



signProx: 1-bit Proximal Algorithm for Nonconvex Stochastic Optimization

Xiaojian Xu¹, Ulugbek Kamilov^{1,2}

**¹ Computer Science and Engineering,
Washington University in St. Louis**

**² Electrical and Systems Engineering,
Washington University in St. Louis**

A common inverse problem

- A common inverse problem model

True signal

$$\boldsymbol{x}^*$$

Forward model

$$\boldsymbol{H}$$

Observation

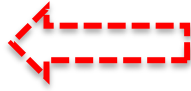
$$\boldsymbol{y} = \boldsymbol{H}\boldsymbol{x}^*$$

A common inverse problem

- A common inverse problem model

True signal

$$\mathbf{x}^*$$



Forward model

$$\mathbf{H}$$

Observation

$$\mathbf{y} = \mathbf{H}\mathbf{x}^*$$

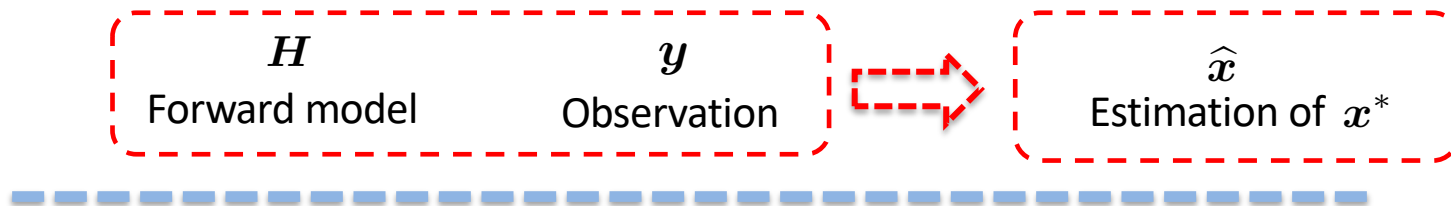
The inverse problem could be formulated as an optimization problem

- A feasible optimization framework to solve the problem



The inverse problem could be formulated as an optimization problem

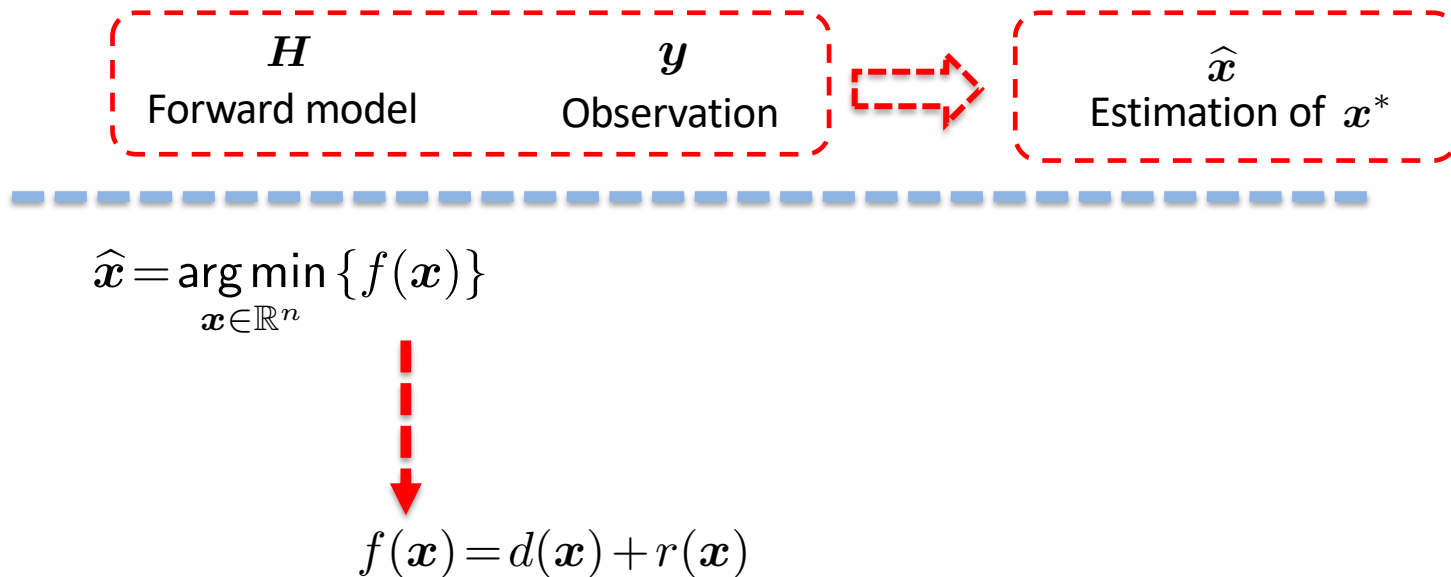
- A feasible optimization framework to solve the problem



$$\hat{x} = \arg \min_{x \in \mathbb{R}^n} \{f(x)\}$$

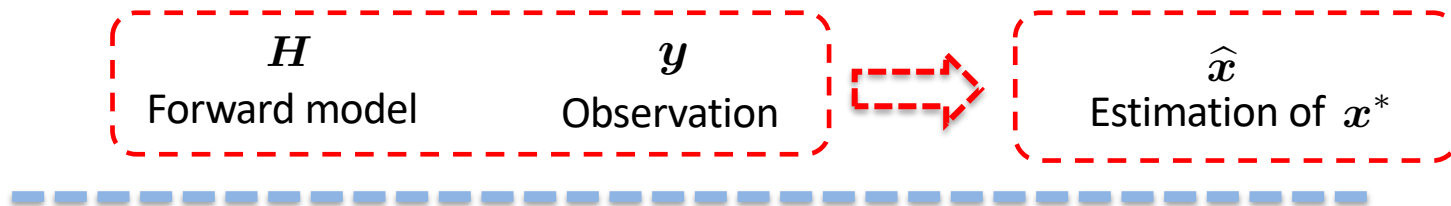
The inverse problem could be formulated as an optimization problem

