

Consistently better than SGD

Best regret bound

Ready to use

Code: <http://www.cs.umb.edu/~yangmu/code/CSGD.zip>

Constrained Stochastic Gradient Descent for Large-scale Least Squares Problem

Yang Mu¹, Wei Ding¹, Tianyi Zhou², Dacheng Tao²

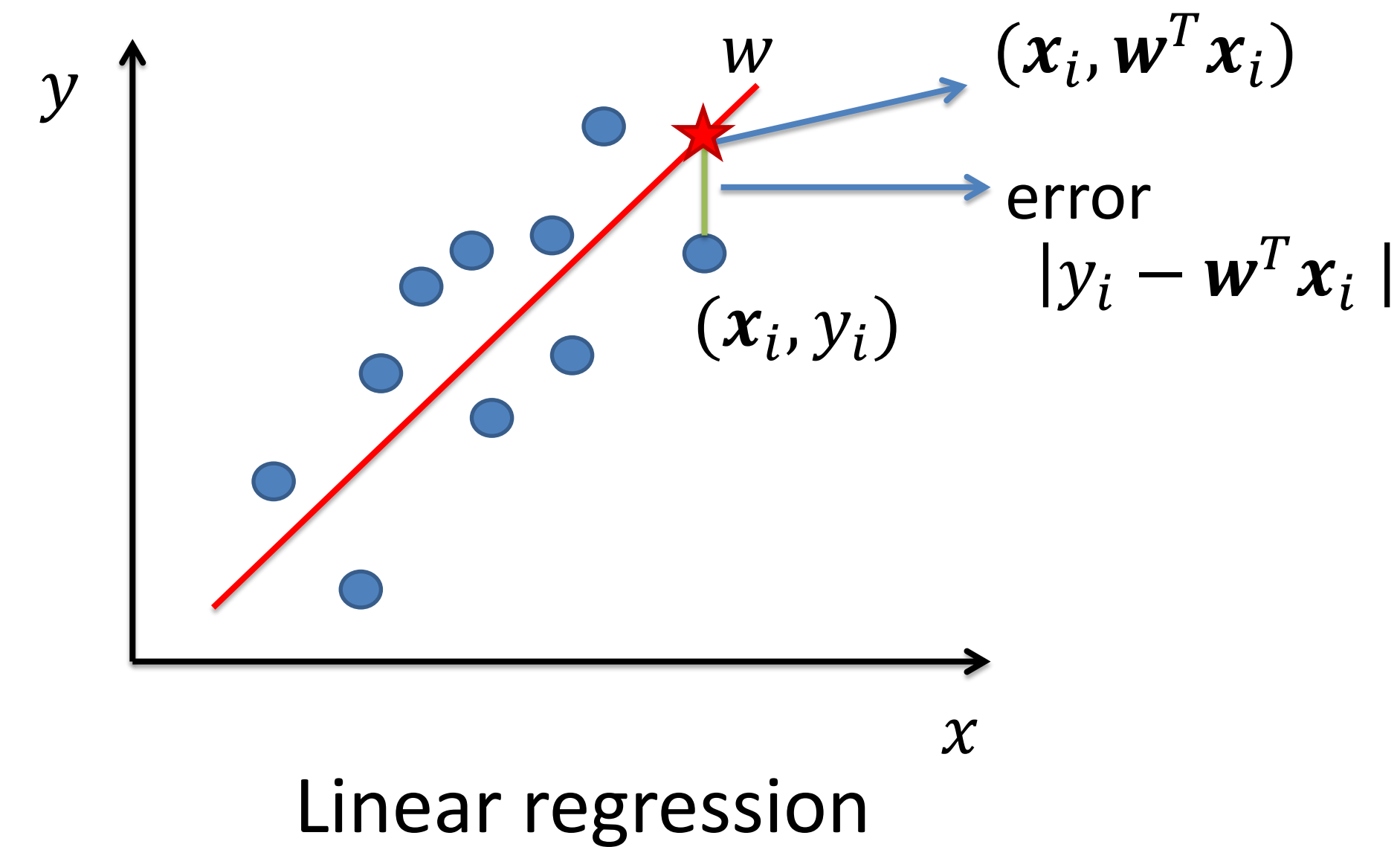
¹University of Massachusetts Boston

²University of Technology, Sydney

1. The Least Squares problem

Objective function:

$$\min_{\mathbf{w}} \sum_i \frac{1}{2} \|\mathbf{y}_i - \mathbf{w}^T \mathbf{x}_i\|_2^2, \text{ where } \mathbf{x}_i \in \mathbb{R}^d, \mathbf{y}_i \in \mathbb{R}^1 \text{ and } \mathbf{w} \in \mathbb{R}^d$$



2. When Least Squares problem meets Large Scale

Use closed form solution: $\mathbf{w} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$?



It's $O(nd^2)$!

Note: Any algorithm with time complexity greater than $O(nd)$ is not applicable in large scale high dimension cases.

Stochastic Gradient Descent (**SGD**) with $O(d)$ time each iteration is an appealing approach, which takes the form

$$\mathbf{w}_{t+1}^* = \arg \min_{\mathbf{w}} \frac{1}{t} \sum_{i=1}^t \frac{1}{2} \|\mathbf{y}_i - \mathbf{w}^T \mathbf{x}_i\|_2^2,$$

where $l(\mathbf{w}, \mathbf{x}_i, \mathbf{y}_i) = \frac{1}{2} \|\mathbf{y}_i - \mathbf{w}^T \mathbf{x}_i\|_2^2$ is the loss function at step i , denoted as $l(\mathbf{w})$ for short.

