

# Dual Averaging Method for Online Graph-Structured Sparsity

Baojian Zhou<sup>1,2</sup>, Feng Chen<sup>1</sup>, and Yiming Ying<sup>2</sup>

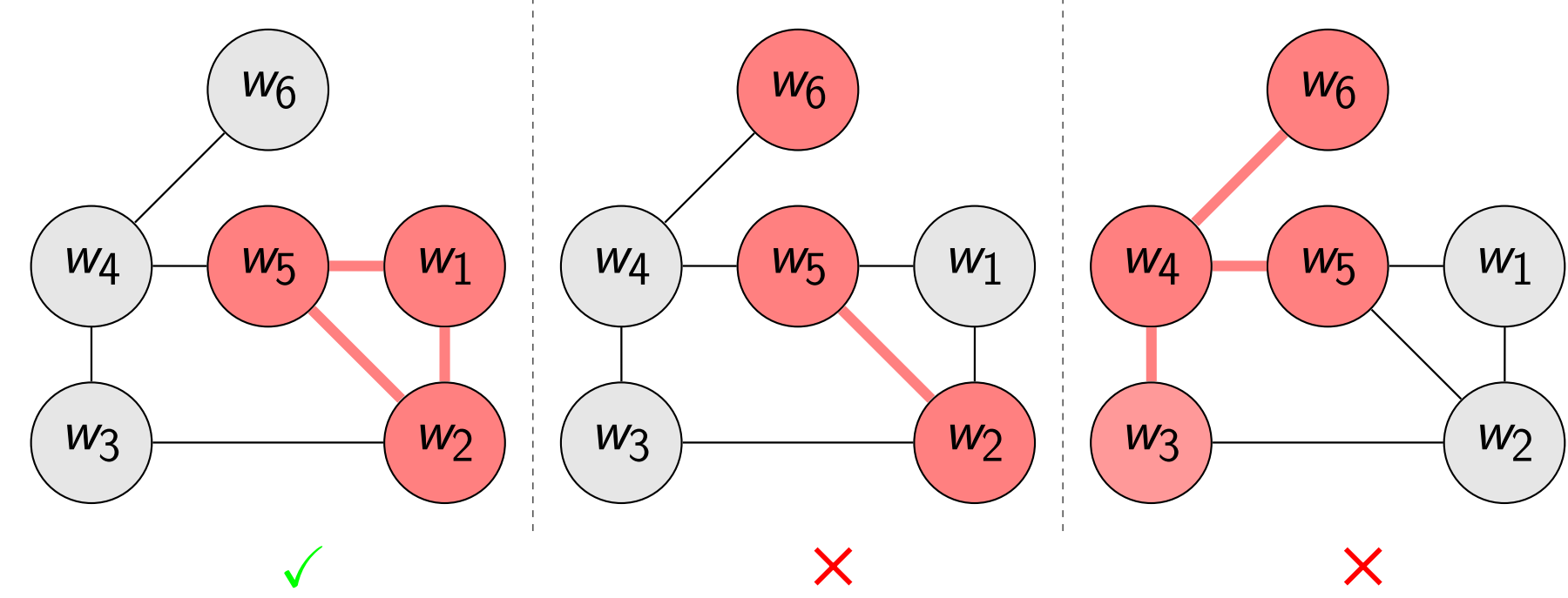
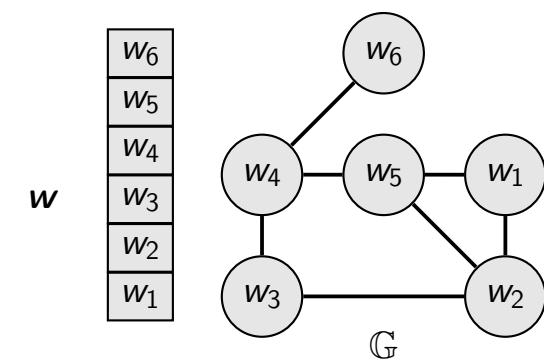
<sup>1</sup>Department of Computer Science, <sup>2</sup>Department of Mathematics and Statistics, University at Albany, NY, USA

## Online Graph-structured Learning

We study an online graph-structured learning problem, which is to minimize the regret as defined in the following

$$R(T, \mathcal{M}(\mathbb{M})) := \sum_{t=1}^T f_t(\mathbf{w}_t, \{\mathbf{x}_t, \mathbf{y}_t\}) - \min_{\mathbf{w} \in \mathcal{M}(\mathbb{M})} \sum_{t=1}^T f_t(\mathbf{w}, \{\mathbf{x}_t, \mathbf{y}_t\}),$$

- ▶  $f_t(\mathbf{w}, \{\mathbf{x}_t, \mathbf{y}_t\})$  is a convex loss
- ▶  $\mathcal{M}(\mathbb{M})$  models structured sparsity such as connected subgraphs, dense subgraphs, and subgraphs isomorphic to a query graph, ...



