

Randomized Proximal Algorithm  
with  
Automatic Dimension Reduction

Dmitry GRISHCHENKO  
**grishchenko.org**

joint work with  
F. IUTZELER, J. MALICK, M.-R. AMINI

Université Grenoble Alpes

Optimization Days 2018, Grenoble

# Distributed setup

one **master** machine

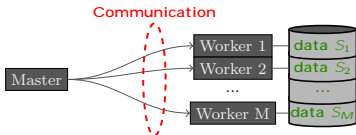
$M$  **worker** machines

data stored locally

on worker machines

communication cost

proportional to sending data size



# Distributed setup

one **master** machine

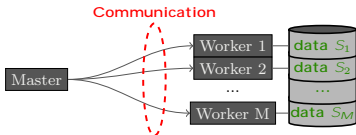
$M$  **worker** machines

data stored locally

on worker machines

communication cost

proportional to sending data size



# Distributed Learning

Global objective:

$$\min_{x \in \mathbb{R}^d} \frac{1}{m} \sum_{j=1}^m \ell_j(x) + g(x)$$

$m$  examples individual losses  $(\ell_j)$  empirical risk minimization regularizer  $g$

# Distributed Learning

Global objective:

$$\min_{x \in \mathbb{R}^d} \frac{1}{m} \sum_{j=1}^m \ell_j(x) + g(x)$$

$m$  examples individual losses  $(\ell_j)$  empirical risk minimization regularizer  $g$

Local data:

$$\min_{x \in \mathbb{R}^d} \sum_{i=1}^M f_i(x) + g(x)$$

$\{f_i\}$  convex, smooth  
 $g$  convex, nonsmooth

$M$  data blocks stored locally local function  $(f_i)$

$$f_i(x) = \frac{1}{|S_i|} \sum_{j \in S_i} \ell_j(x)$$

proportion  $|S_i| = m$  at  $i$

# Distributed Learning

Global objective:

$$\min_{x \in \mathbb{R}^d} \frac{1}{m} \sum_{j=1}^m \ell_j(x) + g(x)$$

$m$  examples individual losses  $(\ell_j)$  empirical risk  
minimization regularizer  $g$

Local data:

$$\min_{x \in \mathbb{R}^d} \sum_{i=1}^M f_i(x) + \begin{cases} g(x) \\ \ell(\{z\}) \end{cases}$$

convex, nonsmooth  
convex, smooth

$M$  data blocks stored locally local function  $(f_i)$   
 $f_i(x) = \frac{1}{|S_i|} \sum_{j \in S_i} \ell_j(x)$   
 proportion  $|S_i| = m$  at  $i$



## Review on Proximal Gradient

Problem:

$$\min_{x \in \mathbb{R}^n} f(x) + g(x);$$

$f(x)$  is differentiable,  $L$  smooth and  $\mu$  strongly convex

$g(x)$  is non-smooth but convex

## Review on Proximal Gradient

**Problem:**

$$\min_{x \in \mathbb{R}^n} f(x) + g(x);$$

$f(x)$  is differentiable,  $L$  smooth and  $\mu$  strongly convex

$g(x)$  is non-smooth but convex

**Algorithm:**

$$x^{k+1} = \operatorname{prox}_g(x^k - \tau \nabla f(x^k));$$

where *proximity operator* of  $g$

$$\operatorname{prox}_g(x) := \operatorname{argmin}_u g(u) + \frac{1}{2} \|x - u\|^2$$



## Review on Proximal Gradient

**Problem:**

$$\min_{x \in \mathbb{R}^n} f(x) + g(x);$$

$f(x)$  is differentiable,  $L$ -smooth and  $\mu$ -strongly convex  
 $g(x)$  is non-smooth but convex

**Algorithm:**

$$x^{k+1} = \text{prox}_g(x^k - \tau \nabla f(x^k));$$

where *proximity operator* of  $g$

$$\text{prox}_g(x) := \underset{u}{\text{argmin}} \quad g(u) + \frac{1}{2} \|x - u\|^2$$

