



Optimization in Machine Learning: Lecture 9

Solving Large-scale Problems II

by Xiaolin Huang

xiaolinhuang@sjtu.edu.cn

SEIEE 2-429

Institute of Image Processing and Pattern Recognition

<http://www.pami.sjtu.edu.cn/>



上海交通大學

SHANGHAI JIAO TONG UNIVERSITY

1

Frank-Wolfe Algorithm

2

ADMM

3

Parallel Computing



1

Frank-Wolfe Algorithm

2

ADMM

3

Parallel Computing



Conditional/reduced gradient



$$\min_x f(x) \text{ s. t. } x \in D$$



- f is convex and contiguously differentiable
- D is a compact convex set
- projected descent method

$$x^{k+1} = \underline{P_D}(x^k - t_k \Delta x_{nt})$$

- “*projection are not always easy*”

Conditional/reduced gradient



$$\min_x f(x) \text{ s. t. } x \in D$$

- f is convex and contiguously differentiable
- D is a compact convex set
- Frank-Wolfe method
 - linearize the objective function
 - solve an LP (if D is polynomial)
 - determine the conditional/reduced gradient
 - determine the stepsize

there is no projection



Frank-Wolfe method

Algorithm 4 Fully-Corrective Variant, Re-Optimizing over all Previous Directions (with $\mathbf{s}^{(0)} := \mathbf{x}^{(0)}$)

... as Algorithm 2, except replacing line **b)** with
b') Update $\mathbf{x}^{(k+1)} := \arg \min_{\mathbf{x} \in \text{conv}(\mathbf{s}^{(0)}, \dots, \mathbf{s}^{(k+1)})} f(\mathbf{x})$

$$\min_x f(x) \text{ s.t. } x \in D$$



$$z^k = \operatorname{argmin}_z z^\top \nabla f(x^k) \text{ s.t. } z \in D$$



set the direction as $\Delta x^k = z^k - x^k$

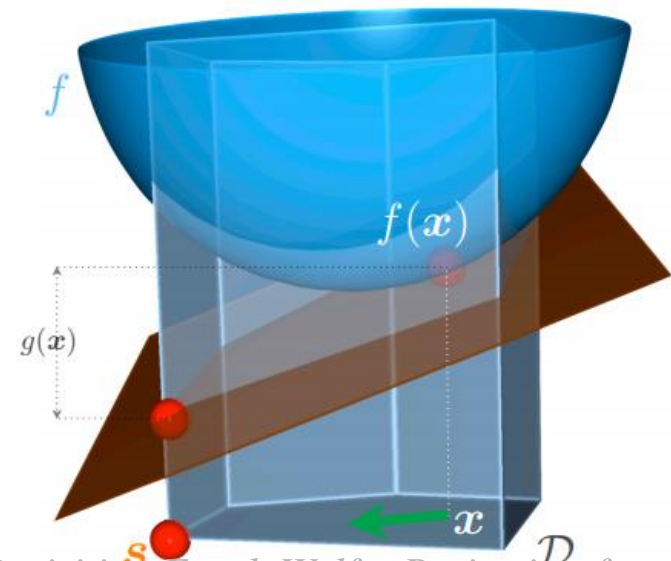


set $t^k = \frac{2}{k+1}$ or do line search

$$x^{k+1} = x^k - t^k \Delta x^k$$

x^k is kept in the feasible domain

there is no projection
 affine invariance



Frank-Wolfe method



$$\min_x f(x) \text{ s.t. } x \in D$$

- for example, when $D = \{x: \|x\| \leq t\}$

$$\begin{aligned} z^k &\in \operatorname{argmin}_z z^\top \nabla f(x^k) \text{ s.t. } \|z\| \leq t \\ &= -t \left(\operatorname{argmax}_{\|z\| \leq 1} z^\top \nabla f(x^k) \right) \\ &= -t \partial \|\nabla f(x^k)\|_* \end{aligned}$$

- this could be often simpler/cheaper than projection onto D

$$D = \{x: \|x\|_1 \leq t\} \quad \rightarrow \quad z^k \in -t \partial \|\nabla f(x^k)\|_\infty$$

$$D = \{x: \|x\|_p \leq t\} \quad \rightarrow \quad z_i^k = -t' \operatorname{sign}(\nabla f(x^k)) |\nabla f_i(x^k)|^{p/q}$$

Intuitive Explanation



$$z^k = \operatorname{argmin}_z z^\top \nabla f(x^k) \text{ s.t. } z \in D$$

- from the optimality of z^k , we have

$$\nabla f(x^k)^\top (z^k - x^k) \geq 0$$

- from convexity $f(y) \geq f(x^k) + \nabla f(x^k)^\top (y - x^k), \forall x^k, y \in D$

$$\begin{aligned} f(x^*) &\geq f(x^k) + \nabla f(x^k)^\top (x^* - x^k) \\ &\geq \min_{z \in D} \{f(x^k) + \nabla f(x^k)^\top (z - x^k)\} \\ &\geq f(x^k) - \nabla f(x^k)^\top x^k + \min_{z \in D} z^\top \nabla f(x^k) \end{aligned}$$