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

## Main Idea

An intuitive way to do this is to use online projected gradient descent Zinkevich (2003) where the algorithm needs to solve the following projection at iteration  $t$ :

$$\mathbf{w}_{t+1} = P(\mathbf{w}_t - \eta_t \nabla f_t(\mathbf{w}_t), \mathcal{M}(\mathbb{M})), \quad (1)$$

where  $\eta_t$  is the learning rate and  $P$  is the projection operator onto  $\mathcal{M}(\mathbb{M})$ , i.e.,  $P(\cdot, \mathcal{M}(\mathbb{M})) : \mathbb{R}^p \rightarrow \mathbb{R}^p$  is defined as

$$P(\mathbf{w}, \mathcal{M}(\mathbb{M})) = \arg \min_{\mathbf{x} \in \mathcal{M}(\mathbb{M})} \|\mathbf{w} - \mathbf{x}\|_2^2. \quad (2)$$

However, there are two essential drawbacks of online PGD

- ▶ The projection in (1) only uses single gradient  $\nabla f_t(\mathbf{w}_t)$  which is too noisy (large variance) to capture the graph-structured information at each iteration;
- ▶ The training samples coming later are less important than these coming earlier due to the decay of learning rate  $\eta_t$ .

Inspired by dual averaging-based methods, at each iteration, our method updates  $\mathbf{w}_t$  by using the following minimization step:

$$\mathbf{w}_{t+1} = \arg \min_{\mathbf{w} \in \mathcal{M}(\mathbb{M})} \left\{ \left\langle \frac{1}{t+1} \sum_{i=0}^t \mathbf{g}_i, \mathbf{w} \right\rangle + \frac{\beta_t}{2t} \|\mathbf{w}\|_2^2 \right\}, \quad (3)$$

where  $\beta_t$  is to control the learning rate implicitly and  $\mathbf{g}_i$  is a subgradient of  $\partial f_i(\mathbf{w}, \{\mathbf{x}_i, \mathbf{y}_i\}) = \{ \mathbf{g} : f_i(\mathbf{z}, \{\mathbf{x}_i, \mathbf{y}_i\}) \geq f_i(\mathbf{w}, \{\mathbf{x}_i, \mathbf{y}_i\}) + \langle \mathbf{g}, \mathbf{z} - \mathbf{w} \rangle, \forall \mathbf{z} \in \mathcal{M}(\mathbb{R}) \}$ . The minimization step (3) has the following equivalent projection problems, specified in the following Theorem.

**Theorem:** Assume  $\beta_t = \gamma \sqrt{t}$ , where  $\gamma > 0$  and denote  $\bar{\mathbf{s}}_{t+1} = \frac{1}{t+1} \sum_{i=0}^t \mathbf{g}_i$ . The minimization step of (3) can be expressed as the following two equivalent optimization problems:

$$\max_{S \in \mathbb{M}} \left\| P\left(-\frac{\sqrt{t} \bar{\mathbf{s}}_{t+1}}{\gamma}, S\right) \right\|_2^2 \quad (4)$$

$$\min_{S \in \mathbb{M}} \left\| -\frac{\sqrt{t} \bar{\mathbf{s}}_{t+1}}{\gamma} - P\left(-\frac{\sqrt{t} \bar{\mathbf{s}}_{t+1}}{\gamma}, S\right) \right\|_2^2, \quad (5)$$

