

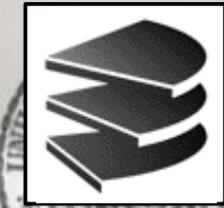
# Combining Near-Optimal Feature Selection with *gSpan*

Marisa Thoma<sup>1</sup>

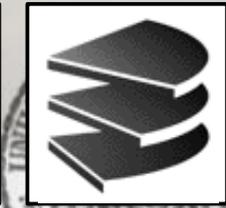
*joint work with*

Karsten Borgwardt<sup>2</sup>, Xifeng Yan<sup>3</sup>, Hong Cheng<sup>6</sup>, Arthur Gretton<sup>4</sup>,  
Le Song<sup>5</sup>, Alex Smola<sup>5</sup>, Jiawei Han<sup>6</sup>, Philip Yu<sup>2</sup>, Hans-Peter Kriegel<sup>1</sup>

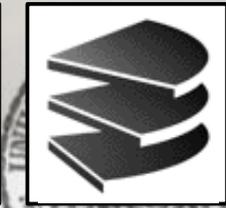
<sup>1</sup> LMU Munich, <sup>2</sup> University of Cambridge, <sup>3</sup> IBM Watson Research  
Center New York, <sup>4</sup> MPI Tübingen, <sup>5</sup> NICTA Canberra, <sup>6</sup> UIUC



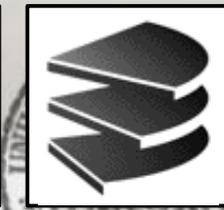
- Why feature selection?
- Submodular selection criteria
- CORK, a novel feature selection criterion
- Inclusion of CORK into *gSpan*
- Experimental validation
- Summary and outlook



- Two-class problem:
  - Collection of graphs
  - Predict class label from graph topology
- Current solutions:
  - Pattern-based learning (subgraphs, paths, circles)
  - Graph kernel approaches (random walk, pattern based approaches)
  - Nested approaches (feature generation adapted to the dataset)



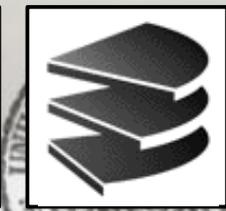
- Frequent Subgraphs
  - Can be efficiently enumerated (*gSpan*, FSG, MoSS)
- They also contain
  - Insignificant features
  - Redundant features, already covered by other substructures
- Potentially exponentially many



- Identify the most discriminative subgraphs
- Set  $\mathcal{J}_{opt} \subseteq \mathcal{F}$  (total feature set), s.t.

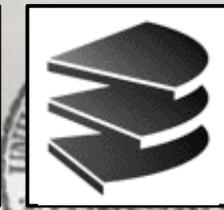
$$\mathcal{J}_{opt} = \operatorname{argmax}_{\mathcal{S} \subseteq \mathcal{F}} f(\mathcal{S})$$

according to information criterion  $f$ .



## Approaches

- Complete subset enumeration and test (Wrapper)
- Ranker selection (assuming feature independence)
- Greedy approaches
  - Backward Elimination vs. Forward Selection
- Nested approaches
  - Feature selection performed in combination with feature generation

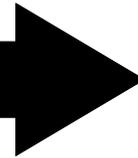


# Proof by Krause and Guestrin

A note on the Budgeted Maximization of Submodular Functions. CMU 2005

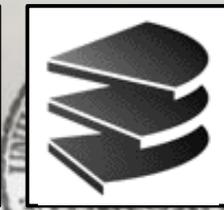
Submodular decision function

**Greedy Forward Selection**



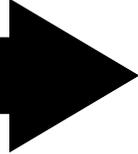
Guarantee of

$$f(\hat{\mathcal{J}}) \geq \left(1 - \frac{1}{e}\right) f(\mathcal{J}_{opt})$$



## Proof by Krause and Guestrin A note on the Budgeted Maximization of Submodular Functions. CMU 2005

Submodular decision function

**Greedy Forward Selection** 

Guarantee of

$$f(\widehat{\mathcal{F}}) \geq \left(1 - \frac{1}{e}\right) f(\mathcal{F}_{opt})$$