Convergence analysis



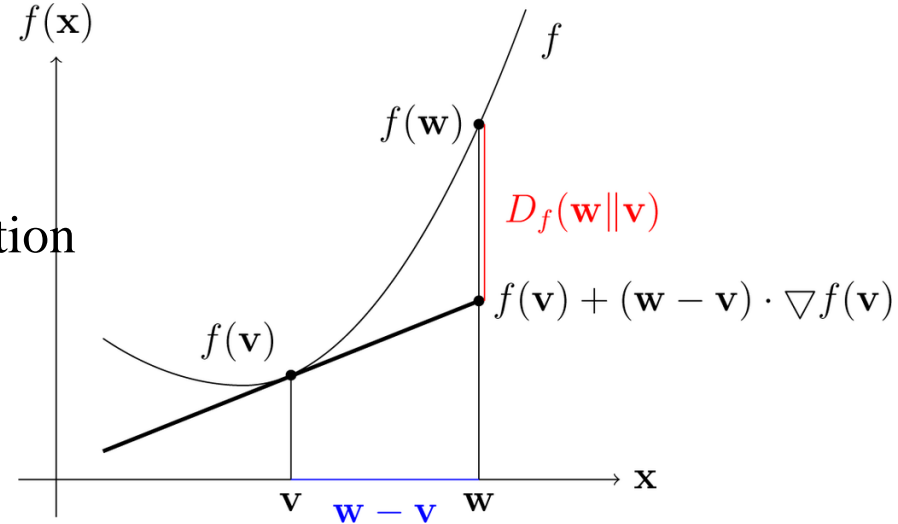
- curvature constant to characterize the *non-linearity* *Bregman divergence*

$$C_f \triangleq \sup_{\substack{x, s \in D \\ \gamma \in [0, 1] \\ y = x + \gamma(s - x)}} \frac{2}{\gamma^2} (f(y) - f(x) - (y - x)^\top \nabla f(x))$$

- C_f is closely related to Lipschitz condition

$$C_f \leq \text{diam}(D)^2 L$$

- additionally with convexity, we have



$$\begin{aligned} f(x^{k+1}) &= f(x^k + t^k(z^k - x^k)) \\ &\leq f(x^k) + t^k(z^k - x^k)^\top \nabla f(x) + \frac{\gamma^2}{2} C_f \end{aligned}$$

Convergence analysis



$$f(x^{k+1}) \leq f(x^k) + t^k(z^k - x^k)^\top \nabla f(x) + \frac{\gamma^2}{2} C_f$$

allowing inaccuracy solution

$$(z^k - x^k)^\top \nabla f(x) \leq \underbrace{\min_{z \in D} z^\top \nabla f(x^k) - x^k \nabla f(x^k)}_{\text{duality gap} - g(x^k)} + 2\delta C_f$$

duality gap $- g(x^k)$

$$f(x^{k+1}) \leq f(x^k) - t^k g(x^k) + \frac{\gamma^2 + \delta}{2} C_f$$

$$\min_x f(x) + I_D(x) \quad \text{indicator function}$$

$$\max_{\alpha} -f^*(\alpha) - I_D^*(-\alpha) \quad \text{supporting function}$$

$$f(x) + I_D(x) + f^*(\alpha) + I_D^*(-\alpha) \geq x^\top \alpha - I_D^*(-\alpha)$$

$$\downarrow \quad \text{at } x = x^k, \alpha = \nabla f(x^k)$$

$$x^k \nabla f(x^k) + \max_{z \in D} -z^\top \nabla f(x^k)$$

Convergence analysis



$$f(x^{k+1}) \leq f(x^k) - t^k g(x^k) + \frac{\gamma^2 + \delta}{2} C_f$$



$$\begin{aligned} f(x^{k+1}) - f^* &\leq f(x^k) - f^* - t^k g(x^k) + \frac{\gamma^2 + \delta}{2} C_f \\ &\leq f(x^k) - f^* - t^k (f(x^k) - f^*) + \frac{\gamma^2 + \delta}{2} C_f \\ &\leq (1 - t^k) (f(x^k) - f^*) + \frac{\gamma^2 + \delta}{2} C_f \end{aligned}$$

- sublinear linear convergence

Theorem: Frank-Wolfe using step sizes $\gamma_k = 2/(k+1)$, $k = 1, 2, 3, \dots$, and inaccuracy parameter $\delta \geq 0$, satisfies

$$f(x^{(k)}) - f^* \leq \frac{2M}{k+1} (1 + \delta)$$

1

Frank-Wolfe Algorithm

2

ADMM

3

Parallel Computing



LASSO



$$\min_x \lambda \|x\|_1 + \|b - Ax\|_2^2$$

- **properties**

- convex
- one part is non-smooth and coordinately separable
- one part is smooth and pointwise separable