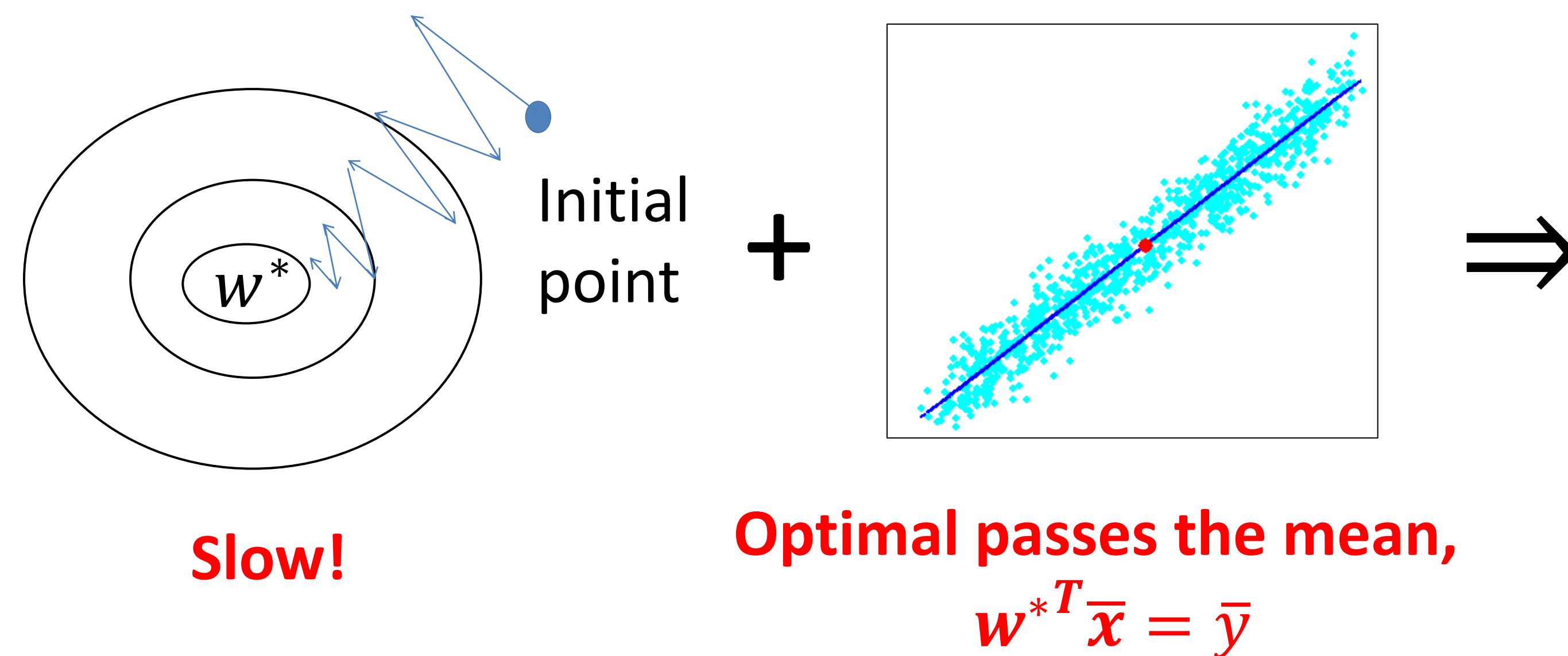
SGD update rule is:

$$\mathbf{w}_{t+1} = \mathbf{w}_t - \eta_t \mathbf{g}_t,$$

where $\mathbf{g}_t = \partial l(\mathbf{w})$.