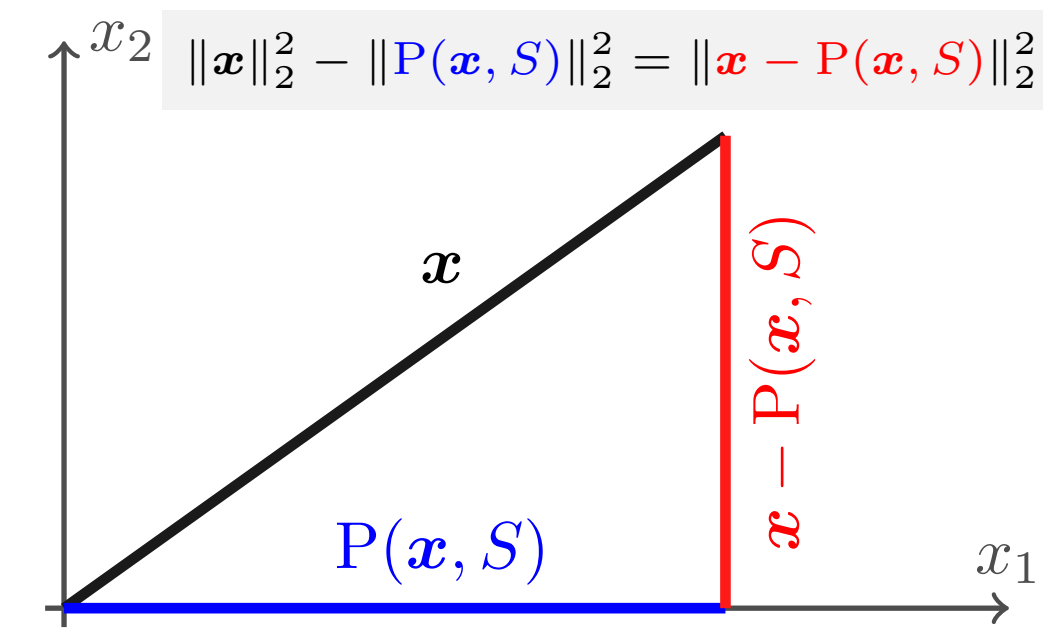
where  $P(s, S)$  is the projection operator that projects  $s$  onto the subspace spanned by  $S$ .

The original minimization problem can be equivalently expressed as

$$\begin{aligned} \mathbf{w}_{t+1} &= \arg \min_{\mathbf{w} \in \mathcal{M}(\mathbb{M})} \left\{ \langle \bar{\mathbf{s}}_{t+1}, \mathbf{w} \rangle + \frac{\gamma}{2\sqrt{t}} \|\mathbf{w}\|_2^2 \right\} \\ &= \arg \min_{\mathbf{w} \in \mathcal{M}(\mathbb{M})} \left\| \mathbf{w} - \left( -\frac{\sqrt{t}}{\gamma} \bar{\mathbf{s}}_{t+1} \right) \right\|_2^2, \end{aligned}$$

Each step is essentially a projection !

## Theorem Insight



Define  $\mathbf{x} := -\sqrt{t} \bar{\mathbf{s}}_{t+1}$  and adding minimization to both sides

$$\min_{S \in \mathbb{M}} \left\{ \|\mathbf{x}\|_2^2 - \|P(\mathbf{x}, S)\|_2^2 \right\} = \min_{S \in \mathbb{M}} \|\mathbf{x} - P(\mathbf{x}, S)\|_2^2.$$

