February 9, 2018

# Homework 1

Sanchit Singhal

CE 395R – Data Mining with Carlos Caldas

Spring 2018

University of Texas at Austin

# Pima Indians Diabetes Database

## a) Data Selected

- Number of Attributes = 8 + class (all numeric)
  - Number of times pregnant
  - Plasma glucose centration
  - Diastolic blood pressure (mm Hg)
  - Triceps skin fold thickness (mm)
  - Serum insulin (mu U/ml)
  - Body Mass Index (kg/m^2)
  - Diabetes pedigree
  - Age (years)
  - Class variable (0 or 1) – Diabetes

- Number of Instances = 768

*I downloaded this dataset in .data and .name file types. I first converted this information into .csv:*
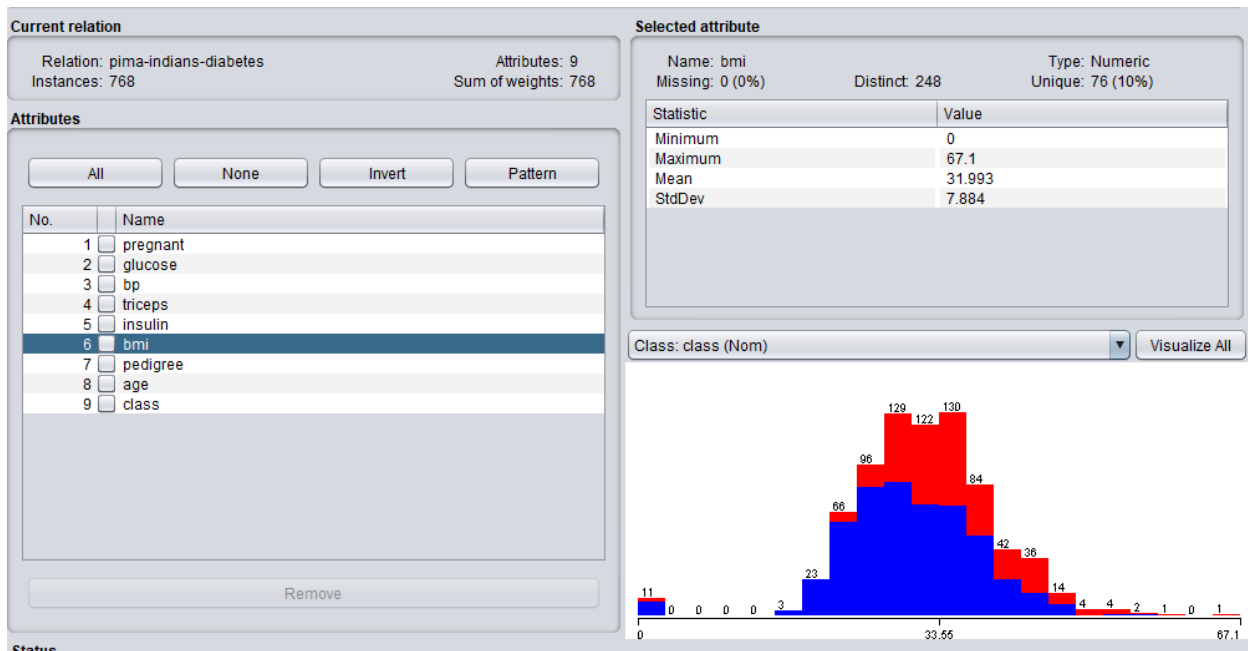
```
     pima-indians-diabetes.csv ⊠
  1  pregnant,glucose,bp,triceps,insulin,bmi,pedigree,age,class
  2  6,148,72,35,0,33.6,0.627,50,1
  3  1,85,66,29,0,26.6,0.351,31,0
  4  8,183,64,0,0,23.3,0.672,32,1
  5  1,89,66,23,94,28.1,0.167,21,0
  6  0,137,40,35,168,43.1,2.288,33,1
  7  5,116,74,0,0,25.6,0.201,30,0
  8  3,78,50,32,88,31.0,0.248,26,1
  9  10,115,0,0,0,35.3,0.134,29,0
 10  2,197,70,45,543,30.5,0.158,53,1
 11  8,125,96,0,0,0.0,0.232,54,1
 12  4,110,92,0,0,37.6,0.191,30,0
 13  10,168,74,0,0,38.0,0.537,34,1
 14  10,139,80,0,0,27.1,1.441,57,0
 15  1,189,60,23,846,30.1,0.398,59,1
 16  5,166,72,19,175,25.8,0.587,51,1
 17  7,100,0,0,0,30.0,0.484,32,1
 18  0,118,84,47,230,45.8,0.551,31,1
```

