



**ADITYA**

COLLEGE OF ENGINEERING & TECHNOLOGY

# DATA WAREHOUSE AND DATA MINING LAB

**ARFF CREATION**

**By**

**Nadella Sunil**

**M.Tech., (PhD) UGC NET & AP SET (Computer Science & Applications)  
TS & AP SET (Mathematical Sciences)**

**Department of Computer Science and Engineering**

**Aditya College of Engineering & Technology, Surampalem, Andhra Pradesh**

# Agenda

- Procedure to create **ARFF** File (Attribute relation File Format) explained in this presentation
- As the name suggests it described a list of instances sharing a set of attributes.
- these files are supported by **WEKA machine Learning tool**,
- The *arff* files are used for the purpose of various operations related to data preprocessing, data cleaning etc.

# Structure of ARFF file



- ARFF file contains 2 sections
  1. Header Section
  2. Data Section
- All the *keywords in ARFF file start with @* symbol.

# Header Section



- This section contains various information related to the dataset like the name of the relation, columns, and type of columns. The header section contains 2 parts
  1. **Table/relation** and
  2. **attribute** part.

# Header Section



- *@relation* :used to give the table name
- *@attribute*: used to give a column name

## *DATA TYPES*

- ***nominal***: represented inside curly brackets (Like constants)
- ***string*** : data type which accepts only string value
- ***numeric***: used to store numbers
- ***date***: used to store date

# Header Section – Syntax



- @relation tablename
- @attribute column\_name type

# Header Section –Example



@relation "employee"

@attribute f\_name string

@attribute l\_name string

@attribute contact\_num numeric

@attribute dept {HR,IT,MANAGEMENT,MAINTAINANCE}

@attribute DOB date dd-mm-yyyy

@attribute city string

*Here dept column is having nominal data type so it can only accept above mentioned types of data only,*



# Data Section



- **Data section** is used to represents the data or entries for available columns. (according to the order in header section data would be inserted).
- **data section starts with @data**, and this section must be added after Header section. only single record can be written in single line.
- **@data**: Used to start data section
- **%**: % sign is used to represent the comment in file.



# Data Section - Syntax



*@data*  
*<record1>*  
*<record2>*  
*•*  
*•*  
*<record N>*

*all the Records must be in the same format as their attributes are defined in Header section Like*

# Data Section - Example



1,naman,N,1234556678,IT,02-08-2000,rjt

2,yash,M,1234556679,HR,04-05-2001,amd

3,kishan,G,1214556678,MANAGEMENT,02-11-2001,pbr

4,?,?,5234556678,IT,03-05-2000,amd

- ✓ We separate values by comma(,) and to
- ✓ represent the empty or missing value for a particular column we use the (?)sign.

# Emp.arff



@relation "employee"