By moving the minimization into the negative term, we obtain

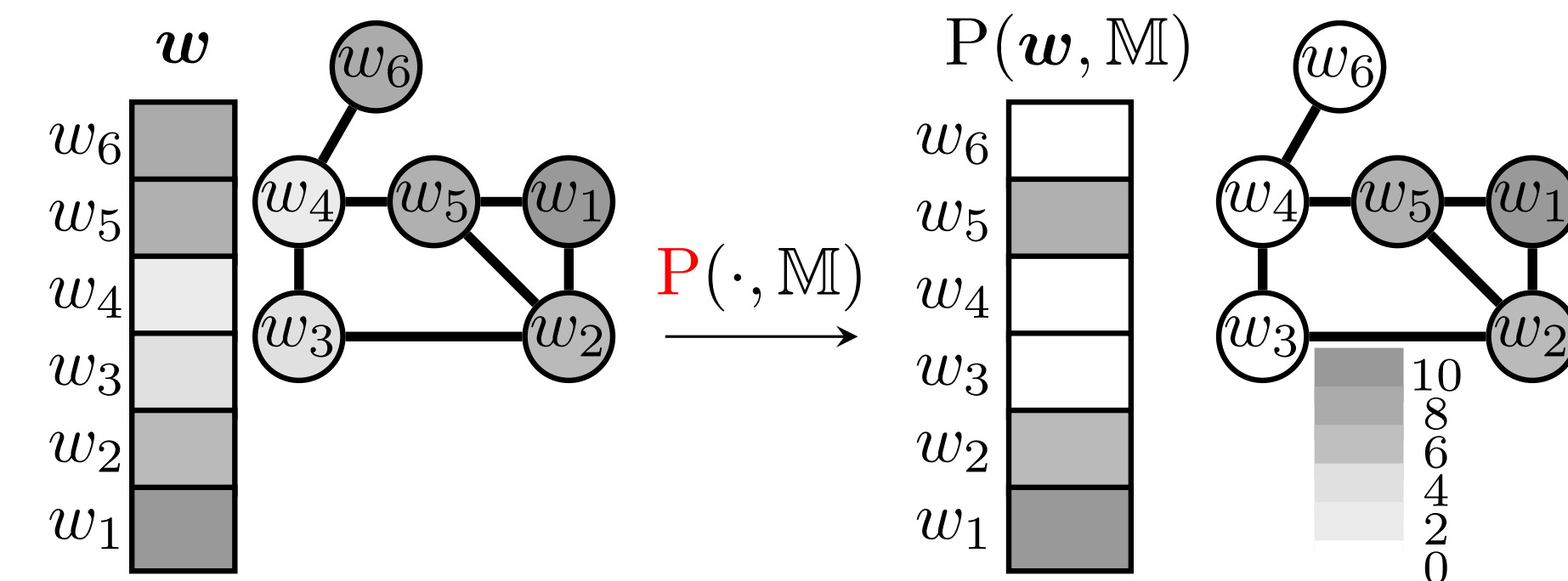
$$\|\mathbf{x}\|_2^2 + \max_{S \in \mathbb{M}} \|P(\mathbf{x}, S)\|_2^2 = \min_{S \in \mathbb{M}} \|\mathbf{x} - P(\mathbf{x}, S)\|_2^2.$$

## Proposed Algorithms

**Algorithm 1** GRAPHDA: Online Graph Dual Averaging Algorithm

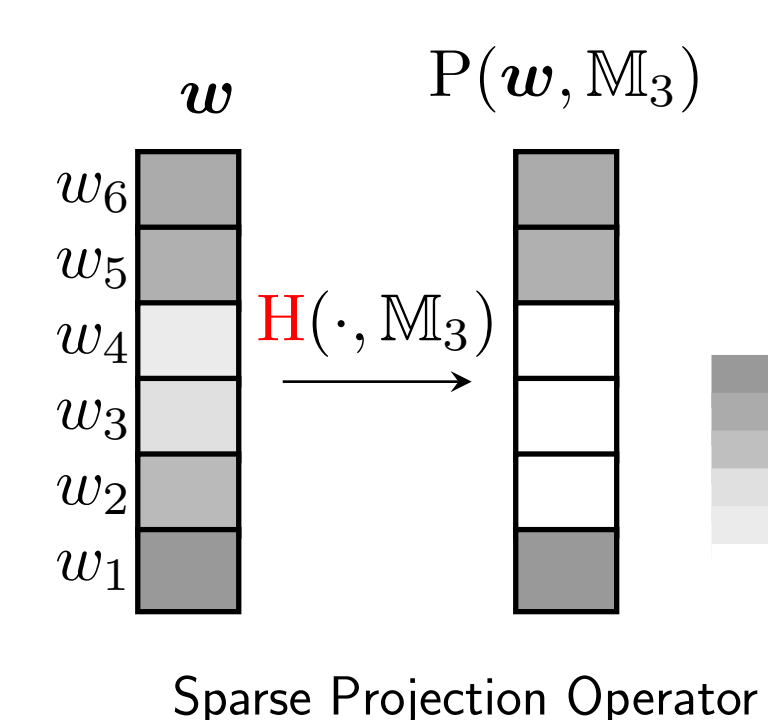
- 1: **Input:**  $\gamma, \mathbb{M}$
- 2:  $\bar{\mathbf{s}}_0 = \mathbf{0}, \mathbf{w}_0 = \mathbf{0}$
- 3: **for**  $t = 0, 1, 2, \dots$  **do**
- 4: receive  $\{\mathbf{x}_t, \mathbf{y}_t\}$  and compute  $\mathbf{g}_t = \nabla f_t(\mathbf{w}_t, \{\mathbf{x}_t, \mathbf{y}_t\})$
- 5:  $\bar{\mathbf{s}}_{t+1} = \bar{\mathbf{s}}_t + \mathbf{g}_t$
- 6:  $\mathbf{b}_{t+1} = P(\bar{\mathbf{s}}_{t+1}, \mathbb{M})$
- 7:  $\mathbf{w}_{t+1} = P(-\frac{\sqrt{t}}{\gamma} \mathbf{b}_{t+1}, \mathbb{M})$
- 8: **end for**

