

Dr. Dale E. Parson, Assignment 3, Using Weka numeric-predicting models to correlate several meteorological and temporal attributes with numeric Hawk Mountain BW migration counts for fall 2017 & fall 2018. Due by 11:59 PM on Wednesday November 13 via make turnitin.

Perform the following steps to set up for this semester's projects and to get my handout. Start out in your login directory on csit (a.k.a. acad).

```
cd $HOME
mkdir DataMine # This should already be there from assignment 1.
cp ~/parson/DataMine/csc458fall2019assn3.problem.zip DataMine/csc458fall2019assn3.problem.zip
cd ./DataMine
unzip csc458fall2019assn3.problem.zip
cd ./csc458fall2019assn3
```

EDIT THE SUPPLIED FILE README.txt when the following questions starting at Q1 below.

Keep with the supplied format, and do not turn in a Word or PDF or other file format. I will deduct 20% for other file formats, because with this many varying assignments being turned in, I need a way to grade these in reasonable time, which for me is a batch edit run on the **vim** editor.

STEP1: From the assignment directory, load file **FilteredCSC458assn3.arff** into Weka. This file is identical to the FilteredCSC458assn2.arff file that we saved in Assignment 2 with two exceptions. First, **HawkYear** is now numeric rather than a string to be converted to a nominal value. Even though discrete years such as HawkYear are typically converted to nominal sets such as {2017, 2018} because they are not continuous numeric ranges, Assignment 3 uses them as numeric data because of an AddExpression filtering step that requires numeric data. Second, **SunYear** is deleted because it is 100% redundant with HawkYear, contributing no new information, but potentially adding complexity to the learning process for some models. There are remaining attributes that are partially redundant with each other. We will inspect the effects of that redundancy.¹ **Q1** through **Q15** are worth 6.66% each.

Q1: Run the **ZeroR** classifier under functions after loading this ARFF file and paste the following results. What is the basis for the “ZeroR predicts class value”?

ZeroR predicts class value: n.n

...

Correlation coefficient	n.n
Mean absolute error	n.n
Root mean squared error	n.n
Relative absolute error	n %
Root relative squared error	n %
Total Number of Instances	2255

What is the basis for the “ZeroR predicts class value”?

¹ <https://faculty.kutztown.edu/parson/fall2019/csc458fall2019legend.html> has the attribute legend for all attributes being studied in this semester's projects.

ZeroR predicts class value: 8.289135254988913

...
Correlation coefficient -0.0719
Mean absolute error 14.8173
Root mean squared error 85.8593
Relative absolute error 100 %
Root relative squared error 100 %
Total Number of Instances 2255

What is the basis for the “ZeroR predicts class value”? 8.289 is the **MEAN** of the BW values.

Q2: Run the **LinearRegression** classifier under **functions** and paste the following results. Are the error measures² for LinearRegression better or worse than those for ZeroR? What do you think accounts for the changes in degree of error going from ZeroR to LinearRegression? Remember that LinearRegression attempts to fit the relationships of non-target attributes → target attribute to a multidimensional line.

Correlation coefficient n.n
Mean absolute error n.n
Root mean squared error n.n
Relative absolute error n %
Root relative squared error n %
Total Number of Instances 2255

Correlation coefficient 0.1035
Mean absolute error 20.656
Root mean squared error 86.0161
Relative absolute error 139.4045 %
Root relative squared error 100.1826 %
Total Number of Instances 2255

Correlation coefficient is slightly better. Errors measures got worse. ZeroR just picks the mean BW count, which in this case largely discards the non-linear outliers of Figure 1 that LinearRegression attempts to fit to a line.

Q3: Apply the **filter unsupervised -> attribute -> Normalize** after setting **ignoreClass to true**. Make sure that BW gets normalized. This step recalibrates each attribute on the scale (value – minValue)/(maxValue – minValue) as a percentage, i.e., [0, 100]% of its range. We are doing this in order to interpret weights (coefficients) in the LinearRegression formula on a normalized scale. Run **LinearRegression** on this filtered dataset and paste both the LinearRegression formula and its results. Have the error measures gotten better, worse, or stayed the same? Ignore Mean absolute error and Root mean squared error because they are no longer in application units. Instead, compare **Correlation coefficient, Relative absolute error, and Root relative squared error** to Q2 results, since those measures are always normalized.

Linear Regression Model

² Definitions for error measures appear in [Chapter 5 slides](#), slides 60 "Evaluating numeric prediction" through 64 ""Which measure?. Recall from my lecture that I usually consult **Mean absolute error** and **Root mean squared error** because they are in application units, BW counts for this dataset. The % measures are simply those values normalized across the range [0, maxErrorValue]. **Mean absolute error** emphasizes the average error, and **Root mean squared error** emphasizes the outliers, so it is much bigger **when there are significant outliers**.

