

# Exercise: Data Fusion

Web Data Integration



# Agenda

1. Exercise Overview
2. Use Case for this Exercise
3. The WInte.r Framework
  1. Load Data for Fusion
  2. Profile Correspondences
  3. Define a Fusion Strategy
  4. Specify Tolerance Ranges
  5. Run Fusion
  6. Evaluate Fusion Result
  7. Adjust Data Fusion Strategy
4. Hands-on: Tasks of the Exercise

# 1. Exercise Overview

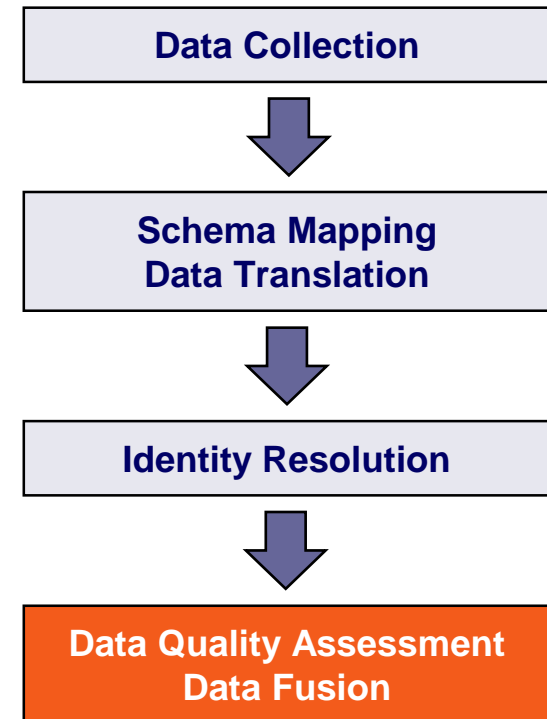
## Learning goal

Learn how to use the Winte.r Framework to:

1. Fuse data and resolve data conflicts
2. Experiment with different conflict resolution functions
3. Measure the density and accuracy of the resulting fused dataset

## Result

1. Fused data set in which each real-world entity is described by only a single record and these records contain no data conflicts
2. A well-founded idea about the quality and profile of the fused data



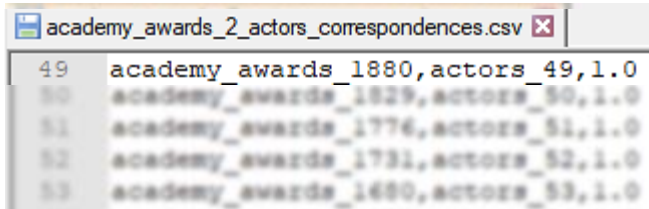
## 2. Use Case for this Exercise

1. Download the .zip of the project from the course page
2. Unzip it and look at the files in \data\ folder
  - .xml input datasets in **input** folder
  - .csv output of Exercise 2 in **correspondences** folder
  - gold.xml in **goldstandard** folder
3. Open the project in a Java IDE (import as maven project)

The project serves as a quick-start for today's tasks and contains implementations for:

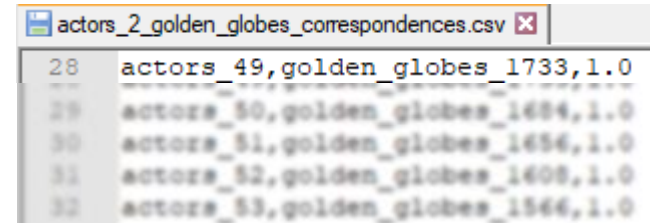
- loading and profiling datasets and correspondences
- various conflict resolution functions
- comparing fusion results to a gold standard (ground truth)
- computing the evaluation metrics density, consistency, and accuracy

# Provided Datasets and Correspondences



id	title	director	actors
49	academy_awards_1880	actors_49	1.0
50	academy_awards_1829	actors_50	1.0
51	academy_awards_1776	actors_51	1.0
52	academy_awards_1731	actors_52	1.0
53	academy_awards_1680	actors_53	1.0

(i) academy awards – actors correspondences



id	title	director	actors
28	actors_49	golden_globes_1733	1.0
29	actors_50	golden_globes_1684	1.0
30	actors_51	golden_globes_1654	1.0
31	actors_52	golden_globes_1608	1.0
32	actors_53	golden_globes_1564	1.0