Let  $\mathbb{M} = \{S : |S| \leq 3, S \text{ is connected}\}$ . Finding a connected subgraph up to 3 nodes.



Graph Projection Operator Hegde et al. (2015b)

What if the graph information is not available ?



**Algorithm 2** DAIHT

- 1: **Input:**  $\gamma, \mathbb{M}$
- 2:  $\bar{\mathbf{s}}_0 = \mathbf{0}, \mathbf{w}_0 = \mathbf{0}$
- 3: **for**  $t = 0, 1, 2, \dots$  **do**
- 4: receive  $\{\mathbf{x}_t, \mathbf{y}_t\}$  and compute  $\mathbf{g}_t = \nabla f_t(\mathbf{w}_t, \{\mathbf{x}_t, \mathbf{y}_t\})$
- 5:  $\bar{\mathbf{s}}_{t+1} = \bar{\mathbf{s}}_t + \mathbf{g}_t$
- 6:  $\mathbf{b}_{t+1} = H(\bar{\mathbf{s}}_{t+1}, \mathbb{M})$
- 7:  $\mathbf{w}_{t+1} = H(-\frac{\sqrt{t}}{\gamma} \mathbf{b}_{t+1}, \mathbb{M})$
- 8: **end for**

## Time Complexity and Regret

The time complexity of GRAPHDA mainly depends on two projections. If we use the weighted-graph model, the per-iteration time cost could be

- ▶ non-sparse graph:  $\mathcal{O}(p + |\mathbb{E}| \log^3(p))$  **Edge-dependent**
- ▶ sparse graph:  $\mathcal{O}(p + p \log^3(p))$  **Nearly-linear !**

If  $\mathcal{M}(\mathbb{M})$  is a **convex set**, then

- ▶ The regret can be bounded as:  $R(T, \mathcal{M}(\mathbb{M})) = C \cdot \mathcal{O}(\sqrt{T})$ , where  $C$  is a constant.

- ▶ If we assume further that the loss is strongly convex, then  $\|\mathbf{w}_T - \mathbf{w}^*\|_2 = \mathcal{O}(\frac{\ln T}{\sqrt{T}})$ .

However,  $\mathcal{M}(\mathbb{M})$  is not convex in our case. We leave the regret bound analysis of this case as a future work.

## Experiments

Method	Proposed in	Dataset	$ \mathbb{V} $	$ \mathbb{E} $
ADAM	Kingma and Ba (2014)	Benchmark	1,089	2,112
$\ell_1$ -RDA	Xiao (2010)	MNIST	786	1,516
DA-GL	Yang et al. (2010)	KEGG	5,372	78,545
DA-SGL	Yang et al. (2010)			
ADAGRAD	Duchi et al. (2011)			
STOHT	Nguyen et al. (2017)			
GRAPHSTOHT	Zhou et al. (2019)			
DA-IHT	<b>This paper</b>			
GRAPHDA	<b>This paper</b>			

Legend:  
■ non-sparse  
■ sparse: convex-based  
■ sparse: nonconvex-based

Aim to answer the following two questions:

- ▶ Can GRAPHDA achieve better classification performance?
- ▶ Can GRAPHDA learn stronger interpretative model through capturing more meaningful graph-structured features?

