

CSC662 ปฏิบัติการ ๑๑

การใช้ผังการไหลของความรู้ Knowledge Flow

เขียนโดย ผศ. ดร. กรุง สีนอกิรมย์สราญ

ภาควิชาคณิตศาสตร์ คณะวิทยาศาสตร์ จุฬาลงกรณ์มหาวิทยาลัย

เนื้อหาที่ครอบคลุม

- แนะนำซอฟต์แวร์ Weka
- ฐานข้อมูลเชิงสัมพันธ์และการเตรียมข้อมูล
- สถิติและการทำเหมืองข้อมูลแบบอธิบาย
- การทำเหมืองข้อมูลแบบกฎเชื่อมโยง
- การทำเหมืองข้อมูลแบบจำแนกประเภท
- การทำเหมืองข้อมูลแบบการวิเคราะห์การเกาะกลุ่ม
- **การใช้ผังการไหลของความรู้ Knowledge flow**

2

ผังการไหล (Knowledge flow)

นิยามผังการไหลของความรู้

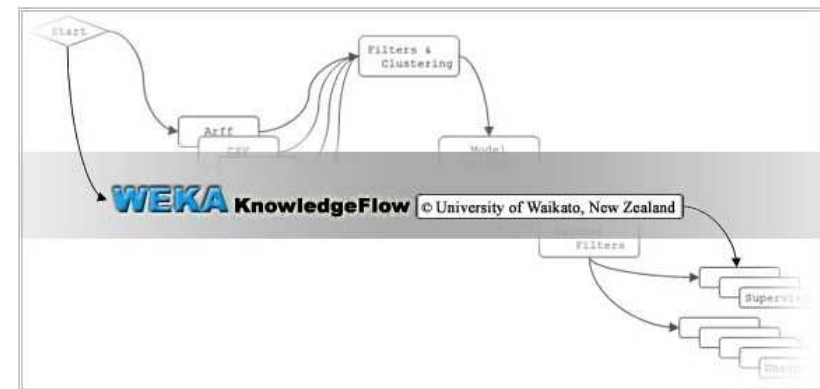
- ผังการไหลของความรู้ คือแผนภาพที่แสดงถึงการได้มาของความรู้ (Knowledge) โดยผ่านกระบวนการ/ขั้นตอนวิธี/การแสดงผลที่ใช้ในการทำเหมืองข้อมูล
- แผนภาพที่สร้างใช้สัญลักษณ์ (Icon) ที่สื่อถึงกระบวนการ/ขั้นตอนวิธี/การแสดงผลหนึ่งลักษณะหรือหนึ่งแบบ
- เส้นที่เชื่อมโยงระหว่างสัญลักษณ์แสดงการไหลของข้อมูล (data) ที่ผ่านกระบวนการ (icon) จนถึงความรู้ที่ได้

ตัวอย่าง การไหลของข้อมูลเพื่อให้ได้ความรู้ DataSources → Filter → Classifier → Evaluator → Visualization

3

ผังการไหล (Knowledge flow)




หน้าจอเริ่มต้นของผังการไหลของความรู้



4

ผังการไหล (Knowledge flow)



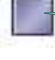
รายการของฟังก์ชันไหลของความรู้

-  **New layout**
 - New layout สร้างฟังก์ชันไหลใหม่
-  **Save layout**
 - Save layout เก็บฟังก์ชันที่สร้างไว้ใน Knowledge Flow Layout บันทึกลง เพิ่มข้อมูลเพื่อนำกลับมาใช้ใหม่
-  **Open layout**
 - Open layout เปิดเพิ่มข้อมูลทีเก็บฟังก์ชันที่สร้างแล้ว เพื่อนำกลับมาใช้ใหม่

5

ฟังก์ชันไหล (Knowledge Flow)

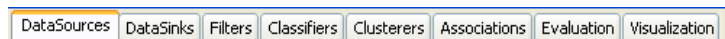
รายการของฟังก์ชันไหลของความรู้

-  **Selection**
 - Selection เปลี่ยนเมาส์ให้เป็นตัวชี้เพื่อเลือกภาพสัญลักษณ์หรือเส้นเชื่อม
-  **Display help**
 - Display help แสดงข้อความอธิบายการใช้เครื่องมือของฟังก์ชันไหลของความรู้
-  **Stop all execution**
 - Stop all execution หยุดการประมวลผลทุกอย่างที่เกิดขึ้น