set function  $f: \mathcal{F} \rightarrow \mathbb{R}$

$f$  is submodular  $\Leftrightarrow$

$S \subset \mathcal{T} \subseteq \mathcal{F}$ ,  $s \in \mathcal{F}$  :

$$f(S \cup \{s\}) - f(S) \geq f(\mathcal{T} \cup \{s\}) - f(\mathcal{T})$$

## Proof by Krause and Guestrin

A note on the Budgeted Maximization of Submodular Functions. CMU 2005

Submodular decision function

**Greedy Forward Selection** →

Guarantee of

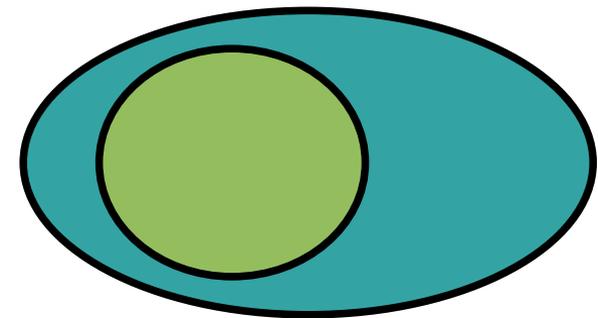
$$f(\widehat{\mathcal{F}}) \geq \left(1 - \frac{1}{e}\right) f(\mathcal{F}_{opt})$$

set function  $f: \mathcal{F} \rightarrow \mathbb{R}$

$f$  is submodular  $\Leftrightarrow$

$S \subset \mathcal{T} \subseteq \mathcal{F}, s \in \mathcal{F} :$

$$f(S \cup \{s\}) - f(S) \geq f(\mathcal{T} \cup \{s\}) - f(\mathcal{T})$$



$$f(S) = \text{area}(S)$$

## Proof by Krause and Guestrin

A note on the Budgeted Maximization of Submodular Functions. CMU 2005

Submodular decision function

**Greedy Forward Selection** →

Guarantee of

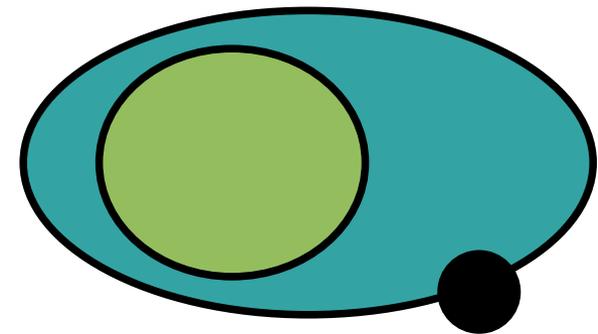
$$f(\widehat{\mathcal{F}}) \geq \left(1 - \frac{1}{e}\right) f(\mathcal{F}_{opt})$$

set function  $f: \mathcal{F} \rightarrow \mathbb{R}$

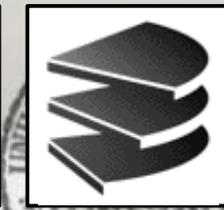
$f$  is submodular  $\Leftrightarrow$

$S \subset \mathcal{T} \subseteq \mathcal{F}, s \in \mathcal{F} :$

$$f(S \cup \{s\}) - f(S) \geq f(\mathcal{T} \cup \{s\}) - f(\mathcal{T})$$



$$f(S) = \text{area}(S)$$



## *Correspondence:*

A pair of instances  $i, j \in \text{Dataset } D$  with  $\text{class}(i) \neq \text{class}(j)$  is called a correspondence in a feature set  $\mathcal{S}$ , if  $i$  and  $j$  have the same values for  $\mathcal{S}$ .

