



# DATA ENGINEERING CHALLENGE

Gokul Kannan

Master of Science, Computer Science



2

# CHALLENGE 1: Skyline Operator

# Skyline Operator : Overview

- **Skyline** query returns a set of points  $P$ , such that any point  $p_i = (x_i; y_i)$  in  $P$  is not **dominated** by any other point in the dataset.
- Different methods to identify Skyline points are
  - Block Nested Loop
  - Divide and Conquer
  - Plane-sweep

# Skyline Operator : Goal

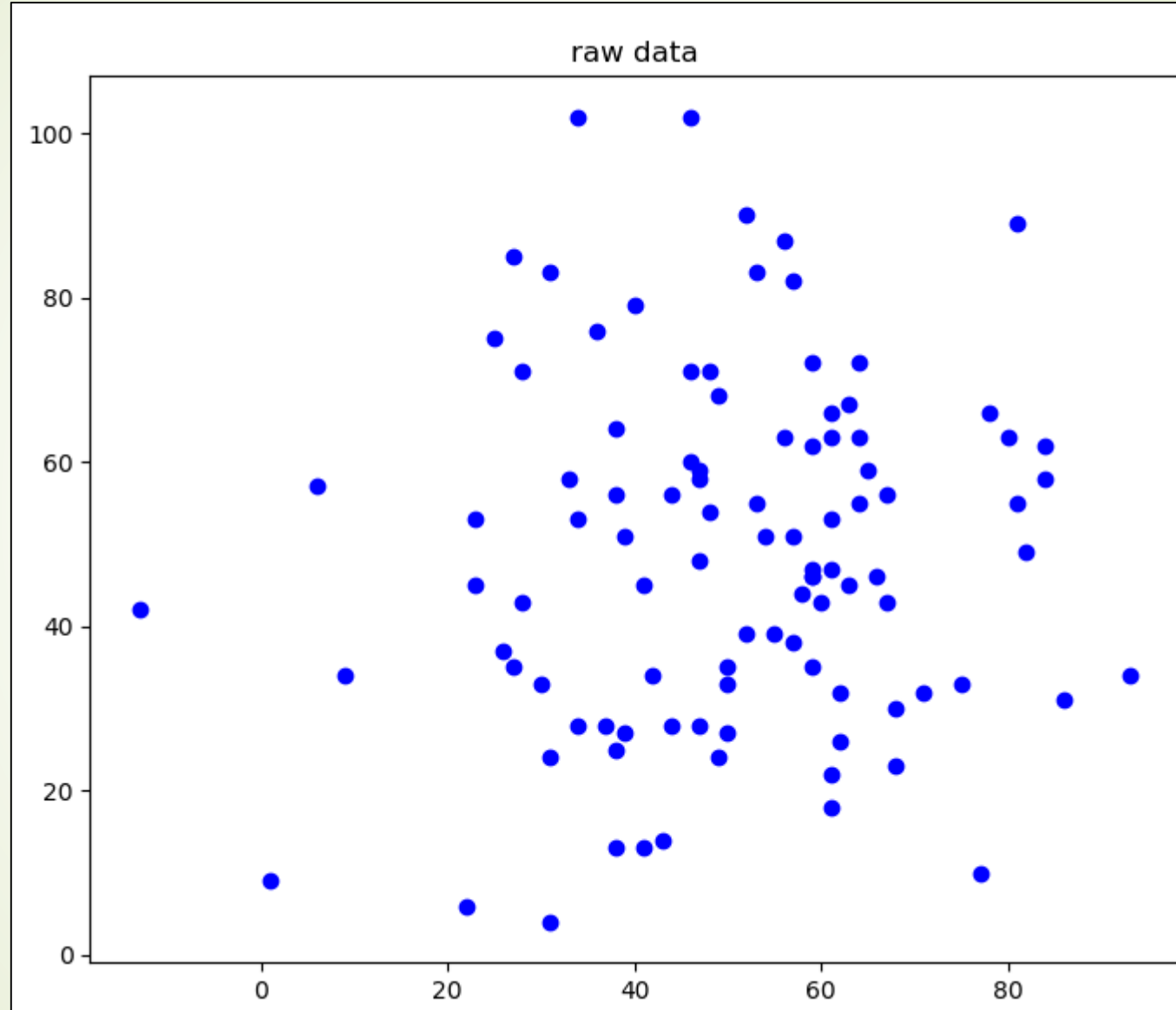
Objective	Identify all the skyline points
Input	skylinepoints.csv file A set of two-dimensional points $p_i = (x_i; y_i)$ as input
Output	<ul style="list-style-type: none"><li>• The set of all skyline points that is not worse than any other point in data set with respect to the dimensions x and y.</li><li>• The program execution time must be calculated.</li></ul>
Assumptions	Higher values are preferred in both dimensions x and y
Constraints	<ol style="list-style-type: none"><li>1. Efficiency : Overall program execution time should be less.</li><li>2. Algorithm shall not be tailored to the provided data set.</li><li>3. Any Algorithm shall be used to identify the skyline points</li></ol>
Programming Language	The algorithm shall be implemented in Java 8

# Skyline Operator : Algorithm

- Sort the input data pair (X, Y) with respect to data point X in **descending order**.
  - If X has same values, sort according to the values of Y.
  - Remove duplicate points
  - The first point in the sorted list is an Skyline point.
  - Store the Y value of this point (Max\_Y).
- Iterate through the sorted array.
- If there is an y value in the dataset, which is greater than the stored y value (Max\_Y),
  - Add this point as Skyline point.
  - Update the Max\_Y value
  - Repeat this process for all the sorted points.

