

# CSC662 ปฏิบัติการ ๖

ตัวกรองที่ใช้ในซอฟต์แวร์ Weka

เขียนโดย ผศ. ดร. กรุง สีนอกภิรมย์สรราชู

ภาควิชาคณิตศาสตร์ คณะวิทยาศาสตร์ จุฬาลงกรณ์มหาวิทยาลัย

## เนื้อหาที่ครอบคลุม

- ตัวกรองในซอฟต์แวร์ Weka
- การกำจัดลักษณะประจำ
- ตัวกรองแบบอัตโนมัติ Supervised filter
  - ลักษณะประจำ
  - ระเบียบ
- ตัวกรองที่ผู้ใช้กำหนดเอง Unsupervised filter
  - ลักษณะประจำ
  - ระเบียบ

2

Weka filter

## ตัวอย่างเพิ่มข้อมูล sample01.csv

ID,SEX,PASS/FAIL,Score,Class

- 1,M,Pass,45.5,B
- 2,F,Pass,56.78,B
- 3,M,Pass,89,A
- 4,F,Pass,77,A
- 5,M,Fail,32,C
- 6,F,Fail,12,D
- 7,M,Fail,35,C
- 8,F,Pass,62,B
- 9,M,Pass,68,B+
- 10,F,Fail,10,D

