

# Stochastic Block BFGS: Squeezing More Curvature out of Data

Robert Mansel Gower

Joint work with Donald Goldfarb and Peter Richtárik



International Conference on Machine Learning, New York, June 2016

# The Problem

$$\min_{w \in \mathbf{R}^d} f(w) \stackrel{\text{def}}{=} \frac{1}{n} \sum_{i=1}^n f_i(w).$$

- ◆ Each  $f_i$  is strongly convex and twice continuously differentiable.
- ◆ Far more data samples than features  $n \gg d$ , access through subsampling

# The Problem

**Motivation** is from stochastic optimization, such as empirical risk minimization

$$\min_{w \in \mathbf{R}^d} f(w) \stackrel{\text{def}}{=} \frac{1}{n} \sum_{i=1}^n f_i(w).$$

- ◆ Each  $f_i$  is strongly convex and twice continuously differentiable.
- ◆ Far more data samples than features  $n \gg d$ , access through subsampling

# The Problem

**Motivation** is from stochastic optimization, such as empirical risk minimization

$$\min_{w \in \mathbf{R}^d} f(w) \stackrel{\text{def}}{=} \frac{1}{n} \sum_{i=1}^n f_i(w).$$

- ◆ Each  $f_i$  is strongly convex and twice continuously differentiable.
- ◆ Far more data samples than features  $n \gg d$ , access through subsampling

$$\nabla f_S(w) \stackrel{\text{def}}{=} \frac{1}{|S|} \sum_{i \in S} \nabla f_i(w)$$



$$S \subset \{1, \dots, n\}$$

$$\nabla^2 f_T(x) \stackrel{\text{def}}{=} \frac{1}{|T|} \sum_{i \in T} \nabla^2 f_i(x)$$



$$T \subset \{1, \dots, n\}$$

# Variable Metric Method

$$w_{t+1} = w_t - \eta H_t g_t$$

# Variable Metric Method

$$w_{t+1} = w_t - \eta H_t g_t$$

stepsize



The diagram illustrates the role of the stepsize parameter in the Variable Metric Method. A yellow rectangular box containing the word "stepsize" is positioned below the equation. A yellow arrow points from the top of this box to the Greek letter eta (η) in the equation  $w_{t+1} = w_t - \eta H_t g_t$ , which is displayed within a light green rectangular box.

# Variable Metric Method

$$w_{t+1} = w_t - \eta H_t g_t$$

stepsize



A yellow box containing the text 'stepsize' has a yellow arrow pointing upwards to the  $\eta$  term in the equation above.

$$\mathbf{E}[g_t] = \nabla f(w_t)$$


A yellow box containing the equation  $\mathbf{E}[g_t] = \nabla f(w_t)$  has a yellow arrow pointing upwards and to the left to the  $g_t$  term in the equation above.

# Variable Metric Method

$$H_t \approx \nabla^2 f(w_t)^{-1}$$

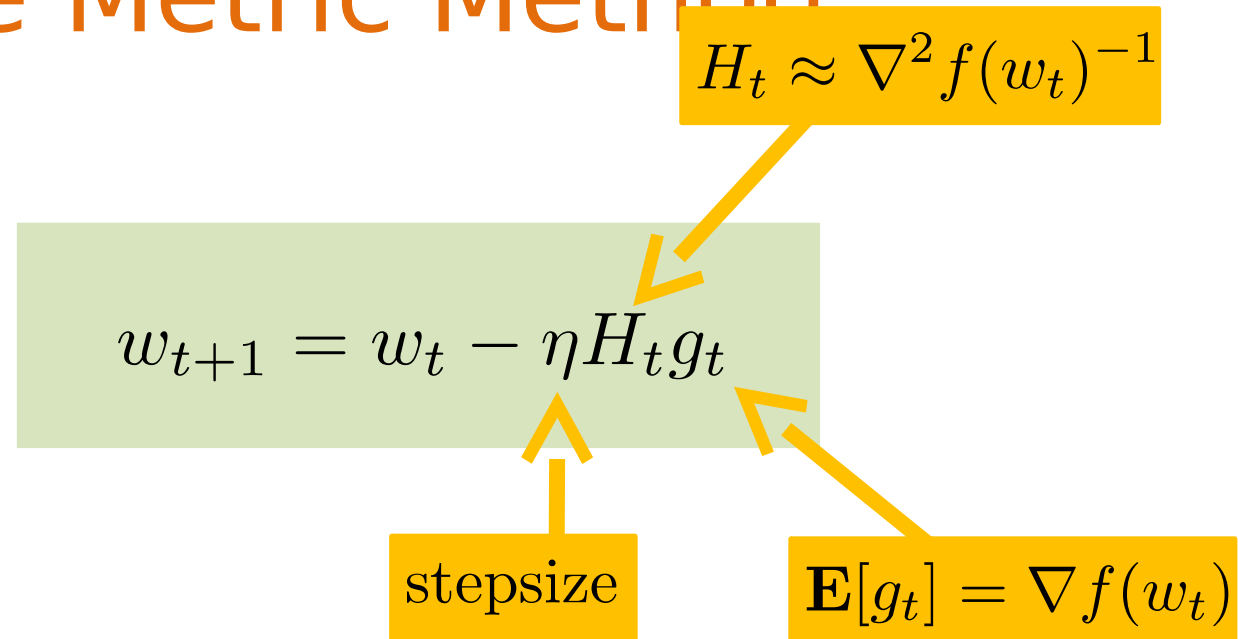
$$w_{t+1} = w_t - \eta H_t g_t$$

stepsize

$$\mathbf{E}[g_t] = \nabla f(w_t)$$



# Variable Metric Method



**Exe:** ♦ Newton's Method

$$w_{t+1} = w_t - \eta \nabla^2 f(w_t)^{-1} \nabla f(w_t)$$

♦ Steepest descent

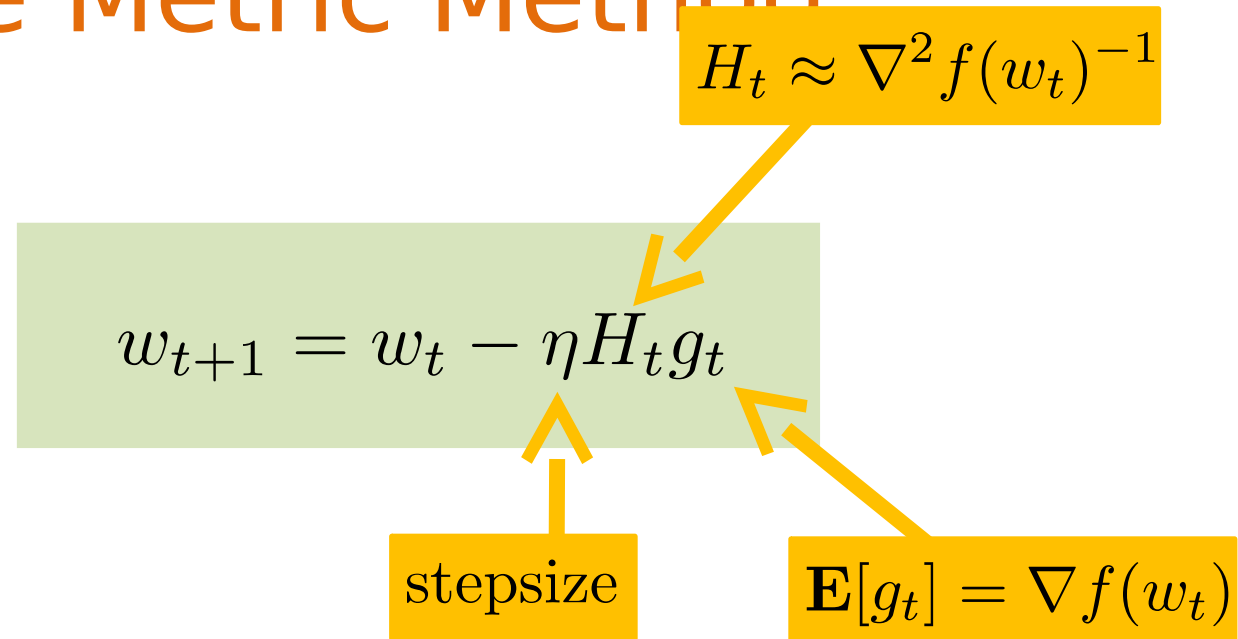
$$w_{t+1} = w_t - \eta \nabla f(w_t)$$

♦ Stochastic gradient descent (SGD)

$$w_{t+1} = w_t - \eta \nabla f_S(w_t)$$

♦ SAG, SVRG, S2GD, ...etc

# Variable Metric Method



**Exe:** ♦ Newton's Method

$$w_{t+1} = w_t - \eta \nabla^2 f(w_t)^{-1} \nabla f(w_t)$$

♦ Steepest descent

$$w_{t+1} = w_t - \eta \nabla f(w_t)$$

♦ Stochastic gradient descent (SGD)

$$w_{t+1} = w_t - \eta \nabla f_S(w_t)$$

♦ SAG, SVRG, S2GD, ...etc

**Challenge** how to construct an effective  $H_t$  that is cheap to calculate?

# Stochastic Second order Methods

$H_t$  is directly estimated from  $\nabla^2 f_T(x_t)$

- ◆ Low rank decomposition (Agarwal, Bullins and Hazan 2016)
- ◆ SVD decomposition (Erdogdu and Montanari 2015)
- ◆ Sketching full Hessian (Pilanci and Wainwright 2015)

$H_t$  is updated using the (L)BFGS update

- ◆ (Schraudolph, Yu and Gunter 2007)
- ◆ (Mokhtari and Ribeiro 2014, 2015)
- ◆ (Byrd, Hansen, Nocedal and Singer 2015)
- ◆ (**MNJ** Moritz, Nishihara, Jordan 2016)

# Hessian Sketching

**Fact:** Calculating a directional derivative of the gradient is cheap

$$\nabla^2 f_T(x_t)v = \left. \frac{d}{d\alpha} \nabla f_T(x_t + \alpha v) \right|_{\alpha=0}$$

