

# Entity Resolution with Active Learning

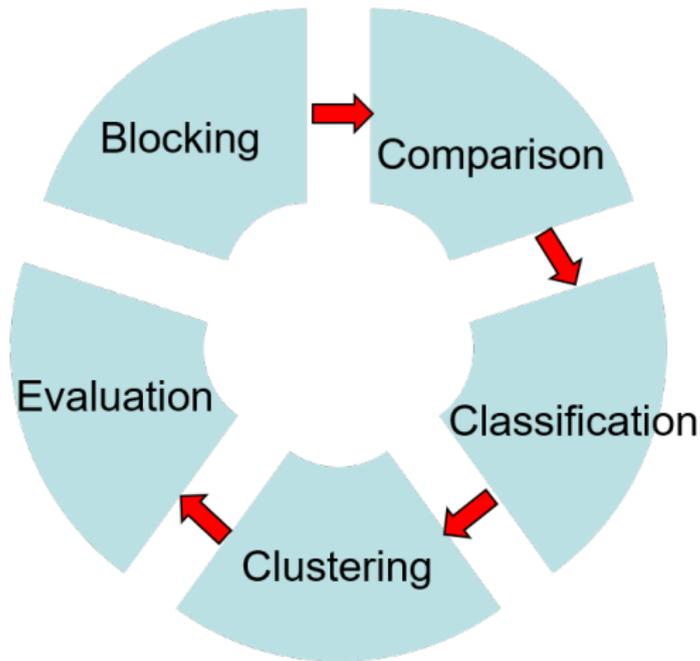
Jingyu Shao

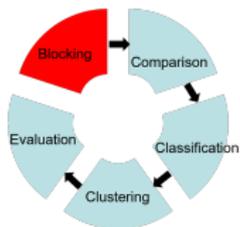
School of Computing  
Australian National University

May 24, 2021

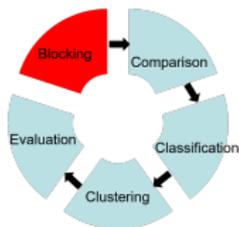
- \* Introduction and challenges
- \* How to build a set of “optimal” blocking schemes efficiently?
- \* How to design an AL approach under various data distributions?
- \* How to alleviate the overfitting problem for powerful models?
- \* Conclusion

The process of identifying **records which represent the same real-world entity** from one or more datasets



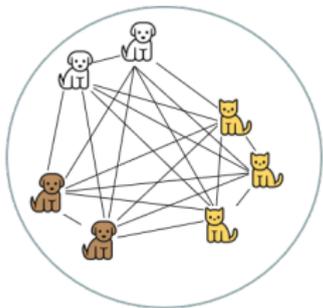


Reduce the number of record pairs to be compared by **grouping potentially matched records** into the same block.  
E.g., millions of pairs in real life.

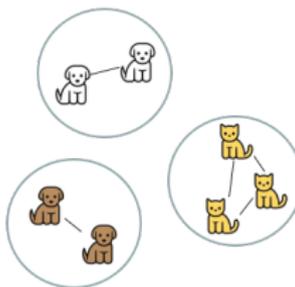


Reduce the number of record pairs to be compared by **grouping potentially matched records** into the same block.  
E.g., millions of pairs in real life.

Without blocking:  
7 records with 21 pairs

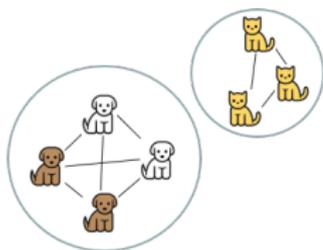


With blocking:  
7 records with 5 pairs

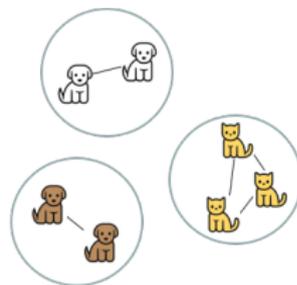


Using blocking schemes: (Which is better?)

Name

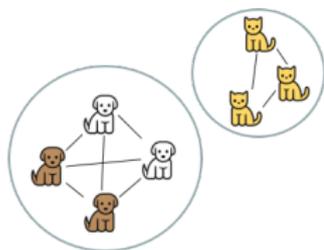


Color

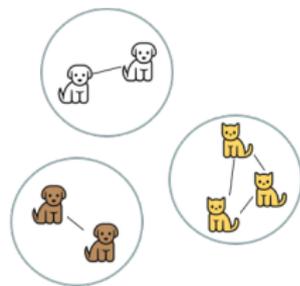


Using blocking schemes: (Which is better?)

Name

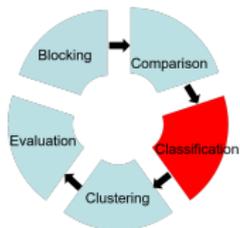


Color

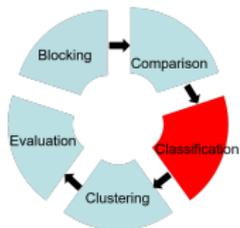


How to learn a good blocking scheme?

- Millions of record pairs, with highly imbalanced labels hard to obtain.
- The search space for all possible blocking schemes is large.

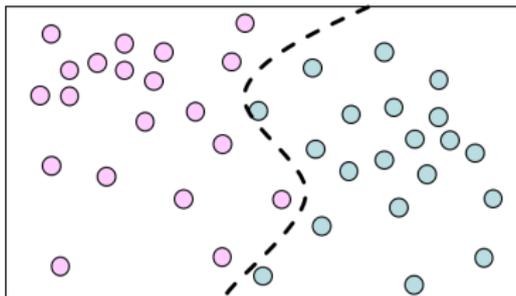


A classifier is used to categorize samples into **matches** and **non-matches**.



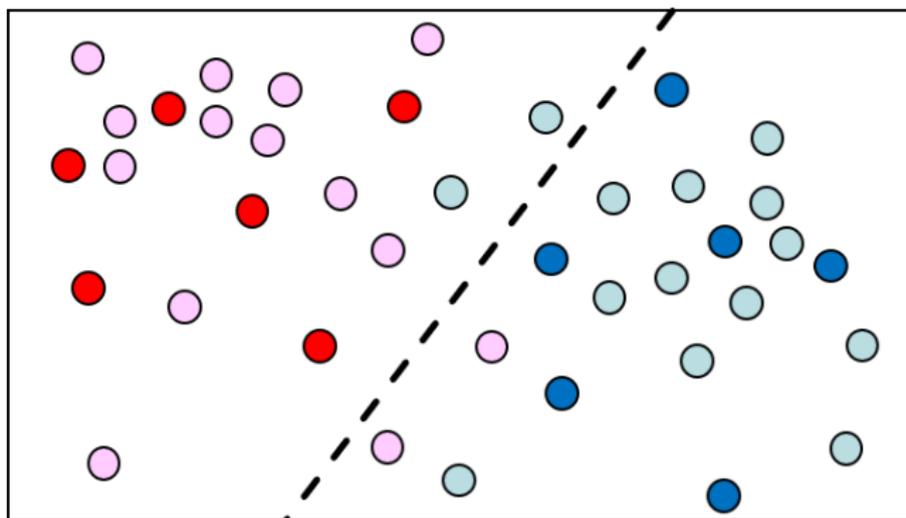
A classifier is used to categorize samples into **matches** and **non-matches**.

Considering we have samples within a block, and they are mapped into a feature space shown as below:



The red and blue points refer to matches and non-matches

Sufficient number of samples are necessary for training, but obtaining their labels for learning is costly.



