# Stochastic Iterative Hard Thresholding for Graph-Structured Sparsity Optimization

Baojian Zhou[1,2], Feng Chen[1], and Yiming Ying[2]

[1]Department of Computer Science,   [2]Department of Mathematics and Statistics, University at Albany, NY, USA
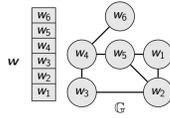
## Structured Sparse Learning

Given $\mathcal{M}(\mathbb{M}) = \{w : \text{supp}(w) \in \mathbb{M}\}$, the structured sparse learning problems can be formulated as

$$\min_{w \in \mathcal{M}(\mathbb{M})} F(w) := \frac{1}{n}\sum_{i=1}^{n} f_i(w), \text{ where}$$

▶ $F(w)$ is a convex loss such as least square, logistic loss, ...
▶ $\mathcal{M}(\mathbb{M})$ models structured sparsity such as connected subgraphs, dense subgraphs, and subgraphs isomophic to a query graph, ...
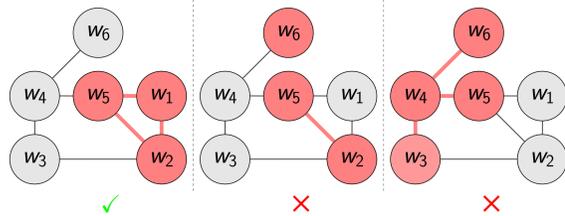
## Previous Work



Figure: Weighted Graph Model $\mathbb{M} = \{S : |S| \leq 3, S \text{ is connected }\}$ Hegde et al. (2015a).

To solve above problem under sparsity constraint, Nguyen et al. (2017) proposed Stochastic Iterative Hard Thresholding (STOIHT) choose $\xi_t$ from $[n]$ with probability $p_{\xi_t}$ and project $w^t$ onto a subspace

$$w^{t+1} = P(w^t - \eta_t \nabla f_{\xi_t}(w^t), \Gamma^t),$$

where the orthogonal projection $P(\cdot, \Gamma)$ is defined as

$$P(w, \Gamma) := \arg\min_{w' \in \mathcal{R}(\Gamma)} \|w - w'\|_2^2.$$

Why stochastic?
▶ More steady
▶ Less computation per-iteration

Two issues of STOIHT
▶ Cannot handle graph-structured constraint
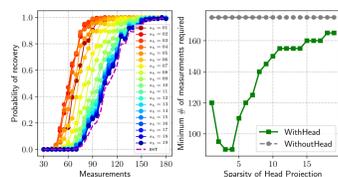▶ Ideally, $\nabla f_{\xi_t}(w^t)$ also needs to be in a subspace

## Our Algorithm

The hybird of Nguyen et al. (2017) and Hegde et al. (2016).

**Algorithm 1** GRAPHSTOIHT

1: **Input**: $\eta_t, F(\cdot), \mathbb{M}_{\mathcal{H}}, \mathbb{M}_{\mathcal{T}}$
2: **Initialize**: $w^0$ and $t = 0$
3: **for** $t = 0, 1, 2, \ldots$ **do**
4:   Choose $\xi_t$ from $[n]$ with prob. $p_{\xi_t}$
5:   $b^t = P(\nabla f_{\xi_t}(w^t), \mathbb{M}_{\mathcal{H}})$
6:   $w^{t+1} = P(w^t - \eta_t b^t, \mathbb{M}_{\mathcal{T}})$
7: **end for**
8: **Return** $w^{t+1}$

Why projection $b^t = P(\nabla f_{\xi_t}(w^t), \mathbb{M}_{\mathcal{H}})$ ?