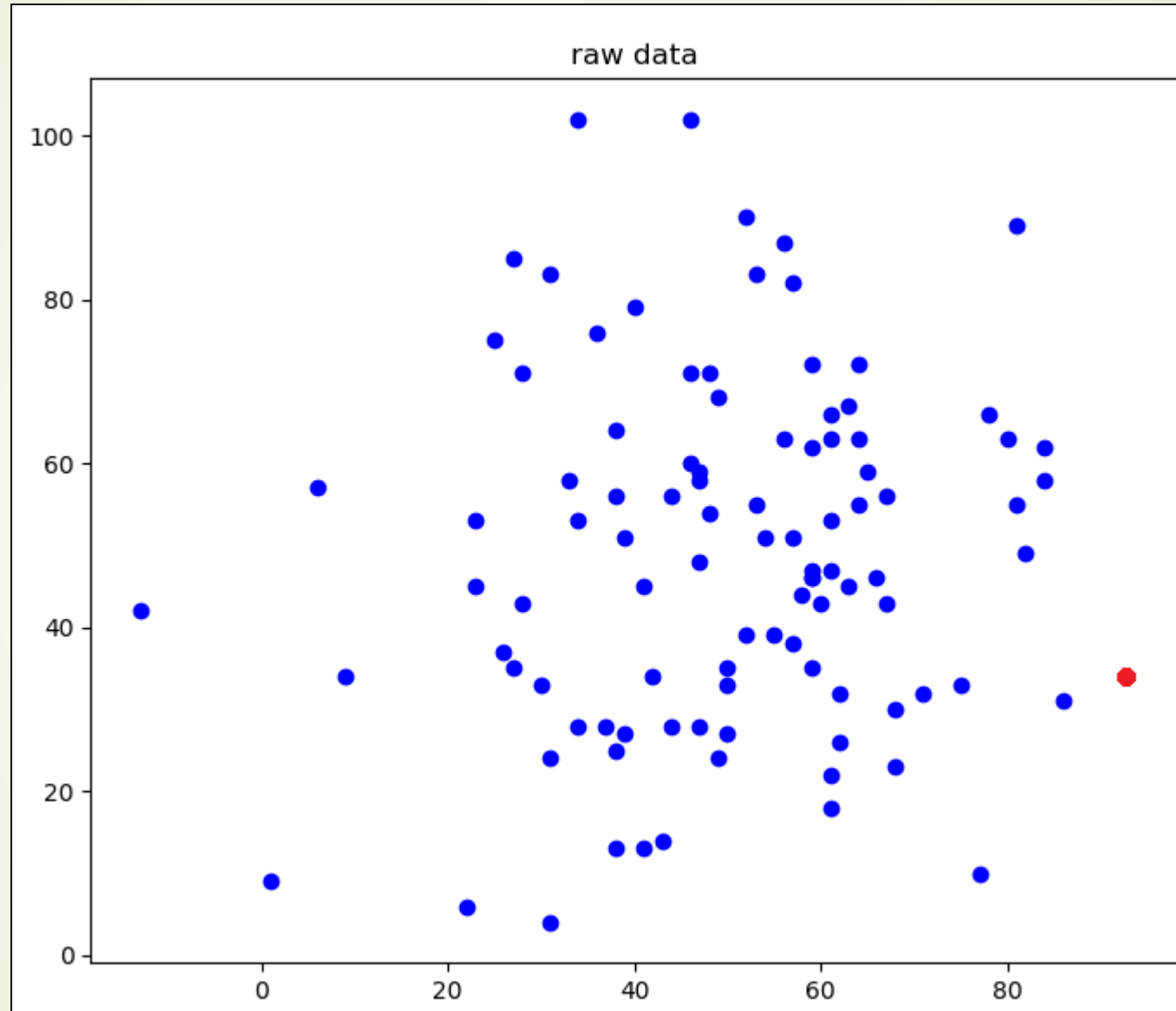
# Skyline Operator : Input Data Points

6



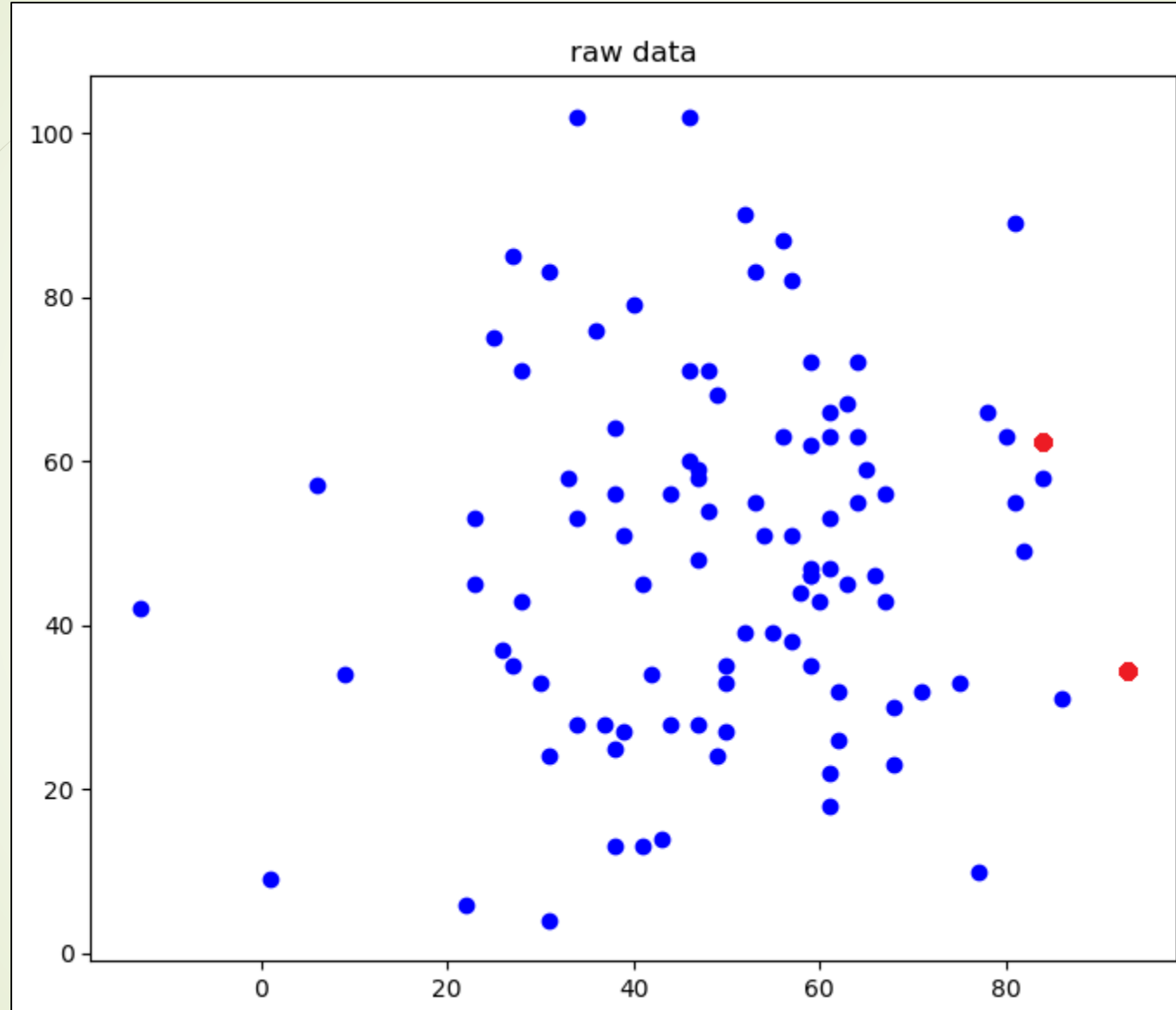
# Skyline Operator : Skyline Point 1

7



# Skyline Operator : Skyline Point 2

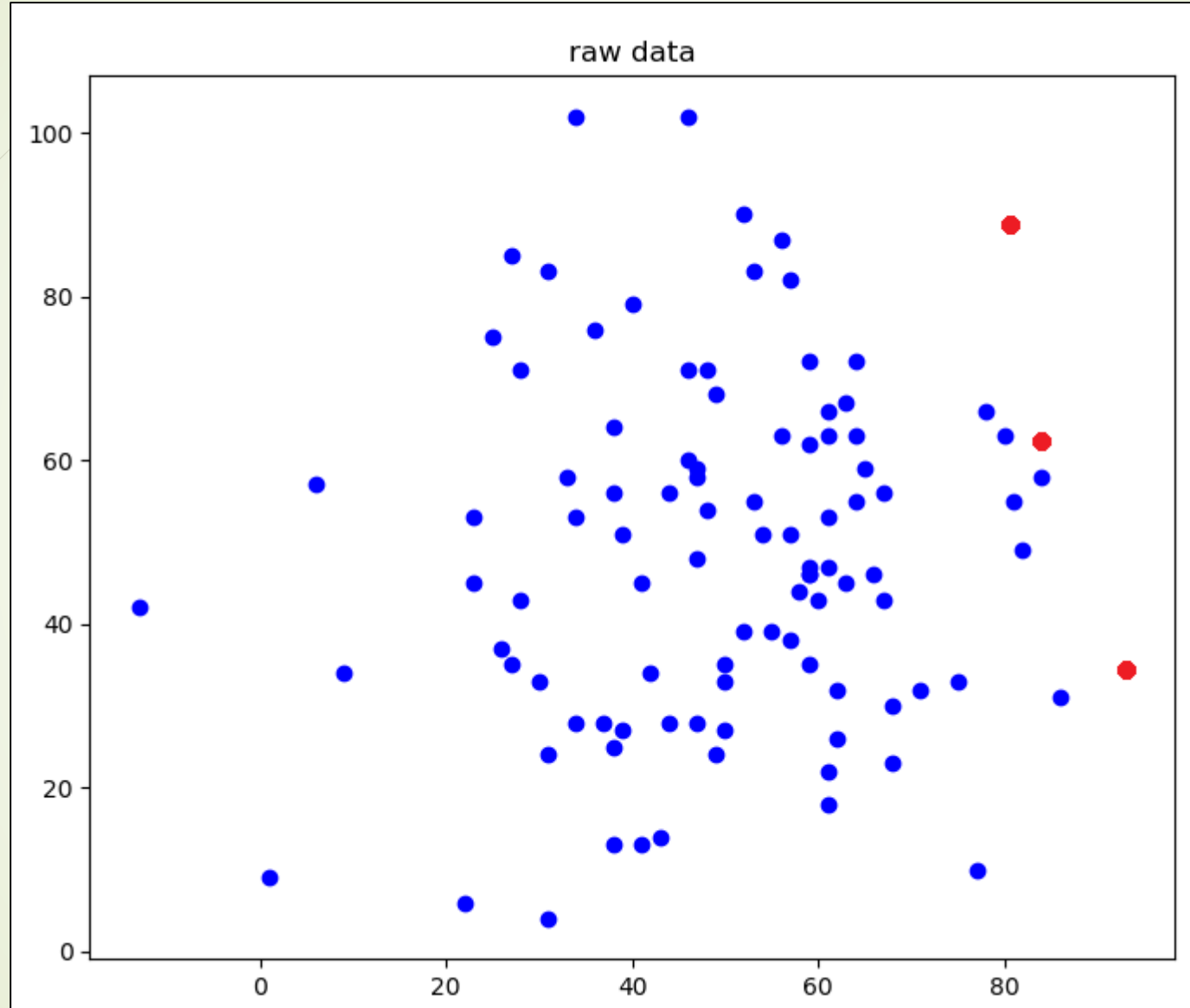
8





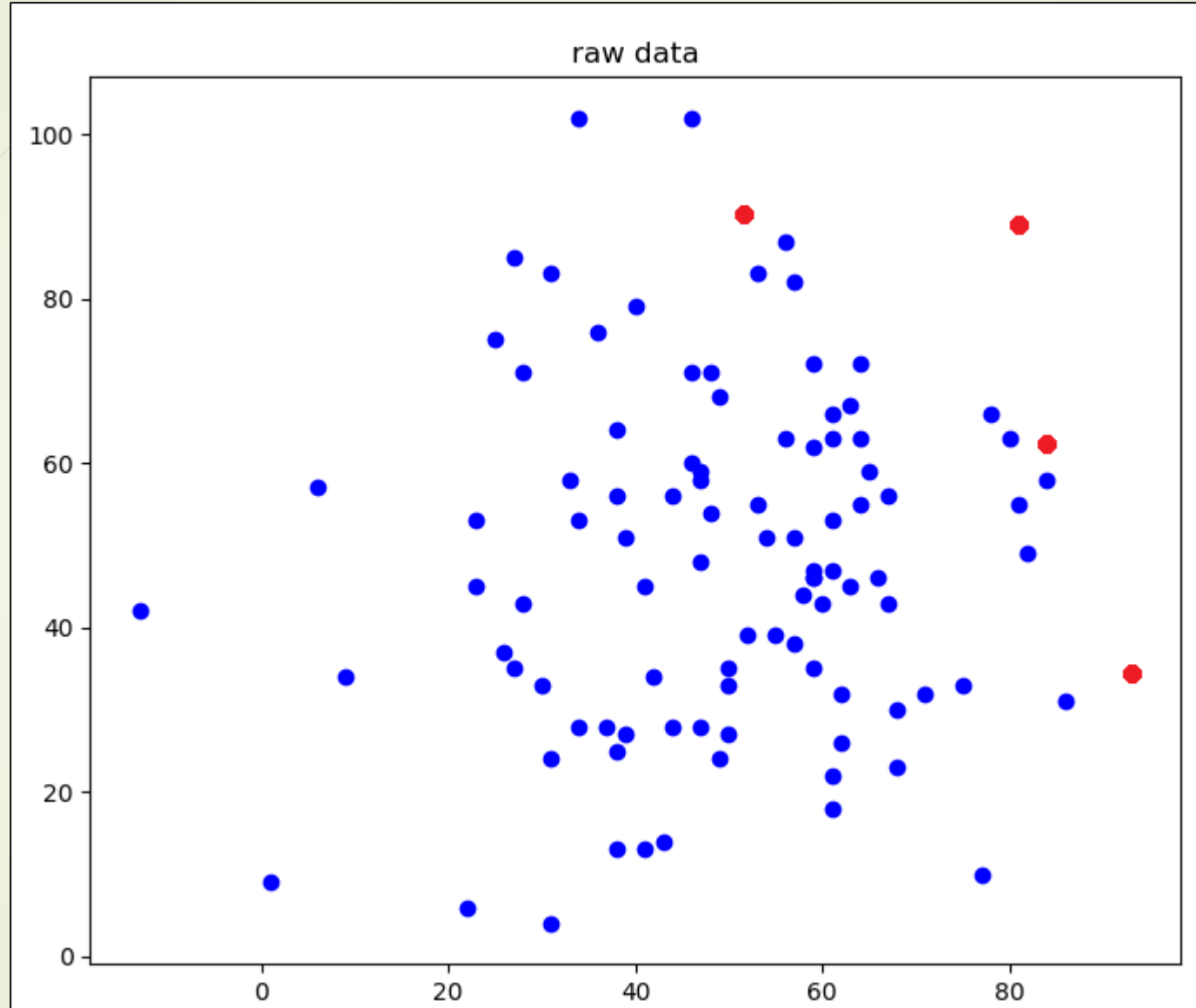
# Skyline Operator : Skyline Point 3

9



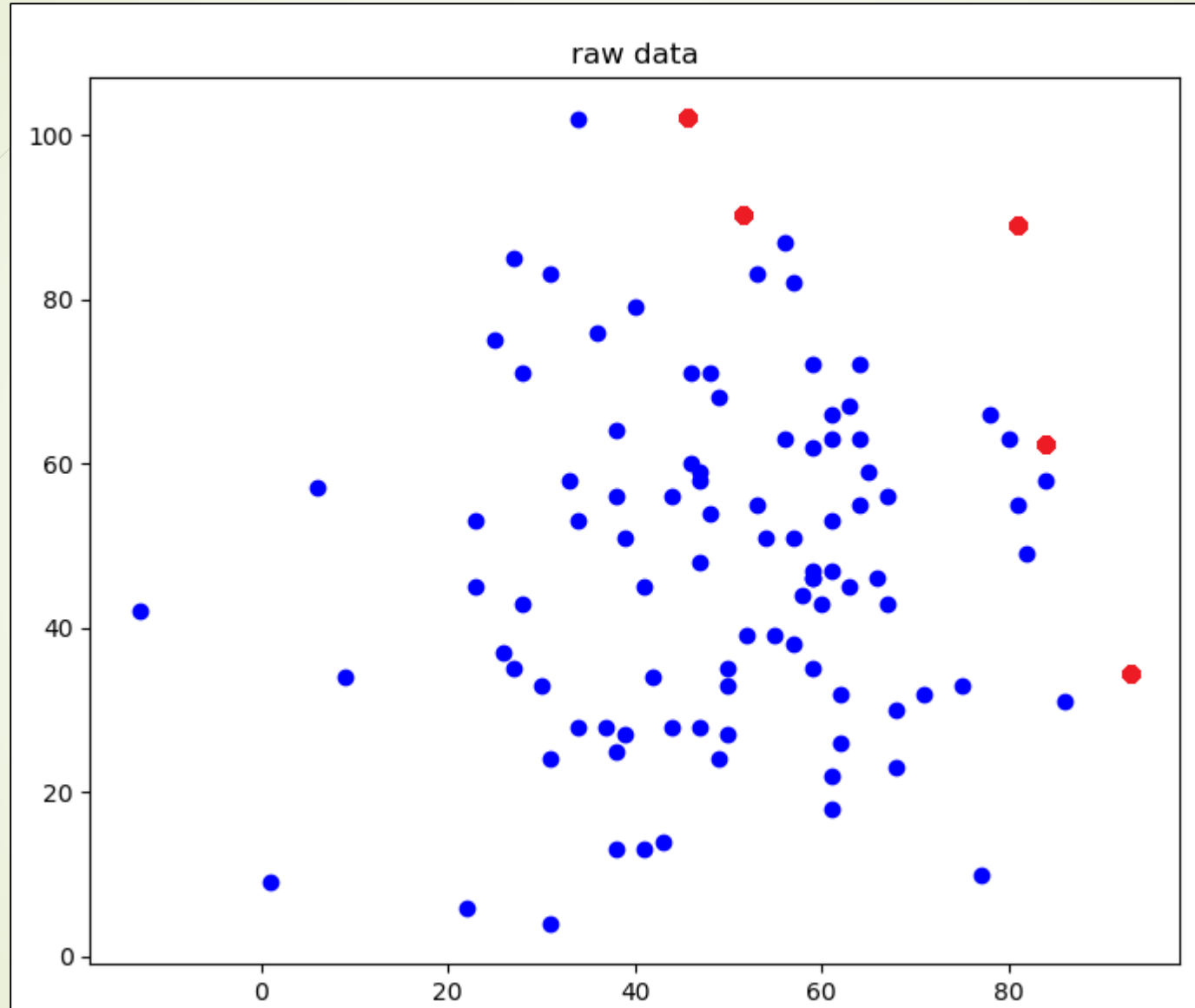
# Skyline Operator : Skyline Point 4

10



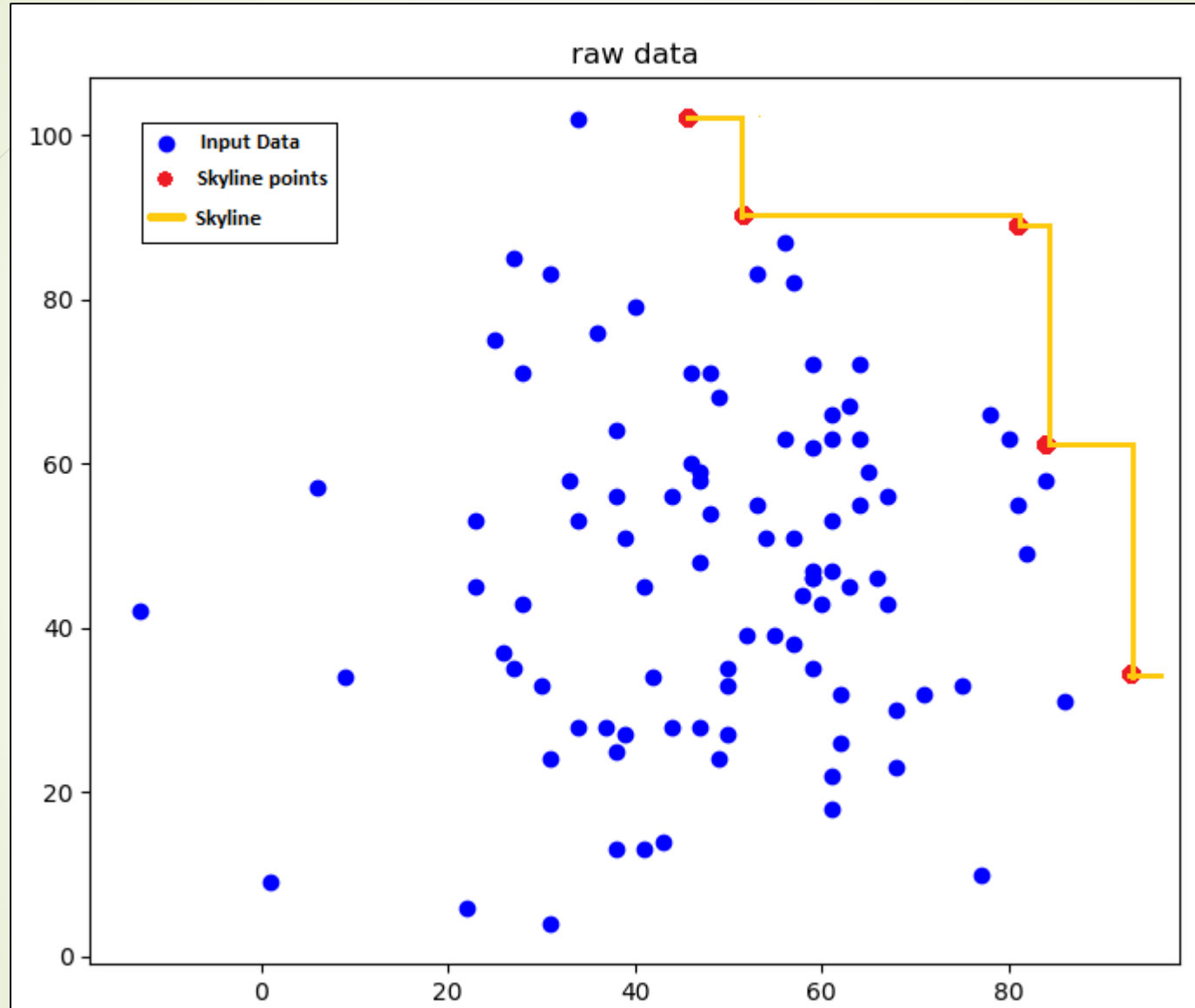
# Skyline Operator : Skyline Point 5

11



# Skyline Operator : Skyline

12



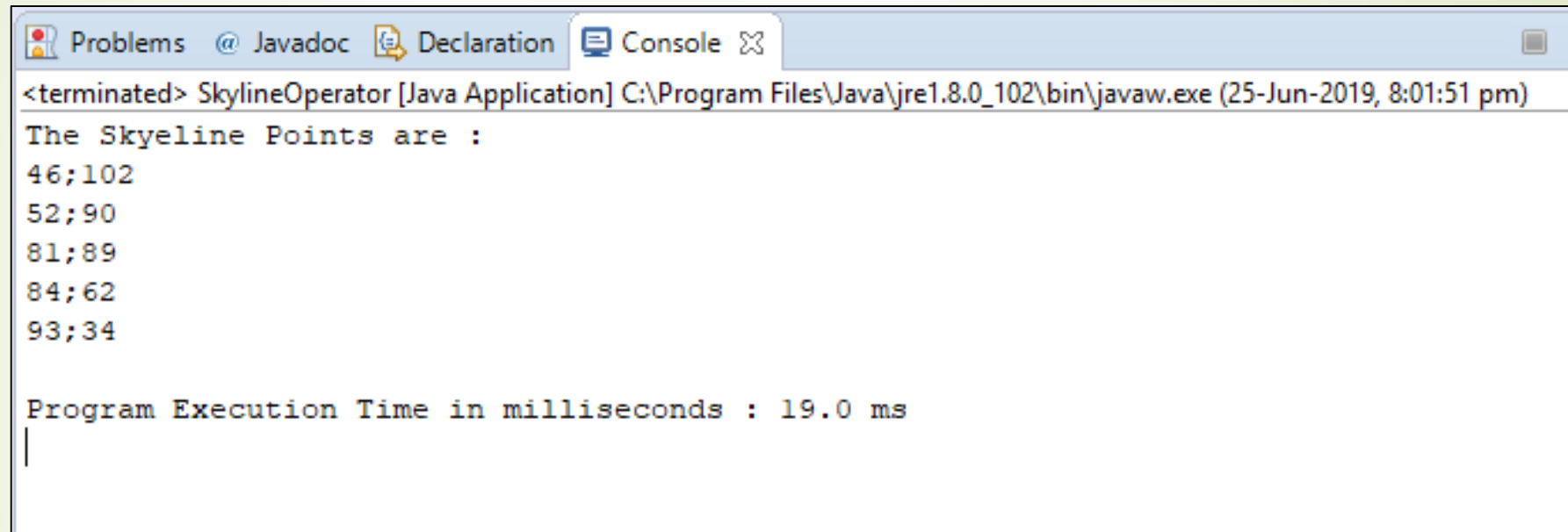
# Skyline Operator : Implementation (Sorting)

```
for(int i = 0; i<Data.size(); i++)
{
    for(int j=i+1;j<Data.size();j++)
    {
        if(Integer.parseInt(Data.get(i).get(0)) == Integer.parseInt(Data.get(j).get(0)))
        {
            // Remove Duplicate Points
            if(Integer.parseInt(Data.get(i).get(1)) == Integer.parseInt(Data.get(j).get(1)))
            {
                System.out.println("Removed Data Point : " + Data.get(j));
                Data.remove(j);
            }
            // If X value is same, Sort using Y value
            else if(Integer.parseInt(Data.get(i).get(1)) < Integer.parseInt(Data.get(j).get(1)))
            {
                tempdata = Data.get(i);
                Data.set(i, Data.get(j));
                Data.set(j, tempdata);
            }
        }
        else if(Integer.parseInt(Data.get(i).get(0)) < Integer.parseInt(Data.get(j).get(0)))
        {
            tempdata = Data.get(i);
            Data.set(i, Data.get(j));
            Data.set(j, tempdata);
        }
    }
}
```