3

Weka filter

หน้าต่างที่ได้จากการเปิดเพิ่ม sample01.csv

ปุ่มเลือกโมดูล Filter เพื่อเตรียมข้อมูล

ลักษณะทางสถิติที่สำคัญของตัวแปรที่เลือก

Label	Count
M	5
F	5

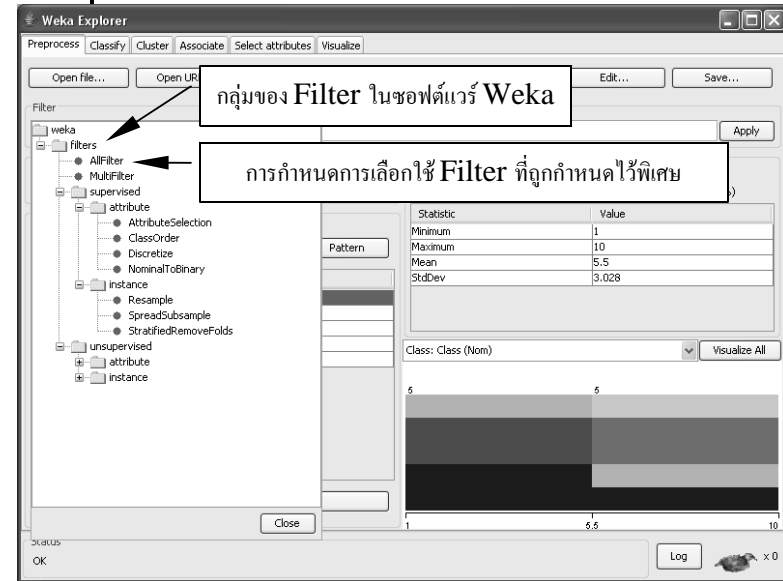
Class: Class (Nom) Visualize All

Status OK Log x 0

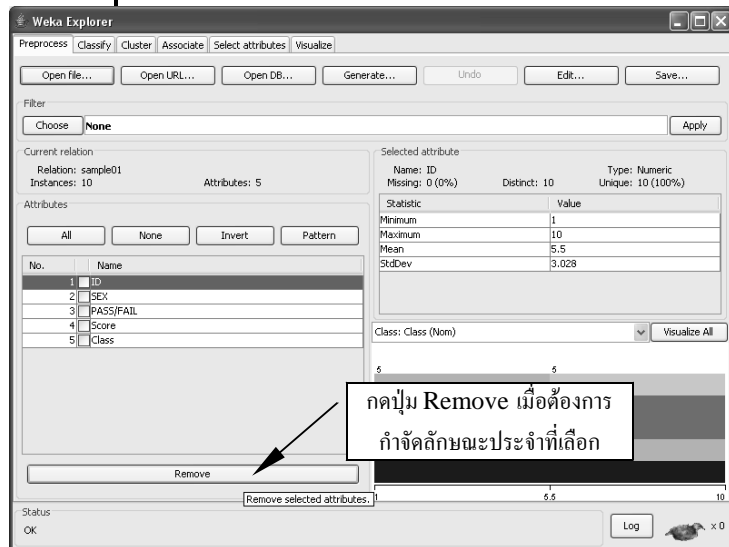
# การใช้ตัวกรองในซอฟต์แวร์ Weka

- ตัวกรอง (Filters) รวบรวมโมดูลในชั้นการเตรียมข้อมูล
- ตัวกรองแบ่งออกเป็นสองลักษณะคือ
  - Supervised รวมโมดูลที่แปลงข้อมูลแบบอัตโนมัติที่มีการควบคุมด้วยพารามิเตอร์ที่ผู้ใช้กำหนด แบ่งเป็นสองหมวดใหญ่คือลักษณะประจำ (attribute) กับข้อมูลแต่ละระเบียน (instance)
  - Unsupervised รวมโมดูลที่แปลงข้อมูลที่ผู้ใช้กำหนดเอง แบ่งเป็นสองหมวดใหญ่คือ ลักษณะประจำ (attribute) กับข้อมูลแต่ละระเบียน (instance)

# ตัวกรอง



# การกำจัดลักษณะประจำ

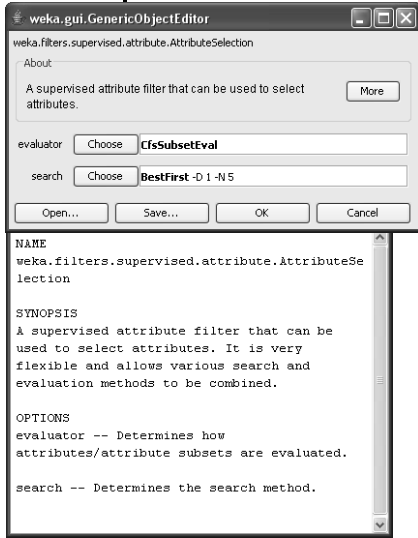


- เรากำจัดลักษณะประจำที่ไม่ต้องการออก โดยทำการสั่งหามาถูกหน้าลักษณะประจำที่ต้องการ แล้วกดปุ่ม Remove

# ตัวกรองแบบอัตโนมัติ Supervised

- ประกอบด้วย
  - ลักษณะประจำ: AttributeSelection, ClassOrder, Discretize, NominalToBinary
  - ระเบียน: Resample, SpreadSubsample, StratifiedRemoveFolds

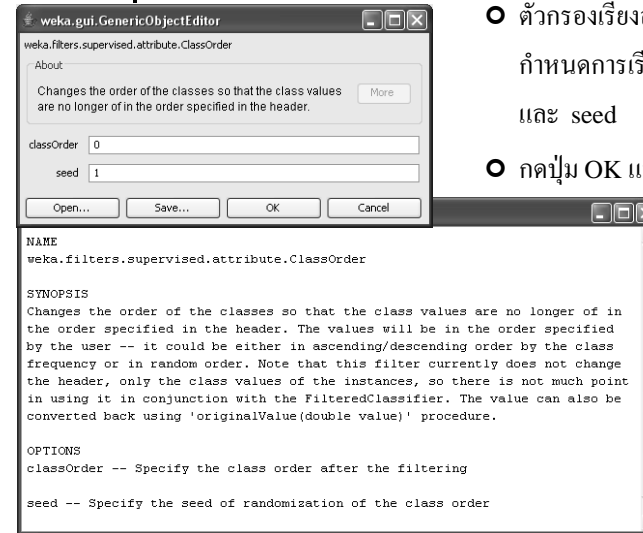
# AttributeSelection



Weka filter

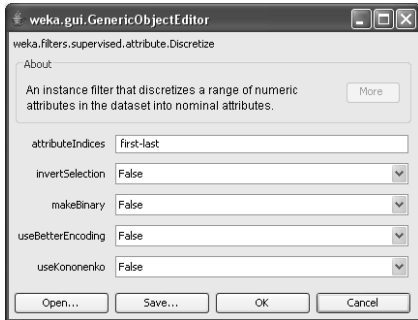
- ตัวกรองที่เลือกลักษณะประจำที่นำมาวิเคราะห์แบบอัตโนมัติ โดยผู้ใช้กำหนดตัวประเมินในกล่อง evaluator และวิธีการค้นในกล่อง search
- กดปุ่ม OK
- แล้วกดปุ่ม Apply

# ClassOrder



- ตัวกรองเรียงลำดับคลาส โดยผู้ใช้กำหนดการเรียงในกล่อง classOrder และ seed
- กดปุ่ม OK แล้วกดปุ่ม Apply

