# Web Pages Classification using Concept Analysis
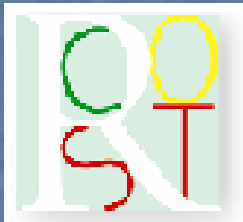
## Porfirio Tramontana

## Anna Rita Fasolino

Dipartimento di Informatica e Sistemistica
*University of Naples Federico II, Italy*

## Giuseppe A. Di Lucca

RCOST – Research Centre on Software Technology

*University of Sannio, Benevento, Italy*

# The Context

- With the diffusion of new paradigms and technological solutions for the Web (RIA and Ajax, SOA, ...), existing Web Applications are rapidly become legacy

- A strategic objective: integrating that existing applications with the new platforms

- An open issue: defining effective approaches for analysing and classifying the Web Applications User Interface
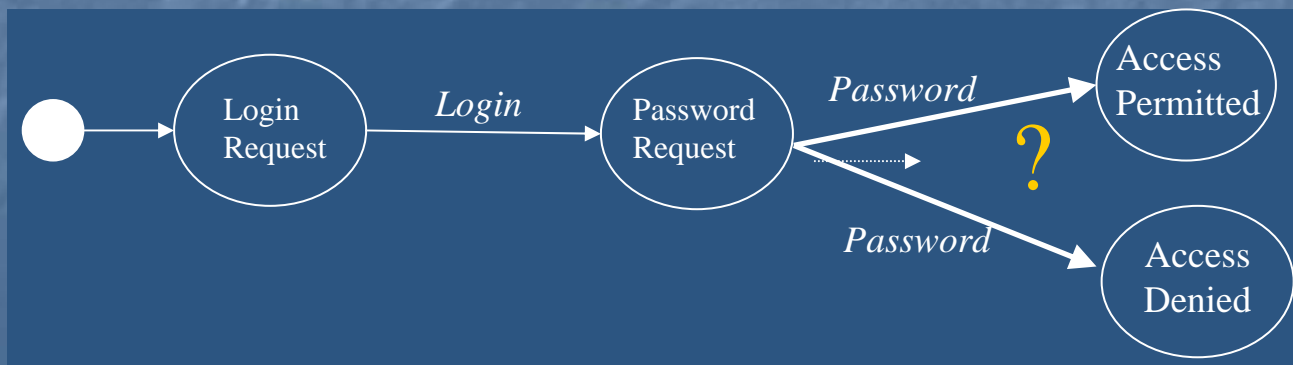
# Automatic Web Pages Classification

- Web Applications interact with users by means of Client Pages that are dynamically generated

- Given a classification of Built Client pages in a set of equivalent classes
  - … classes may correspond to the different reached Web Application scenarios
  - … classes may correspond to the different results of the execution of test cases

- The problem:
  - To propose a technique for the automatic identification of the equivalence class a Web page belongs to

# A possible application scenario

- Designing wrappers that encapsulate the original UI with the aim of exporting a renewed interface, such as a Web Service one
  - The possible reached scenarios represent the classes to recognise.
  - The Wrapper has to automatically identify the class in order to know the state of the interaction with the UI and the actions to be performed

G. Di Lorenzo, A. R. Fasolino, L. Melcarne, P. Tramontana, V. Vittorini, "*Turning Web Applications into Web Services by Wrapping Techniques*" in the 14th Working Conference on Reverse Engineering WCRE 2007, Vancouver, BC, Canada
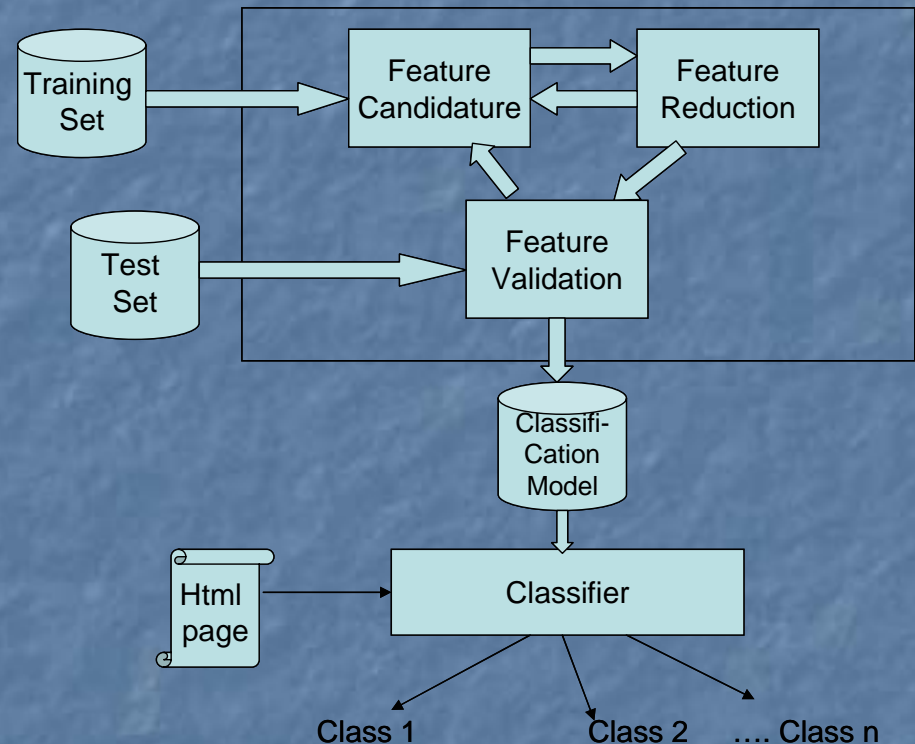
# The proposed solution