# Skyline Operator : Implementation

```
Max_Y = Integer.parseInt(Data.get(0).get(1));
Output.add(Data.get(0));
for(int i = 1; i<Data.size(); i++)
{
    if(Max_Y < Integer.parseInt(Data.get(i).get(1))) {
        Max_Y = Integer.parseInt(Data.get(i).get(1));
        Output.add(Data.get(i));
    }
}
System.out.println(Data.size());
return Output;
```

# Skyline Operator : Output



The screenshot shows a Java IDE console window with the following content:

```
<terminated> SkylineOperator [Java Application] C:\Program Files\Java\jre1.8.0_102\bin\javaw.exe (25-Jun-2019, 8:01:51 pm)
The Skyline Points are :
46;102
52;90
81;89
84;62
93;34

Program Execution Time in milliseconds : 19.0 ms
|
```

## CHALLENGE 2: Entity Resolution



# Entity Resolution : Overview

- Entity resolution is used to identify duplicates within a single or among two datasets using similarity functions

# Entity Resolution : Goal

Objective	Perform entity resolution on the CORA data set. The program must identify and return duplicate pairs of publication, author and venue entities.
Input	CORA Dataset in XML format <a href="https://hpi.de/fileadmin/user_upload/fachgebiete/naumann/projekte/repeatability/CORA/cora-all-id.xml">https://hpi.de/fileadmin/user_upload/fachgebiete/naumann/projekte/repeatability/CORA/cora-all-id.xml</a>
Output	duplicates_publications.csv – List of all publications duplicates duplicates_authors.csv – List of all authors duplicates duplicates_venues.csv – List of all venues duplicates Ex: ahlskog1994a;ahlskog1994a
Constraints	1. Id fields shall not be used to solve this task 2. Efficiency : accuracy of dataset in terms of the F-measure.
Programming Language	Python 3.6

# Entity Resolution : Input CORA dataset

19

```
<coraRADD>
  <publication id="ahlskog1994a">
    <author id="199">M. Ahlskog</author>
    <author id="74"> J. Paloheimo</author>
    <author id="64"> H. Stubb</author>
    <author id="103"> P. Dyreklev</author>
    <author id="54"> M. Fahlman</author>
    <title>Inganas</title>
    <title>and</title>
    <title>M.R.</title>
  <venue>
    <venue pubid="ahlskog1994a" id="1">
      <name>Andersson</name>
      <name> J Appl. Phys.</name>
      <vol>76</vol>
      <date> (1994). </date>
    </venue>
  </venue>
</publication>
  <publication id="ahlskog1994a">
    <author id="199">M. Ahlskog</author>
    <author id="74"> J. Paloheimo</author>
    <author id="64"> H. Stubb</author>
    <author id="103"> P. Dyreklev</author>
    <author id="54"> M. Fahlman</author>
    <author id="101"> O. Inganas and M.R. Andersson</author>
  <venue>
    <venue pubid="ahlskog1994a" id="1">
      <name>J Appl. Phys.</name>
      <vol>76</vol>
      <date> (1994). </date>
    </venue>
  </venue>
</publication>
  <publication id="ahlskog1994a">
```

# Entity Resolution : Algorithm

- Data Collection – Get the structured CORA data from the provided URL link (XML file)
- Data Extraction – Extract the XML data into list as an usable format
- Data Merging – Merge meaningful data which are separated into multiple tags

```
<title>Provably</title>  
<title>correct</title>  
<title>compiler</title>  
<title>development</title>  
<title>and</title>  
<title>implementation.</title>
```



“Provably correct compiler development and implementation.”

```
<venue pubid="ahlskog1994a" id="1">  
  <name>J Appl. Phys.</name>  
  <vol>76</vol>  
  <date> (1994). </date>  
</venue>  
</venue>
```



“J Appl. Phys. 76 (1994).”

- Identify the similarity between the each elements.
  - Levenshtein Distance (Edit Distance)
- Add the elements which have high similarity to the corresponding CSV files

# Entity Resolution : Data Collection

```
def GetWebData(WebURL) :
    file = urllib.request.urlopen(WebURL)
    data = file.read()
    file.close()

    data = xmltodict.parse(data)
    return data
```

```
▼<coraRADD>
  ▼<publication id="ahlskog1994a">
    <author id="199">M. Ahlskog</author>
    <author id="74"> J. Paloheimo</author>
    <author id="64"> H. Stubb</author>
    <author id="103"> P. Dyreklev</author>
    <author id="54"> M. Fahlman</author>
    <title>Inganas</title>
    <title>and</title>
    <title>M.R.</title>
  ▼<venue>
    ▼<venue pubid="ahlskog1994a" id="1">
      <name>Andersson</name>
      <name> J Appl. Phys.</name>
      <vol>76</vol>
      <date> (1994). </date>
    </venue>
  </venue>
</publication>
```

```
OrderedDict([('@id', 'ahlskog1994a'),
('author', [OrderedDict([('@id', '199'), ('#text',
'M. Ahlskog')]), OrderedDict([('@id', '74'),
('#text', 'J. Paloheimo')]),
OrderedDict([('@id', '64'), ('#text', 'H.
Stubb')]), OrderedDict([('@id', '103'), ('#text',
'P. Dyreklev')]), OrderedDict([('@id', '54'),
('#text', 'M. Fahlman')])]), ('title', ['Inganas',
'and', 'M.R.']), ('venue',
OrderedDict([('venue',
OrderedDict([('@pubid', 'ahlskog1994a'),
('@id', '1'), ('name', ['Andersson', 'J Appl.
Phys.']), ('vol', '76'), ('date', '(1994).')]))]))])
```

# Entity Resolution : Data Extraction and Merging

```
for Publication in InputData['coraRADD']['publication']:

    tempAut, tempID = GetInstanceAuthor(Publication, 'author', '#text', '@id')
    PubIdList.append(Publication['@id'])
    TempTitleList.append(GetInstanceTitle(Publication, 'title'))
    VenNameList.append(GetInstanceVenue(Publication, 'venue', 'name'))
    VenIdList.append(GetInstanceVenue(Publication, 'venue', '@id'))
    VenDateList.append(GetInstanceVenue(Publication, 'venue', 'date'))

    AuthorList.append(tempAut)
    AutID.append(tempID)

VenIdList, VenNameList, VenDateList = MergeVenueName(VenIdList, VenNameList, VenDateList)
```



# Entity Resolution : Title List

```
def GetInstanceTitle(Pub, Param):
    OutputList = []
    if Param in Pub.keys():
        for Title in Pub[Param]:
            OutputList.append(Title)
    else:
        OutputList.append('NAN')

    return OutputList
```

```
for Title in TempTitleList:
    if Title != ['NAN']:
        tempstr = str()
        for Entry in Title:
            tempstr = tempstr + Entry + " "
        TitleList.append(tempstr)
    else:
        TitleList.append('NAN')
```



```
Inganas and M.R.
NAN
NAN
NAN
NAN
NAN
NAN
NAN
Robots and Manufacturing Automation.
A spatial model of interaction in large virtual environments.
Viewpoints, Actionpoints and Spatial Frames for Collaborative User Interfaces,
User Embodiment in Collaborative Virtual Environments,
User Embodiment in Collaborative Virtual Environments.
User Embodiment in Collaborative Virtual Environments.
Networked Virtual reality and Cooperative Work.
Actress: an action semantics directed compiler generator,
Actress: an action semantics directed compiler generator,
an action semantics directed compiler generator.
Provably correct compiler development and implementation.
Provably Correct Compiler Implementation,
```

# Entity Resolution : Author List

```
def GetInstanceAuthor(Pub, Param, Key, Idkey):
    OutputList = []
    TempList = []
    tempstr = str()
    KeyList = []
    if Param in Pub.keys():
        if(isinstance(Pub[Param],list) == False):
            OutputList = Pub[Param][Key]
            KeyList = Pub[Param][Idkey]
        else:
            for Aut in Pub[Param]:
                if(isinstance(Aut,str) == False):
                    OutputList = Aut[Key]
                    KeyList = Aut[Idkey]
            else:
                OutputList = 'NAN'
                KeyList = 'NAN'

    return OutputList,KeyList
```



M. Fahlman  
 O. Inganas and M.R. Andersson  
 O. Inganas and M.R. Andersson  
 O. Inganas and M.R. Andersson  
 O. Inganas and M.R. Andersson  
 O. Inganas and M.R. Andersson  
 O. Inganas and M.R. Andersson  
 O. Inganas and M.R. Andersson  
 C. Ray Asfahl.  
 Steve Benford and Lennart E. Fahlen.  
 and Fahn  
 and Snowdon  
 Snowdon  
 Snowdon  
 and Tom Rodden.  
 H. and Watt  
 H. and Watt  
 and D. A. Watt. Actress:  
 and M. Muller-Olm.  
 B. Buth et. al.



