabla_y H(x^{k+1}, y^k) \rangle + \frac{d_k}{2} \|y - y^k\|^2 + g(y), \quad (d_k > 0)$$

## PALM: Proximal Alternating Linearized Minimization

Initialization: start with any  $(x^0, y^0) \in \mathbb{R}^n \times \mathbb{R}^m$ .

For each  $k = 0, 1, \dots$  generate a sequence  $\{(x^k, y^k)\}_{k \in \mathbb{N}}$  as follows:

Take  $\gamma_1 > 1$ , set  $c_k = \gamma_1 L_1(y^k)$  and compute

$$x^{k+1} \in \text{prox}_{c_k}^f \left( x^k - \frac{1}{c_k} \nabla_x H(x^k, y^k) \right).$$

Take  $\gamma_2 > 1$ , set  $d_k = \gamma_2 L_2(x^{k+1})$  and compute

$$y^{k+1} \in \text{prox}_{d_k}^g \left( y^k - \frac{1}{d_k} \nabla_y H(x^{k+1}, y^k) \right).$$

- When there is no  $y$  term, PALM reduces to PFB
- The gradient of the smooth part  $H$  has to be globally Lipschitz continuous
- PALM converges to a critical point of  $\Psi$ .
- Convergence results for PALM exist for the case that the function  $\Psi$  to be minimized satisfies the so-called Kurdyka-Lojasiewicz (KL) property (sharpness of the gradient near critical points)

## Example: PALM for BSS

$$\min_{\mathbf{A}, \boldsymbol{\alpha}} \frac{1}{2} \|\mathbf{X} - \mathbf{A}\boldsymbol{\alpha}\Phi^*\|_F^2 + \lambda \|\boldsymbol{\alpha}\|_1 + \iota_{\mathcal{C}}(\mathbf{A})$$

$$\mathbf{S} = \boldsymbol{\alpha}\Phi^* \quad \boldsymbol{\alpha} \in \mathbb{R}^{N \times K}$$

[Feng & Kowalski, 2018]

---

### Algorithm 1: BSS-PALM

---

Initialization :  $\boldsymbol{\alpha}^{(1)} \in \mathbb{R}^{N \times K}$ ,  $\mathbf{A}^{(1)} \in \mathbb{R}^{M \times N}$ ,  $L^{1,(1)} = \|\mathbf{A}^{(1)}\|_2^2$ ,  $L^{2,(1)} = \|\boldsymbol{\alpha}^{(1)}\Phi^*\|_2^2$ ,  
 $j = 1$ ;

**repeat**

1.  $\nabla_{\boldsymbol{\alpha}} Q(\boldsymbol{\alpha}^{(j)}, \mathbf{A}^{(j)}) = -\mathbf{A}^{(j)T} (\mathbf{X} - \mathbf{A}^{(j)}\boldsymbol{\alpha}^{(j)}\Phi^*)\Phi$ ;
2.  $\boldsymbol{\alpha}^{(j+1)} = \mathcal{S}_{\lambda/L^{1,(j)}}(\boldsymbol{\alpha}^{(j)} - \frac{1}{L^{1,(j)}} \nabla_{\boldsymbol{\alpha}} Q(\boldsymbol{\alpha}^{(j)}, \mathbf{A}^{(j)}))$ ;
3.  $\nabla_{\mathbf{A}} Q(\boldsymbol{\alpha}^{(j+1)}, \mathbf{A}^{(j)}) = -(\mathbf{X} - \mathbf{A}^{(j)}\boldsymbol{\alpha}^{(j+1)}\Phi^*)\Phi\boldsymbol{\alpha}^{(j+1)H}$ ;
4.  $\mathbf{A}^{(j+1)} = \mathcal{P}_{\mathcal{C}}(\mathbf{A}^{(j)} - \frac{1}{L^{2,(j)}} \nabla_{\mathbf{A}} Q(\boldsymbol{\alpha}^{(j+1)}, \mathbf{A}^{(j)}))$ ;
5.  $L^{1,(j+1)} = \|\mathbf{A}^{(j+1)}\|_2^2$ ;
6.  $L^{2,(j+1)} = \|\boldsymbol{\alpha}^{(j+1)}\Phi^*\|_2^2$ ;
7.  $j = j + 1$ ;