## Correspondence-based Quality Criterion (*CORK*):

$$q(\mathcal{S}) = (-1) * \text{“number of correspondences in } \mathcal{S}\text{”}$$

From now:

use CORK for binary feature values  $\{0, 1\}$ .



Correspondence-based Quality Criterion (*CORK*):

$$q(S) = (-1) * \text{“number of correspondences in } S\text{”}$$

$A_{S_0}$ : # “feature  $S = 0$  in class  $A$ “ (# mis-matches)

$A_{S_1}$ : # “feature  $S = 1$  in class  $A$ “ (# matches)

$$q(\{S\}) = -(A_{S_0} \cdot B_{S_0} + A_{S_1} \cdot B_{S_1})$$

2-class dataset

		<b>A</b>			<b>B</b>		
		$a_1$	$a_2$	$a_3$	$b_1$	$b_2$	$b_3$
Features	$S_1$	1	1	0	1	0	0
	$S_2$	1	0	0	1	0	0

Correspondence-based Quality Criterion (*CORK*):

$$q(S) = (-1) * \text{“number of correspondences in } S\text{”}$$

$A_{S_0}$ : # “feature  $S = 0$  in class  $A$ “ (# mis-matches)

$A_{S_1}$ : # “feature  $S = 1$  in class  $A$ “ (# matches)

$$q(\{S\}) = -(A_{S_0} \cdot B_{S_0} + A_{S_1} \cdot B_{S_1})$$

$$q(\{s_1\}) = -(1 \cdot 2 + 2 \cdot 1) = -4$$

2-class dataset

	<b>A</b>			<b>B</b>		
	$a_1$	$a_2$	$a_3$	$b_1$	$b_2$	$b_3$
Features $s_1$	1	1	0	1	0	0
Features $s_2$	1	0	0	1	0	0

Correspondence-based Quality Criterion (*CORK*):

$$q(S) = (-1) * \text{“number of correspondences in } S\text{”}$$

$A_{S_0}$ : # “feature  $S = 0$  in class  $A$ “ (# mis-matches)

$A_{S_1}$ : # “feature  $S = 1$  in class  $A$ “ (# matches)

$$q(\{S\}) = -(A_{S_0} \cdot B_{S_0} + A_{S_1} \cdot B_{S_1})$$

$$q(\{s_1\}) = -(1 \cdot 2 + 2 \cdot 1) = -4$$

$$q(\{s_2\}) = -(2 \cdot 2 + 1 \cdot 1) = -5$$

2-class dataset

	<b>A</b>			<b>B</b>		
	$a_1$	$a_2$	$a_3$	$b_1$	$b_2$	$b_3$
Features $s_1$	1	1	0	1	0	0
Features $s_2$	1	0	0	1	0	0

$$q(\{S\}) = -(A_{S_0} \cdot B_{S_0} + A_{S_1} \cdot B_{S_1})$$

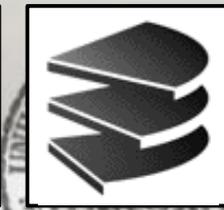
$$q(\{S, T\}) = -\left(\sum_{i,j=0}^1 A_{S_i, T_j} \cdot B_{S_i, T_j}\right)$$

	A			B		
	$a_1$	$a_2$	$a_3$	$b_1$	$b_2$	$b_3$
$s_1$	1	1	0	1	0	0
$s_2$	1	0	0	1	0	0

Histogram over *equivalence classes* for the possible feature combinations

For  $M$  features:  $2^M$  possible equivalence classes

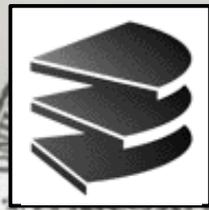
$$q(\{s_1, s_2\}) = -\left(\begin{matrix} 0,0 & 0,1 & 1,0 & 1,1 \\ 1 \cdot 2 + 0 \cdot 0 + 1 \cdot 0 + 1 \cdot 1 \end{matrix}\right) = -3$$