6

ฟังก์ชันไหล (Knowledge Flow)

ส่วนประกอบหลักของฟังก์ชันไหลของความรู้

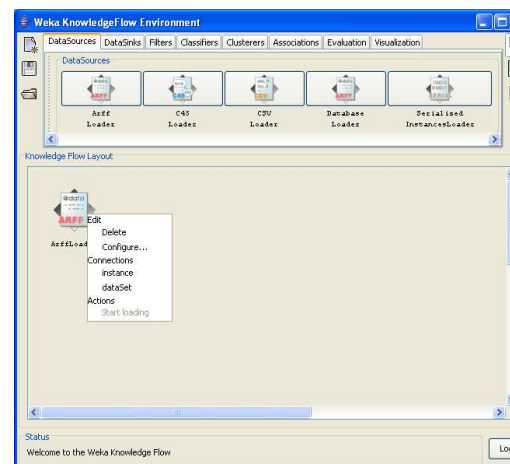


- DataSources: กำหนดแหล่งข้อมูลที่อ่านเข้าฟังก์ชัน
- DataSinks: กำหนดการบันทึกข้อมูลหรือจุดสุดท้ายของกระบวนการ
- Filters: ขั้นตอนการจัดการเตรียมข้อมูล
- Classifiers: การสร้างตัวแบบและวิธีการในการจัดจำแนกประเภท
- Clusterers: การใช้ขั้นตอนวิธีการวิเคราะห์การเกาะกลุ่ม
- Associations: การใช้ขั้นตอนวิธีการหาความสัมพันธ์
- Evaluation: ประเมินและแบ่งชุดข้อมูลออกเป็นส่วน ๆ
- Visualization: สำหรับแสดงผลลัพธ์ด้วยภาพนามธรรม

7

ฟังก์ชันไหล (Knowledge Flow)

ตัวอย่างการสร้างฟังก์ชันอ่านข้อมูล

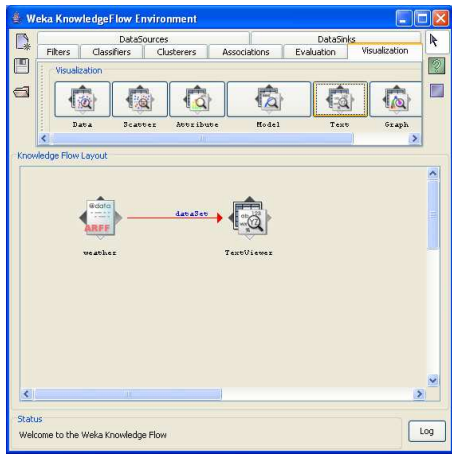


- เริ่มจากเลือกแถบ DataSources
- เลือกสัญลักษณ์ ArffLoader เมาส์เปลี่ยนเป็นเครื่องหมายกากบาท
- กดเมาส์ในบริเวณ Knowledge Flow Layout
- กดเมาส์ปุ่มขวาที่ ArffLoader เลือก Configure...
- เลือกแฟ้มที่ชื่อ weather.arff

8

ฟังก์ชันไหล (Knowledge Flow)

ตัวอย่างผังการไหลที่แสดงข้อความของข้อมูล

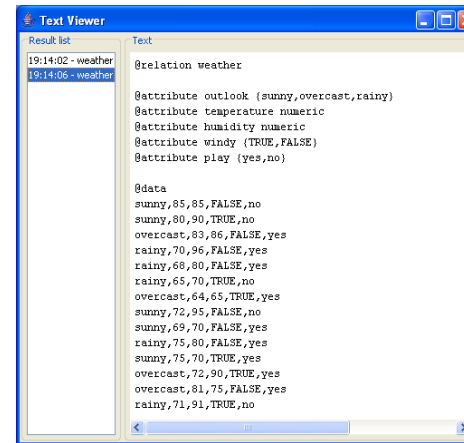


- เลือกแถบ Visualization
- กดเมาส์ที่ Text Viewer เมาส์ เปลี่ยนเป็นเครื่องหมายกากบาท
- กดเมาส์ปุ่มขวาเพื่อเชื่อม ArffLoader ไปยัง TextViewer โดยเลือก dataSet บนเมนูของ weather.arff
- กดเมาส์ปุ่มขวาที่ ArffLoader โดยเลือก Start loading ได้รายการ Action