**until** *convergence*;

---

# The stochastic asynchronous Proximal Alternating Linearized Minimization (PALM) algorithm

$$\underset{(x_1, \dots, x_m) \in \mathcal{H}_1 \times \dots \times \mathcal{H}_m}{\text{minimize}} \quad f(x_1, \dots, x_m) + \sum_{j=1}^m r_j(x_j)$$

- $f$  is a smooth,  $C^1$  function, it can be the data fidelity term
- $r$  is nonsmooth, it can be a structural regularizer of the solution

The paper's idea facilitates parallel computation:

- assign the computation of each  $x_i$  to a different processor

---

**Algorithm 1** SAPALM [Local view]

---

**Input:**  $x \in \mathcal{H}$

- 1: All processors in parallel do
  - 2: **loop**
  - 3:     Randomly select a coordinate block  $j \in \{1, \dots, m\}$
  - 4:     Read  $x$  from shared memory
  - 5:     Compute  $g = \nabla_j f(x) + \nu_j$
  - 6:     Choose stepsize  $\gamma_j \in \mathbb{R}_{++}$
  - 7:      $x_j \leftarrow \mathbf{prox}_{\gamma_j r_j}(x_j - \gamma_j g)$
- 

Main features:

- Inconsistent iterates. Other processors may write updates to  $x$  in the time required to read  $x$  from memory.
- Coordinate blocks. When the coordinate blocks  $x_j$  are low dimensional, it reduces the likelihood that one update will be immediately erased by another, simultaneous update.
- Noise. The noise  $\nu \in \mathcal{H}$  is a random variable that we use to model injected noise. It can be set to 0, or chosen to accelerate each iteration, or to avoid saddle points

---

**Algorithm 2** SAPALM [Global view]

---

**Input:**  $x^0 \in \mathcal{H}$

```
1: for  $k \in \mathbb{N}$  do
2:   Randomly select a coordinate block  $j_k \in \{1, \dots, m\}$ 
3:   Read  $x^{k-d_k} = (x_1^{k-d_k,1}, \dots, x_m^{k-d_k,m})$  from shared memory
4:   Compute  $g^k = \nabla_{j_k} f(x^{k-d_k}) + \nu_{j_k}^k$ 
5:   Choose stepsize  $\gamma_{j_k}^k \in \mathbb{R}_{++}$ 
6:   for  $j = 1, \dots, m$  do
7:     if  $j = j_k$  then
8:        $x_{j_k}^{k+1} \leftarrow \text{prox}_{\gamma_{j_k}^k r_{j_k}}(x_{j_k}^k - \gamma_{j_k}^k g^k)$ 
9:     else
10:       $x_j^{k+1} \leftarrow x_j^k$ 
```

---

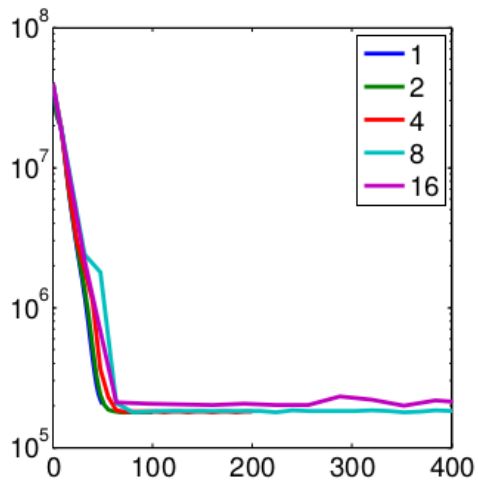
**Stochastic Gradients:** Noise due to stochastic approximations of delayed gradients.