- **decomposition into two parts**

$$\min_x \lambda \|x\|_1 + \|b - Az\|_2^2$$

$$\text{s. t. } x = z$$

Properties



■ dual ascent



Lagrangian: $L(x, z; \alpha) = \lambda \|x\|_1 + \|b - Az\|_2^2 + \underline{\alpha^\top (x - z)}$

dual function: $g(\alpha) = \inf_{x, z} L(x, z; \alpha)$

dual problem: $\alpha^* = \operatorname{argmax} g(\alpha)$

primal-dual: $(x^*, z^*) = \operatorname{argmin}_{x, z} L(x, z; \alpha^*)$

■ gradient method for dual ascent



update $\left\{ \begin{array}{l} \alpha^{k+1} = \alpha^k + t^k \nabla g(\alpha^k) \\ \nabla g(\alpha^k) = \tilde{x} - \tilde{z} \quad \text{where} \quad (\tilde{x}, \tilde{z}) = \operatorname{argmin}_{x, z} L(x, z; \alpha^k) \end{array} \right.$

separable

Dual ascent



- maximize the dual problem $g(\lambda, v)$, which is always convex

$$g(\lambda, v) = \min_x L(x, \lambda, v) = \min_x f_0(x) + \lambda^\top f(x) + v^\top h(x)$$

$$\leq f_0(x^k) + \lambda^\top f(x^k) + v^\top h(x^k)$$

$$= f_0(x^k) + \lambda^k{}^\top f(x^k) + (\lambda - \lambda^k)^\top f(x^k)$$

$$+ v^k{}^\top h(x^k) + (v - v^k)^\top h(x^k)$$

x^k is the minimizer of the Lagrangian at λ^k, v^k

$$= g(\lambda^k, v^k) + (\lambda - \lambda^k)^\top f(x^k) + (v - v^k)^\top h(x^k)$$



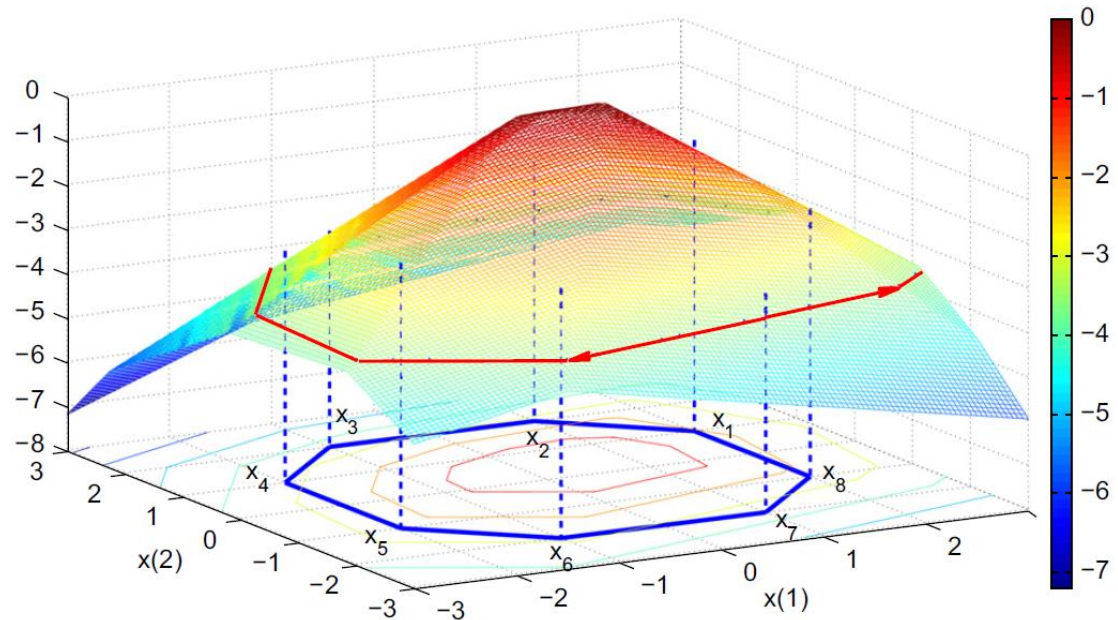
$$\begin{bmatrix} f(x^k) \\ h(x^k) \end{bmatrix} \in \begin{bmatrix} \frac{\partial g(\lambda^k, v^k)}{\partial \lambda} \\ \frac{\partial g(\lambda^k, v^k)}{\partial v} \end{bmatrix}$$



Dual ascent

- dual ascent method

- $x^{k+1} = \min_x f_0(x) -$
- $\lambda^{k+1} = \max\{\lambda^k + t$
- $v^{k+1} = v^k + t^k h(x$