- To find, for any class, a combination of Web page features for any class that are satisfied just by pages belonging to that class (Classification model)

- We are interested to a solution that is:
  - Accurate
    - The classification approach should be very reliable, in order to support automatic tasks
  - Efficient
    - The classification approach should support a fast classification of dynamic built client pages

# The proposed approach

## Three step iterative process:

1. Feature Candidature
   - Generation of a set of features, representing page characteristics

2. Feature Reduction
   - Proposition of a combination of the candidate features allowing the identification of the class of any Web page

3. Feature Validation
   - Evaluation of the accurateness of the proposed classification model (i.e. the set of discriminating expressions)



The resulting classification model allows a Classifier module to automatically determine the class any Web page belongs to

# 1 - Feature Candidature

- A **feature** is a characteristic retrievable in a Web page interface
  - It is retrievable by analysing its HTML code
    - A feature can be associated to a XPath query (if the HTML code is XHTML-compatible)

- A set of features that could be useful for discriminating a class are proposed
  - Ideally speaking, we search for features that are retrievable just in the pages belonging to a class

//a/text()="Activities"

//h1/img

/html/body/div[4]/form

7

# 2 - Feature Reduction

- The goal of the Feature Reduction step is to obtain a **Discriminating Expression** for each class, i.e. a combination of features that is true for each Web page belonging to the class, false elsewhere

- The Feature Reduction is executed on the set of candidate features
  - A training set of pre-classified pages is adopted

# Evaluation Matrix

M(i,j) is the value of the Feature j for the Web Page i

Candidate Features

Training
Set
Pages

| | Login | Logout | Incorrect | AdminLog... | LinkToHome |
|------|-------|--------|-----------|-------------|------------|
| P1 | X | | | | X |
| P2 | | X | | X | X |
| P3 | X | | | | X |
| P4 | | X | | | X |
| P5 | X | | X | | X |
| P6 | X | | X | | X |
| P7 | X | | X | | X |

Class C1={P1, P3}
Class C2={P2, P4}
Class C3={P5, P6, P7}

The features with all true (or false) values are deleted: they surely are not useful to discriminate Classes

9

# Concept Analysis

• Concept Analysis is adopted to find features that characterise classes of equivalent pages

• A Concept Lattice may be automatically built on the basis of the Evaluation Matrix values (evaluated on the Training Set)



Attributes ≡ Candidate Features

Objects ≡ Training Set Pages

# Feature Classification

According to Eisenbarth, Koschke and Simon*, features are classified, with respect to classes as:

- **Specific,** if it is satisfied by all and only the pages of the class C;

- **Relevant**: if it is satisfied by all the pages of the class C, but also by other pages from other classes;

- **CSPC** (i.e., Conditionally Specific): if it is satisfied by a subset of pages of the class C and by no other page from other classes;

- **Shared**: if it is satisfied by a subset of pages of the class C and by other pages from other classes too;

- **Irrelevant**: if it is not satisfied by any page of the classes C.

