# fast effective molecular feature mining by local optimization

Albrecht Zimmermann, Björn Bringmann, Ulrich Rückert

# fast effective molecular feature mining by local optimization

Albrecht Zimmermann, Björn Bringmann, Ulrich Rückert



#### simple approach: **QSA** dump a load of patterns into svm



You are here

(110...10)



cancer



svm



#### simple approach: **CSar** dump a load of patterns into svm You are here (110...10) better pattern set mining with cancer a) post-processing b) iterative mining svm

model

### a tale of two paradigms

#### **Supervised**

pattern set to distinguishes classes

correspondences, ig, ...

Bringmann & Zimmermann '05 Cheng et al. '08 Thoma et al. '09

#### Unsupervised

pattern set to induce fine, balanced partition

joint entropy

Knobbe & Ho '06 Bringmann & Zimmermann '07

### a tale of two paradigms

#### Supervised

pattern set to distinguishes classes

correspondences, ig, ...

Bringmann & Zimmermann '05 Cheng et al. '08 Thoma et al. '09

Trade-off

#### Unsupervised

pattern set to induce fine, balanced partition

joint entropy

Knobbe & Ho '06 Bringmann & Zimmermann '07

Trade-off

#### **Combined Scores**

Cheng et al. '07 Rückert & Kramer '07

class-correlated dispersion score

+++ 888 +++ 888 +++

#### **Supervised**













which paradigm is justified?

#### some empirical exploration

algorithmpropertiesbaseline500 most general frequent graphspicker\*cond. pattern prob.post processingrückert '07class-correlated disp-scoreiterativefcorkcorrespondence-basediterativedtminformation gainiterative

relation to AUC?

#### some empirical exploration

algorithmpropertiesbaseline500 most general frequent graphspicker\*cond. pattern prob.Unsupervisedrückert '07class-correlated disp-scoreCombinedfcorkcorrespondence-basedSuperviseddtminformation gainSupervised

relation to AUC?

characteristic	correlates w/auc	expected?
# features		$\checkmark$
joint entropy	╉╸╉╸╉╸	7
# equivalence classes	╉╸╉╸	
high correspondence		
low correspondence		





characteristic	correlates w/auc	expected?	
# features	**	$\checkmark$	
joint entropy	╉╋╋	Unsupervised	
# equivalence classes	╺╋╸╋		
high correspondence		$\checkmark$	
low correspondence			





characteristic	correlates w/auc	expected?
# features	**	$\checkmark$
joint entropy	╉╋╋	Unsupervised
# equivalence classes	╺╬╸╬	
high correspondence		$\checkmark$
low correspondence		





characteristic	correlates w/auc	expected?	
# features	**	$\checkmark$	
joint entropy	╉╋╋	Unsupervised	
# equivalence classes	╺╬╸╬		
high correspondence		Supervised	
low correspondence			





characteristic	correlates w/auc	expected?	
# features	**	$\checkmark$	
joint entropy	╉╋╋	Unsupervised	
# equivalence classes	╺╋╸╋		
high correspondence		Supervised	
low correspondence	**		









### what we found - algorithms

characteristic	result
auc	dtm >> fcork
# correspondences	dtm ~ fcork
# equivalence classes	dtm > fcork
# features	dtm >> fcork
running times	dtm >> fcork

### what we found - algorithms

characteristic	result
auc	dtm >> fcork
# correspondences	dtm ~ fcork
# equivalence classes	dtm > fcork
# features	dtm >> fcork
running times	dtm >> fcork

Vit's in the way they work!

#### another tale of two paradigms mine global, use global

8 8 8 + + + + ()

mine global, use global

 $p_1$ 

mine global, use global

 $p_1$ 

mine global, use global

 $p_1$ 

mine global, use global

 $p_1 p_2$ 





mine global, use global

 $p_1 p_2$ 



+ + +

patterns mined on all data
patterns applied to all data

#### another tale of two paradigms mine global, use global

 $p_1 p_2$ 





patterns mined on all data
patterns applied to all data





![](_page_37_Figure_0.jpeg)

![](_page_38_Figure_0.jpeg)

![](_page_39_Figure_0.jpeg)

patterns applied to all data

patterns applied its subset

![](_page_40_Figure_0.jpeg)

![](_page_41_Figure_0.jpeg)

![](_page_42_Figure_0.jpeg)

![](_page_43_Figure_0.jpeg)

![](_page_44_Figure_0.jpeg)

![](_page_45_Figure_0.jpeg)

### experimental evaluation

measure	expected	actual
# equivalence classes	remine ~ dtm	remine ~ dtm 💉
#footuros (procelido)	remine < dtm	remine < dtm 💉
# reatures (prec silue)	remine ? fcork	remine > fcork
# correspondences	remine ~ dtm	remine < dtm 🛕
auc	remine ~ dtm	remine ~ dtm 🎸
rupping timos	remine < dtm	remine < dtm 🗸
running umes	remine ? fcork	remine < fcork 🤐

measure	evaluation
# correspondences	remine < dtm
# equivalence classes	remine < dtm
auc	remine < dtm

measure	evaluation	exploration
# correspondences	remine < dtm	dtm ~ fcork
# equivalence classes	remine < dtm	dtm > fcork
auc	remine < dtm	dtm >> fcork

measure	evaluation	exploration
# correspondences	remine < dtm	dtm ~ fcork
# equivalence classes	remine < dtm	dtm > fcork
auc	remine < dtm	dtm >> fcork
$\Rightarrow$	Unsupervised	

measure	evaluation	exploration
# correspondences	remine < dtm	dtm ~ fcork
# equivalence classes	remine < dtm	dtm > fcork
auc	remine < dtm	dtm >> fcork
$\rightarrow$	Unsupervised	
+ features	Unsupervised fcork < remine	< dtm

measure	evaluation	exploration
# correspondences	remine < dtm	dtm ~ fcork
# equivalence classes	remine < dtm	dtm > fcork
auc	remine < dtm	dtm >> fcork
$\Rightarrow$	Unsupervised	
⇒ # features/auc	Unsupervised	< dtm

measure	evaluation	exploration
# correspondences	remine < dtm	dtm ~ fcork
# equivalence classes	remine < dtm	dtm > fcork
auc	remine < dtm	dtm >> fcork
$\Rightarrow$	Unsupervised	
# features/auc	fcork < remine < dtm	
$\Rightarrow$	overfitting?!	

![](_page_53_Picture_1.jpeg)

#### preliminary evidence that

Unsupervised > Supervised

![](_page_54_Picture_3.jpeg)

#### preliminary evidence that

Unsupervised

>

Supervised

#### & class information can improve partition-scores

#### preliminary evidence that

Unsupervised

S

Supervised

#### & class information can improve partition-scores

![](_page_56_Figure_6.jpeg)

![](_page_57_Figure_1.jpeg)

#### **Thank You for Your Attention**