- DC (difference of convex functions) algorithm

$$\min_x g(x) - h(x)$$



$$\max_y g^*(y) - h^*(y)$$

both are convex

concave optimization

$$\min_x t - h(x), \text{ s.t. } g(x) \leq t$$

- update: $y^k = \underset{y}{\operatorname{argmin}} h^*(y) - g^*(y^{k-1}) - x^{k\top}(y - y^{k-1})$
- $x^{k+1} = \underset{x}{\operatorname{argmax}} x^\top y^k - g(x)$

Alternating direction method of multipliers



- **Dual decomposition: Lagrangian**

$$L(x, z; \alpha) = \lambda \|x\|_1 + \|b - Az\|_2^2 + \alpha^T (x - z)$$

- **Method of multipliers: augmented Lagrangian**

$$L_\rho(x, z; \alpha) = \lambda \|x\|_1 + \|b - Az\|_2^2 + \alpha^T (x - z) + \frac{\rho}{2} \|x - z\|_2^2$$



- **ADMM**

$$x^{k+1} = \operatorname{argmin}_x L_\rho(x, v^k; \alpha^k)$$

$$z^{k+1} = \operatorname{argmin}_z L_\rho(x^{k+1}, z; \alpha^k)$$

$$\alpha^{k+1} = \alpha^k + \rho(x^{k+1} - z^{k+1})$$



strong duality
condition number adjustment

ADMM for LASSO



$$\min_x \lambda \|x\|_1 + \|b - Az\|_2^2$$

- **augmented Lagrangian**

$$L_\rho(x, z; \alpha) = \lambda \|x\|_1 + \|b - Az\|_2^2 + \alpha^T(x - z) + \frac{\rho}{2} \|x - z\|_2^2$$



- **x-update**

$$\begin{aligned} x^{k+1} &= \operatorname{argmin}_x L_\rho(x, z^k; \alpha^k) \\ &= \operatorname{argmin}_x \lambda \|x\|_1 + x^\top \alpha^k + \frac{\rho}{2} \|x - z^k\|_2^2 \\ &= S_{\lambda/\rho}(z^k + \alpha^k/\rho) \end{aligned}$$

ADMM for LASSO



$$\min_x \lambda \|x\|_1 + \|b - Az\|_2^2$$

- **augmented Lagrangian**

$$L_\rho(x, z; \alpha) = \lambda \|x\|_1 + \|b - Az\|_2^2 + \alpha^T(x - z) + \frac{\rho}{2} \|x - z\|_2^2$$

- **z-update**

$$\begin{aligned} z^{k+1} &= \operatorname{argmin}_z L_\rho(x^{k+1}, z; \alpha^k) \\ &= \operatorname{argmin}_z \|b - Az\|_2^2 + z^T \alpha^k + \frac{\rho}{2} \|x^{k+1} - z\|_2^2 \\ &= (A^T A + \rho I)^{-1} (A^T b + \rho x^{k+1} - \alpha^k) \end{aligned}$$

ADMM for LASSO



- **Dual decomposition: Lagrangian**

$$L(x, z; \alpha) = \lambda \|x\|_1 + \|b - Az\|_2^2 + \alpha^T (x - z)$$

- **Method of multipliers: augmented Lagrangian**

$$L_\rho(x, z; \alpha) = \lambda \|x\|_1 + \|b - Az\|_2^2 + \alpha^T (x - z) + \frac{\rho}{2} \|x - z\|_2^2$$