**Convergence result:**

Let each  $f$  be  $L$ -smooth and  $\mu$ -strongly convex. Then, for  $\tau \in (0; 2/(L + \mu)]$ ,

$$\|x^k - x^*\| \leq (1 - \tau \mu)^k \|x^0 - x^*\| + \frac{\tau L}{2} \|x^0 - x^*\|^2;$$

for  $\tau = 2/(L + \mu)$

## Distributed Proximal Gradient

Problem:

$$\min_{x \in \mathbb{R}^d} \sum_{i=1}^n f_i(x) + g(x)$$

$\underbrace{\hspace{10em}}_{F(x)}$

Gradient property:

$$\nabla F(x) = \sum_{i=1}^n \nabla f_i(x)$$

## Distributed Proximal Gradient

Problem:

$$\min_{x \in \mathbb{R}^d} \sum_{i=1}^M f_i(x) + g(x)$$

$\underbrace{\qquad\qquad\qquad}_{F(x)}$

Gradient property:

$$\nabla F(x) = \sum_{i=1}^M \nabla f_i(x)$$

Algorithm: on each iteration:

Master gathering of the local variables

$$x^{k+1} = \text{P} \sum_{i=1}^M x_i^{k+1/2} = x^k - \gamma \nabla F(x)$$

Master performs a proximity operation

$$x_1^{k+1} = \dots = x_M^{k+1} = \text{prox}_g x^{k+1}$$

COMMUNICATION

Worker  $i$  update on local variable

$$x_i^{k+1/2} = x_i^k - \gamma \nabla f_i(x_i^k)$$

for all  $i = 1; \dots; M$

## Distributed Proximal Gradient

Problem:

$$\min_{x \in \mathbb{R}^d} \sum_{i=1}^M f_i(x) + g(x)$$

$\underbrace{\hspace{10em}}_{F(x)} \{z\}$

Gradient property:

$$\nabla F(x) = \sum_{i=1}^M \nabla f_i(x)$$

Algorithm: on each iteration:

Master gathering of the local variables

$$x^{k+1} = \text{P} \sum_{i=1}^M x_i^{k+1/2} = x^k \quad \nabla F(x)$$

Master performs a proximity operation

$$x_1^{k+1} = \dots = x_M^{k+1} = \text{prox}_g x^{k+1}$$

COMMUNICATION

Worker  $i$  update on local variable

$$x_i^{k+1/2} = x_i^k \quad \nabla f_i(x_i^k)$$

for all  $i = 1; \dots; M$

It's exactly proximal gradient descent

## Distributed Proximal Gradient

Problem:

$$\min_{x \in \mathbb{R}^d} \sum_{i=1}^M f_i(x) + g(x)$$

$\left\{ \begin{array}{l} \text{---} \\ F(x) \end{array} \right\}$

Gradient property:

$$\nabla F(x) = \sum_{i=1}^M \nabla f_i(x)$$

Algorithm: on each iteration:

Master gathering of the local variables

$$x^{k+1} = \frac{1}{M} \sum_{i=1}^M x_i^{k+1} = x^k \quad \nabla F(x)$$

Master performs a proximity operation

$$x_1^{k+1} = \dots = x_M^{k+1} = \text{prox}_g x^{k+1}$$

COMMUNICATION

Worker  $i$  update on local variable

$$x_i^{k+1} = x_i^k - \eta \nabla f_i(x_i^k)$$

for all  $i = 1; \dots; M$

It's exactly proximal gradient descent

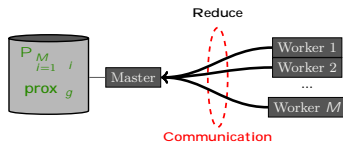
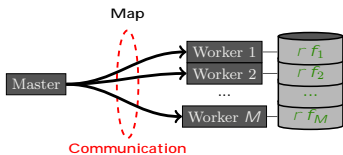
$k$  = number of master updates

Convergence rate:

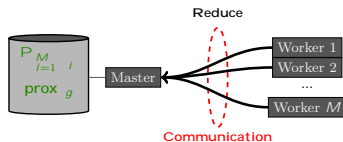
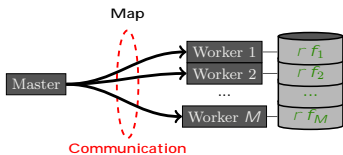
Let each  $f_i$  be  $L_i$ -smooth and  $\mu_i$ -strongly convex. Then, for  $\eta \in (0; 2/(L + \mu)]$  and  $L = \max_i L_i; \mu = \min_i \mu_i$ ,

$$\|x^k - x^*\| \leq \frac{1}{2} \left( \frac{L - \mu}{L + \mu} \right)^k \|x^0 - x^*\|$$

# Communication Problem



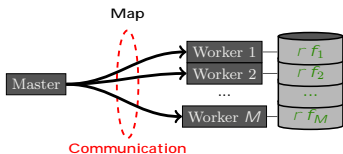
# Communication Problem



## Question:

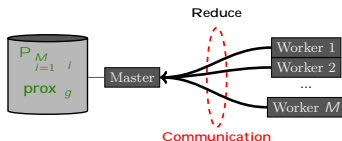
what if dimension  $d$  is extremely high?