*From my csv file I converted the file into .arff for ease of use with Weka 3.8:*



```
    pima-indians-diabetes.arff ⊠
  1   @RELATION pima-indians-diabetes
  2
  3   @ATTRIBUTE pregnant REAL
  4   @ATTRIBUTE glucose REAL
  5   @ATTRIBUTE bp REAL
  6   @ATTRIBUTE triceps REAL
  7   @ATTRIBUTE insulin REAL
  8   @ATTRIBUTE bmi REAL
  9   @ATTRIBUTE pedigree REAL
 10   @ATTRIBUTE age REAL
 11   @ATTRIBUTE class {0,1}
 12
 13   @DATA
 14   6,148,72,35,0,33.6,0.627,50,1
 15   1,85,66,29,0,26.6,0.351,31,0
 16   8,183,64,0,0,23.3,0.672,32,1
 17   1,89,66,23,94,28.1,0.167,21,0
 18   0,137,40,35,168,43.1,2.288,33,1
 19   5,116,74,0,0,25.6,0.201,30,0
 20   3,78,50,32,88,31.0,0.248,26,1
 21   10,115,0,0,0,35.3,0.134,29,0
 22   2,197,70,45,543,30.5,0.158,53,1
 23   8,125,96,0,0,0.0,0.232,54,1
```

## b) Data Quality Problem Identified

*After reviewing the data in the Explorer, I realized that the BMI attribute had a minimum of 0 and there were 11 such instances. This makes no sense as nobody can have zero mass. This tells me that this is actually just missing data:*

### c) Proposed Solution

*The first step to clean this data would be to identify these instances as having missing values. I did this using the unsupervised.attribute.NumericalCleaner filter. I set the attributeIndicies to 6 to specify the bmi attribute and changed the minThreshold to slightly above 0. The minDefault was changed to NaN as that is unknown:*



*After applying the filter, I was able to identify 11 attribute values as missing which were formally set to 0 :*

*Now that I have the missing data identified, I must decide on a strategy to handle them. Two obvious methods would be to either remove them or impute them with the mean. Because we are dealing with population data, I think it would make sense to replace the missing values with the mean of the bmi attribute. This will ensure the results will not be affected much by the artificial values we impute. I did this using the unsupervised.attribute.ReplaceMissingValues:*



*After applying the filter, I was able impute the mean for the missing values for bmi attribute:*



*I exported this cleaned dataset into a new CSV.*

# Adult Income Database

## a) Data Selected

- Number of Attributes = 14 + class
  - Age (years)
  - Workclass (categorical)
  - Final Weight (continuous)
  - Education (categorical)
  - Education-num (continuous)
  - Martial-status (categorical)
  - Occupation (categorical)
  - Relationship (categorical)
  - Race (categorical)
  - Sex (categorical)
  - Capital-gain (continuous)
  - Capital-loss (continuous)
  - Hours-per-week (continuous)
  - Native-country (categorical)
  - Class (>50K, <=50K)