**Ideally**  $H_t$  should satisfy the *inverse equation*

$$H_t \nabla^2 f_T(x_t) = I$$

Solving the **sketched inverse equation** is easier

$$H_t \nabla^2 f_T(x_t) D_t = D_t$$

# Hessian Sketching

**Fact:** Calculating a directional derivative of the of  $O(\text{eval}(f_T(x)))$  with Automatic Differentiation.

$$\nabla^2 f_T(x_t)v = \left. \frac{d}{d\alpha} \nabla f_T(x_t + \alpha v) \right|_{\alpha=0}$$

**Ideally**  $H_t$  should satisfy the *inverse equation*

$$H_t \nabla^2 f_T(x_t) = I$$

Solving the **sketched inverse equation** is easier

$$H_t \nabla^2 f_T(x_t) D_t = D_t$$

# Hessian Sketching

**Fact:** Calculating a directional derivative of the of  $O(\text{eval}(f_T(x)))$  with Automatic Differentiation.


$$\nabla^2 f_T(x_t)v = \left. \frac{d}{d\alpha} \nabla f_T(x_t + \alpha v) \right|_{\alpha=0}$$

**Ideally**  $H_t$  should satisfy the *inverse equation*

$$H_t \nabla^2 f_T(x_t) = I$$

Solving the **sketched inverse equation** is easier

$$H_t \nabla^2 f_T(x_t) D_t = D_t$$


$$D_t \in \mathbf{R}^{d \times q}, q \ll d$$

# Hessian Sketching

**Fact:** Calculating a directional derivative of the of  $O(\text{eval}(f_T(x)))$  with Automatic Differentiation.

$$\nabla^2 f_T(x_t)v = \left. \frac{d}{d\alpha} \nabla f_T(x_t + \alpha v) \right|_{\alpha=0}$$

**Ideally**  $H_t$  should satisfy the *inverse equation*

$$H_t \nabla^2 f_T(x_t) = I$$

Solving the **sketched inverse equation** is easier

$$H_t \nabla^2 f_T(x_t) D_t = D_t$$

Cost of evaluating  $\nabla^2 f_T(x_T) D_t$  is  $q \times O(\text{eval}(f_T(x)))$

$$D_t \in \mathbf{R}^{d \times q}, q \ll d$$

# Block BFGS: Least change formulation

$$H_t = \arg \min_{H \in \mathbb{R}^{d \times d}} ||H - H_{t-1}||_t^2$$

subject to  $H \nabla^2 f_T(x_t) D_t = D_t, \quad H = H^T$

where  $||H||_t^2 \stackrel{\text{def}}{=} \text{Tr} (H \nabla^2 f_T(x_t) H^T \nabla^2 f_T(x_t))$ .

**The constraint** serves as a fidelity term, enforcing that a sketch of the inverse equation be satisfied

**The objectives** serves as a regularizer, enforcing a low rank update



Goldfarb, D. (1970). **A Family of Variable-Metric Methods Derived by Variational Means**. Mathematics of Computation, 24(109), 23.



# Block BFGS: Random update formulation

**Cost of update:**

$$O(d^2 \times q)$$

$$H_t = D_t \Delta_t D_t^T + (I - D_t \Delta_t Y_t^T) H_{t-1} (I - Y_t \Delta_t D_t^T),$$

where  $Y_t = \nabla^2 f_T(x_t) D_t$  and  $\Delta_t = (D_t^T Y_t)^{-1}$



RMG and Peter Richtárik (2016). **Randomized quasi-Newton updates are linearly convergent matrix inversion algorithms.** arXiv:1602.01768

# Stochastic Block BFGS Method

**Initialize**  $H_{-1} = I, w_0 \in \mathbf{R}^d$ , stepsize  $\eta > 0$

**For**  $t = 0, 1, \dots$ ,

- 1 Calculate  $g_t$
2. Sample  $T_t \subseteq [n]$ , independently
3. Form  $D_t \in \mathbf{R}^{d \times q}$
4. Compute sketch  $Y_t = \nabla^2 f_{T_t}(w_t) D_t$
5.  $H_t = D_t \Delta_t D_t^T$   
 $\quad + (I - D_t \Delta_t Y_t^T) H_{t-1} (I - Y_t \Delta_t D_t)$
6.  $d_t = H_t g_t$
7.  $w_{t+1} = w_t - \eta d_t$

**Output**  $w_{t+1}$

# Stochastic Block BFGS Method

**Initialize**  $H_{-1} = I, w_0 \in \mathbf{R}^d$ , stepsize  $\eta > 0$

**For**  $t = 0, 1, \dots$ ,

- 1 Calculate  $g_t$
2. Sample  $T_t \subseteq [n]$ , independently
3. Form  $D_t \in \mathbf{R}^{d \times q}$
4. Compute sketch  $Y_t = \nabla^2 f_{T_t}(w_t) D_t$
5.  $H_t = D_t \Delta_t D_t^T$   
 $\quad + (I - D_t \Delta_t Y_t^T) H_{t-1} (I - Y_t \Delta_t D_t)$
6.  $d_t = H_t g_t$
7.  $w_{t+1} = w_t - \eta d_t$

**Output**  $w_{t+1}$

**How to choose**  $D_t$ ?

# Stochastic Block BFGS Method

**Initialize**  $H_{-1} = I, w_0 \in \mathbf{R}^d$ , stepsize  $\eta > 0$

**For**  $t = 0, 1, \dots$ ,

- 1 Calculate  $g_t$
2. Sample  $T_t \subseteq [n]$ , independently
3. Form  $D_t \in \mathbf{R}^{d \times q}$
4. Compute sketch  $Y_t = \nabla^2 f_{T_t}(w_t) D_t$
5.  $H_t = D_t \Delta_t D_t^T + (I - D_t \Delta_t Y_t^T) H_{t-1} (I - Y_t \Delta_t D_t)$
6.  $d_t = H_t g_t$
7.  $w_{t+1} = w_t - \eta d_t$

**Output**  $w_{t+1}$

**How to choose**  $D_t$ ?

**Do we need to store**  $H_t$ ?

# Choosing the sketch matrix

$$H_t \nabla^2 f_T(x_t) D_t = D_t$$

We employ one of three strategies

- ◆ **gauss**:  $D_t \sim \mathcal{N}(0, I)$  has Gaussian entries sampled i.i.d at each iteration
- ◆ **prev** (**p**revious search directions delayed) : Let  $d_t = H_t g_t$ . Store  $q$  previous search directions  $D_t = [d_{t-q}, \dots, d_{t-1}]$ , update  $H_t$  once every  $q$  iterations
- ◆ **fact** (**f**actorized self-conditioning) : Sample the columns of a factored form  $L_t$  of  $H_t$  (*i.e.*  $L_t L_t^T = H_t$ ) uniformly at random. Fortunately we can maintain and update  $L_t$  efficiently.

# Limited Memory Block BFGS

Expanding  $M \in \mathbb{N}$  block BFGS updates gives

$$\begin{aligned} H_t &= (I - D_t \Delta_t Y_t^T) H_{t-1} (I - Y_t \Delta_t D_t^T) + D_t \Delta_t D_t^T \\ &\vdots \\ &= \text{FUNCTION} (H_{t-M}, D_t, Y_t, \Delta_t, \dots, D_{t+1-M}, Y_{t+1-M}, \Delta_{t+1-M}) \end{aligned}$$

# Limited Memory Block BFGS

Expanding  $M \in \mathbb{N}$  block BFGS updates gives

$$\begin{aligned} H_t &= (I - D_t \Delta_t Y_t^T) H_{t-1} (I - Y_t \Delta_t D_t^T) + D_t \Delta_t D_t^T \\ &\vdots \\ &= \text{FUNCTION} (H_{t-M}, D_t, Y_t, \Delta_t, \dots, D_{t+1-M}, Y_{t+1-M}, \Delta_{t+1-M}) \end{aligned}$$

$H_t$  is a function of  $H_{t+1-M}$  and  $(D_{t+1-i}, Y_{t+1-i}, \Delta_{t+1-i})$  for  $i = 1, \dots, M$

# Limited Memory Block BFGS

Expanding  $M \in \mathbb{N}$  block BFGS updates gives

$$H_t = (I - D_t \Delta_t Y_t^T) H_{t-1} (I - Y_t \Delta_t D_t^T) + D_t \Delta_t D_t^T$$

$\vdots$

$$= \text{FUNCTION} (H_{t-M}, D_t, Y_t, \Delta_t, \dots, D_{t+1-M}, Y_{t+1-M}, \Delta_{t+1-M})$$

$$= \text{FUNCTION} (D_t, Y_t, \Delta_t, \dots, D_{t+1-M}, Y_{t+1-M}, \Delta_{t+1-M})$$

$H_t$  is a function of  $H_{t+1-M}$  and  $(D_{t+1-i}, Y_{t+1-i}, \Delta_{t+1-i})$  for  $i = 1, \dots, M$

To simplify  $H_{t-M} = I$



# Limited Memory Block BFGS

**Store** the M block triples

$$(D_t, Y_t, \Delta_t), \dots, (D_{t+1-M}, Y_{t+1-M}, \Delta_{t+1-M})$$

# Limited Memory Block BFGS

**Store** the  $M$  block triples

Store  $M(2qd + q^2)$  doubles

$$(D_t, Y_t, \Delta_t), \dots, (D_{t+1-M}, Y_{t+1-M}, \Delta_{t+1-M})$$

# Limited Memory Block BFGS

**Store** the  $M$  block triples

Store  $M(2qd + q^2)$  doubles

$$(D_t, Y_t, \Delta_t), \dots, (D_{t+1-M}, Y_{t+1-M}, \Delta_{t+1-M})$$

**Calculate**  $H_t g_t$  using the following algorithm

**Two-loop recursion**

**inputs**  $g_t \in \mathbf{R}^d$ ,  $D_i, Y_i \in \mathbf{R}^{d \times q}$  and  $\Delta_i \in \mathbf{R}^{q \times q}$