# Communication Problem



## Question:

what if dimension  $d$  is extremely high?



## Answer:

sparsify data before sending!



# Identification

[Malick-Fadili-Peyré' 18]

Let  $(u^k)$  be a sequence converging to  $u^?$ , verifying

$$x^k := \underset{g}{\text{prox}}(u^k) \quad ! \quad x^?$$

where  $x^?$  is the unique minimizer of the  $\min_x \sum_{i=1}^M f_i(x) + g(x)$ .

Then, there is  $K < \infty$  such that:

$$g(x) = \|x\|_1.$$

$$\text{supp}(x^?) \subseteq \text{supp}(x^k) \subseteq \text{supp}(y_i^?) \quad \text{for all } k \leq K;$$

where  $\text{supp}(x) = \{i \in [1; n] \mid x_i \neq 0\}$

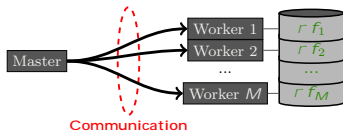
$$g(x) = \text{1-dimensional TV}(x) = \sum_{i=1}^{n-1} |x_{i+1} - x_i|$$

$$\text{jumps}(x^?) \subseteq \text{jumps}(x^k) \subseteq \text{jumps}(y_i^?) \quad \text{for all } k \leq K$$

where  $\text{jumps}(x) = \{i \in [1; n-1] \mid x_i \neq x_{i+1}\}$

where  $y_i^? = \underset{(1-\epsilon)g}{\text{prox}}(u^? - \epsilon x^?)$  for any  $\epsilon > 0$ .

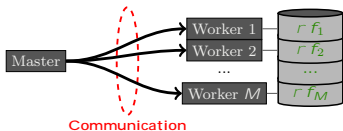
## Rightwards Sparsification



### QUESTION:

What identification gives to us?

## Rightwards Sparsification



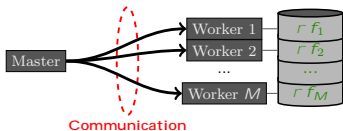
### QUESTION:

What identification gives to us?

### ANSWER:

For some regularizers proximal gradient points become sparse in some meaning:

## Rightwards Sparsification



### QUESTION:

What identification gives to us?

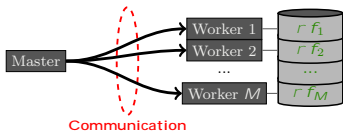
### ANSWER:

For some regularizers proximal gradient points become sparse in some meaning:

for  $\ell_1$  regularizer - coordinate sparsity (small amount of nonzero coordinates)

for  $\mathbf{TV}$  regularizer - block sparsity (small amount of jumps)

## Rightwards Sparsification



### QUESTION:

What identification gives to us?

### ANSWER:

For some regularizers proximal gradient points become sparse in some meaning:

for  $\ell_1$  regularizer - coordinate sparsity (small amount of nonzero coordinates)

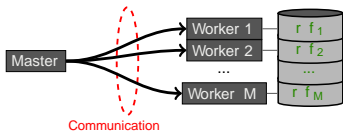
for **TV** regularizer - block sparsity (small amount of jumps)

### CONCLUSION:

master sends  $\text{prox}_g$  which is “sparse”

rightwards communications are “sparse”

## Rightwards Sparsification



### QUESTION:

What identification gives to us?

### ANSWER:

For some regularizers proximal gradient points become sparse in some meaning:

for  $\ell_1$  regularizer - coordinate sparsity (small amount of nonzero coordinates)

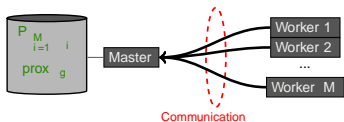
for TV regularizer - block sparsity (small amount of jumps)

### CONCLUSION:

master sends prox  $g$  which is "sparse"

rightwards communications are "sparse"

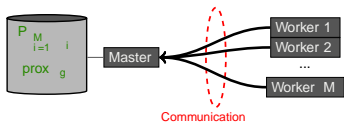
## Leftwards Sparsification



### Ideas of sparsification:

$\text{prox}_g x_i^k$  is not an option to send  $\{ P_{i=1}^M \}$   $\text{prox}_g x_i^k$  leads to nothing!  
 master knows  $x^k$  { we can send only gradient from slave!