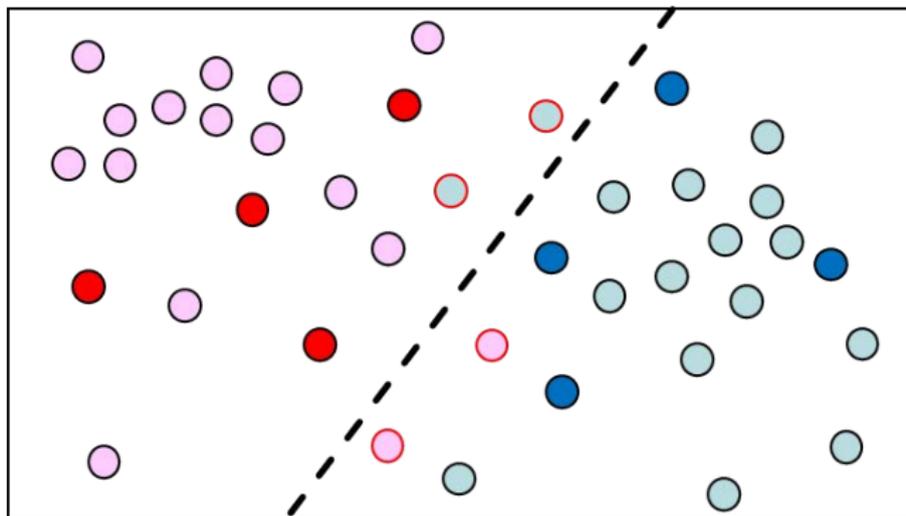
Accuracy: 90%

\* Red: 6/20

\* Blue: 6/20

Initialization: random seed samples

Select the most uncertain instances



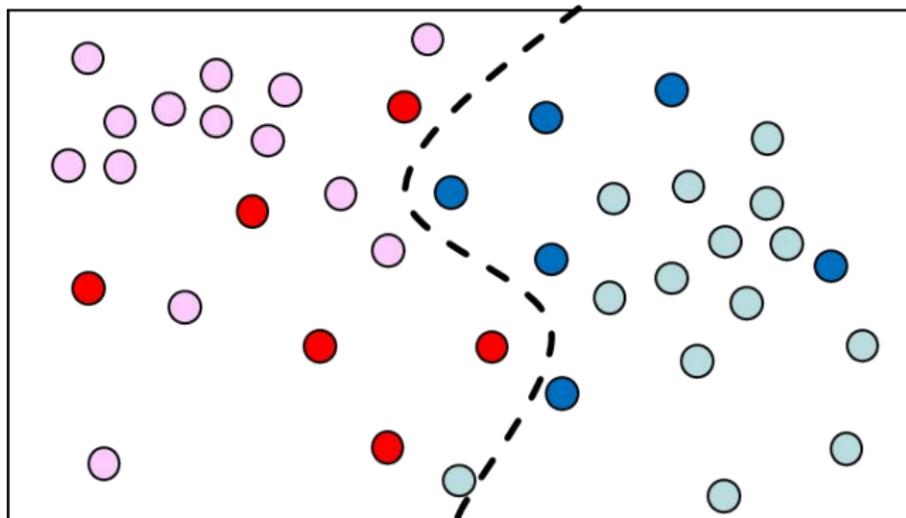
Accuracy: 90%

\* Red: 4/20

\* Blue: 4/20 <sup>1</sup>

<sup>1</sup>B. Settles, Active learning literature survey, 2010

4 more samples are labeled



Accuracy: 97.5%

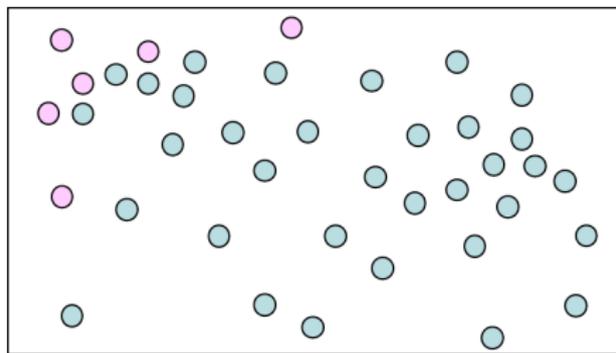
\* Red: 6/20

\* Blue: 6/20

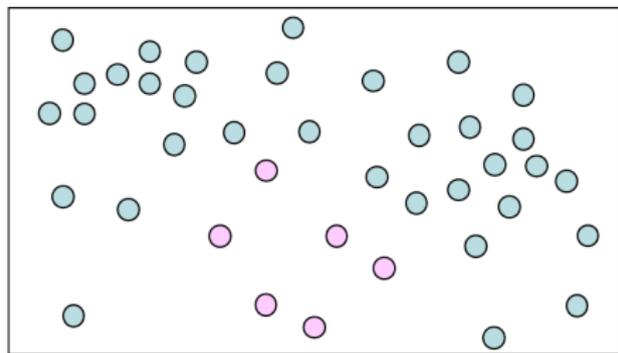
The distribution of matches and non-matches is highly imbalanced.

A small number of samples are labeled.

- Various strategies: different datasets



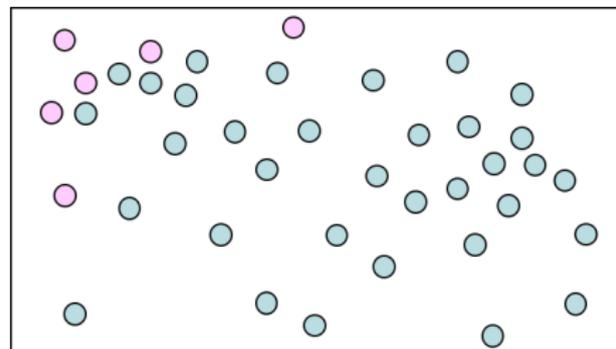
Sample distribution



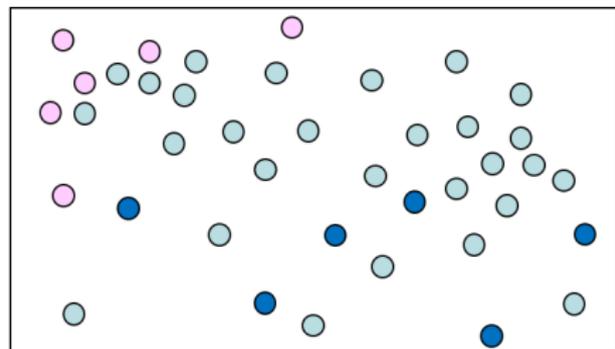
Sample distribution

The distribution of matches and non-matches is highly imbalanced.  
A small number of samples are labeled.

- Various strategies: different datasets
- Cold start: imbalanced ER data distribution



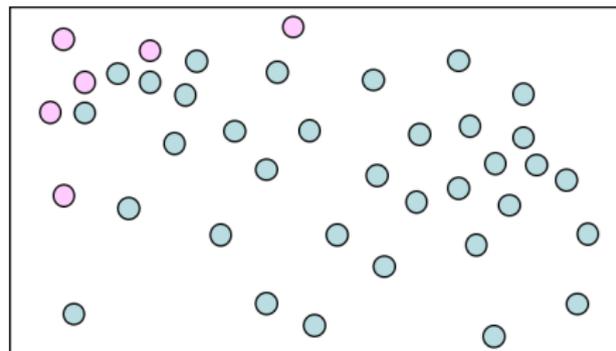
Sample distribution



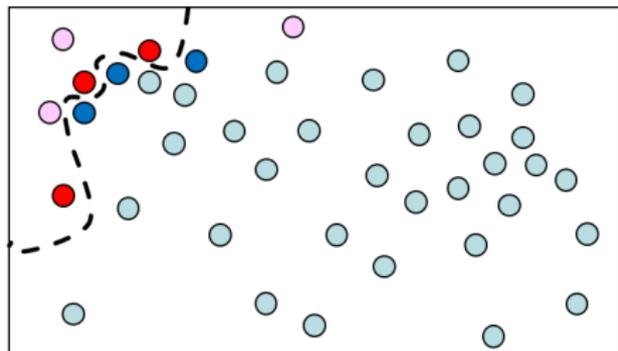
Cold Start

The distribution of matches and non-matches is highly imbalanced.  
A small number of samples are labeled.

- Various strategies: different datasets
- Cold start: imbalanced ER data distribution
- Overfitting: powerful models



Sample distribution



Overfitting

- \* Introduction and challenges
- \* How to build a set of “optimal” blocking schemes efficiently?
  - Active scheme learning and scheme skyline learning<sup>12</sup>
- \* How to design an AL approach under various data distributions?
- \* How to alleviate the overfitting problem for powerful models?
- \* Conclusion

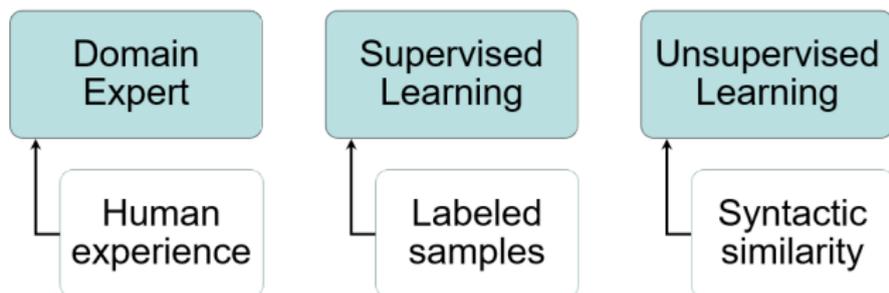
---

<sup>1</sup>J. Shao and Q. Wang. Active Blocking Scheme Learning for Entity Resolution. PAKDD'18.

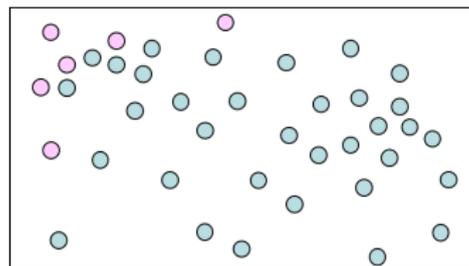
<sup>2</sup>J. Shao, Q. Wang and Y. Lin. Skyblocking for Entity Resolution. IS'19.

Disjunction of conjunction of attributes

Blocking schemes are built from:



Class Imbalance Problem:



Sample distribution

Large Search Space  $2^{\binom{n}{2}}$ :

Possible  
schemes for  
attributes  
 $\{A, B, C, D\}$

$A \wedge B, A \vee B$   
 $A \wedge B \vee C$   
 $(A \wedge B) \vee (C \wedge D)$   
...

Our observation: similar attribute values –  $>$  matches

How to select attributes to build schemes?

Some values are frequent but useless, e.g. year.

Balanced samples for all possible attributes and select!

**Balance Rate**  $\gamma(s, X)$  describes the balance degree for a given scheme  $s$  under a sample set  $X$

Our observation: similar attribute values –  $>$  matches

How to select attributes to build schemes?

Some values are frequent but useless, e.g. year.

Balanced samples for all possible attributes and select!

**Balance Rate**  $\gamma(s, X)$  describes the balance degree for a given scheme  $s$  under a sample set  $X$

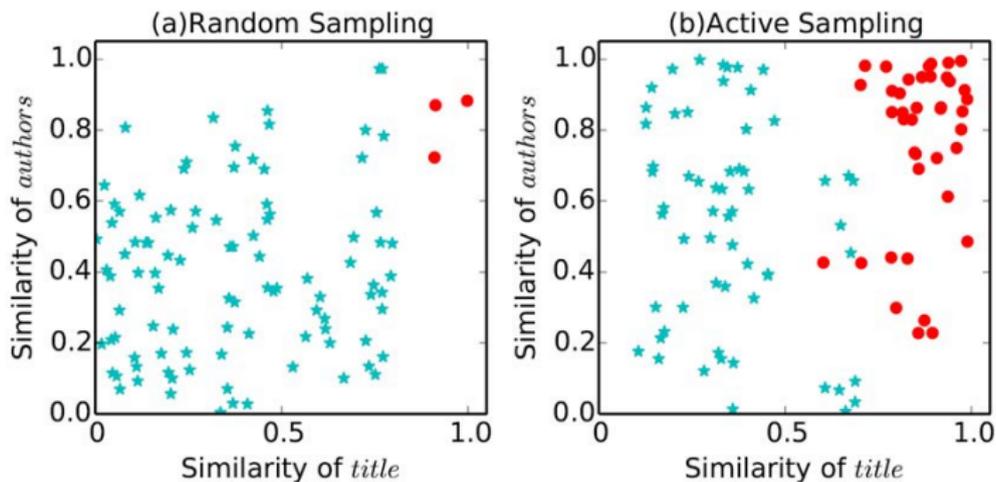
E.g. if  $s = A \wedge B$ ,  $X = \{x_1, x_2\}$ , then  $s(x_1) = \text{true}$ ,  $s(x_2) = \text{false}$

Thus  $\gamma(s, X) = \frac{1(\#true) - 1(\#false)}{2} = 0$  (balanced):

	A	B	C	D
$x_1$	1	1	0	1
$x_2$	0	1	0	1

Select samples to minimize the balance rate for a given set of schemes:

$$\text{minimize } \sum_{s_i \in S} \gamma(s_i, X)^2$$



Reduce the search space by extending “proper” schemes w.r.t. a specific criterion, e.g. Pair Completeness (Recall) and Pair Quality (Precision).