@attribute id numeric

@attribute f\_name string

@attribute l\_name string

@attribute contact\_num numeric

@attribute dept {HR,IT,MANAGEMENT,MAINTAINANCE}

@attribute DOB date dd-mm-yyyy

@attribute city string

@data

1,naman,N,1234556678,IT,02-08-2000,rjt

2,yash,M,1234556679,HR,04-05-2001,amd

3,kishan,G,1214556678,MANAGEMENT,02-11-2001,pbr

4,?,?,5234556678,IT,03-05-2000,amd

# How to Create and open arff file - Procedure



**Step 1:** Open any text editor and paste the above code.

**Step 2:** Save the file with emp\_dm.arff file extension

**Step 3:** Open weka tool

**Step 4:** Click on Explorer

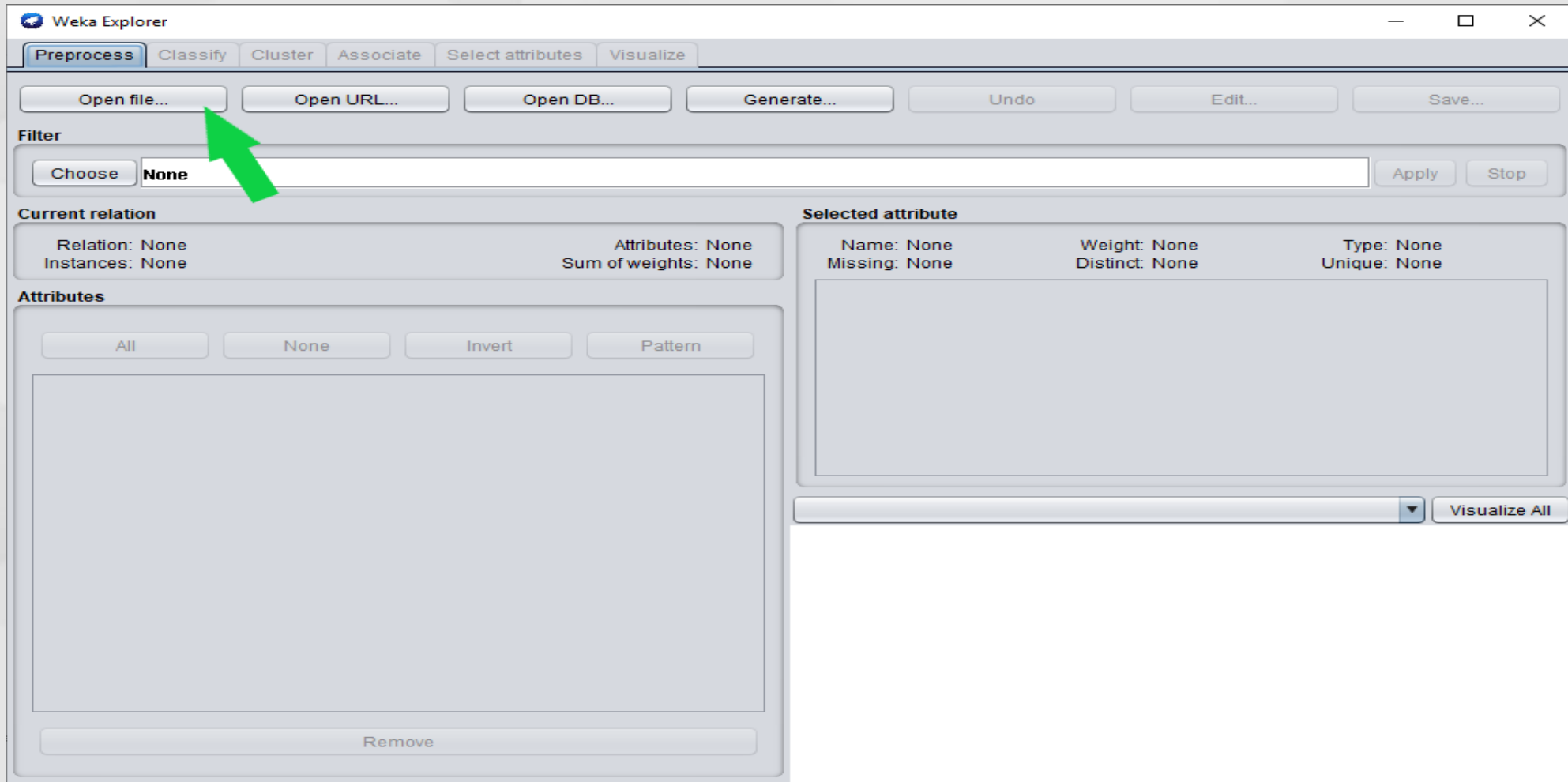
# Open WEKA tool and select Explore



**ADITYA**  
COLLEGE OF ENGINEERING & TECHNOLOGY

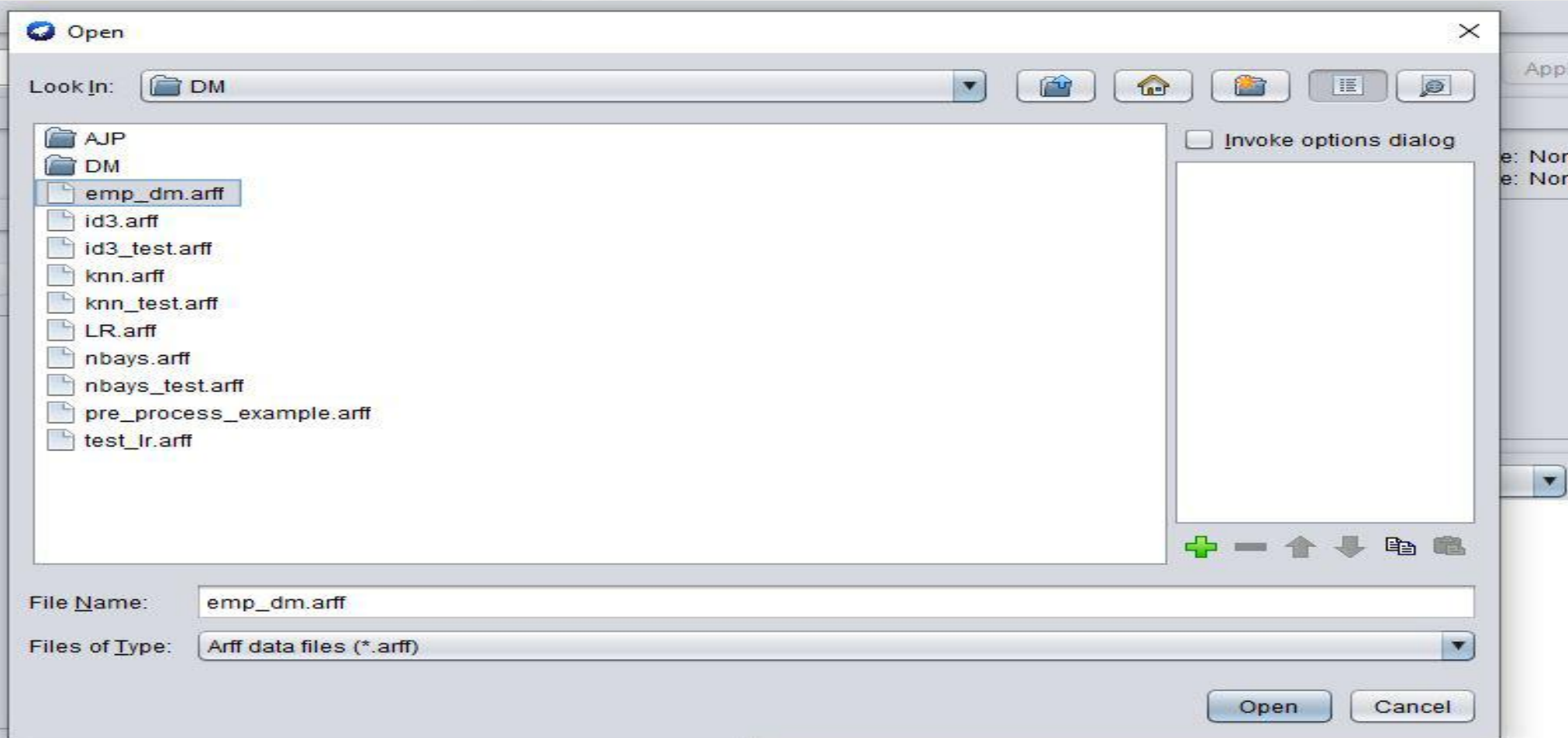


# Step – 5 : Then click on Open file



## Step – 6 :

Select/Locate arff file from disk then click On Open.





## Step – 7 :

file is now Loaded now click on Edit from Preprocess Tab



**ADITYA**  
COLLEGE OF ENGINEERING & TECHNOLOGY

Weka Explorer

Preprocess Classify Cluster Associate Select attributes Visualize

Open file... Open URL... Open DB... Generate... Undo Edit... Save...

Filter  
