

# Semi-Supervised Feature Selection for Graph Classification

Xiangnan Kong, Philip S. Yu



Department of Computer Science  
University of Illinois at Chicago



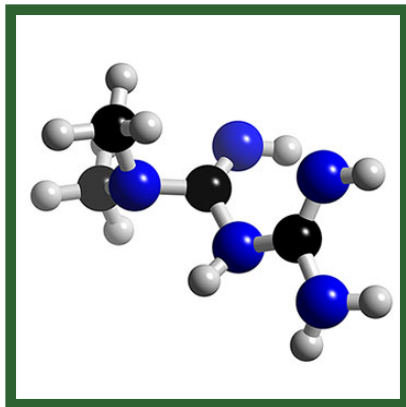
# Graph Classification

- why should we care?

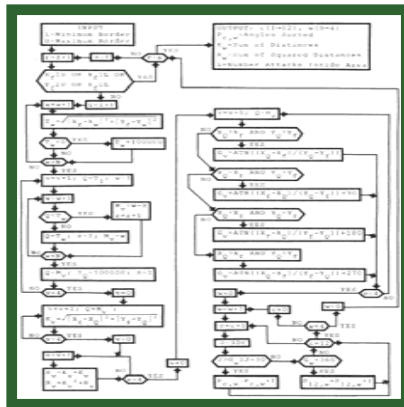
❑ Conventional data mining and machine learning approaches assume data are represented as feature vectors. E.g.  $(x_1, x_2, \dots, x_d) - y$

❑ In real apps, data are not directly represented as feature vectors, but **graphs** with complex structures. E.g.  $G(V, E, I) - y$

Chemical Compounds



Program Flows



XML Docs



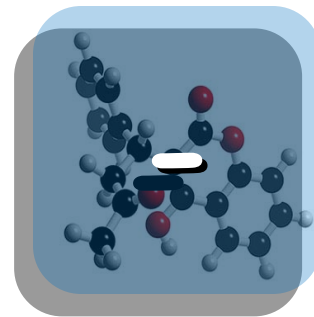
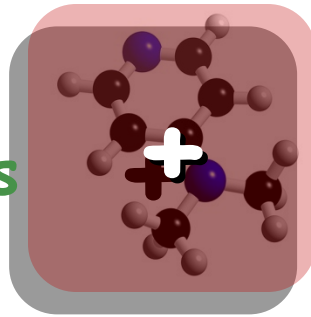
# Example: Graph Classification

- Drug activity prediction problem

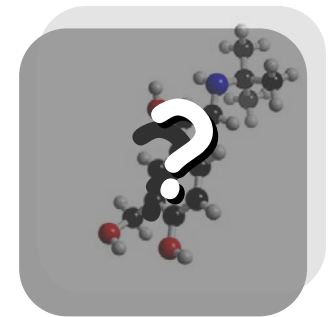
- Given a set of chemical compounds labeled with activities

- Predict the activities of testing molecules

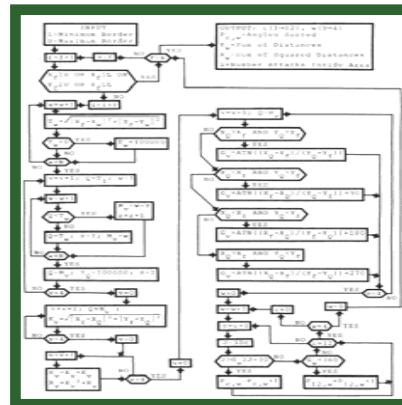
Training Graphs



Testing Graph



Program Flows



XML Docs



# Subgraph-based Graph Classification

## Subgraph Patterns

$g_1$

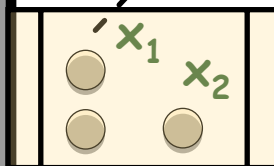
$g_2$

$g_3$



How to mine a set of **subgraph** patterns in order to **effectively** perform graph classification?

Classifier



$x_2$

Feature Vectors

Graph Objects  $\longrightarrow$  Feature Vectors  $\longrightarrow$  Classifiers

# Conventional Methods

## - Two Components

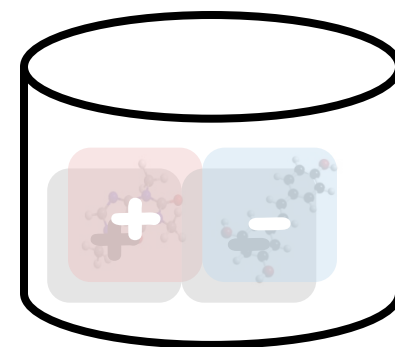
Two Components:

1. Evaluation (effective)

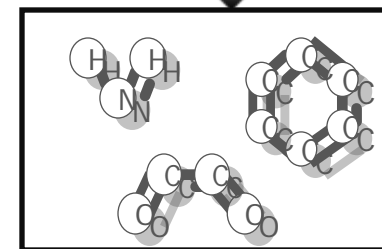
whether a subgraph feature is relevant to graph classification?

2. Search space pruning (efficient)

how to avoid enumerating all subgraph features?