# Entity Resolution : Venue Instance

```
def GetInstanceVenue (Pub, Param, Key) :
    OutputList = []
    if Param in Pub.keys() :
        if Key in Pub[Param][Param].keys() :
            OutputList.append(Pub[Param][Param][Key])
        else:
            OutputList.append('NAN')
    return OutputList
```

```
['(1994).']
['(1994).']
['(1994).']
['(1994).']
['(1994).']
['(1994).']
['(1994).']
['(1994).']
['(1994).']
['1992.']
[['September', '1993.']]
[['(1994),', 'June 1994,']]
['May 7-11, 1995,']
['(1995).']
[]
['1995.']
['(1992b),']
['(1992b),']
['October 1992.']
['1992.']
['1992,']
```

```
[['Andersson', 'J Appl. Phys.']]
['J Appl. Phys.']
['J Appl. Phys.']
['J Appl. Phys.']
['J Appl. Phys.']
['J Appl. Phys.']
['J Appl. Phys.']
['J Appl. Phys.']
['Journal of Applied Physics']
[]
["In Proceedings of ECSCW'93"]
['6th ERCIM workshop']
['in Proc. ACM Conference on Human Factors in Computing Systems (CHI95)']
['In Proceedings of CHI95']
['In Proceedings of CHI95']
['Presence']
["`Proceedings of the International Workshop on Compiler Construction (CC-92)'"']
["`Proceedings of the International Workshop on Compiler Construction (CC-92)'"']
["Proceedings of the 4th International Conference on Compiler Construction (CC'92)"]
['Compiler Construction']
['Compiler Construction']
```

# Entity Resolution : Venue Name and Date

26

```
def MergeVenueName(VenIdList, VenNameList, VenDateList):
    TempId = []
    TempName = []
    TempDate = []
    |
    for Id in VenIdList:
        for Entry in Id:
            TempId.append(int(Entry))

    VenNameList = ProcessNameDate(VenNameList)
    VenDateList = ProcessNameDate(VenDateList)

    return TempId, VenNameList, VenDateList
```

```
def ProcessNameDate(VenNameList):
    TempName = []
    TempDate = []
    TempStr = 'NAN'
    TempStrDate = 'NAN'
    for Name in VenNameList:
        if(len(Name) == 0):
            TempName.append('NAN')
        else:
            for NameEntry in Name:
                if isinstance(NameEntry, str) == True:
                    TempName.append(NameEntry)
                else:
                    TempStr = str()
                    for Index in NameEntry:
                        TempStr = TempStr + Index + " "
                    TempName.append(TempStr)

    return TempName
```

```
Andersson J Appl. Phys.
J Appl. Phys.
J Appl. Phys.
J Appl. Phys.
J Appl. Phys.
J Appl. Phys.
J Appl. Phys.
Journal of Applied Physics
NAN
In Proceedings of ECSCW'93
6th ERCIM workshop
in Proc. ACM Conference on Human Factors in Computing Systems (CHI95)
In Proceedings of CHI95
In Proceedings of CHI95
Presence
'Proceedings of the International Workshop on Compiler Construction (CC-92)'
'Proceedings of the International Workshop on Compiler Construction (CC-92)'
Proceedings of the 4th International Conference on Compiler Construction (CC'92)
Compiler Construction
Compiler Construction
```

```
(1994).
(1994).
(1994).
(1994).
(1994).
(1994).
(1994).
1992.
September 1993.
(1994), June 1994,
May 7-11, 1995,
(1995).
NAN
1995.
(1992b),
(1992b),
October 1992.
1992.
1992,
```

# Entity Resolution : Venue Name and Date

```
def MergeVenueData(VenNameList, VenDateList):
    for i in range(0, len(VenNameList)):
        if VenNameList[i] == 'NAN':
            VenNameList[i] = VenDateList[i]
    return VenNameList
```

```
Andersson J Appl. Phys.
J Appl. Phys.
J Appl. Phys.
J Appl. Phys.
J Appl. Phys.
J Appl. Phys.
J Appl. Phys.
Journal of Applied Physics
1992.
In Proceedings of ECSCW'93
6th ERCIM workshop
in Proc. ACM Conference on Human Factors in Computing Systems (CHI95)
In Proceedings of CHI95
In Proceedings of CHI95
Presence
`Proceedings of the International Workshop on Compiler Construction (CC-92)'
`Proceedings of the International Workshop on Compiler Construction (CC-92)'
Proceedings of the 4th International Conference on Compiler Construction (CC'92)
Compiler Construction
Compiler Construction
```

# Entity Resolution : Similarity measure

```
def SimilarityMeasure(InputStr, ID, AvgLen):
    MatchList = []
    for i in range(0, len(InputStr)):
        for j in range(i+1, len(InputStr)):
            if InputStr[i] != 'NAN' and InputStr[j] != 'NAN':
                Sim = SequenceMatcher(None, InputStr[i], InputStr[j]).ratio()
                AverageLen = (len(InputStr[i]) + len(InputStr[j])) / 2
                if AverageLen < AvgLen:
                    if Sim > 0.9: # Similarity threshold is set as 0.9 for smaller strings
                        MatchList.append([ID[i], ID[j], InputStr[i], InputStr[j], Sim])
                else:
                    if Sim > 0.7: # Similarity threshold is set as 0.7 for Larger strings
                        MatchList.append([ID[i], ID[j], InputStr[i], InputStr[j], Sim])
            elif InputStr[i] == 'NAN' and InputStr[j] == 'NAN':
                MatchList.append([ID[i], ID[j], InputStr[i], InputStr[j], Sim])
            else:
                # Do nothing
                None
    return MatchList
```

```
def LD(s, t):
    if s == "":
        return len(t)
    if t == "":
        return len(s)
    if s[-1] == t[-1]:
        cost = 0
    else:
        cost = 1

    res = min([LD(s[:-1], t)+1,
               LD(s, t[:-1])+1,
               LD(s[:-1], t[:-1]) + cost])
    return res
```

# Entity Resolution : Duplicate Publications

fahlman1988b	fahlman1988b	An empirical study of learning speed in backpropagation networks.	An empirical study of learning speed in back-propagation networks.	0.992481
--------------	--------------	---	--	----------

brodley1992	brodley1992	Multivariate Versus Univariate Decision Trees.	Multivariate Versus Univariate Decision Trees.	1
brodley1992	brodley1992b	Multivariate Versus Univariate Decision Trees.	Multivariate decision trees.	0.710526
brodley1992	brodley1992b	Multivariate Versus Univariate Decision Trees.	Multivariate decision trees.	0.710526
brodley1992	brodley1992b	Multivariate Versus Univariate Decision Trees.	Multivariate decision trees.	0.710526
brodley1992	brodley1992b	Multivariate Versus Univariate Decision Trees.	Multivariate decision trees.	0.710526

utgoff1982aaai	utgoff1982aaai	Acquisition of appropriate bias for inductive concept learning.	Acquisition of appropriate bias for inductive concept learning.	1
utgoff1982aaai	utgoff1984phd	Acquisition of appropriate bias for inductive concept learning.	Shift of bias for inductive concept learning.	0.781818
utgoff1982aaai	utgoff1984phd	Acquisition of appropriate bias for inductive concept learning.	Shift of bias for inductive concept learning.	0.781818

# Entity Resolution : Duplicate Authors

54	54	M. Fahlman	S. Fahlman	0.9
54	54	M. Fahlman	S. Fahlman	0.9
54	54	M. Fahlman	Fahlman	0.823529
54	54	M. Fahlman	Fahlman	0.823529
54	54	M. Fahlman	S. E. Fahlman.	0.75
54	54	M. Fahlman	S. E. Fahlman.	0.75
54	54	M. Fahlman	S.E. Fahlman.	0.782609

117	21	S. and Lebiere	C. Lebiere:	0.72
117	117	S. and Lebiere	S. Fahlman and C. Lebiere	0.717949
117	117	S. and Lebiere	S. Fahlman and C. Lebiere	0.717949
117	117	S. and Lebiere	S. Fahlman and C. Lebiere	0.717949

173	173	D. Aha and D. Kibler.	David W. Aha and Dennis Kibler.	0.769231
173	173	D. Aha and D. Kibler.	D.W. Aha and D. Kibler.	0.954545
173	28	D. Aha and D. Kibler.	D. and Kibler	0.764706
173	173	D. Aha and D. Kibler.	D. Aha and D. Kibler.	1
173	173	D. Aha and D. Kibler.	D. Aha and D. Kibler.	1



# Entity Resolution : Duplicate Venues

1	1	Andersson J Appl. Phys.	J Appl. Phys.	0.702703
1	1	Andersson J Appl. Phys.	J Appl. Phys.	0.702703
1	1	Andersson J Appl. Phys.	J Appl. Phys.	0.702703
1	1	Andersson J Appl. Phys.	J Appl. Phys.	0.702703

3	4	In Proceedings of ECSCW'93	In Proceedings of CHI95	0.816327
3	4	In Proceedings of ECSCW'93	In Proceedings of CHI95	0.816327
3	20	In Proceedings of ECSCW'93	In Proceedings of Interchi '93	0.75
3	20	In Proceedings of ECSCW'93	In Proceedings of Interchi '93	0.75
3	20	In Proceedings of ECSCW'93	In Proceedings of Interchi '93	0.75
3	20	In Proceedings of ECSCW'93	In Proceedings of Interchi '93	0.75

6	6	'Proceedings of the International Workshop on Compiler Construction (CC-92)'	in Proceedings of the Fourth International Conference on Compiler Construction	0.779221
6	6	'Proceedings of the International Workshop on Compiler Construction (CC-92)'	in Proceedings of the Fourth International Conference on Compiler Construction	0.779221
6	6	'Proceedings of the International Workshop on Compiler Construction (CC-92)'	in Proceedings of the Fourth International Conference on Compiler Construction	0.779221

13	14	Proceedings of the Royal Society London B (in press).	Proceedings of the Royal Society of London B	0.857143
13	14	Proceedings of the Royal Society London B (in press).	Proceedings of the Royal Society of London B	0.857143

22	22	1979	1979	1
22	22	1979	1979	1

# Entity Resolution : Current and Future Adaptations

- Similarity Threshold was modified to based on the average length of the comparing string to solve below mentioned issue.
  - 1992 and 1994 - 75% similarity
  - Hence the similarity threshold was modified to 90% for strings of length less than 5.
- The following characters and words shall be removed before comparing the strings to get a better comparison measure.
  - The special characters - { “ , ‘ ( ) } { [ ] } .
  - Words like The, of , and , is
  - Initials, Salutation, Degrees
- The characters in uppercase shall be converted to lowercase
- Some Special processing is needed depending data.
  - Proceedings of the **Eighth** International Workshop on Machine Learning
  - Proceedings of the **Ninth** International Workshop on Machine Learning
  - In Proceedings of the **4 th** International Workshop on Machine Learning





33

## CHALLENGE 3: Data Cleaning

# Data Cleaning : Overview

- During data cleaning, errors in data such as illegal or wrong values are resolved to improve the quality of the given data.
- **OpenRefine** is a power tool for working with messy data. Use it to improve data consistency, link it to data registries like Wikidata, augment it with data from other sources, transform it into different formats for other tools to consume, and contribute it back to the original sources.
- OpenRefine is not a web service but a desktop app that runs on your own computer, so you can process sensitive data with privacy.
- OpenRefine was originally developed as "Freebase Gridworks" by Metaweb Technologies, Inc.. Metaweb was acquired by Google in July 2010 and they renamed the product Google Refine.
- In October, 2012, the product was renamed OpenRefine as it transitioned to a community supported project.