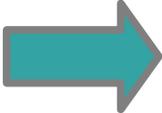
$\mathcal{S} \subset \mathcal{T} \subseteq \mathcal{F}, s \in \mathcal{F}:$

Improvement of  $\mathcal{T}$   must also improve  $\mathcal{S}$

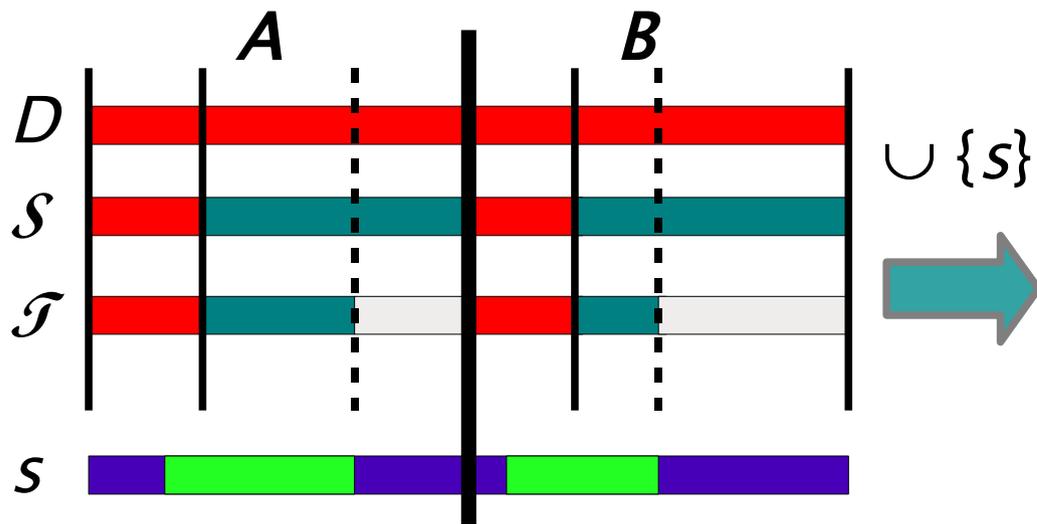
*Improvement*  $\triangleq$  remove a correspondence



$\mathcal{S} \subset \mathcal{T} \subseteq \mathcal{F}, s \in \mathcal{F}$ :