Labeled Graphs



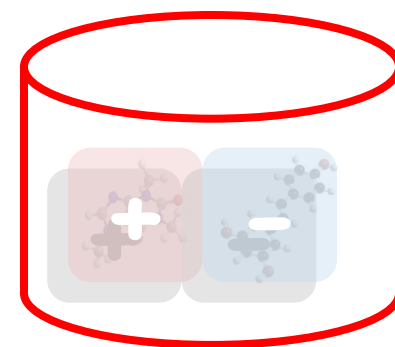
Discriminative Subgraphs

# One Problem

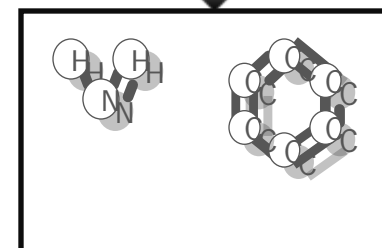
- Supervised Settings
  - Require a large set of labeled training graphs
- However...



Labeling  
a graph  
is hard !



Labeled Graphs

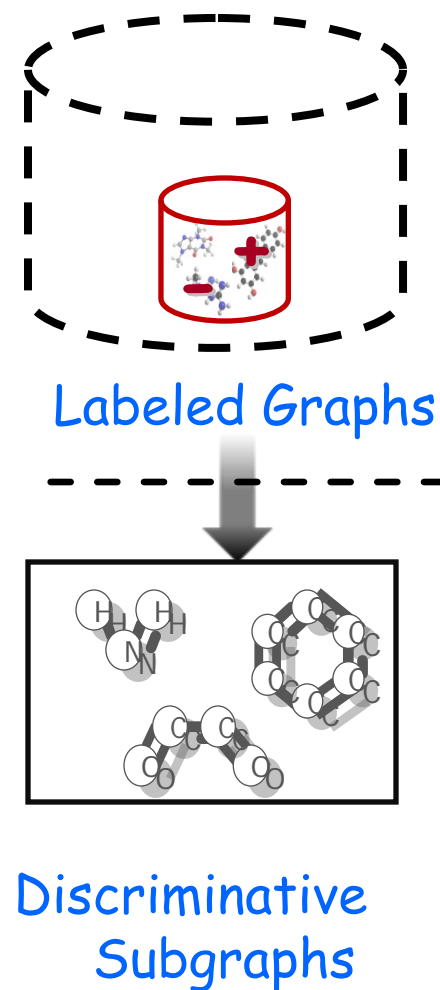


Discriminative  
Subgraphs

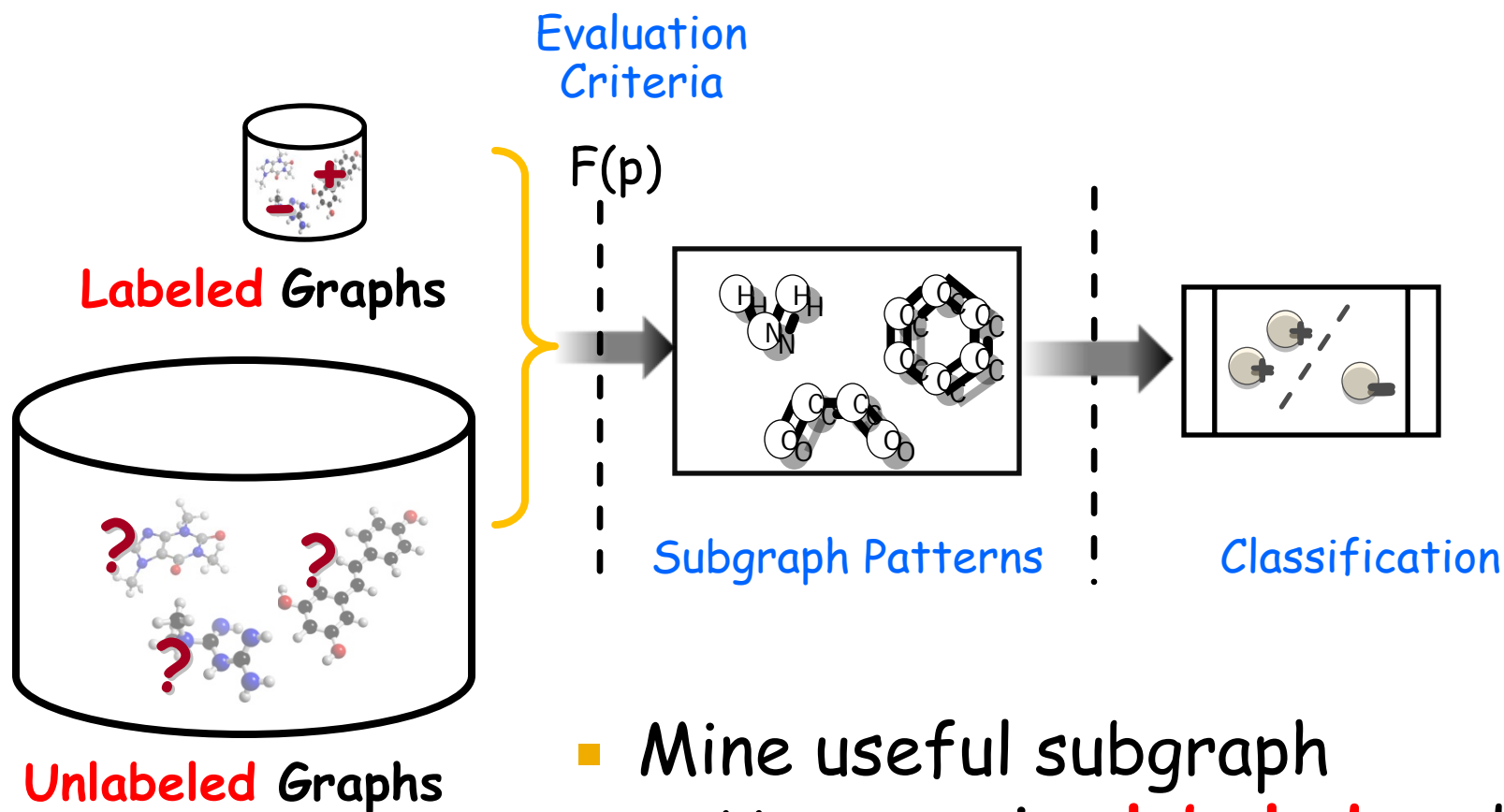
# Lack of labels -> problems

## Supervised Methods:

1. Evaluation *effective?*  
require large amount of label information
2. Search space pruning *efficient?*  
pruning performances rely on large amount of label information