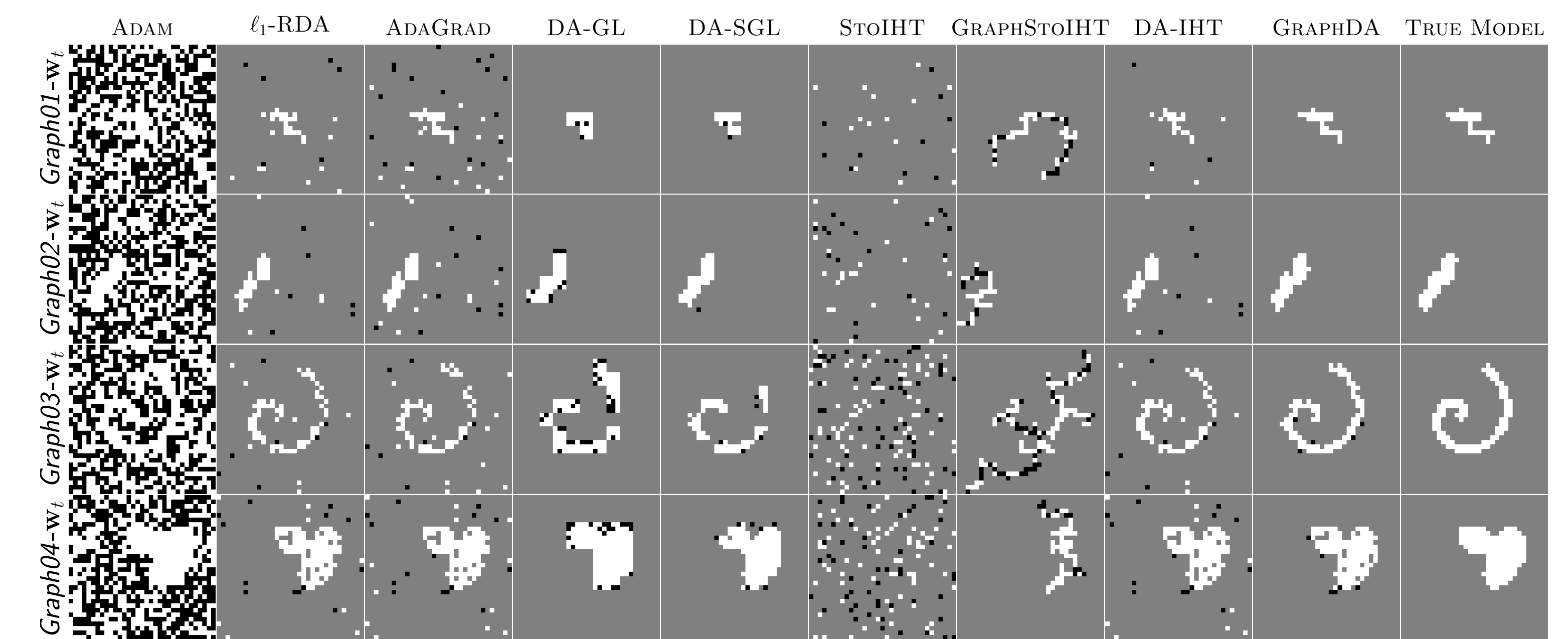
## Evaluation Metric

- ▶ Area Under the ROC Curve(AUC) score
- ▶ Classification Accuracy(Acc)
- ▶ Nonzero Ratio