**For**  $i = t, \dots, t - M + 1$

$$\alpha_i \leftarrow \Delta_i D_i^T v$$

$$v \leftarrow v - Y_i \alpha_i$$

**For**  $i = t - M + 1, \dots, t$

$$\beta_i \leftarrow \Delta_i Y_i^T v$$

$$v \leftarrow v + D_i(\alpha_i - \beta_i)$$



**output**  $H_t g_t \leftarrow v$

Costs  $Mq(4d + 2q)$  to apply

# Stochastic Block BFGS Method

**Initialize**  $H_{-1} = I, w_0 \in \mathbf{R}^d$ , stepsize  $\eta > 0$

**For**  $t = 0, 1, \dots$ ,

1. Calculate  $g_t$  
2. Sample  $T_t \subseteq [n]$ , independently
3. Form  $D_t \in \mathbf{R}^{d \times q}$
4. Compute sketch  $Y_t = \nabla^2 f_T(w_t) D_t$
5.  $d_t = H_t g_t$  
6.  $w_{t+1} = w_t - \eta d_t$

**Output**  $w_{t+1}$

Using SVRG

Two-loop recursion



Full Algorithm

# Stochastic Block BFGS Method

**Initialize**  $H_{-1} = I, w_0 \in \mathbf{R}^d$ , stepsize  $\eta > 0$

**For**  $t = 0, 1, \dots$ ,

1. Compute  $\mu = \nabla f(w_t)$

2. Set  $x_0 = w_t$

**For**  $k = 0, 1, \dots, m - 1$

3. Sample  $S_k, T_k \subseteq [n]$ , independently

4.  $g_k = \nabla f_{S_k}(x_k) - \nabla f_{S_k}(w_t) + \mu$

5. Form  $D_k \in \mathbf{R}^{d \times q}$

6.  $x_{k+1} = x_k - \eta H_k g_k$

7. **Option I:** Set  $w_{t+1} = x_m$

8. **Option I:** Set  $w_{t+1} = x_i$ , where  $i$  is selected uniformly at random from  $[m] = \{1, 2, \dots, m\}$

**Output**  $w_{t+1}$



Back

# Experiments

# Logistic regression with L2 regularizer

## Test problem

$$\min_w \sum_{i=1}^n \ln (1 + \exp(-y_i a^i, w)) + \frac{1}{n} ||w||_2^2,$$

where  $[a^1, \dots, a^n] \in \mathbf{R}^{d \times n}$  and  $y \in \{0, 1\}^n$  are the given data.

Data taken from LIBSVM

# Key to methods

SVRG

Johnson Zhang (2013)

MNJ

Moritz, Nishihara Jordan (2013)

gauss\_q\_M

Gaussian elements

fact\_q\_M

Self-conditioning factorized sampling

prev\_q\_M

Previous search directions delayed



# Key to methods

SVRG

Johnson Zhang (2013)

MNJ

Moritz, Nishihara Jordan (2013)

gauss\_q\_M

Gaussian elements

fact\_q\_M

Self-conditioning factorized sampling

prev\_q\_M

Previous search directions delayed

Stochastic Block BFGS Methods

# Key to methods

SVRG

Johnson Zhang (2013)

MNJ

Moritz, Nishihara Jordan (2013)

