

PREGUNTA 2

De un dataset de su tarea anterior en WEKA, realice tres algoritmos de pre procesamiento.

NORMALIZACION

Usa una fórmula o un algoritmo para transformar las variables medidas en diferentes escalas en una escala común para que puedan ser comparables (manzanas con manzanas) o analizadas en un modelo estadístico elegido. Un típico ejemplo es calcular el logaritmo de las variables para hacer una distribución sesgada normal (por ejemplo, desplegado en un gráfico como una curva normal). En weka realizamos la normalización llevando a una escala de 0 a 1 para todos los atributos.

Weka Explorer

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

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

Filter: Choose **None** [Apply] [Stop]

Current relation
 Relation: train
 Instances: 2000
 Attributes: 21
 Sum of weights: 2000

Attributes
 All | None | Invert | Pattern

No.	Name
1	<input checked="" type="checkbox"/> battery_power
2	<input type="checkbox"/> blue
3	<input type="checkbox"/> clock_speed
4	<input type="checkbox"/> dual_sim
5	<input type="checkbox"/> fc
6	<input type="checkbox"/> four_g
7	<input type="checkbox"/> int_memory
8	<input type="checkbox"/> m_dep
9	<input type="checkbox"/> mobile_wt
10	<input type="checkbox"/> n_cores
11	<input type="checkbox"/> pc
12	<input type="checkbox"/> px_height
13	<input type="checkbox"/> px_width
14	<input type="checkbox"/> ram
15	<input type="checkbox"/> sc_h
16	<input type="checkbox"/> sc_w
17	<input type="checkbox"/> talk_time
18	<input type="checkbox"/> theta_0

[Remove]

Selected attribute
 Name: battery_power
 Missing: 0 (0%)
 Distinct: 1094
 Type: Numeric
 Unique: 491 (25%)

Statistic	Value
Minimum	501
Maximum	1998
Mean	1238.518
StdDev	439.418

Class: price_range (Num) [Visualize All]

Status
 OK [Log] x 0

Histogram of battery_power

Bin Range	Frequency
501 - 600	191
600 - 700	168
700 - 800	178
800 - 900	160
900 - 1000	159
1000 - 1100	165
1100 - 1200	161
1200 - 1300	156
1300 - 1400	172
1400 - 1500	160
1500 - 1600	161
1600 - 1700	169

Weka Explorer

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

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

Filter

- ☐ ClusterMembership
- ☐ Copy
- ☐ DateToNumeric
- ☐ Discretize
- ☐ FirstOrder
- ☐ FixedDictionaryStringToWordVector
- ☐ InterquartileRange
- ☐ KernelFilter
- ☐ MakeIndicator
- ☐ MathExpression
- ☐ MergeInfrequentNominalValues
- ☐ MergeManyValues
- ☐ MergeTwoValues
- ☐ NominalToBinary
- ☐ NominalToString
- ☒ Normalize
- ☐ NumericCleaner
- ☐ NumericToBinary
- ☐ NumericToDate
- ☐ NumericToNominal
- ☐ NumericTransform
- ☐ Obfuscate
- ☐ OrdinalToNumeric
- ☐ PartitionedMultiFilter
- ☐ PKIDiscretize
- ☐ PrincipalComponents

Attributes: 21
Sum of weights: 2000

Invert Pattern

Remove

Selected attribute

Name: battery_power
Missing: 0 (0%)
Distinct: 1094
Type: Numeric
Unique: 491 (25%)

Statistic	Value
Minimum	501
Maximum	1998
Mean	1238.518
StdDev	439.418

Class: price_range (Num)

Visualize All

Status
OK

Weka Explorer

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

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

Filter

Choose Normalize -S 1.0 -T 0.0

Apply Stop

Current relation

Relation: train
Instances: 2000

Attributes: 21
Sum of weights: 2000

Attributes

All None Invert Pattern

No.	Name
1	<input checked="" type="checkbox"/> battery_power
2	<input type="checkbox"/> blue
3	<input type="checkbox"/> clock_speed
4	<input type="checkbox"/> dual_sim
5	<input type="checkbox"/> fc
6	<input type="checkbox"/> four_g
7	<input type="checkbox"/> int_memory
8	<input type="checkbox"/> m_dep
9	<input type="checkbox"/> mobile_wt
10	<input type="checkbox"/> n_cores
11	<input type="checkbox"/> pc
12	<input type="checkbox"/> px_height
13	<input type="checkbox"/> px_width
14	<input type="checkbox"/> ram
15	<input type="checkbox"/> sc_h
16	<input type="checkbox"/> sc_w
17	<input type="checkbox"/> talk_time
18	<input type="checkbox"/> three_g
19	<input type="checkbox"/> touch_screen
20	<input type="checkbox"/> wifi
21	<input type="checkbox"/> price_range