- A feasible optimization framework to solve the problem



The inverse problem could be formulated as an optimization problem

- A feasible optimization framework to solve the problem



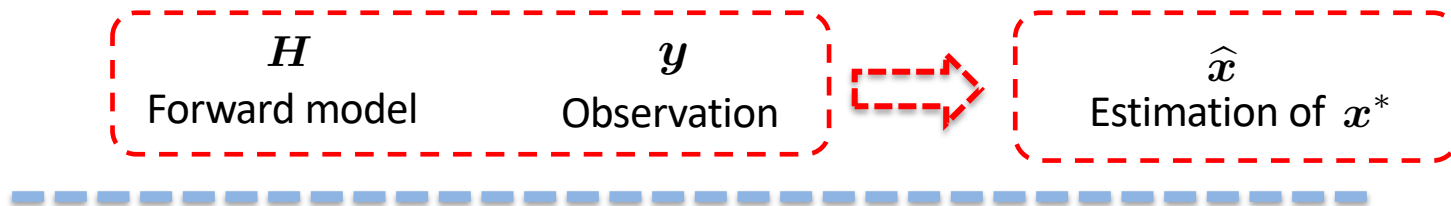
$$\hat{x} = \arg \min_{x \in \mathbb{R}^n} \{f(x)\}$$

A diagram showing the decomposition of the objective function $f(x)$. A red dashed arrow points from the x in the minimization problem above to the x in the equation $f(x) = d(x) + r(x)$. A red dashed arrow labeled "Data-fidelity term" points from the text to the $d(x)$ term in the equation. The equation is followed by a curly brace and the definition of $d(x)$.

$$f(x) = d(x) + r(x) \left\{ \begin{array}{l} d(x) = \frac{1}{2} \|y - Hx\|_2^2 \quad \dots \end{array} \right.$$

The inverse problem could be formulated as an optimization problem

- A feasible optimization framework to solve the problem



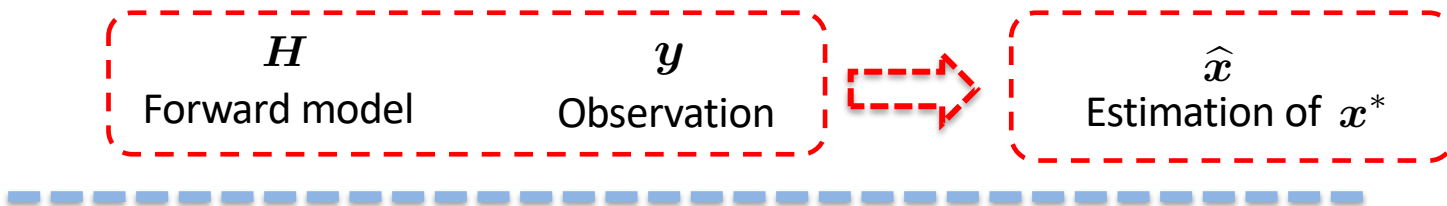
$$\hat{x} = \arg \min_{x \in \mathbb{R}^n} \{f(x)\}$$

A diagram showing the decomposition of the cost function $f(x)$. A red dashed arrow points from the $f(x)$ in the equation above to the $f(x)$ in the equation below. A red dashed arrow labeled "Data-fidelity term" points from the $d(x)$ term to its definition. Another red dashed arrow labeled "Regularizer term" points from the $r(x)$ term to its definition. The equation is:

$$f(x) = d(x) + r(x) \quad \left\{ \begin{array}{l} d(x) = \frac{1}{2} \|y - Hx\|_2^2 \quad \dots \\ r(x) = \end{array} \right.$$

The inverse problem could be formulated as an optimization problem

- A feasible optimization framework to solve the problem



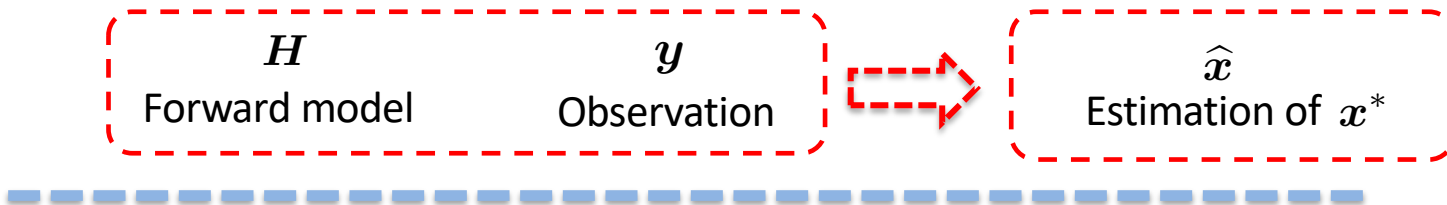
$$\hat{x} = \arg \min_{x \in \mathbb{R}^n} \{f(x)\}$$

A diagram showing the decomposition of the objective function $f(x)$. A red dashed arrow points from the $f(x)$ in the optimization problem above to the $f(x)$ in the equation below. The equation is $f(x) = d(x) + r(x)$. A red dashed arrow labeled "Data-fidelity term" points from the text to the $d(x)$ term. Another red dashed arrow labeled "Regularizer term" points from the text to the $r(x)$ term. The terms are further defined as $d(x) = \frac{1}{2} \|y - Hx\|_2^2 \dots$ and $r(x) = \|x\|_1$.

$$f(x) = d(x) + r(x) \begin{cases} d(x) = \frac{1}{2} \|y - Hx\|_2^2 \dots \\ r(x) = \|x\|_1 \end{cases}$$