▶ Both of them solve the same projection problem
▶ Sparsity is both in primal space and dual space
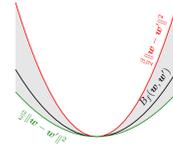▶ Remove some noisy directions at the first stage



## Convergence Analysis

Define the Bregman divergence of $f$ as

$$B_f(w, w') = f(w) - f(w') - \langle \nabla f(w'), w - w' \rangle$$

Two assumptions in $\mathcal{M}(\mathbb{M})$:
▶ $f_i(w)$: $\beta$-Restricted Strong Smoothness
  $F(w)$: $\alpha$-Restricted Strong Convexity
▶ Efficient Approximated projections:
  ● $P(\cdot, \mathbb{M}_{\mathcal{H}})$ with approximation factor $c_{\mathcal{H}}$
  ● $P(\cdot, \mathbb{M}_{\mathcal{T}})$ with approximation factor $c_{\mathcal{T}}$

**Theorem 1 (Linear Convergence)** Let $w^0$ be the start point and choose $\eta_t = \eta$, then $w^{t+1}$ of Algorithm 1 satisfies

$$\mathbb{E}_{\xi_{[t]}}\|w^{t+1} - w^*\| \leq \kappa^{t+1}\|w^0 - w^*\| + \frac{\sigma}{1-\kappa},$$

where $\eta, \tau \in (0, 2/\beta)$ and

$$\kappa = (1 + c_{\mathcal{T}})\left(\sqrt{\alpha\beta\eta^2 - 2\alpha\eta + 1} + \sqrt{1 - \alpha_0^2}\right),$$
$$\alpha_0 = c_{\mathcal{H}}\alpha\tau - \sqrt{\alpha\beta\tau^2 - 2\alpha\tau + 1}, \quad \beta_0 = (1 + c_{\mathcal{H}})\tau,$$
$$\sigma = \left(\frac{\beta_0}{\alpha_0} + \frac{\alpha_0\beta_0}{\sqrt{1-\alpha_0^2}}\right)\mathbb{E}_{\xi_t}\|\nabla_I f_{\xi_t}(w^*)\| + \eta\mathbb{E}_{\xi_t}\|\nabla_I f_{\xi_t}(w^*)\|.$$

## Graph Sparse Linear Regression

Given a design matrix $X \in \mathbb{R}^{m \times p}$ and corresponding observed noisy vector $y \in \mathbb{R}^m$ that are linked via the linear relationship

$$y = Xw^* + \epsilon,$$

where $\epsilon \sim \mathcal{N}(0, \sigma^2 I_m)$. To estimate $w^*$, consider the least square loss and formulate it as

$$\arg\min_{\text{supp}(w) \in \mathcal{M}(\mathbb{M})} F(w) := \frac{1}{n}\sum_{i=1}^{n} \frac{n}{2m}\|X_{B_i}w - y_{B_i}\|^2,$$

where $m$ observations have been partitioned into $n$ blocks, $B_1, B_2, \ldots, B_n$. Let $\alpha = 1 - \delta, \beta = 1 + \delta$.

| Algorithm | $\kappa$ |
|---|---|
| GRAPHIHT | $(1 + c_{\mathcal{T}})\left(\sqrt{\delta} + 2\sqrt{1-\delta}\right)\sqrt{\delta}$ |
| GRAPHSTOIHT | $(1 + c_{\mathcal{T}})\left(\sqrt{\frac{2}{1+\delta}} + \frac{2\sqrt{2(1-\delta)}}{1+\delta}\right)\sqrt{\delta}$ |

$\kappa$ of GRAPHIHT is controlled by $\mathcal{O}(\sqrt{\delta} \cdot 2(1 + c_{\mathcal{T}}))$ while for GRAPHSTOIHT, $\kappa$ is controlled by $\mathcal{O}(\sqrt{\delta} \cdot 3\sqrt{2}(1 + c_{\mathcal{T}}))$. To obtain $\kappa < 1$, $\delta \leq 0.0527$ for GRAPHIHT while $\delta \leq 0.0142$ for GRAPHSTOIHT. The gap between the two $\kappa$ is mainly due to the randomness introduced in our algorithm.