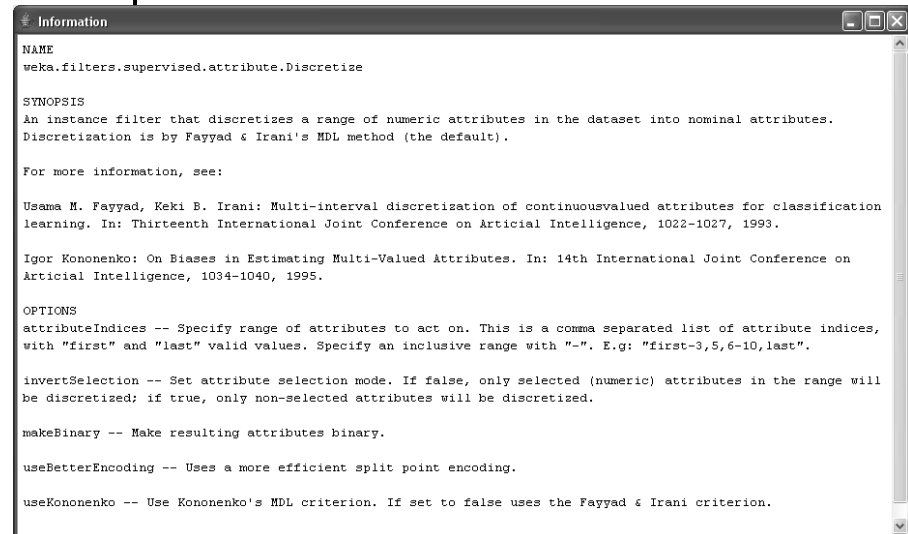
# Discretize



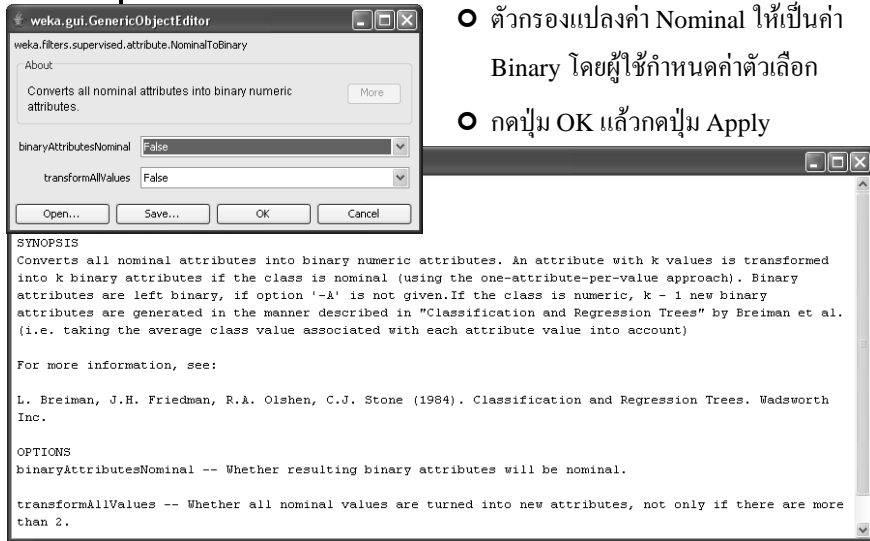
Weka filter

- ตัวกรองแปลงค่าต่อเนื่องให้เป็นค่าไม่ต่อเนื่อง โดยผู้ใช้เลือกลักษณะประจำที่ต้องการเปลี่ยนในกล่อง attributeIndices และผู้ใช้กำหนดตัวเลือกโดยดู Help ในหน้าต่างไป
- กดปุ่ม OK
- แล้วกดปุ่ม Apply