# Semi-Supervised Feature Selection for Graph Classification



- Mine useful subgraph patterns using **labeled** and **unlabeled** graphs



# Two Key Questions to Address

- Evaluation: How to evaluate a set of subgraph features with both **labeled** and **unlabeled** graphs? **(effective)** 
- Search Space Pruning: How to prune the subgraph search space using both **labeled** and **unlabeled** graphs? **(efficient)**

# What is a good feature?

## Cannot-Link

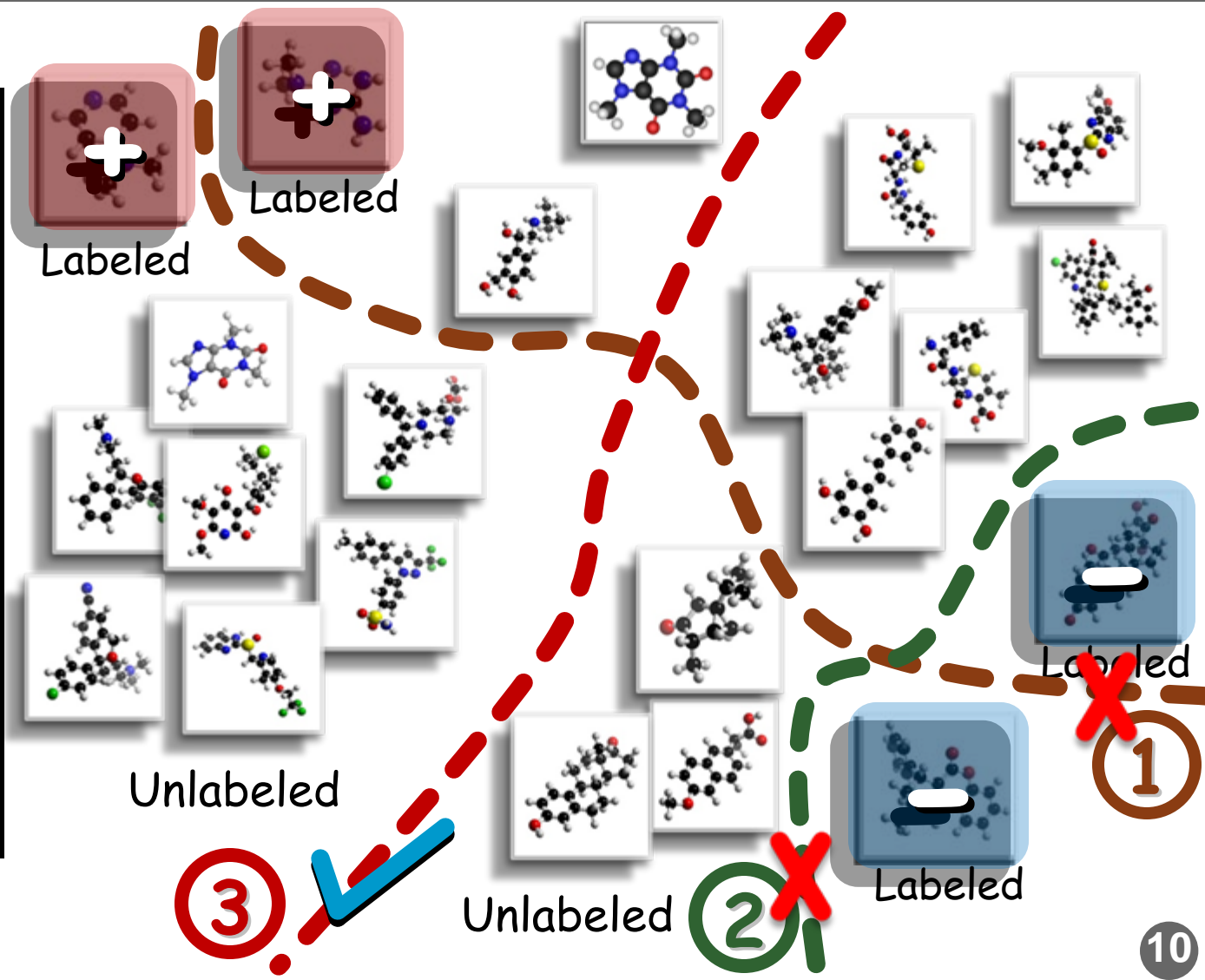
Graphs in different classes should be far away