(ii) actors – golden globes correspondences

```
<movie>
  <id>academy_awards_1880</id>
  <title>One Flew over the Cuckoo's Nest</title>
  <director>
    <name>Milos Forman</name>
  </director>
  <actors>
    <actor>
      <name>Jack Nicholson</name>
    </actor>
    <actor>
      <name>Brad Dourif</name>
    </actor>
    <actor>
      <name>Louise Fletcher</name>
    </actor>
  </actors>
  <date>1975-01-01</date>
  <oscar>yes</oscar>
</movie>
```

```
<movie>
  <id>golden_globes_1733</id>
  <title>One Flew Over The Cuckoo's Nest</title>
  <director>
    <name>Milo Forman</name>
  </director>
  <actors>
    <actor>
      <name>Louise Fletcher</name>
    </actor>
    <actor>
      <name>Jack Nicholson</name>
    </actor>
  </actors>
  <date>1976-01-01</date>
  <globe>yes</globe>
</movie>
```

```
<movie>
  <id>actors_49</id>
  <title>One Flew Over The Cuckoo's Nes</title>
  <actors>
    <actor>
      <name>Louise Fletcher</name>
      <birthday>1934-01-01</birthday>
      <birthplace>Alabama</birthplace>
    </actor>
  </actors>
  <date>1976-01-01</date>
</movie>
```

(i) academy awards dataset

metadata:

dataset\_score: 0.4  
dataset\_date: 2015

(ii) golden globes dataset

metadata:

dataset\_score: 0.5  
dataset\_date: 2018

(iii) actors dataset

metadata:

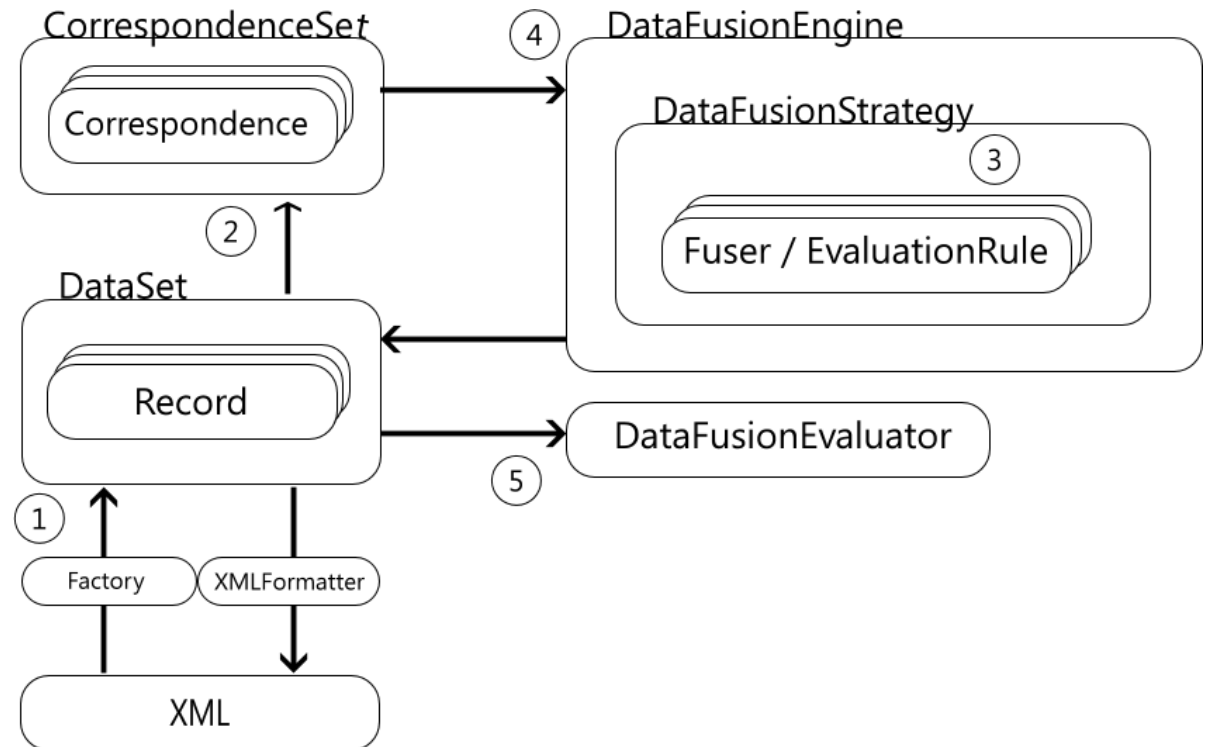
dataset\_score: 0.6  
dataset\_date: 2018

### 3. The WInte.r Framework

- The **W**eb Data **I**ntegration Framework (WInte.r) provides methods for end-to-end data integration
- Implements methods for
  - Data Pre-Processing
  - Schema Matching
  - Identity Resolution
  - **Data Fusion**
  - Evaluation
- Open Source under Apache 2.0 License
- <https://github.com/olehmborg/winter>