# Data Cleaning : Goal

35

Objective	<p>Clean the data in the file.</p> <p>During the collection of the data, wrong, incomplete, and duplicate data have been introduced which should be resolved by according to the given cleaning requirements.</p>
Input	bibliography.csv
Output	Cleaned OpenRefine Project Zip folder.
Constraints	<ol style="list-style-type: none"><li>1. The cleaned data set should contain only correct publication entries.</li><li>2. all false, nonsense or redundant information should be removed from the data set.</li></ol> <p><b>Article:</b> A peer-reviewed article from a journal or magazine</p> <ul style="list-style-type: none"><li>• Author, title, journal, year</li></ul> <p><b>Book:</b> A printed and bound book with explicit publisher</p> <ul style="list-style-type: none"><li>• Title, publisher, year, author, editor</li></ul> <p><b>Incollection</b> : Published as part of a printed and bound book with own title</p> <ul style="list-style-type: none"><li>• Author, title, booktitle, year</li></ul> <p><b>Inproceedings</b> : An article published in a conference proceedings</p> <ul style="list-style-type: none"><li>• Author, title, booktitle, year</li></ul> <p><b>Misc:</b> Entry type that can be used if no other entry type fits.</p>
Tool	OpenRefine ( <a href="http://openrefine.org">http://openrefine.org</a> )



OpenRefine A power tool for working with messy data.

Create Project « Start Over Configure Parsing Options Project name bibliography csv Tags Create Project »

BibliographyType,ISBN,Identifier,Author,Title,Journal,Volume,Number,Month,Pages,Year,Address,Note,URL,Booktitle,Chapter,Edition,Series,Editor,Publisher,ReportType,Howpublished,Institution,Organ

- Article,"1","Sarah Holland and Ahmed Hosny and Sarah Newman and Joshua Joseph and Kasia Chmielinski","The Dataset Nutrition Label: A Framework To Drive Higher Data Quality Standards","May",2018,"incomplete or biased data will often exhibit problematic outcomes. Current methods of data analysis, particularly before model development, are costly and not standardized. The Dataset Nutrition Label (the Label) is a analysis by providing a distilled yet comprehensive overview of dataset "ingredients" before AI model development. Building a Label that can be applied across domains and data types requires that the framework itse diverse qualitative and quantitative modules generated through multiple statistical and probabilistic modelling backends, but displayed in a standardized format. To demonstrate and advance this concept, we generate modules on the ProPublica Dollars for Docs dataset. The benefits of the Label are manifold. For data specialists, the Label will drive more robust data analysis practices, provide an efficient way to select the best data models as a result of more robust training datasets and the ability to check for issues at the time of model development. For those building and publishing datasets, the Label creates an expectation of explanation, whi limitations of the Label, including the challenges of generalizing across diverse datasets, and the risk of using "ground truth" data as a comparison dataset. We discuss ways to move forward given the limitations identi Label project, including research and public policy agendas to further advance consideration of the concept.",,"cs.DB, cs.CY","http://arxiv.org/pdf/1805.03677v1",
- Article,"1050601491","10","Mancuhan, Koray and Clifton, Chris","Combating discrimination using Bayesian networks","Artificial Intelligence and Law",22,2,"211-238",2014,"Discrimination in deci etc{ldots}), but often present in historical decisions. Use of such discriminatory historical decision making as training data can perpetuate discrimi-nation, even if the protected attributes are not directly present in the d and preventing discrimination in classification. First, we propose a discrimination discovery method based on mod-eling the probability distribution of a class using Bayesian networks. This measures the effect of a prot estimated probability distribution (via a Bayesian network). Second, we propose a classification method that corrects for the discovered discrimination without using protected attributes in the decision process. We eva approaches on two different datasets. The empirical results show that a substantial amount of discrimination identified in instances is prevented in future decisions." "Bayesian network Data mining Discrimination disc

Parse data as

CSV / TSV / separator-based files

Line-based text files

Fixed-width field text files

PC-Axis text files

JSON files

MARC files

JSON-LD files

RDF/N3 files

RDF/N-Triples files

RDF/Turtle files

Character encoding

Columns are separated by

☒ commas (CSV)

☐ tabs (TSV)

☐ custom: \t

Escape special characters with \

☐ Column names (comma separated):

☐ Ignore first 0 line(s) at beginning of file

☒ Parse next 1 line(s) as column headers

☐ Discard initial 0 row(s) of data

☐ Load at most 0 row(s) of data

☒ Use character " to enclose cells containing column separators

☐ Parse cell text into numbers, dates, ...

☒ Store blank rows

☒ Store blank cells as nulls

☐ Store file source (file names, URLs) in each row

Update Preview

Version 3.1 [b90e413]

Preferences

Help

About

**OpenRefine** *A power tool for working with messy data.*

Create Project « Start Over Configure Parsing Options Project name  Tags  Create Project »

Open Project  
Import Project  
Language Settings

	Bibliography type	ISBN	Identifier	Author	Title	Journal	Volume	Number	Month	Pages	Year	Address	Note	URL
1.	Article		1	Sarah Holland and Ahmed Hosny and Sarah Newman and Joshua Joseph and Kasia Chmielinski	The Dataset Nutrition Label: A Framework To Drive Higher Data Quality Standards				May		2018			

Parse data as

**CSV / TSV / separator-based files**

- Line-based text files
- Fixed-width field text files
- PC-Axis text files
- JSON files
- MARC files
- JSON-LD files
- RDF/N3 files
- RDF/N-Triples files
- RDF/Turtle files

Character encoding

Columns are separated by

- ☒ commas (CSV)
- ☐ tabs (TSV)
- ☐ custom: \t

Escape special characters with \

☐ Column names (comma separated):

☐ Ignore first 0 line(s) at beginning of file

☒ Parse next 1 line(s) as column headers

☐ Discard initial 0 row(s) of data

☐ Load at most 0 row(s) of data

☒ Use character " to enclose cells containing column separators

☐ Parse cell text into numbers, dates, ...

☒ Store blank rows

☒ Store blank cells as nulls

☐ Store file source (file names, URLs) in each row

Update Preview

Version 3.1 [b90e413]

Preferences  
Help  
About



**OpenRefine** bibliographycsv [Permalink](#) Open... Export ▾ Help

Facet / Filter Undo / Redo 0 / 0 112 rows Extensions: Wikidata ▾

Show as: rows records Show: 5 10 25 50 rows « first < previous 1 - 10 next > last »

**Using facets and filters**

Use facets and filters to select subsets of your data to act on. Choose facet and filter methods from the menus at the top of each data column.

Not sure how to get started?  
[Watch these screencasts](#)

▼ All	▼ BibliographyType	▼ ISBN	▼ Identifier	▼ Author	▼ Title	▼ Journal	▼ Volume	▼ Number	▼ Month	▼ Pages	▼ Year	▼ Address	▼ No
☆	1.	Article	1	Sarah Holland and Ahmed Hosny and Sarah Newman and Joshua Joseph and Kasia Chmielinski	The Dataset Nutrition Label: A Framework To Drive Higher Data Quality Standards				May		2018		

127.0.0.1:3333

OpenRefine interface showing a project with 112 rows. The interface includes a toolbar with buttons for Open..., Export, and Help. The main table displays columns: Journal, Volume, Number, Month, Pages, Year, Address, Note, and URL. A context menu is open over the 'Month' column, showing options like Facet, Text filter, Edit cells, Edit column, Transpose, Sort..., View, and Reconcile. The 'Facet' option is selected, and a sub-menu is visible with options: Text facet, Numeric facet, Timeline facet, Scatterplot facet, Custom text facet..., Custom Numeric Facet..., and Customized facets. The table content includes rows with data such as 'Artificial Intelligence and Law' and 'http://arxiv.org/abs/1412.3756'.

	Journal	Volume	Number	Month	Pages	Year	Address	Note	URL
ing nation yesian s	Artificial Intelligence and Law	22	2		211--238	2014			
g and g e impact					1--28	2014			<a href="http://arxiv.org/abs/1412.3756">http://arxiv.org/abs/1412.3756</a>



Facet / Filter
Undo / Redo 0 / 0

Refresh
Reset All
Remove All

Month
change

17 choices
Sort by: name count
Cluster

23--24 Feb 4  
apr 1  
aug 1  
August 1  
dec 6  
December 1  
jan 1  
jul 1  
jun 1  
mar 2  
March 1  
May 3  
may 3  
nov 1  
oct 2  
October 2  
September 1  
(blank) 80  
Facet by choice counts

112 rows

Show as: rows records
Show: 5 10 25 50 rows

	Journal	Volume	Number	Month	Page
asset Label: A ork To gher ality ds				May	
ing nation ayesian s	Artificial Intelligence and Law	22	2		211--238
g and g e impact					1--28

Google x bibliographycsv - OpenRefine x Bibliography - OpenRefine

127.0.0.1:3333/project?project=1988305065945

**OpenRefine** bibliographycsv [Permalink](#)

Facet / Filter Undo / Redo 0 / 0

Refresh Reset All Remove All

112 rows

Show as: rows records Show: 5 10 25 50 rows

Month change

17 choices Sort by: name count Cluster

23--24 Feb 4

apr 1

aug 1

August 1

dec 6 [edit](#) [include](#)