gauss\_q\_M

Gaussian elements

fact\_q\_M

Self-conditioning factorized sampling

prev\_q\_M

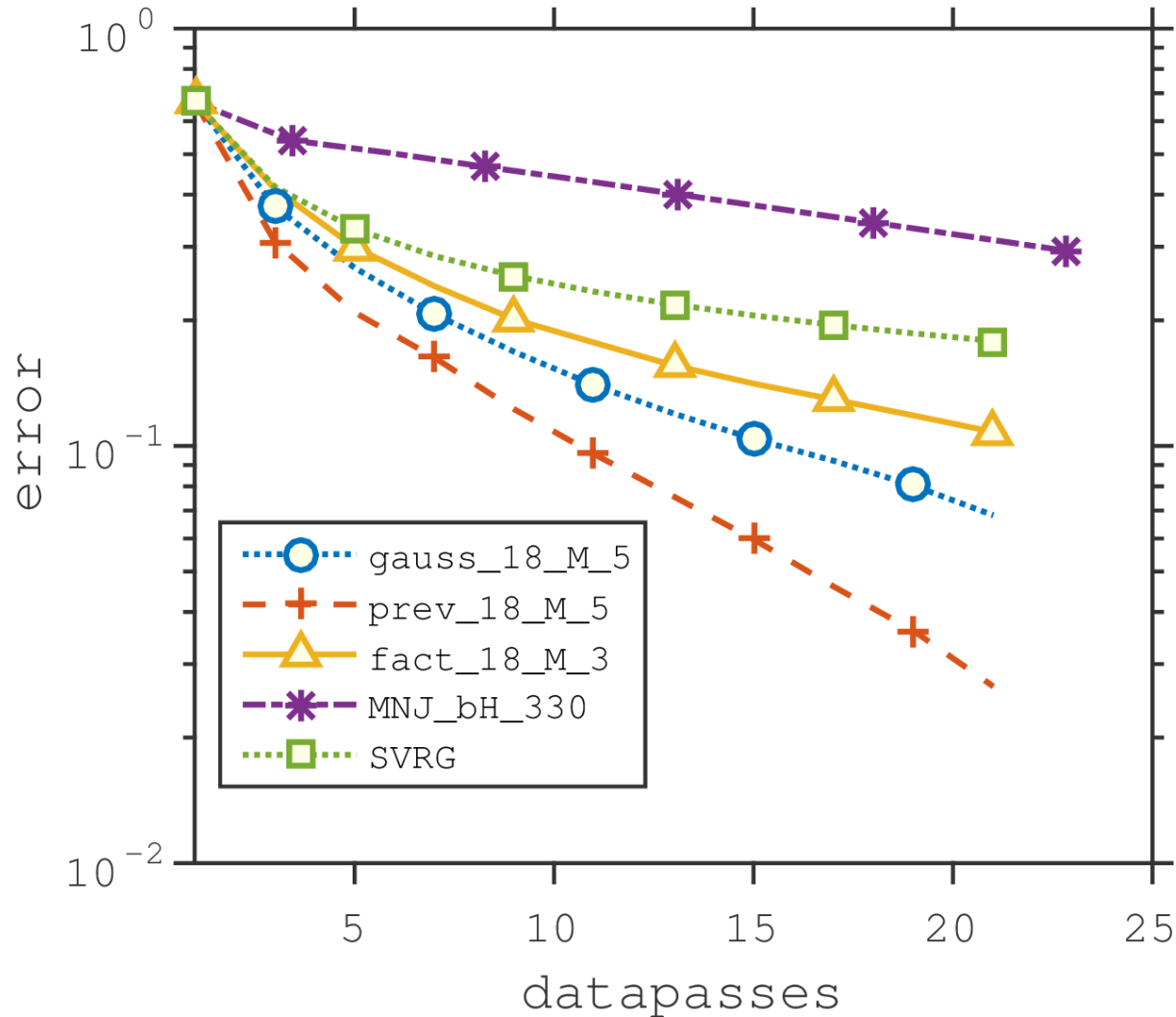
Previous search directions delayed

Stochastic Block BFGS Methods

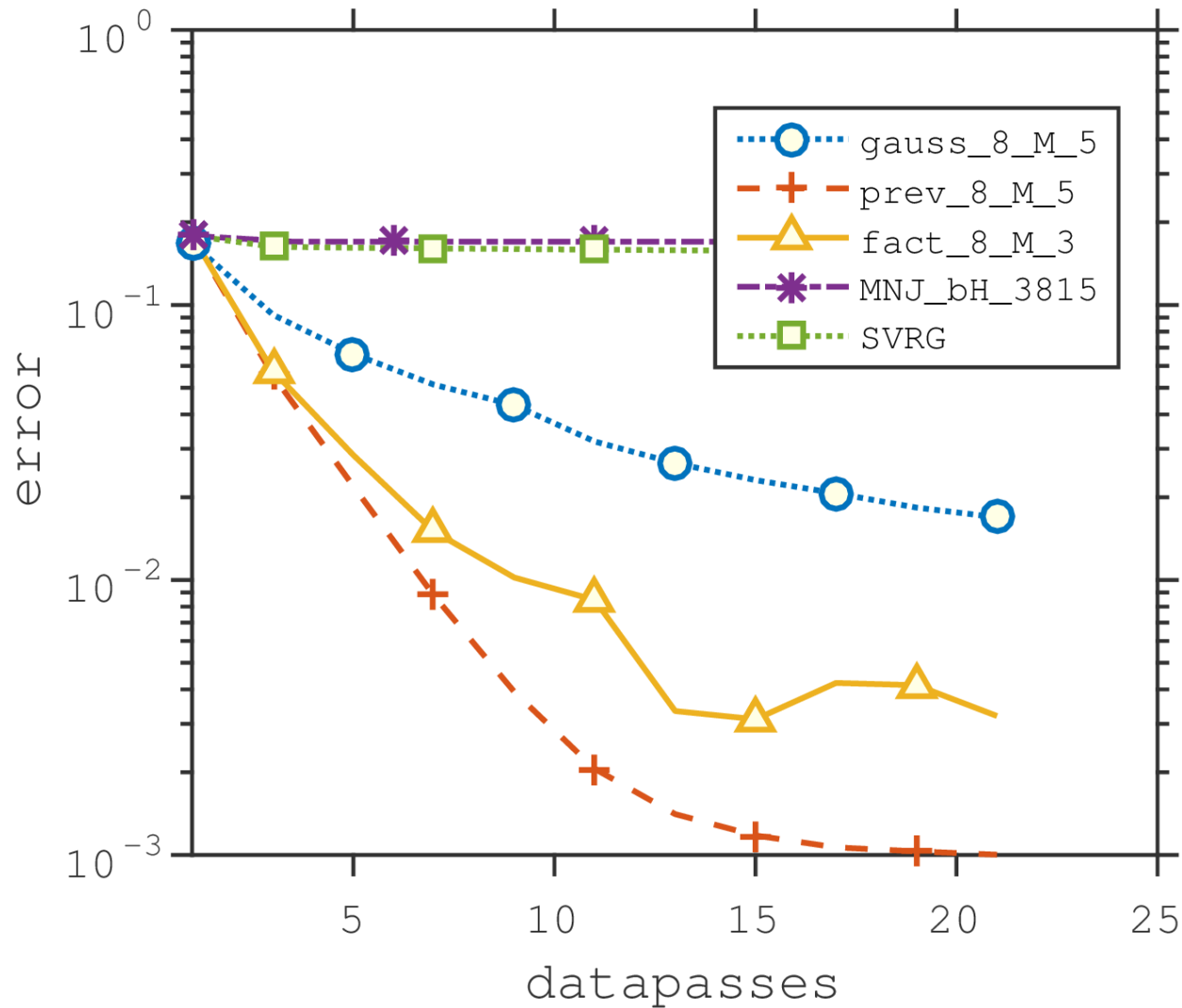
$$D_t \in \mathbf{R}^{n \times q}$$

$M \in \mathbf{N}$  number of block triples stored

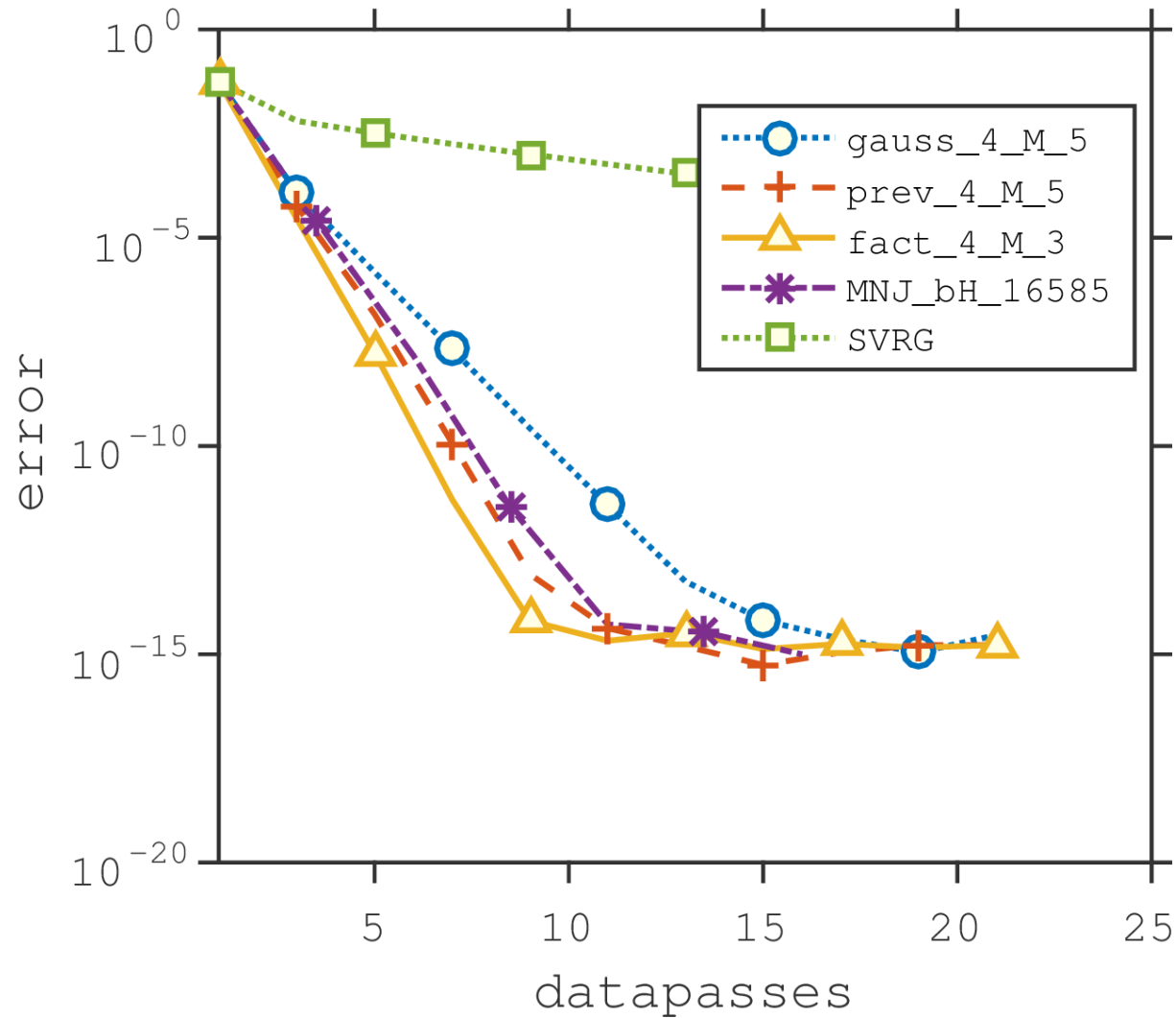
gisette,  $n = 6,000$ ,  $d = 5,000$



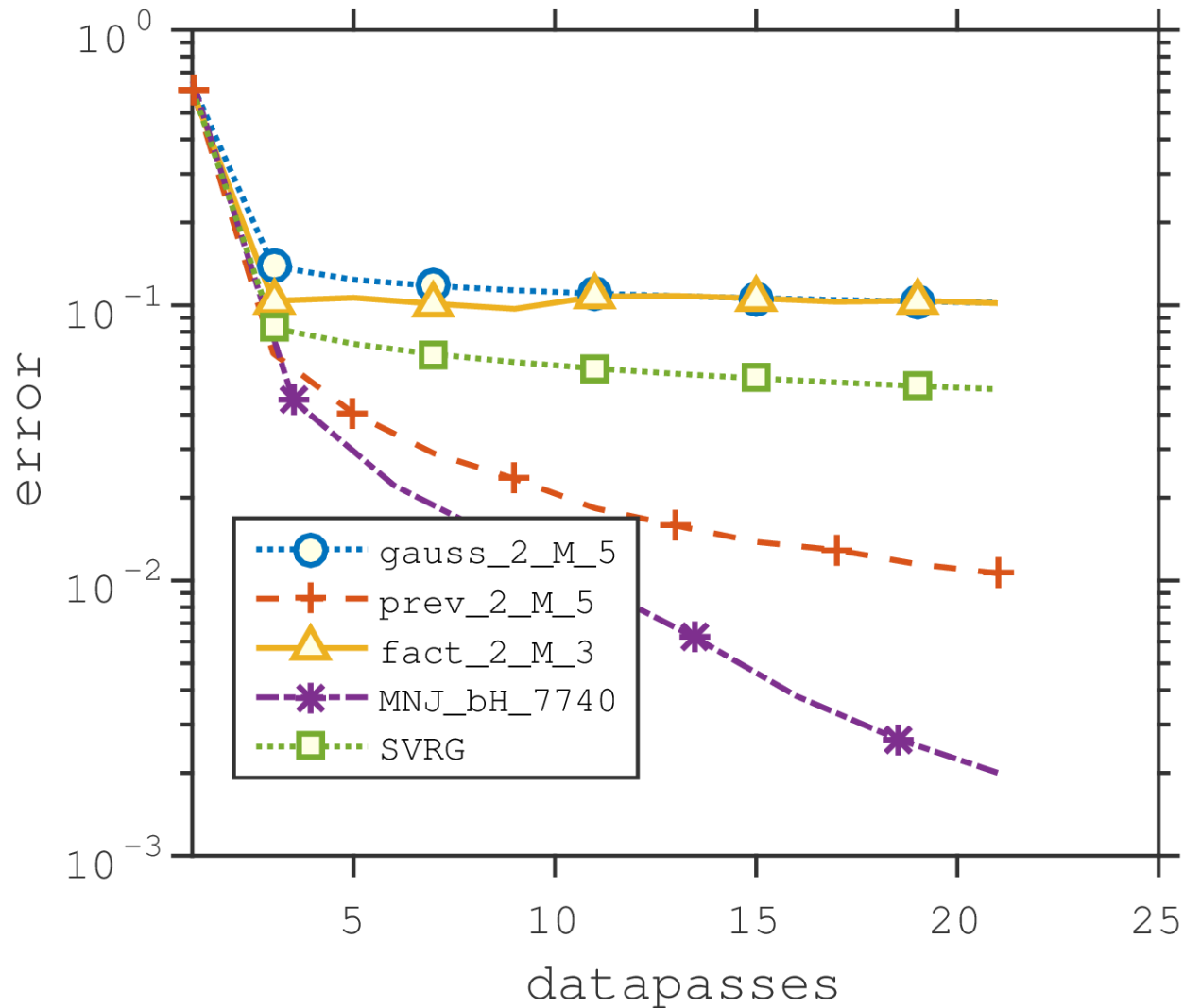
covtype.binary, n= 581,012, d= 54



Higgs,  $n=11,000,000$  ,  $d= 28$



url-combined  $n = 2,396,130$ ,  $d = 3,231,961$



# Conclusions

- ◆ New metric learning framework. A block BFGS framework for gradually learning the metric of the underlying function using sketches of subsampled Hessian matrices
- ◆ New limited memory block BFGS method. May also be of interest for non-stochastic optimization
- ◆ Several matrix sketching possibilities.
- ◆ More reasonable bounds on eigenvalues of  $H_k$  which lead to more reasonable conditions for step size, as compared to MNJ



More Numerics



Convergence results



R. Johnson and T. Zhang (2013).  
**Accelerating stochastic gradient descent using predictive variance reduction.** NIPS.



P. Moritz, R. Nishihara and M. I. Jordan  
(2016). **A Linearly-Convergent Stochastic L-BFGS Algorithm**, AISTATS.



RMG and Peter Richtárik (2016)  
**Randomized Quasi-Newton Updates are Linearly Convergent Matrix Inversion Algorithms**  
arXiv:1602.01768



# Convergence

# Convergence

## Assumption

There exists  $0 < \lambda \leq \Lambda$  such that

$$\lambda I \preceq \nabla^2 f_T(x) \preceq \Lambda I$$

For all  $x \in \mathbf{R}^d$  and all  $T \subseteq [n]$ .