*      T. Eisenbarth, R. Koschke, D. Simon, "Locating features in source code", IEEE Transactions on Software Engineering, Volume 29, Issue 3, IEEE CS Press, March 2003, pp.210 – 224

For the class
Successful Login

• Logout is Specific

• Not(Incorrect) is Relevant

• AdminLogged is CSPC

• Not(AdminLogged) is Shared

• Login is Irrelevant

# Feature Classification Evaluation

■ The Classification of the features may be evaluated on the basis of the Concept Analysis results:

$$Specific(C) = \bigcap_{p \in C} Intent(OC(p)) - \bigcup_{p \in TS-C} Intent(OC(p))$$

$$\mathrm{Re}\,levant(C) = \bigcap_{p \in C} Intent(OC(p)) - Specific(C)$$

$$CSPC(C) = \bigcup_{p \in C} Intent(OC(p)) - \bigcup_{p \in TS-C} Intent(OC(p)) - Specific(C)$$

$$Shared(C) = \bigcup_{p \in C} Intent(OC(p)) - Specific(C) - \mathrm{Re}\,levant(C) - CSPC(C)$$

$$Irrelevant(C) = F - Specific(C) - \mathrm{Re}\,levant(C) - CSPC(C) - Shared(C)$$

# Rules for generating Discriminating Expressions

• Our goal is to obtain combinations of features (*Discriminating Expressions*) that are *specific* for a single class, i.e. that are true for any page belonging to the class and false elsewhere

1. If the class C includes at least one specific feature, any of these specific features can be a candidate discriminating expression of the class.

2. If there are not specific features, then a candidate discriminating expression can be obtained by considering the logic and of the relevant features
   • If there is a couple of features f1 and f2 so that f1==> f2 (i.e. Extent(AC(f1))⊆ Extent(AC(f2))), then it is possible to simplify the expression by discarding the feature f2 from it .

3. If there are neither specific nor relevant features, then a candidate discriminating expression can be obtained by considering the logic or of the CSPC features
   • If there is a couple of features f1 and f2 so that f1==>f2, it is possible to simplify the expression by discarding the feature f1 from it .

14

# 3 - Discriminating Expression Validation

■ **Recall And Precision:**

$$\mathrm{Re}call = \frac{\#(pagesAttributedToClassC)}{\#(AnalysedPagesBelongingToClassC)} \qquad \mathrm{Pr}ecision = \frac{\#(pagesCorrectlyAttributedToClassC)}{\#(pagesAttributedToClassC)}$$

■ *Failed Classification*: the page class has not been identified, because no discriminating feature (or more than one) has been satisfied by the page;

■ *Incorrect Classification*: the page has been attributed to an incorrect class;

■ *Correct Classification*: the page has been attributed to the correct class

# A Classification Example

- Web Pages related to the AdminBooks page of the open source Web application Bookstore has been considered

- Five classes of Web pages

- 25 Web pages training set

- 65 Web pages test set

Empty List



Single Page



Last Page



First Page



Central Page

# A Classification Example



Last Page

Simple Page

Central Page

First Page

# Candidate features

**Candidate Features
-First iteration-**

| Id | Description | XPath Expression |
|---|---|---|
| f1 | Search Form (absolute) | /html/body/table[1]/tbody/tr[1]/td[1]/form[1] |
| f2 | Search Form (relative) | //form[1] |
| f3 | Table with Books (absolute) | /html/body/table[2]/tbody/tr[1]/td[1]/table[1]/tbody/tr[1]/td[1]/a[1]/font/text()="Books" |
| f4 | Table with Books (relative) | //table//text()="Books" |
| f5 | Word "Next" (absolute) | /html/body/table[2]/tbody/tr[1]/td[1]/table[1]/tbody/tr[23]/td[1]/font/a[3]/font/text()="Next" |
| f6 | Word "Next" (relative) | //text()="Next" |
| f7 | Word "Previous" (absolute) | /html/body/table[2]/tbody/tr[1]/td[1]/table[1]/tbody/tr[23]/td[1]/font/a[2]/font/text()="Previous" |
| f8 | Word "Previous" (relative) | //text()="Previous" |

• 8 candidate features have been considered

• f1, f2, f3, f4 have been discarded (always true)

• Negated features have been introduced

**Classification of Candidate features**

| | First Page | Central Page | Last Page | Single page | Empty List |
|---|---|---|---|---|---|
| Specific | | | | | |
| Relevant | $f5, f6, f7, f8$ | $f5, f6, f7, f8$ | $f5, f6, \overline{f7}, f8$ | $\overline{f5}, \overline{f6}, \overline{f7}, \overline{f8}$ | $\overline{f5}, \overline{f6}, \overline{f7}, \overline{f8}$ |
| CSPC | | | | | |
| Shared | | | | | |
| Irrelevant | $\overline{f5}, \overline{f6}, \overline{f7}, \overline{f8}$ | $\overline{f5}, \overline{f6}, \overline{f7}, \overline{f8}$ | $\overline{f5}, \overline{f6}, f7, \overline{f8}$ | $f5, f6, f7, f8$ | $f5, f6, f7, f8$ |

