

# Retrieving and Working with Datasets

**Prof. Pietro Ducange**



# Where to retrieve interesting datasets

- UCI Machine Learning Repository

<https://archive.ics.uci.edu/ml/datasets.html>

- Keel Dataset Repository

<http://sci2s.ugr.es/keel/datasets.php>

- WEKA

<http://www.cs.waikato.ac.nz/ml/weka/datasets.html>

- ChemDB

<http://www.cs.ox.ac.uk/activities/machinelearning/applications.html>



# Where to retrieve interesting datasets: challenges

- Kaggle

<https://www.kaggle.com/datasets>

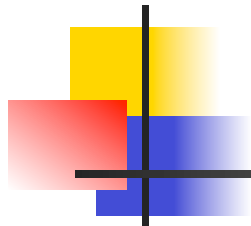
- TIM big data Challenge

<https://dandelion.eu/datamine/open-big-data/>

- TunedIT

<http://tunedit.org/challenges/>





# PIMA Indian Diabetes

---

- From the UCI repository
- The class attribute specifies whether patient shows or not signs of diabetes according to World Health Organization criteria
- 2 classes, 8 attributes, 768 instances, 500 (65.1%) negative, and 268 (34.9%) positive tests for diabetes
- All patients were females at least 21 years old of Pima Indian heritage

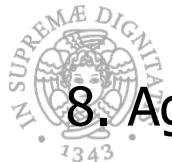


# PIMA Indian Diabetes

Attributes:

1. Number of times pregnant
2. Plasma glucose concentration
3. Diastolic blood pressure (mm Hg)
4. Triceps skin fold thickness (mm)
5. 2-Hour serum insulin ( $\mu$ U/ml)
6. Body mass index ( $\text{weight in kg}/(\text{height in m})^2$ )
7. Diabetes pedigree function ( a sort of ancestor's history)

