

# Σ-Optimality for Active Learning on Gaussian Random Fields

Yifei Ma, Roman Garnett, Jeff Schneider.

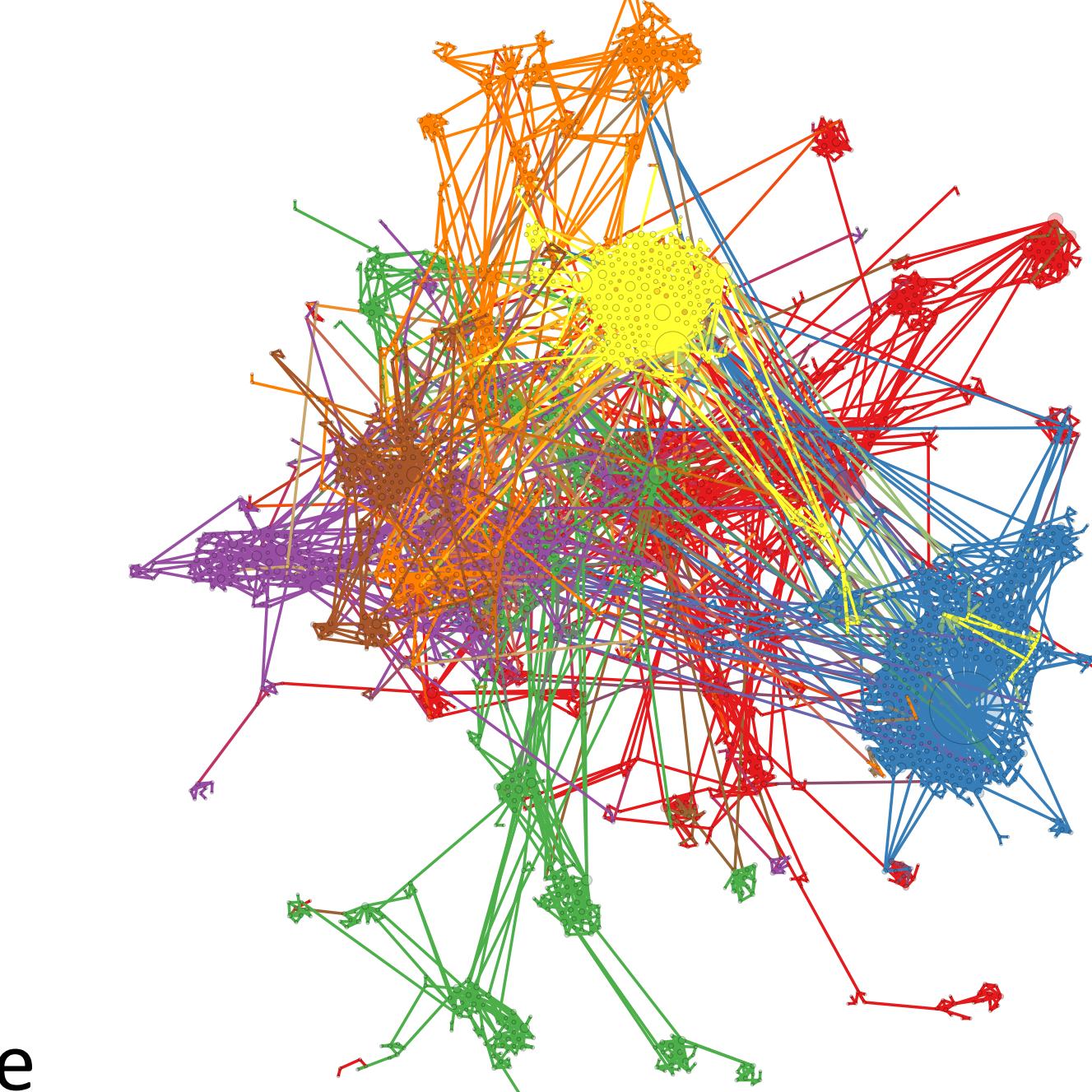
## Motivation

We consider the problem of designing a good **active learning strategy** that, under labeling budget constraints, selects which instances to query for labels that are most helpful for **classification/surveying on a graph-represented database**.

