

A Toolbox for Software Mining

Rahul Premraj Tom Zimmermann

Project #1 Quality of bug reports

What makes a good bug report?

How do they differ from bad ones?

Related: Cool Paper Title [Your names here, Big SE conference]





How long will a bug exist?

What influences its time to fix?

Related: Predicting Eclipse Bug Lifetimes [Panjer, Mining Challenge 2007]



Project #3 Detecting duplicates of bug reports

Is a new bug report a duplicate of another?

Related: Detection of Duplicate Defect Reports Using Natural Language Processing [Runeson et al., ICSE 2007]





How can we assign bugs to developers?

Related: Who Should Fix This Bug? [Anvik et al., ICSE 2006]



Project #5 Visualizing bug reports (networks)

How can we visualize individual bugs?

How can we visualize all bugs?

Related:

Software Bugs and Evolution: A Visual Approach to Uncover Their Relationships [D'Ambros et al., CSMR 2006]



Project #6 Predicting defects with spam filters

Which files have bugs? Which ones don't?

Related: Spam Filter Based Approach for Finding Fault-Prone Software Modules [Mizuno et al., MSR 2007]



Pattern mining





Mining Version Histories to Guide Software Changes (Lecture 3)



Xelopes Library





The Call relation





The Call relation





The Call relation





The Call relation





The Call relation



Many possible patterns



E,

Concept analysis



Concept analysis



Concept Analysis computes all blocks (patterns)

Concept analysis



Concept Analysis computes all blocks (patterns)

Concept analysis



Concept Analysis computes all blocks (patterns)





Support decreases monotonically from top to bottom.



Support decreases monotonically from top to bottom.

Support can be used as a cut-off criterion (support \geq 3).



Support decreases monotonically from top to bottom.

Support can be used as a cut-off criterion (support \geq 3).



Support decreases monotonically from top to bottom.

Support can be used as a cut-off criterion (support \geq 3).

Size increases monotonically from top to bottom. E,

Example patterns

Subject	Pattern			
Ruby CVS	va_end va_start			
Linux 2.6 (s)	mutex_lock mutex_unlock			
Linux 2.6 (xs)	kmem_cache_alloc kmem_cache_free			
Python SVN	PyErr_SetString PyExc_TypeError			
Ruby CVS	rb_fix2int rb_num2int			

Decision trees



Decision trees

Play golf dataset

	Dep. var			
OUTLOOK	TEMPERATURE	HUMIDITY	WINDY	PLAY
sunny	85	85	FALSE	Don't Play
sunny	80	90	TRUE	Don't Play
overcast	83	78	FALSE	Play
rain	70	96	FALSE	Play
rain	68	80	FALSE	Play
rain	65	70	TRUE	Don't Play
overcast	64	65	TRUE	Play
sunny	72	95	FALSE	Don't Play
sunny	69	70	FALSE	Play
rain	75	80	FALSE	Play
sunny	75	70	TRUE	Play
overcast	72	90	TRUE	Play
overcast	81	75	FALSE	Play
rain	71	80	TRUE	Don't Play



Decision trees

Dependent variable: PLAY



Weka



Practical Machine Learning Tools and Techniques





http://www.cs.waikato.ac.nz/ml/weka/

WinMine Toolkit



http://research.microsoft.com/~dmax/winmine/tooldoc.htm

Classification with SVM



The black box: SVMs



SVM: Maximal margins



SVMs: Predictions



SVMs: Non-linear data



SVMs: Non-linear data


Precision and Recall



High precision = returned elements are relevant High recall = relevant elements are returned

Precision and Recall



Recall = TP / (TP + FN) Precision = TP / (TP + FP) Accuracy = (TP + TN) / (TP + FP + FN + TN)









Precision



Precision



Recall



Recall



Visualization



http://www.inf.unisi.ch/phd/dambros/tools/

yEd - Java Graph Editor



http://www.yworks.com/en/products_yed_about.htm

Regression Analysis



What describes the two?



What describes the two?

taste colour size smell shape

texture









Similar problems have similar solutions...

- The car mechanic can infer that if the ignition doesn't turn on, there might be a problem with the spark plug.
- Court case A is similar to my case B. Can I reuse it's solution to win mine?
- Code Reuse ;-)















What about prediction problems?

What about prediction problems?



What about prediction problems?

Recall the Desharnais data set from the Exercises...



New Project

What about prediction problems?



What about prediction problems?



What about prediction problems?



What about prediction problems?

Recall the Desharnais data set from the Exercises...



Average the solution!



Project Length

Evaluation of Results

Evaluation of Results

Actual Effort

Evaluation of Results

Actual Effort Predicted Effort
Actual Effort Predicted Effort Overestimation









Average of Absolute Residuals









Pred(x) % of predictions that lie within x% of Actual Effort



4 of 6 predictions made lie within x% of Actual Effort

Pred(x) = 66.67%















Boneless



Holuid

Congratulations! You won the lottery...



Holuid



Congratulations! You won the lottery...

Congratulations! You won the competition...



\$23 million in your custody

Nigeria



Holuid



Boneless