9

ฝึกการไหล (Knowledge flow)

การแสดงผลของที่สัญลักษณ์ภาพนามธรรม

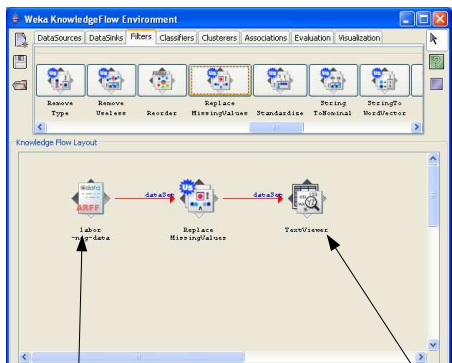


- หลังจาก เลือก Start loading ได้รายการ Action ของ ArffLoader
- ชื่อแฟ้มที่อ่านจะปรากฏได้ภาพ ArffLoader
- แสดงข้อความโดยเลือก Show results ภายได้รายการใน TextViewer โดยกดเมาส์ปุ่มขวา
- ผลลัพธ์ที่ได้แสดงทางภาพด้านซ้าย

10

ฝึกการไหล (Knowledge flow)

การเพิ่มขั้นตอนในการกรอง Missing value

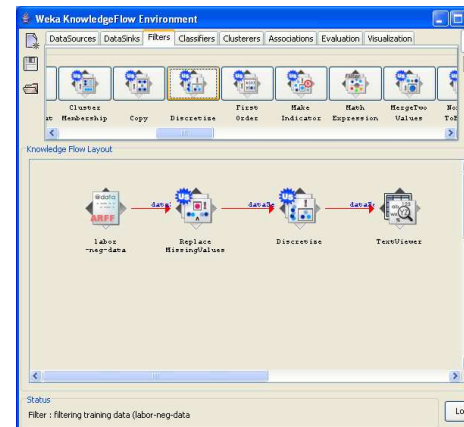


- เริ่มจาก DataSources โดยใช้ ArffLoader
- เลือก Configure... แล้วเลือกแฟ้ม labor.arff
- เลือกแถบ Filters ที่เรียก Replace Missing Values เพื่อเติมค่าที่หายไป
- เลือกแถบ Visualization แล้วเลือก TextViewer เพื่อแสดงผลลัพธ์

กดเมาส์ปุ่มขวาที่ ArffLoader แล้วเลือก Start loading

กดเมาส์ปุ่มขวาที่ TextViewer แล้วเลือก Show results

การเพิ่มขั้นตอน Discretization

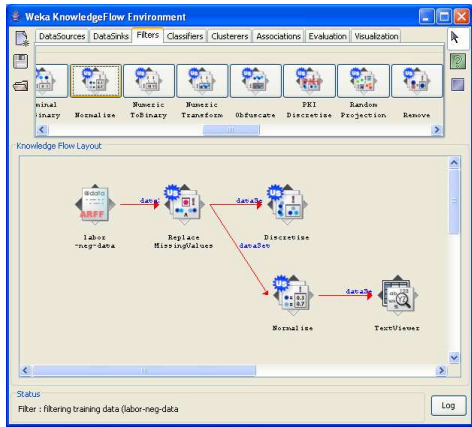


- เริ่มจาก DataSources โดยใช้ ArffLoader
- เลือก Configure... แล้วเลือกแฟ้ม labor.arff
- เลือกแถบ Filters แล้วเลือก Discretize เพื่อเปลี่ยนตัวแปรที่มีค่าต่อเนื่องเป็นตัวแปรที่มีค่าไม่ต่อเนื่อง
- เลือก TextViewer จากแถบ Visualization

12

ฝึกการไหล (Knowledge flow)