Improvement of  $\mathcal{T}$   must also improve  $\mathcal{S}$

*Improvement*  $\triangleq$  remove a correspondence

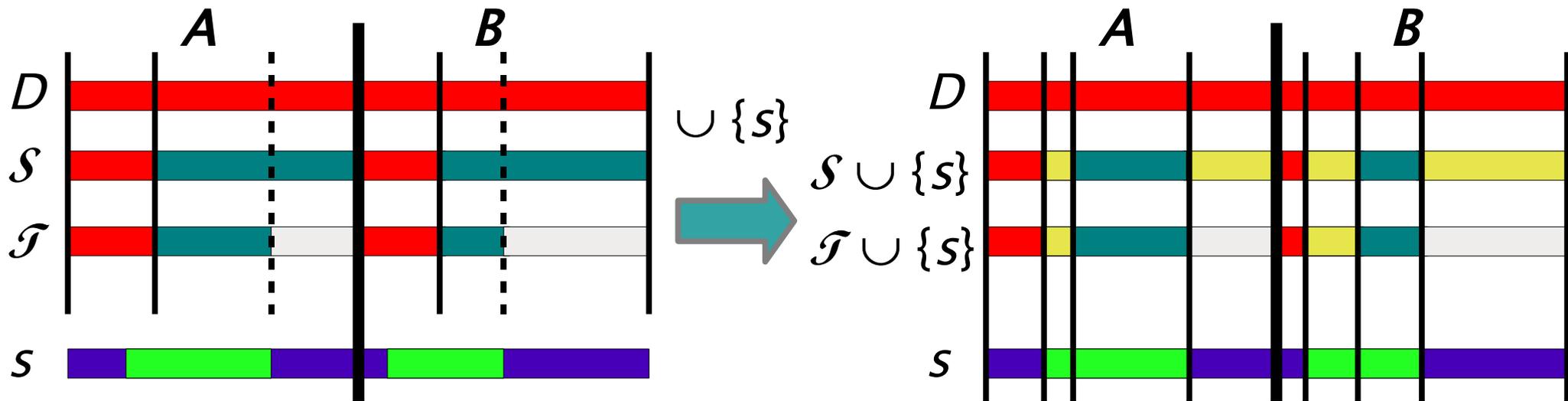


$$f(\mathcal{S} \cup \{s\}) - f(\mathcal{S}) \geq f(\mathcal{T} \cup \{s\}) - f(\mathcal{T})$$

$\mathcal{S} \subset \mathcal{T} \subseteq \mathcal{F}, s \in \mathcal{F}$ :

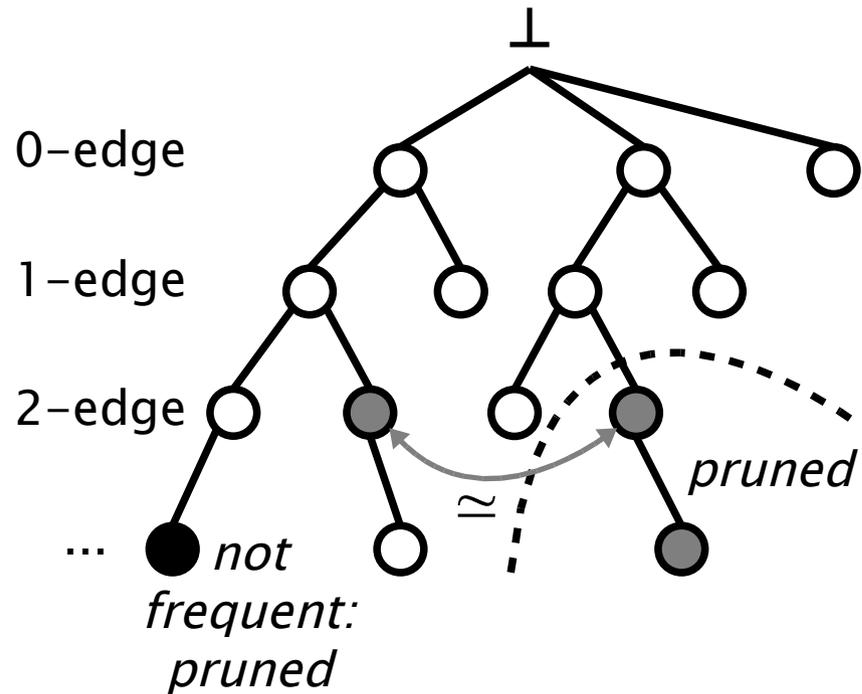
Improvement of  $\mathcal{T}$  must also improve  $\mathcal{S}$