# Data Fusion Walkthrough: Movie Use Case

1. Data Sets
2. Correspondences
3. Fusion Strategy
4. Fusion Engine
5. Evaluation



## 3.1 Load Data for Fusion: Data Model

- We need to load:
  1. The individual movie data sets
  2. The correspondences files
- Data model for fusion needs to
  - Extend the *AbstractRecord* class which implements both *Matchable* and *Fusible*

```
public class Movie
    extends AbstractRecord<Attribute> {

    public Movie(
        String identifier,
        String provenance) {
        super(identifier, provenance);
        actors = new LinkedList<>();
    }

    private String title;
    private String director;
    private LocalDateTime date;
    private List<Actor> actors;

    public String getTitle() {
        return title;
    }

    public void setTitle(String title) {
        this.title = title;
    }

    ...
}
```



# Load Data for Fusion: Data Model

- *Custom XMLReader* which implements the FusibleFactory and will create the fused objects

```
public class MovieXMLReader extends XMLMatchableReader<Movie, Attribute> implements
FusibleFactory<Movie, Attribute> {
    @Override
    public Movie createInstanceForFusion(RecordGroup<Movie, Attribute> cluster) {

        List<String> ids = new LinkedList<>();

        // collect the ids of all records that are fused in this group
        for (Movie m : cluster.getRecords()) {
            ids.add(m.getIdentifier());
        }

        // sort and merge the ids to create an id for the fused record
        Collections.sort(ids);

        String mergedId = StringUtils.join(ids, '+');

        // create the fused record
        return new Movie(mergedId, "fused");
    }
}
```

# Load Data for Fusion: Create Fusible Data sets

- Load data using the FusibleDataSet class
  - Extends the DataSet class


```
// Load the Data into FusibleDataSet
FusibleDataSet<Movie, Attribute> ds1 = new FusibleHashedDataSet<>();
new MovieXMLReader().loadFromXML(new File("data/input/academy_awards.xml"), "/movies/movie", ds1);
```

- Allows you to add data set meta data

```
ds1.setScore(3.0);
ds1.setDate(LocalDate.parse("2012-01-01", formatter));
```

- Calculates density report

```
ds1.printDataSetDensityReport();
```



```
DataSet density:                0,58
Attributes densities:
    Actors:                     0,23
    Date:                       1,00
    Title:                      1,00
    Director:                   0,09
```

# Load Data for Fusion: Load Correspondences

- Use the *CorrespondenceSet* class to load all correspondences found in data/correspondences
  - Call *loadCorrespondences* multiple times!

```
// load the correspondences
CorrespondenceSet<FusableMovie> correspondences = new CorrespondenceSet<>();

correspondences.loadCorrespondences( new File("correspondences.csv"),
                                     ds1,
                                     ds2);
```

- The order of the data sets (2<sup>nd</sup> and 3<sup>rd</sup> parameter) must match the order of Ids in the correspondence file!

## 3.2 Profile Correspondences: Group Size Distribution

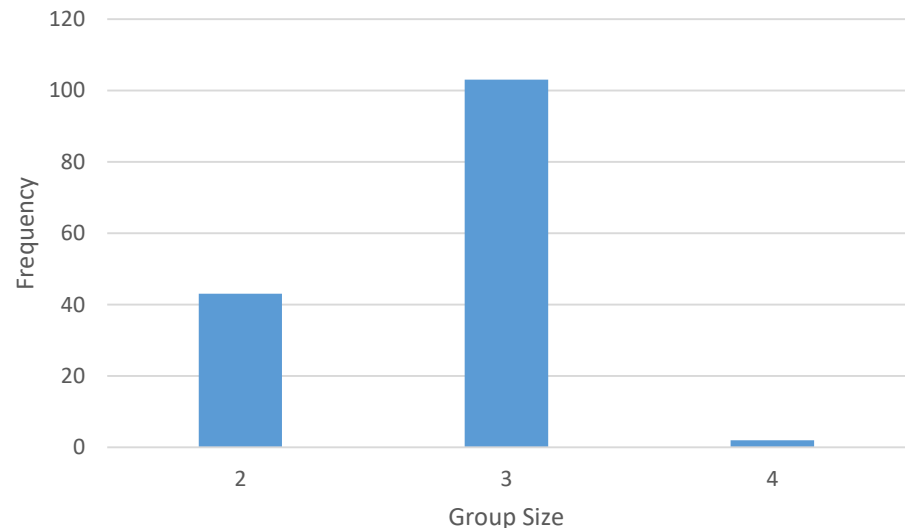
- Check group size distribution
  - All correspondences from the identity resolution are combined
  - All records that are believed to describe the same real-world entity end up in a group (transitive!)

```
// write group size distribution
correspondences.printGroupSizeDistribution();
```