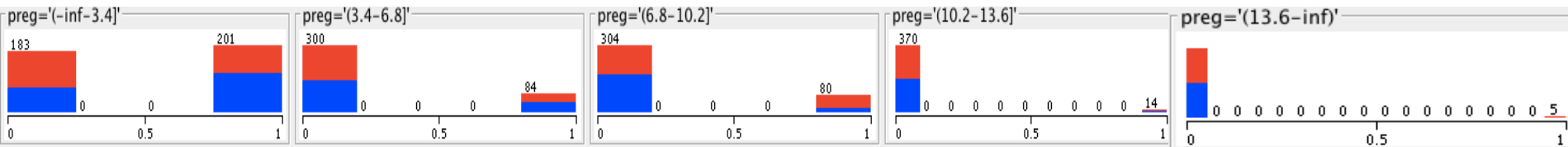
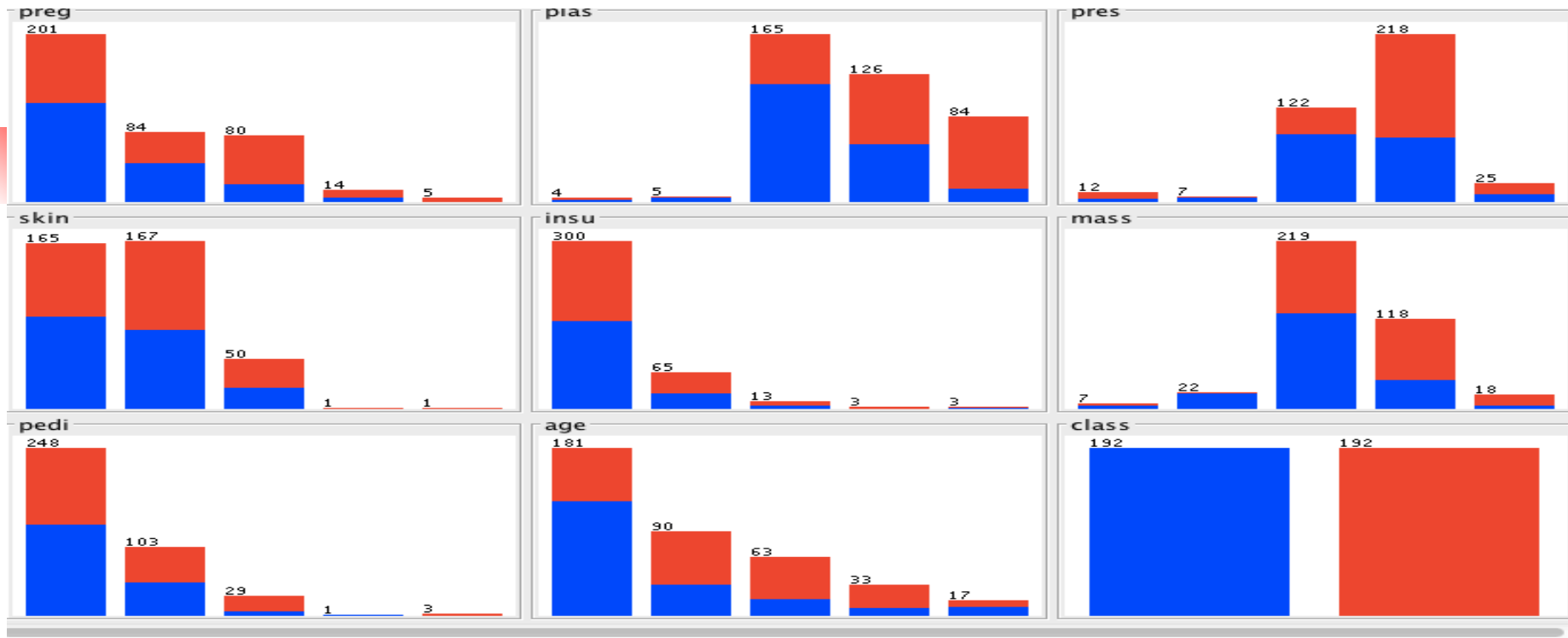
8. Age (years)



# Exercise

1. Load the diabetes.arff dataset
2. Apply a supervised resampling considering just 50% of the original instances (with default parameters)
3. Apply a supervised resampling considering just 50% of the original instances (with the first parameter set to 1), save the dataset and analyze the differences with the results of 2)
4. Discretize the saved dataset by using 5 bins (unsupervised filter)
5. Apply the NominalToBinary Filter to the discretized dataset and comment the results





# Hepatitis Data Set

**Type** Classification

**Origin** Real world **Features** 19

**(Real / Integer / Nominal)** (2 / 17 / 0)

**Classes** 2 **Missing values?** Yes

**Total instances** 155

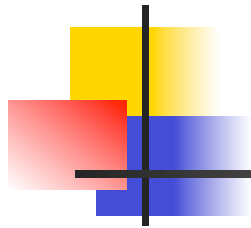
UCI URL: <https://archive.ics.uci.edu/ml/datasets/Hepatitis>

*The task is to predict if patients will die (1) or survive (2).*

Class Distribution: DIE: 32 LIVE: 123







# Hepatitis Data Set

## Attribute Information:

1. Class: DIE, LIVE
2. AGE: 10, 20, 30, 40, 50, 60, 70, 80
3. SEX: male, female
4. STEROID: no, yes
5. ANTIVIRALS: no, yes
6. FATIGUE: no, yes
7. MALAISE: no, yes
8. ANOREXIA: no, yes
9. LIVER BIG: no, yes
10. LIVER FIRM: no, yes
11. SPLEEN PALPABLE: no, yes
12. SPIDERS: no, yes
13. ASCITES: no, yes
14. VARICES: no, yes
15. BILIRUBIN: 0.39, 0.80, 1.20, 2.00, 3.00, 4.00
16. ALK PHOSPHATE: 33, 80, 120, 160, 200, 250
17. SGOT: 13, 100, 200, 300, 400, 500,
18. ALBUMIN: 2.1, 3.0, 3.8, 4.5, 5.0, 6.0
19. PROTINE: 10, 20, 30, 40, 50, 60, 70, 80, 90
20. HISTOLOGY: no, yes