## Graph Sparse Logistic Regression

Given a dataset $\{x_i, y_i\}_{i=1}^m$, the graph logistic regression is formulated as the following problem

$$\arg\min_{\text{supp}(w) \in \mathcal{M}(\mathbb{M})} F(w) := \frac{1}{n}\sum_{i=1}^{n}\frac{n}{m}\sum_{j=1}^{m/n} h(w, i_j) + \frac{\lambda}{2}\|w\|^2,$$

where $h(w, i_j) = \log(1 + \exp(-y_{i_j} \cdot \langle x_{i_j}, w \rangle))$. Problem above has an important application on gene pathway analysis. If each sample $a_i$ is normalized, then $F(x)$ satisfies $\lambda$-RSC and each $f_i(x)$ satisfies $(\alpha + (1 + \nu)\theta_{max})$-RSS. The condition of $\kappa < 1$ is
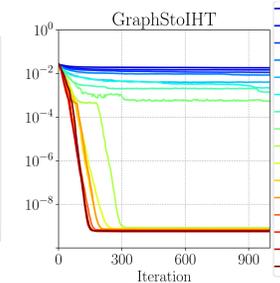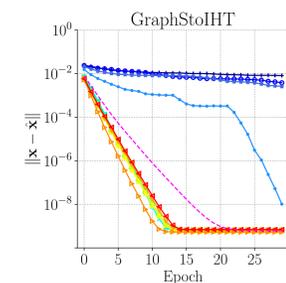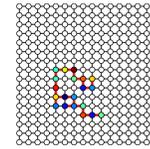
$$\frac{\lambda}{\lambda + n(1+\nu)\theta_{max}/4m} \geq \frac{243}{250},$$

with probability $1 - p\exp(-\theta_{max}\nu/4)$, where $\theta_{max} = \lambda_{max}(\sum_{j=1}^{m/n}\mathbb{E}[a_{i_j}a_{i_j}^T])$ and $\nu \geq 1$.

## Experiments

**Simulation Dataset:**
▶ Each entry $\sqrt{m}X_{ij} \sim \mathcal{N}(0, 1)$
▶ Supp($w^*$) is generated by random walk
▶ Entries of $w^*$ from $\mathcal{N}(0, 1)$
▶ Weighted Graph Model



▶ **Upper Right:** Probability of recovery on synthetic dataset. The probability of recovery is a function of number of observations $m$.
▶ **Lower Left:** The left part illustrates the *estimation error* as a function of epochs for different choices of $b$. When $b = 180$, it degenerates to GRAPHIHT (the dashed line). The right part shows the *estimation error* as a function of iterations for different choices of $\eta$.
▶ **Lower Right:** Robustness to noise $\epsilon$. The number of observations required is a function of different block sizes.





**Real image dataset:**
▶ IHT (Blumensath and Davies, 2009)
▶ STOIHT (Nguyen et al., 2017)
▶ NIHT (Blumensath and Davies, 2010)
▶ COSAMP (Needell and Tropp, 2009)
▶ GRAPHIHT (Hegde et al., 2016) + WGM
▶ GRAPHCOSAMP (Hegde et al., 2015b)

**Experimental settings:**
▶ Resized real images (Hegde et al., 2015b)
▶ $\eta$ of IHT-based in $\{0.2, 0.4, 0.6, 0.8\}$
▶ $b$ of STOIHT-based in $\{m/5, m/10\}$
▶ Tune $b$ and $\eta$ on 100 observations.
▶ $A$ used here is Gaussian matrix

