#### Incremental gradients, parallel methods (Optml++ Meeting 5)

#### Suvrit Sra

#### Massachusetts Institute of Technology

#### OPTML++, Fall 2015



## Outline

- Lect 1: Recap on convexity
- Lect 1: Recap on duality, optimality
- Lect 2: First-order optimization algorithms
- Lect 3: Operator splitting
- Lect 4: Stochastic and incremental methods
- Lect 5: Parallel / sparse-data methods



Optimization for ML and beyond: OPTML++



## **Stochastic gradient**

Method	Assumptions	Full	Stochastic
Subgradient	convex	$O(1/\sqrt{k})$	$O(1/\sqrt{k})$
Subgradient	strongly cvx	O(1/k)	O(1/k)

So using stochastic subgradient, solve *n* times faster.

Method	Assumptions	Full	Stochastic
Gradient	convex	O(1/k)	$O(1/\sqrt{k})$
Gradient	strongly cvx	$O((1-\mu/L)^k)$	O(1/k)

- For smooth problems, stochastic gradient needs more iterations

- Widely used in ML, rapid initial convergence

- Several speedup techniques studied, but worst case remains same

Method	Assumptions	Rate
Gradient	convex	O(1/k)
Gradient	strongly cvx	$O((1-\mu/L)^k)$
Stochastic	strongly cvx	O(1/k)
SAG	strongly convex	$O((1 - \min\{\frac{\mu}{n}, \frac{1}{8n}\})^k)$

This speedup also observed in practice

#### Complicated convergence analysis

Similar rates for many other methods

- stochastic dual coordinate (SDCA); [Shalev-Shwartz, Zhang, 2013]
- stochastic variance reduced gradient (SVRG); [Johnson, Zhang, 2013]
- proximal SVRG [Xiao, Zhang, 2014]
- hybrid of SAG and SVRG, SAGA (also proximal); [Defazio et al, 2014]
- accelerated versions [Lin, Mairal, Harchoui; 2015]
- asynchronous hybrid SVRG [Reddi et al. 2015]
- incremental Newton method, S2SGD and MS2GD, ...

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

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

Incremental gradient methods

$$x_{k+1} = x_k - \frac{\eta_k}{n} \nabla f_{i(k)}(x_k), \quad k \geq 0.$$

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

Incremental gradient methods

$$x_{k+1} = x_k - \frac{\eta_k}{n} \nabla f_{i(k)}(x_k), \quad k \geq 0.$$

View as gradient-descent with perturbed gradients

$$x_{k+1} = x_k - \frac{\eta_k}{n} (\nabla F(x_k) + \boldsymbol{e}_k)$$

Optimization for ML and beyond: OPTML++

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

Incremental gradient methods

$$x_{k+1} = x_k - \frac{\eta_k}{n} \nabla f_{i(k)}(x_k), \quad k \geq 0.$$

View as gradient-descent with perturbed gradients

$$x_{k+1} = x_k - \frac{\eta_k}{n} (\nabla F(x_k) + \boldsymbol{e_k})$$

► Perturbation slows down rate of convergence. Typically  $\eta_k = O(1/k)$ ; convergence rate also O(1/k) (sublinear).

Suvrit Sra (MIT)

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

Incremental gradient methods

$$x_{k+1} = x_k - \frac{\eta_k}{n} \nabla f_{i(k)}(x_k), \quad k \geq 0.$$

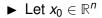
View as gradient-descent with perturbed gradients

$$x_{k+1} = x_k - \frac{\eta_k}{n} (\nabla F(x_k) + \boldsymbol{e}_k)$$

- ► Perturbation slows down rate of convergence. Typically  $\eta_k = O(1/k)$ ; convergence rate also O(1/k) (sublinear).
- Can we reduce impact of perturbation to speed up?

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

The incremental gradient method (IGM)



► For *k* ≥ 0

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

#### The incremental gradient method (IGM)

▶ Let 
$$x_0 \in \mathbb{R}^n$$

- For  $k \ge 0$ 
  - Pick  $i(k) \in \{1, 2, \dots, n\}$  uniformly at random

2 
$$x_{k+1} = x_k - \eta_k \nabla f_{i(k)}(x_k)$$

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

#### The incremental gradient method (IGM)

- ▶ Let  $x_0 \in \mathbb{R}^n$
- For  $k \ge 0$ 
  - 1 Pick  $i(k) \in \{1, 2, ..., n\}$  uniformly at random

2 
$$x_{k+1} = x_k - \eta_k \nabla f_{i(k)}(x_k)$$

 $g \equiv \nabla f_{i(k)}$  may be viewed as a stochastic gradient



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

#### The incremental gradient method (IGM)

- ▶ Let  $x_0 \in \mathbb{R}^n$
- For  $k \ge 0$ 
  - 1 Pick  $i(k) \in \{1, 2, ..., n\}$  uniformly at random 2  $x_{k+1} = x_k - \eta_k \nabla f_{i(k)}(x_k)$