The inverse problem could be formulated as an optimization problem

- A feasible optimization framework to solve the problem



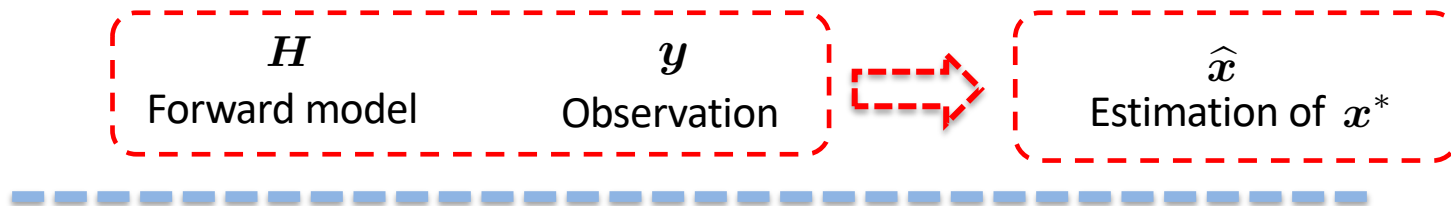
$$\hat{x} = \arg \min_{x \in \mathbb{R}^n} \{f(x)\}$$

A diagram showing the decomposition of the objective function $f(x)$. A red dashed arrow points from the $f(x)$ in the optimization problem above to the $f(x)$ in the equation $f(x) = d(x) + r(x)$. A red dashed arrow labeled "Data-fidelity term" points from the $d(x)$ term to the equation $d(x) = \frac{1}{2} \|y - Hx\|_2^2 \dots$. Another red dashed arrow labeled "Regularizer term" points from the $r(x)$ term to the equation $r(x) = \|x\|_1 \quad \|x\|_2$.

$$f(x) = d(x) + r(x) \begin{cases} d(x) = \frac{1}{2} \|y - Hx\|_2^2 & \dots \\ r(x) = \|x\|_1 \quad \|x\|_2 \end{cases}$$

The inverse problem could be formulated as an optimization problem

- A feasible optimization framework to solve the problem



$$\hat{x} = \arg \min_{x \in \mathbb{R}^n} \{ f(x) \}$$

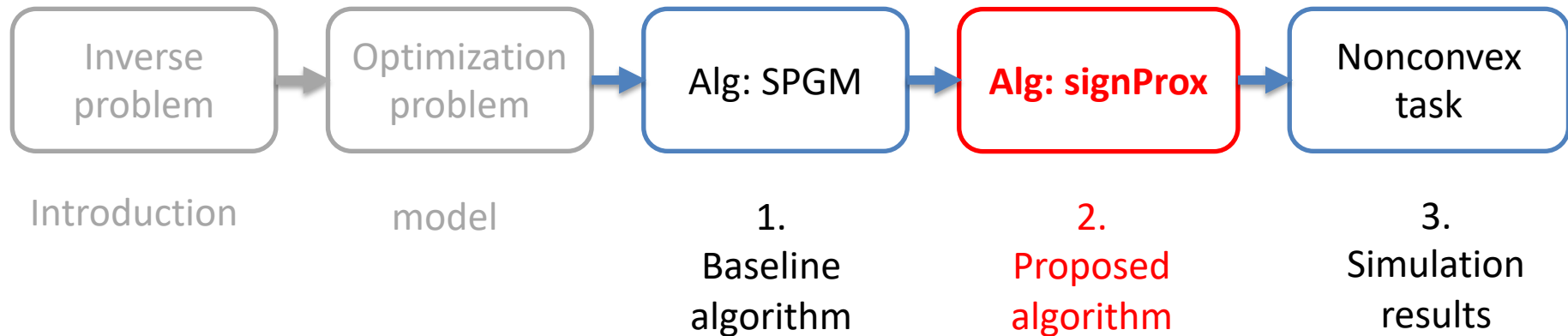
A diagram showing the decomposition of the objective function $f(x)$. A red dashed arrow points from the $f(x)$ in the equation above to the $f(x)$ in the equation below. The equation is $f(x) = d(x) + r(x)$. A red dashed arrow points from the text "Data-fidelity term" to the $d(x)$ term. Another red dashed arrow points from the text "Regularizer term" to the $r(x)$ term. The terms are defined as follows:

$$\left\{ \begin{array}{l} d(x) = \frac{1}{2} \|y - Hx\|_2^2 \quad \dots \\ r(x) = \|x\|_1 + \|x\|_2 \quad \dots \end{array} \right.$$

Guideline of the talk

- Large scale optimization problem

$$\hat{\mathbf{x}} = \arg \min_{\mathbf{x} \in \mathbb{R}^n} \{f(\mathbf{x})\} \quad \text{with} \quad f(\mathbf{x}) = d(\mathbf{x}) + \frac{1}{K} \{r_1(\mathbf{x}) + \cdots + r_k(\mathbf{x}) + \cdots + r_K(\mathbf{x})\}$$



Proximal gradient method could minimize the non-differentiable functions

- Modeled optimization problem

$$\hat{\boldsymbol{x}} = \arg \min_{\boldsymbol{x} \in \mathbb{R}^n} \{f(\boldsymbol{x})\} \quad \text{with} \quad f(\boldsymbol{x}) = d(\boldsymbol{x}) + r(\boldsymbol{x})$$

Proximal gradient method could minimize the non-differentiable functions

- Modeled optimization problem

$$\hat{\boldsymbol{x}} = \arg \min_{\boldsymbol{x} \in \mathbb{R}^n} \{f(\boldsymbol{x})\} \quad \text{with} \quad f(\boldsymbol{x}) = d(\boldsymbol{x}) + r(\boldsymbol{x})$$



- Proximal gradient method (PGM)

$$\boldsymbol{x}^t \leftarrow \text{prox}_{\gamma r}(\boldsymbol{x}^{t-1} - \gamma \nabla d(\boldsymbol{x}^{t-1}))$$

Proximal gradient mapping

Proximal gradient method could minimize the non-differentiable functions

- Modeled optimization problem

$$\hat{\mathbf{x}} = \arg \min_{\mathbf{x} \in \mathbb{R}^n} \{f(\mathbf{x})\} \quad \text{with} \quad f(\mathbf{x}) = d(\mathbf{x}) + r(\mathbf{x})$$