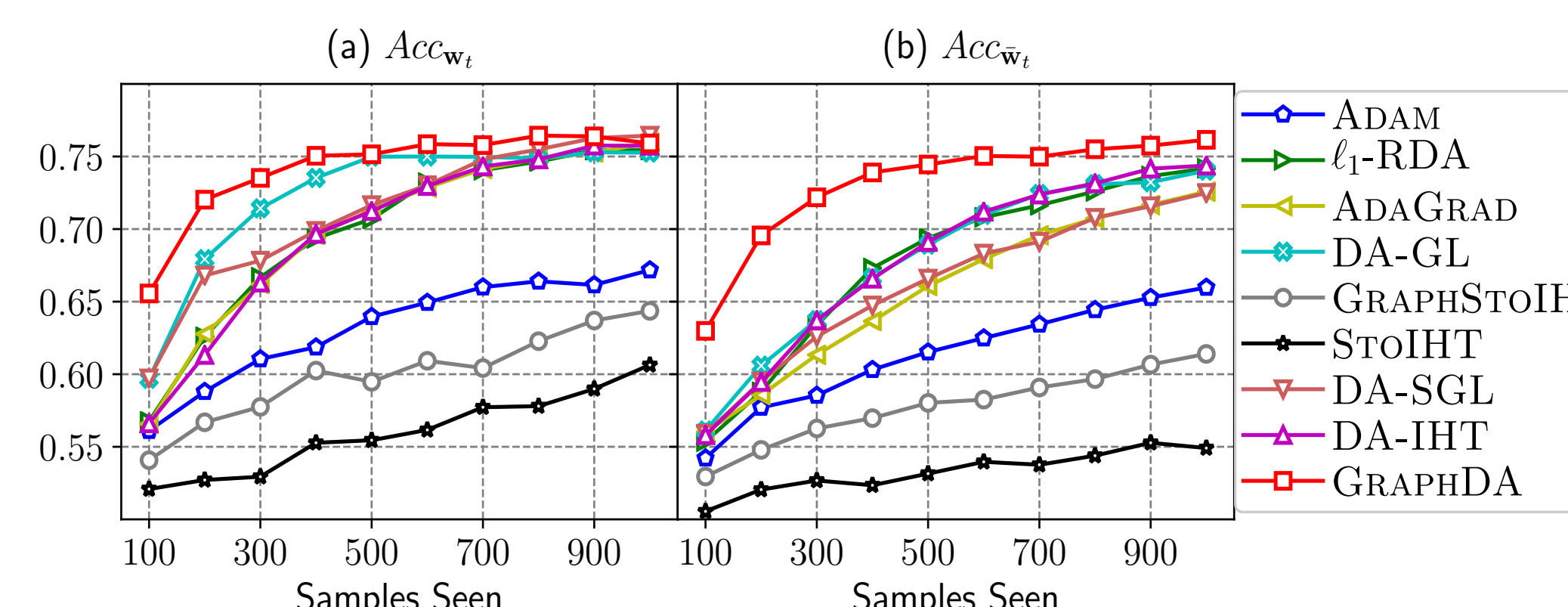
$$NR_{\mathbf{w}} = \frac{|\text{supp}(\mathbf{w})|}{p}$$

Method	$AUC_{w, \pi}$	$Acc_{w, \pi}$	$NR_{w, \pi}$
ADAM	(0.618, 0.603)	(0.619, 0.603)	(100.0%, 100.0%)
$\ell_1$ -RDA	(0.693, 0.672)	(0.694, 0.673)	(11.58%, 83.60%)
ADAGRAD	(0.696, 0.636)	(0.696, 0.637)	(11.33%, 100.0%)
DA-GL	(0.735, 0.666)	(0.735, 0.667)	(15.99%, 100.0%)
DA-SGL	(0.699, 0.647)	(0.699, 0.647)	(25.54%, 100.0%)
STOHT	(0.552, 0.523)	(0.553, 0.523)	(7.79%, 40.62%)
GRAPHSTOHT	(0.603, 0.570)	(0.602, 0.570)	(7.84%, 22.06%)
DA-IHT	(0.697, 0.666)	(0.697, 0.666)	(4.35%, 39.50%)
GRAPHDA	<b>(0.749, 0.739)</b>	<b>(0.749, 0.739)</b>	<b>(2.56%, 32.12%)</b>

- ▶ GRAPHDA has better classification performance.
- ▶ GRAPHDA has stronger model interpretability.