### □ Examples of graphs

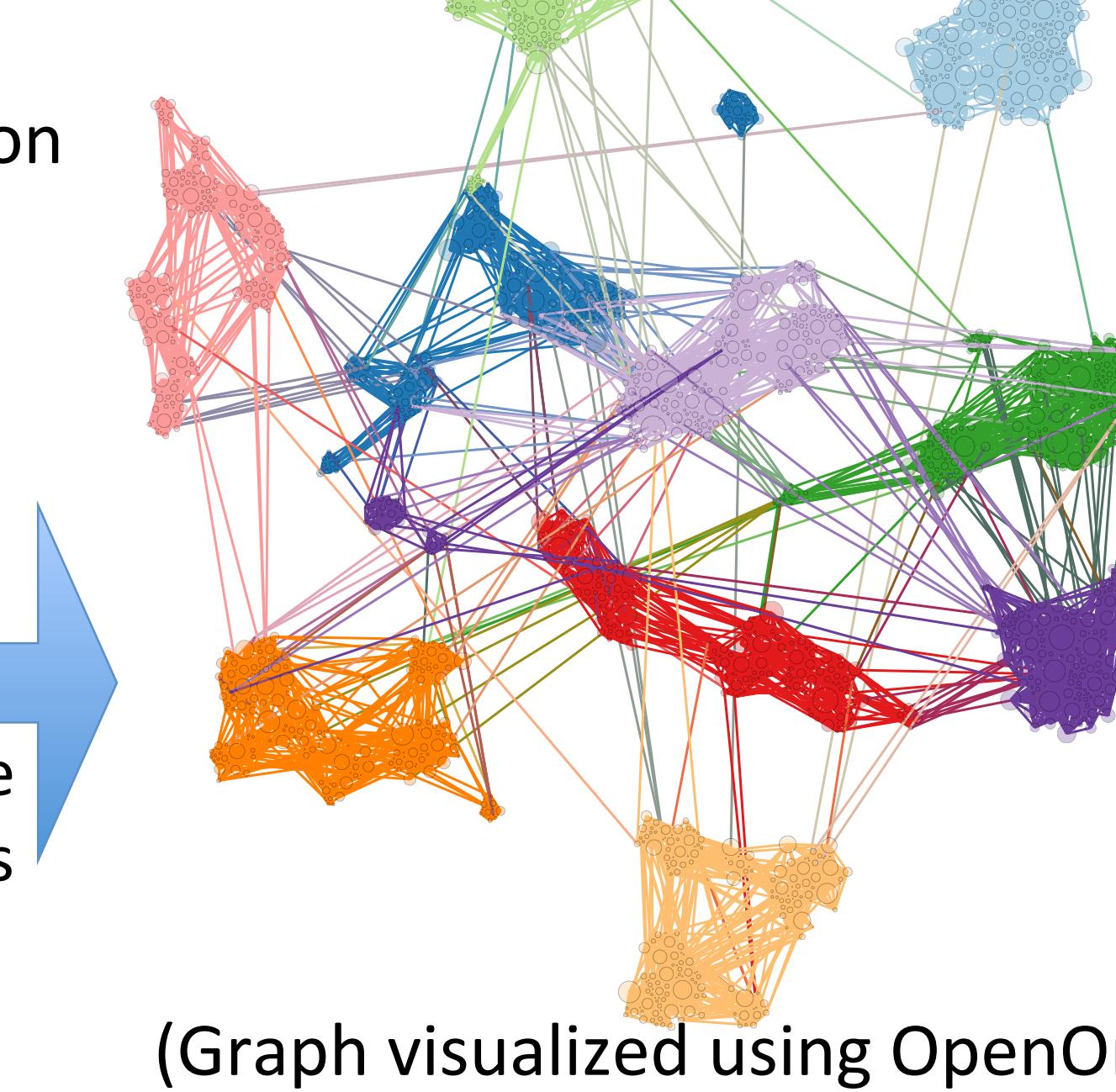
#### ◊ Citation graph

**Node:** a paper.  
**Node label (y):** topic of the paper.  
**Edge ( $A_{ij}$ ):** presence of a citation.



#### ◊ K-nn graph on feature space

**Node:** a hand-written digit image.  
**Node label (y):** 0-9, actual digit.  
**Edge ( $A_{ij}$ ):** 4-nearest-neighbors on Euclidean distance.