- Proximal gradient method (PGM)

$$\mathbf{x}^t \leftarrow \text{prox}_{\gamma r}(\mathbf{x}^{t-1} - \gamma \nabla d(\mathbf{x}^{t-1})) \quad \left\{ \begin{array}{l} \mathbf{s} = \mathbf{x}^{t-1} - \gamma \nabla d(\mathbf{x}^{t-1}) \end{array} \right.$$

Proximal gradient mapping

Proximal gradient method could minimize the non-differentiable functions

- Modeled optimization problem

$$\hat{\mathbf{x}} = \arg \min_{\mathbf{x} \in \mathbb{R}^n} \{f(\mathbf{x})\} \quad \text{with} \quad f(\mathbf{x}) = d(\mathbf{x}) + r(\mathbf{x})$$



- Proximal gradient method (PGM)

$$\mathbf{x}^t \leftarrow \text{prox}_{\gamma r}(\mathbf{x}^{t-1} - \gamma \nabla d(\mathbf{x}^{t-1})) \quad \left\{ \begin{array}{l} \mathbf{s} = \mathbf{x}^{t-1} - \gamma \nabla d(\mathbf{x}^{t-1}) \\ \mathbf{x}^t = \text{prox}_{\gamma r}(\mathbf{s}) \triangleq \arg \min_{\mathbf{x} \in \mathbb{R}^n} \left\{ \frac{1}{2} \|\mathbf{x} - \mathbf{s}\|_2^2 + \gamma r(\mathbf{x}) \right\} \end{array} \right.$$

Proximal gradient mapping

Proximal operation

Proximal average is a good estimation of the true proximal gradient mapping

- Large scale optimization problem

$$\hat{\mathbf{x}} = \arg \min_{\mathbf{x} \in \mathbb{R}^n} \{f(\mathbf{x})\} \quad \text{with} \quad \begin{cases} f(\mathbf{x}) = \frac{1}{K} \sum_{k=1}^K f_k(\mathbf{x}) \\ f_k(\mathbf{x}) = d(\mathbf{x}) + r_k(\mathbf{x}) \end{cases}$$

- Bauschke, Heinz H., et al. "The proximal average: basic theory." SIAM Journal on Optimization 19.2 (2008): 766-785.
- Yu, Yao-Liang. "Better approximation and faster algorithm using the proximal average." Advances in neural information processing systems. 2013.

Proximal average is a good estimation of the true proximal gradient mapping

- Large scale optimization problem

$$\hat{\mathbf{x}} = \arg \min_{\mathbf{x} \in \mathbb{R}^n} \{f(\mathbf{x})\} \quad \text{with} \quad \begin{cases} f(\mathbf{x}) = \frac{1}{K} \sum_{k=1}^K f_k(\mathbf{x}) \\ f_k(\mathbf{x}) = d(\mathbf{x}) + r_k(\mathbf{x}) \end{cases}$$




- Proximal average

$$P(\mathbf{x}) \triangleq \frac{1}{K} \sum_{k=1}^K P_k(\mathbf{x}) \quad \text{with} \quad P_k(\mathbf{x}) \triangleq \text{prox}_{\gamma r_k}(\mathbf{x} - \gamma \nabla d(\mathbf{x})), \quad k \in [1, \dots, K]$$

- Bauschke, Heinz H., et al. "The proximal average: basic theory." SIAM Journal on Optimization 19.2 (2008): 766-785.
- Yu, Yao-Liang. "Better approximation and faster algorithm using the proximal average." Advances in neural information processing systems. 2013.

Stochastic proximal gradient method could approximate the proximal average

- Stochastic proximal gradient method

$$P(x) \triangleq \frac{1}{K} \sum_{k=1}^K P_k(x)$$


Could be time consuming and
computational resource
demanding for K components

- H. Robbins and S. Monro, "A stochastic approximation method," The Annals of Mathematical Statistics, vol. 22, no. 3, pp. 400–407, September 1951.

Stochastic proximal gradient method could approximate the proximal average

- Stochastic proximal gradient method

$$P(\mathbf{x}) \triangleq \frac{1}{K} \sum_{k=1}^K P_k(\mathbf{x}) \quad \approx \quad \hat{P}(\mathbf{x}) \triangleq \frac{1}{B} \sum_{b=1}^B P_{k_b}(\mathbf{x})$$

Could be time consuming and
computational resource
demanding for K components

$$B \ll K$$

- H. Robbins and S. Monro, "A stochastic approximation method," The Annals of Mathematical Statistics, vol. 22, no. 3, pp. 400–407, September 1951.

Stochastic proximal gradient method could approximate the proximal average

- Stochastic proximal gradient method

$$\boxed{P(x)} \triangleq \frac{1}{K} \sum_{k=1}^K P_k(x) \approx \boxed{\hat{P}(x)} \triangleq \frac{1}{B} \sum_{b=1}^B P_{k_b}(x)$$

Could be time consuming and computational resource demanding for K components

$B \ll K$

- H. Robbins and S. Monro, "A stochastic approximation method," The Annals of Mathematical Statistics, vol. 22, no. 3, pp. 400–407, September 1951.

Stochastic proximal gradient method could approximate the proximal average

- Stochastic proximal gradient method

$$\boxed{P(x)} \triangleq \frac{1}{K} \sum_{k=1}^K P_k(x) \approx \boxed{\hat{P}(x)} \triangleq \frac{1}{B} \sum_{b=1}^B P_{k_b}(x)$$

Could be time consuming and computational resource demanding for K components

$B \ll K$

$$x^t \leftarrow \hat{P}(x^{t-1})$$

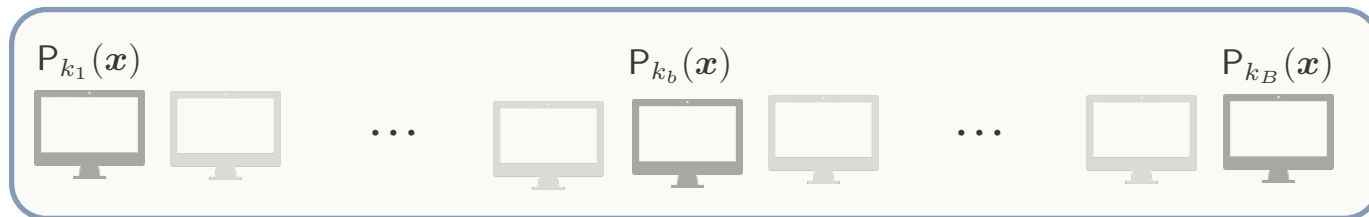
Stochastic proximal gradient method (SPGM)

- H. Robbins and S. Monro, "A stochastic approximation method," The Annals of Mathematical Statistics, vol. 22, no. 3, pp. 400–407, September 1951.