December 1

jan 1

jul 1

jun 1

mar 2

March 1

May 3

may 3

nov 1

oct 2

October 2

September 1

(blank) 80

Facet by choice counts

	Journal	Volume	Number	Month	
asset Label: A ork To gher ality ds				May	
ing nation yesian s	Artificial Intelligence and Law	22	2		211
g and g e impact					1--2

javascript:()

← → ↻ 127.0.0.1:3333/project?project=1988305065945

**OpenRefine** bibliographycsv [Permalink](#)

Facet / Filter Undo / Redo 0 / 0

Refresh Reset All Remove All

112 rows

Show as: rows records Show: 5 10 25 50 rows

Journal Volume Number Month Pages

Month change

17 choices Sort by: name count Cluster

23--24 Feb 4

apr 1

aug 1

August 1

dec 6

December 1

jan 1

jul 1

jun 1

mar 2

March 1

May 3

may 3

nov 1

oct 2

October 2

September 1

(blank) 80

Facet by choice counts

dec

Apply Cancel

Enter Esc

Label: A ork To gher ality ds				May	
ing nation ayesian s	Artificial Intelligence and Law	22	2		211--238
g and g e impact					1--28

Month change

11 choices Sort by: name count Cluster

April 1

August 2

December 7

February 3

January 1

July 1

March 3

May 6

November 1

October 4

September 1

(blank) 75

Facet by choice counts

26 choices Sort by: name count Cluster

2017	8
2011	7
2013	5
2008	4
2010	4
2009	3
2015	3
14	2
16	2
17	2
18	2
2007	2
06	1
10	1
13	1
1986	1
1988	1
2003	1
2019	1
88	1
95	1
97	1
(blank)	10

Facet by choice counts

Google x bibliographycsv - OpenRefine x Bibliography - OpenRefine x ILIAS für Lehre & Lernen – Univer x +

127.0.0.1:3333/project?project=1988305065945

**OpenRefine** bibliographycsv [Permalink](#)

Facet / Filter Undo / Redo 0 / 0

Refresh Reset All Remove All Show as: rows records Show: 5 10 25 50 rows

26 choices Sort by: name count Cluster

	Journal	Volume	Number	Month	Pages	Year	Address	Note	URL
aset Label: A ork To gher ality ds				May					
ing nation ayesian s	Artificial Intelligence and Law	22	2		211--238	2014			
g and g e impact					1--28	2014			<a href="http://arxiv.org/abs/1412">http://arxiv.org/abs/1412</a>

Facet by choice counts

javascript:[]

Context menu for 'Edit cells':

- Facet
- Text filter
- Edit cells**
  - Transform...**
  - Common transforms
    - Fill down
    - Blank down
    - Split multi-valued cells...
    - Join multi-valued cells...
  - Cluster and edit...
  - Replace
- Edit column
- Transpose
- Sort...
- View
- Reconcile

```
if value == None:  
    return value  
elif int(value) >=0 and int(value) <20:  
    value = str(int(value)+2000)  
    return value  
elif int(value) >19 and int(value) <99:  
    value = str(int(value)+1900)  
    return value  
else:  
    return value
```

## Custom text transform on column Year

Expression

Language Python / Jython

```
if value == None:  
    return value  
elif int(value) >=0 and int(value) <20:  
    value = str(int(value)+2000)  
    return value
```

No syntax error.

Preview

History

Starred

Help

37.	2013	2013
38.	2003	2003
39.	97	1997
40.	2016	2016
41.	1988	1988
42.	88	1988
43.	10	2010
44.	null	null

On error

- ☒ keep original
- ☐ set to blank
- ☐ store error

☐ Re-transform up to 10 times until no change

OK

Cancel



26 choices	Sort by: name	count	Cluster
2017	8		
2011	7		
2013	5		
2008	4		
2010	4		
2009	3		
2015	3		
14	2		
16	2		
17	2		
18	2		
2007	2		
06	1		
10	1		
13	1		
1986	1		
1988	1		
2003	1		
2019	1		
88	1		
95	1		
97	1		
(blank)	10		
Facet by choice counts			

19 choices	Sort by: name	count	Cluster
2018	18		
2012	13		
2016	12		
2014	11		
2017	10		
2011	7		
2013	6		
2010	5		
2008	4		
2009	3		
2015	3		
1988	2		
2007	2		
1986	1		
1995	1		
1997	1		
2003	1		
2006	1		
2019	1		
(blank)	10		
Facet by choice counts			

This feature helps you find groups of different cell values that might be alternative representations of the same thing. For example, the two strings "New York" and "new york" are very likely to refer to the same concept and just have capitalization differences, and "Gödel" and "Godel" probably refer to the same person. [Find out more ...](#)

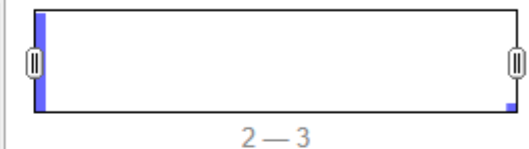
Method key collision

Keying Function metaphone3

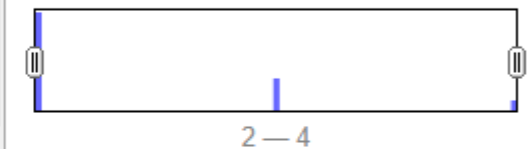
13 clusters found

Cluster Size	Row Count	Values in Cluster	Merge?	New Cell Value
3	4	<ul style="list-style-type: none"> <li>Kamiran, Faisal and Calders, Toon (2 rows)</li> <li>Kamiran, Faisal and Karim, Asim and Zhang, Xiangliang (1 rows)</li> <li>Kamiran, Faisal and {v{Z}}liobaite, Indre and Calders, Toon (1 rows)</li> </ul>	<input checked="" type="checkbox"/>	Kamiran, Faisal and Calders, Toon
2	3	<ul style="list-style-type: none"> <li>Sara Hajian and Josep Domingo-Ferrer (2 rows)</li> <li>Sara Hajian and Josep Domingo-Ferrer and Antoni Mart{v{i}}nez-Ballest{v{e}} (1 rows)</li> </ul>	<input checked="" type="checkbox"/>	Sara Hajian and Josep Domingo-Ferr
2	2	<ul style="list-style-type: none"> <li>Ruggieri, Salvatore (1 rows)</li> <li>Ruggieri, Salvatore and Pedreschi, Dino and Turini, Franco (1 rows)</li> </ul>	<input type="checkbox"/>	Ruggieri, Salvatore
2	2	<ul style="list-style-type: none"> <li>Kamishima, Toshihiro and Akaho, Shotaro (1 rows)</li> <li>Kamishima, Toshihiro and Akaho, Shotaro and Asoh, Hideki and Sakuma, Jun (1 rows)</li> </ul>	<input type="checkbox"/>	Kamishima, Toshihiro and Akaho, S
2	2	<ul style="list-style-type: none"> <li>Hajian, Sara and Domingo-Ferrer, Josep and Martinez-Balleste, Antoni (1 rows)</li> <li>Hajian, Sara and Domingo-ferrer, Josep (1 rows)</li> </ul>	<input type="checkbox"/>	Hajian, Sara and Domingo-Ferrer, J

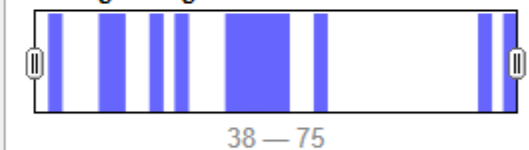
# Choices in Cluster



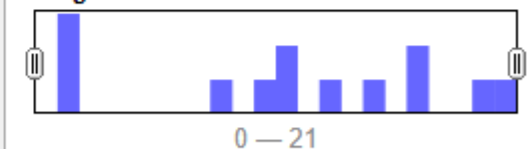
# Rows in Cluster



Average Length of Choices



Length Variance of Choices



Select All Unselect All

Export Clusters

Merge Selected & Re-Cluster

Merge Selected & Close

Close

ISBN		change
18 choices Sort by: name count		Cluster
1011501104	1	
1011501205848	1	
1050601491	1	
9780769539027	1	
9780769544083	1	
9780769544090	1	
9780769549057	1	
9780769549255	3	
9781424433148	1	
9781424499069	1	
9781450300322	1	
9781450308137	1	

								▼ Booktitle	▼ Chapter	▼ Edition	▼ Series	▼ Editor	▼ Publisher
☆	🗨	87.	Book	1011501104	6	Kamiran, Faisal and Calders, Toon	Data preprocessing techniques for classification without discrimination	Knowledge and Information Systems					

1 matching rows (105 total)

Extensions: Wikidata

Show as: rows records

Show: 5 10 25 50 rows

« first ‹ previous 1 - 1 next › last

<input type="checkbox"/> All	<input type="checkbox"/> BibliographyType	<input type="checkbox"/> ISBN	<input type="checkbox"/> Identifier	<input type="checkbox"/> Author	<input type="checkbox"/> Title	<input type="checkbox"/> Journal	<input type="checkbox"/> Volume	<input type="checkbox"/> Number	<input type="checkbox"/> Month	<input type="checkbox"/> Pages	<input type="checkbox"/> Year	<input type="checkbox"/> Address	<input type="checkbox"/> Notes	
		52.	InCollection	0001011501104	6	Kamiran, Faisal and Calders, Toon	Data preprocessing techniques for classification without discrimination		33	1		1--33	2012	

### Custom text transform on column BibliographyType

Expression Language **General Refine Expression Language (GREL)**

```
if(cells["Year"].value == null,"Misc",value)
```

No syntax error.

**Preview** History Starred Help

row	value	if(cells["Year"].value == null ...
1.	Article	Article
2.	Article	Article
3.	Article	Article
4.	Article	Article
5.	Article	Article
6.	Article	Article
7.	Article	Article

On error ☒ keep original ☐ set to blank ☐ store error ☐ Re-transform up to  times until no change

OK Cancel

OpenRefine bibliographycsv Permalink

Facet / Filter Undo / Redo 3 / 4

Extract... Apply...