## Must-Link

Graphs in the same class should be close

## Separability

Unlabeled graphs are able to be separated from each other



# Optimization

$$\mathcal{T}^* = \operatorname{argmax}_{\mathcal{T} \subseteq \mathcal{S}} J(\mathcal{T}) \quad \text{s.t.} \quad |\mathcal{T}| \leq t$$

Evaluation Function:

$$\begin{aligned} & \frac{\alpha}{2|\mathcal{C}|} \sum_{y_i y_j = -1} (D_{\mathcal{T}} \mathbf{x}_i - D_{\mathcal{T}} \mathbf{x}_j)^2 \\ & - \frac{\beta}{2|\mathcal{M}|} \sum_{y_i y_j = 1} (D_{\mathcal{T}} \mathbf{x}_i - D_{\mathcal{T}} \mathbf{x}_j)^2 \\ & + \frac{1}{2|\mathcal{D}_u|^2} \sum_{G_i, G_j \in \mathcal{D}_u} (D_{\mathcal{T}} \mathbf{x}_i - D_{\mathcal{T}} \mathbf{x}_j)^2 \end{aligned}$$

## Cannot-Link

Graphs in different classes should be far away

## Must-Link

Graphs in the same class should be close

## Separability

Unlabeled graphs are able to be separated from each other

# Evaluation: gSemi Criterion

In matrix form:

$$J(\mathcal{T}) = \frac{1}{2} \sum_{i,j} (D_{\mathcal{T}} \mathbf{x}_i - D_{\mathcal{T}} \mathbf{x}_j)^2 W_{ij}$$

$$W_{ij} = \begin{cases} \frac{\alpha}{|C|} & \text{if } y_i y_j = -1 \\ -\frac{\beta}{|\mathcal{M}|} & \text{if } y_i y_j = 1 \\ \frac{1}{|\mathcal{D}_u|^2} & \text{if } G_i, G_j \in \mathcal{D}_u \\ 0 & \text{otherwise} \end{cases}$$

$$L = D - W$$

$$= \text{tr}(D_{\mathcal{T}}^{\top} X (D - W) X^{\top} D_{\mathcal{T}})$$

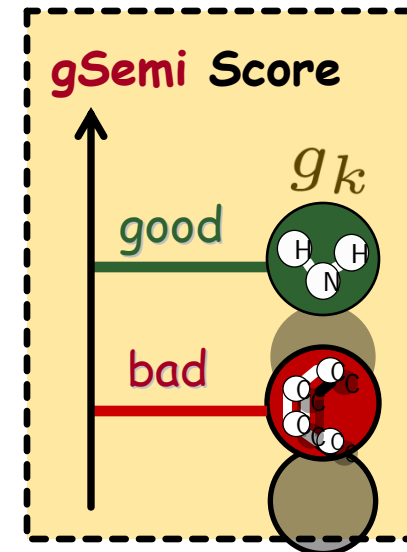
$$= \text{tr}(D_{\mathcal{T}}^{\top} X L X^{\top} D_{\mathcal{T}})$$

$$= \sum_{g_k \in \mathcal{T}} (\mathbf{f}_k^{\top} L \mathbf{f}_k) \quad (\text{the sum over all selected features})$$

■ **gSemi** Score:

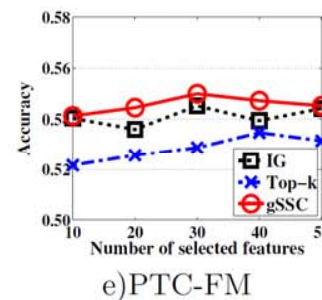
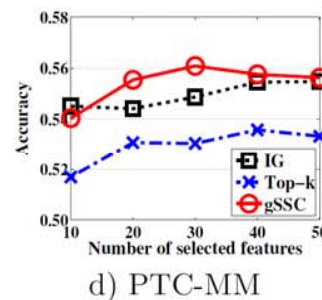
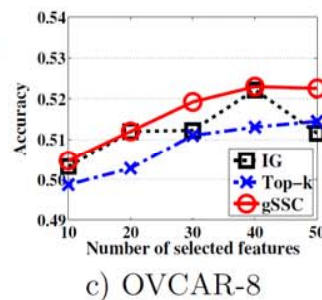
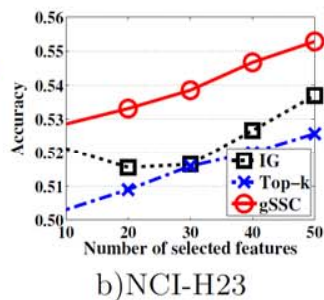
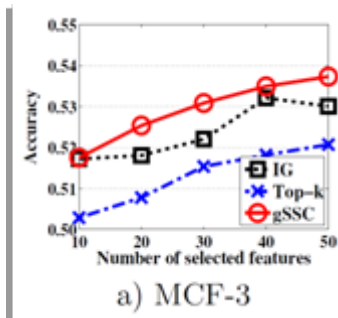
$$h(g_k, L) = \mathbf{f}_k^{\top} L \mathbf{f}_k$$