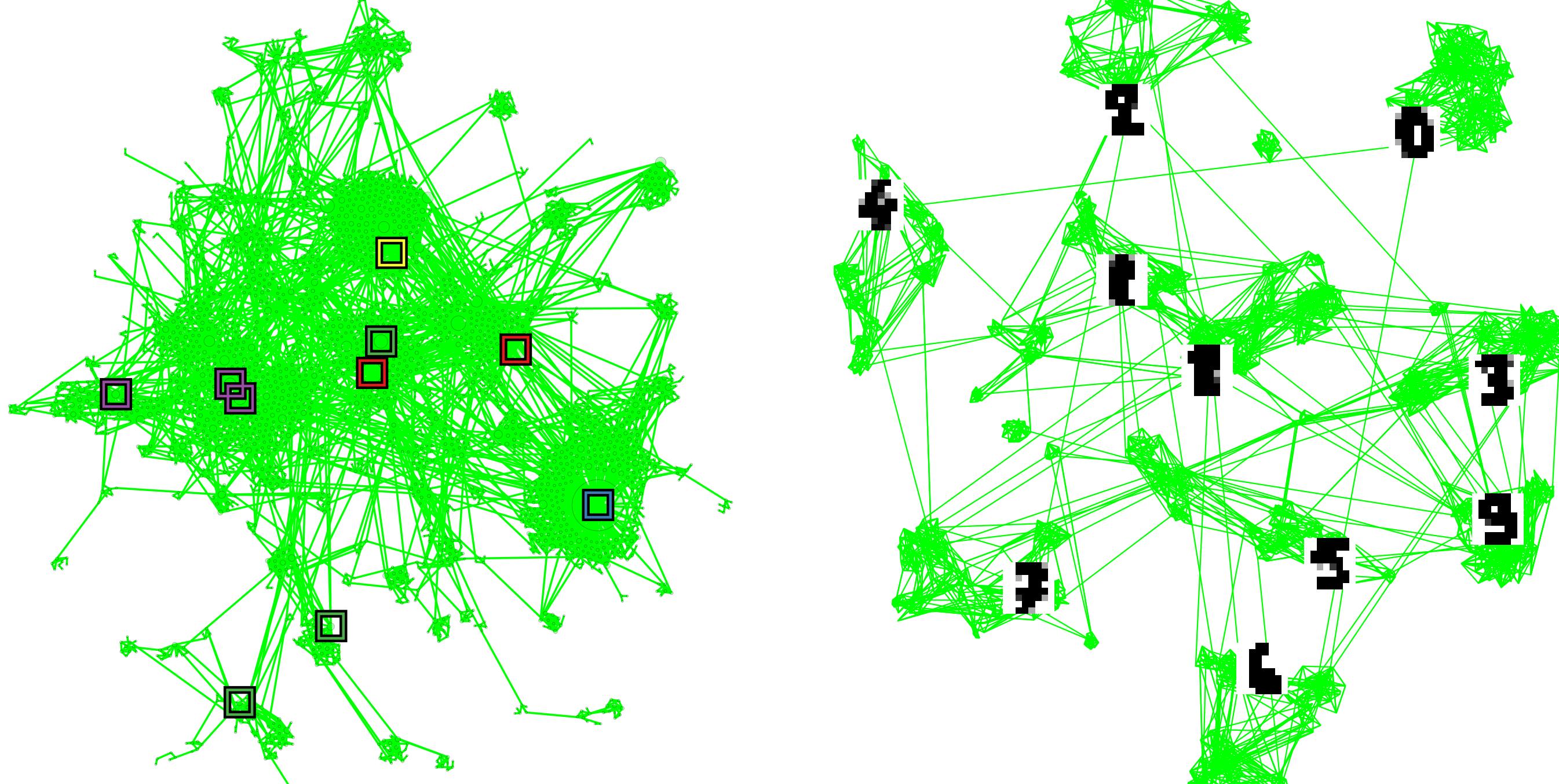
### □ Assume labels are connected to graph structure

#### ◊ Specifically, assume Gaussian random field model (GRF)

Bayesian generative model: where  $L$  is graph Laplacian:

$$P(y) \propto \exp \left\{ -\frac{1}{2} \sum_{i,j} A_{ij}(y_i - y_j)^2 \right\} = \exp \left\{ -\frac{1}{2} y^T L y \right\} \quad L = \text{diag}(\sum_j A_{ij}) - A$$

### □ The goal: query decisions (start from no label)



## Approach

For GRFs, the distribution of labels on unlabeled nodes is the **conditional Gaussian given known labels**. We minimize its **Bayesian risk**, which is also the **predictive variance**. Σ-optimality can be viewed as a variant of V-optimality.

### □ Distribution of labels on unlabeled nodes

◊ Let  $\ell$ : labeled,  $u$ : unlabeled.  $(u, \ell)$ : complementary. Split  $y, L$ .

$$P(y_u | y_\ell) \propto \mathcal{N}(y_u; \hat{y}_u, L_u^{-1}), \quad \hat{y}_u = -L_u^{-1} L_{u\ell} y_\ell$$

### □ Bayesian risk minimization

◊ Greedy L2 risk minimization: **V-optimality**

$$R_V(\ell) = \mathbb{E} \left[ \mathbb{E} \left[ \sum_{u_i \in u} (y_{u_i} - \hat{y}_{u_i})^2 \mid y_\ell \right] \right] = \text{tr}(L_u^{-1})$$