Zig-Zag property of SGD

The beauty of Least Squares

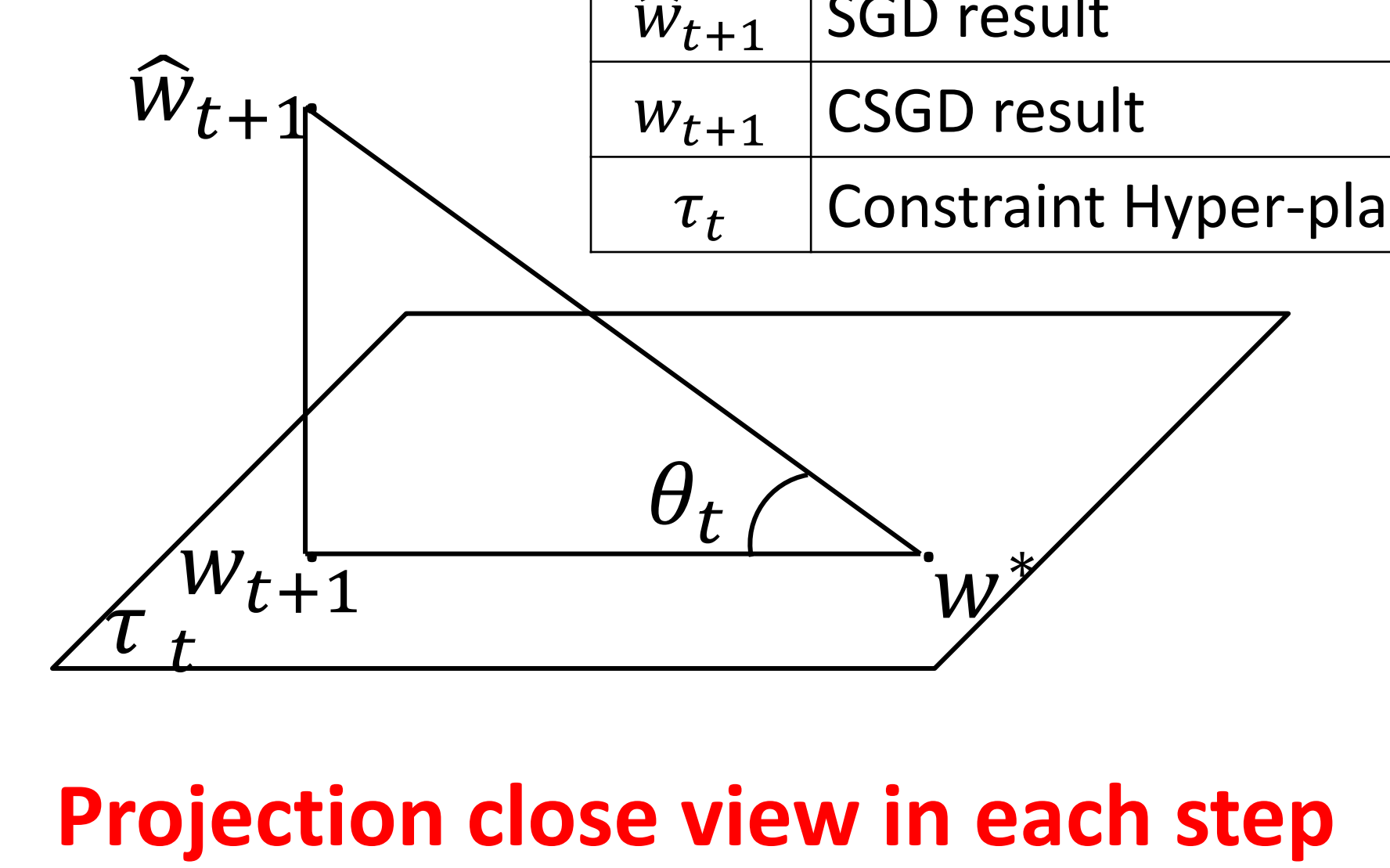
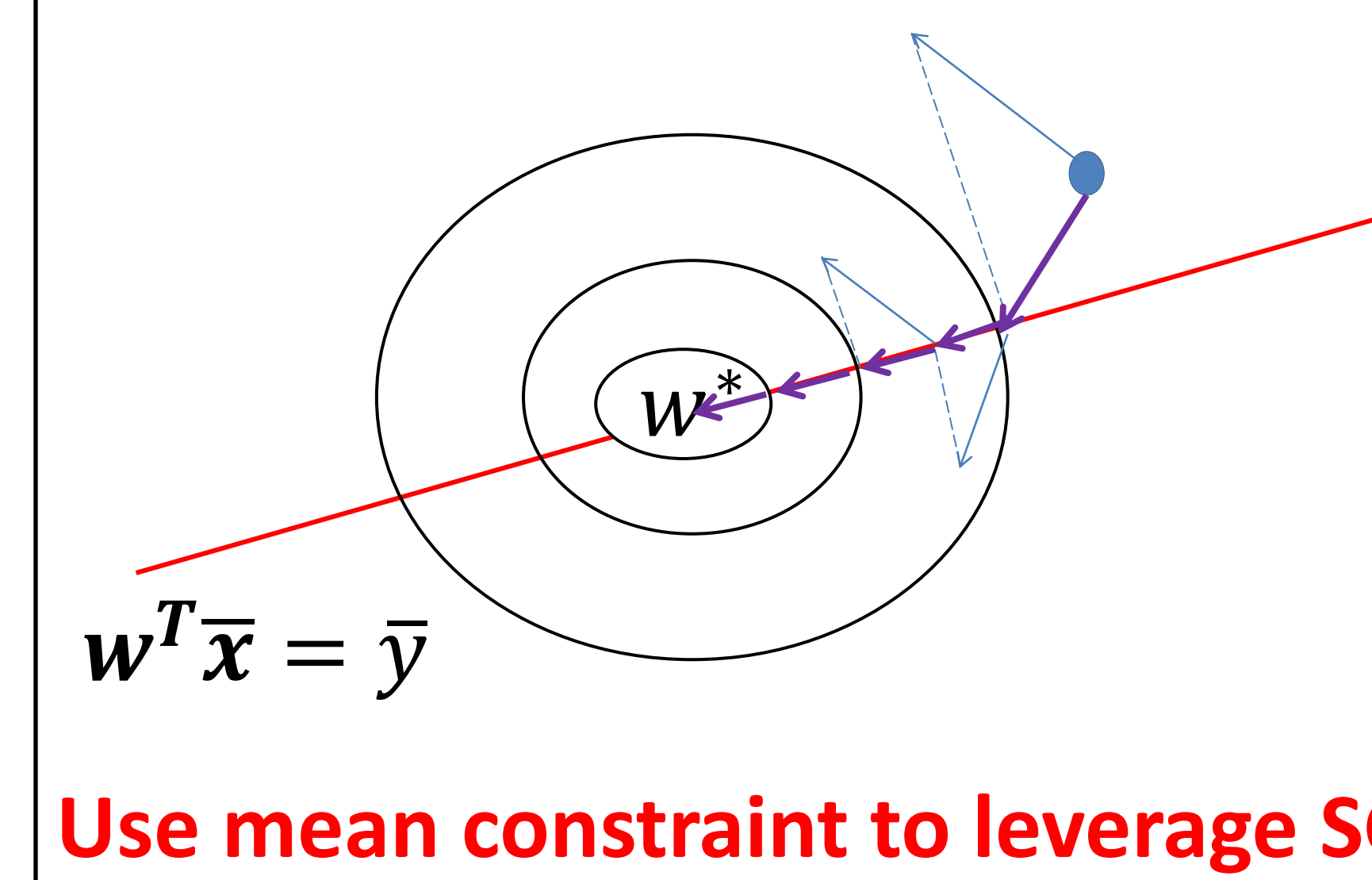


Remarks:

1. performs **consistently better** than SGD in terms of batch optimal. 2. has $O(\log T)$ regret bound $R_G(T) \leq \frac{G^2}{2H} (1 + \log T)$ 3. **extensible**.