## Leftwards Sparsification



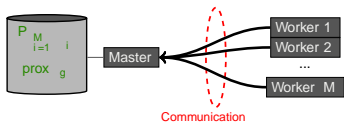
### Ideas of sparsification:

$\text{prox}_g x_i^k$  is not an option to send  $\{P_{i=1}^M\}$   $\text{prox}_g x_i^k$  leads to nothing!  
master knows  $x^k$  { we can send only gradient from slave!

**QUESTION:** How to sparsify gradient?



## Leftwards Sparsification



### Ideas of sparsification:

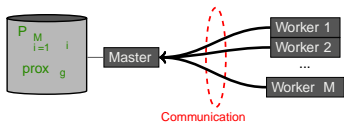
$\text{prox}_g x_i^k$  is not an option to send  $\{ \sum_{i=1}^M \text{prox}_g x_i^k \}$  leads to nothing!  
master knows  $x^k$  { we can send only gradient from slave!

**QUESTION:** How to sparsify gradient?

Option 1: [Tong Zhang' 17]

Use stochastic gradient against real one

## Leftwards Sparsification



### Ideas of sparsification:

$\text{prox}_g x_i^k$  is not an option to send  $\{P_{i=1}^M\}$   
 $\text{prox}_g x_i^k$  leads to nothing!  
master knows  $x^k$  { we can send only gradient from slave!

**QUESTION:** How to sparsify gradient?

Option 1: [Tong Zhang' 17]

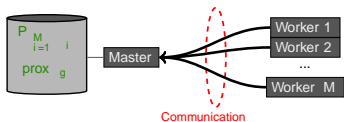
Use stochastic gradient against real one

**Drawback:**

decreasing stepsize

full gradient computation

## Leftwards Sparsification



### Ideas of sparsification:

$\text{prox}_g x_i^k$  is not an option to send  $\{ \text{prox}_g x_i^k \}_{i=1}^M$  leads to nothing!  
master knows  $x^k$  { we can send only gradient from slave!

**QUESTION:** How to sparsify gradient?

**Option I:**[Tong Zhang' 17]

Use stochastic gradient against real one

**Option II:**[Peter Richtárik' 16]

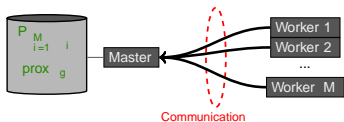
Use parallel coordinate descent

**Drawback:**

decreasing stepsize

full gradient computation

## Leftwards Sparsification



### Ideas of sparsification:

$\text{prox}_g x_i^k$  is not an option to send  $\{ \text{prox}_g x_i^k \}_{i=1}^P$  leads to nothing!  
master knows  $x^k$  { we can send only gradient from slave!

**QUESTION:** How to sparsify gradient?

**Option I:**[Tong Zhang' 17]

Use stochastic gradient against real one

**Option II:**[Peter Richtárik' 16]

Use parallel coordinate descent

**Drawback:**

decreasing stepsize

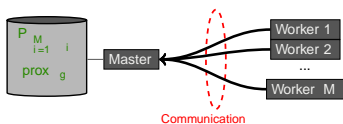
full gradient computation

**Drawback:**

block-separability

shared memory

## Leftwards Sparsification



### Ideas of sparsification:

$\text{prox}_g x_i^k$  is not an option to send  $\{ \text{prox}_g x_i^k \}$  leads to nothing!  
master knows  $x^k$  { we can send only gradient from slave!

**QUESTION:** How to sparsify gradient?

**Option I:**[Tong Zhang' 17]

Use stochastic gradient against real one

**Option II:**[Peter Richtárik' 16]

Use parallel coordinate descent

**Drawback:**

decreasing stepsize

full gradient computation

**Drawback:**

block-separability

shared memory

**Our option:** Use coordinate descent based algorithm taking into account sparsity structure of final solution

## Some Notations

### Projections:

Let  $P$  be a set of orthogonal projections  $\{P_i\}$  such that:

$P_i$  is linear operator

$$(\exists i : P_i(z^?) = P_i(y^?)) , \quad z^? = y^?$$

## Some Notations

### Projections:

Let  $\mathcal{P}$  be a set of orthogonal projections  $\{P_i\}$  such that:

$P_i$  is linear operator

$$(\exists i : P_i(z^?) = P_i(y^?)) , \quad z^? = y^?$$

### Expectation:

We select  $P \in \mathcal{P}$  random with the same probabilities

Let us denote by  $P = \mathbb{E}P$

Also let  $Q = P^{\frac{1}{2}}$

## Some Notations

### Projections:

Let  $\mathcal{P}$  be a set of orthogonal projections  $\{P_i\}$  such that:

$P_i$  is linear operator

$$(\exists i : P_i(z^?) = P_i(y^?)) , \quad z^? = y^?$$

### Expectation:

We select  $P \in \mathcal{P}$  random with the same probabilities

Let us denote by  $P = \mathbb{E}P$

Also let  $Q = P^{-\frac{1}{2}}$

### Examples:

Subspaces with sparsity equal to  $s$ :

$\ell_1$   $s$  dimensional subspace with  $\ell_1$  supp of size  $s$

$\text{TV}$   $s$  dimensional subspace with  $\ell_1$  jumps of size  $s - 1$



## Some Notations

### Projections:

Let  $\mathcal{P}$  be a set of orthogonal projections  $\{P_i\}$  such that:

$P_i$  is linear operator

$$(\exists i : P_i(z^?) = P_i(y^?)) , \quad z^? = y^?$$

### Expectation:

We select  $P \in \mathcal{P}$  random with the same probabilities

Let us denote by  $P = \mathbb{E}P$

Also let  $Q = P^{-\frac{1}{2}}$

### Examples:

Subspaces with sparsity equal to  $s$ :

$\ell_1$   $s$  dimensional subspace with  $\ell_1$  supp of size  $s$

$\text{TV}$   $s$  dimensional subspace with  $\ell_1$  jumps of size  $s - 1$

## Some Notations

### Projections:

Let  $\mathcal{P}$  be a set of orthogonal projections  $\{P_i\}$  such that:

$P_i$  is linear operator

$$(P_i : P_i(z) = P_i(y)) , \quad z = y$$

### Expectation:

We select  $P \in \mathcal{P}$  random with the same probabilities

Let us denote by  $P = E\mathcal{P}$

Also let  $Q = P^{-\frac{1}{2}}$

### Examples:

Subspaces with sparsity equal to  $s$ :

$\ell_1$   $s$  dimensional subspace with  $\text{card}(\text{supp})$  of size  $s$

TV  $s$  dimensional subspace with  $\text{card}(\text{jumps})$  of size  $s - 1$

Projections  $\mathcal{P}$ :

$\ell_1$  set of diagonal matrices with  $s$  ones and all other zeros

TV set of projections, each projection is block-diagonal matrix with  $s$  blocks; each blocks is fully filled with values equal to inverse of block's size

## Some Notations

### Projections:

Let  $\mathcal{P}$  be a set of orthogonal projections  $\{P_i\}$  such that:

$P_i$  is linear operator

$$(P_i : P_i(z) = P_i(y)) , \quad z = y$$

### Expectation:

We select  $P \in \mathcal{P}$  random with the same probabilities

Let us denote by  $P = E\mathcal{P}$

Also let  $Q = P^{-\frac{1}{2}}$

### Examples:

Subspaces with sparsity equal to  $s$ :

$\ell_1$   $s$  dimensional subspace with  $\text{card}(\text{supp})$  of size  $s$

TV  $s$  dimensional subspace with  $\text{card}(\text{jumps})$  of size  $s - 1$

Projections  $\mathcal{P}$ :

$\ell_1$  set of diagonal matrices with  $s$  ones and all other zeros

TV set of projections, each projection is block-diagonal matrix with  $s$  blocks; each blocks is fully filled with values equal to inverse of block's size