◊ Greedy surveying risk minimization: **Σ-optimality**

$$R_\Sigma(\ell) = \mathbb{E} \left[ \mathbb{E} \left[ \left( \frac{y_u \cdot \mathbf{1}}{n} - \frac{\hat{y}_u \cdot \mathbf{1}}{n} \right)^2 \mid y_\ell \right] \right] = \frac{1}{n^2} \mathbf{1}^T L_u^{-1} \mathbf{1}$$

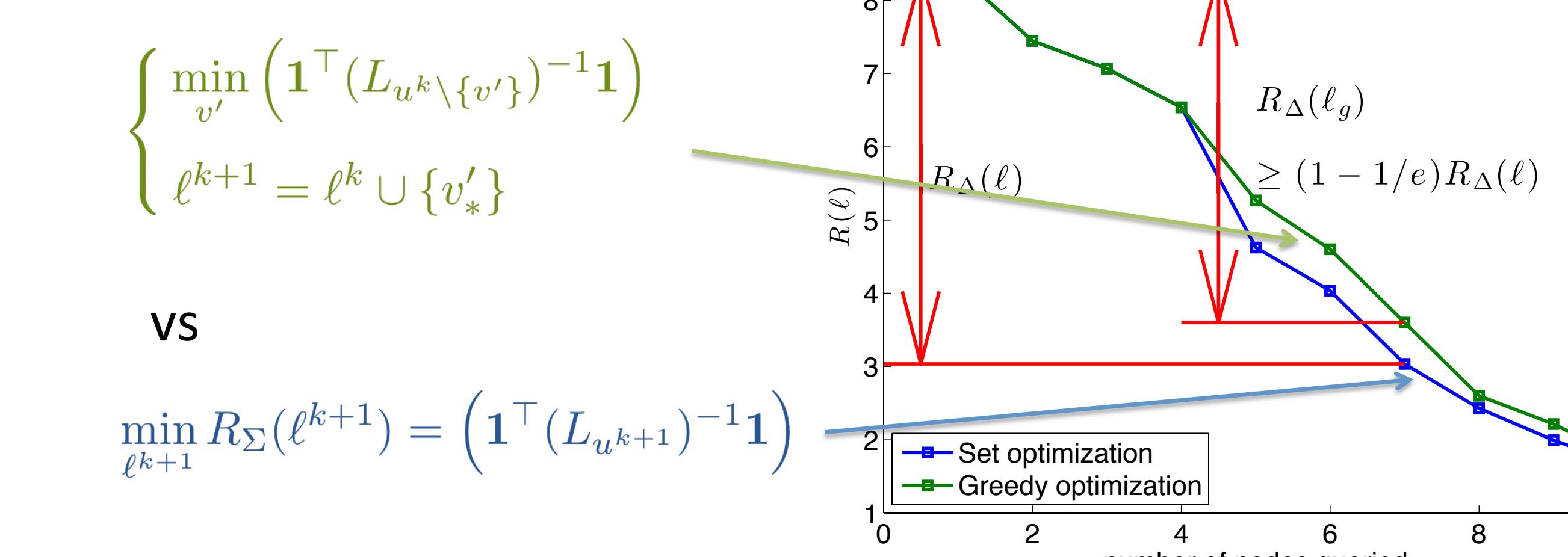
**Surveying** is to predict class proportions.

### □ All active learning strategies considered (at step $k$ with $u^k$ unlabeled)

◊ <b>Σ-Optimality</b>	$\min_{v'} (\mathbf{1}^T (L_{u^k \setminus \{v'\}})^{-1} \mathbf{1})$	ours
◊ <b>V-Optimality<sup>2</sup></b>	$\min_{v'} \text{tr}((L_{u^k \setminus \{v'\}})^{-1})$	Bayesian risk min
◊ <b>Info Gain (IG)<sup>3</sup></b>	$\max_{v'} (L_{u^k \setminus \{v'\}}^{-1})_{v', v'}$	(same as det-opt)
◊ <b>Mutual (MIG)<sup>3</sup></b>	$\max_{v'} (L_{u^k \setminus \{v'\}}^{-1})_{v', v'}/((L_{\ell^k \setminus \{v'\}})^{-1})_{v', v'}$	
◊ <b>Uncertainty<sup>4</sup></b>	$\min_{v'}  \hat{y}_{v'} $	
◊ <b>E Error (EER)<sup>4</sup></b>	$\max_{v'} \mathbb{E}_{y_{v'}} \left[ \left( \sum_{u_i \in u}  \hat{y}_{u_i}  \mid y_{v'} \right) \mid y_{\ell^k} \right]$	

## Theoretical Properties

- ◊ Assume the graph Laplacian  $L$  is connected and diagonal dominant, then the risk reduction  $R_\Delta(\ell) = R(\ell_0) - R(\ell_0 \cup \ell)$  is positive, monotone, and submodular.
- ◊ Greedy application of Bayesian risk minimization has bounded approximation ratio (1-1/e).