# Exercise

---

- Download the rough Hepatitis dataset from UCI
- Prepare it for loading in WEKA (if needed)
- Load the dataset in WEKA
- Handle missing values
- Rebalance the dataset with the resampling filter



## Some indications for relabeling Hepatitis dataset

Relabeled values in attribute SEX

% From: 2 To: male

% From: 1 To: female

% Relabeled values in attribute STEROID

% From: 1 To: no

% From: 2 To: yes

% Relabeled values in attribute

ANTIVIRALS

% From: 2 To: no

% From: 1 To: yes

% Relabeled values in attribute FATIGUE

% From: 2 To: no

% From: 1 To: yes

% Relabeled values in attribute MALAISE

% From: 2 To: no

% From: 1 To: yes

% Relabeled values in attribute ANOREXIA

% From: 2 To: no

% From: 1 To: yes

% Relabeled values in attribute LIVER\_BIG

% From: 1 To: no

% From: 2 To: yes

% Relabeled values in attribute LIVER\_FIRM

% From: 2 To: no

% From: 1 To: yes

% Relabeled values in attribute SPLEEN\_PALPABLE

% From: 2 To: no

% From: 1 To: yes

% Relabeled values in attribute SPIDERS

% From: 2 To: no

% From: 1 To: yes

% Relabeled values in attribute ASCITES

% From: 2 To: no

% From: 1 To: yes

% Relabeled values in attribute VARICES

% From: 2 To: no

% From: 1 To: yes

% Relabeled values in attribute HISTOLOGY

% From: 1 To: no

%<sub>1</sub> From: 2 To: yes



# Rought CSV file (1)

2,30,2,1,2,2,2,2,1,2,2,2,2,2,1.00,85,18,4.0,?,1

....

2,50,1,1,2,1,2,2,1,2,2,2,2,2,0.90,135,42,3.5,?,1

2,31,1,?,1,2,2,2,2,2,2,2,2,2,0.70,46,52,4.0,80,1

2,34,1,2,2,2,2,2,2,2,2,2,2,2,1.00,?,200,4.0,?,1

2,34,1,2,2,2,2,2,2,2,2,2,2,2,0.90,95,28,4.0,75,1

1,51,1,1,2,1,2,1,2,2,1,1,2,2,?,?,?,?,?,1

....

***To load in WEKA the dataset, just open the file paste the string of the name of each attribute, separeted by a comma, as the first line of the file and save it as a .CSV file->***

## Rought CSV file (2)

CLASS, AGE, SEX, STEROID, ANTIVIRALS, FATIGUE, MALAISE,  
ANOREXIA, LIVER\_BIG, LIVER\_FIRM, SLEEN\_PALPABLE,  
SPIDERS, ASCITES, VARICES, BILIRUBIN, ALK\_PHOSPHATE,  
SGOT, ALBUMIN, PROTINE, HISTOLOGY

2,30,2,1,2,2,2,2,1,2,2,2,2,2,1.00,85,18,4.0,?,1

....