∧ reduce block size, increase PQ

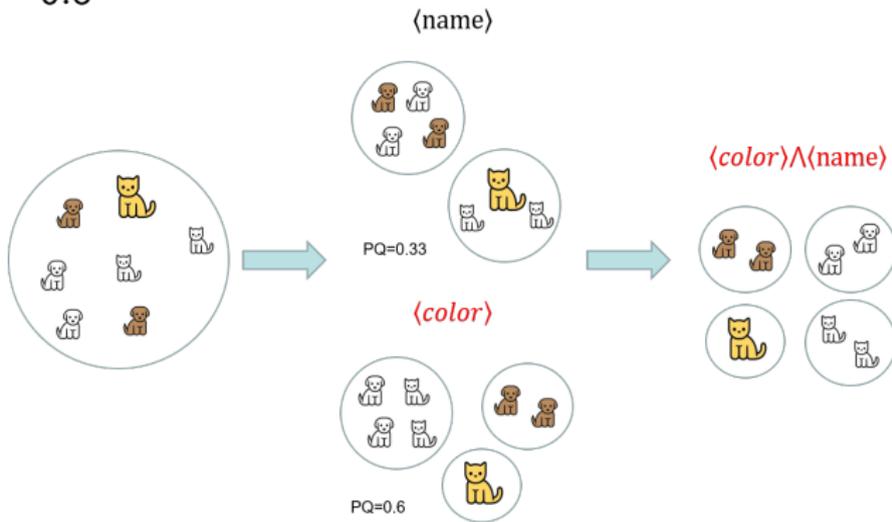
∨ increase block size, increase PC

Reduce the search space by extending “proper” schemes w.r.t. a specific criterion, e.g. Pair Completeness (Recall) and Pair Quality (Precision).

∧ reduce block size, increase PQ

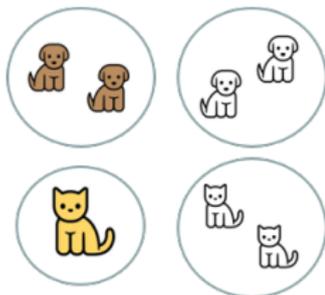
∨ increase block size, increase PC

Example: to learn a blocking scheme with two attributes:  $\langle name \rangle$ ,  $\langle color \rangle$ , w.r.t.  $PQ = 0.8$

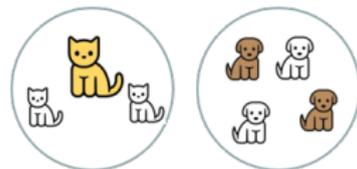


- The higher, the better (PC and PQ values)
- More records in one block (high PC threshold)
- Less records in one block (high PQ threshold)

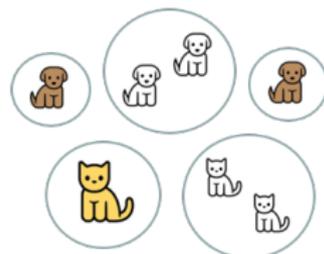
## Ideal Blocks



## Possible Blocks



High PC

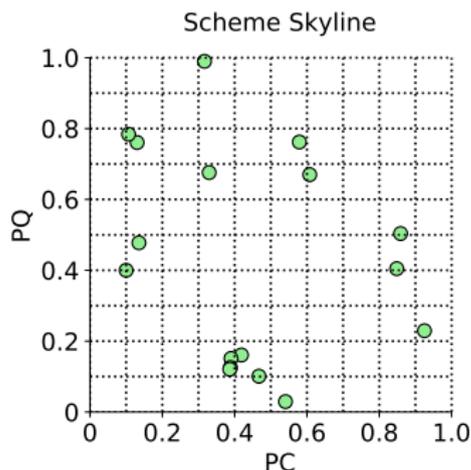


High PQ

Skyline queries under a set of blocking schemes:

Map schemes into a measure space

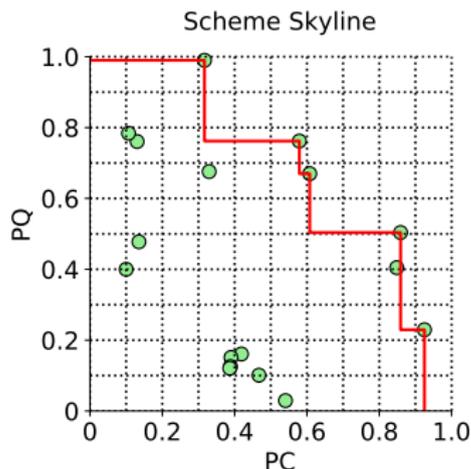
Blocking scheme	PC	PQ
$s_1$	0.13	0.76
$s_2$	0.31	0.99
$s_3$	0.58	0.76
$s_4$	0.84	0.40
$s_5$	0.86	0.50
...	...	...



Skyline queries under a set of blocking schemes:

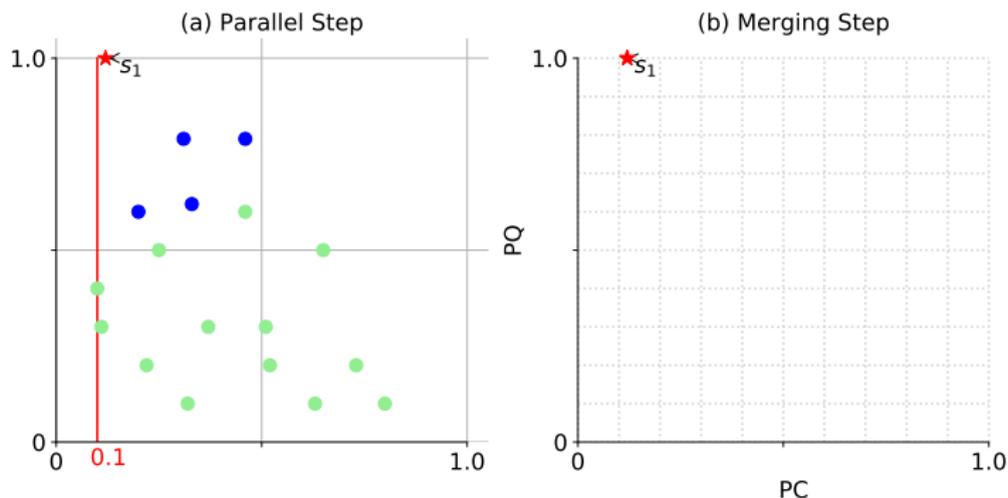
Dominated VS Dominating schemes

Blocking scheme	PC	PQ
$s_1$	0.13	0.76
$s_2$	0.31	0.99
$s_3$	0.58	0.76
$s_4$	0.84	0.40
$s_5$	0.86	0.50
...	...	...



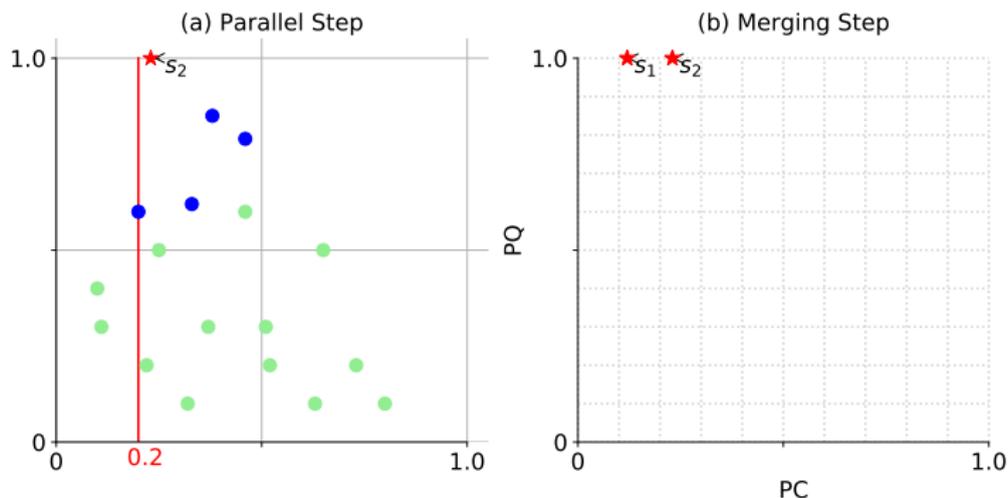
A naïve way to learn scheme skyline:

- Learn “optimal” schemes w.r.t. different thresholds



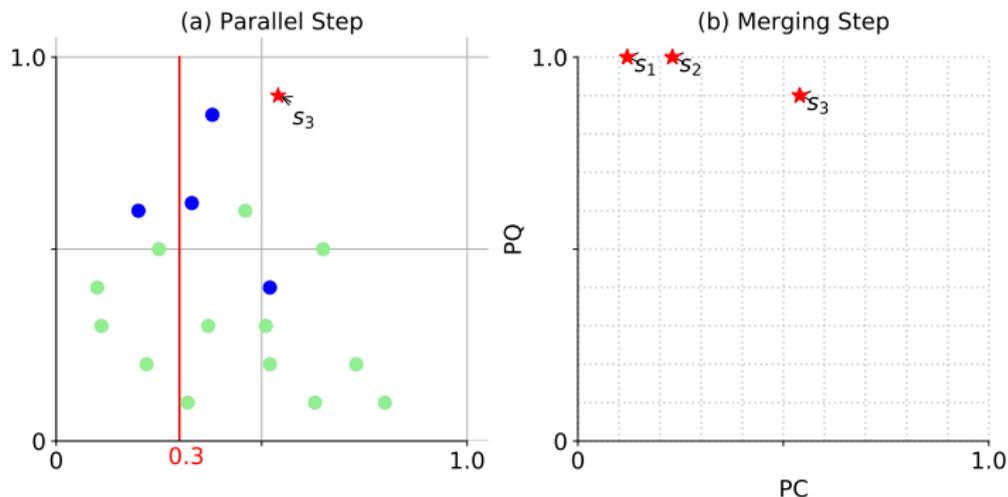
A naïve way to learn scheme skyline:

- New threshold: current one plus a threshold interval, e.g.  $\Delta = 0.1$



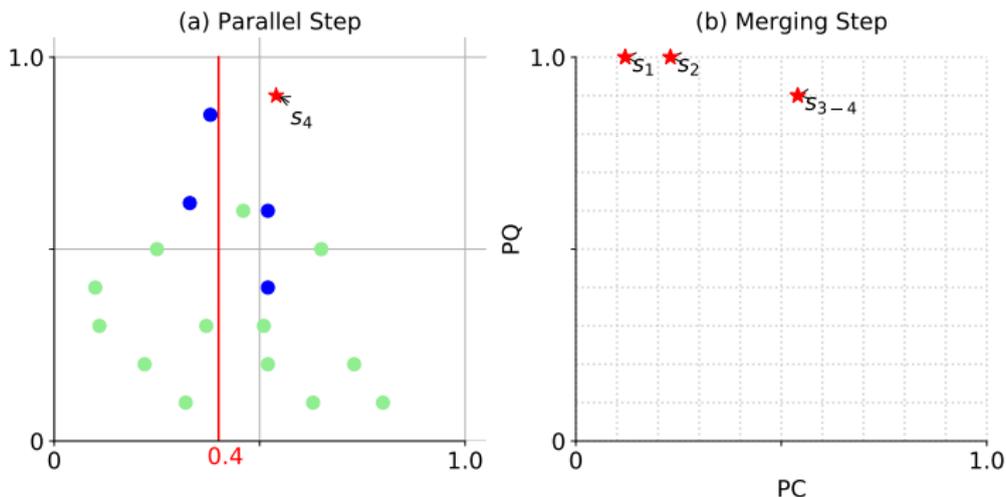
A naïve way to learn scheme skyline:

- New threshold: current one plus a threshold interval, e.g.  $\Delta = 0.1$



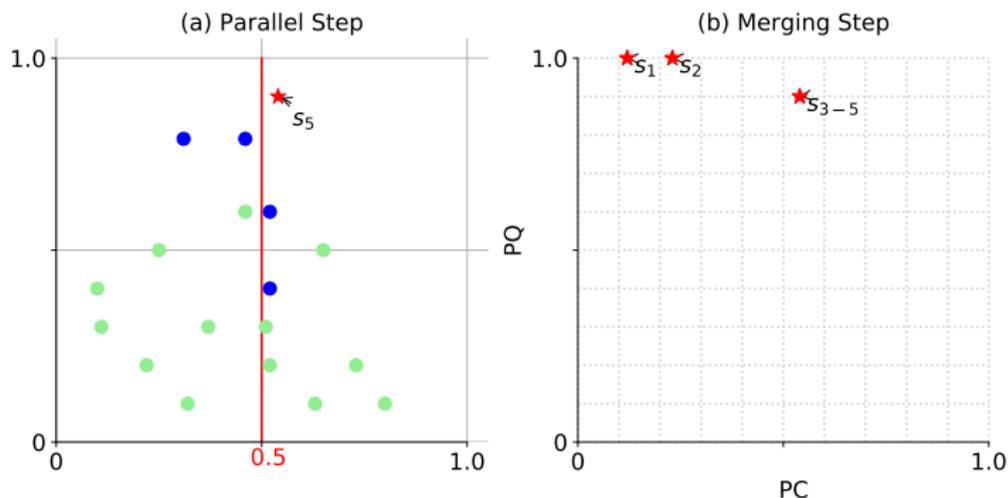
A naïve way to learn scheme skyline:

- New threshold: current one plus a threshold interval, e.g.  $\Delta = 0.1$



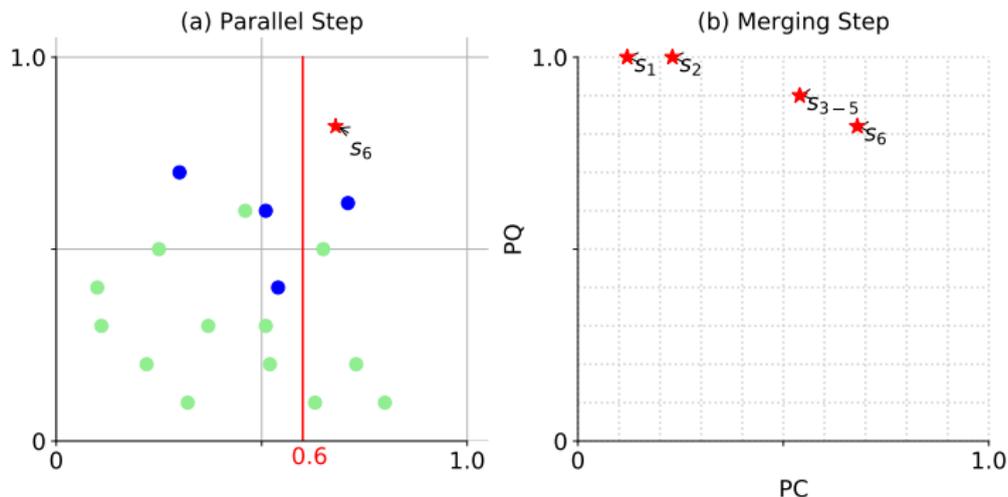
A naïve way to learn scheme skyline:

- New threshold: current one plus a threshold interval, e.g.  $\Delta = 0.1$



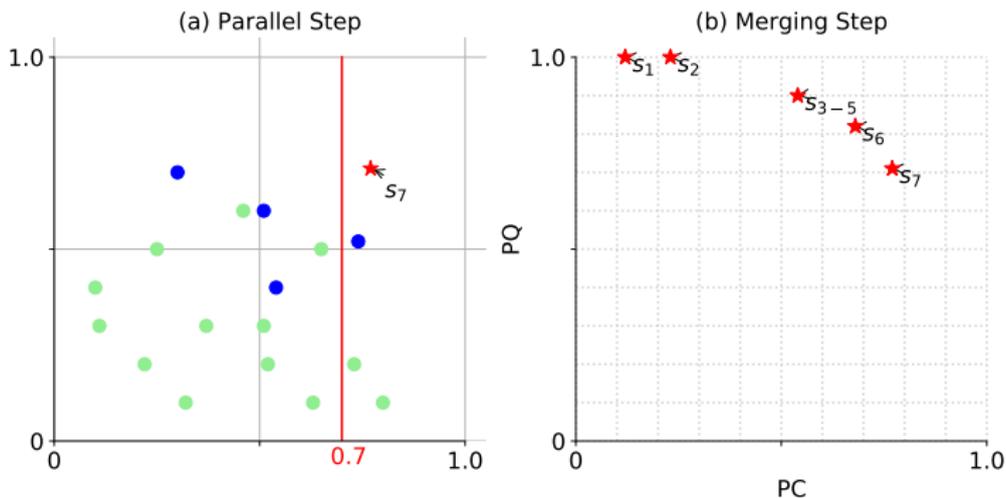
A naïve way to learn scheme skyline:

- New threshold: current one plus a threshold interval, e.g.  $\Delta = 0.1$



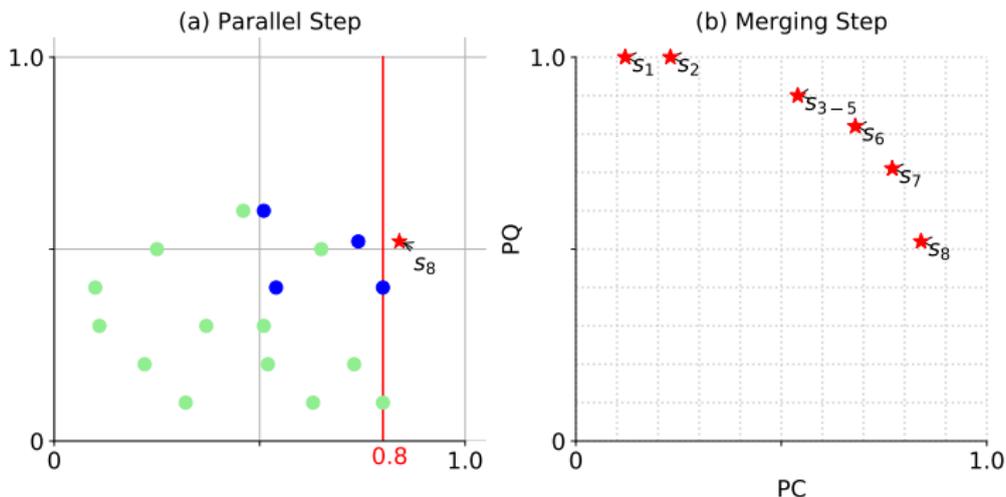
A naïve way to learn scheme skyline:

- New threshold: current one plus a threshold interval, e.g.  $\Delta = 0.1$



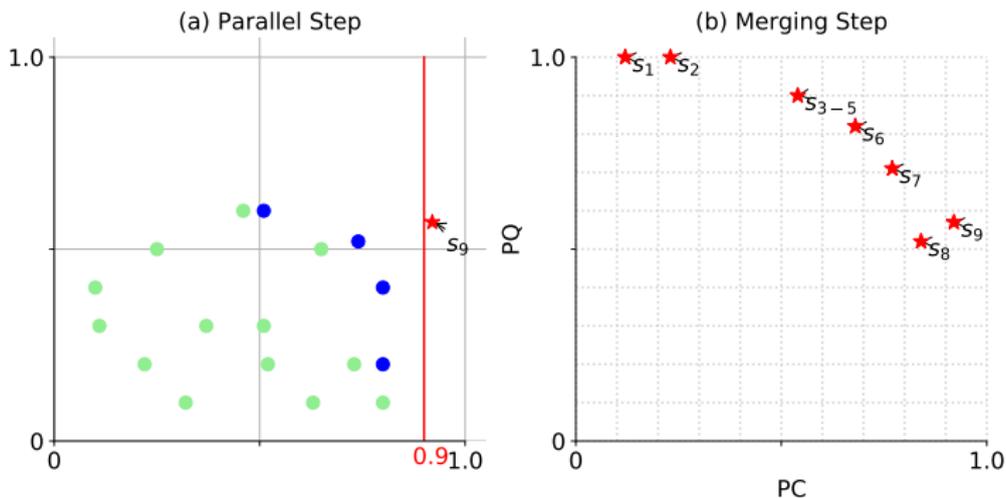
A naïve way to learn scheme skyline:

