# Supervised Learning

**Problem statement:** We are given car data with different car attribute. The data set is giving us the acceptance rate of cars on the basis of these attributes. We are told to do different classifiers and test which one gives better result out of the classifiers.

## Attribute Values:

Buying (Buying price)    :    v-high, high, med, low
Maint (Maintenance price)    :    v-high, high, med, low
Doors (Number of doors)    :    2, 3, 4, 5-more
Persons (accommodated person)    :    2, 4, more
Lug_boot (luggage boot)    :    small, med, big
safety    :    low, med, high

## Missing Attribute Values: No missing value

## Class Distribution (number of instances per class)

| class  | N    | N [%]       |
|--------|------|-------------|
| unacc  | 1210 | (70.023 %)  |
| acc    | 384  | (22.222 %)  |
| good   | 69   | ( 3.993 %)  |
| v-good | 65   | ( 3.762 %)  |

order to evaluate the car data, set we had to convert the data file into csv format. After opening the data.csv file we notice that some of the instances in had string value or combination of numeric value in the instances. So in order to data set in Weka we had to assume the strings as presumable numeric values. We had to leave the class attribute as string because the testing in weka depends on the type of attribute we are trying to predict. After taking the class as the attribute we want to predict we choose five classifiers and got their confusion matrix with TP and FP rate.

## Bayes Net classifier:

## Naïve Bayes:



J48:

## ZeroR:



## Kstar:

**Weka Explorer**

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

**Classifier**

Choose | KStar -B 20 -M a

**Test options**

- Use training set
- Supplied test set — Set...
- Cross-validation Folds 10
- Percentage split % 66

More options...

(Nom) class

Start | Stop

**Result list (right-click for options)**

02:07:46 - bayes.BayesNet
02:09:57 - bayes.NaiveBayes
02:14:16 - trees.J48
02:15:00 - rules.ZeroR
02:15:20 - lazy.KStar

**Classifier output**

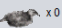```
=== Summary ===

Correctly Classified Instances        1644          95.1389 %
Incorrectly Classified Instances        84           4.8611 %
Kappa statistic                          0.8911
Mean absolute error                      0.1189
Root mean squared error                  0.1923
Relative absolute error                 51.931  %
Root relative squared error             56.8814 %
Total Number of Instances             1728

=== Detailed Accuracy By Class ===

                 TP Rate  FP Rate  Precision  Recall  F-Measure  MCC    ROC Area  PRC Area  Class
                 0.998    0.037    0.985      0.998   0.991      0.971  0.999     1.000     unacc
                 0.948    0.047    0.852      0.948   0.898      0.868  0.995     0.984     acc
                 0.662    0.001    0.977      0.662   0.789      0.798  1.000     0.995     vgood
                 0.420    0.001    0.967      0.420   0.586      0.629  0.996     0.909     good
Weighted Avg.    0.951    0.036    0.954      0.951   0.947      0.928  0.998     0.992