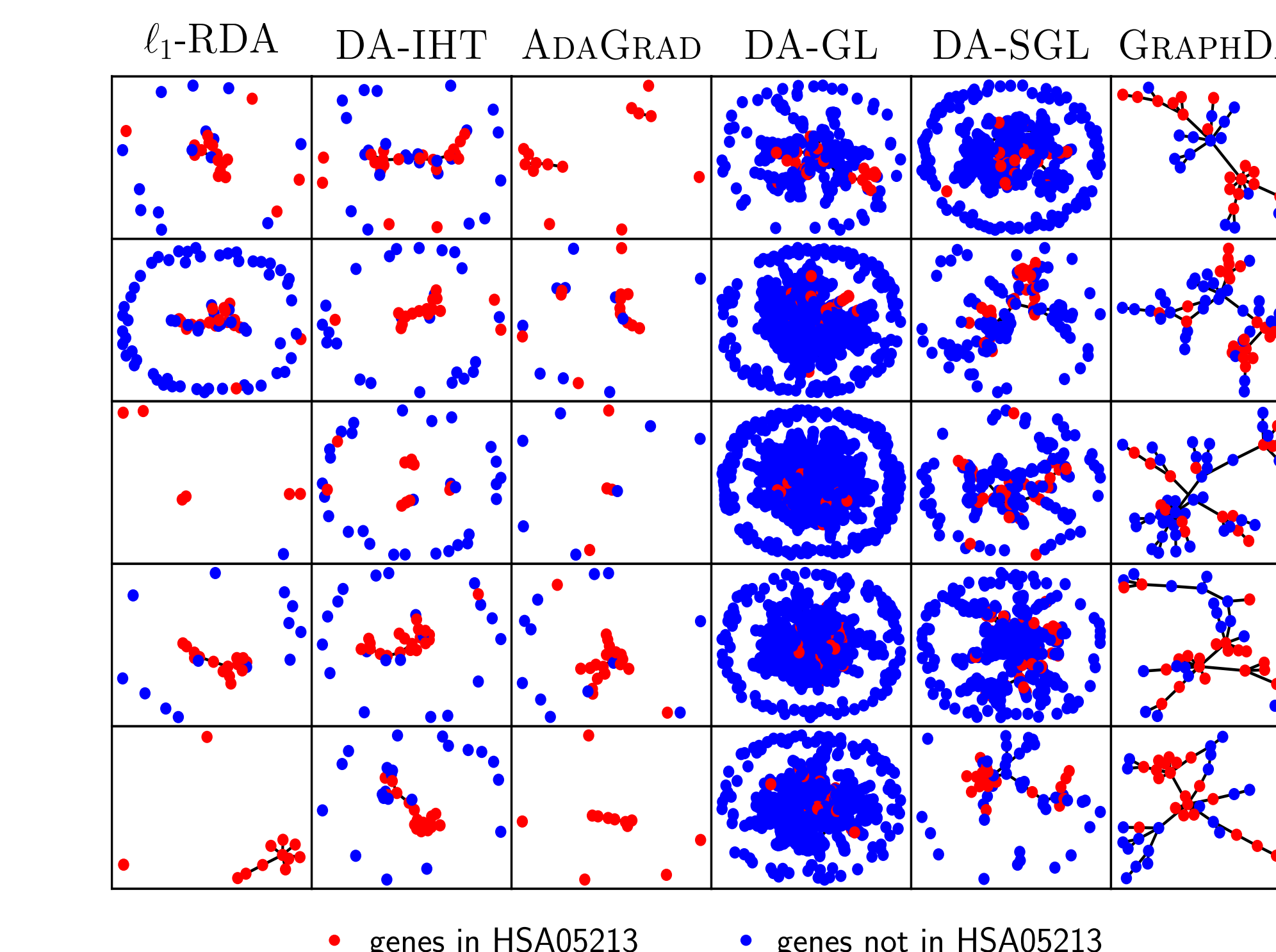


The learned models  $\mathbf{w}_t$  of four benchmark graphs. For each pixel  $i$ , black stands for  $(w_t)_i < 0$ , gray  $(w_t)_i = 0$ , and white  $(w_t)_i > 0$ .



The classification accuracy as a function of number of training samples seen

- ▶ Online PGD-based: STOHT and GRAPHSTOHT **do not work!**
- ▶ Online DA-based:  $\ell_1$ -RDA, DA-GL, DA-SGL, DA-IHT work well.
- ▶ GRAPHDA outperforms other DA-based with the help of graph priors.



• genes in HSA05213 • genes not in HSA05213

## Conclusion and Future Work

### Conclusion

- ▶ We propose a dual averaging-based method, GRAPHDA, for online graph-structured sparsity constraint problems.
- ▶ We prove that the minimization problem in the dual averaging step can be formulated as two equivalent optimization problems.
- ▶ GRAPHDA achieves better classification performance and stronger interpretability.

### Future work

- ▶ Does GRAPHDA have non-regret bound under some proper assumption ?
- ▶ What if true structure of features are time evolving ?

## Code & Datasets

- ▶ Code & Datasets can be found at GitHub: <https://github.com/baojianzhou/graph-da>
- ▶ Email: [bzhou6@albany.edu](mailto:bzhou6@albany.edu)
- ▶ **Baojian Zhou is open to postdoc positions.**