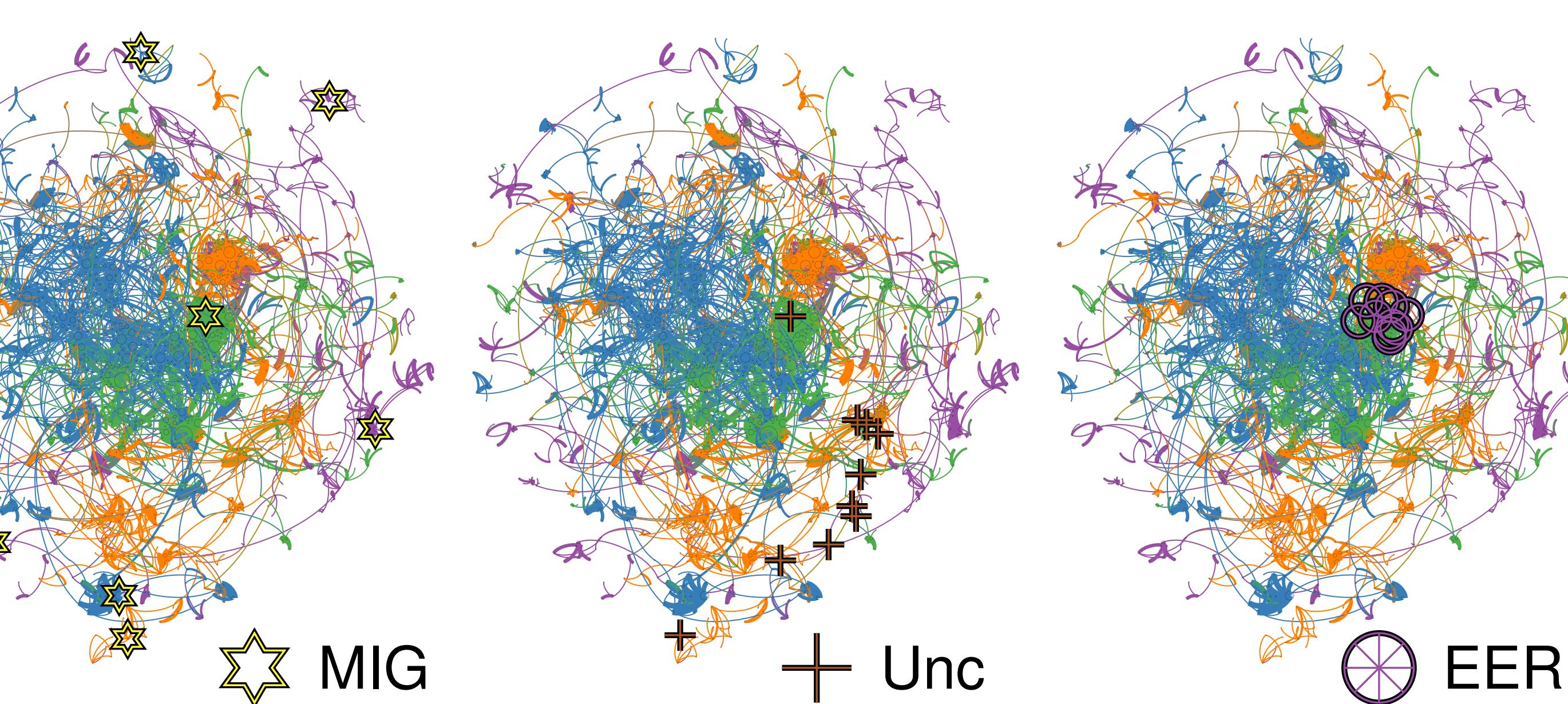
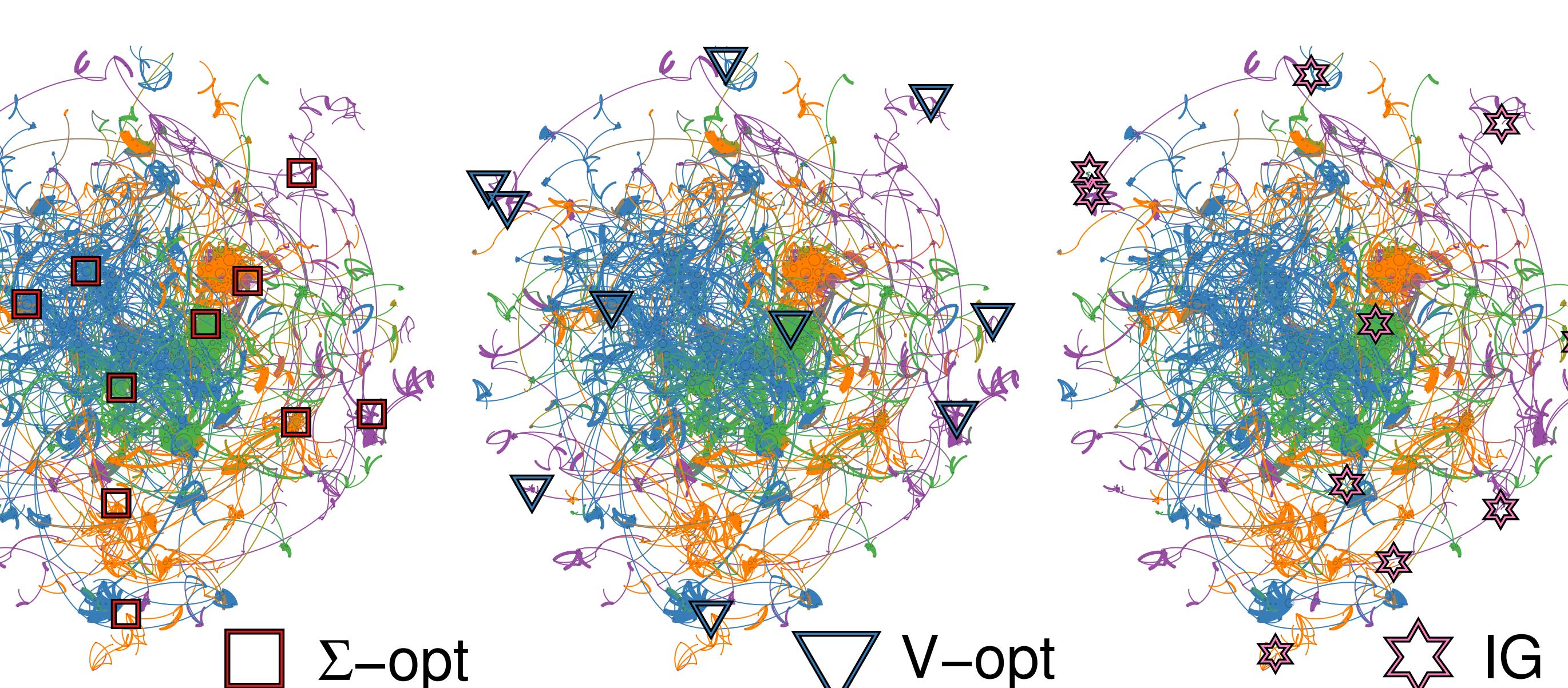
## Observations