$\mathbf{f}_k \in \{0, 1\}^n$  represents the  $k$ -th subgraph feature

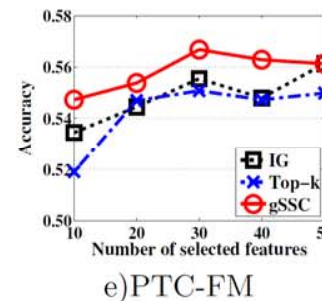
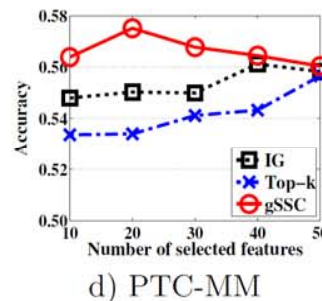
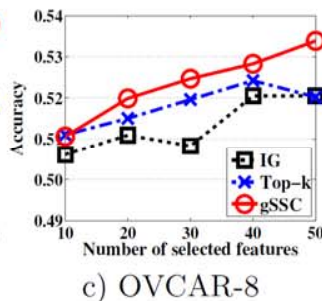
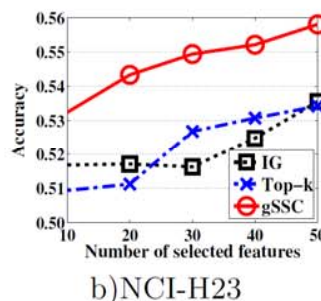
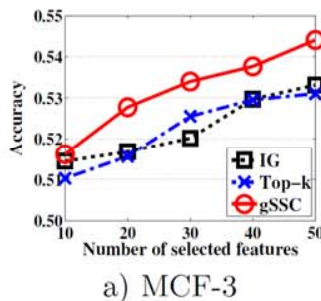


# Experiment Results

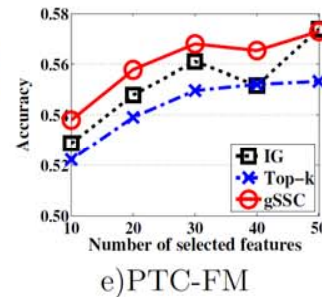
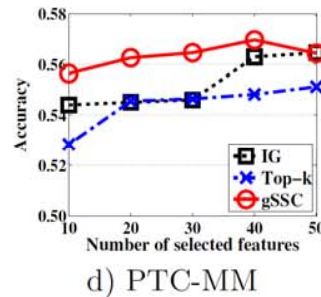
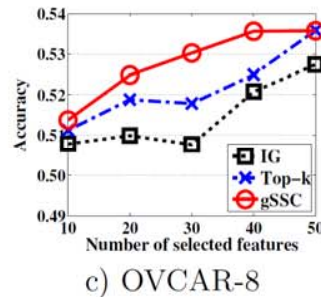
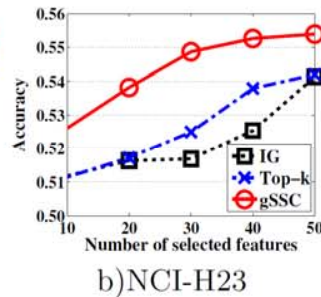
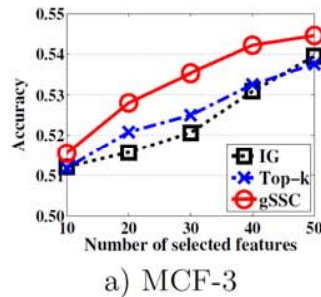
#labeled  
Graphs  
=30