**Two experimental conclusions:**
▶ SGD-based methods are more stable
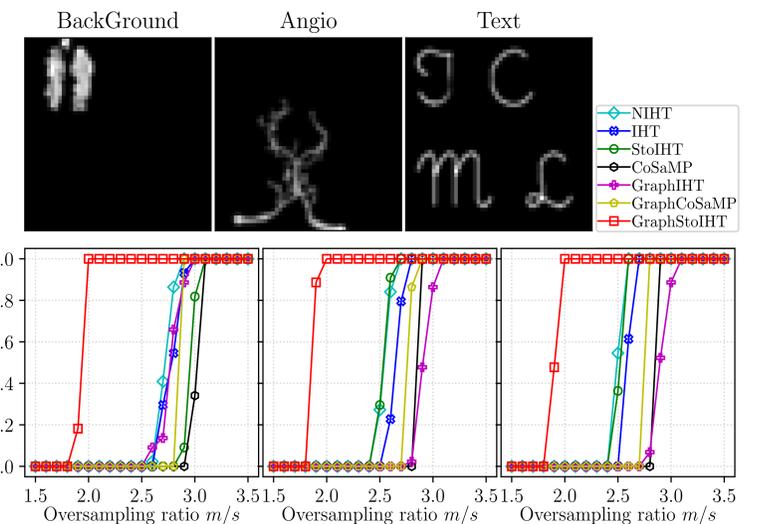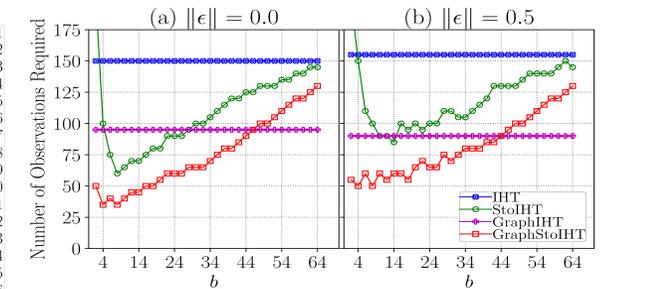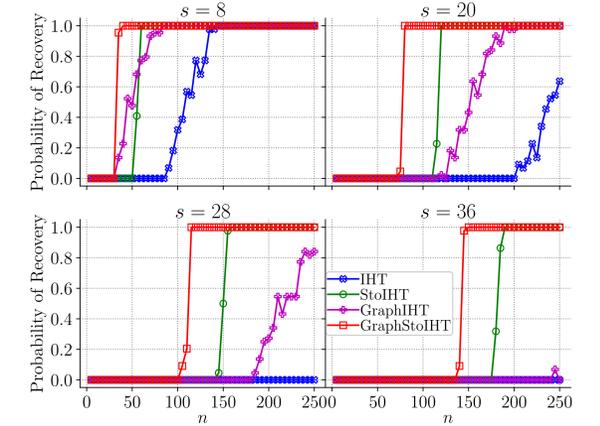▶ Capture the graph-structured sparsity





**Breast Cancer Dataset:**
▶ 295 samples with 78 positives (metastatic) and 217 negatives (non-metastatic) provided in Van De Vijver et al. (2002).
▶ PPI network with 637 pathways is provided in Jacob et al. (2009).

**Four $\ell^1/\ell^2$ mixed norm-based algorithms:**
▶ $\ell^1$-PATHWAY uses pathways as groups
▶ $\ell^1/\ell^2$-PATHWAY uses pathways as groups
▶ $\ell^1$-EDGE uses use edges as groups
▶ $\ell^1/\ell^2$-EDGE uses use edges as groups

| Algorithm | Cancer related genes | $\|w^t\|_0$ | AUC |
|---|---|---|---|
| GRAPHSTOIHT | BRCA2, CCND2, CDKN1A, ATM, AR, TOP2A | 051.7 | **0.715** |
| GRAPHIHT | ATM, CDKN1A, BRCA2, AR, TOP2A | 055.2 | 0.714 |
| $\ell^1$-PATH | BRCA1, CDKN1A, ATM, DSC2 | 061.2 | 0.675 |
| STOIHT | MKI67, NAT1, AR, TOP2A | 059.6 | 0.708 |
| $\ell^1/\ell^2$-PATH | CCND3, ATM, CDH3 | 051.4 | 0.705 |
| $\ell^1$-EDGE | CCND3, AR, CDH3 | 039.9 | 0.698 |
| $\ell^1/\ell^2$-PATH | BRCA1, CDKN1A | 147.6 | 0.705 |
| IHT | NAT1, TOP2A | 067.9 | 0.707 |

## Conclusion and Future Work

▶ We proposed GRAPHSTOIHT.
▶ It enjoys a linear convergence property.
▶ Two real-world applications.

In future, it would be interesting to see if one can apply the variance reduction techniques such as SAGA (Defazio et al., 2014) and SVRG (Johnson and Zhang, 2013) to GRAPHSTOIHT.

## Code & Datasets

▶ Code & Datasets can be found at GitHub:
  https://github.com/baojianzhou/graph-sto-iht
▶ Email: bzhou6@albany.edu
▶ **Baojian Zhou is open to postdoc positions.**