Remove

Selected attribute

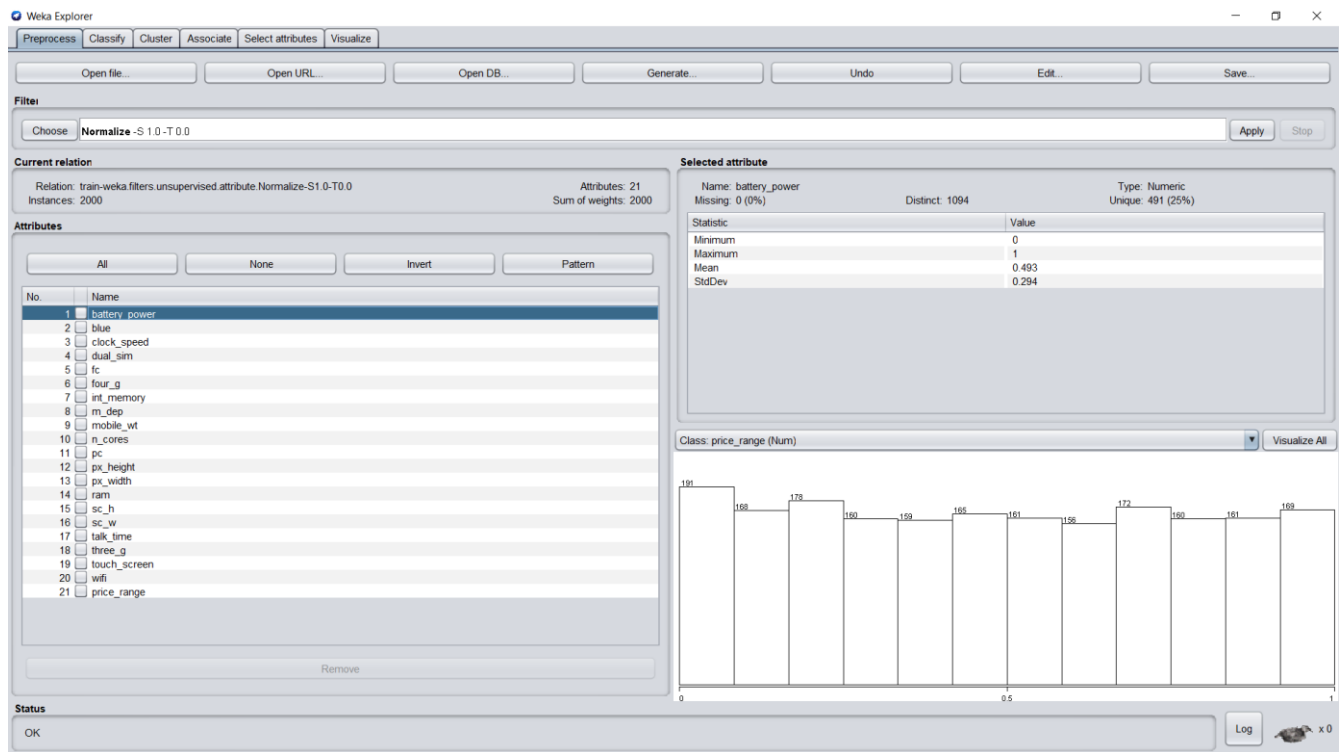
Name: battery_power
Missing: 0 (0%)
Distinct: 1094
Type: Numeric
Unique: 491 (25%)

Statistic	Value
Minimum	501
Maximum	1998
Mean	1238.518
StdDev	439.418

Class: price_range (Num)