#labeled  
Graphs  
=50



#labeled  
Graphs  
=70



MCF-3  
dataset

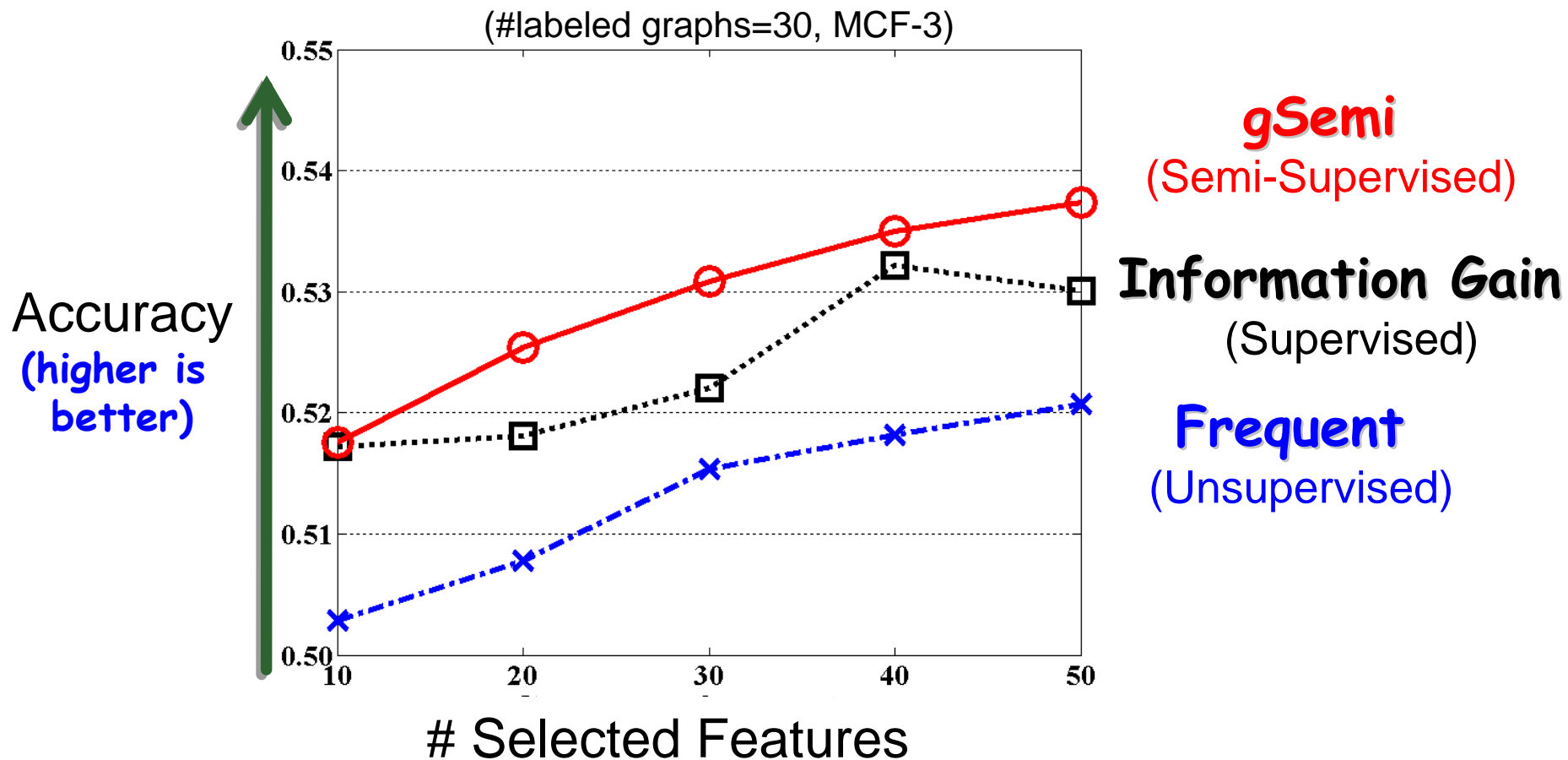
NCI-H23  
dataset

OVCAR-8  
dataset

PTC-MM  
dataset

PTC-FM  
dataset

# Experiment Results



- Our approach performed **best** at NCI and PTC datasets