Filter:

112 rows

Show as: rows records Show: 5 10 25 50 rows

0. Create project

1. Text transform on 14 cells in column Year:  
 jython:if value == None: return value elif  
 int(value) >=0 and int(value) <20: value =  
 str(int(value)+2000) return value elif  
 int(value) >19 and int(value) <99: value =  
 str(int(value)+1900) return value else: return  
 value

2. Mass edit 14 cells in column Author

3. Star row 87

4. Star 111 rows

Transform

Facet

Edit rows

Edit columns

View

Facet by star

Facet by flag

All	BibliographyTyp	ISBN	Identifier	Author	Title	
1				Sarah Holland and Ahmed Hosny and Sarah Newman and Joshua Joseph and Kasia Chmielinski	The Dataset Nutrition Label: A Framework To Drive Higher Data Quality Standards	
★	2.	Article	1050601491	10	Mancuhan, Koray and Clifton, Chris	Combating discrimination using Bayesian networks
★	3.	Article		100	Feldman, Michael and Friedler, Sorelle and Moeller, John and Scheidegger, Carlos and Voelkswagen	Certifying and removing disparate impact

Data javascript:[]

The screenshot shows the OpenRefine web interface in a browser. The address bar displays the URL `127.0.0.1:3333/project?project=1988305065945`. The page title is "OpenRefine bibliographycsv". The interface includes a "Facet / Filter" sidebar on the left with a "Using facets and filters" section. The main area shows "112 rows" and a table with columns: "All", "BibliographyType", "ISBN", "Identifier", and "A". A context menu is open over the "ISBN" column, with "Edit rows" selected. A sub-menu is open for "Edit rows", showing options: "Star rows", "Unstar rows", "Flag rows", "Unflag rows", and "Remove all matching rows". The "Remove all matching rows" option is highlighted with a green box.

Google x bibliographycsv - OpenRefine x Bibliography - OpenRefine

← → ↻ ⓘ 127.0.0.1:3333/project?project=1988305065945

**OpenRefine** bibliographycsv Permalink

Facet / Filter Undo / Redo 3 / 3

**Using facets and filters**

Use facets and filters to select subsets of your data to act on. Choose facet and filter methods from the menus at the top of each data column.

Not sure how to get started?  
[Watch these screencasts](#)

112 rows

Show as: rows records Show: 5 10 25 50 rows

▼ All	▼ BibliographyType	▼ ISBN	▼ Identifier	▼ A
Transform				
Facet		1011501104	6	Kamir Calde
Edit rows				
Edit columns				
View				

Star rows

Unstar rows

Flag rows

Unflag rows

Remove all matching rows

javascript:()



OpenRefine interface showing a bibliography dataset with 112 rows. The interface includes a header bar with tabs for 'Bibliography - OpenRefine' and 'ILIAS für Lehre & Lernen - Universität Stuttgart'. The main area displays a table of bibliographic records with columns: All, BibliographyType, ISBN, Identifier, Author, Title, Journal, Volume, Number, and Month. A context menu is open over the first row, showing options like Facet, Text filter, Edit cells, Edit column, Transpose, Sort..., View, and Reconcile. The 'Edit cells' option is selected, and a sub-menu is open showing various transformation options. The 'Common transforms' option is highlighted, and a further sub-menu is open showing options like Trim leading and trailing whitespace, Collapse consecutive whitespace, Unescape HTML entities, To titlecase, To uppercase, To lowercase, To number, To date, To text, To null, and To empty string. The 'To number' option is highlighted.

All	BibliographyType	ISBN	Identifier	Author	Title	Journal	Volume	Number	Month
☆	1.	Article							May
☆	2.	Article	1050601491	10	Mancu...				
☆	3.	Article		100	Feldman, Michael and Friedler, Sorelle and Moeller, John and Scheidegger, Carlos and Venkatasubramanian	Certifying and removing disparate impact			

**OpenRefine** Bibliography [Permalink](#)

Facet / Filter **Undo / Redo 54 / 54**

Extract... Apply...

Filter:

- Create project
- Mass edit 6 cells in column Month
- Mass edit 2 cells in column Month
- Mass edit 1 cells in column Month
- Mass edit 1 cells in column Month
- Mass edit 6 cells in column Month
- Mass edit 1 cells in column Month
- Mass edit 1 cells in column Month
- Mass edit 1 cells in column Month
- Mass edit 2 cells in column Month
- Mass edit 1 cells in column Month
- Mass edit 4 cells in column Month
- Mass edit 1 cells in column Year
- Text transform on 13 cells in column Year: `jython:if value == None: return value if value == None: return value elif int(value) >=0 and int(value) <20: value = str(int(value)+2000) return value elif int(value) >19 and int(value) <99: value = str(int(value)+1900) return value else:`

27.0.0.1333

**OpenRefine** Bibliography [Permalink](#)

Facet / Filter **Undo / Redo 47 / 54**

Extract... Apply...

Filter:

- Edit single cell on row 93, column BibliographyType
- Edit single cell on row 92, column BibliographyType
- Edit single cell on row 99, column BibliographyType
- Text transform on 5 cells in column BibliographyType: `grel:if(cells["Year"].value == null,"Misc",value)`
- Text transform on 24 cells in column BibliographyType: `grel:value.replace("Article","Misc")`
- Mass edit 5 cells in column BibliographyType
- Remove 1 rows
- Edit single cell on row 77, column Publisher
- Remove 1 rows
- Text transform on 105 cells in column Identifier: `value.toNumber()`
- Reorder rows

The screenshot shows a web application interface with a table of data. The table has columns: Identifier, Author, Title, Journal, Volume, and Number. The first row of data is: 1491, 10, Mancuhan, Koray and Clifton, Chris, Combating discrimination using Bayesian networks, Artificial Intelligence and Law, 22, 2. The table indicates 25 rows are visible out of 50 total rows.

An 'Export' button is highlighted with a red box. A dropdown menu is open, showing various export options. The 'Export project' option is also highlighted with a red box. The dropdown menu includes the following options:

- Export project
- Project data package
- Tab-separated value
- Comma-separated value
- HTML table
- Excel (.xls)
- Excel 2007+ (.xlsx)
- ODF spreadsheet
- Custom tabular exporter...
- SQL Exporter...
- Templating...
- Upload edits to Wikidata
- Export to QuickStatements
- Export schema

	Identifier	Author	Title	Journal	Volume	Number
1491	10	Mancuhan, Koray and Clifton, Chris	Combating discrimination using Bayesian networks	Artificial Intelligence and Law	22	2

Google OpenRefine


127.0.0.1:3333

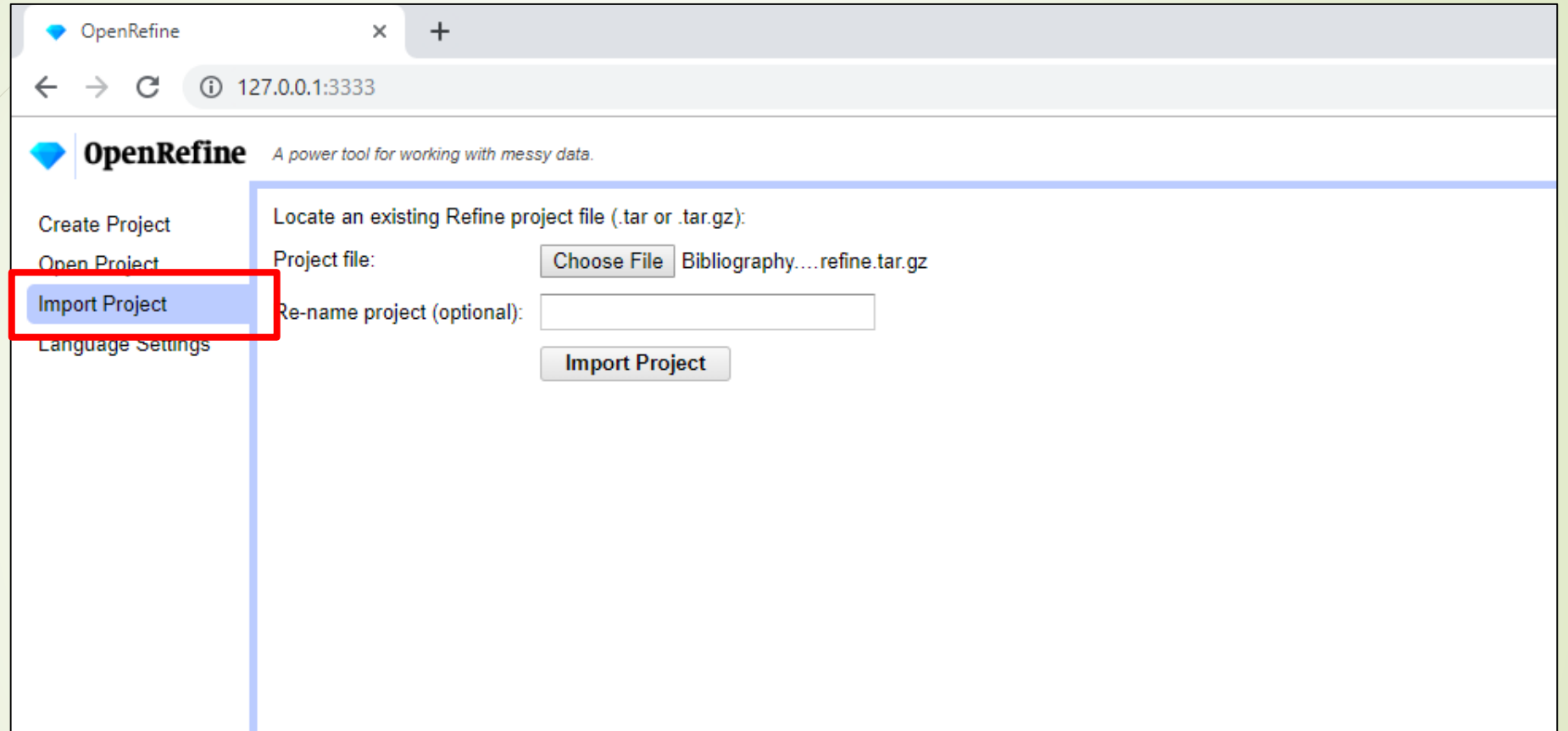
**OpenRefine** *A power tool for working with messy data.*

Create Project  
**Open Project**  
Import Project  
Language Settings

All

		Last modified	Name	Tags	Subject	Description
X	About	2019-06-25 23:22 PM	Bibliography			
X	About	2019-06-25 23:21 PM	bibliographycsv			
X	About	2019-06-19 18:41 PM	Bibliography			
X	About	2019-06-19 18:41 PM	Bibliography			
X	About	2019-06-19 18:32 PM	Bibliography			
X	About	2019-06-19 18:20 PM	Bibliography			
X	About	2019-06-19 18:19 PM	Bibliography			

  
Version 3.1 [b90e413]





60

# Thank You