BW =

PASTE THE FULL FORMULA HERE.

Correlation coefficient	n.n
Mean absolute error	n.n
Root mean squared error	n.n
Relative absolute error	n %
Root relative squared error	n %
Total Number of Instances	2255

Linear Regression Model

BW =

0.0026 * WindSpd=2: 6-11 km/h (4-7 mph),8: 62-74 km/h (39-48 mph),0: less than 1km/h (Calm) +
-0.0047 * WindDir=NNE,N,SW,SE,NE,E,Variable,ENE +
0.0049 * WindDir=N,SW,SE,NE,E,Variable,ENE +
0.0044 * WindDir=SE,NE,E,Variable,ENE +
0.0145 * Temp +
-0.0059 * CloudCover +
0.0081 * FlightDIR=WSW +
0.0047 * FlightHT=2: Unaided eye,(none),3: At limit of unaided vision,4: Binoculars (to 10X) +
0.0177 * FlightHT=4: Binoculars (to 10X) +
0.0047 * SkyCode=3: Overcast,0: Clear +
1.1805 * msnyHstart +
-0.1432 * msnyHend +
0.0058 * humidity +
0.0045 * winddir=East,NE,WNW,North +
-0.8995 * msnyWeath +
-0.1033 * msmnSunrise +
0.0109 * msmnSunset +
0.0071 * tempPrev24 +
-0.0363

Correlation coefficient	0.1035
Mean absolute error	0.0071
Root mean squared error	0.0296
Relative absolute error	139.3816 %
Root relative squared error	100.1824 %
Total Number of Instances	2255

Errors measured stayed essentially the same. Insignificant improvements in **Relative absolute error** and **Root relative squared error** may appear in student answer.

UNDO the **Normalize** filter changes after completing Q3, and verify that attributes including BW have returned to their unnormalized ranges.

Q4. What are the top six attributes in the Q3 LinearRegression formula, in terms of coefficient weights? Ignore numeric signs; compare on basis of numeric magnitude. Paste both the coefficient multiplier and its attribute from the formula, one per line as in the Weka output. Are any of these attributes redundant with each other, i.e., they have almost identical meaning or very close correlation with each other because they measure essentially the same thing? Explain.

Note that when a pair of redundant attributes have opposite signs – one positive, the other negative – then the one of smaller magnitude is used to make an adjustment on the weight of the other. When they both have the same sign, they reinforce each other.

1.1805 * msnyHstart +
-0.8995 * msnyWeath +
-0.1432 * msnyHend +
-0.1033 * msmnSunrise +
0.0177 * FlightHT=4: Binoculars (to 10X) +
0.0145 * Temp +

msnyHstart, msnyWeath, msnyHend, and msmnSunrise are redundant because they all give time of year; **msmnSunrise** correlates with time of year, although short days earlier in the year will have close **msmnSunrise** values to short days later in the year. However, we are analyzing only data from later in year.

Q5: In Weka's Select Attributes tab hit Start with the default configuration parameters and paste the Selected attributes below, including the list of actual attributes. Also, click CfsSubsetEval to inspect this Attribute Evaluators documentation paragraph, and paste that below. Which attributes from your top six in Q4 appear in Q5's results. What accounts for the difference in important attributes?

Selected attributes: Indices of selected non-target attributes

Actual list of most important non-target attributes, one per line, from Weka output

CfsSubsetEval :

Evaluates ... (complete this sentence using paste).

Selected attributes: 1,2,3,6,7,8,9,17,19,22,30 : 11

WindSpd
WindDir
Temp
FlightDIR
FlightHT
SkyCode
HawkYear
winddir
windgust
precipaccum
tempPrev24

CfsSubsetEval :

Evaluates the worth of a subset of attributes by considering the individual predictive ability of each feature along with the degree of redundancy between them.

FlightHT and Temp appear in both Q4 and Q5. The time attributes were eliminated because “along with the degree of redundancy between them” of CfsSubsetEval. Also, the time attributes in Q4 partially canceled due to numeric sign. (1.1805 * msnyHstart + -0.8995 * msnyWeath + -0.1432 * msnyHend) reduces to 0.1378 * MinuteOfYear. Finally, LinearRegression is restricted to fitting attributes onto a line, while CfsSubsetEval is not.