Group Size Distribution of 148 groups:  
Group Size | Frequency

2		43
3		103
4		2

Due to matching errors  
or duplicates in sources



## 3.3 Define Data Fusion Strategy

- Use the *DataFusionStrategy* class to define how each attribute is fused
- For each attribute, you have to add a *Fuser* and an *EvaluationRule*
  - Fusers use a conflict resolution function to fuse the values for an attribute
  - EvaluationRules define tolerance range for considering two values as equal

```
// define the fusion strategy
DataFusionStrategy<Movie,Attribute> strategy =
    new DataFusionStrategy<>(new MovieXMLReader ());

// add attribute fusers
// Note: The attribute name is only used for printing the reports
strategy.addAttributeFuser( "Title",
                           new TitleFuser(),
                           new TitleEvaluationRule());
```

# Define Fusion Strategy: Fuser

- A fuser defines a conflict resolution function
  - That decides which of many values to choose
- And three additional functions
  - If a record has a value that can be fused
  - How to get the value for a record
  - How to assign the fused value to the fused record

```
public class TitleFuser extends AttributeValueFuser<String, Movie, Attribute> {  
    public TitleFuser() {  
        super(new LongestString<Movie>());  
    }  
  
    @Override  
    public boolean hasValue(Movie record) {  
        return record.hasValue(Movie.TITLE);  
    }  
  
    @Override  
    public String getValue(Movie record) {  
        return record.getTitle();  
    }  
}
```

Define the conflict resolution function

Does *record* have a value for this fuser?

Get *record*'s value for this fuser?

# Define Fusion Strategy: Fuser

- A fuser defines a conflict resolution function
  - That decides which of many values to choose
- And three additional functions
  - If a record has a value that can be fused
  - How to get the value for a record
  - How to assign the fused value to the fused record

```
@Override
public void fuse (RecordCluster<Movie> group,
    Movie fusedRecord) {

    FusedValue<String, Movie> fused = getFusedValue (group);

    fusedRecord.setTitle (fused.getValue());

    fusedRecord.setAttributeProvenance (Movie.TITLE, fused.getOriginalIds());
}
```

Get the fused value

Assign the fused value to the fused record

Add provenance info to the fused record

# Implemented Conflict Resolution Functions

## Numeric Functions

Average	Calculates the average of all values
Median	Calculates the median of all values

## String Functions

Longest String	Chooses the longest string
Shortest String	Chooses the shortest string

## List Functions

Union	Creates the union of all values of the lists
Intersection	Creates the intersection of all values of the lists

## Functions that use Provenance Metadata

FavourSources	Chooses the value from the data source with the highest score
MostRecent	Chooses the most up-to-date value

## Data Type Independent Functions

Voting	Chooses the most frequent value
ClusteredVote	Chooses the centroid of the largest value cluster



## 3.4 Specify Tolerance Range

- Evaluation rules are used to calculate consistency and accuracy
- Implement one class per attribute specific evaluation rule which
  - extends the *EvaluationRule* class
  - defines which values should be considered as equal
- It might be OK to tolerate
  - +/- 2% for numeric data like temperature, population
  - edit distance 1 for people names but likely not for movie names
- Ideal case: Complete equality, no tolerance needed

```
public class TitleEvaluationRule extends EvaluationRule<Movie, Attribute> {  
  
    SimilarityMeasure<String> sim = new TokenizingJaccardSimilarity();  
  
    @Override  
    public boolean isEqual(Movie record1, Movie record2) {  
        // the title is correct if all tokens are there, but the order does not matter  
        return sim.calculate(record1.getTitle(), record2.getTitle()) == 1.0;  
    }  
}
```

## 3.5 Run Fusion

- Use the DataFusionEngine class to
  - calculate value consistency

```
// create the fusion engine
DataFusionEngine<Movie,Attribute> engine = new DataFusionEngine<>(strategy);

// calculate attribute consistency
engine.printClusterConsistencyReport(correspondences, null);
```

```
Attribute Consistencies:
Actors: 0,43
Date: 0,00
Title: 0,94
Director: 1,00
```

- and fuse the data sets

```
// run the fusion
FusableDataSet<Movie,Attribute> fusedDataSet = engine.run(correspondences, null);
```