## b) Number of Instances = 32561

*I downloaded this dataset in .data and .name file types. I first converted this information into .csv:*
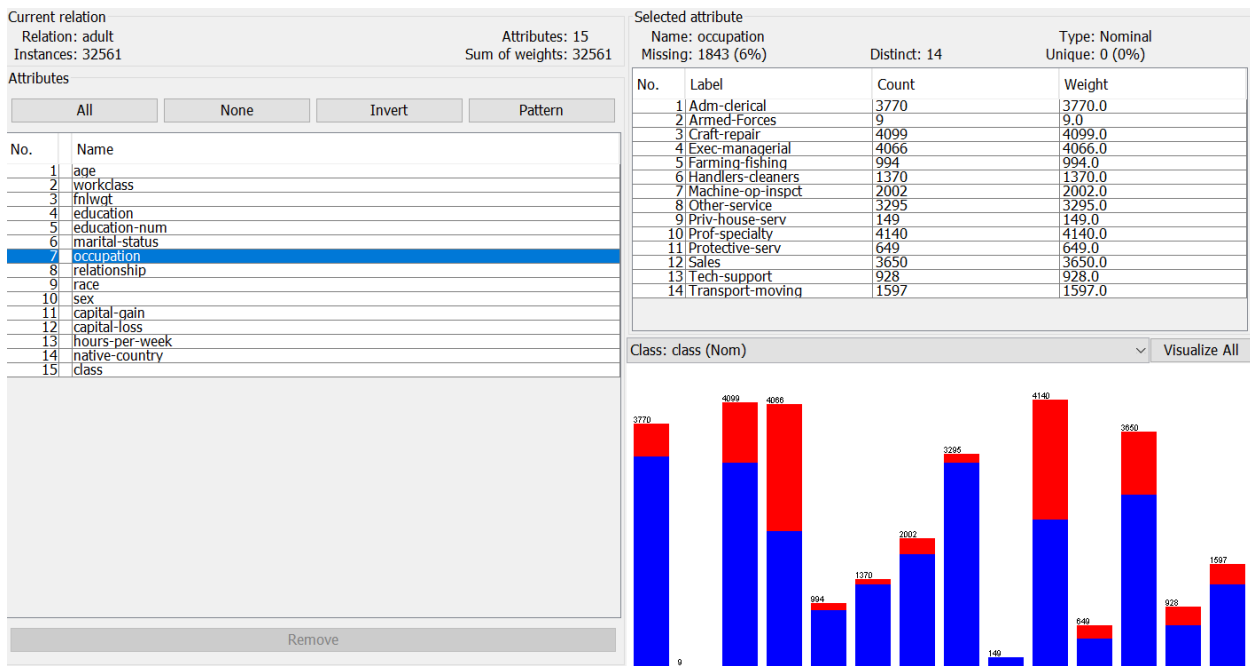


```
   adult.csv

  1  age,workclass,fnlwgt,education,education-num,marital-status,occupation,relationship,race,sex,capital-gain,capital-loss,hours-per-week,native-country,class
  2  39, State-gov, 77516, Bachelors, 13, Never-married, Adm-clerical, Not-in-family, White, Male, 2174, 0, 40, United-States, <=50K
  3  50, Self-emp-not-inc, 83311, Bachelors, 13, Married-civ-spouse, Exec-managerial, Husband, White, Male, 0, 0, 13, United-States, <=50K
  4  38, Private, 215646, HS-grad, 9, Divorced, Handlers-cleaners, Not-in-family, White, Male, 0, 0, 40, United-States, <=50K
  5  53, Private, 234721, 11th, 7, Married-civ-spouse, Handlers-cleaners, Husband, Black, Male, 0, 0, 40, United-States, <=50K
  6  28, Private, 338409, Bachelors, 13, Married-civ-spouse, Prof-specialty, Wife, Black, Female, 0, 0, 40, Cuba, <=50K
  7  37, Private, 284582, Masters, 14, Married-civ-spouse, Exec-managerial, Wife, White, Female, 0, 0, 40, United-States, <=50K
  8  49, Private, 160187, 9th, 5, Married-spouse-absent, Other-service, Not-in-family, Black, Female, 0, 0, 16, Jamaica, <=50K
  9  52, Self-emp-not-inc, 209642, HS-grad, 9, Married-civ-spouse, Exec-managerial, Husband, White, Male, 0, 0, 45, United-States, >50K
 10  31, Private, 45781, Masters, 14, Never-married, Prof-specialty, Not-in-family, White, Female, 14084, 0, 50, United-States, >50K
 11  42, Private, 159449, Bachelors, 13, Married-civ-spouse, Exec-managerial, Husband, White, Male, 5178, 0, 40, United-States, >50K
 12  37, Private, 280464, Some-college, 10, Married-civ-spouse, Exec-managerial, Husband, Black, Male, 0, 0, 80, United-States, >50K
 13  30, State-gov, 141297, Bachelors, 13, Married-civ-spouse, Prof-specialty, Husband, Asian-Pac-Islander, Male, 0, 0, 40, India, >50K
 14  23, Private, 122272, Bachelors, 13, Never-married, Adm-clerical, Own-child, White, Female, 0, 0, 30, United-States, <=50K
 15  32, Private, 205019, Assoc-acdm, 12, Never-married, Sales, Not-in-family, Black, Male, 0, 0, 50, United-States, <=50K
 16  40, Private, 121772, Assoc-voc, 11, Married-civ-spouse, Craft-repair, Husband, Asian-Pac-Islander, Male, 0, 0, 40, ?, >50K
 17  34, Private, 245487, 7th-8th, 4, Married-civ-spouse, Transport-moving, Husband, Amer-Indian-Eskimo, Male, 0, 0, 45, Mexico, <=50K
 18  25, Self-emp-not-inc, 176756, HS-grad, 9, Never-married, Farming-fishing, Own-child, White, Male, 0, 0, 35, United-States, <=50K
 19  32, Private, 186824, HS-grad, 9, Never-married, Machine-op-inspct, Unmarried, White, Male, 0, 0, 40, United-States, <=50K
 20  38, Private, 28887, 11th, 7, Married-civ-spouse, Sales, Husband, White, Male, 0, 0, 50, United-States, <=50K
 21  43, Self-emp-not-inc, 292175, Masters, 14, Divorced, Exec-managerial, Unmarried, White, Female, 0, 0, 45, United-States, >50K
 22  40, Private, 193524, Doctorate, 16, Married-civ-spouse, Prof-specialty, Husband, White, Male, 0, 0, 60, United-States, >50K
 23  54, Private, 302146, HS-grad, 9, Separated, Other-service, Unmarried, Black, Female, 0, 0, 20, United-States, <=50K
 24  35, Federal-gov, 76845, 9th, 5, Married-civ-spouse, Farming-fishing, Husband, Black, Male, 0, 0, 40, United-States, <=50K
 25  43, Private, 117037, 11th, 7, Married-civ-spouse, Transport-moving, Husband, White, Male, 0, 2042, 40, United-States, <=50K
 26  59, Private, 109015, HS-grad, 9, Divorced, Tech-support, Unmarried, White, Female, 0, 0, 40, United-States, <=50K
 27  56, Local-gov, 216851, Bachelors, 13, Married-civ-spouse, Tech-support, Husband, White, Male, 0, 0, 40, United-States, >50K
 28  19, Private, 168294, HS-grad, 9, Never-married, Craft-repair, Own-child, White, Male, 0, 0, 40, United-States, <=50K
```