*Improvement*  $\triangleq$  remove a correspondence

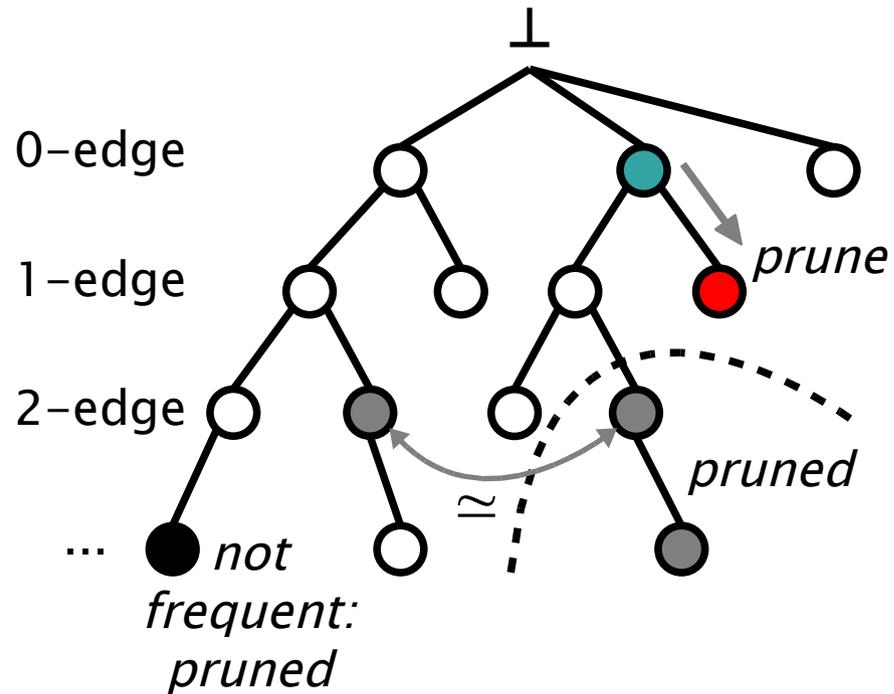


$$f(\mathcal{S} \cup \{s\}) - f(\mathcal{S}) \geq f(\mathcal{T} \cup \{s\}) - f(\mathcal{T})$$

Yan, X. & Han, J.: *gSpan*: Graph-based sub-structure pattern mining. *Proc. ICDM 2002*



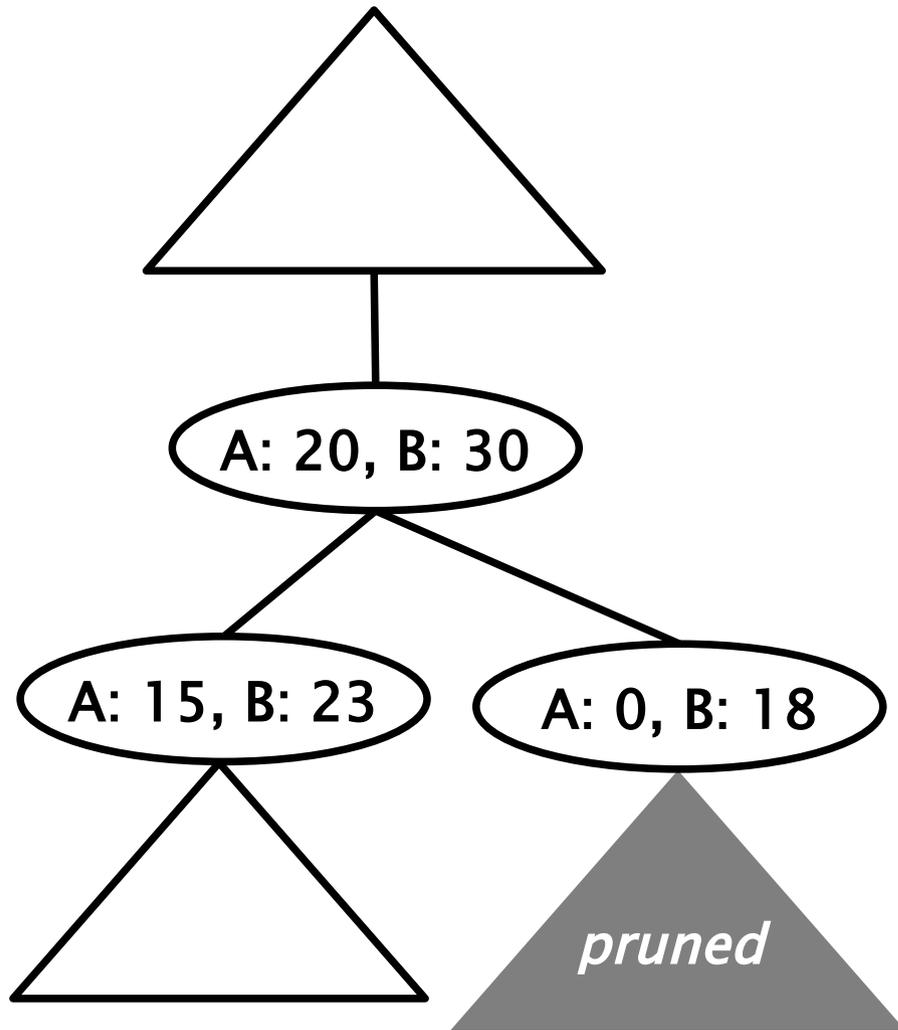
- DFS Code Tree
- Efficient branch-and-bound
  - Frequency bound
  - Minimality bound of DFS Code
- Worst-case runtime: exponential
- For any parent-child relationship in the search tree:  $parent \sqsubseteq child$



