고려대학교
KOREA UNIVERSITY

# Weka 실습 가이드

Human
Intelligence

Machine
Intelligence

# Index

# Data Set

http://archive.ics.uci.edu/ml/

# Data Set

고려대학교 KOREA UNIVERSITY

Data는 csv 로 저장

**UCI Machine Learning Repository**
Center for Machine Learning and Intelligent Systems

About  Citation Policy  Donate a Data Set  Contact

Repository  Web  Search

View ALL Data Sets

Browse Through: **60** Data Sets

**Default Task - Undo**
Classification (60)
Regression (19)
Clustering (13)
Other (4)

**Attribute Type**
Categorical (3)
Numerical (47)
Mixed (8)

**Data Type - Undo**
Multivariate (60)
Univariate (3)
Sequential (11)
Time-Series (13)
Text (3)
Domain-Theory (2)
Other (2)

**Area**
Life Sciences (11)
Physical Sciences (8)
CS / Engineering (18)
Social Sciences (4)
Business (4)
Game (3)
Other (12)

**# Attributes - Undo**
Less than 10 (17)
10 to 100 (60)
Greater than 100 (26)

**# Instances - Undo**
Less than 100 (3)
100 to 1000 (49)
Greater than 1000 (60)

**Format Type**
Matrix (51)
Non-Matrix (9)

| Name | Data Types | Default Task | |
|---|---|---|---|
| Adult | Multivariate | Classification | |
| Australian Sign Language signs | Multivariate, Time-Series | Classification | |
| Australian Sign Language signs (High Quality) | Multivariate, Time-Series | Classification | |
| Bank Marketing | Multivariate | Classification | |
| Buzz in social media | Time-Series, Multivariate | Regression, Classification | |
| Cardiotocography | Multivariate | Classification | |
| Census Income | Multivariate | Classification | |
| Census-Income (KDD) | Multivariate | Classification | |
| Chess (King-Rook vs. King-Pawn) | Multivariate | Classification | |
| Connect-4 | Multivariate, Spatial | Classification | |

**UCI Machine Learning Repository**
Center for Machine Learning and Intelligent Systems

About  Citation Policy  Donate a Data Set  Contact

Repository  Web

View ALL Data Sets

Browse Through: **49** Data Sets

Table View  List View

**Default Task - Undo**
Classification (49)
Regression (10)
Clustering (5)
Other (3)

**Attribute Type**
Categorical (6)
Numerical (25)
Mixed (16)

**Data Type - Undo**
Multivariate (49)
Univariate (1)
Sequential (2)
Time-Series (4)
Text (1)
Domain-Theory (1)
Other (1)

**Area**
Life Sciences (22)
Physical Sciences (8)
CS / Engineering (4)
Social Sciences (2)
Business (3)
Game (0)
Other (10)

**# Attributes - Undo**
Less than 10 (25)
10 to 100 (49)
Greater than 100 (13)

**# Instances - Undo**
Less than 100 (3)
100 to 1000 (49)
Greater than 1000 (60)

**Format Type**
Matrix (30)
Non-Matrix (20)

| Name | Data Types | Default Task | Attribute Types | # Instances | # Attributes | Year |
|---|---|---|---|---|---|---|
| Annealing | Multivariate | Classification | Categorical, Integer, Real | 798 | 38 | |
| Audiology (Standardized) | Multivariate | Classification | Categorical | 226 | 69 | 1992 |
| Breast Cancer Wisconsin (Diagnostic) | Multivariate | Classification | Real | 569 | 32 | 1995 |
| Breast Cancer Wisconsin (Original) | Multivariate | Classification | Integer | 699 | 10 | 1992 |
| Breast Cancer Wisconsin (Prognostic) | Multivariate | Classification, Regression | Real | 198 | 34 | 1995 |
| Breast Tissue | Multivariate | Classification | Real | 106 | 10 | 2010 |
| Chronic_Kidney_Disease | Multivariate | Classification | Real | 400 | 25 | 2015 |
| Climate Model Simulation Crashes | Multivariate | Classification | Real | 540 | 18 | 2013 |
| Congressional Voting Records | Multivariate | Classification | Categorical | 435 | 16 | 1987 |
| Connectionist Bench (Sonar, Mines vs. Rocks) | Multivariate | Classification | Real | 208 | 60 | |