*From my csv file I converted the file into .arff for ease of use with Weka 3.8:*
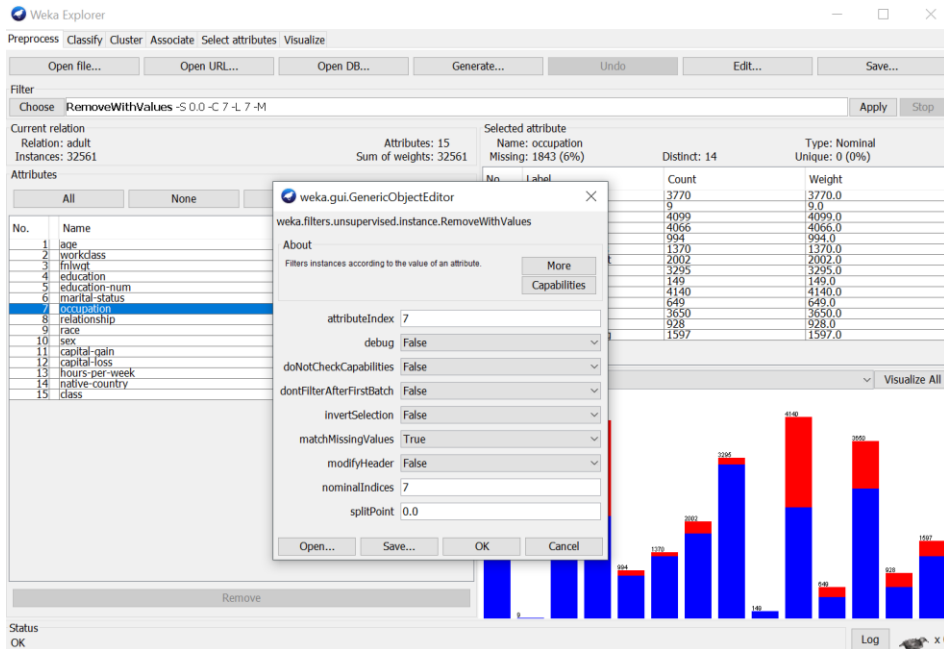