 $g \equiv \nabla f_{i(k)}$  may be viewed as a stochastic gradient

$$g := g^{true} + e$$
, where *e* is mean-zero noise:  $\mathbb{E}[e] = 0$ 

- Index i(k) chosen uniformly from  $\{1, \ldots, n\}$
- ► Thus, in expectation:

 $\mathbb{E}[g] =$ 

- Index i(k) chosen uniformly from  $\{1, \ldots, n\}$
- ► Thus, in expectation:

 $\mathbb{E}[g] = \mathbb{E}_i[\nabla f_i(x)]$ 

- Index i(k) chosen uniformly from  $\{1, \ldots, n\}$
- ► Thus, in expectation:

$$\mathbb{E}[g] = \mathbb{E}_i[\nabla f_i(x)] = \sum_i \frac{1}{n} \nabla f_i(x) =$$

- Index i(k) chosen uniformly from  $\{1, \ldots, n\}$
- ► Thus, in expectation:

$$\mathbb{E}[g] = \mathbb{E}_i[\nabla f_i(x)] = \sum_{i=1}^{\infty} \frac{1}{n} \nabla f_i(x) = \nabla F(x)$$

- Alternatively,  $\mathbb{E}[g g^{true}] = \mathbb{E}[e] = 0$ .
- ▶ We call g an **unbiased estimate** of the gradient
- ► Here, we **obtained** *g* in a two step process:
  - **Sample:** pick an index i(k) unif. at random
  - Oracle: Compute a random gradient based on i(k)

- Index i(k) chosen uniformly from  $\{1, \ldots, n\}$
- ► Thus, in expectation:

$$\mathbb{E}[g] = \mathbb{E}_i[\nabla f_i(x)] = \sum_{i=1}^{\infty} \frac{1}{n} \nabla f_i(x) = \nabla F(x)$$

- Alternatively,  $\mathbb{E}[g g^{true}] = \mathbb{E}[e] = 0.$
- ▶ We call g an **unbiased estimate** of the gradient
- ► Here, we **obtained** *g* in a two step process:
  - **Sample:** pick an index i(k) unif. at random
  - Oracle: Compute a random gradient based on i(k)
- ▶ Individual  $g_k$  values can vary a lot
- ► Variance ( $\mathbb{E}[||g g^{true}||^2]$ ) influences convergence rate

Instead of using g<sub>k</sub> = ∇f<sub>i(k)</sub>(x<sub>k</sub>), correct it by using true gradient every m ≥ n steps (recall: F = <sup>1</sup>/<sub>n</sub> ∑<sup>n</sup><sub>i=1</sub> f<sub>i</sub>(x))

- Instead of using g<sub>k</sub> = ∇f<sub>i(k)</sub>(x<sub>k</sub>), correct it by using true gradient every m ≥ n steps (recall: F = <sup>1</sup>/<sub>n</sub> ∑<sup>n</sup><sub>i=1</sub> f<sub>i</sub>(x))
- ► Reduces variance of  $g_k(x_k, \xi_k)$ ; speeds up convergence

- Instead of using g<sub>k</sub> = ∇f<sub>i(k)</sub>(x<sub>k</sub>), correct it by using true gradient every m ≥ n steps (recall: F = <sup>1</sup>/<sub>n</sub> ∑<sup>n</sup><sub>i=1</sub> f<sub>i</sub>(x))
- ► Reduces variance of  $g_k(x_k, \xi_k)$ ; speeds up convergence

$$\nabla F(\bar{x}) = \frac{1}{m} \sum_{i} f_i(\bar{x})$$
  
$$x_{k+1} = x_k - \eta_k [\underbrace{\nabla f_{i(k)}(x_k) - \nabla f_{i(k)}(\bar{x}) + \nabla F(\bar{x})}_{g_k(x_k,\xi_k)}]$$

- Instead of using g<sub>k</sub> = ∇f<sub>i(k)</sub>(x<sub>k</sub>), correct it by using true gradient every m ≥ n steps (recall: F = <sup>1</sup>/<sub>n</sub> ∑<sup>n</sup><sub>i=1</sub> f<sub>i</sub>(x))
- ► Reduces variance of  $g_k(x_k, \xi_k)$ ; speeds up convergence

$$\nabla F(\bar{x}) = \frac{1}{m} \sum_{i} f_i(\bar{x})$$
  
$$x_{k+1} = x_k - \eta_k [\underbrace{\nabla f_{i(k)}(x_k) - \nabla f_{i(k)}(\bar{x}) + \nabla F(\bar{x})}_{g_k(x_k,\xi_k)}]$$

► Thus, with  $\xi_k = i(k)$ ,  $\mathbb{E}_{\xi}[g_k|x_k] = \nabla F(x_k)$ But with lower variance!