Let  $(L_{u^k}^{-1})_{ij} = \Sigma_{ij} = \rho_{ij}\sigma_i\sigma_j$  be the conditional covariance.

Using rank one update, we have an equivalent selection rule:

$$v^{k+1} = \arg \max_v \left\{ \begin{array}{l} \left( \sum_{t \in u} \rho_{vt} \sigma_t \right)^2 \\ \sum_{t \in u} (\rho_{vt} \sigma_t)^2 \end{array} \right\} \quad \begin{array}{l} \square \quad (1\text{-norm}) \\ \triangle \quad (2\text{-norm}) \end{array}$$

The idea: 1-norm might be less susceptible to outliers.



### ◊ Proof:

► Suppose  $L = D - W$  diagonal dominant  $\Rightarrow$

$$(L_{uu}^{-1} L_{u\ell})^{-1} - \begin{pmatrix} L_{uu}^{-1} & 0 \\ 0 & 0 \end{pmatrix} = \begin{pmatrix} -L_{uu}^{-1} L_{u\ell} \\ 0 \end{pmatrix} (L_{\ell\ell} - L_{\ell u} L_{uu}^{-1} L_{u\ell})^{-1} (-L_{uu}^{-1}, \mathbf{1}) \geq 0$$

► Monotonicity of  $\ell \subset \ell'$   $\Rightarrow$

$$L_u^{-1} \supset L_{u'}^{-1} \Rightarrow \text{tr}(L_u^{-1}) \geq \text{tr}(L_{u'}^{-1}) \Rightarrow R_\Delta(\ell) \leq R_\Delta(\ell')$$