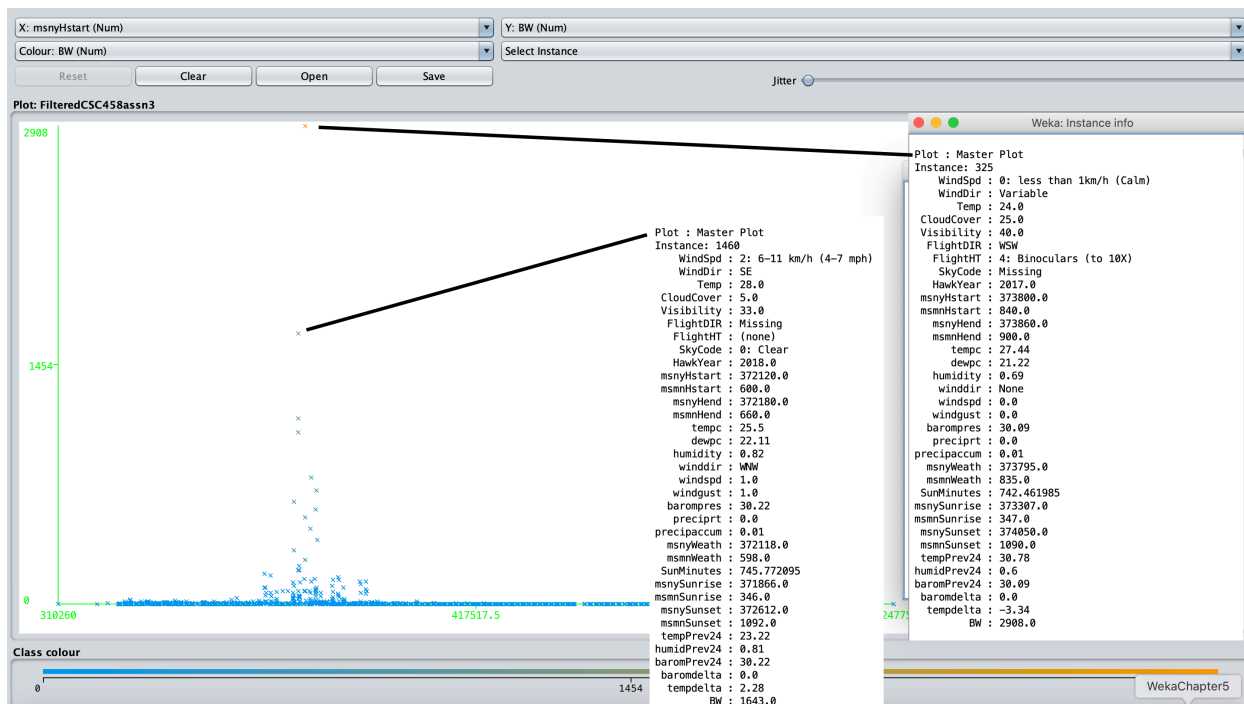


Figure 1: BW (Y axis) count as a function of msnyHstart (X axis) with max hour for 2017 & 2018 Consult HawkYear in the pasted Instance Info panels to determine which is which.

Q6: Based on Figure 1, why might LinearRegression have a hard time modeling this data?

The peak-count hours are far removed from the average line.

Q7: Run the **M5P** tree classifier and paste the following results. How do its error measures compare with ZeroR in Q1 and LinearRegression in Q2?

Number of Rules : N
 Correlation coefficient n.n
 Mean absolute error n.n
 Root mean squared error n.n
 Relative absolute error n %
 Root relative squared error n %
 Total Number of Instances 2255

M5P:
 Number of Rules : 10
 Correlation coefficient 0.2673
 Mean absolute error 12.6235
 Root mean squared error 84.7107
 Relative absolute error 85.1941 %
 Root relative squared error 98.6623 %
 Total Number of Instances 2255

LinearRegression from Q2:

Correlation coefficient	0.1035
Mean absolute error	20.656
Root mean squared error	86.0161
Relative absolute error	139.4045 %
Root relative squared error	100.1826 %
Total Number of Instances	2255

ZeroR from Q1:

ZeroR predicts class value: 8.289135254988913

...

Correlation coefficient	-0.0719
Mean absolute error	14.8173
Root mean squared error	85.8593
Relative absolute error	100 %
Root relative squared error	100 %
Total Number of Instances	2255

M5P is best so far, although error measures are not much better than ZeroR. M5P is still not doing a great job of fitting this data.

After our field trip to Hawk Mountain, a photo and summary appeared in *The Daily Brief*, the KU email-based newsletter for faculty and staff. Our Provost and a former biology professor, Dr. Anne Zayaitz, wrote me to congratulate. We got into a discussion about this dataset, and she wrote, “It’s not really surprising that the data re migration are non linear—what’s probably a bit more interesting to a biologist are those birds that are widely trailing the main group or are coming much earlier than most. We like ‘average’ but it’s the outliers that peak curiosity and questions!” We will take this approach shortly, after a few preliminary steps.