- New threshold: current one plus a threshold interval, e.g.  $\Delta = 0.1$



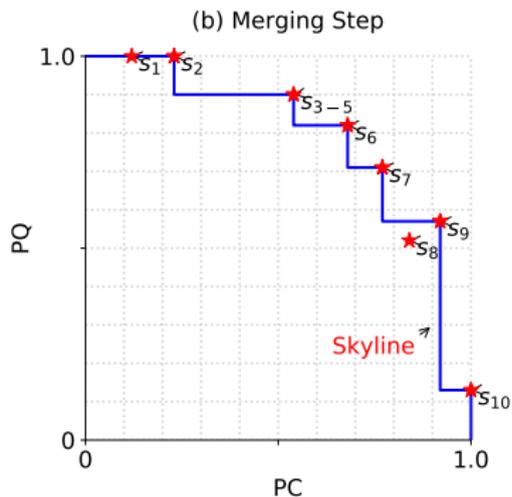
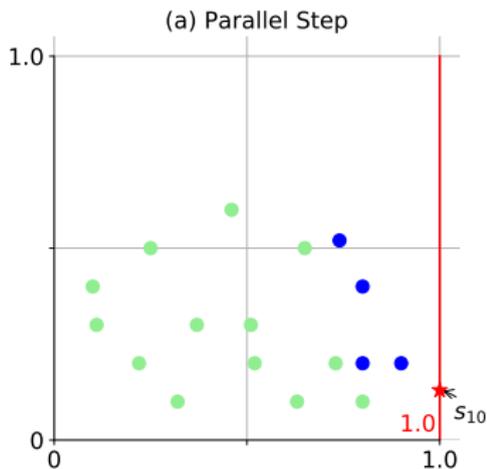
A naïve way to learn scheme skyline:

- New threshold: current one plus a threshold interval, e.g.  $\Delta = 0.1$



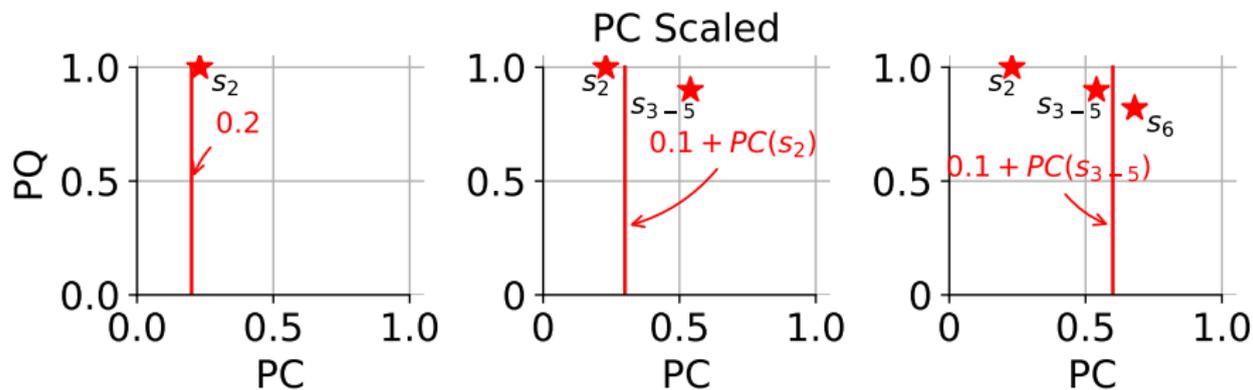
A naïve way to learn scheme skyline:

- Merge them for skyline



Observation: some are redundant under different thresholds: e.g.  $s_{3-5}$

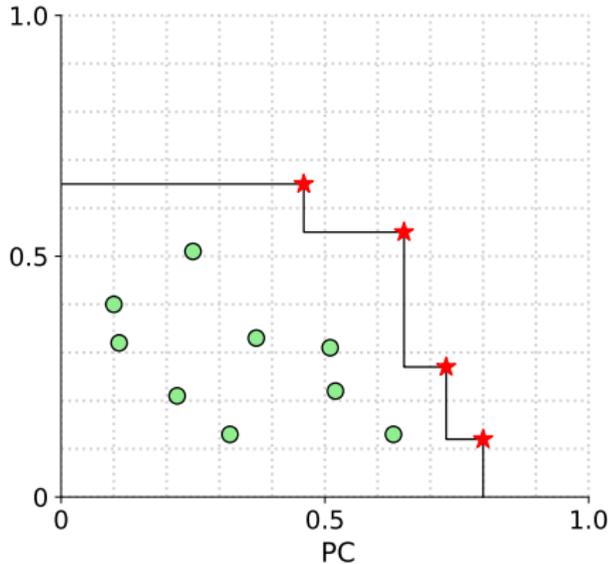
New threshold: PC/PQ value of current scheme plus a threshold interval



Unnecessary label cost in Adap-Sky: samples are independently selected and may be duplicated under different thresholds.

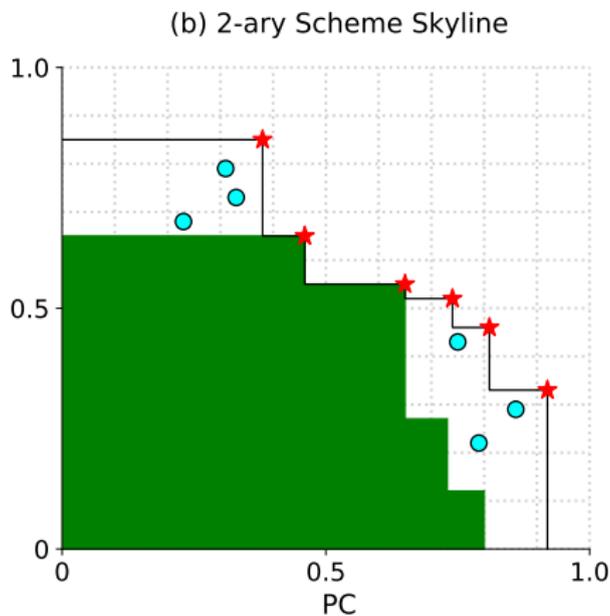
Pro-Sky with scheme extension:

(a) 1-ary Scheme Skyline



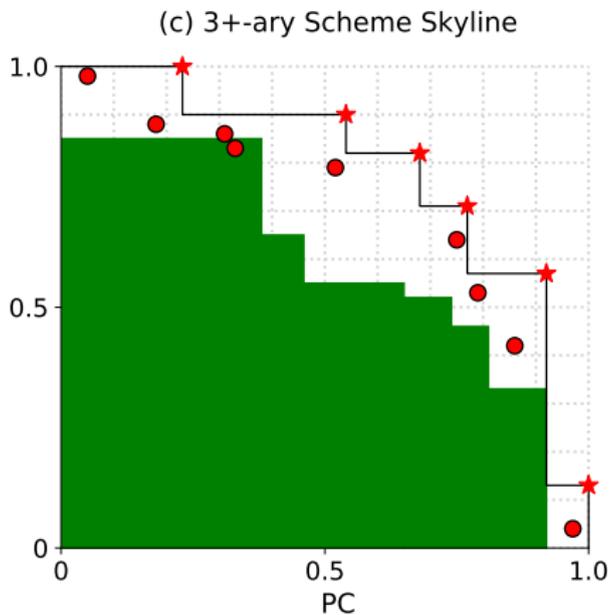
Unnecessary label cost in Adap-Sky: samples are independently selected and may be duplicated under different thresholds.

Pro-Sky with scheme extension:



Unnecessary label cost in Adap-Sky: samples are independently selected and may be duplicated under different thresholds.

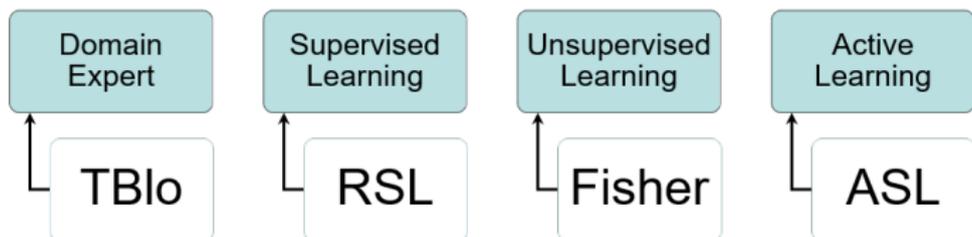
Pro-Sky with scheme extension:



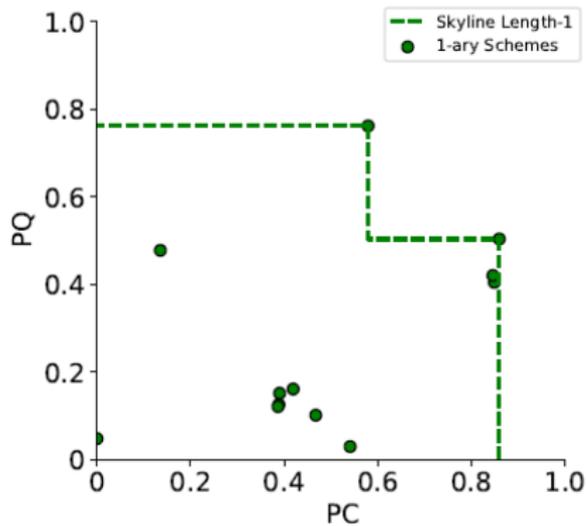
## Datasets

Dataset	# of Attributes	# of Records	Class Imbalance Ratio
Cora	4	1,295	1:49
DBLP - ACM	4	2,616/2,294	1:1,117
DBLP - Scholar	4	2,616:64,263	1:31,440
NCVR	18	267,716/278,262	1:2,692