- Instead of using g<sub>k</sub> = ∇f<sub>i(k)</sub>(x<sub>k</sub>), correct it by using true gradient every m ≥ n steps (recall: F = <sup>1</sup>/<sub>n</sub> ∑<sup>n</sup><sub>i=1</sub> f<sub>i</sub>(x))
- ► Reduces variance of  $g_k(x_k, \xi_k)$ ; speeds up convergence

$$\nabla F(\bar{x}) = \frac{1}{m} \sum_{i} f_i(\bar{x})$$
  
$$x_{k+1} = x_k - \eta_k [\underbrace{\nabla f_{i(k)}(x_k) - \nabla f_{i(k)}(\bar{x}) + \nabla F(\bar{x})}_{g_k(x_k,\xi_k)}]$$

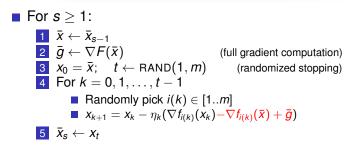
► Thus, with  $\xi_k = i(k)$ ,  $\mathbb{E}_{\xi}[g_k|x_k] = \nabla F(x_k)$ But with lower variance!

Say  $\bar{x}, x_k \to x^*$ . Then  $\nabla F(\bar{x}) \to 0$ . Thus, if  $\nabla f_i(\bar{x}) \to \nabla f_i(x^*)$ , then  $\nabla f_i(x_k) - \nabla f_i(\bar{x}) + \nabla F(\bar{x}) \to \nabla f_i(x_k) - \nabla f_i(x^*) \to 0$ .

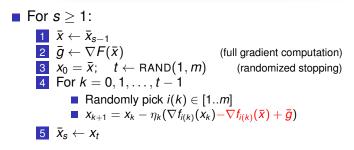
Suvrit Sra (MIT)

Optimization for ML and beyond: OPTML++

## **SVRG**



## **SVRG**



**Theorem** Assume each  $f_i(x)$  is smooth, and F(x) stronglyconvex. Then, for sufficiently large *n*, there is  $\alpha < 1$  s.t.

$$\mathbb{E}[F(\bar{\mathbf{x}}_{\mathbf{s}}) - F(\mathbf{x}^*)] \le \alpha^{\mathbf{s}}[F(\bar{\mathbf{x}}_0) - F(\mathbf{x}^*)]$$

Optimization for ML and beyond: OPTML++

# **Coordinate descent**

Optimization for ML and beyond: OPTML++



min 
$$f(x) = f(x_1, \ldots, x_K)$$
, where  $x_i \in \mathbb{R}^{n_i}$ 

#### min $f(x) = f(x_1, \ldots, x_K)$ , where $x_i \in \mathbb{R}^{n_i}$

Assumption: Gradient of block *i* is Lipschitz continuous

$$\|\nabla_i f(x+E_ih)-\nabla_i f(x)\|\leq L_i\|h\|$$

Block gradient  $\nabla_i f(x)$  is projection of full grad:  $E_i^T \nabla f(x)$ 

min 
$$f(x) = f(x_1, \ldots, x_K)$$
, where  $x_i \in \mathbb{R}^{n_i}$ 

Assumption: Gradient of block *i* is Lipschitz continuous

$$\|\nabla_i f(x+E_ih)-\nabla_i f(x)\| \leq L_i \|h\|$$

Block gradient  $\nabla_i f(x)$  is projection of full grad:  $E_i^T \nabla f(x)$ 

#### **Block Coordinate "Gradient" Descent**

• Using lemma:  $f(y) \le f(x) + \langle \nabla f(x), y - x \rangle + \frac{L}{2} ||y - x||^2$ , we get

$$f(x + E_i h) \leq f(x) + \langle \nabla_i f(x), h \rangle + \frac{L_i}{2} \|h\|^2$$
, for  $i = 1, \dots, n$ .

Optimization for ML and beyond: OPTML++

### min $f(x) = f(x_1, \ldots, x_K)$ , where $x_i \in \mathbb{R}^{n_i}$

Assumption: Gradient of block *i* is Lipschitz continuous

$$\|\nabla_i f(x+E_ih)-\nabla_i f(x)\| \leq L_i \|h\|$$

Block gradient  $\nabla_i f(x)$  is projection of full grad:  $E_i^T \nabla f(x)$ 

#### **Block Coordinate "Gradient" Descent**

► Using lemma:  $f(y) \le f(x) + \langle \nabla f(x), y - x \rangle + \frac{L}{2} ||y - x||^2$ , we get

$$f(x + E_i h) \leq f(x) + \langle \nabla_i f(x), h \rangle + \frac{L_i}{2} ||h||^2$$
, for  $i = 1, \dots, n$ .

- ► BCD algorithm:
  - 1 repeatedly go through blocks in "some" order
  - 2 minimize these upper bounds

For  $k \ge 0$  (no init. of x necessary)

- For  $k \ge 0$  (no init. of x necessary)
- Pick a block *i* from [*n*] with probability  $p_i > 0$

- For  $k \ge 0$  (no init. of x necessary)
- Pick a block *i* from [*n*] with probability  $p_i > 0$
- Optimize upper bound (partial gradient step) for block i

$$h = \underset{h}{\operatorname{argmin}} f(x_k) + \langle \nabla_i f(x_k), h \rangle + \frac{L_i}{2} \|h\|^2$$
$$h = -\frac{1}{L_i} \nabla_i f(x_k)$$