# Writing the Fused Dataset to an XML File

- Use XMLFormatter to write your dataset to an XML file

```
public class MovieXMLFormatter extends XMLFormatter<Movie> {
    @Override
    public Element createRootElement(Document doc) {
        return doc.createElement("movies");
    }

    @Override
    public Element createElementFromRecord(FusableMovie record, Document doc) {
        Element movie = doc.createElement("movie");

        movie.appendChild(createTextElement("id",
            record.getIdentifier(), doc));

        movie.appendChild(createTextElementWithProvenance("title",
            record.getTitle(),
            record.getMergedAttributeProvenance(Movie.TITLE), doc));
        ...
        return movie;
    }
}
```

# Example Record in the Fused Dataset

```
<?xml version="1.0" encoding="UTF-8"?>
<movies>
  <movie>
    <id>academy_awards_2334+actors_39+golden_globes_2000</id>
    <title provenance="golden_globes_2000">Who's Afraid Of Virginia Woolf</title>
    <director provenance="academy_awards_2334">Mike Nichols</director>
    <date provenance="golden_globes_2000+actors_39">1967-01-01T00:00</date>
    <actors provenance="academy_awards_2334+golden_globes_2000+actors_39">
      <actor>
        <name>George Segal</name>
      </actor>
      <actor>
        <name>Richard Burton</name>
      </actor>
      <actor>
        <name>Elizabeth Taylor</name>
      </actor>
      <actor>
        <name>Sandy Dennis</name>
      </actor>
    </actors>
  </movie>

```

## 3.6 Evaluate Fusion Result

- Use the DataFusionEvaluator class
  - needs your fusion strategy!
  - accepts a DataSet as gold standard (ground truth)
  - returns the accuracy of the dataset compared to the gold standard

```
// load the gold standard
DataSet<Movie, Attribute> gs = new FusibleHashedDataSet<>();

new MovieXMLReader().loadFromXML(        new File("data/goldstandard/fused.xml"),
                                          "/movies/movie", gs);

// evaluate
DataFusionEvaluator<Movie, Attribute> evaluator = new DataFusionEvaluator<>(strategy);

double accuracy = evaluator.evaluate(fusedDataSet, gs, null);
```

# Evaluate Fusion: an example

## Record generated by data fusion

```
<movie>
  <id>academy_awards_1880+actors_126+actors_49+golden_globes_1733</id>
  <title provenance="golden_globes_1733">One Flew Over The Cuckoo's Nest</title>
  <director provenance="academy_awards_1880">Milos Forman</director>
  <date provenance="academy_awards_1880">1975-01-01T00:00</date>
  <actors provenance="actors_126+academy_awards_1880+golden_globes_1733+actors_49">
    <actor>
      <name>Brad Dourif</name>
    </actor>
    <actor>
      <name>Louise Fletcher</name>
    </actor>
    <actor>
      <name>Jack Nicholson</name>
    </actor>
  </actors>
</movie>
```

title with the longest string

director from most trusted source

actors list from union

## Gold standard record

```
<movie>
  <id>academy_awards_1880</id>
  <title>One Flew over the Cuckoo's Nest</title>
  <director>
    <name>Milos Forman</name>
  </director>
  <actors>
    <actor>
      <name>Jack Nicholson</name>
    </actor>
    <actor>
      <name>Brad Dourif</name>
    </actor>
    <actor>
      <name>Louise Fletcher</name>
    </actor>
  </actors>
  <date>1975-01-01</date>
  <oscar>yes</oscar>
</movie>
```

Compare the fused values with the ones in the gold standard:

- Did your conflict resolution functions work?
- Should the value comparison be strict or more relaxed?
  - „Cuckoo“s Nest“ vs „Cuckoo’s Nest“

## 3.7 Adjust the Data Fusion Strategy

- Winte.r supports detailed event logging which can help you adjust your fusion strategy and improve your results
- **Default** logging level: dataset densities, group size distributions, gold standard elements, attributes consistencies, fusion accuracy
- **Trace** (tracefile): default + fusion errors and corresponding correct values from the gold standard, attribute specific accuracy

```
private static final Logger logger = WinterLogManager.activateLogger("trace");
```

- Activate the fusion debug report after initializing the DataFusionStrategy object

```
// define the fusion strategy
DataFusionStrategy<Movie, Attribute> strategy = new DataFusionStrategy<>(new MovieXMLReader());
// write debug results to file
strategy.activateDebugReport("data/output/debugResultsDatafusion.csv", -1, gs);
```