SPGM could work parallelly

$$f(\boldsymbol{x}) = \frac{1}{K} \{ f_1(\boldsymbol{x}) + f_2(\boldsymbol{x}) + \cdots + f_K(\boldsymbol{x}) \}$$

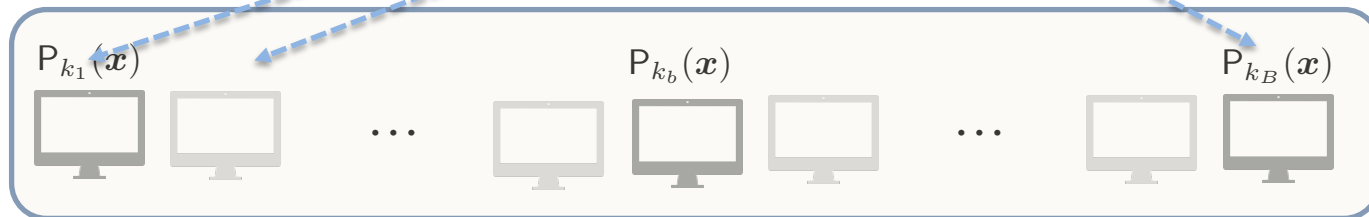
...



SPGM could work parallelly

$$f(x) = \frac{1}{K} \{ \boxed{f_1(x)} + \boxed{f_2(x)} + \cdots + \boxed{f_K(x)} \}$$

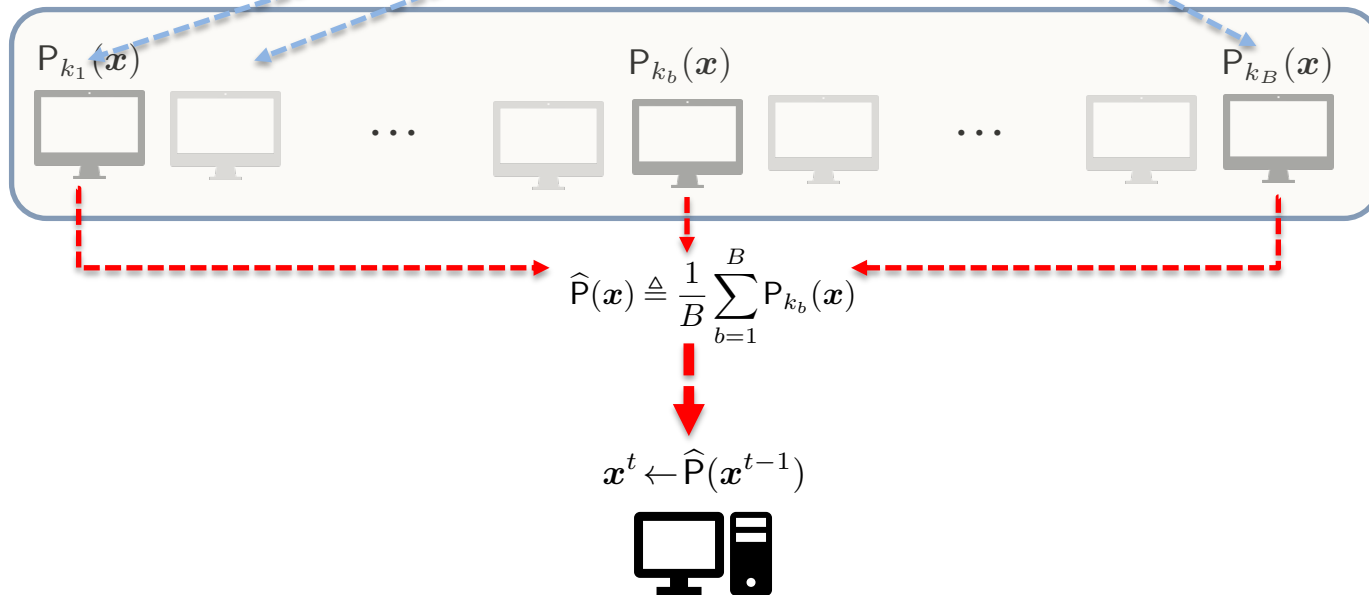
...



SPGM could work parallelly

$$f(x) = \frac{1}{K} \{ \boxed{f_1(x)} + \boxed{f_2(x)} + \cdots + \boxed{f_K(x)} \}$$

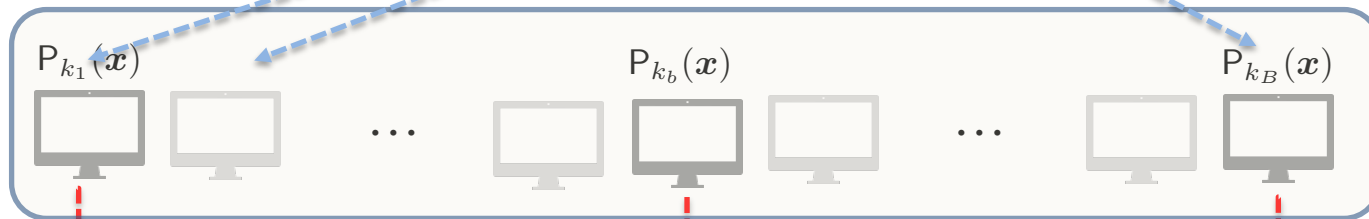
...