- For  $k \ge 0$  (no init. of x necessary)
- Pick a block *i* from [*n*] with probability  $p_i > 0$
- Optimize upper bound (partial gradient step) for block i

$$h = \underset{h}{\operatorname{argmin}} f(x_k) + \langle \nabla_i f(x_k), h \rangle + \frac{L_i}{2} \|h\|^2$$
$$h = -\frac{1}{L_i} \nabla_i f(x_k)$$

► Update the impacted coordinates of *x*, formally

- For  $k \ge 0$  (no init. of x necessary)
- Pick a block *i* from [*n*] with probability  $p_i > 0$
- Optimize upper bound (partial gradient step) for block i

$$h = \underset{h}{\operatorname{argmin}} f(x_k) + \langle \nabla_i f(x_k), h \rangle + \frac{L_i}{2} \|h\|^2$$
$$h = -\frac{1}{L_i} \nabla_i f(x_k)$$

► Update the impacted coordinates of *x*, formally

$$\begin{aligned} x_{k+1}^{(i)} \leftarrow x_k^{(i)} + h \\ x_{k+1} \leftarrow x_k - \frac{1}{L_i} E_i \nabla_i f(x_k) \end{aligned}$$

Optimization for ML and beyond: OPTML++

# **Randomized BCD**

- For  $k \ge 0$  (no init. of x necessary)
- Pick a block *i* from [*n*] with probability  $p_i > 0$
- Optimize upper bound (partial gradient step) for block i

$$h = \underset{h}{\operatorname{argmin}} f(x_k) + \langle \nabla_i f(x_k), h \rangle + \frac{L_i}{2} \|h\|^2$$
$$h = -\frac{1}{L_i} \nabla_i f(x_k)$$

► Update the impacted coordinates of *x*, formally

$$\begin{aligned} x_{k+1}^{(i)} \leftarrow x_k^{(i)} + h \\ x_{k+1} \leftarrow x_k - \frac{1}{L_i} E_i \nabla_i f(x_k) \end{aligned}$$

**Notice:** Original BCD had:  $x_k^{(i)} = \operatorname{argmin}_h f(\dots, \underbrace{h}_{block i}, \dots)$ We'll call this BCM (**Block Coordinate Minimization**)

Suvrit Sra (MIT)

Optimization for ML and beyond: OPTML++

$$h \leftarrow \operatorname{argmin}_{h} f(x_{k}) + \langle \nabla_{i} f(x_{k}), h \rangle + \frac{L_{i}}{2} \|h\|^{2}$$

$$h \leftarrow \operatorname{argmin}_{h} f(x_{k}) + \langle \nabla_{i} f(x_{k}), h \rangle + \frac{L_{i}}{2} \|h\|^{2}$$

#### **Descent:**

$$\begin{array}{rcl} x_{k+1} &=& x_k + E_i h \\ f(x_{k+1}) &\leq& f(x_k) + \langle \nabla_i f(x_k), h \rangle + \frac{L_i}{2} \|h\|^2 \end{array}$$

$$h \leftarrow \operatorname{argmin}_{h} f(x_{k}) + \langle \nabla_{i} f(x_{k}), h \rangle + \frac{L_{i}}{2} \|h\|^{2}$$

#### **Descent:**

$$\begin{array}{rcl} x_{k+1} &=& x_k + E_i h \\ f(x_{k+1}) &\leq& f(x_k) + \langle \nabla_i f(x_k), h \rangle + \frac{L_i}{2} \|h\|^2 \\ x_{k+1} &=& x_k - \frac{1}{L_i} E_i \nabla_i f(x_k) \end{array}$$

$$h \leftarrow \operatorname{argmin}_{h} f(x_{k}) + \langle \nabla_{i} f(x_{k}), h \rangle + \frac{L_{i}}{2} \|h\|^{2}$$

#### **Descent:**

$$\begin{array}{lll} x_{k+1} &=& x_k + E_i h \\ f(x_{k+1}) &\leq& f(x_k) + \langle \nabla_i f(x_k), h \rangle + \frac{L_i}{2} \|h\|^2 \\ x_{k+1} &=& x_k - \frac{1}{L_i} E_i \nabla_i f(x_k) \\ f(x_{k+1}) &\leq& f(x_k) - \frac{1}{L_i} \|\nabla_i f(x_k)\|^2 + \frac{L_i}{2} \left\| -\frac{1}{L_i} \nabla_i f(x_k) \right\|^2 \end{array}$$

Optimization for ML and beyond: OPTML++

$$h \leftarrow \operatorname{argmin}_{h} f(x_{k}) + \langle \nabla_{i} f(x_{k}), h \rangle + \frac{L_{i}}{2} \|h\|^{2}$$