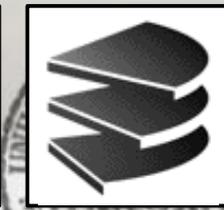
- Let *parent*  $S \sqsubseteq$  *child*  $T$  be frequent subgraphs
- From  $S$  to  $T$ , we can only loose matching subgraph embeddings but never gain additional matches.  $(A_{S_1} \geq A_{T_1} \wedge B_{S_1} \geq B_{T_1})$

CORK:  $q(\{S\}) = -(A_{S_0} \cdot B_{S_0} + A_{S_1} \cdot B_{S_1})$

- The maximal value of  $q(\{T\})$  is reached, when either all matches in  $A$  or all matches in  $B$  are lost.



- Let *parent*  $S \sqsubseteq$  *child*  $T$  be frequent subgraphs
  - From  $S$  to  $T$ , we can only loose matching subgraph embeddings but never gain additional matches.  $(A_{S_1} \geq A_{T_1} \wedge B_{S_1} \geq B_{T_1})$
- CORK:  $q(\{S\}) = -(A_{S_0} \cdot B_{S_0} + A_{S_1} \cdot B_{S_1})$
- The maximal value of  $q(\{T\})$  is reached, when either all matches in  $A$  or all matches in  $B$  are lost.

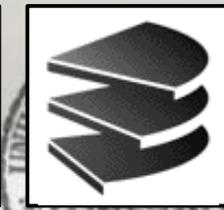


For  $S \sqsubseteq T$  we can derive:

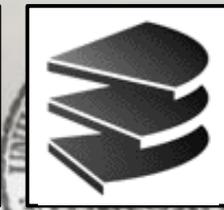
$$q(\{T\}) \leq q(\{S\}) + \max \begin{cases} (A_{S_1} - A_{S_0}) \cdot B_{S_1} \\ A_{S_1} \cdot (B_{S_1} - B_{S_0}) \\ 0 \end{cases}$$

Enables earlier pruning:

- No supergraph of  $T$  can exceed this upper bound.
- Thus, if we have seen a better subgraph, this branch can be pruned.



- CORK bound is extendible to Feature Sets
- For the enumeration of significant features:
  1.  $\mathcal{F} = \emptyset$
  2. *gSpan*  $\Rightarrow$  best subgraph  $S$  according to  $q(\{S\} \cup \mathcal{F})$
  3. If  $q(\{S\} \cup \mathcal{F}) > q(\mathcal{F})$  :
  4.  $\mathcal{F} = \mathcal{F} \cup S$
  5. GoTo 2.
  6. return  $\mathcal{F}$

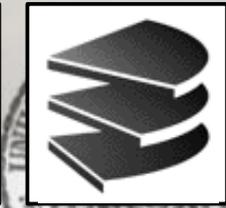


- NCI Dataset:
  - 6 chemical structure sets of variable size
  - Mapped to label “effective against cancer” (Y / N)
- Comparison to 2 feature ranking methods:
  - Sequential Cover via confidence score
  - Pearson Correlation
- Comparison to wrapper approach LAR–LASSO