Visualize All

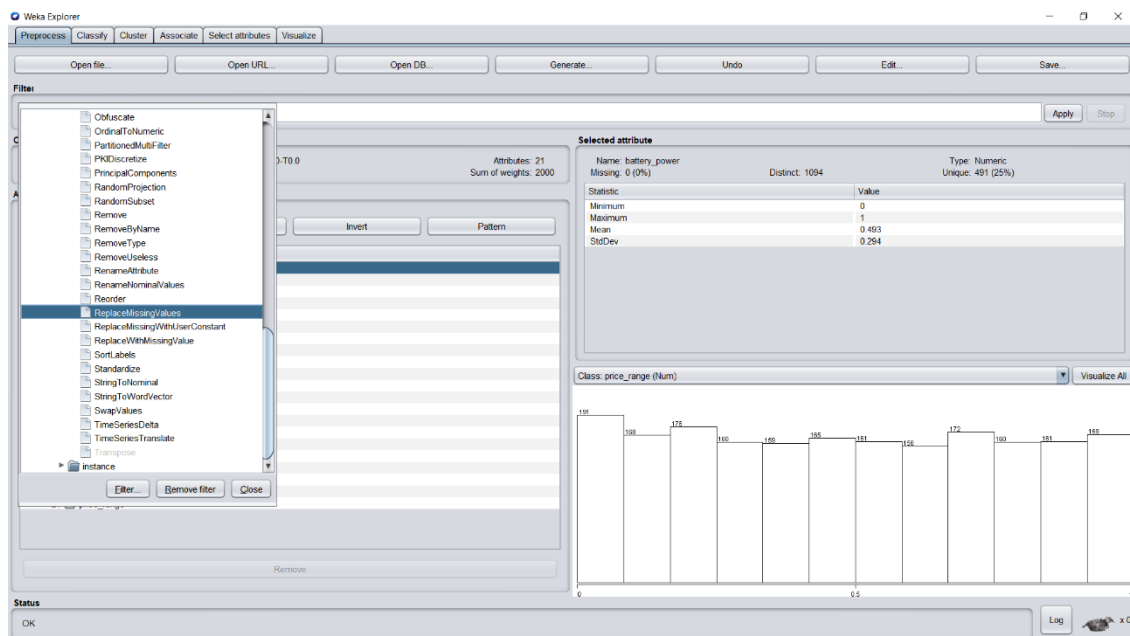
Status
OK



IMPUTACION (REEMPLAZAR DATOS FALTANTES)

La **imputación** es la sustitución de valores no informados en una observación por otros.

A veces es un paso necesario para poder tratar los datos con determinadas técnicas estadísticas de análisis. Idealmente, este análisis debería tener en cuenta el hecho de que algunos de los datos no son observados, sino que han sido imputados. Para este proceso podemos reemplazar los datos faltantes por la media, moda u otro estadístico.



Weka Explorer

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

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

Filter: Choose **ReplaceMissingValues** Apply Stop

Current relation: train-weka.filters.unsupervised.attribute.Normalize-S1.0-T0.0-weka.filters.unsupervised.attribute.Replac... Attributes: 21 Instances: 2000 Sum of weights: 2000

Selected attribute: Name: battery_power Missing: 0 (0%) Distinct: 1094 Type: Numeric Unique: 491 (25%)

Attributes

No.	Name
1	battery_power
2	blue
3	clock_speed
4	dual_sim
5	fc
6	four_g
7	int_memory
8	m_dep
9	mobile_wt
10	n_cores
11	pc
12	px_height
13	px_width
14	ram
15	sc_h
16	sc_w
17	talk_time
18	three_g
19	touch_screen
20	wifi
21	price_range

Remove

Status: OK Log

weka.gui.GenericObjectEditor

weka.filters.unsupervised.attribute.ReplaceMissingValues

About

Replaces all missing values for nominal and numeric attributes in a dataset with the modes and means from the training data.

debug: False

doNotCheckCapabilities: False

ignoreClass: False

Open... Save... OK Cancel

Information

NAME: weka.filters.unsupervised.attribute.ReplaceMissingValues

SYNOPSIS: Replaces all missing values for nominal and numeric attributes in a dataset with the modes and means from the training data. The class attribute is skipped by default.

OPTIONS: debug -- If set to true, filter may output additional info to the console. doNotCheckCapabilities -- If set, the filter's capabilities are not checked before it is built. (Use with caution to reduce runtime.) ignoreClass -- The class index will be unset temporarily before the filter is applied.

CONVERTIR STRINGS A NUMEROS (String to Nominal)

Convierte datos nominales a numeros.

Weka Explorer

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

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

Filter: Choose **StringToNominal** Apply Stop

Current relation: train-weka.filters.unsupervised.attribute.Replac... Attributes: 21 Instances: 2000 Sum of weights: 2000

Selected attribute: Name: four_g Missing: 0 (0%) Distinct: 2 Type: Numeric Unique: 0 (0%)

Statistics:

Statistic	Value
Minimum	0
Maximum	1
Mean	0.521
StdDev	0.5

Class: price_range (Num) Visualize All

Status: OK Log

StringToNominal

Obfuscate

OrdinalToNumeric

PartitionedMultiFilter

PKIDiscretize

PrincipalComponents

RandomProjection

RandomSubset

Remove

RemoveByName

RemoveByType

RemoveUseless

RenameAttribute

RenameNominalValues

Reorder

ReplaceMissingValues

ReplaceMissingWithUserConstant

ReplaceWithMissingValue

SortLabels

Standardize

StringToNominal

StringToWordVector

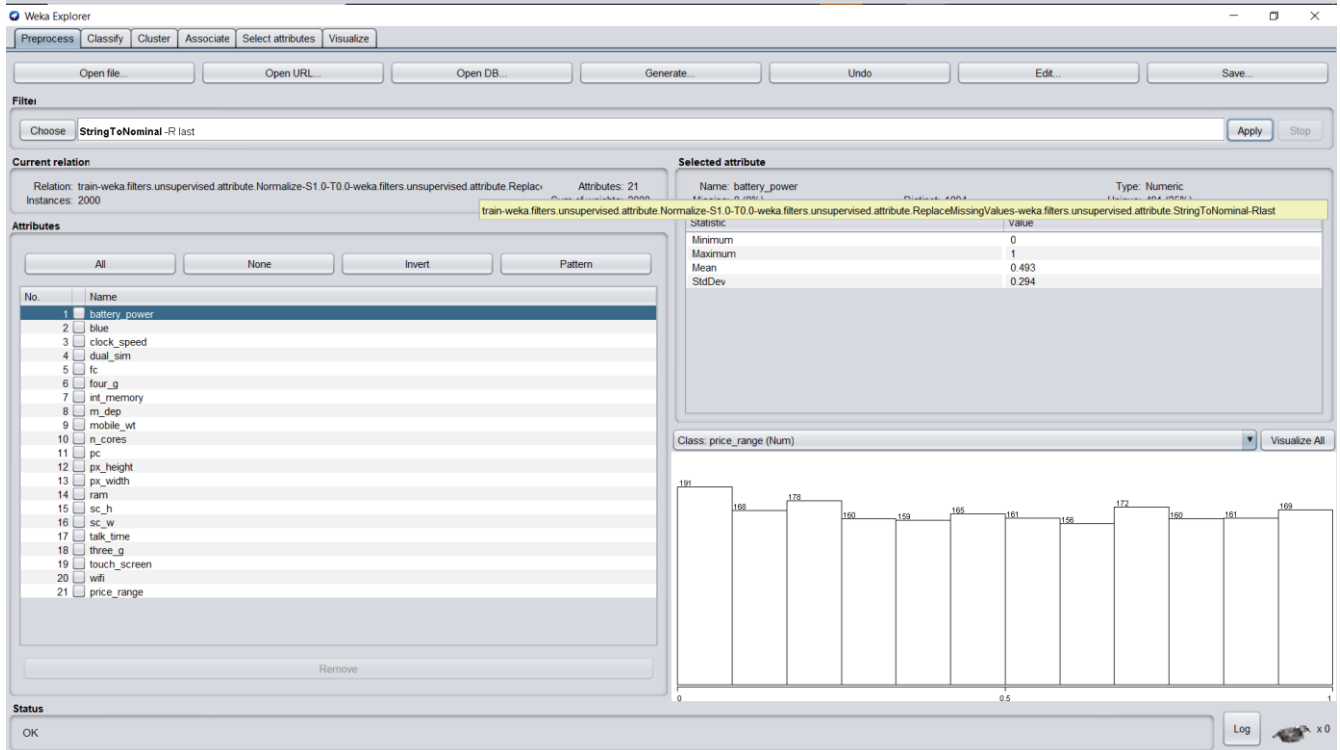
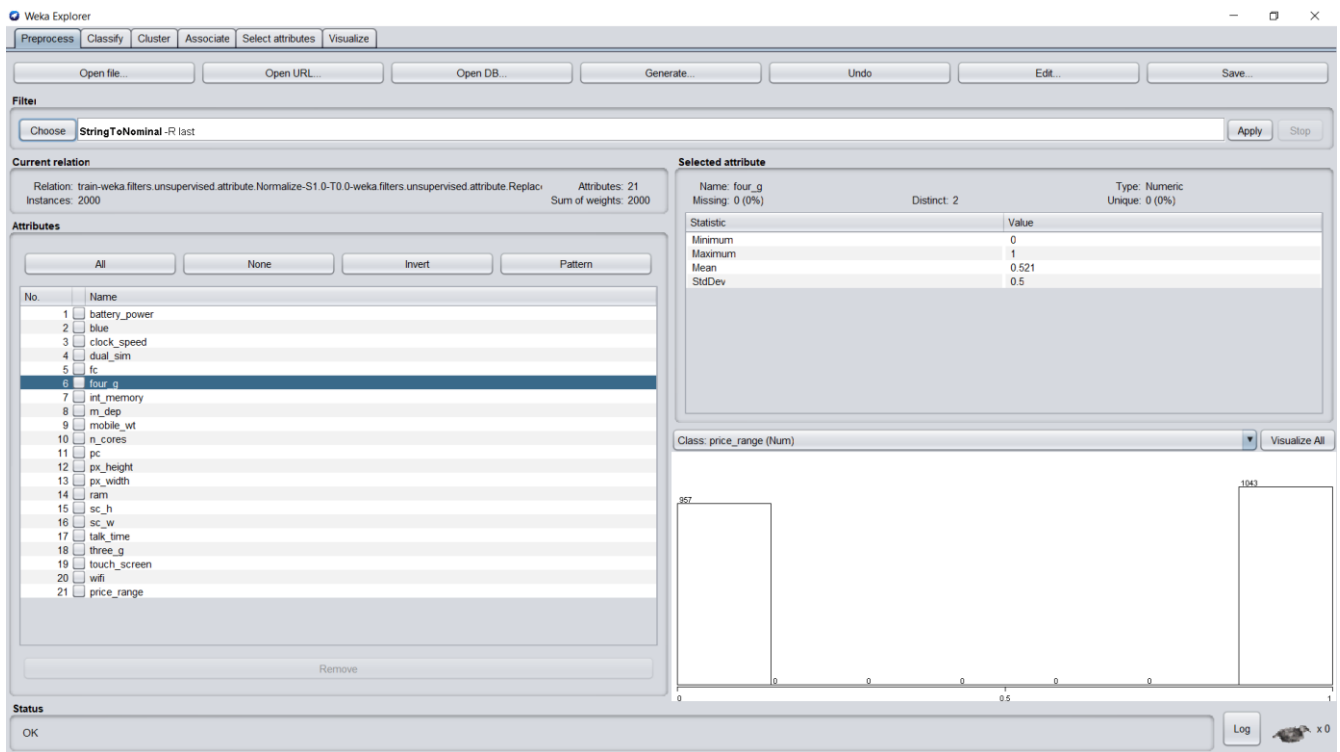
SwapValues