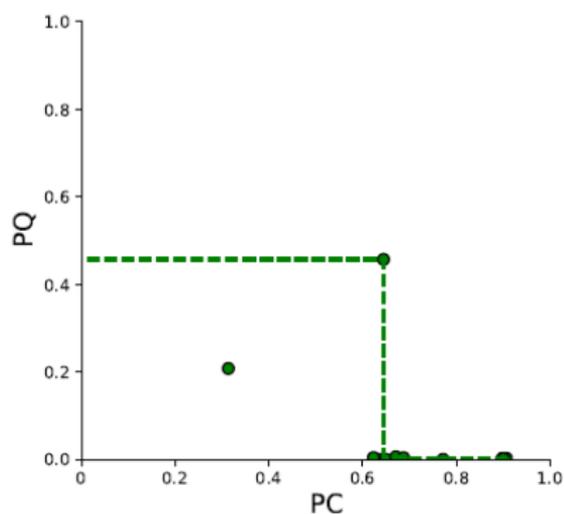
## Baselines

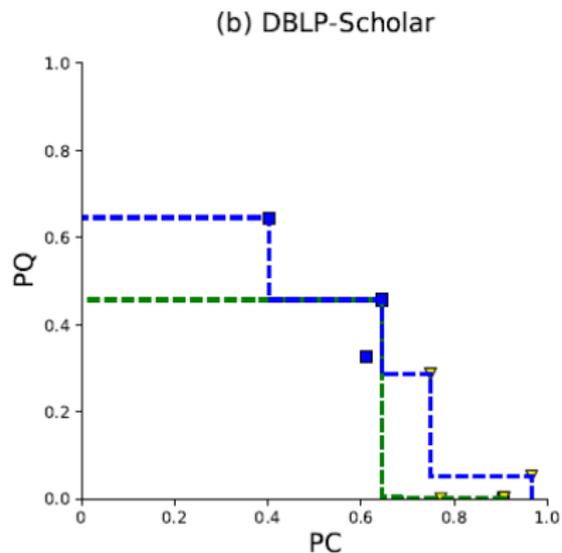
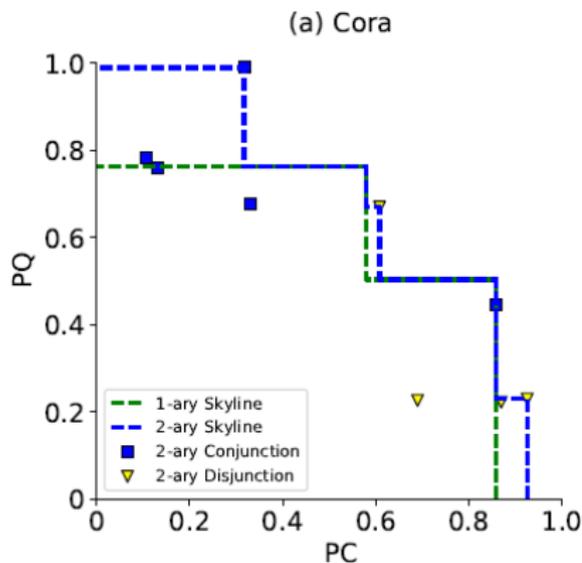


(a) Cora

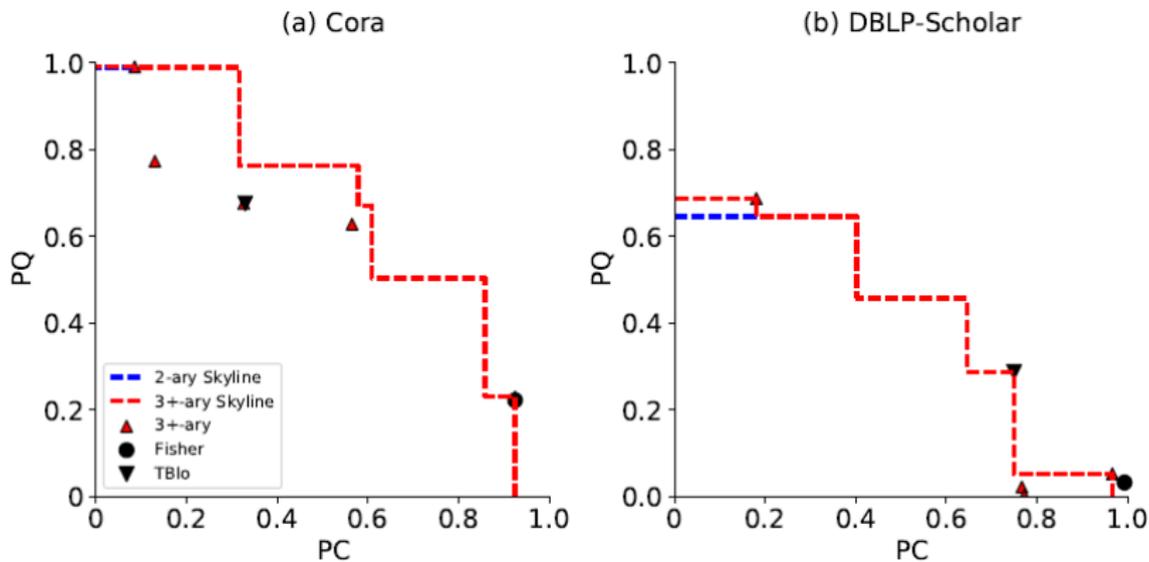


(b) DBLP-Scholar





# Scheme Skylines (Pro-Sky) in Experiments

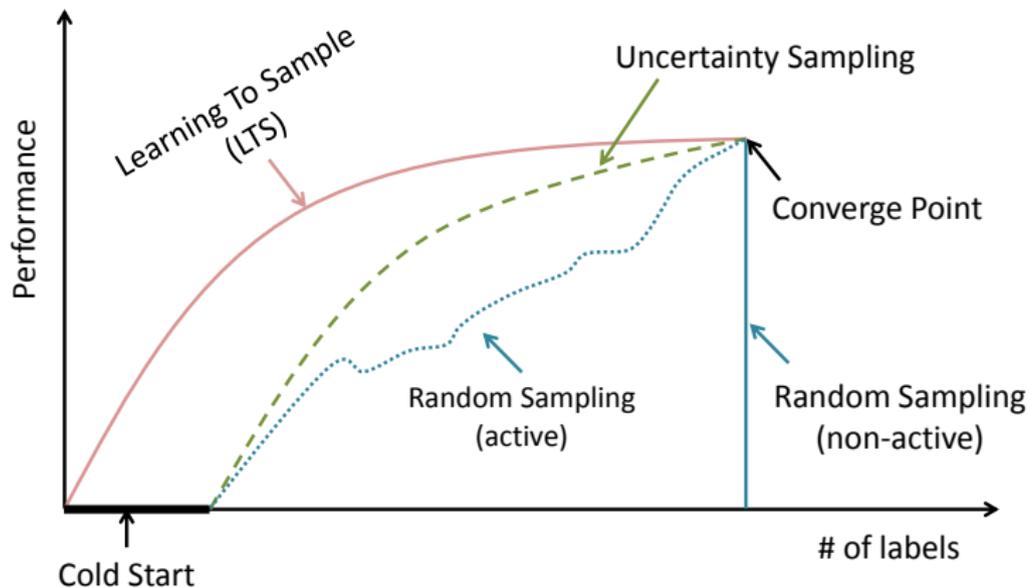


- \* Introduction and challenges
- \* How to build a set of “optimal” blocking schemes efficiently?
- \* **How to design an AL approach under various data distributions?**
  - Learning based active learning for ER <sup>1</sup>
- \* How to alleviate the overfitting problem for powerful models?
- \* Conclusion

---

<sup>1</sup>J. Shao, Q. Wang and F. Liu. Learning To Sample: an Active Learning Framework. ICDM'19.

To build an active learning framework:



## Challenges

- \* No one-fit-all: the “best” active learning strategy varies due to different datasets and machine learning models.
- \* Cold start problem: occurs under limited highly imbalanced samples.

## Challenges

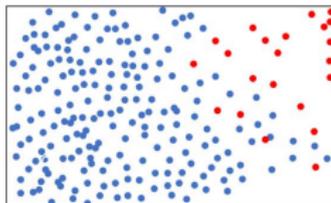
- \* No one-fit-all: the “best” active learning strategy varies due to different datasets and machine learning models.
- \* Cold start problem: occurs under limited highly imbalanced samples.

## Solution

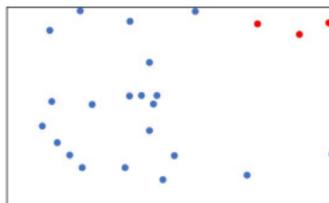
- Dynamical estimation of model performance (learning-based)
- Uncertainty and diversity of samples

Uncertainty sampling: function-based uncertainty measures

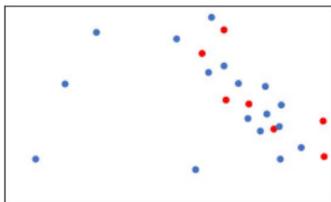
Diversity sampling: considering sample distribution (feature values)