## Lemma [GGR'16]

There exists  $0 < \gamma \leq \Gamma$  such that

$$\gamma I \preceq H_t \preceq \Gamma I, \quad \forall t$$

Furthermore

$$\frac{1}{1 + M\Lambda} \leq \gamma \leq \Gamma \leq (1 + \sqrt{\kappa})^{2M} \left(1 + \frac{1}{\lambda(2\sqrt{\kappa} + \kappa)}\right)$$

where  $\kappa = \Lambda/\lambda$

# Complexity / Convergence

## Theorem [GGR'16]

If

$$m \geq \frac{1}{2\eta (\gamma\lambda - \eta\Gamma^2\Lambda(2\Lambda - \lambda))} \quad \eta < \gamma\lambda/(2\Gamma^2\Lambda^2)$$

# Complexity / Convergence

## Theorem [GGR'16]

If

$$m \geq \frac{1}{2\eta (\gamma\lambda - \eta\Gamma^2\Lambda(2\Lambda - \lambda))} \quad \eta < \gamma\lambda/(2\Gamma^2\Lambda^2)$$

# Complexity / Convergence

## Theorem [GGR'16]

If

$$m \geq \frac{1}{2\eta (\gamma\lambda - \eta\Gamma^2\Lambda(2\Lambda - \lambda))} \quad \eta < \gamma\lambda/(2\Gamma^2\Lambda^2)$$



Inner iterations  
of SVRG

# Complexity / Convergence

## Theorem [GGR'16]

If

$$m \geq \frac{1}{2\eta (\gamma\lambda - \eta\Gamma^2\Lambda(2\Lambda - \lambda))}$$

Inner iterations  
of SVRG

Step size

$$\eta < \gamma\lambda / (2\Gamma^2\Lambda^2)$$

# Complexity / Convergence

## Theorem [GGR'16]

If

$$m \geq \frac{1}{2\eta(\gamma\lambda - \eta\Gamma^2\Lambda(2\Lambda - \lambda))}$$

Inner iterations  
of SVRG

Step size

$$\eta < \gamma\lambda / (2\Gamma^2\Lambda^2)$$

$$\mathbf{E}[f(w_t) - f(w_*)] \leq \rho^t (f(w_0) - f(w_*))$$

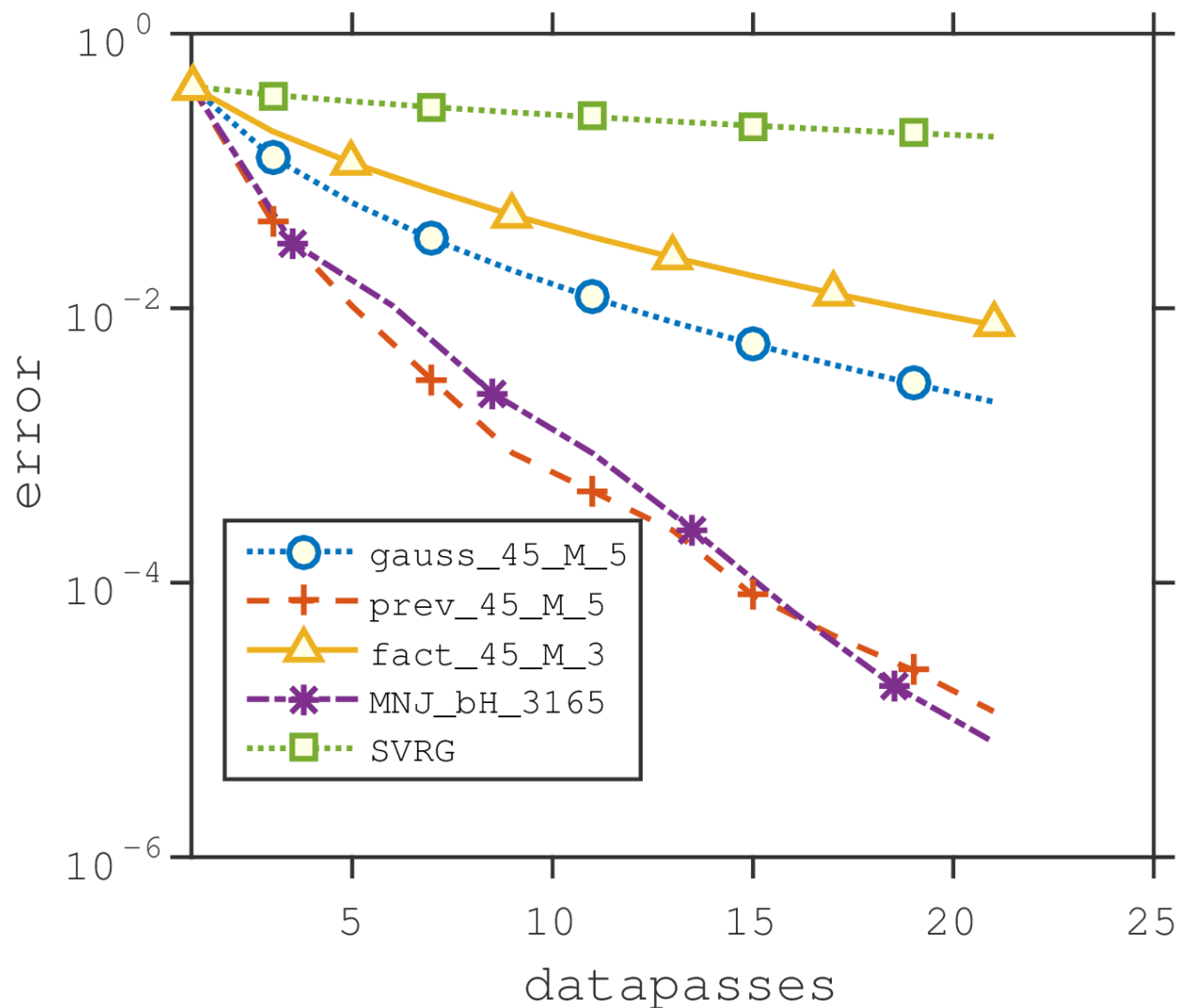
Where,

$$\rho = \frac{1/2m\eta + \eta\Gamma^2\Lambda(\Lambda - \lambda)}{\gamma\lambda - \eta\Gamma^2\Lambda^2} < 1$$

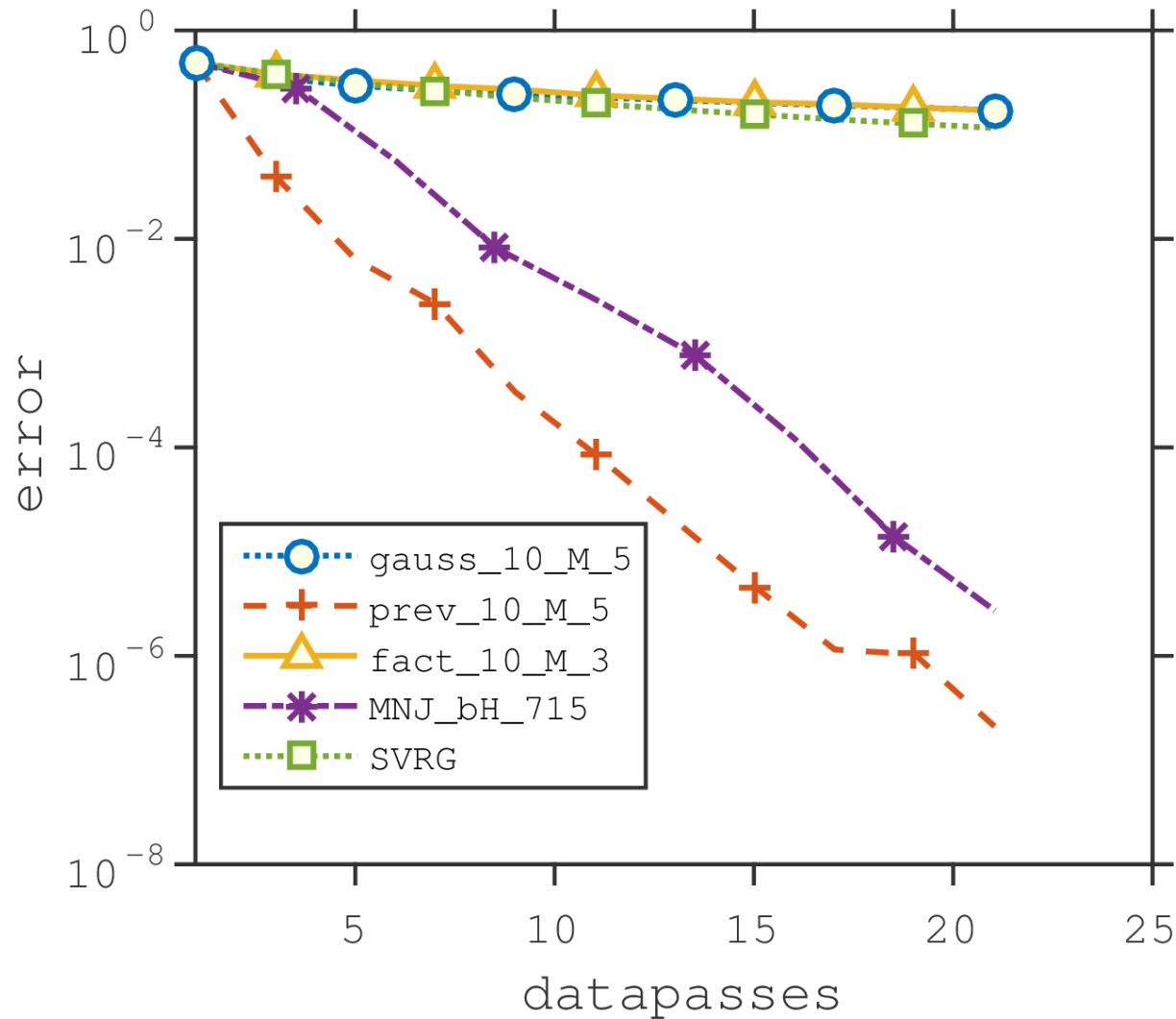
Experimental results  
error  $\times$  datapasses