Q8: Use Weka’s EDIT window to sort instances on BW in descending order. Note the msnyHstart time of year for the maximum-BW 2017 observation hour, and separately for the maximum-BW 2018 observation hour. What are those msnyHstart values and their corresponding BW counts by year? While you are in the Edit window, scroll down and look at HawkYear-msnyHstart-BW values for the top six or seven lines. Note how the years interleave. Close the Edit window after you answer this question.

2017	msnyHstart	N	BW	N
2018	msnyHstart	N	BW	N
2017	msnyHstart	373800	BW	2908
2018	msnyHstart	372120	BW	1643

STEP1.5: In experimenting with Weka, I discovered this week that the **AddExpression** filter that we used last assignment supports the ifelse(), conditional expression construct, meaning that we can avoid using the mutating MathExpression filter. Apply AddExpression as it appears in Figure 2 to create derived attribute **msToYearPeak**, substituting the appropriate aN attribute number (like **a100** for the 100th attribute) for HawkYear and msnyHstart, and the peak BW msnyHstart values from Q8 for PeakTime2017 and PeakTime2018. After applying, re-open the Edit window, sort in descending order on BW, and make sure that the peak 2017 and 2018 BW instances have the new **msToYearPeak** derived attribute set to 0.

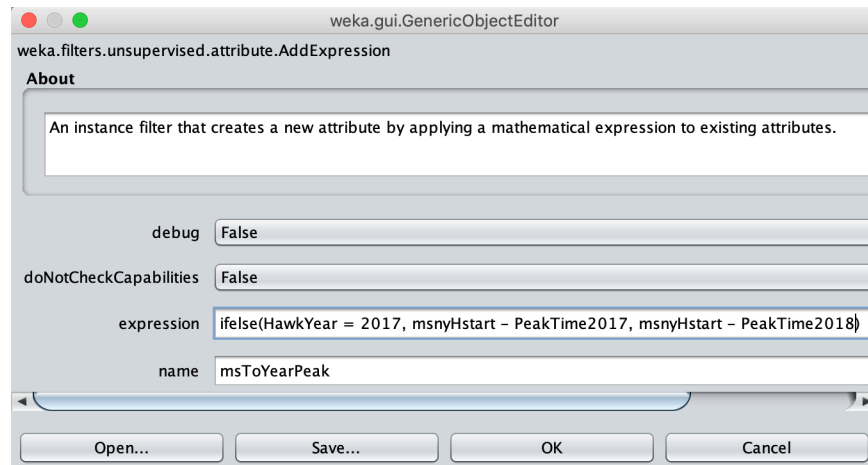


Figure 2: AddExpression for determining minutes from current msnyHstart, to msnyHstart for peak BW, by year

ifelse(a9 = 2017,a10-373800,a10-372120)

STEP2. Apply the **unsupervised -> attribute -> Reorder** filter to place BW last without disturbing the relative order of the other attributes. Remove **ALL** of the msny* prefixed attributes because they are redundant with **msToYearPeak**. Our day-of-year time base **msToYearPeak** measures the observation time's minutes relative to the peak BW hour. That helps to normalize time across years. KEEP msmnHstart, but remove **ALL OF THE OTHER** msmn* prefixed attributes, because they are redundant to some degree with msmnHstart. Keep **msToYearPeak**, of course. Remove SunMinutes because the length-of-daylight correlates with the day-of-year. There would be aliasing (roughly equal SunMinutes) between days before the start of summer, which is the longest day of the year, and corresponding days after start of summer, but since these counts begin in August, SunMinutes correlates exactly with the date, so we Remove SunMinutes. There should be 26 attributes and 2255 instances at this point, with BW last and msToYearPeak just before it. Temporarily Remove attribute HawkYear all by itself (do not Remove it with any other attribute with it – you need to UNDO this removal shortly).

Q9: Run LinearRegression and M5P on this reduced attribute, and copy & paste these values below. Also paste M5P's entire decision tree, and its **LM num: 1** linear regression formula; do not paste the other LM formulas. How do the error measures compare with LinearRegression from Q2 and M5P from Q7? How do LinearRegression versus M5P error measures within Q9 compare for dealing with average BW counts versus outliers? We will discuss the M5P decision tree when I go over the assignment's solution.