SPGM could work parallelly

$$f(x) = \frac{1}{K} \{ f_1(x) + f_2(x) + \cdots + f_K(x) \}$$

...



$$\hat{P}(x) \triangleq \frac{1}{B} \sum_{b=1}^B P_{k_b}(x)$$

Transmission of \hat{P}
is frequent

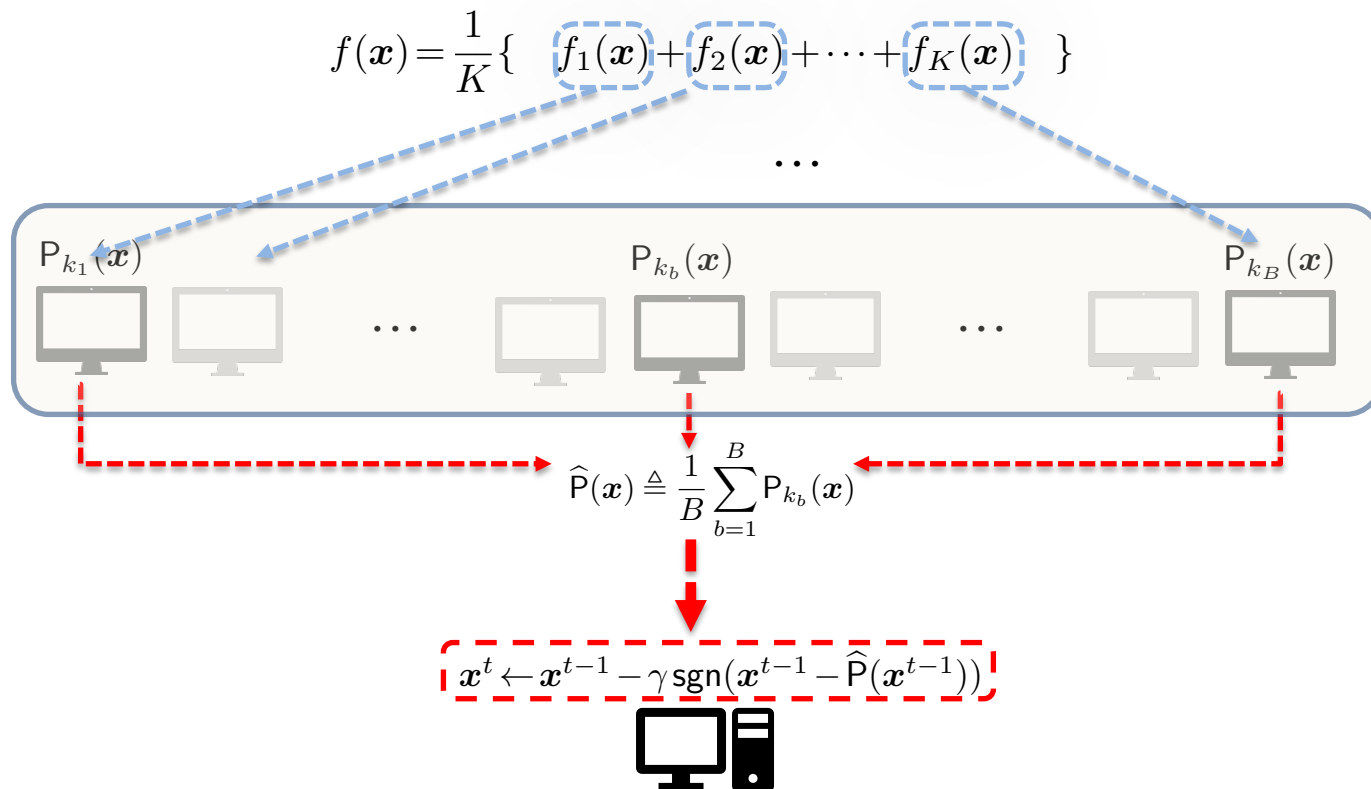
Dimension of \hat{P}
is high

$$x^t \leftarrow \hat{P}(x^{t-1})$$



Significant communication cost !!!

signProx is more efficient than SPGM



- Seide, Frank, et al. "1-bit stochastic gradient descent and its application to data-parallel distributed training of speech dnns." Fifteenth Annual Conference of the International Speech Communication Association. 2014.
- J. Bernstein, Y.-X. Wang, K. Azizzadenesheli, and A. Anandkumar, "signSGN: Compressed optimization for non-convex problems," in Proc. 35th Int. Conf. Machine Learning (ICML), Stockholm, Sweden, July 2018.

SPGM uses the true direction to update while signProx only uses the sign

- Stochastic proximal gradient method (SPGM)

» Update rule:

$$\mathbf{x}^t \leftarrow \hat{\mathbf{P}}(\mathbf{x}^{t-1})$$

-
- 1-bit stochastic proximal gradient method (signProx)

» Update rule:

$$\mathbf{x}^t \leftarrow \mathbf{x}^{t-1} - \gamma \operatorname{sgn}(\mathbf{x}^{t-1} - \hat{\mathbf{P}}(\mathbf{x}^{t-1}))$$

signProx could achieve the comparable performance with SPGM

- Convergence rate tells how fast you can reduce your loss function

» Convergence conclusion for **signProx**:

$$\mathbb{E} \left[\frac{1}{T} \sum_{t=1}^T \|\mathbf{x}^{t-1} - \mathbf{P}(\mathbf{x}^{t-1})\|_1 \right] \leq \frac{1}{\sqrt{T}} C_1 = \mathcal{O}\left(\frac{1}{\sqrt{T}}\right)$$

» Convergence conclusion for **SPGM**:

$$\mathbb{E} \left[\frac{1}{T} \sum_{t=1}^T \|\mathbf{x}^{t-1} - \mathbf{P}(\mathbf{x}^{t-1})\|_2^2 \right] \leq \frac{1}{\sqrt{T}} C_2 = \mathcal{O}\left(\frac{1}{\sqrt{T}}\right)$$