epsilon\_normalized n= 400,000 , d=2,000

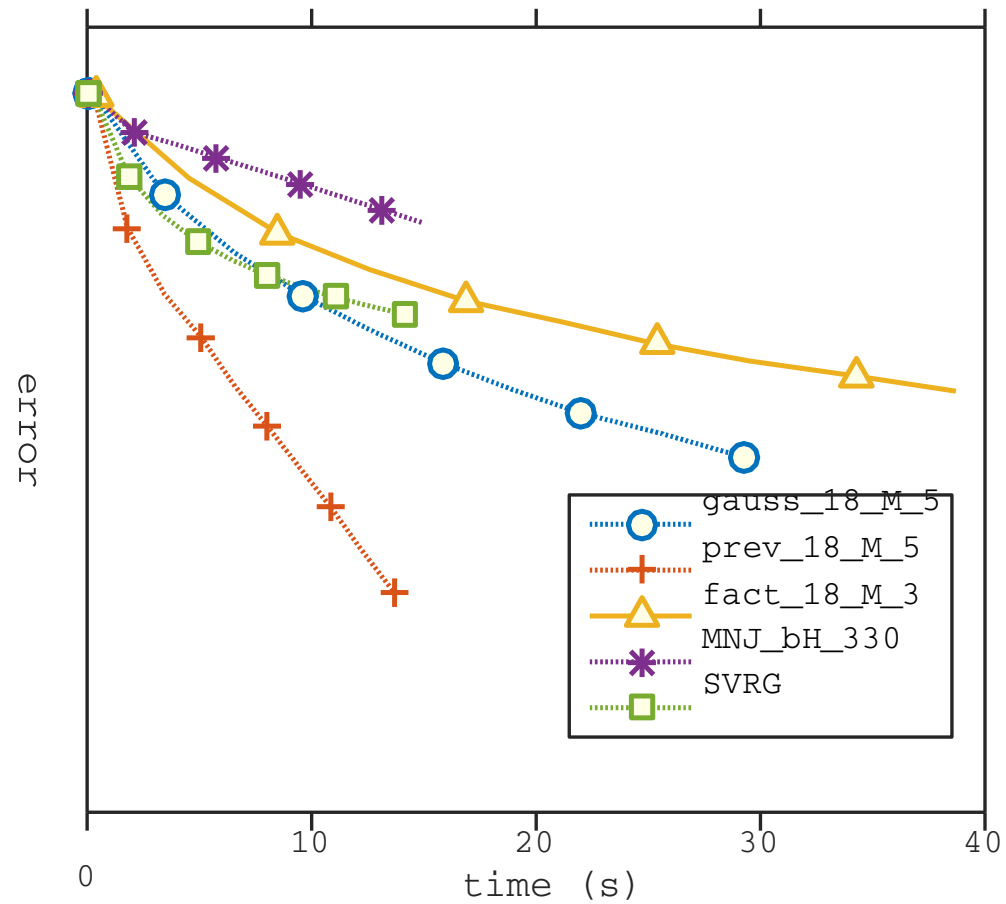


rcv1-training  $n = 20,242$ ,  $d = 47,236$

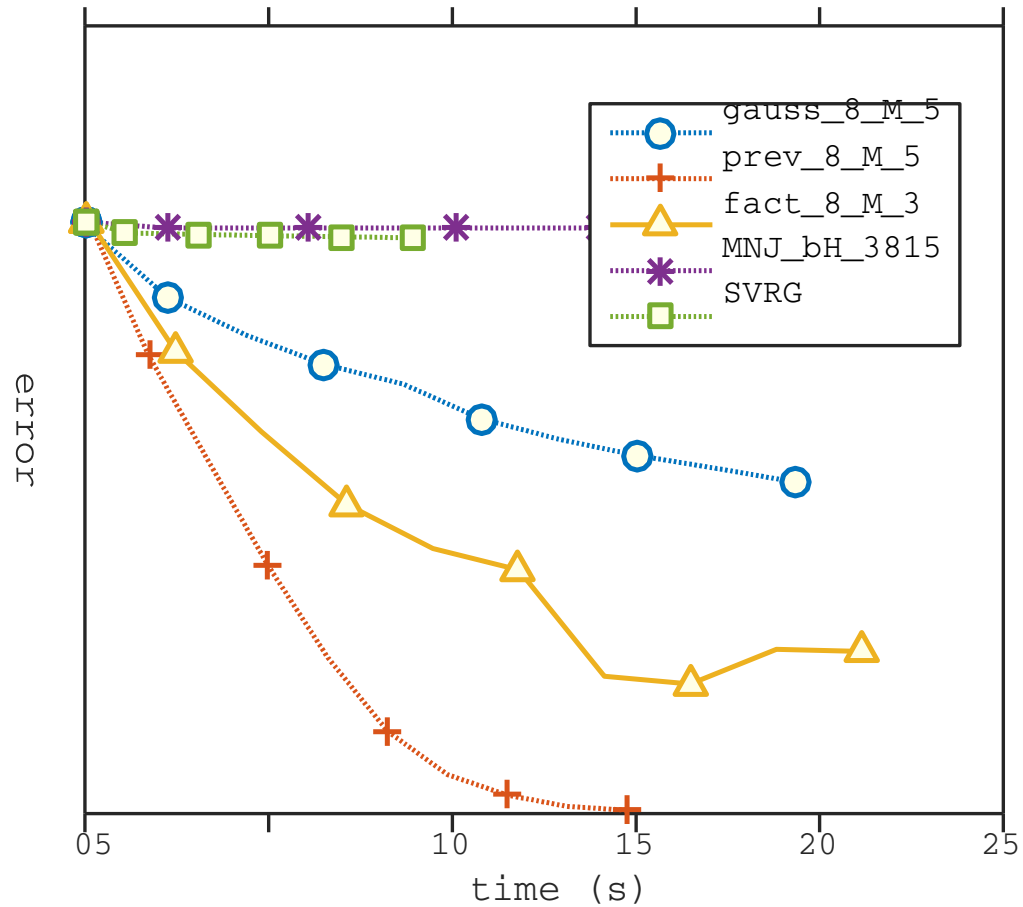


Experimental results  
error X time

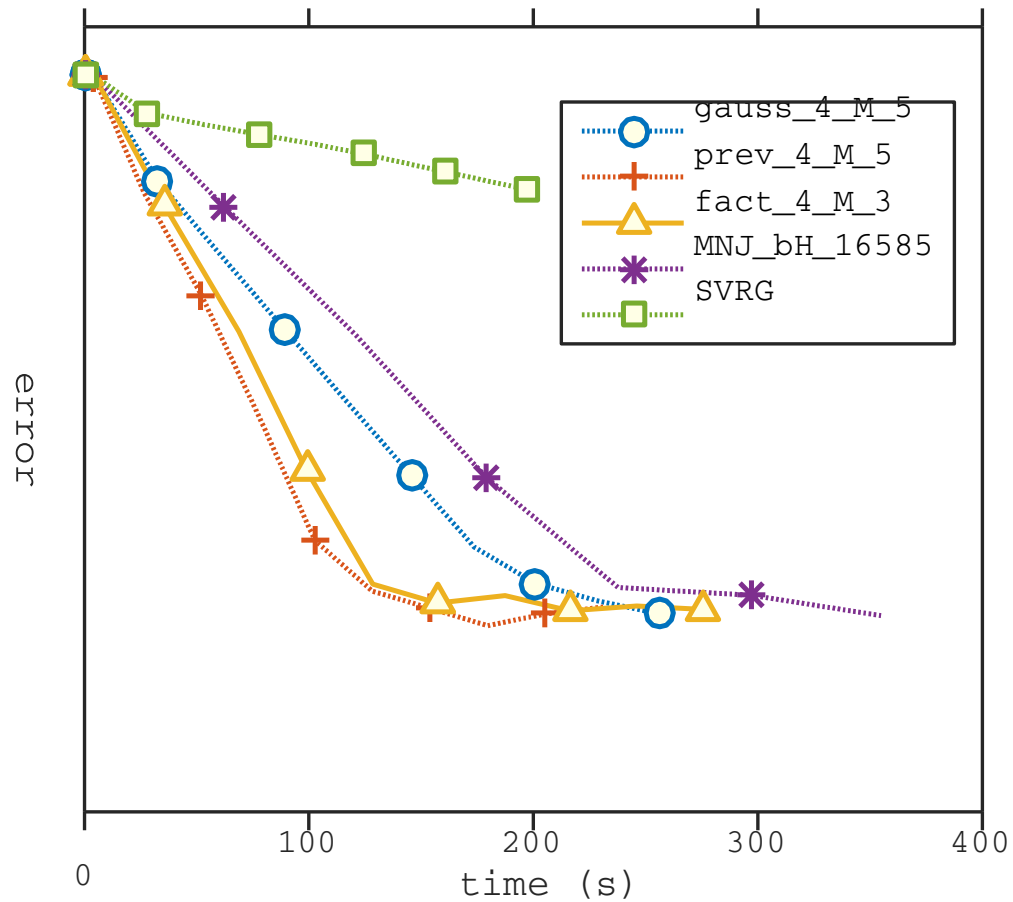
gisette,  $n = 6,000$ ,  $d = 5,000$



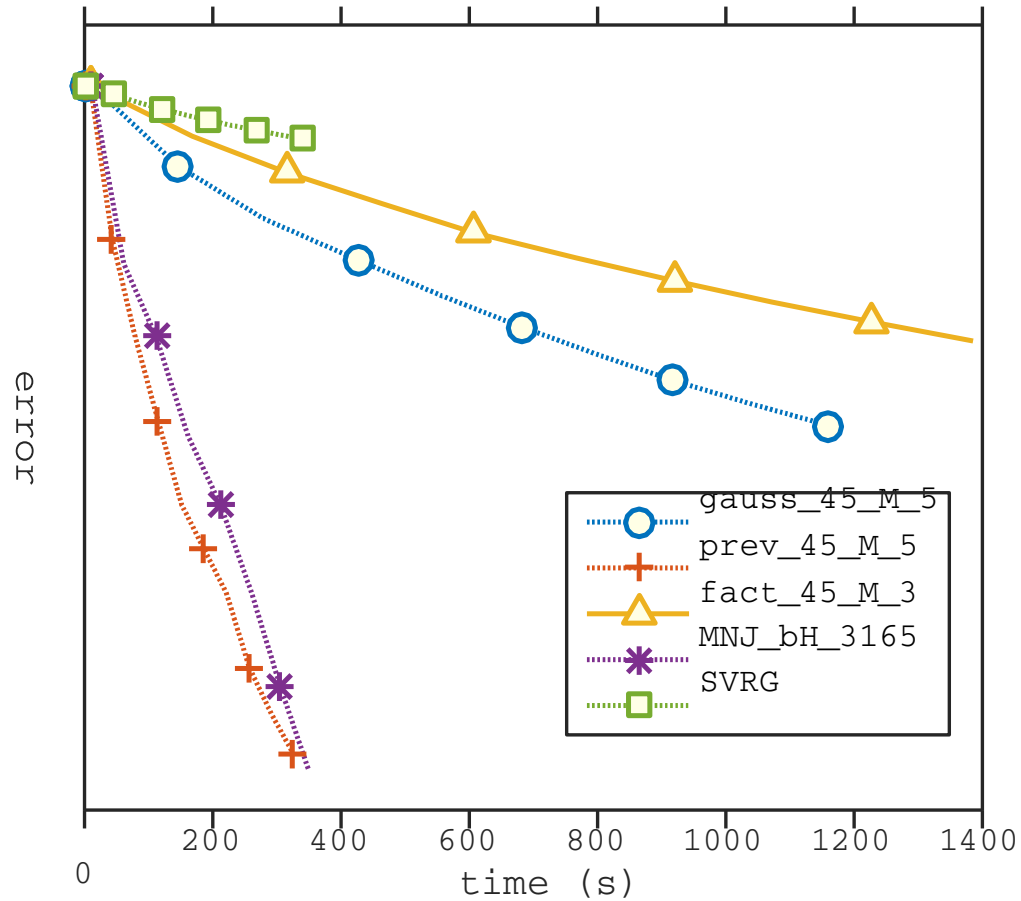