► Diminishing marginal reward of a singleton  $\ell$

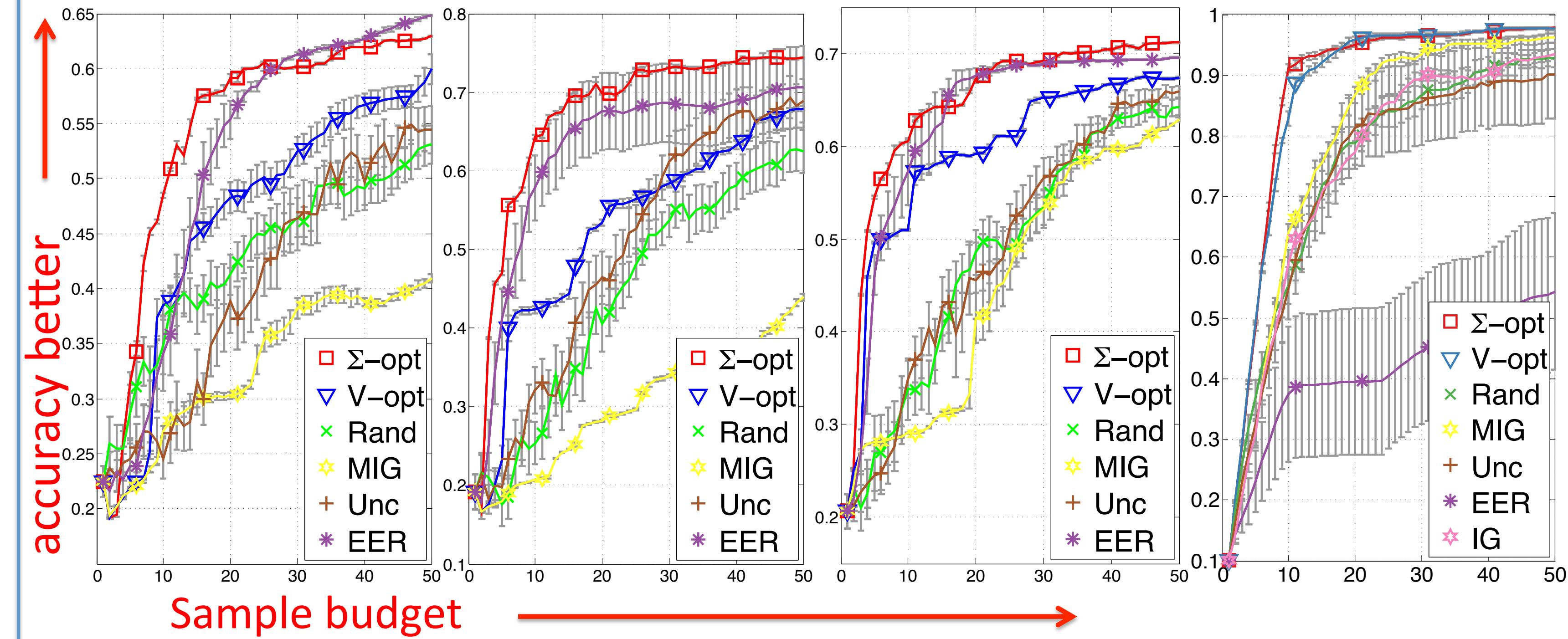
$$\Delta(\ell) = \text{tr}(L^{-1}) - \text{tr}(L_{\ell\ell}^{-1}) = \frac{\|(-L_{\ell u} L_{uu}^{-1}, \mathbf{1})^T\|^2}{L_{\ell\ell} - (-L_{\ell u} L_{uu}^{-1})^T (-L_{\ell u})}$$

It decreases as  $L$  shrinks in size nestingly, i.e., as  $\ell$  is queried later.

◊ Bonus: The GRF is a subclass of suppressor-free GPs. Knowing more decouples the unknown. It extends the conditional independence idea. Suppressor-free is desirable.