2,50,1,1,2,1,2,2,1,2,2,2,2,2,0.90,135,42,3.5,?,1

2,31,1,?,1,2,2,2,2,2,2,2,2,2,0.70,46,52,4.0,80,1

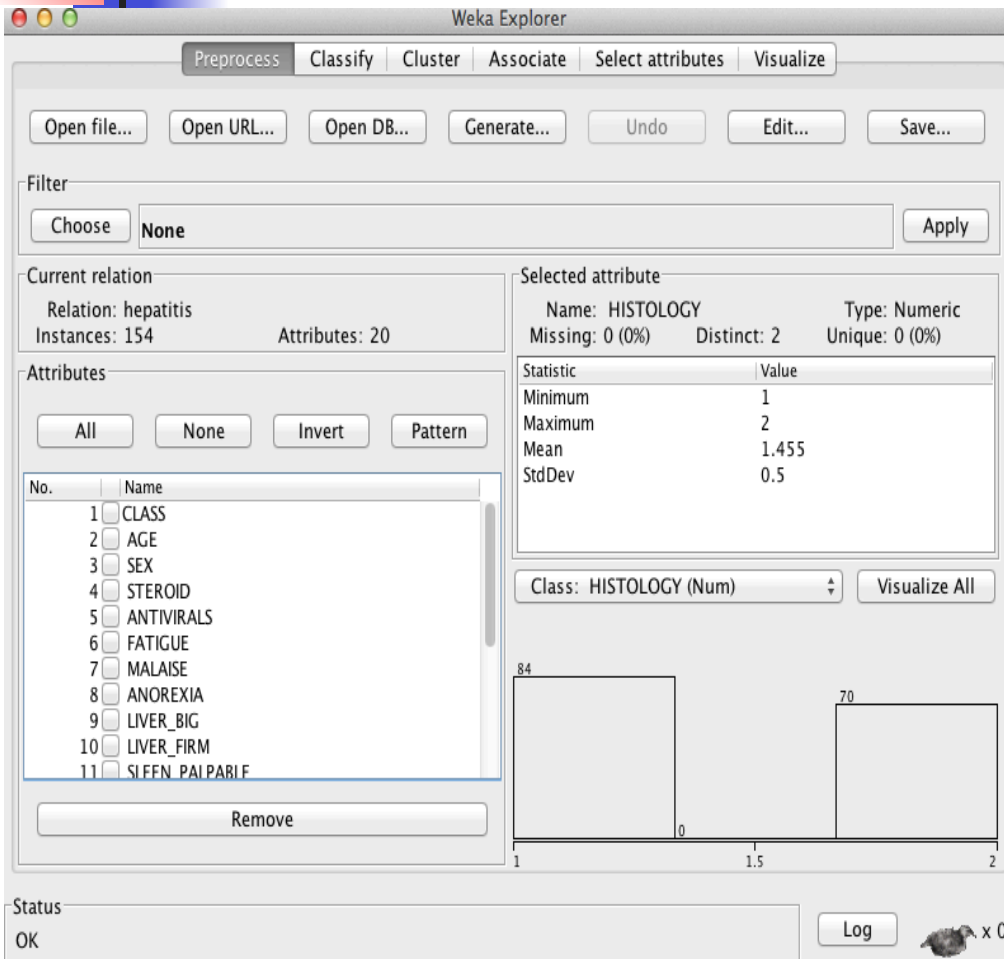
2,34,1,2,2,2,2,2,2,2,2,2,2,1.00,?,200,4.0,?,1

2,34,1,2,2,2,2,2,2,2,2,2,2,0.90,95,28,4.0,75,1

1,51,1,1,2,1,2,1,2,2,1,1,2,2,?,?,?,?,1



# Loading the CSV file for the dataset in WEKA



Open the dataset editing (edit button) and set the CLASS attribute as a class.

Using the save button, you can save the dataset in .arff format.



If you open the arff file with a text editor you will find:

---

@relation hepatitis

@attribute 'AGE' numeric

@attribute 'SEX' numeric

@attribute 'STEROID' numeric

@attribute 'ANTIVIRALS' numeric

@attribute 'FATIGUE' numeric

@attribute 'MALAISE' numeric

@attribute 'ANOREXIA' numeric

@attribute 'LIVER\_BIG' numeric

@attribute 'LIVER\_FIRM' numeric

@attribute 'SLEEN\_PALPABLE' numeric

@attribute 'SPIDERS' numeric

@attribute 'ASCITES' numeric

@attribute 'VARICES' numeric

@attribute 'BILIRUBIN' numeric

@attribute 'ALK\_PHOSPHATE' numeric

@attribute 'SGOT' numeric

@attribute 'ALBUMIN' numeric

@attribute 'PROTIME' numeric

@attribute 'HISTOLOGY' numeric

@attribute CLASS numeric



You can relabel by hand the arff file (pay attention to the class attribute).