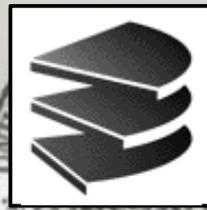
Tsuda, K. : Entire regularization paths for graph data. *ICML* 2007

- 10 repetitions of 10-fold cross validation
- on *gSpan* enumeration with freq. threshold 10%
- validated via a linear SVM

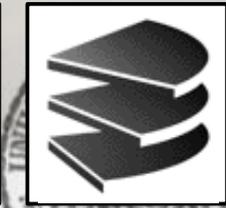
Dataset	#Features	Filter						Wrapper	
		SC		PC		CORK		LAR	
		Acc.	Std.	Acc.	Std.	Acc.	Std.	Acc.	Std.
NCI1	57	66.98	2.31	65.43	3.82	70.98	2.31	73.08	2.06
NCI33	53	66.50	2.57	64.15	3.46	70.08	2.76	72.81	2.51
NCI41	49	70.20	3.23	65.37	4.27	70.38	2.72	72.39	2.58
NCI47	56	67.04	2.35	67.00	3.45	71.42	2.22	72.62	2.07
NCI81	64	69.04	2.17	64.27	5.01	70.76	2.21	72.58	1.88



- Improve runtime
  - Save minimal DFS Codes
  - Sharpen the bound for later iterations
  - Combine with other bounding criteria
    - Subgraph mining procedure can be applied without giving a frequency threshold
- Exploit tree structure for decision tree learning



Thank you.



**Input:** Set of features  $\mathcal{F}$ , FS criterion  $f$

$\mathcal{J} := \emptyset$

$\mathcal{S} := \mathcal{F}$

$s = \mathbf{argmax} \{f(\{s^*\}) : s^* \in \mathcal{S}\}$

**while**  $f(\mathcal{J}) < f(\mathcal{J} \cup \{s\})$  **do**

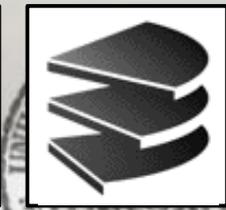
$\mathcal{J} := \mathcal{J} \cup \{s\}$

$\mathcal{S} := \mathcal{S} \setminus \{s\}$

$s = \mathbf{argmax} \{f(\mathcal{J} \cup \{s^*\}) : s^* \in \mathcal{S}\}$

**end**

**Output:** selected feature set  $\mathcal{J}$



$$S \sqsubseteq T, (A_{S_1} \geq A_{T_1} \wedge B_{S_1} \geq B_{T_1}) \quad q(\{S\}) = -(A_{S_0} \cdot B_{S_0} + A_{S_1} \cdot B_{S_1})$$

- $T$  can loose all matches in  $B$ :

$$q(\{T\}) \leq -(A_{S_0} \cdot (B_{S_0} + B_{S_1}) + A_{S_1} \cdot 0) = -A_{S_0} \cdot (B_{S_0} + B_{S_1})$$

- $T$  can loose all matches in  $A$ :

$$q(\{T\}) \leq -((A_{S_0} + A_{S_1}) \cdot B_{S_0} + 0 \cdot B_{S_1}) = -(A_{S_0} + A_{S_1}) \cdot B_{S_0}$$

- The maximal CORK value of  $T$  is thus

$$q(\{T\}) \leq \max \begin{Bmatrix} -A_{S_0} \cdot |B| \\ -|A| \cdot B_{S_0} \\ q(\{S\}) \end{Bmatrix} = q(\{S\}) + \max \begin{Bmatrix} (A_{S_1} - A_{S_0}) \cdot B_{S_1} \\ A_{S_1} \cdot (B_{S_1} - B_{S_0}) \\ 0 \end{Bmatrix}$$