# Discretize Help

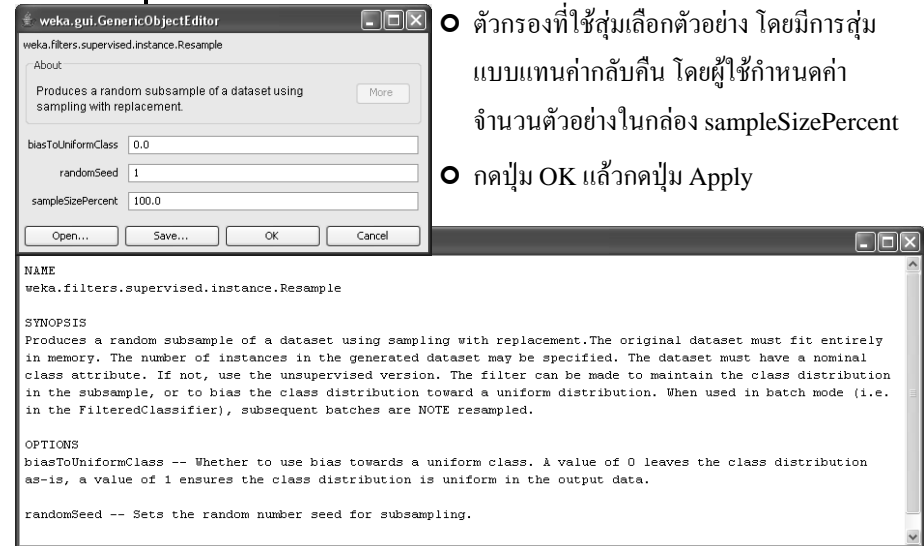


# NominalToBinary



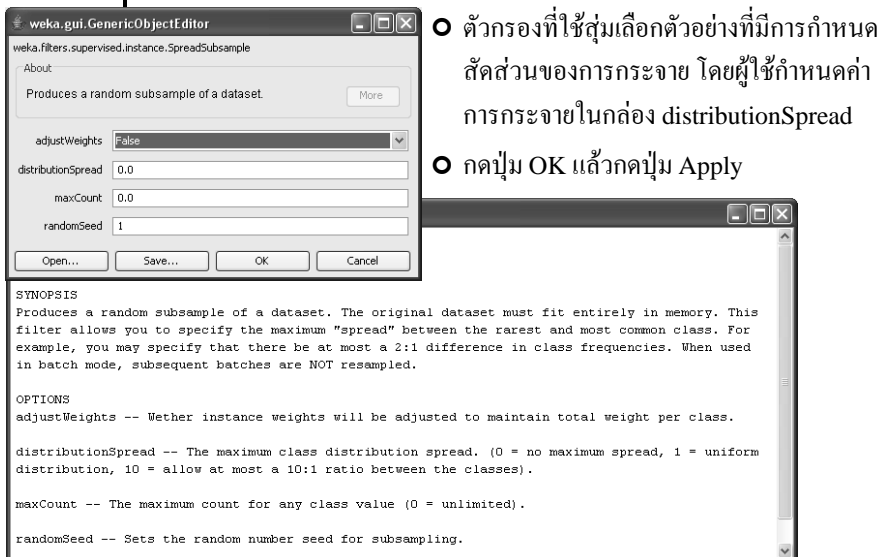
- ตัวกรองแปลงค่า Nominal ให้เป็นค่า Binary โดยผู้ใช้กำหนดค่าตัวเลือก
- กดปุ่ม OK แล้วกดปุ่ม Apply

# Resample



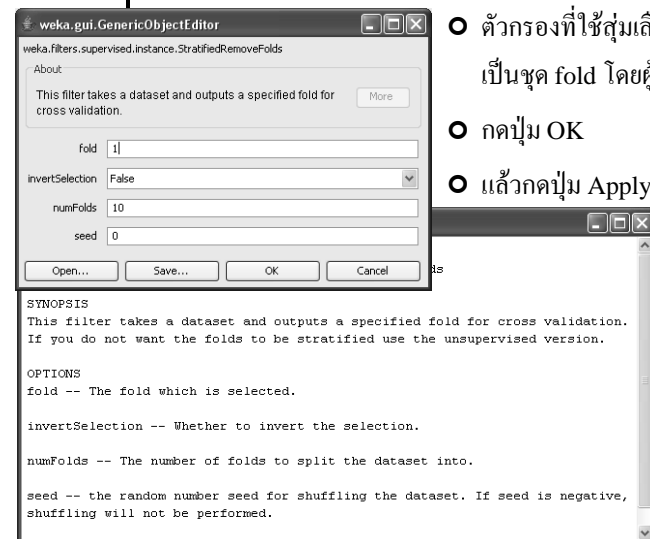
- ตัวกรองที่ใช้สุ่มเลือกตัวอย่าง โดยมีการสุ่มแบบแทนค่ากลับคืน โดยผู้ใช้กำหนดค่าจำนวนตัวอย่างในกล่อง sampleSizePercent
- กดปุ่ม OK แล้วกดปุ่ม Apply

# SpreadSubsample



- ตัวกรองที่ใช้สุ่มเลือกตัวอย่างที่มีการกำหนดสัดส่วนของการกระจาย โดยผู้ใช้กำหนดค่าการกระจายในกล่อง distributionSpread
- กดปุ่ม OK แล้วกดปุ่ม Apply