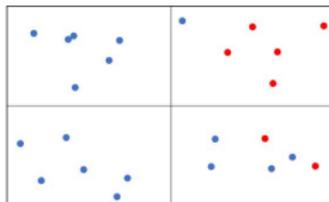
(a) Entire Data Distribution



(b) Random Sampling



(c) Uncertainty Sampling

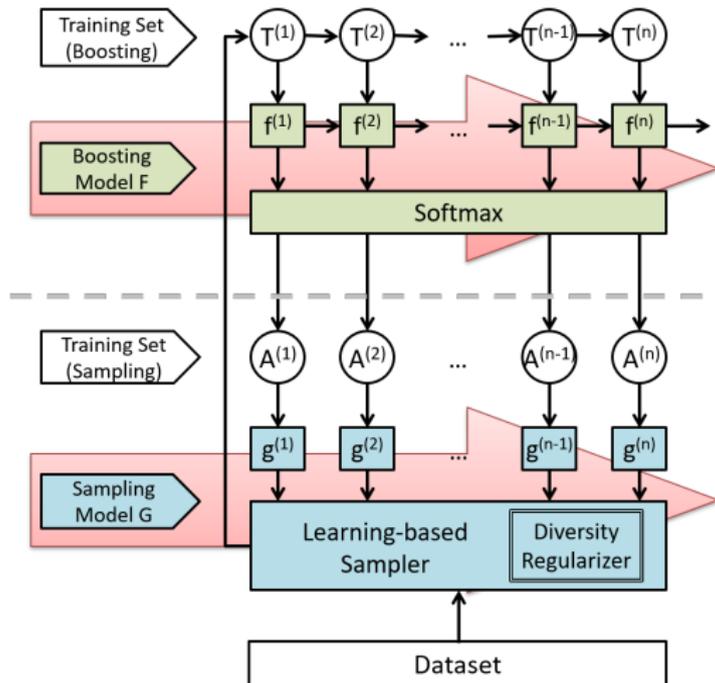


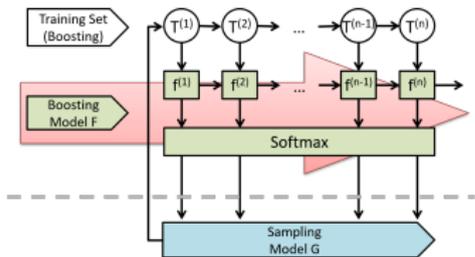
(d) Diversity Sampling  
(4 groups)

# Framework: Learning to Sample (LTS)



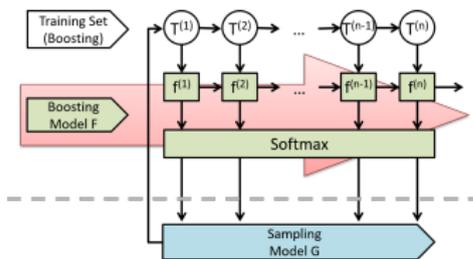
Two models dynamically learn from each other in iterations for performance improvement.





The boosting model  $F$  is a set of classifiers  $\langle f^{(1)}, \dots, f^{(n)} \rangle$ .

A classifier  $f^{(t)} \in F$  at the  $t$ -th iteration is trained by minimizing:



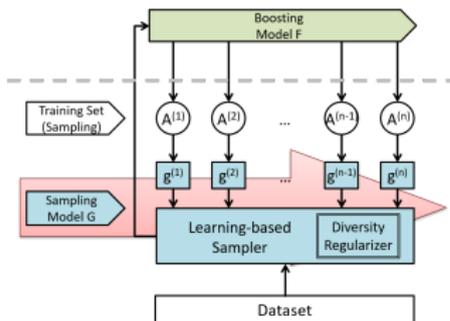
The boosting model  $F$  is a set of classifiers  $\langle f^{(1)}, \dots, f^{(n)} \rangle$ .

A classifier  $f^{(t)} \in F$  at the  $t$ -th iteration is trained by minimizing:

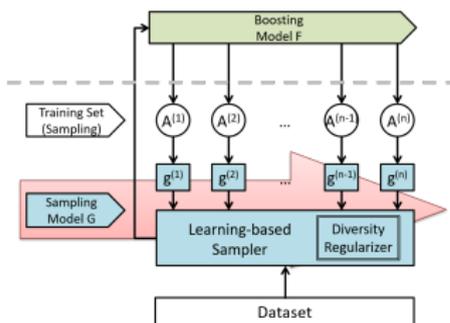
$$\sum_{(x_i, y_i) \in T^{(t)}} \ell_1(\hat{y}_i^{(t-1)} + f^{(t)}(x_i), y_i) + \Omega_1(f^{(t)})$$

where:

- $T^{(t)}$ : training set;
- $\hat{y}_i^{(t-1)} = \sum_{k=1}^{t-1} f^{(k)}(x_i)$ : predicted label of  $x_i$ ;
- $\ell_1$ : a differentiable loss function;
- $\Omega_1(f^{(t)})$ : the complexity penalty for  $f^{(t)}$ .



The sampling model  $G$  actively selects a set  $\Delta^{(t)}$  of uncertainty and diversity samples at the  $t$ -th iteration by:

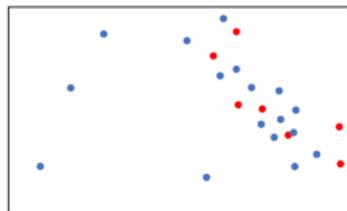


The sampling model  $G$  actively selects a set  $\Delta^{(t)}$  of uncertainty and diversity samples at the  $t$ -th iteration by:

$$\begin{aligned} &\text{maximize} && \sum_{i=1}^k v_i g^{(t)}(x_i) + \alpha \times \Gamma(\mathbf{v}) \\ &\text{subject to} && \|\mathbf{v}\|_1 = |\Delta^{(t)}| \end{aligned}$$

where  $\mathbf{v} = (v_1, \dots, v_k)^T \in \{0, 1\}^k$ ,  $k$  is the number of samples, and  $\alpha$  is a parameter.

- A regressor  $g^{(t)}(x_i)$  for uncertainty sampling
- A regularizer  $\Gamma(\mathbf{v})$  for diversity sampling.



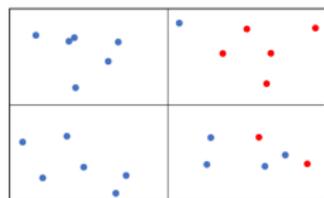
Uncertainty Sampling

A regressor is trained to predict the uncertainty of samples by minimizing:

$$\sum_{(x_i, z_i^{(t)}) \in A^{(t)}} w_i^{(t)} \ell_2(g^{(t)}(x_i), z_i^{(t)}) + \Omega_2(g^{(t)})$$

where:

- $A^{(t)} = \{(x_i, z_i^{(t)}) | x_i \in \mathcal{T}^{(t)}, z_i^{(t)} \in [0, 1]\}$ : uncertainty sample set;
- $z_i^{(t)}$ : the uncertainty of  $x_i$ ;
- $w_i^{(t)}$ : the weights of  $x_i$ ;
- $\ell_2$ : a differentiable loss function;
- $\Omega_2(g^{(t)})$ : the complexity penalty for  $g^{(t)}$ .



Diversity Sampling  
(4 groups)

The diversity  $\Gamma(\mathbf{v})$  is defined using a  $l_{2,1}$ -norm function:

$$\Gamma(\mathbf{v}) = \|\mathbf{v}\|_{2,1} = \sum_{j=1}^b \|\mathbf{v}_j\|_2$$

where:

- The sample space  $\mathbf{v}$  with  $b$  groups  $\{\mathbf{v}_1, \dots, \mathbf{v}_b\}$ ;
- The vector  $\mathbf{v}_j \in \{0, 1\}^m$  indicates samples selected in a group;
- Sample size  $m = |X_j^{(t)}|$  in a group.

# Results under Different Label Budgets



Dataset	Label Budget $\zeta$ (% of $ X $ )	CART	XG	XG+RS	XG + US $\alpha = 0$	XG+LTS $\alpha = 1$	XG + DS $\alpha \rightarrow \infty$
Cora	0.01	0	0	0	0	0.857	<b>0.878</b>
	0.05	0.741	0.763	0.750	0.827	0.864	<b>0.885</b>
	0.1	0.788	0.796	0.787	0.823	0.862	0.886
	0.5	0.848	0.835	0.835	0.873	<b>0.900</b>	0.893
	1	0.868	0.878	0.880	0.870	<b>0.902</b>	0.894
	5	0.878	0.897	0.892	0.907	<b>0.915</b>	0.898
NCVoter	0.01	0	0	0	0	0.324	<b>0.875</b>
	0.05	0	0	0	0	0.954	0.991
	0.1	0	0	0	0	<b>0.994</b>	0.993
	0.5	0	0	0	0	<b>0.994</b>	0.991
	1	0.334	0.379	0.398	0	0.993	<b>0.994</b>
	5	0.993	0.993	0.994	0.993	<b>0.997</b>	0.993
DBLP- ACM	0.1	0	0	0	0	0	<b>0.397</b>
	0.5	0	0	0	0	0.702	0.632
	1	0.348	0.347	0.279	0	<b>0.878</b>	0.721 3
	2	0.599	0.767	0.680	0.403	<b>0.884</b>	0.783
	5	0.870	0.850	0.803	0.874	<b>0.931</b>	0.833
	10	0.903	0.911	0.890	0.926	<b>0.981</b>	0.899
DBLP- Scholar	0.1	0	0	0	0	0.723	<b>0.731</b>
	0.5	0.378	0.54	0.498	0.555	0.773	<b>0.780</b>
	1	0.562	0.669	0.659	0.738	0.804	<b>0.792</b>
	2	0.772	0.806	0.771	0.807	<b>0.815</b>	0.801
	5	0.773	0.822	0.803	<b>0.836</b>	<b>0.836</b>	0.818 8
	10	0.808	0.835	0.830	<b>0.865</b>	0.851	0.829

- \* Introduction and challenges
- \* How to build a set of “optimal” blocking schemes efficiently?
- \* How to design an AL approach under various data distributions?
- \* **How to alleviate the overfitting problem for powerful models?**
  - A generative model with adversarial nets<sup>1</sup>
- \* Conclusion

---

<sup>1</sup>J. Shao, Q. Wang, A. Wijesinghe and E. Rahm. ERGAN: Generative Adversarial Networks for Entity Resolution. ICDM'20.