Check the performance of SPGM and signProx on a phase retrieval model

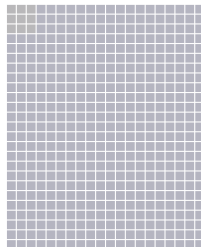
- **Nonconvex** phase retrieval task

$$\mathbf{x}^* \in \mathbb{R}^{2500}$$



True image

$$\mathbf{H} \in \mathbb{R}^{3000 \times 2500}$$



Forward model

$$\mathbf{y} = |\mathbf{H}\mathbf{x}^*|^2 \in \mathbb{R}^{3000}$$



Observation

Check the performance of SPGM and signProx on a phase retrieval model

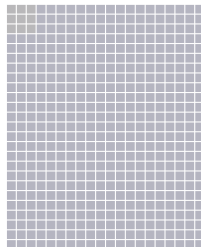
- **Nonconvex** phase retrieval task

$$\mathbf{x}^* \in \mathbb{R}^{2500}$$



True image

$$\mathbf{H} \in \mathbb{R}^{3000 \times 2500}$$



Forward model

$$\mathbf{y} = |\mathbf{H}\mathbf{x}^*|^2 \in \mathbb{R}^{3000}$$



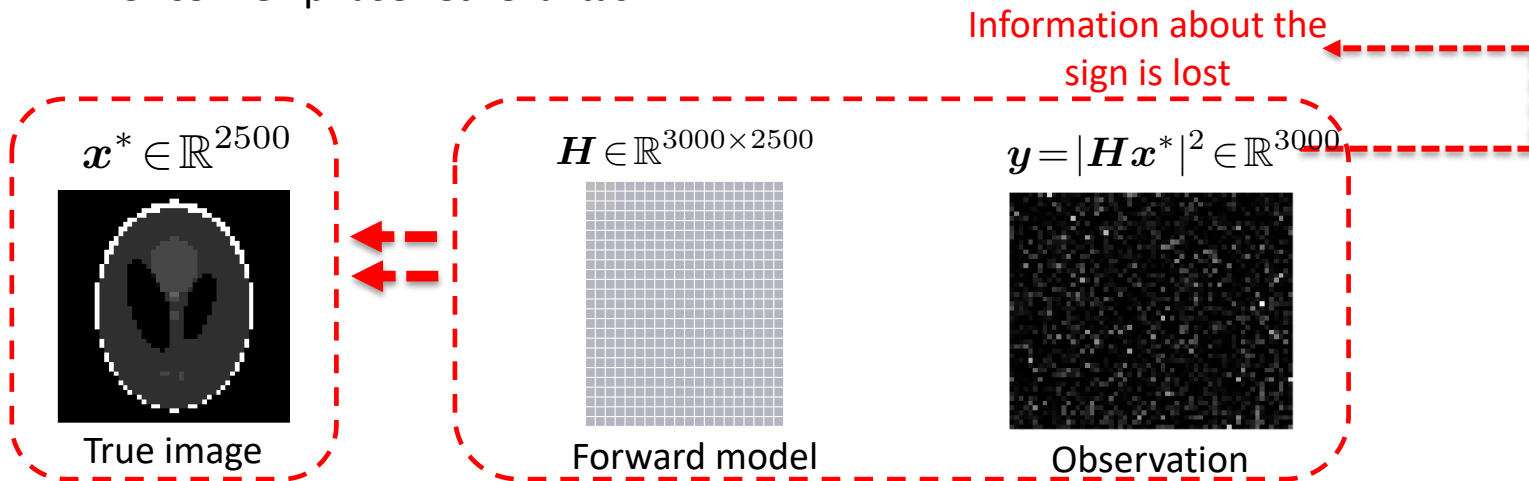
Observation

Information about the
sign is lost



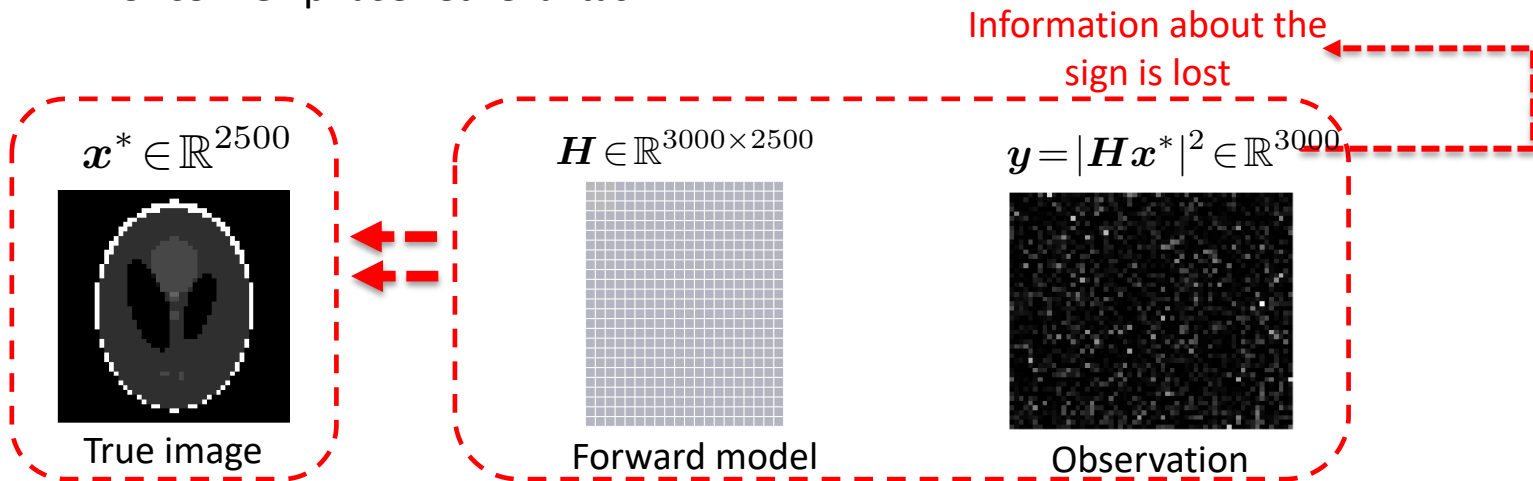
Check the performance of SPGM and signProx on a phase retrieval model