Choose None Apply

Current relation  
Relation: employee  
Instances: 4  
Attributes: 7  
Sum of weights: 4

Attributes  
All None Invert Pattern

No.	Name
1	<input checked="" type="checkbox"/> id
2	<input type="checkbox"/> f_name
3	<input type="checkbox"/> l_name
4	<input type="checkbox"/> contact_num
5	<input type="checkbox"/> dept
6	<input type="checkbox"/> DOB
7	<input type="checkbox"/> city

Selected attribute

Name: id  
Missing: 0 (0%)  
Distinct: 4  
Type: Numeric  
Unique: 4 (100%)

Statistic	Value
Minimum	1
Maximum	4
Mean	2.5
StdDev	1.291

Class: city (Str) Vis

# Step – 8 : viewing the emp.arff dataset



**ADITYA**  
COLLEGE OF ENGINEERING & TECHNOLOGY

Weka Explorer

Preprocess | Classify | Cluster | Associate | Select attributes | Visualize

Open file... | Open URL... | Open DB... | Generate... | Undo | Edit... | Save...

Filter

Viewer

Relation: employee

No.	1: id	2: f_name	3: l_name	4: contact_num	5: dept	6: DOB	7: city
	Numeric	String	String	Numeric	Nominal	Date	String
1	1.0	naman	N	1.23455667...	IT	02-0...	rjt
2	2.0	yash	M	1.23455667...	HR	04-0...	amd
3	3.0	kishan	G	1.21455667...	MAN...	02-1...	pbr
4	4.0			5.23455667...	IT	03-0...	amd

No.

Thank You.