@relation hepatitis

@attribute 'AGE' numeric

@attribute 'SEX' numeric

@attribute 'STEROID' numeric

@attribute 'ANTIVIRALS' numeric

@attribute 'FATIGUE' numeric

@attribute 'MALAISE' numeric

@attribute 'ANOREXIA' numeric

@attribute 'LIVER\_BIG' numeric

@attribute 'LIVER\_FIRM' numeric

@attribute 'SLEEN\_PALPABLE' numeric

@attribute 'SPIDERS' numeric

@attribute 'ASCITES' numeric

@attribute 'VARICES' numeric

@attribute 'BILIRUBIN' numeric

@attribute 'ALK\_PHOSPHATE' numeric

@attribute 'SGOT' numeric

@attribute 'ALBUMIN' numeric

@attribute 'PROTIME' numeric

@attribute 'HISTOLOGY' numeric

@attribute 'class' { 1, 2 }





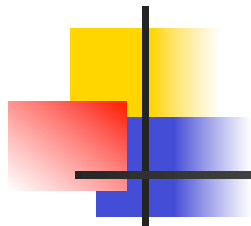
# If you load the new arff dataset in WEKA

The screenshot shows the Weka Explorer window with the following details:

- Filter:** Discretize -B 2 -M -1.0 -R last
- Current relation:** Relation: hepa-weka.filters.u... Attributes: 20 Instances: 155 Sum of weights: 155
- Attributes:** A list of 20 attributes is shown, with 'CLASS' selected at index 20. Other attributes include SLEEN\_PALPABLE, SPIDERS, ASCITES, VARICES, BILIRUBIN, ALK\_PHOSPHATE, SGOT, ALBUMIN, PROTINE, and HISTOLOGY.
- Selected attribute:** Name: CLASS, Type: Nominal, Missing: 0 (0%), Distinct: 2, Unique: 0 (0%).
- Table of Selected Attribute:**

No.	Label	Count	Weight
1	'(-inf-1.5]'	32	32.0
2	'(1.5-inf)'	123	123.0
- Bar Chart:** A bar chart visualizes the distribution of the 'CLASS' attribute, with a blue bar for the first class (count 32) and a red bar for the second class (count 123).
- Status:** OK





After applying the ReplaceMissingValues filter  
(weka.filters.unsupervised.attribute.ReplaceMissingValues),  
you can apply the supervised resample filter with 100% and  
the first parameter set to 1



# The rebalanced Hepatitis Dataset in WEKA

Weka Explorer

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

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

Filter: Choose **Resample -B 1.0 -S 1 -Z 100.0** Apply

Current relation: Relation: hepa-weka.filters.u... Attributes: 20 Instances: 154 Sum of weights: 154

Attributes: All | None | Invert | Pattern

No.	Name
10	SLEEN_PALPABLE
11	SPIDERS
12	ASCITES
13	VARICES
14	BILIRUBIN
15	ALK_PHOSPHATE
16	SGOT
17	ALBUMIN
18	PROTIME
19	HISTOLOGY
20	CLASS

Remove

Selected attribute: Name: CLASS Missing: 0 (0%) Distinct: 2 Type: Nominal Unique: 0 (0%)

No.	Label	Count	Weight
1	'(-inf-1.5]'	77	77.0
2	'(1.5-inf)'	77	77.0

Class: CLASS (Nom) Visualize All

77 77

Status: OK Log x 0