TimeSeriesDelta

TimeSeriesTranslate

Transpose

Instance



Weka Explorer

PreprocessClassifyClusterAssociateSelect attributesVisualize

Open file...Open URL...Open DB...Generate...UndoEdit...Save...

Filter

ChooseStringToNominal-R lastApplyStop

Current relation

Relation: train-weka.filters.unsupervised.attribute.Normalize-S1.0-T0.0-weka.filters.unsupervised.attribute.ReplaceAttributes: 21Instances: 2000Sum of weights: 2000

Selected attribute

Name: battery_powerMissing: 0 (0%)Distinct: 1094Type: NumericUnique: 491 (25%)StatisticValue

Attributes

AllNoneInvert

No.	Name
1	<input checked="" type="checkbox"/> battery_power
2	<input type="checkbox"/> blue
3	<input type="checkbox"/> clock_speed
4	<input type="checkbox"/> dual_sim
5	<input type="checkbox"/> fc
6	<input type="checkbox"/> four_g
7	<input type="checkbox"/> int_memory
8	<input type="checkbox"/> m_dep
9	<input type="checkbox"/> mobile_wt
10	<input type="checkbox"/> n_cores
11	<input type="checkbox"/> pc
12	<input type="checkbox"/> px_height
13	<input type="checkbox"/> px_width
14	<input type="checkbox"/> ram
15	<input type="checkbox"/> sc_h
16	<input type="checkbox"/> sc_w
17	<input type="checkbox"/> talk_time
18	<input type="checkbox"/> three_g
19	<input type="checkbox"/> touch_screen
20	<input type="checkbox"/> wifi
21	<input type="checkbox"/> price_range

Remove

weka.gui.GenericObjectEditor

weka.filters.unsupervised.attribute.StringToNominal

About

Converts a range of string attributes (unspecified number of values) to nominal (set number of values).MoreCapabilities

attributeRange: lastdebug: False (dropdown)doNotCheckCapabilities: False (dropdown)Open...Save...OKCancel

Information

NAMEweka.filters.unsupervised.attribute.StringToNominalSYNOPSISConverts a range of string attributes (unspecified number of values) to nominal (set number of values). You should ensure that all string values that will appear are represented in the first batch of the data.OPTIONSdebug -- If set to true, filter may output additional info to the console.attributeRange -- Sets which attributes to process ("first" and "last" are valid values and ranges and lists can also be used).doNotCheckCapabilities -- If set, the filter's capabilities are not checked before it is built. (Use with caution to reduce runtime.)

Status

OKLogx 0