Data 선택조건

1. Classification    3. # of Attribute 10 to 1000

2. Multivariate     4. # of Instance (100 to 1000, Greater than 1,000)

# Install

# Install

# Explorer

# Explorer

**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: iris
Instances: 150
Attributes: 5
Sum of weights: 150

**Selected attribute**

Name: sepallength
Missing: 0 (0%)
Distinct: 35
Type: Numeric
Unique: 9 (6%)

| Statistic | Value |
|-----------|-------|
| Minimum | 4.3 |
| Maximum | 7.9 |
| Mean | 5.843 |
| StdDev | 0.828 |

**Attributes**

All | None | Invert | Pattern

| No. | Name |
|-----|------|
| 1 | sepallength |
| 2 | sepalwidth |
| 3 | petallength |
| 4 | petalwidth |
| 5 | class |

Remove

Class: class (Nom) | Visualize All

**Status**

OK | Log

---

View in main window
View in separate window
Save result buffer
Delete result buffer

Load model
Save model
Re-evaluate model on current test set
Re-apply this model's configuration

Visualize classifier errors
**Visualize tree**
Visualize margin curve
Visualize threshold curve ▶
Cost/Benefit analysis ▶
Visualize cost curve ▶

---

**Weka Explorer** ①

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

**Classifier**

② Choose | J48 -C 0.25 -M 2

**Test options**

③
- Use training set
- Supplied test set | Set...
- Cross-validation Folds 10
- Percentage split % 66

More options...

(Nom) class

④ Start | Stop

**Result list (right-click for options)**

⑤ 15:37:02 - trees.J48

**Classifier output**

```
=== Classifier model (full training set) ===

J48 pruned tree
------------------

petalwidth <= 0.6: Iris-setosa (50.0)
petalwidth > 0.6
|   petalwidth <= 1.7
|   |   petallength <= 4.9: Iris-versicolor (48.0/1.0)
|   |   petallength > 4.9
|   |   |   petalwidth <= 1.5: Iris-virginica (3.0)
|   |   |   petalwidth > 1.5: Iris-versicolor (3.0/1.0)
|   petalwidth > 1.7: Iris-virginica (46.0/1.0)

Number of Leaves  :     5

Size of the tree :     9
```

**Status**

OK | Log | x 0

---

**Weka Classifier Tree Visualizer: 15:37:02 - trees.J48 (iris)**

**Tree View**

petalwidth
- <= 0.6 → Iris-setosa (50.0)
- > 0.6 → petalwidth
  - <= 1.7 → petallength
    - <= 4.9 → Iris-versicolor (48.0/1.0)
    - > 4.9 → petalwidth
      - <= 1.5 → Iris-virginica (3.0)
      - > 1.5 → Iris-versicolor (3.0/1.0)
  - > 1.7 → Iris-virginica (46.0/1.0)

# Experimenter

# Experimenter

# Knowledge flow



Double click

# Knowledge flow

# Knowledge flow

# Knowledge flow

# Knowledge flow

# Knowledge flow

# Knowledge flow

# Knowledge flow

# Weka Reference

http://statweb.stanford.edu/~lpekelis/13_datafest_cart/WekaManual-3-7-8.pdf

http://software.ucv.ro/~eganea/AIR/KnowledgeFlowTutorial-3-5-8.pdf

https://www.youtube.com/watch?v=bPrTeUAS6_I&list=PLJbE6j2EG1pZnBhOg3_Rb63WLCprtyJag&index=26