

Introduction to Weka

by Prof. Weijia Jia (賈維嘉) E11-4007

Email: jiawj@umac.mo

Introduction

- WEKA的全名是怀卡托智能分析环境（Waikato Environment for Knowledge Analysis），它的源代码：
- <http://www.cs.waikato.ac.nz/ml/weka>
- 同时weka也是新西兰的一种鸟名，而WEKA的主要开发者来自新西兰。



Introduction

- WEKA作为一个公开的数据挖掘工作平台，集合了大量能承担数据挖掘任务的机器学习算法，包括对数据进行预处理，分类，回归、聚类、关联规则以及在新的交互式界面上的可视化。
- 2005年8月，在第11届ACM SIGKDD国际会议上，怀卡托大学的Weka小组荣获了数据挖掘和知识探索领域的最高服务奖。

Data Format

- WEKA所处理的数据集是图中那样的一个二维的表格
- 表里的一个横行称作一个实例 (Instance)，相当于统计学中的一个样本，或者数据库中的一条记录
- 竖行称作一个属性 (Attribute)，相当于统计学中的一个变量
- 这样一个表格 (数据集)，在WEKA看来，呈现了属性之间的一种关系 (Relation)

| weather.arff | | | | | |
|-------------------|----------------------|--------------------------|-----------------------|--------------------|-------------------|
| Relation: weather | | | | | |
| No. | 1.outlook Nominal | 2.temperature Numeric | 3.humidity Numeric | 4.windy Nominal | 5.play Nominal |
| 1 | sunny | 85.0 | 85.0 | FALSE | no |
| 2 | sunny | 80.0 | 90.0 | TRUE | no |
| 3 | overcast | 83.0 | 86.0 | FALSE | yes |
| 4 | rainy | 70.0 | 96.0 | FALSE | yes |
| 5 | rainy | 68.0 | 80.0 | FALSE | yes |
| 6 | rainy | 65.0 | 70.0 | TRUE | no |
| 7 | overcast | 64.0 | 65.0 | TRUE | yes |
| 8 | sunny | 72.0 | 95.0 | FALSE | no |
| 9 | sunny | 69.0 | 70.0 | FALSE | yes |
| 10 | rainy | 75.0 | 80.0 | FALSE | yes |
| 11 | sunny | 75.0 | 70.0 | TRUE | yes |
| 12 | overcast | 72.0 | 90.0 | TRUE | yes |
| 13 | overcast | 81.0 | 75.0 | FALSE | yes |
| 14 | rainy | 71.0 | 91.0 | TRUE | no |

Data Format

- WEKA存储数据的格式是ARFF (Attribute-Relation File Format) 文件, 这是一种ASCII文本文件
- 图中所示的二维表格存储在ARFF文件中

```
% ARFF file for the weather data with some numeric features ↓
% ↓
@relation weather ↓
↓
@attribute outlook {sunny, overcast, rainy} ↓
@attribute temperature real ↓
@attribute humidity real ↓
@attribute windy {TRUE, FALSE} ↓
@attribute play {yes, no} ↓
↓
@data ↓
% ↓
% 14 instances ↓
% ↓
sunny,85,85,FALSE,no ↓
sunny,80,90,TRUE,no ↓
overcast,83,86,FALSE,yes ↓
rainy,70,96,FALSE,yes ↓
rainy,68,80,FALSE,yes ↓
rainy,65,70,TRUE,no ↓
overcast,64,65,TRUE,yes ↓
sunny,72,95,FALSE,no ↓
sunny,69,70,FALSE,yes ↓
rainy,75,80,FALSE,yes ↓
sunny,75,70,TRUE,yes ↓
overcast,72,90,TRUE,yes ↓
overcast,81,75,FALSE,yes ↓
rainy,71,91,TRUE,no ↵
```

Data Format

- 关系声明
关系名称在ARFF文件的第一个有效行来定义， 格式为
@relation <relation-name>

```
% ARFF file for the weather data with some numeric features ↓
% ↓
@relation weather ↓
↓
@attribute outlook {sunny, overcast, rainy} ↓
@attribute temperature real ↓
@attribute humidity real ↓
@attribute windy {TRUE, FALSE} ↓
@attribute play {yes, no} ↓
↓
@data ↓
% ↓
% 14 instances ↓
% ↓
sunny,85,85,FALSE,no ↓
sunny,80,90,TRUE,no ↓
overcast,83,86,FALSE,yes ↓
rainy,70,96,FALSE,yes ↓
rainy,68,80,FALSE,yes ↓
rainy,65,70,TRUE,no ↓
overcast,64,65,TRUE,yes ↓
sunny,72,95,FALSE,no ↓
sunny,69,70,FALSE,yes ↓
rainy,75,80,FALSE,yes ↓
sunny,75,70,TRUE,yes ↓
overcast,72,90,TRUE,yes ↓
overcast,81,75,FALSE,yes ↓
rainy,71,91,TRUE,no ↵
```