3. Motivation

Constrained Stochastic Gradient Descent (CSGD)



w^*	batch optimal
\hat{w}_{t+1}	SGD result
w_{t+1}	CSGD result
τ_t	Constraint Hyper-plane

4. Algorithm

CSGD:

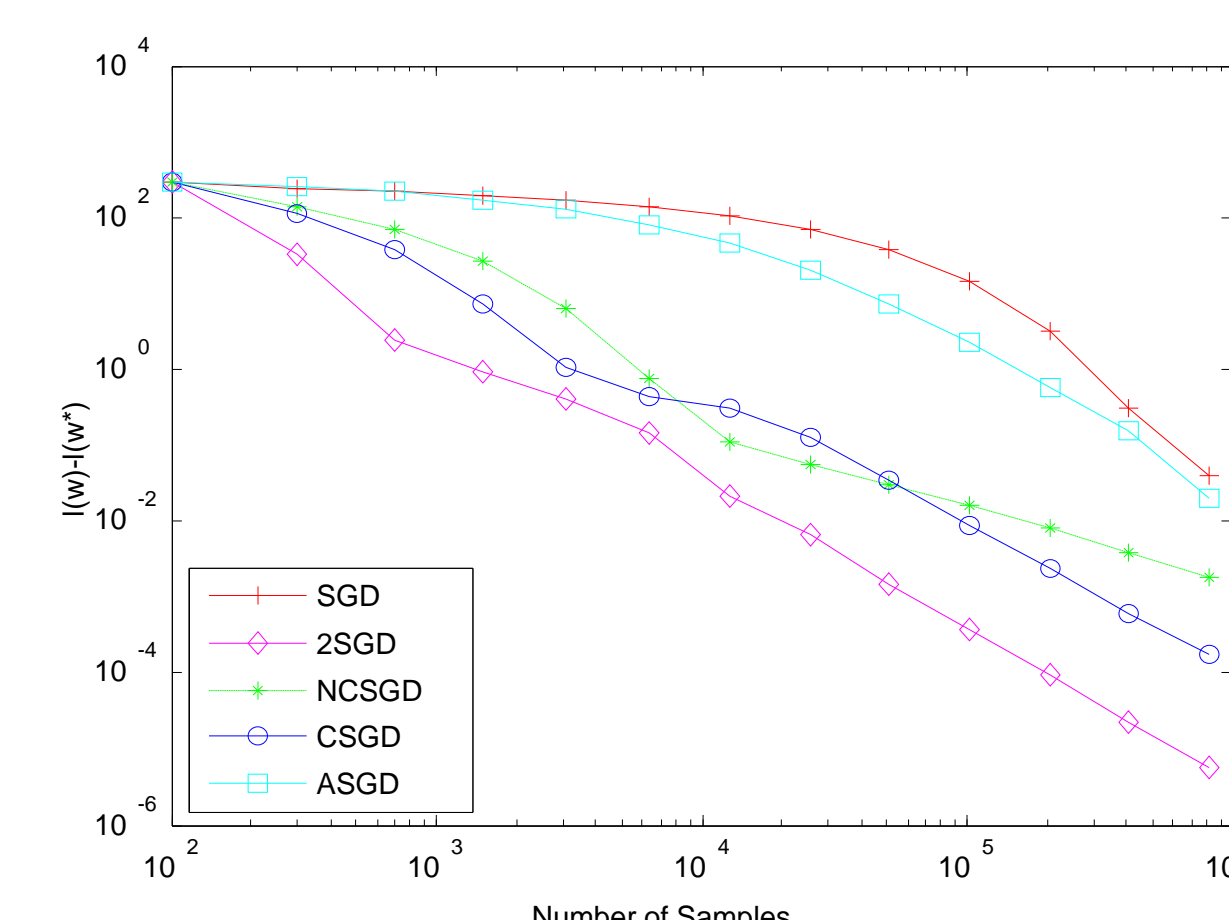
$$\mathbf{w}_{t+1}^* = \arg \min_{\mathbf{w}} \frac{1}{t} \sum_{i=1}^t \frac{1}{2} \|\mathbf{y}_i - \mathbf{w}^T \mathbf{x}_i\|_2^2, \text{ s.t. } \mathbf{w}^T \bar{\mathbf{x}}_t = \bar{\mathbf{y}}_t. \text{ The update rule } \mathbf{w}_{t+1} = \mathbf{P}_t(\mathbf{w}_t - \eta_t \mathbf{g}_t) + \mathbf{r}_t, \text{ where } \mathbf{P}_t = \mathbf{I} - \frac{\bar{\mathbf{x}}_t \bar{\mathbf{x}}_t^T}{\|\bar{\mathbf{x}}_t\|_2^2}, \mathbf{r}_t = \frac{\bar{\mathbf{y}}_t}{\|\bar{\mathbf{x}}_t\|_2^2} \bar{\mathbf{x}}_t$$

WAIT! Is $\mathbf{P}_t \in \mathbb{R}^{d \times d}$? **YES**, but it is rank ONE!