- it allows us to prove convergence for a stochastic- or minibatch-gradient version of APALM, rather than requiring processors to compute a full (delayed) gradient.
- Stochastic gradients can be computed faster than their batch counterparts, allowing more frequent updates

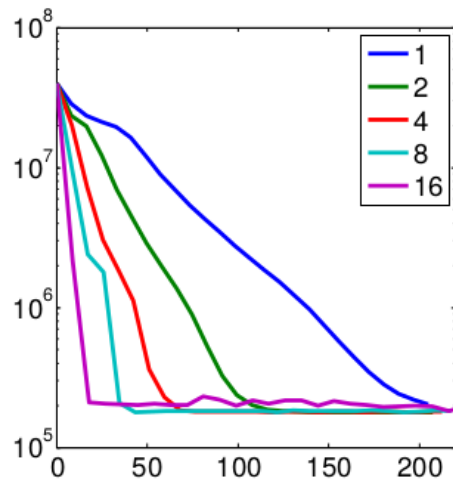
Convergence theorem proves that the SAPALM sequence is summable and  $\alpha$ -diminishing (expected error is below a threshold)



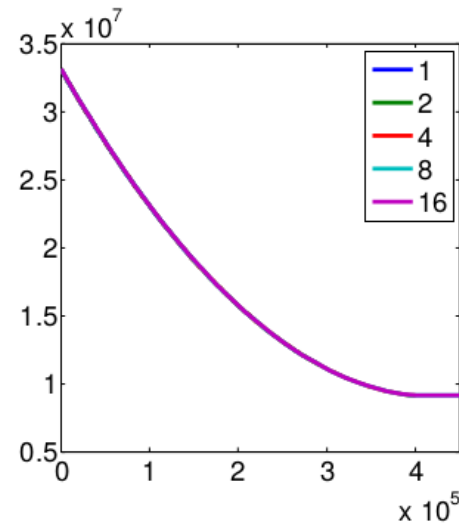
# Numerical experiments



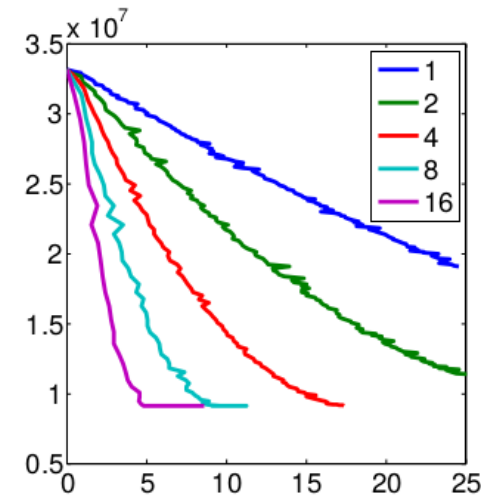
(a) Iterates vs objective



(b) Time (s) vs. objective



(c) Iterates vs. objective



(d) Time (s) vs. objective

- Sparse PCA

$$\operatorname{argmin}_{X,Y} \frac{1}{2} \|A - X^T Y\|_F^2 + \lambda \|X\|_1 + \lambda \|Y\|_1$$

- Quadratically Regularized Firm Thresholding PCA with Asynchronous Stochastic Gradients

$$\operatorname{argmin}_{X,Y} \frac{1}{2} \|A - X^T Y\|_F^2 + \lambda (\|X\|_{\text{Firm}} + \|Y\|_{\text{Firm}}) + \frac{\mu}{2} (\|X\|_F^2 + \|Y\|_F^2)$$



*That's all Folks!*

Proximal operator over  $\sigma$

$$\text{prox}_t^\sigma(x) := \operatorname{argmin} \left\{ \sigma(u) + \frac{t}{2} \|u - x\|^2 : u \in \mathbb{R}^d \right\}$$

Moreau proximal envelope associated to  $\sigma$

$$m^\sigma(x, t) := \inf \left\{ \sigma(u) + \frac{1}{2t} \|u - x\|^2 : u \in \mathbb{R}^d \right\}$$