Data Format

- 属性声明
属性声明用一系列以 “@attribute” 开头的语句表示。数据集中的每一个属性都有它对应的 “@attribute” 语句，来定义它的属性名称和数据类型。这些声明语句的顺序很重要。首先它表明了该项属性在数据部分的位置。

```
% ARFF file for the weather data with some numeric features ↓
% ↓
@relation weather ↓
↓
@attribute outlook {sunny, overcast, rainy} ↓
@attribute temperature real ↓
@attribute humidity real ↓
@attribute windy {TRUE, FALSE} ↓
@attribute play {yes, no} ↓
↓
@data ↓
% ↓
% 14 instances ↓
% ↓
sunny,85,85,FALSE,no ↓
sunny,80,90,TRUE,no ↓
overcast,83,86,FALSE,yes ↓
rainy,70,96,FALSE,yes ↓
rainy,68,80,FALSE,yes ↓
rainy,65,70,TRUE,no ↓
overcast,64,65,TRUE,yes ↓
sunny,72,95,FALSE,no ↓
sunny,69,70,FALSE,yes ↓
rainy,75,80,FALSE,yes ↓
sunny,75,70,TRUE,yes ↓
overcast,72,90,TRUE,yes ↓
overcast,81,75,FALSE,yes ↓
rainy,71,91,TRUE,no ↵
```

Data Format

- 属性声明的格式为
@attribute <attribute-name>
<datatype>
- WEKA支持的<datatype>有四种:
 - numeric(real/integer)-数值型
 - <nominal-specification>-分类
(nominal) 型
 - string-字符串型
 - date [<date-format>]-日期和时间
型

```
% ARFF file for the weather data with some numeric features ↓
% ↓
@relation weather ↓
↓
@attribute outlook {sunny, overcast, rainy} ↓
@attribute temperature real ↓
@attribute humidity real ↓
@attribute windy {TRUE, FALSE} ↓
@attribute play {yes, no} ↓
↓
@data ↓
% ↓
% 14 instances ↓
% ↓
sunny,85,85,FALSE,no ↓
sunny,80,90,TRUE,no ↓
overcast,83,86,FALSE,yes ↓
rainy,70,96,FALSE,yes ↓
rainy,68,80,FALSE,yes ↓
rainy,65,70,TRUE,no ↓
overcast,64,65,TRUE,yes ↓
sunny,72,95,FALSE,no ↓
sunny,69,70,FALSE,yes ↓
rainy,75,80,FALSE,yes ↓
sunny,75,70,TRUE,yes ↓
overcast,72,90,TRUE,yes ↓
overcast,81,75,FALSE,yes ↓
rainy,71,91,TRUE,no ↵
```


Data Format

- 数据信息
数据信息中 “@data” 标记独占一行，剩下的是各个实例的数据。
- 每个实例占一行。实例的各属性值用逗号 “,” 隔开。如果某个属性的值是缺失值（missing value），用问号 “?” 表示，且这个问号不能省略。

```
% ARFF file for the weather data with some numeric features ↓
% ↓
@relation weather ↓
↓
@attribute outlook {sunny, overcast, rainy} ↓
@attribute temperature real ↓
@attribute humidity real ↓
@attribute windy {TRUE, FALSE} ↓
@attribute play {yes, no} ↓
↓
@data ↓
% ↓
% 14 instances ↓
% ↓
sunny,85,85,FALSE,no ↓
sunny,80,90,TRUE,no ↓
overcast,83,86,FALSE,yes ↓
rainy,70,96,FALSE,yes ↓
rainy,68,80,FALSE,yes ↓
rainy,65,70,TRUE,no ↓
overcast,64,65,TRUE,yes ↓
sunny,72,95,FALSE,no ↓
sunny,69,70,FALSE,yes ↓
rainy,75,80,FALSE,yes ↓
sunny,75,70,TRUE,yes ↓
overcast,72,90,TRUE,yes ↓
overcast,81,75,FALSE,yes ↓
rainy,71,91,TRUE,no ↵
```