การเพิ่มขั้นตอน Normalization

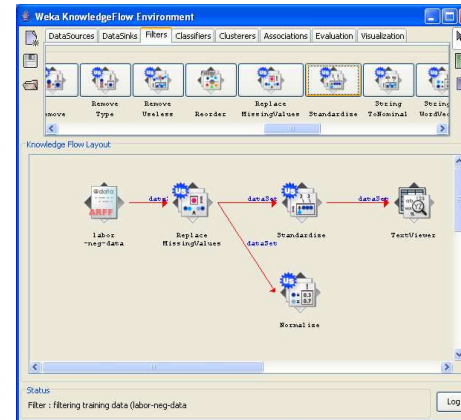


- เริ่มจาก DataSources โดยใช้ ArffLoader
- เลือก Configure... แล้วเลือกเพิ่ม labor.arff
- เลือกแถบ Filters ที่เรียก Normalize เพื่อแปลงตัวแปรที่มีค่าต่อเนื่องให้อยู่ในช่วง [0, 1]
- เลือกแถบ Visualization แล้วเลือก Text Viewer เพื่อแสดงผลลัพธ์

13

ฝึกการไหล (Knowledge flow)

การเพิ่มขั้นตอน Standardize

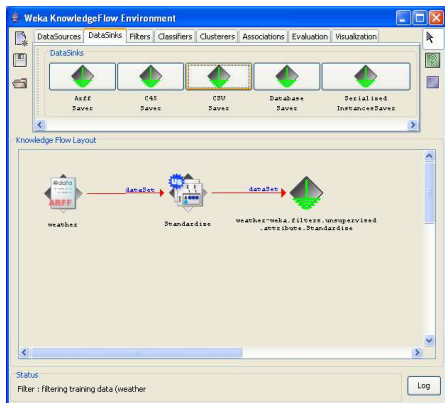


- เริ่มจาก DataSources โดยใช้ ArffLoader
- เลือก Configure... แล้วเลือกเพิ่ม labor.arff
- เลือกแถบ Filters ที่เรียก Standardize เพื่อแปลงตัวแปรให้มีค่าตกอยู่ในช่วงของการกระจายแบบปรกติมาตรฐาน
- เลือกแถบ Visualization แล้วเลือก Text Viewer เพื่อแสดงผลลัพธ์

14

ฝึกการไหล (Knowledge flow)

การบันทึกข้อมูลลงเพิ่ม CSV

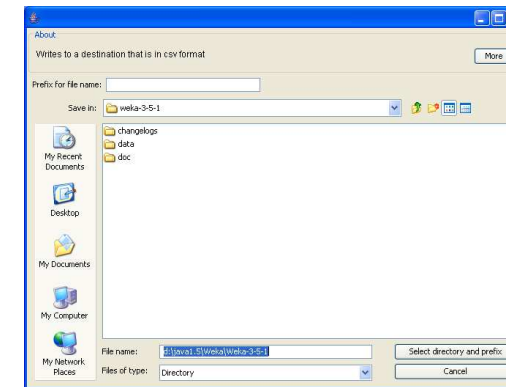


- ซอฟต์แวร์ Weka สามารถแปลงข้อมูลให้อยู่ในรูปแบบ csv เพื่อนำไปใช้กับซอฟต์แวร์อื่น เช่น Calc
- เริ่มจากเลือก ArffLoader ใน DataSources
- แปลงข้อมูลให้เหมาะสม
- เลือกแถบ DataSinks แล้วเลือก CSVsaver
- เลือกเพิ่มข้อมูล Arff ที่ต้องการ แล้วเลือก Start Loading

15

ฝึกการไหล (Knowledge flow)

การบันทึกลงเพิ่ม CSV ต่อ



- เลือก Configure... ในเมนูของ CSVsaver
- เปลี่ยนสถานที่ที่ต้องการเก็บไปตำแหน่งที่ต้องการเก็บ โดยเพิ่ม prefix ให้กับชื่อเพิ่มที่ต้องการ
- เก็บข้อมูลโดยเลือก Start loading ใน ArffLoader

16

ฝึกการไหล (Knowledge flow)