- **Nonconvex** phase retrieval task



Check the performance of SPGM and signProx on a phase retrieval model

- **Nonconvex** phase retrieval task



- Optimization objective function

$$\hat{x} = \arg \min_{x \in \mathbb{R}^n} \left\{ \frac{1}{2} \|y - |Hx|^2\|_2^2 + \text{TV}(x) \right\}$$

Simulate the stochasticity of the proximal gradient mapping by adding noise

- Stochastic simulation of proximal gradient mapping

$$\hat{P} = \underbrace{P}_{\text{True proximal mapping}} + \overset{\text{Noise}}{p(e)}$$

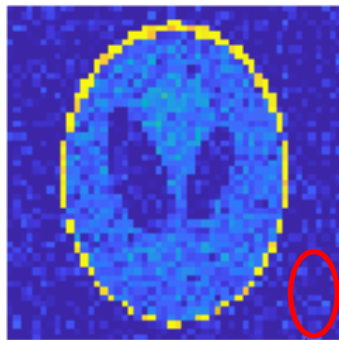
True proximal mapping

P



Dense stochasticity

\hat{P}_1



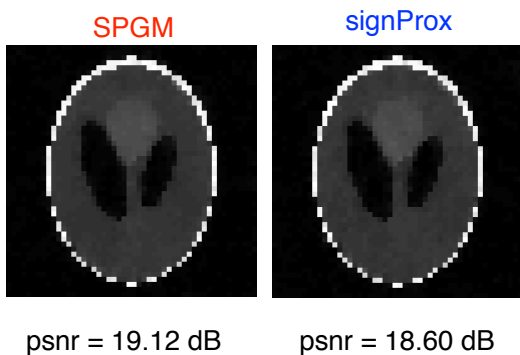
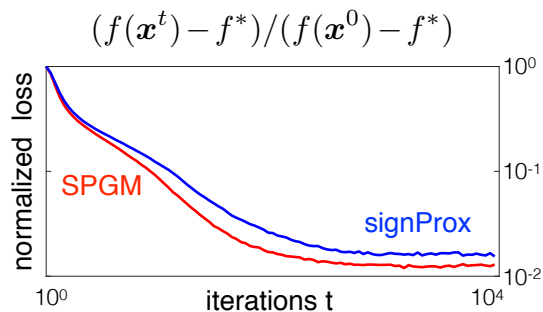
Sparse stochasticity

\hat{P}_2



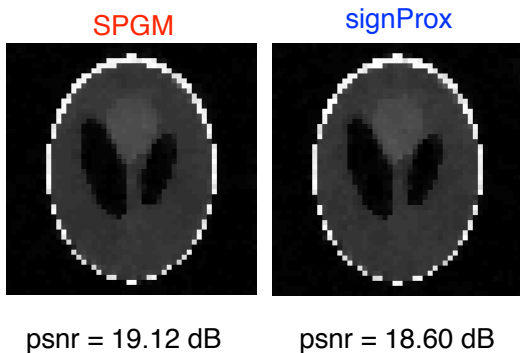
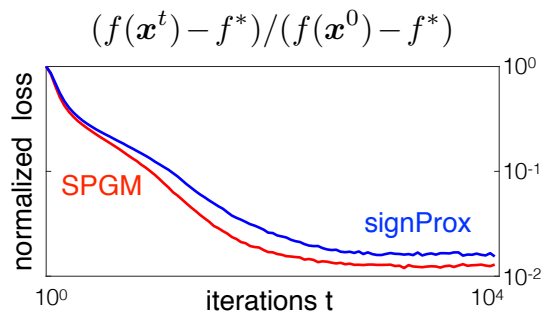
signProx outperforms SPGM in the some sparse stochasticity scenario

Dense stochasticity scenario

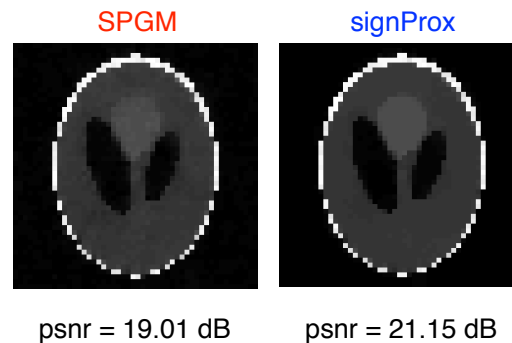
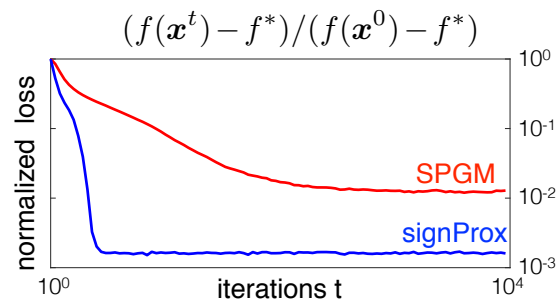


signProx outperforms SPGM in the some sparse stochasticity scenario

Dense stochasticity scenario



Sparse stochasticity scenario

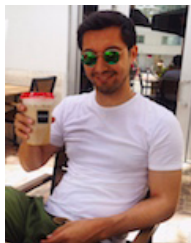


Conclusions

- » Proposed a compressed proximal gradient method signProx to solve the low efficiency problem of SPGM in a large scale optimization scenario.
- » Proved the convergence of the signProx under nonconvex assumption and showed it achieves the comparable theoretical performance with SPGM.
- » Simulated a phase retrieval problem and showed signProx has a comparable performance with SPGM and in some scenario it even outperforms SPGM.

Acknowledgements

- Support
 - » National Science Foundation under Grant No. 1813910.
- Lab homepage
 - » <https://cigroup.wustl.edu/>
- Follow us
 - » <https://twitter.com/wustlcig>



Ulugbek Kamilov



Xiaojian Xu



Yu Sun



Jiaming Liu



Guangxiao Song

Thanks & questions?