LinearRegression for Q9:

Correlation coefficient	n.n
Mean absolute error	n.n
Root mean squared error	n.n
Relative absolute error	n %
Root relative squared error	n %
Total Number of Instances	2255

M5P for Q9:

Number of Rules : N	
Correlation coefficient	n.n
Mean absolute error	n.n
Root mean squared error	n.n

Relative absolute error	n	%
Root relative squared error	n	%
Total Number of Instances	2255	

M5 pruned model tree:
(using smoothed linear models)
PASTE THE ENTIRE TREE

LM num: 1
BW = PASTE FORMULA LM1

LinearRegression for Q9:
Correlation coefficient 0.1026
Mean absolute error 20.5162
Root mean squared error 86.0106
Relative absolute error 138.4606 %
Root relative squared error 100.1762 %
Total Number of Instances 2255

M5P for Q9:
Number of Rules : 15
Correlation coefficient 0.22
Mean absolute error 13.7605
Root mean squared error 86.3709
Relative absolute error 92.8675 %
Root relative squared error 100.5958 %
Total Number of Instances 2255

LinearRegression from Q2:
Correlation coefficient 0.1035
Mean absolute error 20.656
Root mean squared error 86.0161
Relative absolute error 139.4045 %
Root relative squared error 100.1826 %
Total Number of Instances 2255

M5P from Q7:
Number of Rules : 10
Correlation coefficient 0.2673
Mean absolute error 12.6235
Root mean squared error 84.7107
Relative absolute error 85.1941 %
Root relative squared error 98.6623 %
Total Number of Instances 2255

M5 pruned model tree:
(using smoothed linear models)

msToYearPeak <= 20430 :


```

| msToYearPeak <= -10470 : LM1 (502/4.011%)
| msToYearPeak > -10470 :
| | FlightHT=2: Unaided eye,(none),3: At limit of unaided vision,4: Binoculars (to 10X) <= 0.5 :
| | | tempc <= 20.195 : LM2 (126/3.88%)
| | | tempc > 20.195 :
| | | | msToYearPeak <= 4470 : LM3 (54/36.824%)
| | | | msToYearPeak > 4470 : LM4 (49/4.203%)
| | FlightHT=2: Unaided eye,(none),3: At limit of unaided vision,4: Binoculars (to 10X) > 0.5 :
| | | msToYearPeak <= 2850 :
| | | | msToYearPeak <= -1770 :
| | | | | CloudCover <= 95 : LM5 (26/55.121%)
| | | | | CloudCover > 95 : LM6 (14/0%)
| | | | msToYearPeak > -1770 :
| | | | | CloudCover <= 27.5 :
| | | | | | tempc <= 28.14 :
| | | | | | | Temp <= 23 : LM7 (4/99.723%)
| | | | | | | Temp > 23 : LM8 (6/417.545%)
| | | | | | tempc > 28.14 : LM9 (6/77.74%)
| | | | | | CloudCover > 27.5 :
| | | | | | | tempdelta <= 1.945 :
| | | | | | | | windgust <= 0.5 : LM10 (5/31.699%)
| | | | | | | | windgust > 0.5 : LM11 (13/250.617%)
| | | | | | | tempdelta > 1.945 : LM12 (9/21.748%)
| | | msToYearPeak > 2850 :
| | | | humidPrev24 <= 0.62 : LM13 (30/2.881%)
| | | | humidPrev24 > 0.62 : LM14 (102/57.922%)
| msToYearPeak > 20430 : LM15 (1309/0.588%)

```

LM num: 1

BW =

```

0.9327 * WindDir=SE,NE,E,Variable,ENE
+ 0.0474 * Temp
- 0.0106 * CloudCover
+ 0.9967 * FlightHT=2: Unaided eye,(none),3: At limit of unaided vision,4: Binoculars (to 10X)
+ 1.1306 * FlightHT=3: At limit of unaided vision,4: Binoculars (to 10X)
+ 0.0016 * msmnHstart
+ 4.5381 * humidity
+ 0.627 * winddir=SE,East,NE,WNW,North
+ 0.2157 * winddir=NE,WNW,North
+ 0.0105 * tempPrev24
+ 0 * msToYearPeak
- 3.7315

```

M5P got a little worse than Q7 with this reduced attribute set. LinearRegression stayed about the same as Q2. The eliminate attributes were not contributing much useful data. Within Q9, M5P did better than

LinearRegression for the Mean absolute error (average BW counts), but no better than LinearRegression for Root mean squared error (outliers).

Q10: Execute UNDO once to restore only the **HawkYear** attribute, undoing its Removal. After consulting Figure 3, write an AddExpression using nested ifelse() and value comparisons on BW to create

7 distinct value ranges for BW as follows. Name this derived attribute **BWbins**. **SAVE** this ARFF 27-attribute data as **HawkData20172018.arff**. TURN this file into me by placing it in the handout directory before running **make turnitin** after you have completed the project. This is the only ARFF file you need to turn in. **Q11**. Also, paste your AddExpression for **BWbins** into README.txt. You should get a BWbins distribution that looks like Figure 4. You should also inspect the BWbins (Y) to msToYearPeak(X) visualization to see that this custom discretization looks correct.