# StratifiedRemoveFolds



- ตัวกรองที่ใช้สุ่มเลือกกลุ่มตัวอย่างออกเป็นชุด fold โดยผู้ใช้กำหนดตัวเลือก
- กดปุ่ม OK
- แล้วกดปุ่ม Apply

## ตัวกรองที่ผู้ใช้กำหนดเอง Unsupervised

- เราจะเลือกอธิบายตัวกรองบางตัวเท่านั้น สำหรับตัวกรองอื่น ผู้ใช้สามารถอ่านได้จาก **Help** ของซอฟต์แวร์ **Weka**
- ลักษณะประจำ: Add, AddCluster, AddExpression, AddNoise, ClusterMembership, Copy, Discretize, FirstOrder, MakeIndicator, MergeTwoValues, NominalToBinary, Normalize, NumericToBinary, NumericTransform, Obfuscate, PKIDiscretize, RandomProjection, Remove, RemoveType, RemoveUseless, ReplaceMissingValues, Standardize, StringToNominal, StringToWordVector, SwapValues, TimeSeriesDelta, TimeSeriesTranslate
- ระเบียบ: Normalize, NonSparseToSparse, Randomize, RemoveFolds, RemoveMisclassified, RemovePercentage, RemoveRange, RemoveWithValues, Resample, SparseToNonSparse

17

Weka filter

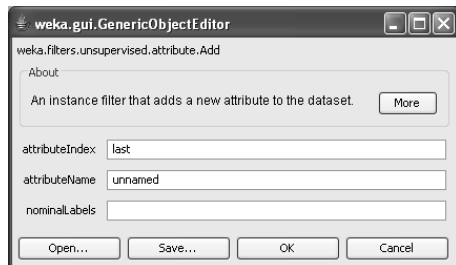
## ตัวกรองที่ผู้ใช้กำหนดเองกับลักษณะประจำ

- Add filter
- AddExpression filter
- NominalToBinary filter
- NumericToBinary filter
- NumericTransform filter
- Remove filter
- ReplaceMissingValues filter
- Standardize filter
- AddCluster filter
- Discretize filter
- Normalize filter
- RemoveType filter

18

Weka filter

## Add filter

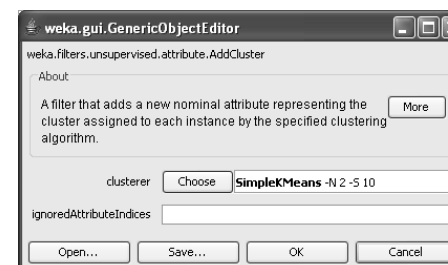


- ตัวกรองเพิ่มลักษณะประจำ เลือก add โดยเพิ่มลักษณะประจำที่มีค่าตั้งต้นคือ missing value
- กดปุ่ม OK
- แล้วกดปุ่ม Apply

19

Weka filter

## AddCluster filter

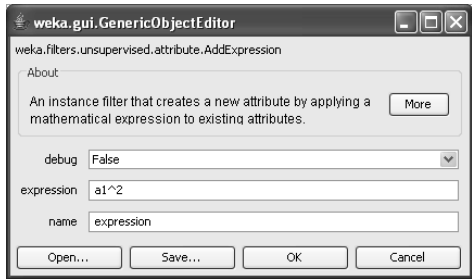


- ตัวกรองเพิ่มลักษณะประจำตามการเกาะกลุ่ม เลือก addCluster เลือกวิธีเกาะกลุ่ม เช่น SimpleKMeans
- กำหนดลักษณะประจำที่ไม่นำมาใช้ในการวิเคราะห์การเกาะกลุ่มใน ignoredAttributeIndices
- กดปุ่ม OK
- แล้วกดปุ่ม Apply

20

Weka filter

# AddExpression filter

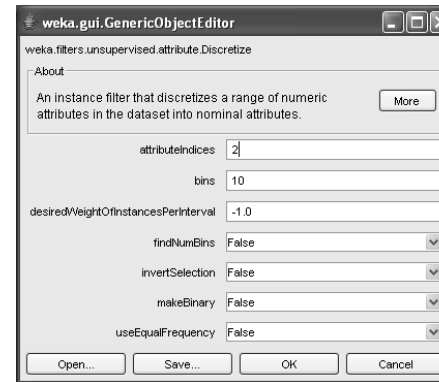


- ตัวกรองเพิ่มลักษณะประจำตามนิพจน์ จากลักษณะประจำที่กำหนด เลือก addExpression พิมพ์นิพจน์ที่ต้องการ สร้างลักษณะประจำใหม่
- กดปุ่ม OK
- แล้วกดปุ่ม Apply