=== Confusion Matrix ===

   a    b    c    d   <-- classified as
1208    2    0    0 |   a = unacc
  19  364    0    1 |   b = acc
   0   22   43    0 |   c = vgood
   0   39    1   29 |   d = good
```

**Status**

OK

---

We are considering the 'unacc' class positive interest.

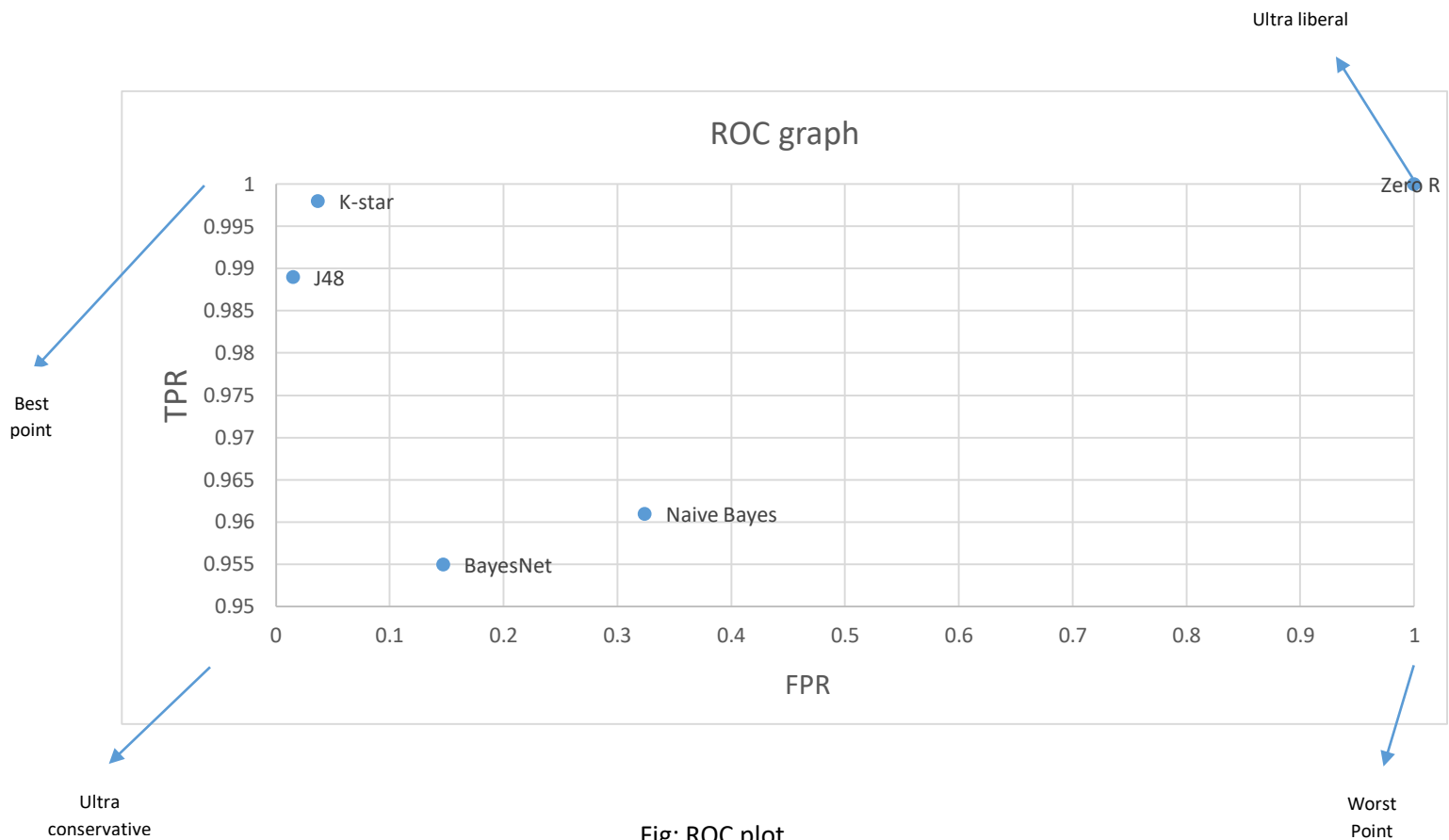| Classifiers | True positive Rate (TPR) | False Positive Rate (FPR) |
|---|---|---|
| Bayes Net | 0.955 | 0.147 |
| Naïve Bayes | 0.961 | 0.324 |
| ZeroR classifier | 1.0 | 1.0 |
| J48 | 0.989 | 0.015 |
| KStar Classifier | 0.998 | 0.037 |

Fig: ROC plot

Evaluation: From the ROC graph I should choose the classifier that gives best result out of these five classifiers. Each classifier's point is labeled in the ROC. The selected point shows the best point, worst, UL, UC points of the ROC plot. This point shows us which classifier is giving us the best relative best classification. Amongst all the classifier, Kstar and J48 classifier is the closes to the best point and farthest from the worst point possible. This worst point is very unreliable. Zero R is situated in the ultra-liberal point which is biased towards giving positive result more often. Naive Bayes and Bayes Net is close to worst point which tends to make a lot of mistake. So amongst all these classifiers the j48 classifier gives better classification of the car dataset. Amongst all these Kstaris closest to being the best point.

Decision: Kstar is the best classifier amongst these selected classifiers for car dataset.