covtype.binary, n= 581,012, d= 54



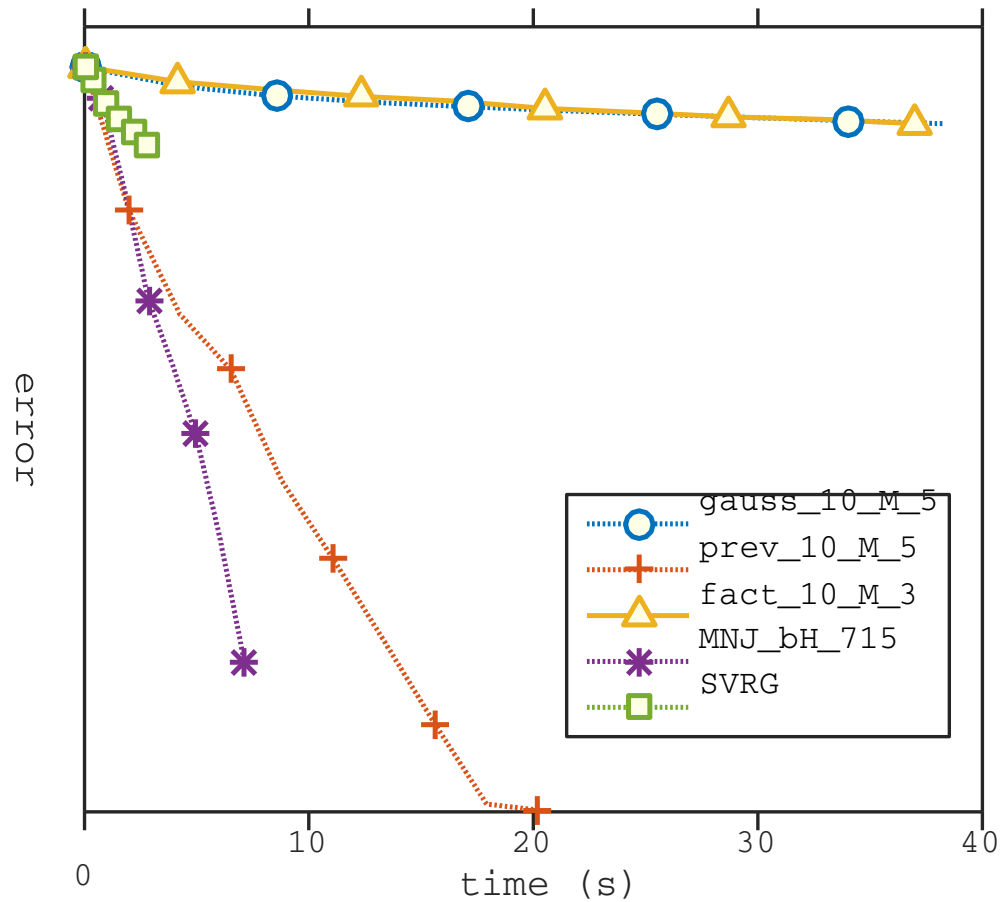
Higgs,  $n=11,000,000$  ,  $d= 28$



epsilon\_normalized n= 400,000 , d=2,000

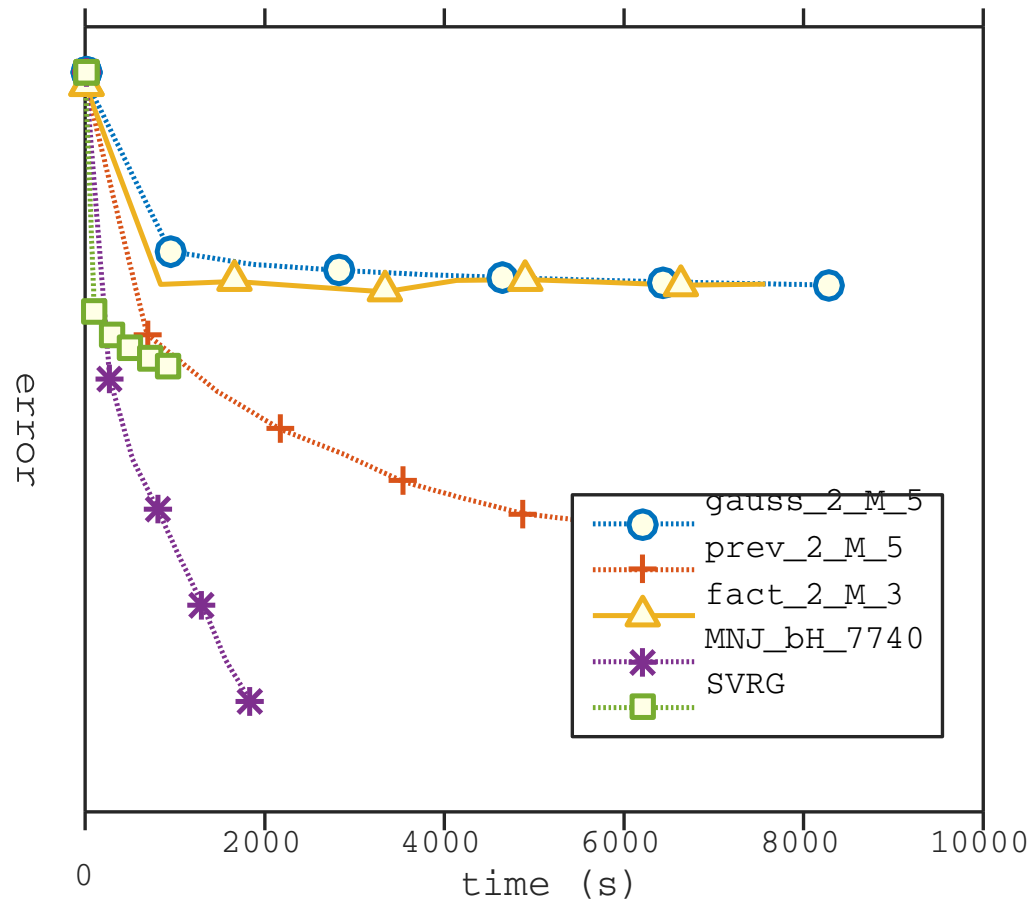


rcv1-training  $n = 20,242$ ,  $d = 47,236$



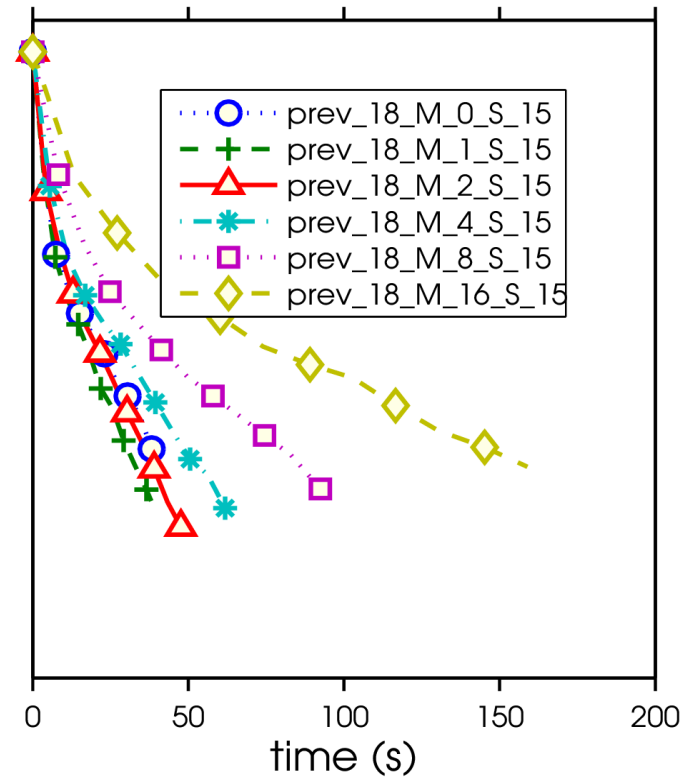
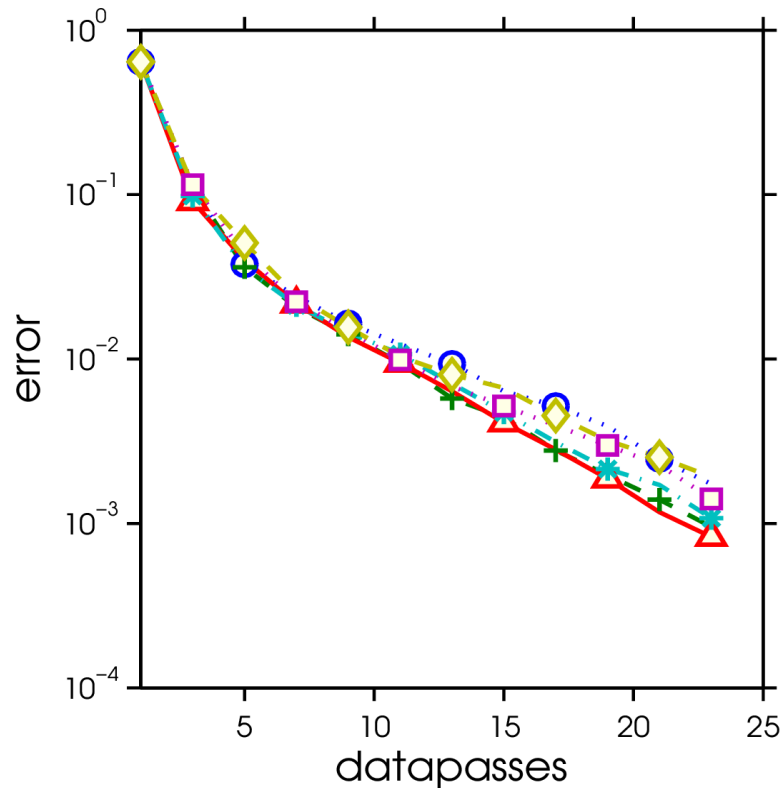