## Empirical Results

### □ Active classification



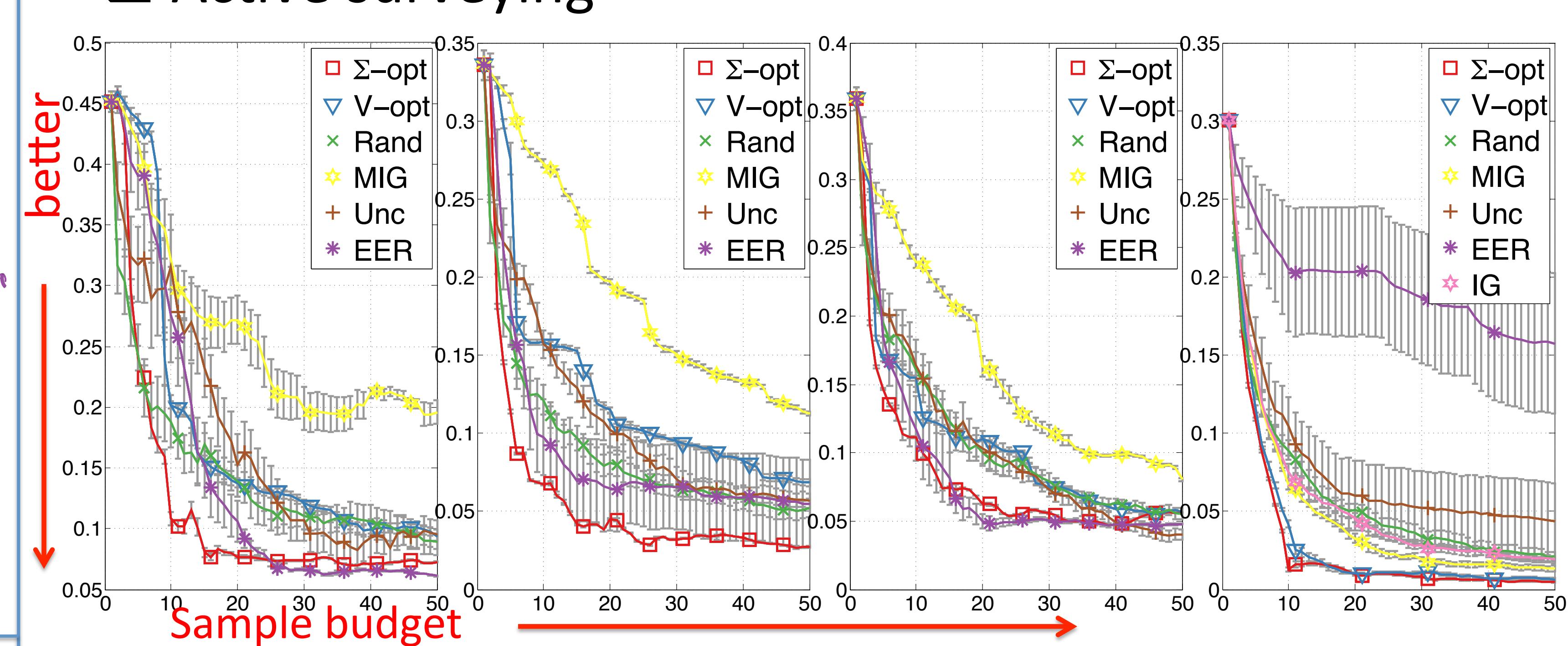
DBLP coauthorship graph  
1711 nodes (scholars). Labels: Machine learning / data mining / information retrieval / database  
2898 edges: Coauthorship  
Used in Ji & Han 2012

Cora citation graph  
2485 nodes (papers). Labels: Case based / Genetic algorithms / Neural networks / Probabilistic methods / Reinforcement learning / Rule learning / Theory  
3665 edges: Undirected citation  
Used in Sen et al 2008

CiteSeer citation graph  
2109 nodes (papers) labels: Agents / AI / DB / IR / ML / HCI  
4-nearest-neighbors graph of Euclidean distances between concatenations of raw pixels.  
Labels: actual digit in images  
Random subsample 70% digits From UCI Repository 1998

4-nn handwritten digits  
1797 images of 8x8 resolution of 4-nearest-neighbors graph of Euclidean distances between concatenations of raw pixels.  
Labels: actual digit in images  
Random subsample 70% digits From UCI Repository 1998

### □ Active surveying



### □ Active regression (omitted): V-opt worked better

## Reference

- Methods
- Zhu, Xiaojin, Lafferty, John, and Ghahramani, Zoubin. Combining active learning and semi-supervised learning using gaussian fields and harmonic functions. In *ICML 2003 workshop on The Continuum from Labeled to Unlabeled Data in Machine Learning and Data Mining*, 2003.
  - Ji, Ming and Han, Jiwei. A variance minimization criterion to active learning on graphs. In *AISTAT*, 2012.
  - Krause, Andreas, Singh, Ajit, and Guestrin, Carlos. Near-optimal sensor placements in gaussian processes: Theory, efficient algorithms and empirical studies. *JMLR*, 2008.
  - Settles, Burr. Active learning literature survey. 2010.
- Theory
- Das, Abhiram and Kempe, David. Algorithms for subset selection in linear regression. *ACM symposium on Theory of computing*, 2008.
  - Friedland, S and Gaubert, S. Submodular spectral functions of principal submatrices of a hermitian matrix, extensions and applications. *Linear Algebra and its Applications*, 2011.
- Analysis
- Garnett, Roman, Krishnamurthy, Yamuna, Xiong, Xuehan, Schneider, Jeff, and Mann, Richard. Bayesian optimal active search and surveying. In *ICML*, 2012.
  - Martin, Shawn, Brown, W Michael, Klavans, Richard, and Boyack, Kevin W. Openord: an open-source toolbox for large graph layout. In *IS&T/SPIE Electronic Imaging*, 2011.