BW range	BWbins value
0	0
1	1
2	2
<30	3
<200	4
<1000	5
Else (≥ 1000)	6

ifelse(a26=0,0,ifelse(a26=1,1,ifelse(a26=2,2,ifelse(a26<30,3,ifelse(a26<200,4,ifelse(a26<1000,5,6))))))

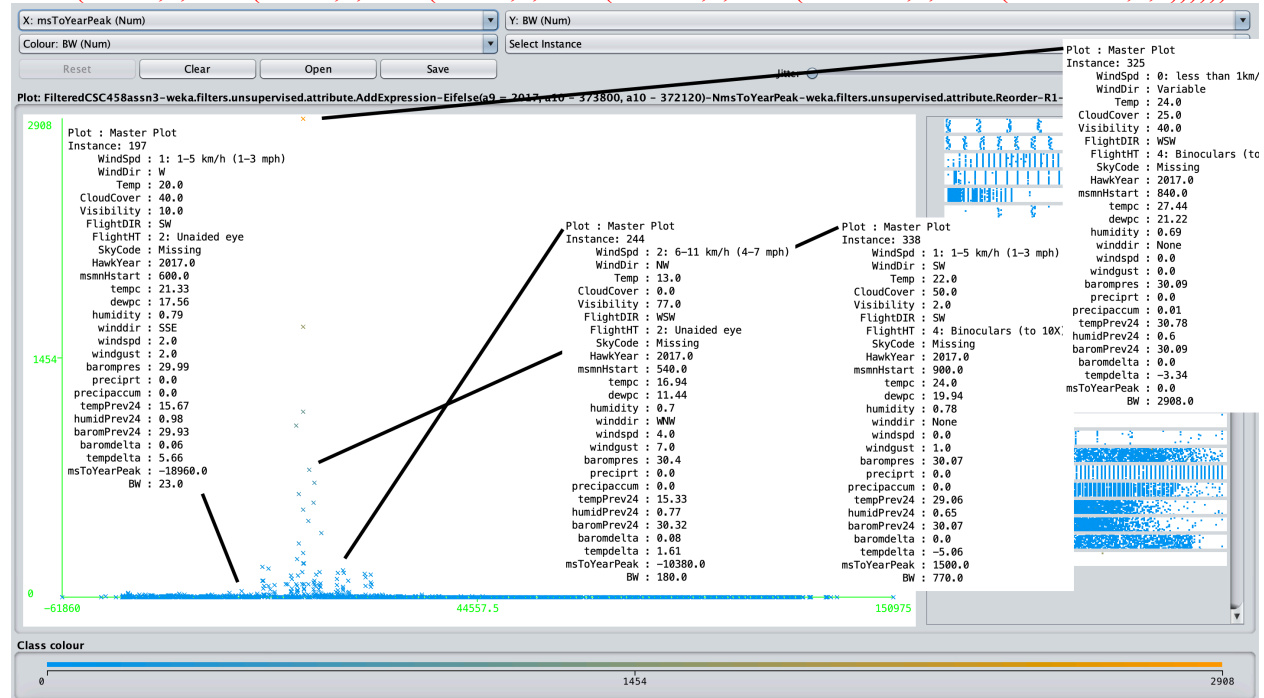


Figure 3: Regions of BW custom discretization via AddExpression: 0, 1, 2, <30, <200, <1000, ≥ 1000

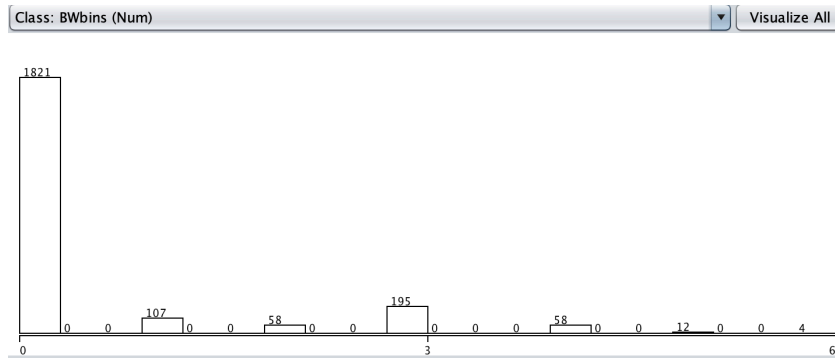


Figure 4: BWbins distribution

Q12: Remove **BW** by itself, since **BWbins** is now the target attribute, and they are non-linearly redundant. Then, temporarily Remove attribute **HawkYear** all by itself (do not Remove it with any other attribute with it – you need to UNDO this removal shortly). Run both LinearRegression and M5P against this dataset with BWbins as the target attribute, and paste indicated results below. How do you account for the changes in Correlation Coefficient and the error measures going from Q9 to Q12? (Note that **Mean absolute error** and **Root mean squared error** are in units of our 0-6 bin scale, not in BW counts, so the **Relative absolute error** and **Root relative squared error** measures are probably more useful.) Discuss the changes from Q9 in **Root relative squared error** brought about by the compression of the outlying peaks in bin 6.