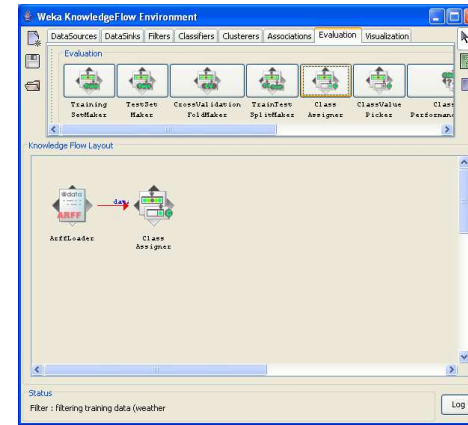
โครงการทำเหมืองข้อมูลโดยใช้ฟังก์ชันไหล

- กำหนดข้อมูลในแฟ้ม iris.arff ให้หาตัวแบบต้นไม้การตัดสินใจที่ดีที่สุด โดยใช้ 5 fold cross-validation กับขั้นตอนวิธี J48 แสดงผลลัพธ์ที่ได้ในรูปแบบต้นไม้
- แนวทางวางฟังก์ชันไหล: DataSources → Evaluation → J48 → Visualization
 - เริ่มจากการอ่านแฟ้ม iris.arff
 - กำหนดลักษณะประจำที่ใช้แทนคลาส
 - แบ่งข้อมูลออกเป็น 5 ส่วนเพื่อทำ cross-validation
 - ใช้ขั้นตอนวิธี J48
 - แสดงผลลัพธ์

17

ฟังก์ชันไหล (Knowledge Flow)

การอ่าน iris.arff

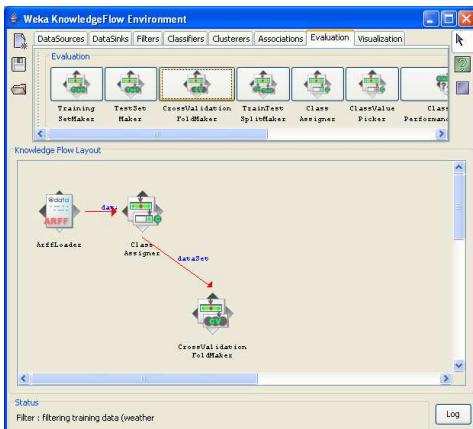


18

ฟังก์ชันไหล (Knowledge Flow)

- เลือก ArffLoader จากแถบ DataSources
- ปรับแต่งให้เลือกแฟ้ม iris.arff จาก Configure... เมนู
- เลือก Class Assignment จากแถบ Evaluation
- เลือกคลาสเป้าหมาย

การแยกออกเป็น k-fold cross validation

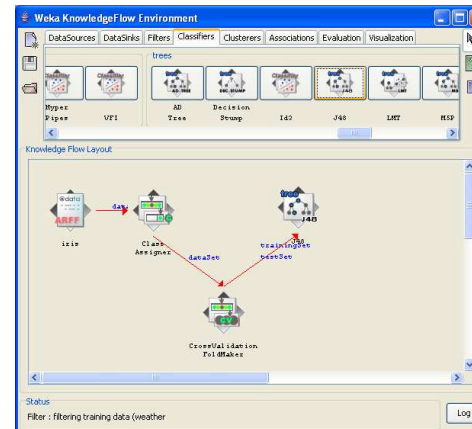


19

ฟังก์ชันไหล (Knowledge Flow)

- เลือก Cross Validation FoldMaker จากแถบ Evaluation
- ปรับแต่งให้มีจำนวน fold เท่ากับ 5
- ส่งข้อมูล DataSet จาก Class Assigment

การเรียกใช้ขั้นตอนวิธี J48

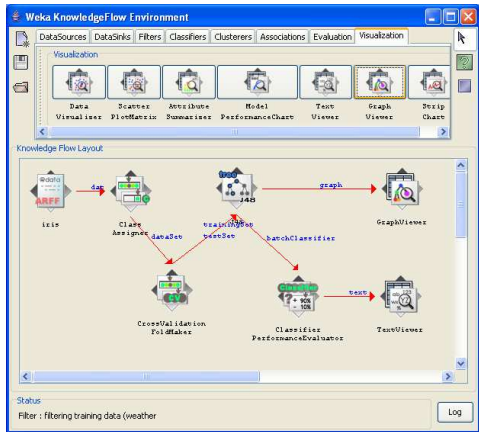


20

ฟังก์ชันไหล (Knowledge Flow)

- เลือก J48 จากแถบ Classifiers
- โยงข้อมูลจาก Cross Validation FoldMaker โดยเลือก training set และ test set โยงไปที่ J48
- สั่งให้ข้อมูลนำเข้า โดยเลือก Start loading จาก ArffLoader