#### **Descent:**

$$\begin{aligned} x_{k+1} &= x_k + E_i h \\ f(x_{k+1}) &\leq f(x_k) + \langle \nabla_i f(x_k), h \rangle + \frac{L_i}{2} \|h\|^2 \\ x_{k+1} &= x_k - \frac{1}{L_i} E_i \nabla_i f(x_k) \\ f(x_{k+1}) &\leq f(x_k) - \frac{1}{L_i} \|\nabla_i f(x_k)\|^2 + \frac{L_i}{2} \left\| -\frac{1}{L_i} \nabla_i f(x_k) \right\|^2 \\ f(x_{k+1}) &\leq f(x_k) - \frac{1}{2L_i} \|\nabla_i f(x_k)\|^2. \end{aligned}$$

$$f(x_k) - f(x_{k+1}) \ge \frac{1}{2L_i} \|\nabla_i f(x_k)\|^2$$

**Expected descent:** 

$$f(x_k) - \mathbb{E}[f(x_{k+1}|x_k)] = \sum_{i=1}^{a} p_i (f(x_k) - f(x_k - \frac{1}{L_i} E_i \nabla_i f(x_k)))$$

-1

#### **Expected descent:**

$$f(x_k) - \mathbb{E}[f(x_{k+1}|x_k)] = \sum_{i=1}^d p_i (f(x_k) - f(x_k - \frac{1}{L_i} E_i \nabla_i f(x_k)))$$
  
$$\geq \sum_{i=1}^d \frac{p_i}{2L_i} \|\nabla_i f(x_k)\|^2$$

#### **Expected descent:**

$$f(x_k) - \mathbb{E}[f(x_{k+1}|x_k)] = \sum_{i=1}^d p_i (f(x_k) - f(x_k - \frac{1}{L_i} E_i \nabla_i f(x_k)))$$
  
$$\geq \sum_{i=1}^d \frac{p_i}{2L_i} \|\nabla_i f(x_k)\|^2$$
  
$$= \frac{1}{2} \|\nabla f(x_k)\|^2_W \quad (\text{suitable } W).$$

٦

# Expected descent: $f(x_k) - \mathbb{E}[f(x_{k+1}|x_k)] = \sum_{i=1}^d p_i (f(x_k) - f(x_k - \frac{1}{L_i} E_i \nabla_i f(x_k)))$ $\geq \sum_{i=1}^d \frac{p_i}{2L_i} \|\nabla_i f(x_k)\|^2$ $= \frac{1}{2} \|\nabla f(x_k)\|^2_W \text{ (suitable W).}$

What is the expected descent with uniform probabilities?

Expected descent:  $f(x_k) - \mathbb{E}[f(x_{k+1}|x_k)] = \sum_{i=1}^d p_i (f(x_k) - f(x_k - \frac{1}{L_i} E_i \nabla_i f(x_k)))$   $\geq \sum_{i=1}^d \frac{p_i}{2L_i} \|\nabla_i f(x_k)\|^2$   $= \frac{1}{2} \|\nabla f(x_k)\|^2_W \text{ (suitable W).}$ 

What is the expected descent with uniform probabilities?

Descent + more notation + some work yields

$$O(\frac{d}{\epsilon}\sum_{i}L_{i}||x_{0}^{(i)}-x_{*}^{(i)}||^{2})$$

as the iteration complexity of obtaining  $\mathbb{E}[f(x_k)] - f^* \leq \epsilon$ 

Suvrit Sra (MIT)

Optimization for ML and beyond: OPTML++

# Exercise

- Recall Lasso problem:  $\min \frac{1}{2} \|Ax b\|^2 + \lambda \|x\|_1$
- ▶ Here  $x \in \mathbb{R}^d$ ; use *d* blocks
- Show what the Randomized BCD iterations look like
- Recall 1D prox operations for  $\lambda | \cdot |$  arise
- Try to implement it as efficiently as you can (do not copy or update vectors / coordinates unless necessary)

#### Exercise – pseudocode

Assuming *d* blocks, each update is scalar valued.

• Let 
$$x_0 = 0$$
;  $y_0 = Ax_0 - b = -b$ 

- For  $k \ge 0$ 
  - Pick random coordinate  $j \in [d]$
  - Compute  $\alpha \leftarrow \langle a_j, y \rangle$  i.e.,  $\nabla_j f(x_k)$
  - Min  $\alpha h + \frac{L_i}{2}h^2 + \lambda |h|$

$$h = \operatorname{prox}_{\lambda|\cdot|}(x_j - \frac{1}{L_j}\alpha)$$
$$h = \operatorname{sgn}(x_j - \frac{1}{L_j}\alpha) \max(|x_j - \frac{1}{L_j}\alpha| - \lambda, 0)$$

- Update:  $x_{k+1} = x_k + he_j$
- Update:  $y_{k+1} \leftarrow y_k + ha_j$

Suvrit Sra (MIT)

Optimization for ML and beyond: OPTML++

# **Parallel Optimization**

Optimization for ML and beyond: OPTML++



 Intuition: degree of separability strongly correlated with degree of parallelism possible

- Intuition: degree of separability strongly correlated with degree of parallelism possible
- ► Not insisting on exact computation allows more parallelism

- Intuition: degree of separability strongly correlated with degree of parallelism possible
- ► Not insisting on exact computation allows more parallelism
- Suppose *f* is the fraction of sequential computation. Then speedup for any number of processors (cores) is ≤ 1/*f*