LinearRegression for Q12:

Correlation coefficient	n.n
Mean absolute error	n.n
Root mean squared error	n.n
Relative absolute error	n %
Root relative squared error	n %
Total Number of Instances	2255

M5P for Q12:

Number of Rules : N	
Correlation coefficient	n.n
Mean absolute error	n.n
Root mean squared error	n.n
Relative absolute error	n %
Root relative squared error	n %
Total Number of Instances	2255

LinearRegression for Q12:

Correlation coefficient	0.5083
Mean absolute error	0.6914
Root mean squared error	0.9792
Relative absolute error	85.8272 %
Root relative squared error	86.214 %
Total Number of Instances	2255

M5P for Q12:

Number of Rules : 42

Correlation coefficient	0.7309
Mean absolute error	0.3905
Root mean squared error	0.7763
Relative absolute error	48.4781 %
Root relative squared error	68.3464 %
Total Number of Instances	2255

LinearRegression for Q9:

Correlation coefficient	0.1026
Mean absolute error	20.5162
Root mean squared error	86.0106
Relative absolute error	138.4606 %
Root relative squared error	100.1762 %
Total Number of Instances	2255

M5P for Q9:

Number of Rules : 15	
Correlation coefficient	0.22
Mean absolute error	13.7605
Root mean squared error	86.3709
Relative absolute error	92.8675 %
Root relative squared error	100.5958 %
Total Number of Instances	2255

The six bins show more linear relationships to non-target attributes, thanks to non-linear compression.

Root relative squared error dropped significantly because the “flash mob” outliers are compressed.

STEP3: The second-last stage of this assignment is to use these BWbins data from 2017 as training data, and the 2018 data as test data. UNDO once to restore the **HawkYear** attribute. Take a note in the lower right of the Preprocess window how many 2017 instances there are, and how many 2018 instances there are. Apply **unsupervised -> instance -> RemoveWithValues** to remove all 2018 instances. Check to make sure that there remain the correct number of values, then save this data as TRAIN2017.arff. Execute UNDO once to restore the 2018 instances, then apply **unsupervised -> instance -> RemoveWithValues** to remove all 2017 instances. Check to make sure that there remain the correct number of values, then save this data as TEST2018.arff. Re-load (Open file) TRAIN2017.arff. Remove attribute **HawkYear** again. (Do not turn in any TRAIN*.arff or TEST*.arff files.)

Q13: Run both **LinearRegression** and **M5P** against this dataset, using 10-fold cross validation (default) with BWbins as the target attribute, and paste indicated results below. How do Correlation Coefficient and the error measure compare to those of Q12 for **LinearRegression** and **M5P**?

LinearRegression for Q13:

Correlation coefficient	n.n
Mean absolute error	n.n
Root mean squared error	n.n
Relative absolute error	n %
Root relative squared error	n %
Total Number of Instances	1141

M5P for Q13:

Number of Rules : N	
Correlation coefficient	n.n
Mean absolute error	n.n
Root mean squared error	n.n
Relative absolute error	n %
Root relative squared error	n %
Total Number of Instances	1141

LinearRegression for Q13:

Correlation coefficient	0.5131
Mean absolute error	0.7142
Root mean squared error	0.9949
Relative absolute error	86.4248 %
Root relative squared error	86.2013 %
Total Number of Instances	1141

M5P for Q13:

Number of Rules : 6	
Correlation coefficient	0.7326
Mean absolute error	0.4121
Root mean squared error	0.7854
Relative absolute error	49.8605 %
Root relative squared error	68.0469 %
Total Number of Instances	1141

LinearRegression for Q12:

Correlation coefficient	0.5083
Mean absolute error	0.6914
Root mean squared error	0.9792
Relative absolute error	85.8272 %
Root relative squared error	86.214 %
Total Number of Instances	2255

M5P for Q12:

Number of Rules : 42	
Correlation coefficient	0.7309
Mean absolute error	0.3905
Root mean squared error	0.7763
Relative absolute error	48.4781 %
Root relative squared error	68.3464 %
Total Number of Instances	2255

About the same. LinearRegression marginally better due to more uniform 2017 data.