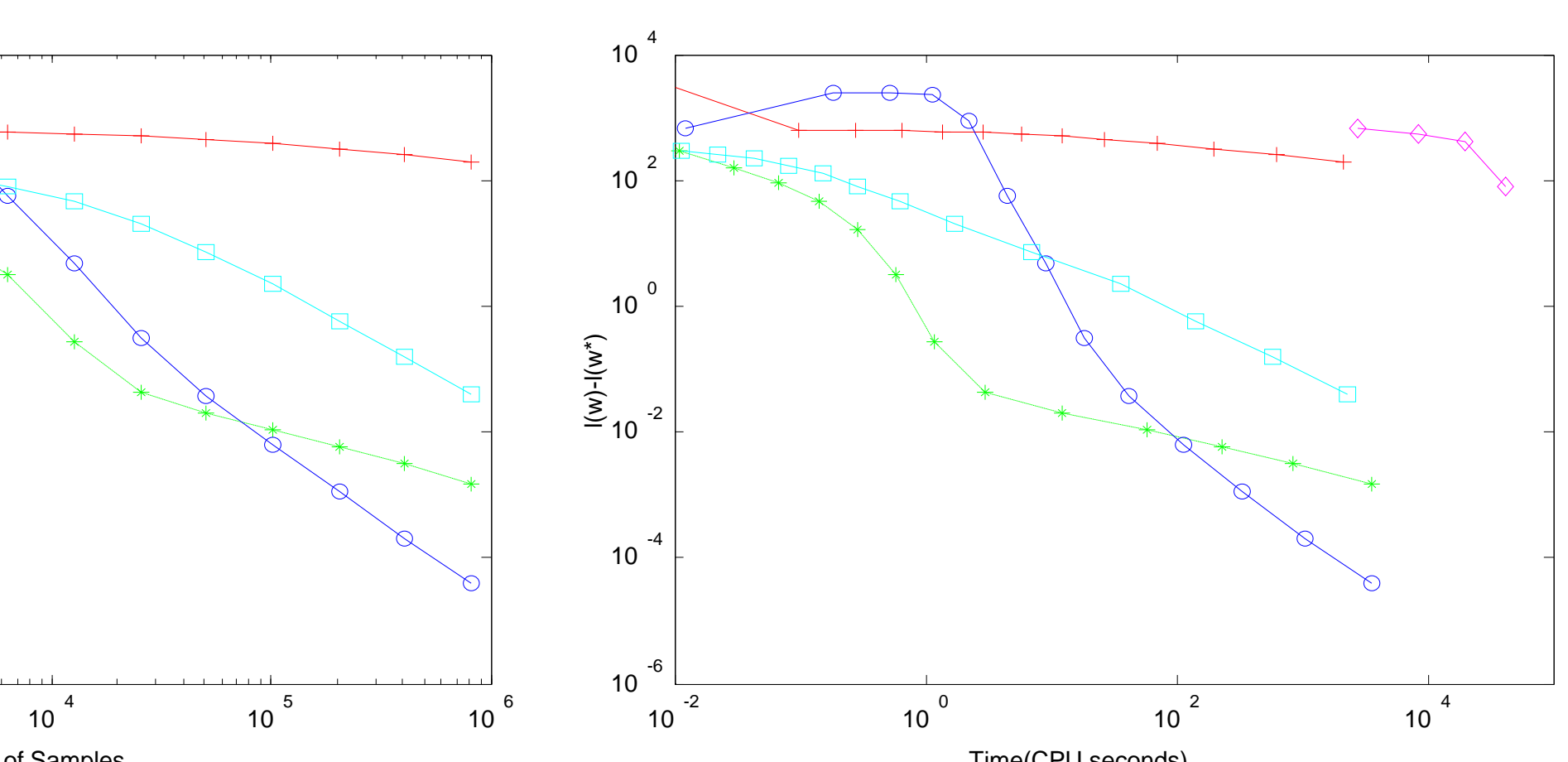
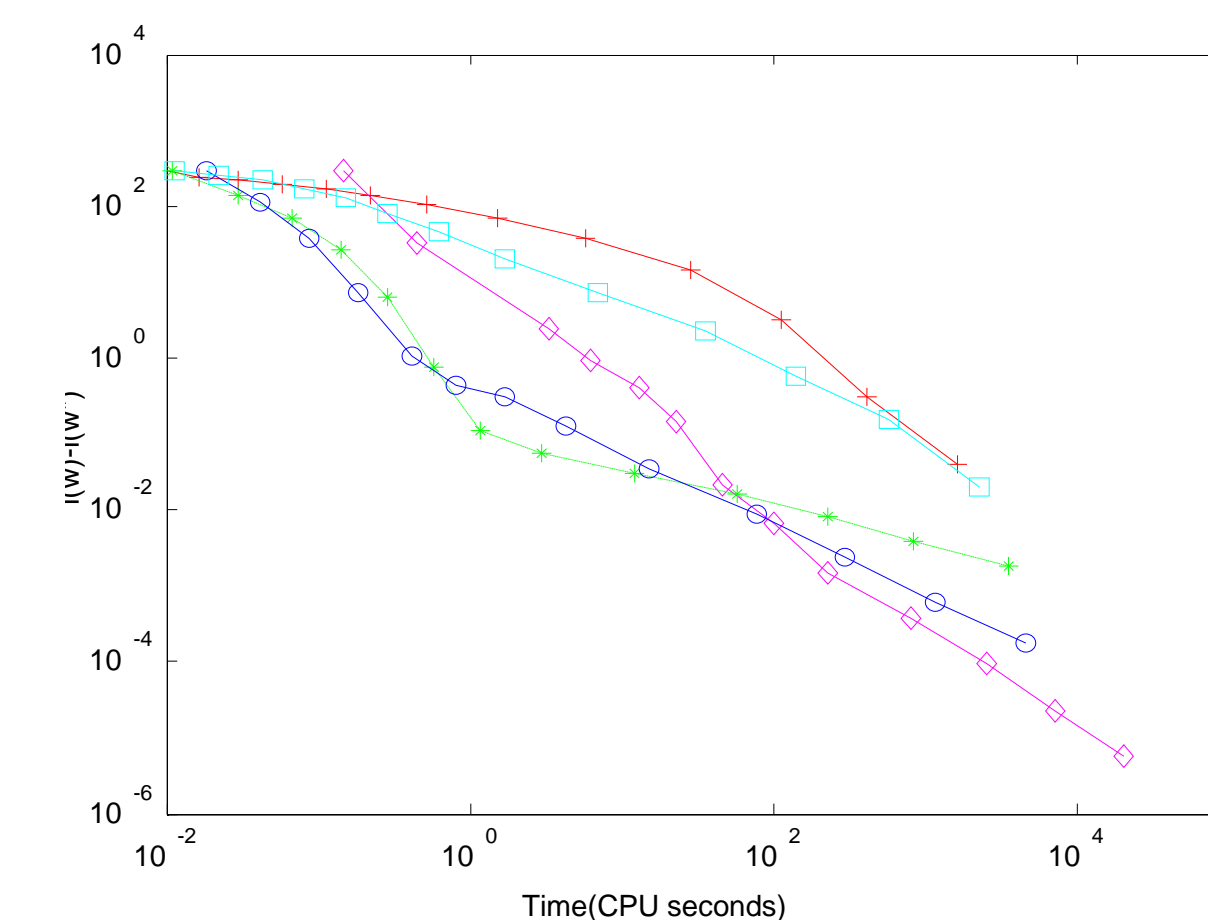
Therefore, we still have a $O(d)$ time complexity algorithm each iteration: $\mathbf{w}_{t+1} = \mathbf{w}_t - \eta_t \mathbf{g}_t - \bar{\mathbf{x}}_t (\bar{\mathbf{x}}_t^T (\mathbf{w}_t - \eta_t \mathbf{g}_t)) / \|\bar{\mathbf{x}}_t\|_2^2 + \mathbf{r}_t$.