## b) Data Quality Problem Identified

*After reviewing the data in the Explorer, I realized that the occupation attribute has 1842 missing values which is roughly 6% of the total data set:*
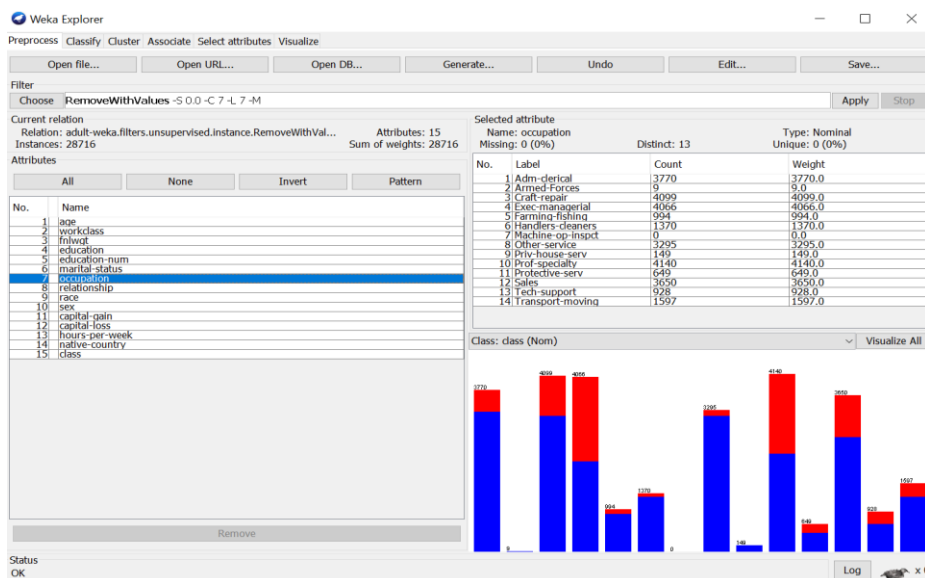
### c) Proposed Solution

*Because this is a nominal attribute, it does not make sense to try and impute them like the other database and therefore I decided to simply remove these records. I did this using the unsupervised.instance.RemoveWithValues filter. I set the attributeIndicies and nomialIndicies to 7, the index of the occupation attribute. I marked 'True' for the matchMissingValues option as we have the instances marked in the dataset already:*



*After applying the filter, I was able to get rid of the missing values for occupation attribute:*



*Like before, I exported this clean dataset into a new csv file. Obviously, this new file will contain 1842 fewer records than the original database since we have removed some records.*