# Adjust the Data Fusion Strategy: an example

- Inspect the debug report

Attribute Name	Consistency	ValueIDS	Values	FusedValue	IsCorrect	CorrectValue
Date	0	Date-{actors_97 academy_awards_3480}	{1947-01-01T00:00 1946-01-01T00:00}	1947-01-01T00:00	FALSE	1946-01-01T00:00
Date	0	Date-{actors_104 academy_awards_3059}	{1954-01-01T00:00 1953-01-01T00:00}	1954-01-01T00:00	FALSE	1953-01-01T00:00
Date	0	Date-{actors_91 academy_awards_4015}	{1941-01-01T00:00 1940-01-01T00:00}	1941-01-01T00:00	FALSE	1940-01-01T00:00

the fused value comes from the *actors* dataset as a result of the **FavourSource** fusion strategy

the fused value is **wrong**

the correct date agrees with the *academy\_awards* date value

- Possible implication if this is the case for many values

Change dataset scores from

**academy awards dataset**

metadata:  
dataset\_score: 0.4

**actors dataset**

metadata:  
dataset\_score: 0.6

to

**academy awards dataset**

metadata:  
**dataset\_score: 0.6**

**actors dataset**

metadata:  
**dataset\_score: 0.4**



# WInte.r Tutorial also covers Data Fusion

- <https://github.com/olehmborg/winter/wiki/WInte.r-Tutorial>

The screenshot shows the GitHub Wiki page for the WInte.r Tutorial. The page title is "WInte.r Tutorial" and it was edited 15 minutes ago. The page content includes an introduction to the WInte.r framework for identity resolution and data fusion, followed by a list of topics covered in the tutorial. A sidebar on the right shows a table of contents with 16 pages.

olehmborg / winter

Unwatch 5 Star 41 Fork 13

Code Issues 0 Pull requests 0 Projects 0 Wiki Insights

## WInte.r Tutorial

olehmborg edited this page 15 minutes ago · 1 revision

Edit New Page

This tutorial gives a step-by-step introduction to using the WInte.r framework for identity resolution and data fusion. The goal of identity resolution (also known as data matching or record linkage) is to identify records in different datasets that describe the same real-world entity. Data fusion methods merge all records which describe the same real-world entity into a single, consolidated record while resolving data conflicts.

The tutorial explains the usage of WInte.r along the use case of integrating data about movies. The goal is to integrate the two datasets [Actors](#) and [Academy awards](#) into a single, duplicate-free dataset containing comprehensive descriptions of all movies. The complete source code of the tutorial is found in the folder [use case / movies](#).

The tutorial is structured as follows:

1. [Overview of the Datasets](#)
2. [Define Data Model and Load Data](#)
3. [Identity Resolution](#)
  - i. [Creating a Matching Rule](#)
  - ii. [Running the Identity Resolution](#)
  - iii. [Creating a Gold Standard for Identity Resolution](#)
  - iv. [Evaluating the Matching Result](#)

Pages 16

Contents

- [Data Model](#)
  - [Data Normalisation](#)
  - [Web Tables](#)
- [Matching](#)
  - [Similarity Measures](#)
  - [Blocking](#)
  - [Schema Matching](#)
  - [Identity Resolution](#)
  - [Learning Matching Rules](#)
    - [RapidMiner Integration](#)
- [Data Fusion](#)
- [Evaluation](#)
- [Event and Result Logging](#)
- [WInte.r Tutorial: Movie Data Integration](#)

# Project Related Information

# Select Data for Fusion Experiments

- Your input is the output of Exercises 1 and 2
  - schema is aligned, unique IDs are in place, identity resolution is done
- What data is suitable for fusion experiments?
  - Select **at least 2 datasets** for which
    1. attribute intersection is big enough
      - you should be able to fuse data for  **$\geq 5$  attributes**
    2. quality of identity resolution is good enough
    3. it makes sense to apply different fusion functions to the attributes
      - compare different functions
      - experiment with using provenance data vs. no provenance data

# Prepare Your Gold Standard

- Your gold standard should contain
  - $\geq 15$  entities
  - $\geq 4$  attributes per entity
- Manually look up the correct values (in an external data source)!
- The Gold Standard uses your target schema
  - Use IDs from one of your data sets, so the evaluator can find the correct records!
- Don't create ambiguities!
  - If records a, b, c are the same according to the identity resolution
  - Only one corresponding record can appear in the gold standard
  - Choose a, b or c as ID in the gold standard

# Requirements for the Final Project Report

- **12 pages** (sharp!) – counted without title page, table of content, literature list
  - Every extra page (including appendix pages) will reduce your mark by 0.33
- Due to **Friday, 29th November 2019, 23:59**
  - Send by **email to Chris, Ralph and Anna**
- You must use the **DWS master thesis layout** (without Chapters)
  - <https://www.uni-mannheim.de/dws/teaching/thesis-guidelines/>
- Also **submit**
  1. **your code** and
  2. (a subset) of **your data**
- Please cite sources properly if you use any
  - Preferred citation style [Author, year]