5. Experiments

Optimization
Study

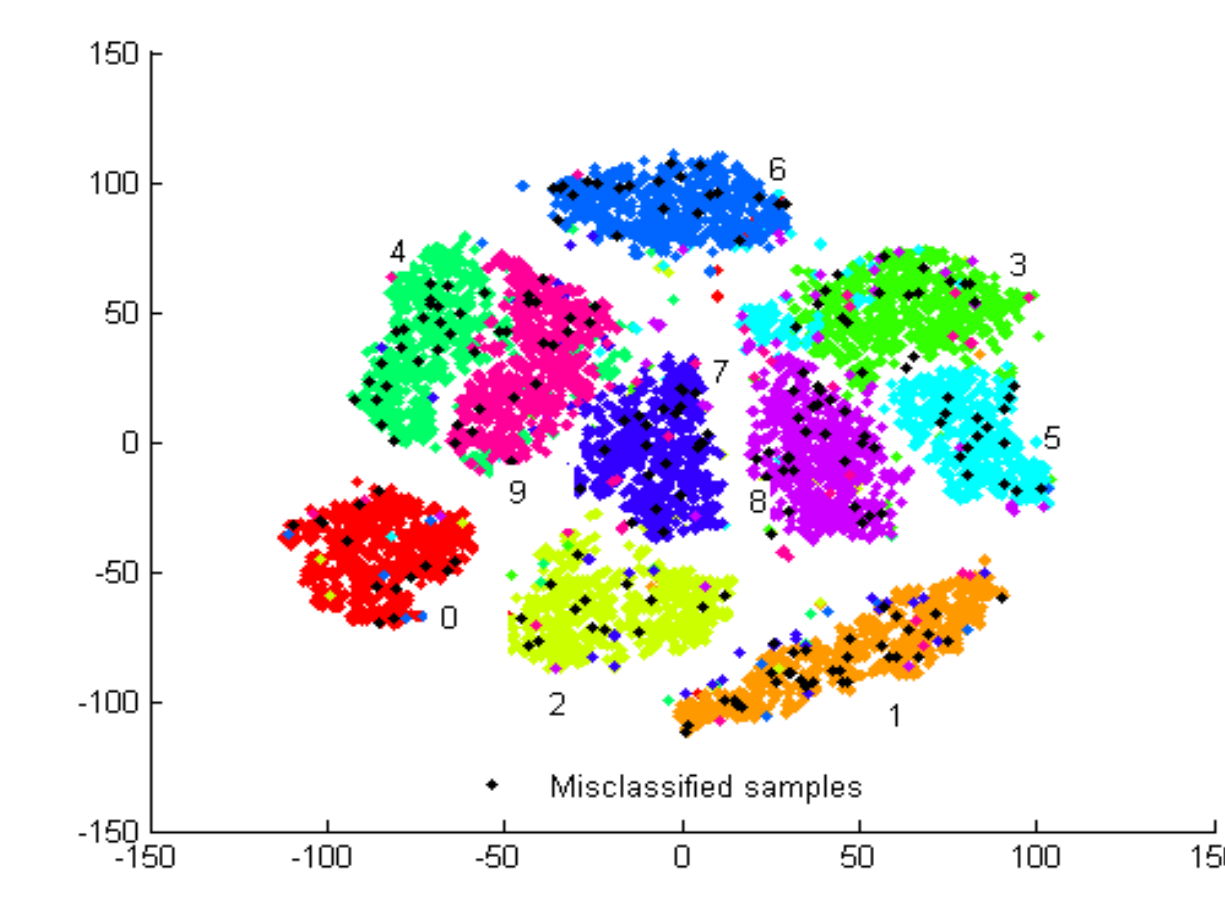