21

Weka filter

# Discretize filter

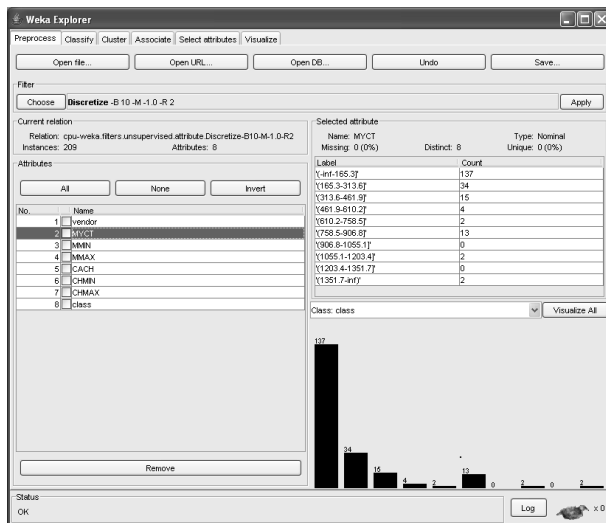


- ตัวกรองแปลงเป็นค่าไม่ต่อเนื่อง ผู้ใช้เลือกลักษณะประจำในช่อง attributeIndices ตามลำดับลักษณะประจำที่กำหนด
- กำหนดจำนวนกล่องที่ต้องการใน bins
- เราสามารถแบ่งแบบ equal width หรือ equal depth โดยปรับเป็น False ที่ useEqualFrequency
- กดปุ่ม OK
- แล้วกดปุ่ม Apply

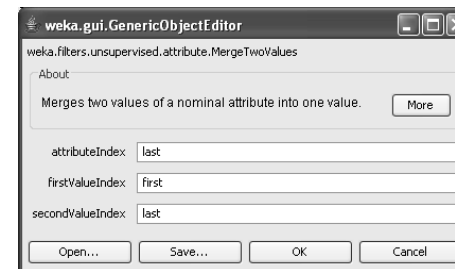
22

Weka filter

# ผลการใช้ตัวกรอง Discretize



# MergeTwoValues filter

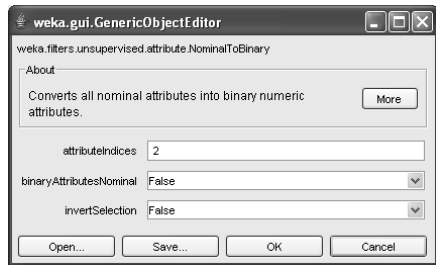


- ตัวกรองรวมค่าสองค่าเป็นหนึ่ง เลือก MergeTwoValues
- กำหนดดัชนีของลักษณะประจำใน attributeIndex
- กำหนดค่าใน firstValueIndex และ secondValueIndex
- กดปุ่ม OK
- แล้วกดปุ่ม Apply

24

Weka filter

# NominalToBinary filter

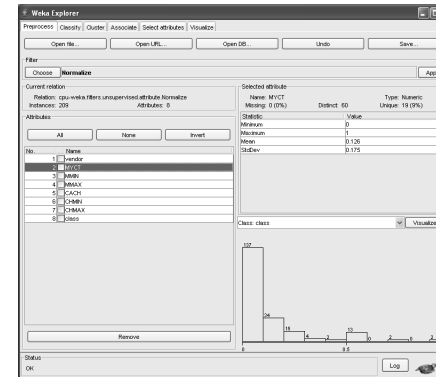


- เลือกตัวกรองแปลงค่าไม่ต่อเนื่องเป็นค่า 0 หรือ 1 เลือก NominalToBinary
- กำหนดครรหระของลักษณะประจำใน attributeIndices ที่ต้องการ
- กดปุ่ม OK
- แล้วกดปุ่ม Apply

25

Weka filter

# Normalize filter



- ตัวกรองเปลี่ยนเป็นค่ามาตรฐานเลือก Normalize เพื่อปรับลักษณะประจำทุกลักษณะประจำเฉพาะลักษณะประจำที่เป็นจำนวน จะถูกแปลงให้มีค่าอยู่ในช่วง 0-1 โดยใช้สูตร