- Intuition: degree of separability strongly correlated with degree of parallelism possible
- ► Not insisting on exact computation allows more parallelism
- Suppose *f* is the fraction of sequential computation. Then speedup for any number of processors (cores) is ≤ 1/*f*
- Parallel optimization on multi-core machines: shared memory architecture. Main penalty: synchronization / atomic operations

- Intuition: degree of separability strongly correlated with degree of parallelism possible
- ► Not insisting on exact computation allows more parallelism
- Suppose *f* is the fraction of sequential computation. Then speedup for any number of processors (cores) is ≤ 1/*f*
- Parallel optimization on multi-core machines: shared memory architecture. Main penalty: synchronization / atomic operations
- Distributed optimization across machines: synchronization and communication biggest burden;

- Intuition: degree of separability strongly correlated with degree of parallelism possible
- ► Not insisting on exact computation allows more parallelism
- Suppose *f* is the fraction of sequential computation. Then speedup for any number of processors (cores) is ≤ 1/*f*
- Parallel optimization on multi-core machines: shared memory architecture. Main penalty: synchronization / atomic operations
- Distributed optimization across machines: synchronization and communication biggest burden; node failure, network failure, load-balancing, etc.

- Intuition: degree of separability strongly correlated with degree of parallelism possible
- ► Not insisting on exact computation allows more parallelism
- Suppose *f* is the fraction of sequential computation. Then speedup for any number of processors (cores) is ≤ 1/*f*
- Parallel optimization on multi-core machines: shared memory architecture. Main penalty: synchronization / atomic operations
- Distributed optimization across machines: synchronization and communication biggest burden; node failure, network failure, load-balancing, etc.
- Synchronous vs. asynchronous computation

min 
$$f(x) := \sum_{i=1}^m f_i(x)$$
  $x \in \mathbb{R}^n$ .

min 
$$f(x) := \sum_{i=1}^m f_i(x)$$
  $x \in \mathbb{R}^n$ .

**Product space trick** 

min 
$$f(x) := \sum_{i=1}^m f_i(x)$$
  $x \in \mathbb{R}^n$ .

#### **Product space trick**

► Introduce (local) variables  $(x_1, ..., x_m)$ 

min 
$$f(x) := \sum_{i=1}^m f_i(x)$$
  $x \in \mathbb{R}^n$ .

#### **Product space trick**

- ► Introduce (local) variables  $(x_1, ..., x_m)$
- ▶ Problem is now over  $\mathcal{H}^m := \mathcal{H} \times \mathcal{H} \times \cdots \times \mathcal{H}$  (*m*-times)

min 
$$f(x) := \sum_{i=1}^m f_i(x)$$
  $x \in \mathbb{R}^n$ .

#### **Product space trick**

- ► Introduce (local) variables  $(x_1, ..., x_m)$
- ▶ Problem is now over  $\mathcal{H}^m := \mathcal{H} \times \mathcal{H} \times \cdots \times \mathcal{H}$  (*m*-times)
- **Consensus** constraint:  $x_1 = x_2 = \ldots = x_m$

$$\min_{\substack{(x_1,\ldots,x_m)\\ \text{s.t.}}} \sum_i f_i(x_i)$$
  
s.t.  $x_1 = x_2 = \cdots = x_m$ .

Optimization for ML and beyond: OPTML++

$$\min_{\boldsymbol{x}} f(\boldsymbol{x}) + \mathbb{1}_{\mathcal{B}}(\boldsymbol{x})$$
where  $\boldsymbol{x} \in \mathcal{H}^m$  and  $\mathcal{B} = \{ \boldsymbol{z} \in \mathcal{H}^m \mid \boldsymbol{z} = (x, x, \dots, x) \}$ 

$$\min_{\boldsymbol{x}} f(\boldsymbol{x}) + \mathbb{1}_{\mathcal{B}}(\boldsymbol{x})$$
where  $\boldsymbol{x} \in \mathcal{H}^m$  and  $\mathcal{B} = \{ \boldsymbol{z} \in \mathcal{H}^m \mid \boldsymbol{z} = (x, x, \dots, x) \}$ 

► Can solve using proximal splitting methods (e.g., DR, ADMM)

$$\min_{\mathbf{x}} f(\mathbf{x}) + \mathbb{1}_{\mathcal{B}}(\mathbf{x})$$

where  $\boldsymbol{x} \in \mathcal{H}^m$  and  $\mathcal{B} = \{ \boldsymbol{z} \in \mathcal{H}^m \mid \boldsymbol{z} = (x, x, \dots, x) \}$ 

- Can solve using proximal splitting methods (e.g., DR, ADMM)
- ► Each component of  $f_i(x_i)$  independently in parallel
- Communicate / synchronize to ensure consensus
- ► Asynchronous versions exist (results from 2014, 2015)

$$\min_{\mathbf{x}} f(\mathbf{x}) + \mathbb{1}_{\mathcal{B}}(\mathbf{x})$$

where  $\boldsymbol{x} \in \mathcal{H}^m$  and  $\mathcal{B} = \{ \boldsymbol{z} \in \mathcal{H}^m \mid \boldsymbol{z} = (x, x, \dots, x) \}$ 