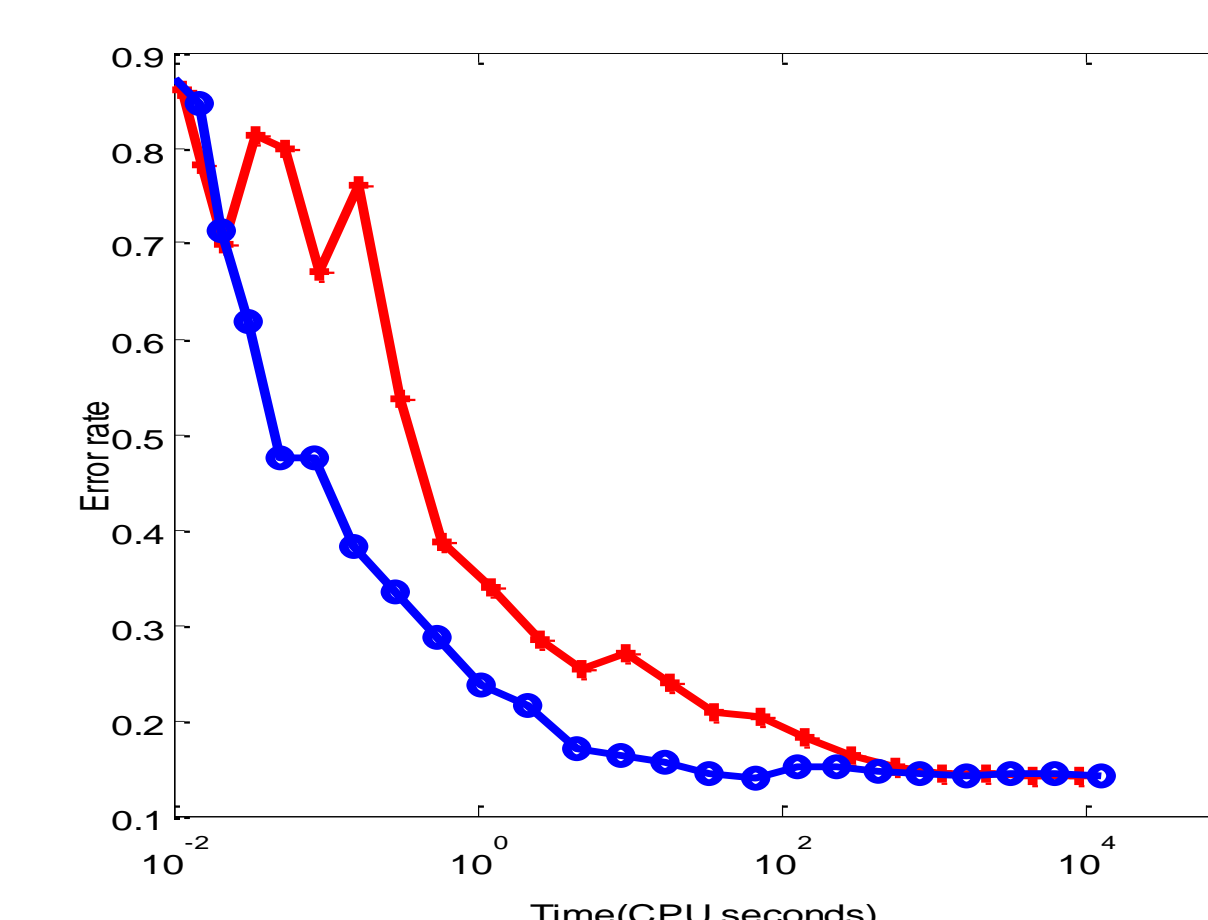
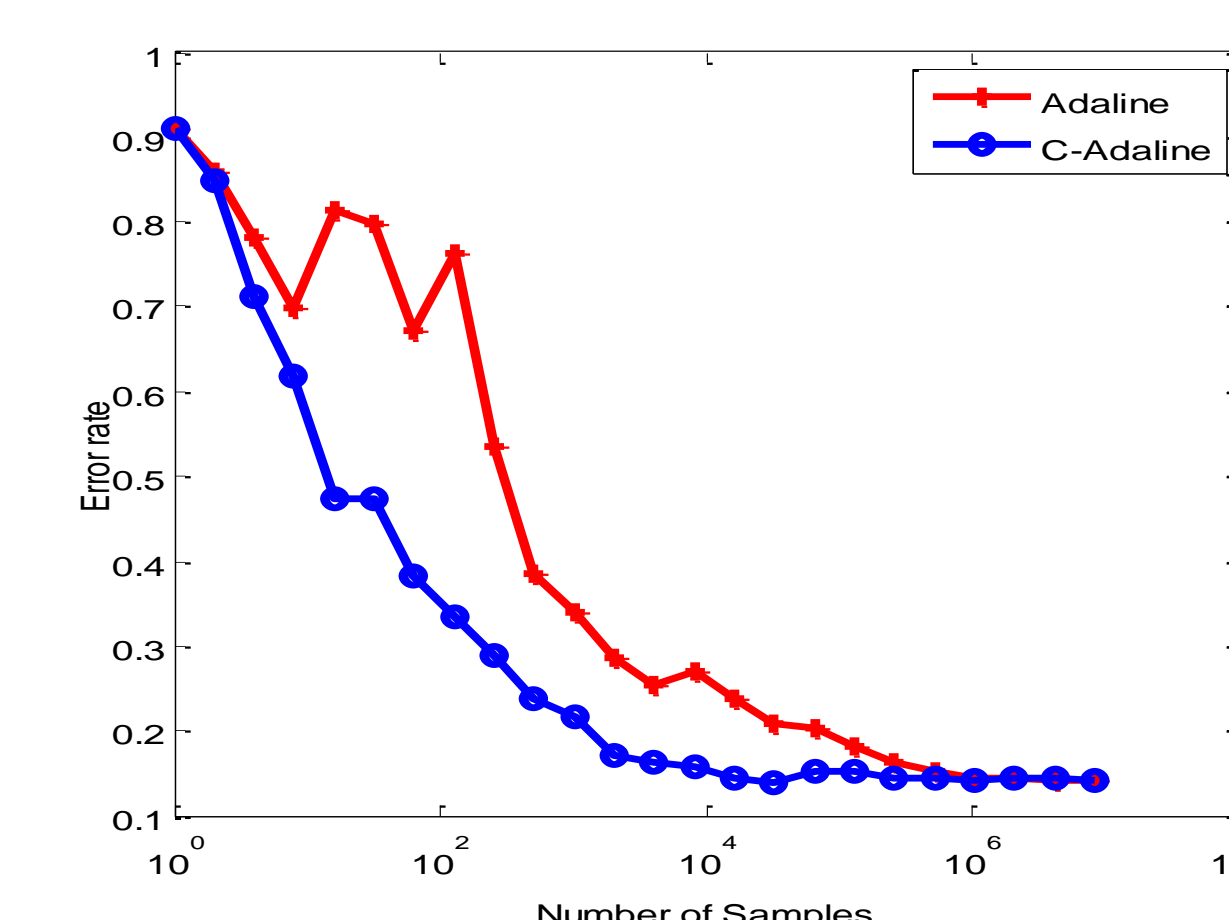