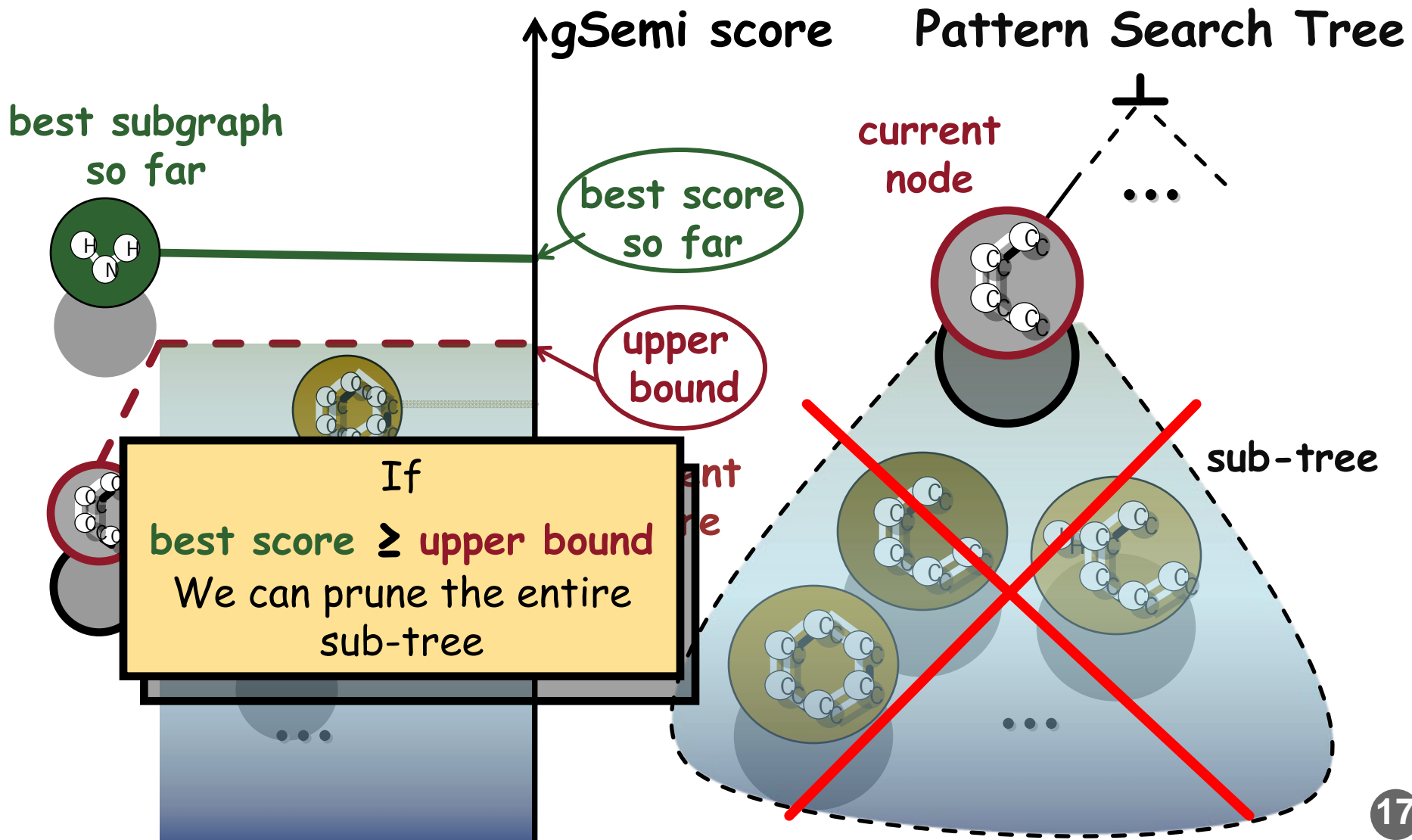
# Two Key Questions to Address

- How to evaluate a set of subgraph features with both labeled and unlabeled graphs? (effective)
- How to prune the subgraph search space using both labeled and unlabeled graphs? (efficient) ←