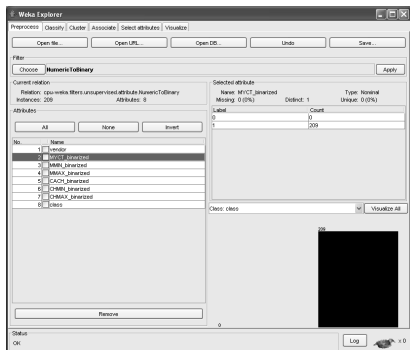
$$v'_i = \frac{v_i - \min_A}{\max_A - \min_A}$$

- กดปุ่ม Apply

26

Weka filter

# NumericToBinary filter

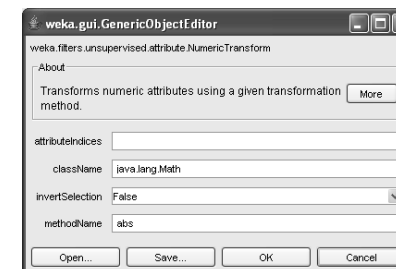


- ตัวกรองแปลงข้อมูลจำนวนให้เป็นค่า 0 หรือ 1 เลือก NumericToBinary โดยเปลี่ยนทุกลักษณะประจำที่เป็นจำนวน ค่าจำนวนที่เป็น 0 จะยังคงค่า 0 แต่ค่าที่ไม่ใช่ 0 จะเปลี่ยนเป็น 1 ทั้งหมด
- กดปุ่ม Apply

27

Weka filter

# NumericTransform filter

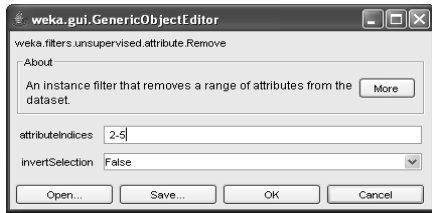


- ตัวกรองแปลงโดยใช้ฟังก์ชันจำนวนเลือก NumericTransform จะแปลงค่าในลักษณะประจำตามฟังก์ชันที่กำหนด เช่น abs
- กดปุ่ม OK
- แล้วกด Apply

28

Weka filter

# Remove filter

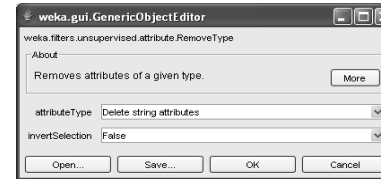


- ตัวกรองกำจัดลักษณะประจำ เลือก Remove โดย attributeIndices
- กดปุ่ม OK
- แล้วกดปุ่ม Apply

29

Weka filter

# RemoveType filter

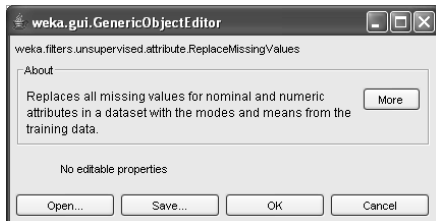


- ตัวกรองกำจัดลักษณะประจำตามชนิดของลักษณะประจำ เลือก RemoveType โดยเลือกชนิดที่ต้องการกำจัดใน attributeType
- กดปุ่ม OK
- แล้วกดปุ่ม Apply

30

Weka filter

# ReplaceMissingValue

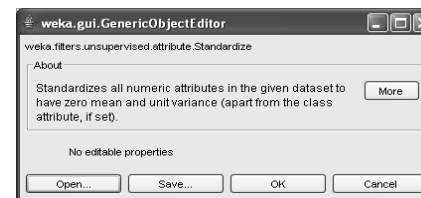


- ตัวกรองการแทนค่าที่ขาดหายไปเลือก ReplaceMissingValue
  - แทนด้วยค่าเฉลี่ยสำหรับลักษณะประจำที่เป็นจำนวน
  - แทนด้วยฐานนิยมสำหรับลักษณะประจำที่เป็นค่าไม่ต่อเนื่อง
- กดปุ่ม OK
- แล้วกดปุ่ม Apply

31

Weka filter

# Standardize filter



- ตัวกรองเปลี่ยนข้อมูลให้อยู่ในรูปที่มีการแจกแจงมาตรฐานโดยใช้ z-score โดยเลือก Standardize

$$v'_i = \frac{v_i - \bar{V}}{std_A}, \bar{V} = \sum_{i=1}^n \frac{v_i}{n}$$

$$std_A = \sum_{i=1}^n \frac{(v_i - \bar{V})^2}{n - 1}$$

- กดปุ่ม OK
- แล้วกดปุ่ม Apply

32

Weka filter



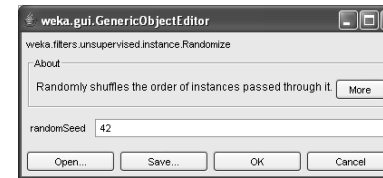
## ตัวกรองที่ตัวกรองที่ผู้ใช้กำหนดเองกับระเบียบ