- **ADMM**

$$x^{k+1} = \operatorname{argmin}_x L_\rho(x, v^k; \alpha^k)$$

analytical solution

$$z^{k+1} = \operatorname{argmin}_z L_\rho(x^{k+1}, z; \alpha^k)$$

analytical solution, inverse only once

$$\alpha^{k+1} = \alpha^k + \rho(x^{k+1} - z^{k+1})$$

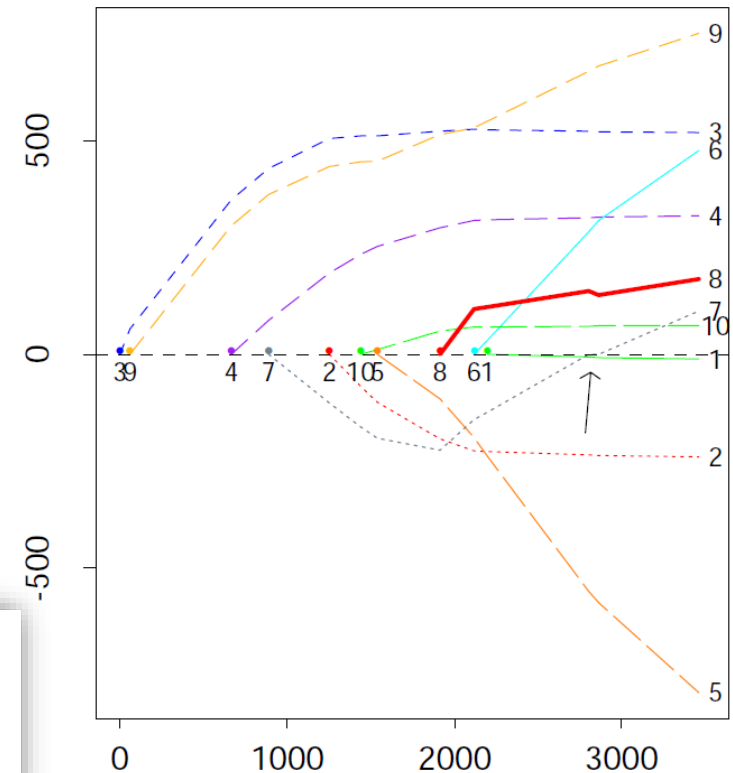
LARS: Least Angle Regression



$$\min_x \lambda \|x\|_1 + \|b - Ax\|_2^2$$

- **another interesting idea**
 - the optimal solution is a function of λ , denoted by $x(\lambda)$
 - $x(\lambda)$ is piecewise linear and thus, trackable
 - $x(\infty)$ is known

Lasso



Least Angle Regression

Bradley Efron, Trevor Hastie, Iain Johnstone and Robert Tibshirani
 Statistics Department, Stanford University

Alternating direction method of multipliers

$$\begin{aligned} \min \quad & f(x) + g(z) \\ \text{s. t.} \quad & Ax + Bz = c \end{aligned}$$

Method of multipliers: augmented Lagrangian

$$L_\rho(x, z; \alpha) = \lambda f(x) + g(z) + \alpha^T (Ax + Bz - c) + \frac{\rho}{2} \|Ax + Bz - c\|_2^2$$

ADMM

$$x^{k+1} = \operatorname{argmin}_x L_\rho(x, z^k; \alpha^k)$$

$$z^{k+1} = \operatorname{argmin}_z L_\rho(x^{k+1}, z; \alpha^k)$$

$$\alpha^{k+1} = \alpha^k + \rho(Ax^{k+1} + Bz^{k+1} - c)$$

$$\begin{aligned} z^{k+1} \quad \text{minimize} \quad & L(x^{k+1}, z, y^k) \\ \Rightarrow 0 = \nabla_z L(x^{k+1}, z, y^k) \\ & = \nabla g(z^{k+1}) + B^T y^k + \rho B^T (Ax^{k+1} + Bz^{k+1} - c) \\ & = \nabla g(z^{k+1}) + B^T (y^k + \rho (Ax^{k+1} + Bz^{k+1} - c)) \\ (x^{k+1}, z^{k+1}, y^{k+1}) & = \nabla g(z^{k+1}) + B^T y^{k+1} \\ & \text{Satisfies dual feasibility.} \end{aligned}$$

Convergence on ADMM



$$\begin{array}{ll} \min & f(x) + g(z) \\ \text{s. t.} & Ax + Bz = c \end{array}$$

- **assumption**

- f and g are convex, closed, proper
- the Lagrangian has a saddle point

- **convergence**

- approach feasibility
- approach optimal solution

More discussions



- in-exact iteration

$$z^{k+1} = (A^\top A + \rho I)^{-1} (A^\top b + \rho x^{k+1} - \alpha^k)$$

- e.g., solve by SGD
- update order
- multiple-block

$$\min_{x_i} \quad \sum_i h_i(x_i) \quad \text{s.t.} \quad \sum_i A_i x_i = c$$

- convergence discussion becomes hard
 - difficult to find hyper-parameters

Related algorithms



- operator splitting methods (Douglas, Bregman, Lions and Mercier, ...)
- proximal methods (Rockafellar, ...)

$$\mathbf{prox}_{\lambda f}(u) = \operatorname{argmin}_x (f(x) + 1/2\lambda \|x - u\|_2^2)$$

$$x^{k+1} = \mathbf{prox}_{\lambda f}(x^k - \lambda^k \nabla f(x^k))$$

proximal gradient method

- Dykstra's alternating projections algorithm
- Singarn's method of partial inverses
- Rockafellar-Wets progressive hedging for scenario-based decomposition

1

Frank-Wolfe Algorithm

2

ADMM

3

Parallel Computing



Parallel computing by ADMM



$$\min_x \Omega(x) + \sum_i l(a_i^\top x)$$



$$\min_x \Omega(x_0) + \sum_i f_i(x_i) \quad \text{s.t.} \quad x_i = x_0, \forall i$$

- augmented Lagrangian:

$$L_\rho(x_0, x_i; \alpha_i) = \Omega(x_0) + \sum_i f_i(x_i) + \sum_i \alpha_i^\top (x_i - x_0) + \frac{\rho}{2} \sum_i \|x_i - x_0\|_2^2$$