# Discriminating Expressions – First iteration

**Candidate Discriminating Expressions -First iteration-**

| Class | Discriminating Expression | Recall | Precision |
|---|---|---|---|
| First Page | $f7$ | 5/5 | 5/10 |
| Central Page | $f7$ | 5/5 | 5/10 |
| Last Page | $\overline{f7}$ AND $f8$ | 5/5 | 5/5 |
| Single Page | $\overline{f7}$ | 5/5 | 5/10 |
| Empty List | $\overline{f7}$ | 5/5 | 5/10 |

| | | | |
|---|---|---|---|
| f9 | Link "Next" (absolute) | /html/body/table[2]/tbody/tr[1]/td[1]/table[1]/tbody/tr[23]/td[1]/font/a[3][@href]/font/text()="Next" |
| f10 | Link "Next" (relative) | //a[@href]//text()="Next" |
| f11 | Link "Previous" (absolute) | /html/body/table[2]/tbody/tr[1]/td[1]/table[1]/tbody/tr[23]/td[1]/font/a[2][@href]/font/text()="Previous" |
| f12 | Link "Previous" (relative) | //a[@href]//text()="Previous" |
| f13 | Text "No records" (absolute) | /html/body/table[2]/tbody/tr[1]/td[1]/table[1]/tbody/tr[3]/td[1]/font/text()="No records" |
| f14 | Text "No records" (relative) | //text()="No records" |

Training Set Validation

There is not enough precision!

A new iteration of feature candidature is needed
New features have been added:
1. We need to distinguish when the words "Next" and "Previous" are links or simple labels
2. We need to distinguish between empty page and single page

19

# Discriminating Expressions – Second Iteration

| Class | Discriminating Expressions | Training Set | | Test Set | |
|---|---|---|---|---|---|
| | | Recall | Prec. | Recall | Prec. |
| First Page | $\overline{f12}\ AND\ f7$ | 5/5 | 5/5 | 16/16 | 16/16 |
| Central Page | $f11$ | 5/5 | 5/5 | 11/11 | 11/12 |
| Last Page | $f12\ AND\ \overline{f7}$ | 5/5 | 5/5 | 14/15 | 14/14 |
| Single Page | $\overline{f14}\ AND\ \overline{f8}$ | 5/5 | 5/5 | 14/14 | 14/14 |
| Empty List | $f14$ | 5/5 | 5/5 | 9/9 | 9/9 |

Training Set Validation is OK but there are some problems in a page belonging to the Test Set (assigned to Central page instead of Last Page)

The incorrect-assigned page has been added to the Training Set and the Candidate Discriminating Expressions have been evaluated

| Class | Discriminating Expression |
|---|---|
| First Page | $\overline{f12}\ AND\ f7$ |
| Central Page | $f11\ AND\ f10$ |
| Last Page | $f12\ AND\ \overline{f10}$ |
| Single Page | $\overline{f14}\ AND\ \overline{f8}$ |
| Empty List | $f14$ |

← Modified Expressions

Recall and Precision are, now, 100% both on the Training Set and on the Test Set

20

# Discussion

- **Process Effort**
  - In the case study, about 100 minutes were needed to obtain the classification model
  - The activities requiring human intervention was the most critical
    - Training and Test Set Collection step
      - The reliability of the classification model depends on the *representativeness* of the Training and Test Sets
      - Web pages resulting from the execution of a Test Suite may be good candidates to fill in the Training and the Test Set
    - Feature Candidature step
      - Heuristic techniques allowing a semi-automatic generation of features are currently under experimentation

# Further application scenarios

- ## Support to testing automation
  - The possible reached pages are arranged in a set of equivalence classes, corresponding to different test results
  - The proposed discriminating expressions make it possible the automatic evaluation of the testing results

- ## Support to dynamic analysis
  - The recognition of the reached interaction scenarios makes it possible the automatic collection of logs of the visited use cases and scenarios, providing also useful information for user profiling

# Conclusions and Future Works

- A process for the definition of a classification model allowing the automatic identification of classes of Built Client pages has been defined and validated

- The classification model allows the run-time identification of the classification of a Built Client page, with short response time and high precision

- Future Works:
    - We are working for the process improvement, in order to reduce the effort needed to human made tasks
    - Further experimentation will be carried out in order to assess the scalability of the approach

Time is over ... Are there any questions?