# Randomized Strata Descent

## Master Initialization

Initialize  $z^0$

Fix "measure of sparsity dimension", generate set  $P$  and calculate  $P; Q$

Compute  $x^0 = \text{prox}_g(Q^{-1} z^0)$

Randomly select  $P_0$  and send  $P_0; x^0; Q$  to workers

# Randomized Strata Descent

## Master Initialization

Initialize  $z^0$

Fix "measure of sparsity dimension", generate set

$P$  and calculate  $P; Q$

Compute  $x^0 = \text{prox}_g(Q^{-1} z^0)$

Randomly select  $P_0$  and send  $P_0; x^0; Q$  to workers

## Master

Initialize

for  $k=1, \dots$  do

Receive  $y_i^{k-1}$  from workers

$z^k = z^{k-1} - P_{k-1}(z^{k-1})$

$+ P_{k-1} Q^{-1} x^{k-1} + \sum_{i=1}^M y_i^{k-1}$

$x^k = \text{prox}_g(Q^{-1} z^k)$

Randomly select  $P_k$

Send  $x^k; P_k$  to workers

end for

C  
O  
M  
M  
U  
N  
I  
C  
A  
T  
I  
O  
N

## Worker i

for  $k=0, \dots$  do

Receive  $x^k; P_k$

$y_i^k =$

$P_k Q^{-1} \text{r f}_i(x^k)$

Send  $y_i^k$  to master

end for

# Randomized Strata Descent

## Master Initialization

Initialize  $z^0$   
 Fix "measure of sparsity dimension", generate set  $P$  and calculate  $Q$ ;  $Q$   
 Compute  $x^0 = \text{prox}_g(Q^{-1} z^0)$   
 Randomly select  $P_0$  and send  $P_0$ ;  $x^0$ ;  $Q$  to workers

## Master

Initialize

for  $k=1, \dots$  do

Receive  $y_i^{k-1}$  from workers

$z^k = z^{k-1} - P_{k-1}(z^{k-1})$

$+ P_{k-1} Q^{-1} x^{k-1} + \sum_{i=1}^M y_i^{k-1}$

$x^k = \text{prox}_g(Q^{-1} z^k)$

Randomly select  $P_k$

Send  $x^k$ ;  $P_k$  to workers

end for

C  
O  
M  
M  
U  
N  
I  
C  
A  
T  
I  
O  
N

## Worker i

for  $k=0, \dots$  do

Receive  $x^k$ ;  $P_k$

$y_i^k =$

$P_k Q^{-1} \text{r f}_i(x^k)$

Send  $y_i^k$  to master

end for

Is it "coordinate descent"?