# Final Report: Content

- Your final report should contain
  1. Results of Phase 1: Data Translation
  2. Results of Phase 2: Identity Resolution
  3. Results of Phase 3: Data Fusion
  4. Summary of the overall result

# Data Translation in the Final Report

- Your report should contain
  1. Profiling results describing your **input data sets**
    - e.g. updated versions of the tables that you created for your project proposal

Dataset	Source (*)	format	class (**)	# of entities	# of attributes	list of attributes (***)
IMDB	Download URL	csv	movie	17,000	10	title, director, year, ...
DBpedia	dbpedia.org/sparql	xml	actor	23,500	8	name, birthDate (MV), activeYears,...
Freebase	Download URL	csv	actor	11,000	14	given_name, surname, spouse (MV)

.....

Class name	Attribute name	Datasets in which attribute is found
movie	name	dataset1, dataset2, dataset3, dataset4
movie	director	dataset1, dataset3
movie	year	dataset2, dataset3, dataset4

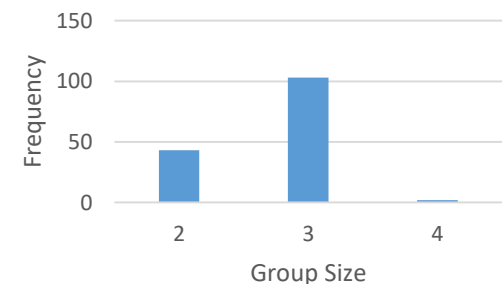
.....

2. Your **consolidated schema** and how you created it
3. Which **transformations** you used and why
  - if there was any information you could not transform

# Identity Resolution in the Final Report

- Your report should contain
  1. Content and size of your **gold standard** and procedure used to create it
  2. Which **matching rules** did you try?
    - Discuss what happened with P/R and F1?
    - Please include a table comparing the results of the different matching rules
  3. Which **blocking methods** did you try?
    - Report and discuss the change in runtime, number of matches, and reduction ratio. How do P/R/F1 change?
    - Please include a table comparing the results of the different blocking methods that you tested
  4. What's the group size distribution of your result?
  5. An **analysis of the errors** that remain when applying your best matching rule.

#	Matching Rule	Blocker	P	R	F1	# Corr	Time
1	Rule1:Title&Year	No Blocking	0.71	0.95	0.82	10.230	90 min
2	Rule1:Title&Year	StandardYear	0.71	0.73	0.72	9.609	18 sec
3	Rule1:Title&Year	SNBYear	0.71	0.89	0.79	10.215	50 sec
4	Rule2:Title&Actors	SNBYear	0.81	0.89	0.83	9.919	19 sec





# Data Fusion in the Final Report

- Your report should contain
  1. Which **datasets** you selected for fusion?
  2. What kind of **provenance data** you added?
  3. What was the **density** of your input and the merged datasets?
  4. How **consistent** were your datasets?
  5. Size and content of your **gold standard** and how you created it
  6. Which **conflict resolution functions** you tried for each attribute
    - Whether you define your own conflict resolution functions
  7. Which **accuracy** did the different conflict resolution functions deliver? What was the best function for each attribute?
    - Please include a table comparing the accuracies reached by the different resolution functions.
  8. An **analysis of the errors** that remain after applying your best resolution function.

# Summary of Overall Results in Final Report

Please conclude your report with a summary answering the following three questions:

1. How many **additional entities** did you add compared to the largest of your input datasets?
2. How much did you **increase the density** of your data compared to the largest of your input datasets?
3. What is the **overall accuracy** of your final dataset according to your gold standard?

# Final Report: Important

- **Balance your content** between the 3 exercises
  - not 10 pages on identity resolution and 2 pages on the rest
- If you have done something cool – **write about it!**
  - it is highly unlikely we dig it out of your code ourselves
- We are strict about the **12 pages limit**. Thus,
  - include lots of tables to show us what you have tried
  - briefly discuss the results of each of your experiments
  - do not repeat theoretical stuff from the slides (e.g. definition of X)
  - we will reduce your mark by 0.33 for each extra page no matter how interesting it is!

# Task

1. Open and run the provided Java project
2. Inspect the profiling results
  1. What is the potential of every dataset to fill missing values within the other datasets?
  2. How many real-world entities are presented by more than one dataset?
  3. How is it possible to have 4 conflicting values given that we have only 3 datasets?
  4. Which attributes do you expect to be more conflicting?
3. Create your own value-based fusion strategy that selects the longest title value
4. Inspect the debug report to find which source is the most trustworthy for the date attribute. Change the dataset metadata so that this source is preferred.