Q14: Run both **LinearRegression** and **M5P** against this dataset, using this 2017 data with HawkYear still removed as the training dataset, and TEST2018.arff as the external supplied test dataset, and paste indicated results below. How do Correlation Coefficient and the error measure compare to those of Q13 for LinearRegression and M5P? A degradation of more than 10% (.10) of Correlation Coefficient indicates over-fitting to the 2017 training data. Is there overfitting?

LinearRegression for Q14:

Correlation coefficient	n.n
Mean absolute error	n.n
Root mean squared error	n.n
Relative absolute error	n %
Root relative squared error	n %
Total Number of Instances	1114

M5P for Q14:

Number of Rules : N	
Correlation coefficient	n.n
Mean absolute error	n.n
Root mean squared error	n.n
Relative absolute error	n %
Root relative squared error	n %
Total Number of Instances	1114

LinearRegression for Q14:

Correlation coefficient	0.4136
Mean absolute error	0.7225
Root mean squared error	1.0351
Relative absolute error	89.9805 %
Root relative squared error	92.6878 %
Total Number of Instances	1114

M5P for Q14:

Correlation coefficient	0.5559
Mean absolute error	0.5341
Root mean squared error	0.9934
Relative absolute error	66.5234 %
Root relative squared error	88.9585 %
Total Number of Instances	1114

Both had degradation > 10%, so overfitting.

Q15: Load your saved file **HawkData20172018.arff**, and **Remove HawkYear** and **BW**, retaining BWbins as the target attribute out of 25 attributes total. Run **unsupervised -> instance -> Randomize** one time to shuffle (stratify) the 2017 and 2018 instances together. Apply **unsupervised -> instance -> RemovePercentage** with a default value of 50%, and note the number of instance that remain. SAVE this dataset as **TRAINHALF.arff**. Execute UNDO one time to restore the total 2255 instances, then run **instance -> RemovePercentage AFTER** setting **invertSelection** to **true** while leaving the percentage at 50%. SAVE this dataset as **TESTHALF.arff**. Load **TRAINHALF.arff** as the training set, then test it using **TESTHALF.arff** as the supplied external test dataset (not cross validation). Run **LinearRegression** and **M5P** and compare these Q15 results with Q14. Also, compare these Q15 results with Q12. What accounts for improvements in Q15 over Q14, given the fact that both use external test datasets of about the same size? What change do you see in over-fitting in going from Q14 to Q15, and why has that change occurred?

LinearRegression for Q15:

Correlation coefficient	n.n
Mean absolute error	n.n

Root mean squared error	n.n
Relative absolute error	n %
Root relative squared error	n %
Total Number of Instances	1128

M5P for Q15:

Number of Rules : N	
Correlation coefficient	n.n
Mean absolute error	n.n
Root mean squared error	n.n
Relative absolute error	n %
Root relative squared error	n %
Total Number of Instances	1128

LinearRegression for Q15:

Correlation coefficient	0.5024
Mean absolute error	0.7019
Root mean squared error	0.9859
Relative absolute error	87.4893 %
Root relative squared error	86.9972 %
Total Number of Instances	1128

M5P for Q15:

Correlation coefficient	0.7243
Mean absolute error	0.3932
Root mean squared error	0.7826
Relative absolute error	49.0062 %
Root relative squared error	69.0531 %
Total Number of Instances	1128

LinearRegression for Q14:

Correlation coefficient	0.4136
Mean absolute error	0.7225
Root mean squared error	1.0351
Relative absolute error	89.9805 %
Root relative squared error	92.6878 %
Total Number of Instances	1114

M5P for Q14:

Correlation coefficient	0.5559
Mean absolute error	0.5341
Root mean squared error	0.9934
Relative absolute error	66.5234 %
Root relative squared error	88.9585 %
Total Number of Instances	1114

LinearRegression for Q12:

Correlation coefficient	0.5083
Mean absolute error	0.6914
Root mean squared error	0.9792

Relative absolute error	85.8272 %
Root relative squared error	86.214 %
Total Number of Instances	2255

M5P for Q12:

Number of Rules : 42	
Correlation coefficient	0.7309
Mean absolute error	0.3905
Root mean squared error	0.7763
Relative absolute error	48.4781 %
Root relative squared error	68.3464 %
Total Number of Instances	2255

We are basically back to Q12 levels. Stratification via Randomize all but eliminated over-fitting. The fact that Randomization helps significantly implies some hidden variables – there are differences in the two years being smoothed out by Randomization. We have not really analyzed 2017 vs. 2018.

When you have completed all of your work and double-checked the assignment requirements, make sure that both **HawkData20172018.arff** saved in a previous step, and your **README.txt** that answers Q1 through Q15, are sitting in your **csc458fall2019assn3/** directory, then run **make turnitin** by the due date. Late assignments lose 10% per day late, and I will not accept an assignment after I go over its solution in class.