yes because we use coordinate selection in gradient

no because we don't need regularizer to be separable

# Experiments for LASSO

Randomized Strata Descent

Synthetic LASSO problem

$$\min \frac{1}{2} \|Ax - b\|_2^2 + \lambda \|x\|_1$$

dimension  $d = 30$ ,  $\lambda = 0 : 1$

10 machines (1CPU, 1GB) in a cluster

Data divided uniformly

# Experiments for LASSO

Randomized Strata Descent

Synthetic LASSO problem

$$\min \frac{1}{2} \|Ax - b\|_2^2 + \lambda \|x\|_1$$

dimension  $d = 30$ ,  $\lambda = 0.1$

10 machines (1CPU, 1GB) in a cluster

Data divided uniformly



# Experiments for LASSO

Randomized Strata Descent

Synthetic LASSO problem

$$\min \frac{1}{2} \|Ax - b\|_2^2 + \lambda \|x\|_1$$

dimension  $d = 30$ ,  $\lambda = 0 : 1$

10 machines (1CPU, 1GB) in a cluster

Data divided uniformly

## Analysis

**positive** Amount of iterations almost proportional to amount of coordinates selected

**positive** Identification works as expected

**negative** There is no relation between mask recognition and algorithm speedup

# Experiments for Least Squares with 1-dTV Regularizer

Randomized Strata Descent

Synthetic Least Squares problem  
with 1-dTV regularizer

$$\min \frac{1}{2} \|Ax - b\|_2^2 + \lambda \sum_{i=1}^{d-1} |x_i - x_{i+1}|$$

dimension  $d = 30$ ,  $\lambda = 0.5$

10 machines (1CPU, 1GB) in a  
cluster

Data divided uniformly

# Experiments for Least Squares with 1-dTV Regularizer

Randomized Strata Descent

Synthetic Least Squares problem  
with 1-dTV regularizer

$$\min \frac{1}{2} \|Ax - b\|_2^2 + \lambda \sum_{i=1}^{d-1} |x_i - x_{i+1}|$$

dimension  $d = 30$ ,  $\lambda = 0.5$

10 machines (1CPU, 1GB) in a  
cluster

Data divided uniformly

# Experiments for Least Squares with 1-dTV Regularizer

Randomized Strata Descent

Synthetic Least Squares problem  
with 1-dTV regularizer

$$\min \frac{1}{2} \|kAx - bk\|_2^2 + \lambda \sum_{i=1}^d |x_i - x_{i+1}|$$

dimension  $d = 30$ ,  $\lambda = 0 : 5$

10 machines (1CPU, 1GB) in a  
cluster

Data divided uniformly

## Analysis

**positive** Identification works as expected

**negative** Extremely big amount iterations for  
sparsified versions, does not  
correlate even with jumps' amount

**negative** There is no relation between mask  
recognition and algorithm speedup

# Randomized Strata Descent with Automatic Dimension Reduction

Master

```

Initialize
for k=1,p+1,.. do
  calculate sparsity structure of  $x^k \{ S_k$ 
  if  $S_k \notin S_{k-p}$  then
    Generate new  $P; P; Q$ 
    w.r.t to  $S_k$  and s-extra
    Send  $S_k$  to slave
  end if
  for l=1,..,p do
    Receive  $y_i^{k+l-1}$  from workers


$$z^{k+l} = z^{k+l-1} - P_{k+l-1} (z^{k+l-1} + P_{k+l-1} Q^{-1} x^{k+l-1} + \sum_{i=1}^M y_i^{k+l-1})$$


 $x^k = \text{prox}_g(Q^{-1} z^k)$ 
    Randomly select  $P_k$ 
    Send  $x^k; P_k$  to workers
  end for
end for
  
```

Worker i

C  
O  
M  
M  
U  
N  
I  
C  
A  
T  
I  
O  
N

```

for k=0,.. do
  if  $S_k$  received then
    Generate new  $P; P; Q$ 
    w.r.t to  $S_k$ 
    and s-extra
  end if
  Receive  $x^k; P_k$ 
   $y_i^k = P_k Q^{-1} r_{f_i}(x^k)$ 
  Send  $y_i^k$  to master
end for
  
```

# Randomized Strata Descent with Automatic Dimension Reduction

```

Master
Initialize
for k=1,p+1,.. do
  calculate sparsity structure of  $x^k$  {  $S_k$ 
  if  $S_k \notin S_{k-p}$  then
    Generate new  $P; P; Q$ 
    w.r.t to  $S_k$  and s-extra
    Send  $S_k$  to slave
  end if
  for l=1,..,p do
    Receive  $y_i^{k+l-1}$  from workers

     $z^{k+l} = z^{k+l-1} - P_{k+l-1}(z^{k+l-1})$ 
     $+ P_{k+l-1} Q^{-1} x^{k+l-1} + \sum_{i=1}^M y_i^{k+l-1}$ 

     $x^k = \text{prox}_g(Q^{-1} z^k)$ 
    Randomly select  $P_k$ 
    Send  $x^k; P_k$  to workers
  end for
end for
  
```

```

Worker i
C
O
M
M
U
N
I
C
A
T
I
O
N

for k=0,.. do
  if  $S_k$  received then
    Generate new
     $P; P; Q$ 
    w.r.t to  $S_k$ 
    and s-extra
  end if
  Receive  $x^k; P_k$ 

   $y_i^k = P_k Q^{-1} r_{f_i}(x^k)$ 
  Send  $y_i^k$  to master
end for
  
```

Is it "coordinate descent"?

no because we use adapted coordinate selection in gradient

no because we don't need regularizer to be separable

# Experiments for Least Squares with 1-d TV Regularizer

Randomized Strata Descent with Automatic Dimension Reduction

Synthetic Least Squares problem  
with 1-d TV regularizer

$$\min \frac{1}{2} \|kAx - bk\|_2^2 + \lambda \sum_{i=1}^d |x_i - x_{i+1}|$$

dimension  $d = 30$ ,  $\lambda = 0.5$

10 machines (1CPU, 1GB) in a  
cluster

Data divided uniformly

# Experiments for Least Squares with 1-d TV Regularizer

Randomized Strata Descent with Automatic Dimension Reduction

Synthetic Least Squares problem  
with 1-d TV regularizer

$$\min \frac{1}{2} \|kAx - bk\|_2^2 + \lambda \sum_{i=1}^d |x_i - x_{i+1}|$$

dimension  $d = 30$ ,  $\lambda = 0.5$

10 machines (1CPU, 1GB) in a  
cluster

Data divided uniformly



# Experiments for Least Squares with 1-d TV Regularizer

Randomized Strata Descent with Automatic Dimension Reduction

Synthetic Least Squares problem  
with 1-d TV regularizer

$$\min \frac{1}{2} \|kAx - bk\|_2^2 + \lambda \sum_{i=1}^d |x_i - x_{i+1}|$$

dimension  $d = 30$ ,  $\lambda = 0.5$

10 machines (1CPU, 1GB) in a  
cluster

Data divided uniformly

## Analysis

**positive** Identification works as expected

**positive** Small amount of iterations

**positive** Mask recognition leads to fast  
convergence

# Convergence Rate

## Randomized Strata Descent with Automatic Dimension Reduction

### Theorem

Let each  $f_i$  be  $L_i$ -smooth and  $\mu_i$ -strongly convex. Then, for  $\alpha \in (0, 2/(L + \mu)]$ , and  $L = \max_i L_i$ ,  $\mu = \min_i \mu_i$