- Randomize
- RemoveFolds
- RemovePercentage
- RemoveRange
- RemoveWithValues
- Resample

33

Weka filter

## Randomize filter

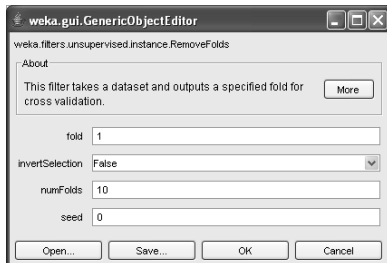


- ตัวกรองสลับสุม เลือก Randomize เพื่อให้มีการเรียงระเบียบแบบสุม
- กดปุ่ม OK
- แล้วกดปุ่ม Apply

34

Weka filter

## RemoveFold filter

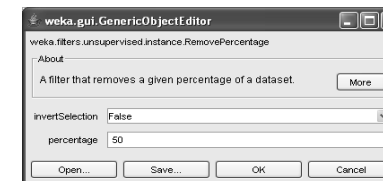


- ตัวกรองกำจัดชุดระเบียบ เลือก RemoveFold เพื่อกำจัดข้อมูลตามจำนวนชุด ตามจำนวนชุดทั้งหมดใน numFolds
- กดปุ่ม Save เพื่อบันทึกชุดระเบียบ
- กดปุ่ม OK
- แล้วกดปุ่ม Apply

35

Weka filter

## RemovePercentage filter

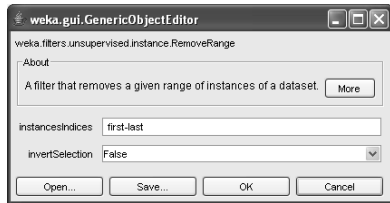


- ตัวกรองกำจัดระเบียบตามเปอร์เซ็นต์ เลือก RemovePercentage เพื่อลดจำนวนข้อมูลโดยเอาออกเท่ากับจำนวนเปอร์เซ็นต์ที่กำหนดใน percentage
- กดปุ่ม OK
- แล้วกดปุ่ม Apply

36

Weka filter

## RemoveRange filter

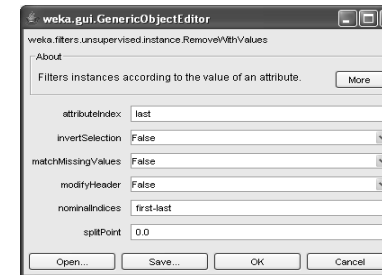


- ตัวกรองกำจัดระเบียบในพิสัยที่กำหนด เลือก RemoveRange เพื่อลดจำนวนข้อมูลที่กำหนดใน instancesIndices
- กดปุ่ม OK
- แล้วกดปุ่ม Apply

37

Weka filter

## RemoveWithValues filter

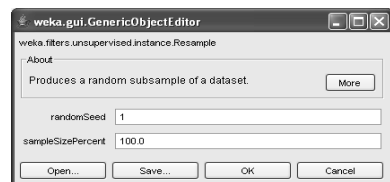


- ตัวกรองกำจัดข้อมูลตามค่า เลือก RemoveWithValues เพื่อลดจำนวนข้อมูลออกโดยใช้ attributeIndex
- ค่าที่ต่ำกว่า splitPoint จะถูกกำจัดทิ้ง
- กดปุ่ม OK
- แล้วกดปุ่ม Apply

38

Weka filter

## Resample filter



- ตัวกรองสุ่มใหม่ เลือก Resample เพื่อให้มีการสุ่มข้อมูลใหม่ โดยกำหนดเป็นเปอร์เซ็นต์ใน sampleSizePercent
- กดปุ่ม save เพื่อบันทึกข้อมูล
- กดปุ่ม OK
- แล้วกดปุ่ม Apply

39

Weka filter

## สรุป

- โมดูลในการเตรียมข้อมูลในซอฟต์แวร์ Weka เรียก ตัวกรอง (Filters) แบ่งออกเป็น
  - Supervised
  - Unsupervised
- นอกจากนี้เราเลือกใช้ตัวกรองกับลักษณะประจำ หรือระเบียบ

40

Weka filter