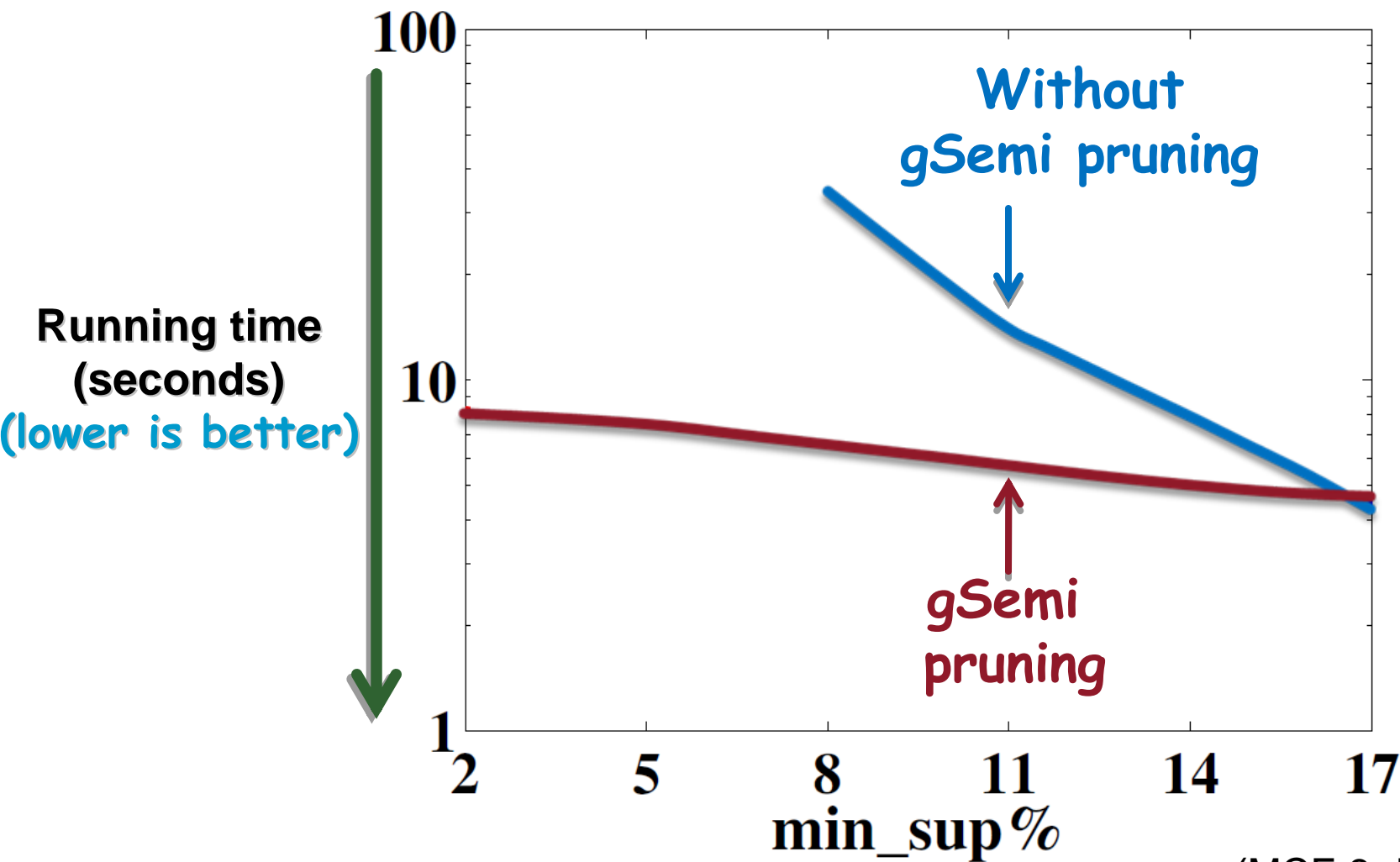




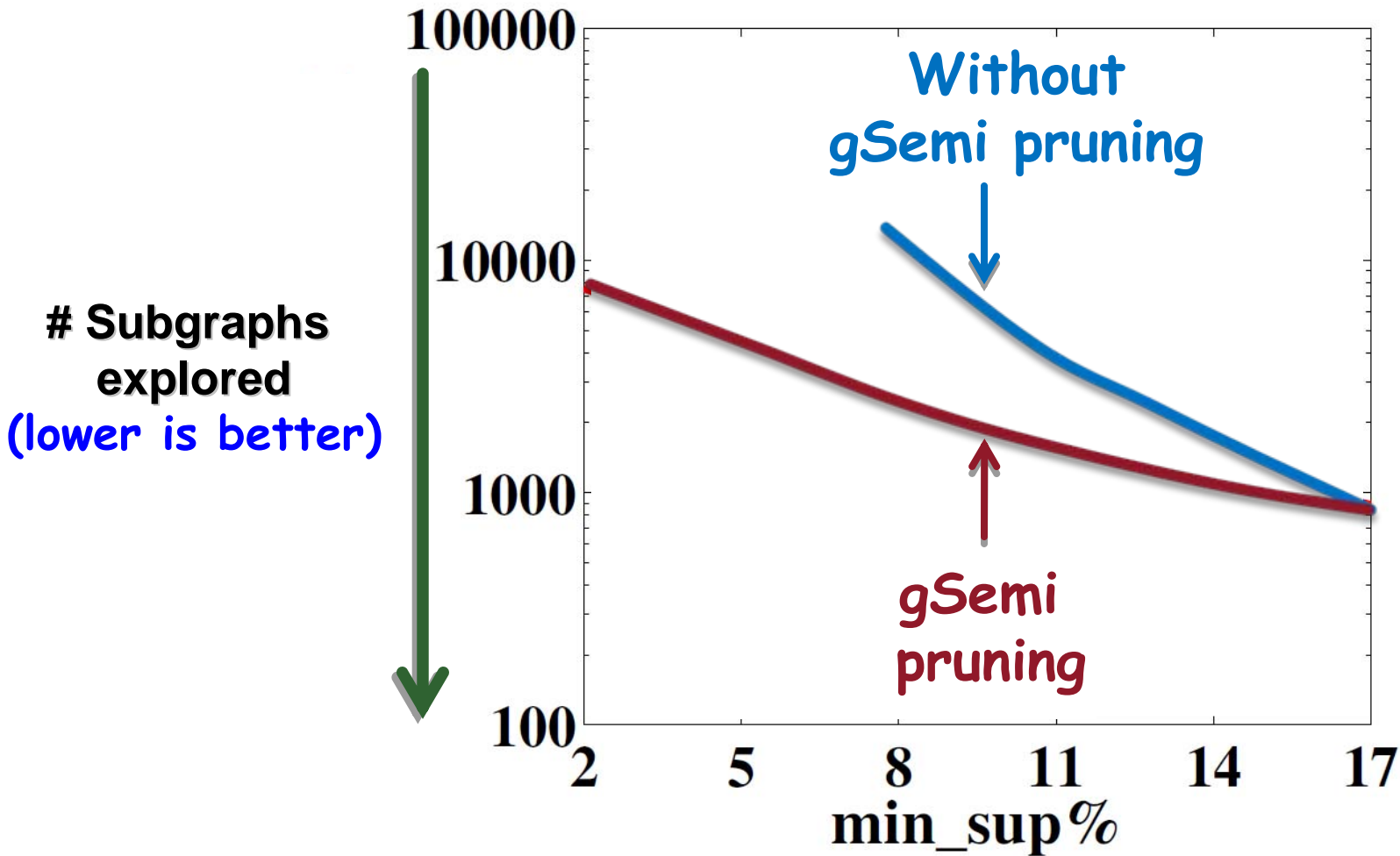
# Pruning Principle



# Pruning Results



# Pruning Results



(MCF-3 dataset)

# Conclusions

- **Semi-Supervised** Feature Selection for Graph Classification
  - Evaluating subgraph features using both labeled and unlabeled graphs (*effective*)
  - Branch&bound pruning the search space using labeled and unlabeled graphs (*efficient*)



Thank you!