$$\mathbb{E} \|x^h - x^*\|_2^2 \leq \frac{2}{\mu} \frac{L}{L + \mu} \lambda_{\min}^{-1} \|x^0 - x^*\|_2^2$$

where  $\lambda_{\min}$  is minimal eigen value of  $\bar{P}$

# Convergence Rate

## Randomized Strata Descent with Automatic Dimension Reduction

### Theorem

Let each  $f_i$  be  $L_i$ -smooth and  $\mu_i$ -strongly convex. Then, for  $\alpha \in (0, 2/(L + \mu)]$ , and  $L = \max_i L_i$ ,  $\mu = \min_i \mu_i$

$$\mathbb{E} \|x^k - x^*\|_2^2 \leq \frac{2}{\mu + L} \lambda_{\min}(\bar{P})^k \|x^0 - x^*\|_2^2;$$

where  $\lambda_{\min}$  is minimal eigen value of  $\bar{P}$

Fixed stepsize same as in standard Proximal Gradient

# Convergence Rate

## Randomized Strata Descent with Automatic Dimension Reduction

### Theorem

Let each  $f_i$  be  $L_i$ -smooth and  $\mu_i$ -strongly convex. Then, for  $\alpha \in (0; 2/(L + \mu)]$ , and  $L = \max_i L_i$ ;  $\mu = \min_i \mu_i$

$$\mathbb{E} \|x^k - x^*\|_2^2 \leq \frac{2}{\mu} \frac{L}{L + \mu} \lambda_{\min}(\bar{P})^k \|x^0 - x^*\|_2^2;$$

where  $\lambda_{\min}$  is minimal eigen value of  $\bar{P}$

Fixed stepsize same as in standard Proximal Gradient

Example:  $\ell_1$  regularizer

$\lambda_{\min} = \rho_{\min}$ , where  $\rho_{\min}$  is minimal probability for coordinate to be chosen

# Convergence Rate

## Randomized Strata Descent with Automatic Dimension Reduction

### Theorem

Let each  $f_i$  be  $L_i$ -smooth and  $\mu_i$ -strongly convex. Then, for  $\alpha \in (0; 2/(L + \mu)]$ , and  $L = \max_i L_i$ ;  $\mu = \min_i \mu_i$

$$\mathbb{E} \|x^k - x^*\|_2^2 \leq \frac{2}{\mu + L} \sum_{i=1}^k \mu_i \|x^0 - x^*\|_2^2;$$

where  $\mu_{\min}$  is minimal eigen value of  $\bar{P}$

Fixed stepsize same as in standard Proximal Gradient

**Example:**  $\ell_1$  regularizer

$\mu = \mu_{\min}$ , where  $\mu_{\min}$  is minimal probability for coordinate to be chosen

$\text{prox}_g$  is separable

$\bar{Q}$  - diagonal matrix

# Convergence Rate

## Randomized Strata Descent with Automatic Dimension Reduction

### Theorem

Let each  $f_i$  be  $L_i$ -smooth and  $\mu_i$ -strongly convex. Then, for  $\alpha \in (0, 2/(L + \mu)]$ , and  $L = \max_i L_i$ ,  $\mu = \min_i \mu_i$

$$\mathbb{E} \|x^k - x^*\|_2^2 \leq \frac{2}{\mu} \frac{L}{L + \mu} \lambda_{\min}(\bar{P})^k \|x^0 - x^*\|_2^2;$$

where  $\lambda_{\min}$  is minimal eigen value of  $\bar{P}$

Fixed stepsize same as in standard Proximal Gradient

**Example:**  $\ell_1$  regularizer

$\mu = \mu_{\min}$ , where  $\mu_{\min}$  is minimal probability for coordinate to be chosen

$\text{prox}_g$  is separable

$\bar{Q}$  - diagonal matrix

$\bar{Q}$  could be skipped in the algorithm

## Conclusion

### Results

Algorithm with automatic dimension reduction

Importance of identification in sparsification

# Conclusion

## Results

Algorithm with automatic dimension reduction  
Importance of identification in sparsification

## Future plans

Asynchronous version  
Approximate computation of  $\bar{Q}$   
Scarse communications  
make less exchanges



## Conclusion

### Results

Algorithm with automatic dimension reduction  
Importance of identification in sparsification

### Future plans

Asynchronous version  
Approximate computation of  $\bar{Q}$   
Scarse communications  
make less exchanges

Thank you!