- update:

$$x_0^{k+1} = \operatorname{argmin}_{x_0} L_\rho(x_0, x_i^k; \alpha_i^k)$$



$$x_i^{k+1} = \operatorname{argmin}_{x_i} L_\rho(x_0^{k+1}, x_i; \alpha_i^k)$$

$$\alpha_i^{k+1} = \alpha_i^k + \rho(x_i^{k+1} - x_0^{k+1})$$

Parallel computing by ADMM



$$\min_x \Omega(x) + \sum_i l(a_i^\top x)$$



$$\min_x \Omega(x_0) + \sum_i f_i(x_i) \quad \text{s.t.} \quad x_i = x_0, \forall i$$

- augmented Lagrangian:

$$L_\rho(x_0, x_i; \alpha_i) = \Omega(x_0) + \sum_i f_i(x_i) + \sum_i \alpha_i^\top (x_i - x_0) + \frac{\rho}{2} \sum_i \|x_i - x_0\|_2^2$$

- update:

$$x_0^{k+1} = \operatorname{argmin}_{x_0} L_\rho(x_0, x_i^k; \alpha_i^k)$$

the master node takes charge of this part

$$x_i^{k+1} = \operatorname{argmin}_{x_i} L_\rho(x_0^{k+1}, x_i; \alpha_i^k)$$

$$\alpha_i^{k+1} = \alpha_i^k + \rho(x_i^{k+1} - x_0^{k+1})$$

} each worker takes charge of this part

Parallel computing for SVM



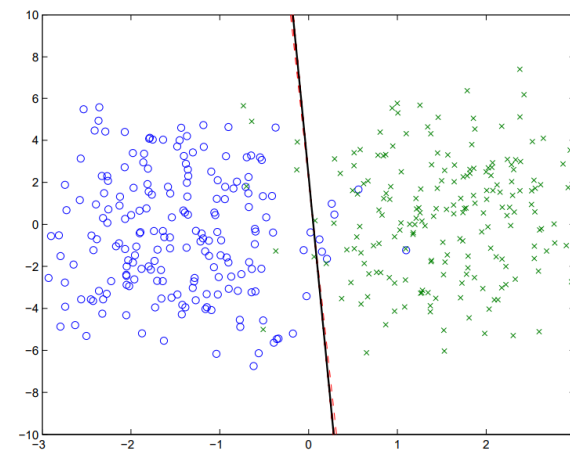
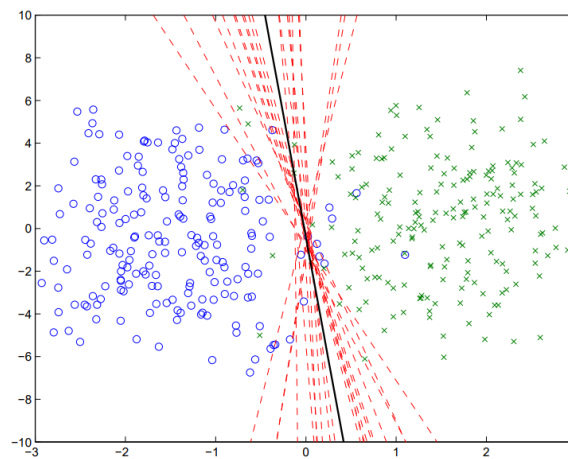
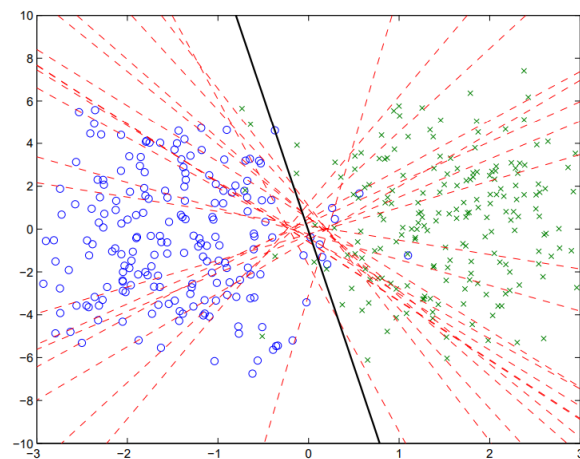
$$\min_x \lambda x^\top x + \sum_i \max\{0, b_i(a_i^\top x)\}$$

- toy example https://stanford.edu/~boyd/papers/pdf/admm_slides.pdf
 - 2 dimension, 400 samples
 - split into 20 groups (worst: each group contains examples in only on class)

Iteration 1

Iteration 5

Iteration 40



Parallel computing by ADMM



$$\min_x \Omega(x) + \sum_i l(a_i^\top x)$$



$$\min_x \Omega(x_0) + \sum_i f_i(x_i) \quad \text{s.t.} \quad x_i = x_0, \forall i$$

- augmented Lagrangian:

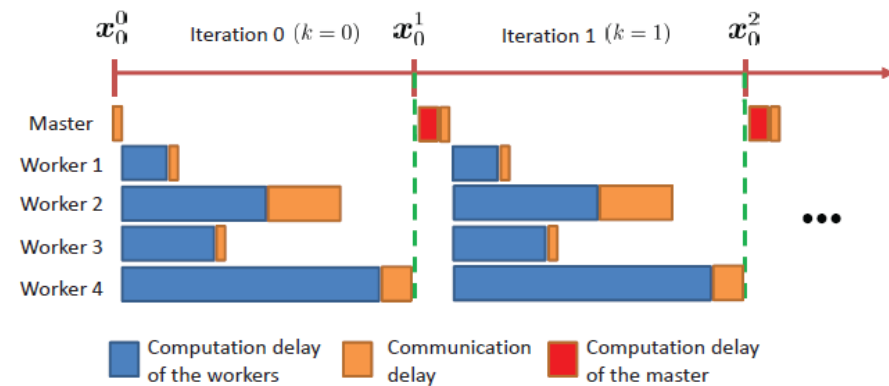
$$L_\rho(x_0, x_i; \alpha_i) = \Omega(x_0) + \sum_i f_i(x_i) + \sum_i \alpha_i^\top (x_i - x_0) + \frac{\rho}{2} \sum_i \|x_i - x_0\|_2^2$$

- update:

$$x_0^{k+1} = \operatorname{argmin}_{x_0} L_\rho(x_0, x_i^k; \alpha_i^k)$$

$$x_i^{k+1} = \operatorname{argmin}_{x_i} L_\rho(x_0^{k+1}, x_i; \alpha_i^k)$$

$$\alpha_i^{k+1} = \alpha_i^k + \rho(x_i^{k+1} - x_0^{k+1})$$



Parallel computing by ADMM



$$\min_x \Omega(x) + \sum_i l(a_i^\top x)$$



$$\min_x \Omega(x_0) + \sum_i f_i(x_i) \quad \text{s.t.} \quad x_i = x_0, \forall i$$

- augmented Lagrangian:

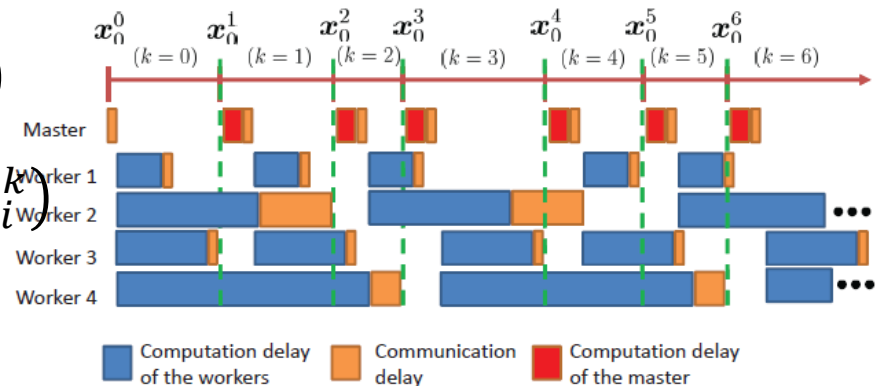
$$L_\rho(x_0, x_i; \alpha_i) = \Omega(x_0) + \sum_i f_i(x_i) + \sum_i \alpha_i^\top (x_i - x_0) + \frac{\rho}{2} \sum_i \|x_i - x_0\|_2^2$$

- update:

$$x_0^{k+1} = \operatorname{argmin}_{x_0} L_\rho(x_0, x_i^k; \alpha_i^k)$$

$$x_i^{k+1} = \operatorname{argmin}_{x_i} L_\rho(x_0^{k+1}, x_i; \alpha_i^k)$$

$$\alpha_i^{k+1} = \alpha_i^k + \rho(x_i^{k+1} - x_0^{k+1})$$



Parallel computing by ADMM



$$\min_x \Omega(x) + \sum_i l(a_i^\top x)$$



$$\min_x \Omega(x_0) + \sum_i f_i(x_i) \quad \text{s.t.} \quad x_i = x_0, \forall i$$

- augmented Lagrangian:

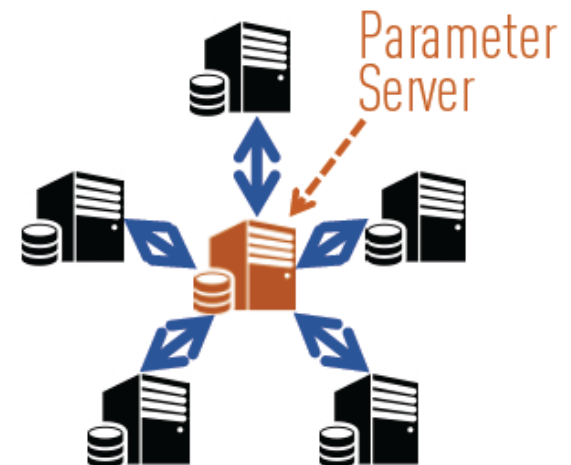
$$L_\rho(x_0, x_i; \alpha_i) = \Omega(x_0) + \sum_i f_i(x_i) + \sum_i \alpha_i^\top (x_i - x_0) + \frac{\rho}{2} \sum_i \|x_i - x_0\|_2^2$$