- Can solve using proximal splitting methods (e.g., DR, ADMM)
- ► Each component of  $f_i(x_i)$  independently in parallel
- ► Communicate / synchronize to ensure consensus
- ► Asynchronous versions exist (results from 2014, 2015)
- ► Alternatively, compute dual and apply || BCD

#### Previously

$$\min f(x) = f(x_1,\ldots,x_d)$$

#### Previously

$$\min f(x) = f(x_1, \ldots, x_d)$$

#### What if?

min 
$$f(x) = \sum_i f_i(x_i)$$

#### Previously

$$\min f(x) = f(x_1, \ldots, x_d)$$

#### What if?

min 
$$f(x) = \sum_i f_i(x_i)$$

- ► Can solve all *d* problems independently in parallel
- ► In theory: *d* times speedup possible compared to serial case

#### Previously

$$\min f(x) = f(x_1, \ldots, x_d)$$

#### What if?

min 
$$f(x) = \sum_i f_i(x_i)$$

- ► Can solve all *d* problems **independently** in **parallel**
- ► In theory: *d* times speedup possible compared to serial case
- if objective "almost separable" we would still expect high speedup, governed by amount of separability
- ▶ Big data problems often have this "almost separable" structure!



# **Partial Separability**

Consider the sparse data matrix

$$\begin{pmatrix} d_{11} & d_{12} & & \\ & d_{22} & d_{23} & \\ & & \ddots & \ddots & \end{pmatrix} \in \mathbb{R}^{m \times n},$$

# Partial Separability

Consider the sparse data matrix

$$\begin{pmatrix} d_{11} & d_{12} & & \\ & d_{22} & d_{23} & \\ & & \ddots & \ddots \end{pmatrix} \in \mathbb{R}^{m \times n},$$

- Objective  $f(x) = \|Dx b\|_2^2 = \sum_{i=1}^m (d_i^T x b_i)^2$  also equals  $(d_{11}x_1 + d_{12}x_2 - b_1)^2 + (d_{22}x_2 + d_{23}x_3 - b_2)^2 + \cdots$
- Each term depends on only 2 coordinates
- ► Formally, we could write this as

$$f(x)=\sum_{J\in\mathcal{J}}f_J(x),$$

where  $\mathcal{J}=\left\{ \left\{ 1,2\right\} ,\left\{ 2,3\right\} ,\cdots\right\}$ 

► Key point:  $f_J(x)$  depends only on  $x_j$  for  $j \in J$ .

# **Partial Separability**

**Def.** Let  $\mathcal{J}$  be a collection of subsets of  $\{1, \ldots, d\}$ . We say *f* has **overlap degree**  $\omega$  if it can be written as

$$f(x) = \sum_{J \in \mathcal{J}} f_J(x),$$

where each  $f_J$  depends only on  $x_j$  for  $j \in J$ , and

$$|\boldsymbol{J}| \leq \boldsymbol{\omega} \quad \forall \boldsymbol{J} \in \mathcal{J}.$$

**Example:** If  $D_{m \times n}$  is a sparse matrix, then  $\omega = \max_{1 \le i \le m} \|d_i^T\|_0$ 

Optimization for ML and beyond: OPTML++

# **Partial Separability**

**Def.** Let  $\mathcal{J}$  be a collection of subsets of  $\{1, \ldots, d\}$ . We say f has **overlap degree**  $\omega$  if it can be written as

$$f(x) = \sum_{J \in \mathcal{J}} f_J(x),$$

where each  $f_J$  depends only on  $x_j$  for  $j \in J$ , and

$$|\boldsymbol{J}| \leq \boldsymbol{\omega} \quad \forall \boldsymbol{J} \in \mathcal{J}.$$

**Example:** If  $D_{m \times n}$  is a sparse matrix, then  $\omega = \max_{1 \le i \le m} ||d_i^T||_0$ **Exercise:** Extend this notion to  $x = (x^{(1)}, \dots, x^{(K)})$ *Hint:* Now,  $f_J$  will depend only on  $x^{(j)}$  for  $j \in J$ 

Optimization for ML and beyond: OPTML++

#### Data sparse ML problems

min 
$$\sum_{i=1}^{n} \ell(\langle a_i, x \rangle) + \lambda \|x\|^2$$

Training data samples  $a_1, \ldots, a_n$  are sparse

#### Data sparse ML problems

min 
$$\sum_{i=1}^{n} \ell(\langle a_i, x \rangle) + \lambda ||x||^2$$

Training data samples  $a_1, \ldots, a_n$  are sparse Rewrite the above problem in the format

min 
$$\sum_{i=1}^{n} \left( \ell(\langle a_i, x \rangle) + \lambda \sum_{u \in J_i} \frac{x_u^2}{d_u} \right),$$

where  $J_i$  are the nonzero coordinates of  $a_i$ ;  $d_u$  is the number of training samples nonzero in coordinate  $u \in [d]$ .