d=100, strong convex

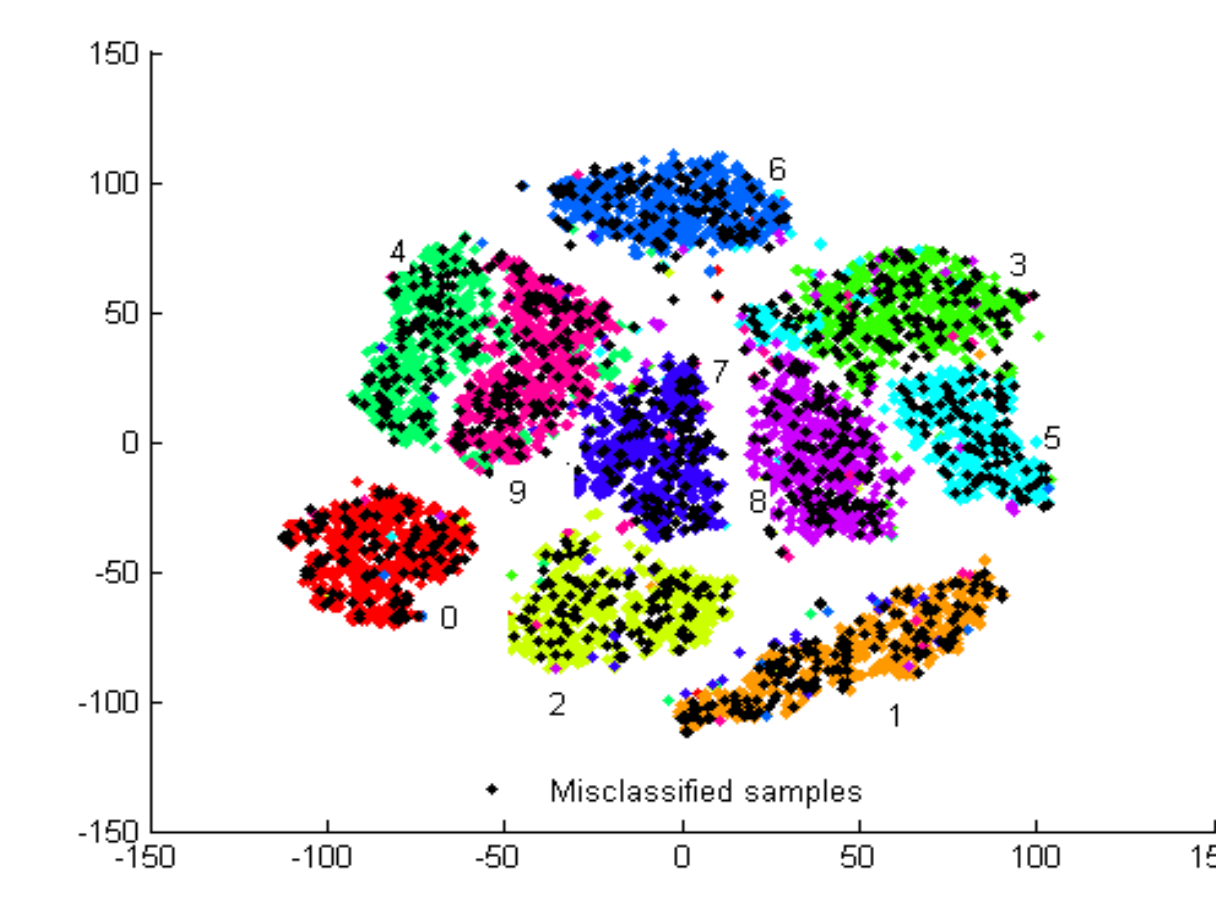


d=5,000, non strong convex

Classification
Study



212 misclassified samples
in 3.38 seconds



1248 misclassified samples
in 112.47 seconds

Synthetic
dataset

MNIST
dataset