การแสดงผลที่ได้จากผังการไหล



21

ผังการไหล (Knowledge Flow)

- เลือก Classifiers PerformanceEvaluators จากแถบ Evaluation
- โยงข้อมูลจาก J48 โดยเลือก batchClassifiers โยงไปที่ Classifiers PerformanceEvaluators
- สร้าง TextViewer และ/หรือ GraphViewer จาก Visualization

ผลลัพธ์ที่ได้ในรูปเนื้อความของ J48

```

Result Set: Text
23:34:52 - J48

=== Evaluation result ===
Scheme: J48
Relation: iris

Correctly Classified Instances      144      96 %
Incorrectly Classified Instances     5       4 %
Kappa statistic                    0.94
Mean absolute error                 0.035
Root mean squared error             0.1582
Relative absolute error             7.8842 %
Root relative squared error        33.5577 %
Total Number of Instances          150

=== Detailed Accuracy By Class ===
TP Rate  FP Rate  Precision  Recall  F-Measure  ROC Area  Class
0.98  0      1      0.98  0.99  0.99  Iris-setosa
0.94  0.03  0.94  0.94  0.94  0.958  Iris-versicolor
0.96  0.03  0.941  0.96  0.95  0.966  Iris-virginica

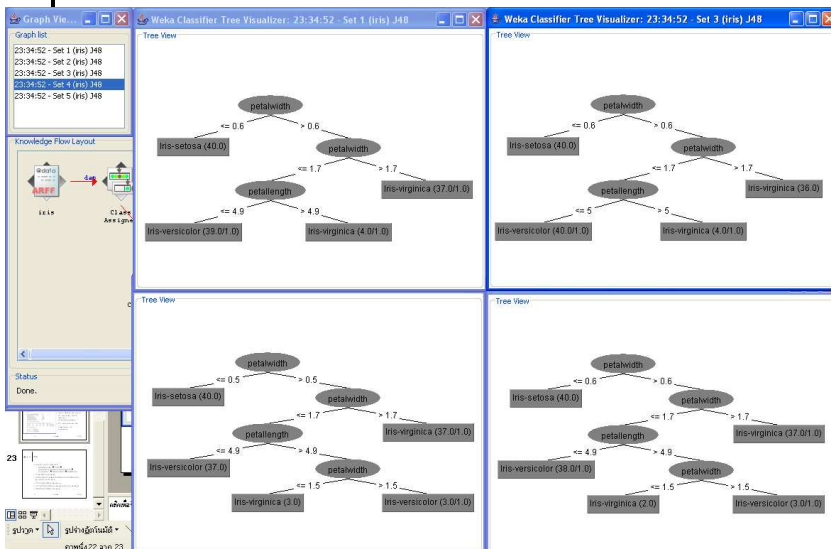
=== Confusion Matrix ===
 a b c <-- classified as
49 1 0 | a = Iris-setosa
 0 47 3 | b = Iris-versicolor
 0 2 40 | c = Iris-virginica
    
```

22

ผังการไหล (Knowledge Flow)

- เลือก Start Loading จากแถบ ArffLoader
- รอจนจบ แล้วเลือก Show results จาก TextViewer
- ผลลัพธ์ที่ได้แสดงดังรูปซ้ายซึ่งให้ค่าที่ถูกต้อง 96%
- ใน Confusion Matrix แสดงผลจากการเปรียบเทียบกับกลุ่มที่สนใจ

ผลลัพธ์ในรูปต้นไม้การตัดสินใจ



23

สรุป

- การออกแบบผังการไหลโดยปรกติ
 - DataSourcees → Filter → Classifier/Clusterers/Associations → Evaluation → Visualization → DataSinks
- Filter ใช้ในการเตรียมข้อมูล
- Classifier/Clusterers/Associations ใช้ในการสร้างตัวแบบในการทำเหมืองข้อมูล
- Evaluation ใช้ในการเลือกตัวแบบ
- Visualization ใช้ในการแสดงผลลัพธ์ของการทำเหมืองข้อมูล
- DataSinks ใช้ในการเก็บผลลัพธ์

24

ผังการไหล (Knowledge Flow)