## Challenges

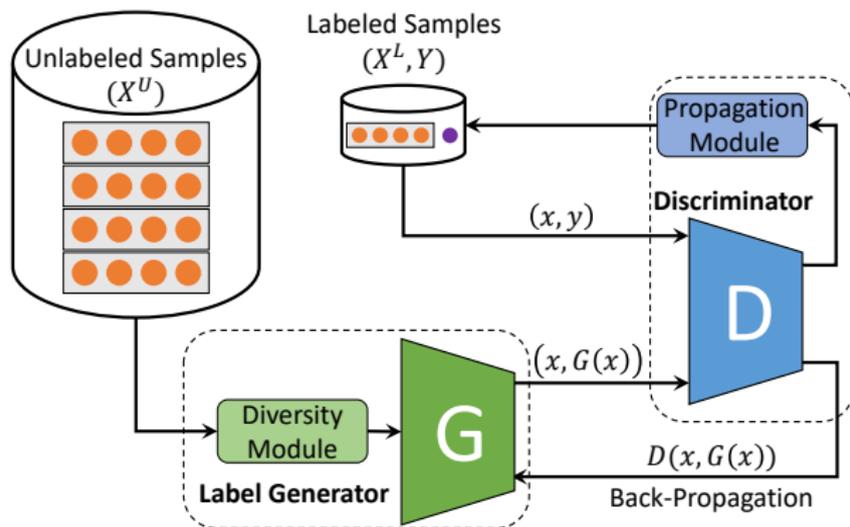
- \* The imbalanced class problem: ER tasks
- \* The overfitting problem: powerful models

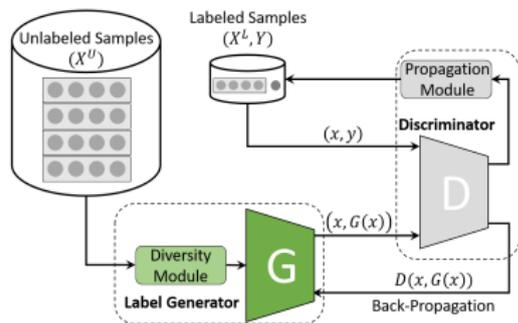
## Challenges

- \* The imbalanced class problem: ER tasks
- \* The overfitting problem: powerful models

## Solution

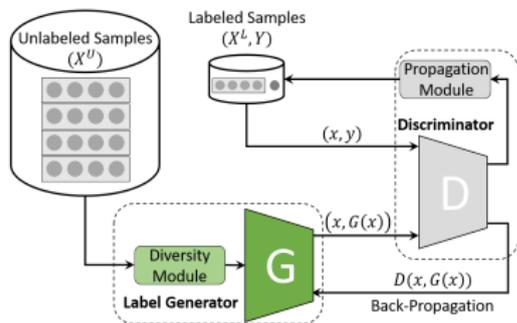
- Label generator  $G$ : only have access to unlabeled samples, consider diverse samples
- Discriminator  $D$ : provide feedback to train  $G$ , limited labels used with propagation





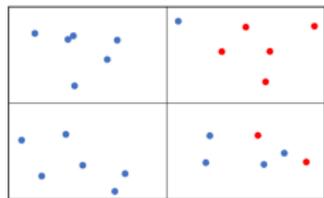
Generate pseudo labels for unlabeled samples

Learn a conditional distribution  
 $p_g(Y|X^U) \approx p(Y|X^U)$



Generate pseudo labels for unlabeled samples

Learn a conditional distribution  $p_g(Y|X^U) \approx p(Y|X^U)$



Diversity Sampling  
(4 groups)

A minibatch of  $m$  samples is selected from  $X^U$  according to the following objective function:

$$\text{maximize} \quad \|\mathbf{v}\|_{2,1} \quad \text{s.t.} \quad \sum_{i,j} v_i^j = m$$

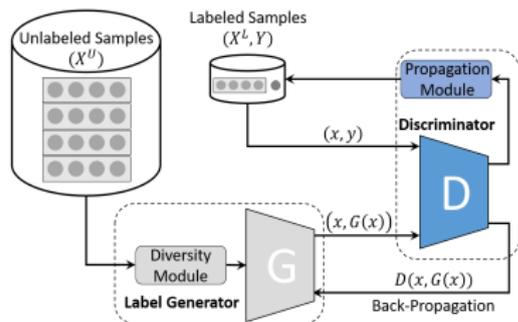
$G$  updates its parameters according to:

$$\mathcal{L}_G = \min_G \mathbb{E}_{x \sim p(x_i^U)} [\log(1 - D(x, G(x)))] \quad (1)$$

where:

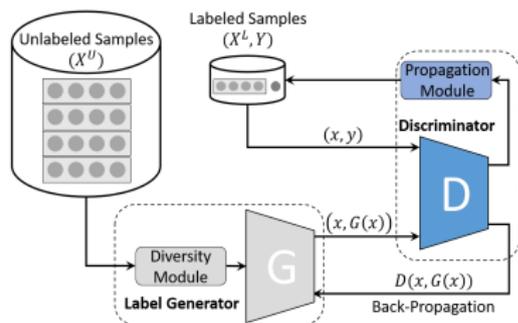
- $G(x_i)$  is the pseudo label of  $x_i$  generated by  $G$ ;
- $(x_i, G(x_i))$  is a pseudo labeled sample sent to the discriminator  $D$ ;
- $D(x, G(x))$  is the feedback from the discriminator  $D$ .

# Discriminator $D$



Distinguish samples with pseudo labels from samples with real labels

Learn a joint distribution  $p(X, Y)$



Distinguish samples with pseudo labels from samples with real labels

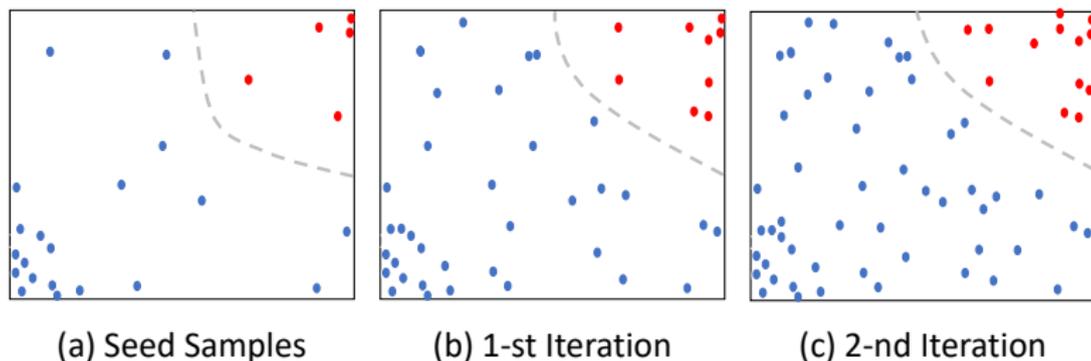
Learn a joint distribution  $p(X, Y)$

The objective function of  $D$  at the  $t$ -th iteration of propagation is:

$$\mathcal{L}_D = \max_D \mathbb{E}_{x \sim p(X_i^U)} \log[(1 - D(x, G(x)))] + \lambda \mathbb{E}_{(x, y) \sim (X^*, Y)^t} \log[D(x, y)] \quad (2)$$

where:

- $\lambda$  refers to a weighted term.
- $(X^*, Y)^t$  refers to the labeled samples in  $t$ -th iteration.



The *propagation module* selects a minibatch of  $|\Delta X^t|$  high-quality pseudo labeled samples for training  $D$ :

$$\operatorname{argmax}_{\Delta X^t \subseteq X^t} \sum_{x \in \Delta X^t} D(x, G(x))$$

- \* *Unsupervised*: **Two-Steps** and **Iterative Term-Entity Ranking and CliqueRank (ITER-CR)**.
- \* *Semi-supervised*: **Semi-supervised Boosted Classifier (SBC)**.
- \* *Fully supervised*: **Magellan** and **eXtreme Gradient boosting (XGboost)**.
- \* *Deep Learning based*: **DeepMatcher (DM)** and **Deep Transfer Active Learning (DTAL)**.
- \* Ablation Study: **ErGAN+WE**, **ErGAN-D**, **ErGAN-P**, and **ErNN**.

Method	Datasets			
	Cora	DBLP- ACM	DBLP- Scholar	NCVoter
2S	62.69	91.43	68.78	98.96
ITER-CR*	89.00	–	–	–
SBC	85.71	97.09	85.47	99.78
SVM	88.95	97.19	85.71	98.48
LR	80.25	95.56	83.84	99.37
XGBoost	91.34	97.20	86.63	<b>100</b>
ERGAN	<b>93.03</b>	<b>98.23</b>	<b>88.32</b>	<b>100</b>
DM	98.58	98.29	94.68	<b>100</b>
DTAL*	98.68 $\pm$ 0.26	98.45 $\pm$ 0.22	92.94 $\pm$ 0.47	–
ERGAN+WE	<b>98.72</b> $\pm$ 0.15	<b>98.51</b> $\pm$ 0.23	<b>94.73</b> $\pm$ 0.35	<b>100</b>

Datasets	Cora				DBLP-ACM			
	0.1%	1%	20%	60%	0.1%	1%	20%	60%
ERNN	84.46	90.67	91.43	92.78	88.05	95.68	98.20	98.22
ERGAN-D	79.87	85.14	91.27	92.97	0	93.30	97.16	98.21
ERGAN-P	85.18	90.76	91.42	<b>93.03</b>	92.67	95.96	98.21	<b>98.23</b>
ERGAN	<b>87.45</b>	<b>91.07</b>	<b>91.54</b>	<b>93.03</b>	<b>96.89</b>	<b>96.93</b>	<b>98.22</b>	<b>98.23</b>
Datasets	DBLP-Scholar				NCVoter			
	0.1%	1%	20%	60%	0.1%	1%	20%	60%
ERNN	82.76	83.17	86.71	87.73	99.39	<b>100</b>	<b>100</b>	<b>100</b>
ERGAN-D	0	78.85	83.43	88.29	0	99.58	<b>100</b>	<b>100</b>
ERGAN-P	83.43	85.34	86.55	<b>88.32</b>	99.39	99.79	<b>100</b>	<b>100</b>
ERGAN	<b>84.23</b>	<b>85.85</b>	<b>86.86</b>	<b>88.32</b>	<b>99.45</b>	<b>100</b>	<b>100</b>	<b>100</b>

In summary, we have proposed four approaches for ER:

- \* ASL: an active scheme learning approach
- \* Skyblocking: scheme skyline learning under different blocking criteria
- \* LST: A learning-based active learning framework
- \* ERGAN: a generative model with adversarial nets

Thank You!

Q & A

Email: [Jingyu.shao@anu.edu.au](mailto:Jingyu.shao@anu.edu.au)