#### Data sparse ML problems

min 
$$\sum_{i=1}^{n} \ell(\langle a_i, x \rangle) + \lambda ||x||^2$$

Training data samples  $a_1, \ldots, a_n$  are sparse Rewrite the above problem in the format

min 
$$\sum_{i=1}^{n} \left( \ell(\langle a_i, x \rangle) + \lambda \sum_{u \in J_i} \frac{x_u^2}{d_u} \right),$$

where  $J_i$  are the nonzero coordinates of  $a_i$ ;  $d_u$  is the number of training samples nonzero in coordinate  $u \in [d]$ .

This is of the form (where  $\mathcal{J} \subset 2^{[d]}$ ),

min 
$$\sum_{J\in\mathcal{J}} f_J(x_J)$$

Degree of overlap  $\omega$ : maximum frequency any given feature

Suvrit Sra (MIT)

Optimization for ML and beyond: OPTML++

Each core runs the computation:

- **1** Sample coordinates J from  $\{1, \ldots, d\}$  (all sets of variables)
- 2 Read current state of  $x_J$  from shared memory
- 3 For each individual coordinate  $j \in J$

 $x_j \leftarrow x_j - \alpha_k [\nabla f_J(x_J)]_j$ 

Each core runs the computation:

- **1** Sample coordinates J from  $\{1, \ldots, d\}$  (all sets of variables)
- 2 Read current state of  $x_J$  from shared memory
- 3 For each individual coordinate  $j \in J$  $x_j \leftarrow x_j - \alpha_k [\nabla f_J(x_J)]_j$
- ► Atomic update only for  $x_j \leftarrow x_j a$  (not for gradient)

Each core runs the computation:

- **1** Sample coordinates J from  $\{1, \ldots, d\}$  (all sets of variables)
- 2 Read current state of  $x_J$  from shared memory
- 3 For each individual coordinate  $j \in J$  $x_j \leftarrow x_j - \alpha_k [\nabla f_J(x_J)]_j$
- ► Atomic update only for  $x_j \leftarrow x_j a$  (not for gradient)
- ► Since the actual coordinate *j* can arise in various *J*, processors can overwrite each others' work.

Each core runs the computation:

- **1** Sample coordinates J from  $\{1, \ldots, d\}$  (all sets of variables)
- 2 Read current state of  $x_J$  from shared memory
- **3** For each individual coordinate  $j \in J$  $x_j \leftarrow x_j - \alpha_k [\nabla f_J(x_J)]_j$
- ► Atomic update only for  $x_j \leftarrow x_j a$  (not for gradient)
- ► Since the actual coordinate *j* can arise in various *J*, processors can overwrite each others' work.
- But if partial overlaps, coordinate *j* does not appear in too many different subsets *J*, method works!
- Several related approaches exist in the literature

1 Choose initial point  $x_0 \in \mathbb{R}^d$ 

- 1 Choose initial point  $x_0 \in \mathbb{R}^d$
- 2 For k ≥ 0
  - Randomly pick (in parallel) a set of blocks  $S_k \subset \{1, \ldots, d\}$

1 Choose initial point  $x_0 \in \mathbb{R}^d$ 

**2** For  $k \ge 0$ 

- Randomly pick (in parallel) a set of blocks  $S_k \subset \{1, \ldots, d\}$
- Perform BCD updates (in parallel) for  $i \in S_k$

$$x_{k+1}^{(i)} \leftarrow x_k^{(i)} - \frac{1}{\beta w_i} \nabla_i f(x_k)$$

1 Choose initial point  $x_0 \in \mathbb{R}^d$ 

2 For k ≥ 0

- Randomly pick (in parallel) a set of blocks  $S_k \subset \{1, \ldots, d\}$
- Perform BCD updates (in parallel) for  $i \in S_k$

$$x_{k+1}^{(i)} \leftarrow x_k^{(i)} - \frac{1}{\beta w_i} \nabla_i f(x_k)$$

- Uniform sampling of blocks (or just coordinates)
- More careful sampling leads to better guarantees

1 Choose initial point  $x_0 \in \mathbb{R}^d$ 

2 For k ≥ 0

- Randomly pick (in parallel) a set of blocks  $S_k \subset \{1, \ldots, d\}$
- Perform BCD updates (in parallel) for  $i \in S_k$

$$x_{k+1}^{(i)} \leftarrow x_k^{(i)} - \frac{1}{\beta w_i} \nabla_i f(x_k)$$

- Uniform sampling of blocks (or just coordinates)
- More careful sampling leads to better guarantees
- Theory requires atomic updates

1 Choose initial point  $x_0 \in \mathbb{R}^d$ 

2 For k ≥ 0

- Randomly pick (in parallel) a set of blocks  $S_k \subset \{1, \ldots, d\}$
- Perform BCD updates (in parallel) for  $i \in S_k$

$$x_{k+1}^{(i)} \leftarrow x_k^{(i)} - \frac{1}{\beta w_i} \nabla_i f(x_k)$$

- Uniform sampling of blocks (or just coordinates)
- More careful sampling leads to better guarantees
- Theory requires atomic updates
- Implement asynchronously (use latest x<sup>(i)</sup> at each core)
- Theory of above method requires guaranteed descent
- Newer asynchronous CD methods also exist (see survey by Wright, 2015); e.g., methods that support inconsistent reads