url-combined  $n = 2,396,130$ ,  $d = 3,231,961$



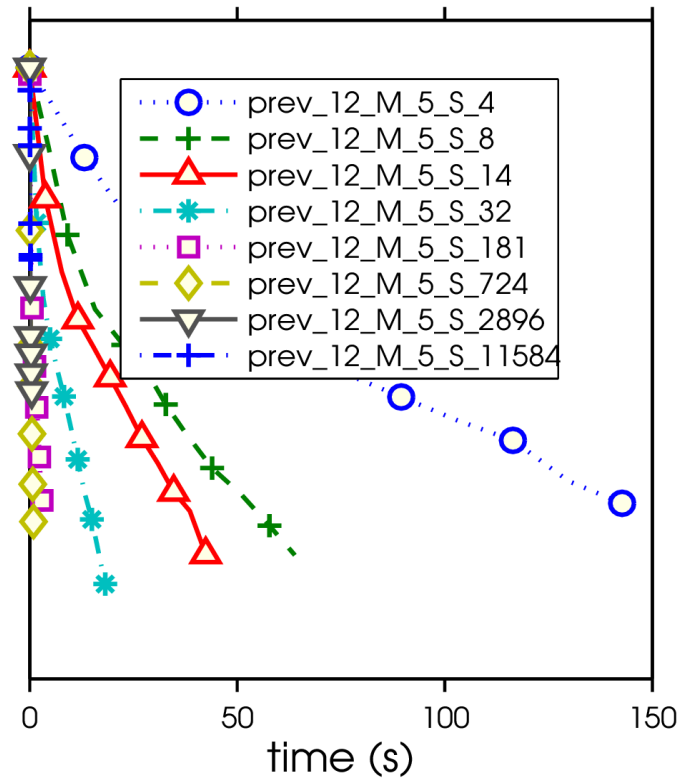
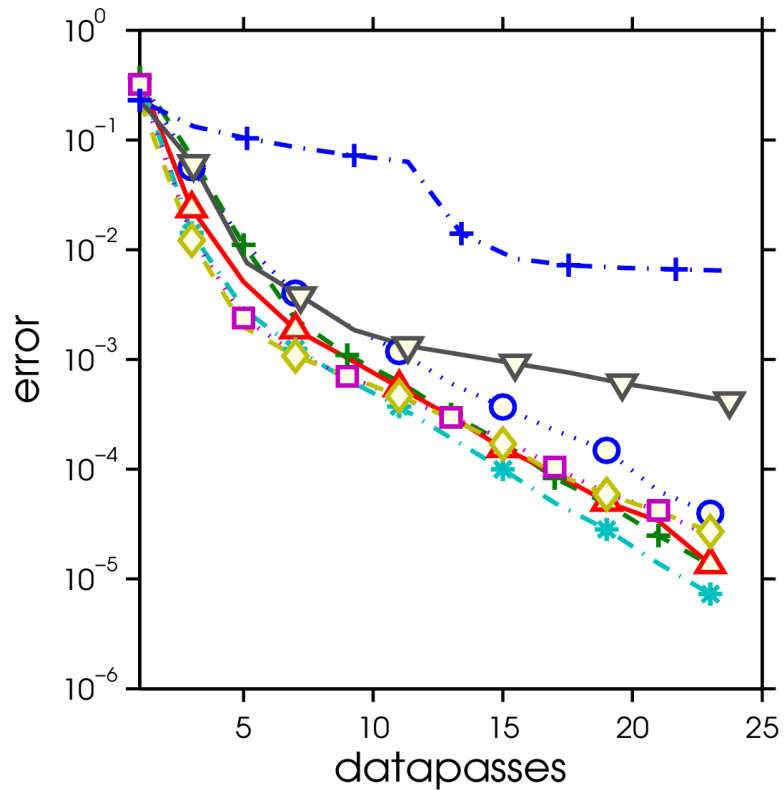
Experimental results  
parameter exploration

w8a  $n = 49,749$ ,  $d = 300$



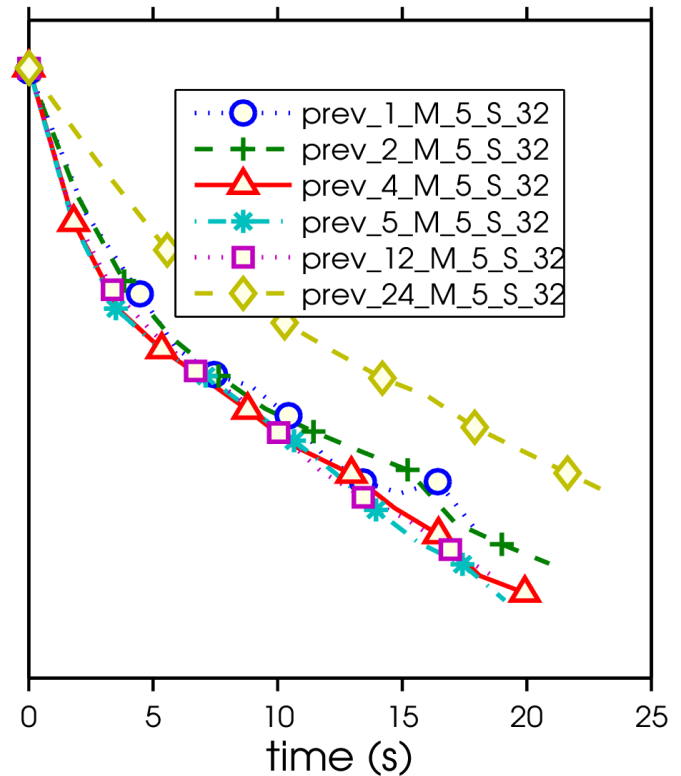
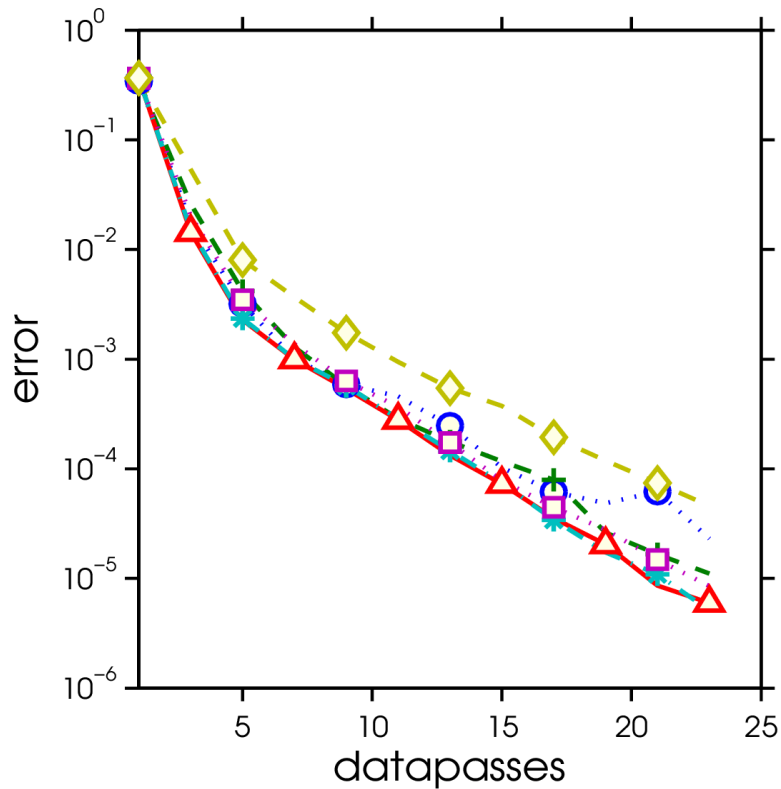
Testing the memory parameter  $M$

a9a  $n = 32,561$ ,  $d = 123$



Testing the subsampling size with  $S=T$

a9a  $n = 32,561$ ,  $d = 123$



Testing update size (delay size)  $q$

**SVRG**

# The Stochastic Variance Reduced Gradient

$$g_t = \nabla f_S(w_t) - \nabla f_S(x_k) + \nabla f(x_k)$$

Where  $x_k$  is a reference point.

$$\begin{aligned}\text{Unbiased : } \mathbf{E}_S[g_t] &= \mathbf{E}[\nabla f_S(w_t)] - \mathbf{E}[\nabla f_S(x_k)] + \nabla f(x_k) \\ &= \nabla f(w_t) + \nabla f(x_k) - \nabla f(x_k) \\ &= \nabla f(w_t)\end{aligned}$$

Maintain  $x_k$  fixed and  
iterate in  $t$  for  $m$  iterations



R. Johnson and T. Zhang (2013). **Accelerating stochastic gradient descent using predictive variance reduction**. NIPS, 1(3), 1-9.