- update:

$$x_0^{k+1} = \operatorname{argmin}_{x_0} L_\rho(x_0, x_i^k; \alpha_i^k)$$

$$x_i^{k+1} = \operatorname{argmin}_{x_i} L_\rho(x_0^{k+1}, x_i; \alpha_i^k)$$

$$\alpha_i^{k+1} = \alpha_i^k + \rho(x_i^{k+1} - x_0^{k+1})$$



Parallel computing by ADMM



$$\min_x \Omega(x) + \sum_i l(a_i^\top x)$$



$$\min_x \Omega(x_0) + \sum_i f_i(x_i) \quad \text{s.t.} \quad x_i = x_0, \forall i$$

- augmented Lagrangian:

$$L_\rho(x_0, x_i; \alpha_i) = \Omega(x_0) + \sum_i f_i(x_i) + \sum_i \alpha_i^\top (x_i - x_0) + \frac{\rho}{2} \sum_i \|x_i - x_0\|_2^2$$

- update:

$$x_0^{k+1} = \operatorname{argmin}_{x_0} L_\rho(x_0, x_i^k; \alpha_i^k)$$

$$x_i^{k+1} = \operatorname{argmin}_{x_i} L_\rho(x_0^{k+1}, x_i; \alpha_i^k)$$

$$\alpha_i^{k+1} = \alpha_i^k + \rho(x_i^{k+1} - x_0^{k+1})$$



Decentralized strategy



$$\min_x \Omega(x) + \sum_i l(a_i^\top x)$$

- Decentralized SGD
 - compute local stochastic gradient
 - compute the neighborhood weighted average

$$x_i^{k+\frac{1}{2}} = \sum_{i \in \text{Nei}} a_i x_i^k$$

- update local variable

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



Can Decentralized Algorithms Outperform Centralized Algorithms? A Case Study for Decentralized Parallel Stochastic Gradient Descent

Xiangru Lian[†], Ce Zhang^{*}, Huan Zhang^{*}, Cho-Jui Hsieh^{*}, Wei Zhang[#], and Ji Liu^{†‡}

Decentralized strategy



$$\min_x \Omega(x) + \sum_i l(a_i^\top x)$$

- Decentralized SGD
 - compute local stochastic gradient
 - compute the neighborhood weighted average

$$x_i^{k+\frac{1}{2}} = \sum_{i \in \text{Nei}} a_i x_i^k$$

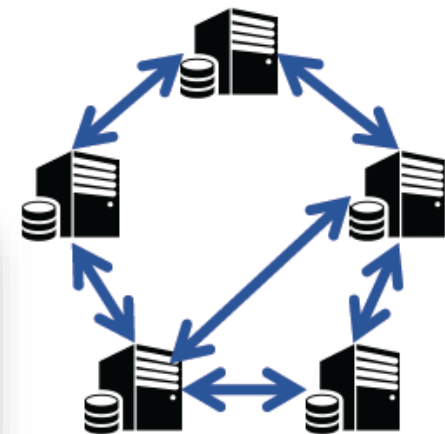
could be asynchronous

- update local variable

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

Asynchronous Decentralized Parallel Stochastic Gradient Descent

Xiangru Lian^{1*} Wei Zhang^{2*} Ce Zhang³ Ji Liu^{4,1}



Decentralized strategy



$$\min_x \Omega(x) + \sum_i l(a_i^\top x)$$

- Decentralized SGD

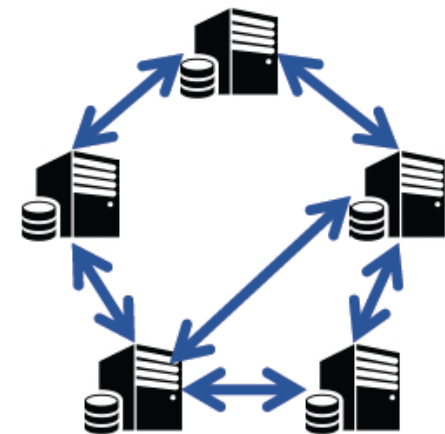
- compute local stochastic gradient
- compute the neighborhood weighted average

$$x_i^{k+\frac{1}{2}} = \sum_{i \in \text{Nei}} a_i x_i^k$$

- update local variable

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

when the data are
really decentralized
(there are variance)



D²: Decentralized Training over Decentralized Data

Hanlin Tang¹ Xiangru Lian¹ Ming Yan^{2,3} Ce Zhang⁴ Ji Liu^{5,1}

THANKS