# Solutions

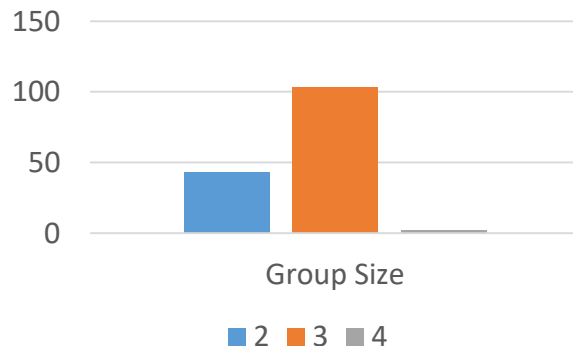
## Interpreting profiling results

Dataset	Number of Elements	Overall Density	Attribute Density			
			Actors	Date	Title	Director
academy_awards.xml	4580	0.58	0.23	1.00	1.00	0.09
actors.xml	151	0.75	1.00	1.00	1.00	0.00
golden_globes.xml	2279	0.78	0.98	1.00	1.00	0.14

Good for missing values

High overlap: More potential for conflicts

Group size distribution of 148 groups



**D1**

```
<movie>
  <title> Crank 1</title>
</movie>
<movie>
  <title> Crank 2</title>
</movie>
```

**D2**

```
<movie>
  <title> Crank </title>
</movie>
```

**Merged**

```
<movie>
  <title_1> Crank 1</title_1>
  <title_1> Crank 2</title_1>
  <title_2> Crank </title_2>
</movie>
```

# Solutions

Create your own value-based fusion strategy that selects the longest title value

```
public class TitleFuserLongestString extends AttributeValueFuser<String, Movie,
Attribute> {

    public TitleFuserLongestString() {
        super(new LongestString<Movie, Attribute>());
    }

    public void fuse(RecordGroup<Movie, Attribute> group, Movie fusedRecord,
Processable<Correspondence<Attribute, Matchable>> schemaCorrespondences,
Attribute schemaElement) {
        ...
    }
    public boolean hasValue(Movie record, Correspondence<Attribute, Matchable>
correspondence) {
        ...
    }
    protected String getValue(Movie record, Correspondence<Attribute, Matchable>
correspondence) {
        ...
    }
}
```

# Solutions

Which source is the most trustworthy as far as the date attribute is concerned?

- Inspect the log messages

```
Error in 'Date': [Movie academy_awards_4080+actors_90: Goodbye, Mr. Chips / Sam Wood / 1940-01-01T00:00] <> [Movie academy_awa
Error in 'Date': [Movie academy_awards_4081+actors_12: Gone with the Wind / Victor Fleming / 1940-01-01T00:00] <> [Movie acade
Error in 'Date': [Movie academy_awards_3480+actors_97: The Best Years of Our Lives / William Wyler / 1947-01-01T00:00] <> [Mov
Error in 'Date': [Movie academy_awards_4270+actors_9: The Great Zeigfeld / Robert Z. Leonard / 1937-01-01T00:00] <> [Movie aca
```

```
ds1.setScore(1.0);
ds2.setScore(2.0);
ds3.setScore(3.0);
```

- The current strategy selects the date value from the *actors* dataset (ds2)

# Solutions

Which source is the most trustworthy as far as the date attribute is concerned?

- Adjust the dataset scores
- Use a conflict resolution function that uses these scores

```
strategy.addAttributeFuser(  
    Movie.DATE,  
    new DateFuserFavourSource(),  
    new DateEvaluationRule());
```

- Check how the accuracy changes

```
ds1.setScore(1.0);  
ds2.setScore(2.0);  
ds3.setScore(3.0);
```

Attribute-specific accuracy: 0.00

```
ds1.setScore(3.0);  
ds2.setScore(1.0);  
ds3.setScore(2.0);
```

Attribute-specific accuracy: 0.95



# Final Presentation and Exam

- Presentation Dates
  - 04.12.2019 and 05.12.2019
- Date and Time of WDI Exam
  - Thursday, 12.12.2019, 8:30 A5 B244
- Format
  - 5-6 open questions that show that you have understood the theory part of the lecture
    - all lecture slide sets including structured data on the Web and data exchange formats + query languages XPath and SPARQL
  - Duration: 60 minutes

# ...and now

- Get the template project and
  - Define your inputs
  - Experiment with creating the merged dataset, and density and consistency evaluation metrics
  - Define your conflict resolution functions
  - Define your gold standard
  - Experiment with data fusion and accuracy evaluation metrics
- **Write your final report**
- **Repeat the theory parts in order to be ready for the final exam**
  - The video recordings of all lectures are online

