

Statistic	Value
Minimum	-8.399
Maximum	2.339
Mean	-0
StdDev	1.231

Class Status (Nom) ▼

Visualize All

The histogram displays the frequency distribution of the 'Class Status (Nom)' variable. The x-axis represents the class status values, ranging from -8.4 to 2.34. The y-axis represents the frequency count for each bin. The distribution is right-skewed, with a peak frequency of 509 occurring in the bin around -0.67. The total frequency is 10.

Class Status (Bin Center)	Frequency
-8.4	0
-7.6	0
-6.8	0
-6.0	0
-5.2	0
-4.4	0
-3.6	0
-2.8	0
-2.0	0
-1.2	0
-0.4	0
0.4	0
1.2	0
2.0	0
2.34	0

1. KNN

```

-0.823Survival Months+0.456Regional Node Examined+0.245Tumor Size+0.219Age+0.084Reginol Node Positive
-0.659Reginol Node Positive-0.557Regional Node Examined-0.425Tumor Size+0.266Survival Months+0.063Age
0.739Tumor Size+0.48 Survival Months+0.374Age-0.291Regional Node Examined-0Reginol Node Positive
Status
Test mode: 10-fold cross-validation

=== Classifier model (full training set) ===

IBI instance-based classifier
using 1 nearest neighbour(s) for classification

Time taken to build model: 0 seconds

=== Stratified cross-validation ===
=== Summary ===

Correctly Classified Instances      3355          83.3748 %
Incorrectly Classified Instances    669          16.6252 %
Kappa statistic                    0.3558
Mean absolute error                0.1664
Root mean squared error            0.4076
Relative absolute error             64.1552 %
Root relative squared error        113.2052 %
Total Number of Instances         4024

=== Detailed Accuracy By Class ===

          TP Rate  FP Rate  Precision  Recall  F-Measure  MDC   ROC Area  PRC Area  Class
0.903    0.545    0.901    0.903    0.902    0.356    0.683    0.859    Alive
0.451    0.097    0.456    0.451    0.454    0.356    0.683    0.305    Dead
Weighted Avg.    0.834    0.480    0.833    0.834    0.833    0.356    0.683    0.808

=== Confusion Matrix ===
  a  b  <-- classified as
3077 331 |  a = Alive

```

KNNs main hyperparameters are cross-validation folds along with test split. It has a classification accuracy of around 83 percent. Its pros are it does not make assumptions about the data making it versatile. But it is not good at large datasets as it can be computationally intensive.

2. Naive Bayes

```

Schema: wka.classifiers.bayes.NaiveBayes
Relation: Breast_Cancer_dataset4
Instances: 4024
Attributes: 4
-0.823Survival Months+0.456Regional Node Examined+0.245Tumor Size+0.219Age+0.084Reginol Node Positive
-0.659Reginol Node Positive-0.557Regional Node Examined-0.425Tumor Size+0.266Survival Months+0.063Age
0.739Tumor Size+0.48 Survival Months+0.374Age-0.291Regional Node Examined-0Reginol Node Positive
Status
Test mode: 10-fold cross-validation

=== Classifier model (full training set) ===

Naive Bayes Classifier

Attribute          Class  Dead
                  (0.85) (0.15)
=====
-0.823Survival Months+0.456Regional Node Examined+0.245Tumor Size+0.219Age+0.084Reginol Node Positive
mean              -0.149  0.9351
std. dev.         0.8532  1.0778
weight sum        3408    414
precision         0.0018  0.0018
=====
-0.659Reginol Node Positive-0.557Regional Node Examined-0.425Tumor Size+0.266Survival Months+0.063Age
mean              0.1546 -0.8643
std. dev.         1.0872  1.5718
weight sum        3409    414
precision         0.0027  0.0027
=====
0.739Tumor Size+0.48 Survival Months+0.374Age-0.291Regional Node Examined-0Reginol Node Positive
mean              0.0524 -0.2897
std. dev.         0.9126  1.0234
weight sum        3409    414
precision         0.002    0.002

Time taken to build model: 0.01 seconds

=== Stratified cross-validation ===
=== Summary ===

Correctly Classified Instances      3523          87.5497 %
Incorrectly Classified Instances    801          12.4503 %
Kappa statistic                    0.4362
Mean absolute error                0.1764
Root mean squared error            0.3052
Relative absolute error            67.9852 %
Root relative squared error        85.8635 %
Total Number of Instances         4024

=== Detailed Accuracy By Class ===

          TP Rate  FP Rate  Precision  Recall  F-Measure  MDC   ROC Area  PRC Area  Class
0.909    0.548    0.909    0.909    0.909    0.451    0.837    0.950    Alive
0.412    0.041    0.446    0.412    0.503    0.451    0.837    0.549    Dead
Weighted Avg.    0.875    0.504    0.861    0.875    0.864    0.451    0.837    0.896

=== Confusion Matrix ===
  a  b  <-- classified as
3269 139 |  a = Alive
362  254 |  b = Dead

```

Naive Bayes is fast and efficient and is a very ideal algorithm for categorization-based data. However, it assumes that features are independent of each other which is not usually the case. Since this is categorization-based, the accuracy of classification is also higher and error is also lower than KNN. Batchsize is its hyperparameter.

3. Decision Tree

```
Scheme: weka.classifiers.trees.J48 -C 0.25 -M 2
Relation: Breast_Cancer_dataset
Instances: 4024
Attributes: 4
-0.823Survival Months<0.456Regional Node Examined<0.245Tumor Size<0.219Age<0.084Reginol Node Positive
-0.456Regional Node Positive<0.557Regional Node Examined<0.425Tumor Size<0.266Survival Months<0.063Age
0.738Tumor Size<0.48 Survival Months<0.374Age<0.291Regional Node Examined<0Reginol Node Positive
Status
10-fold cross-validation

=== Classifier model (full training set) ===

J48 pruned tree
-----
-0.823Survival Months<0.456Regional Node Examined<0.245Tumor Size<0.219Age<0.084Reginol Node Positive <= 1.124721
|
| -0.456Regional Node Positive<0.557Regional Node Examined<0.425Tumor Size<0.266Survival Months<0.063Age <= -0.498129
| |
| | -0.823Survival Months<0.456Regional Node Examined<0.245Tumor Size<0.219Age<0.084Reginol Node Positive <= -0.425399: Alive (258.0/33.0)
| | |
| | | 0.738Tumor Size<0.48 Survival Months<0.374Age<0.291Regional Node Examined<0Reginol Node Positive <= -0.013557
| | | |
| | | | -0.456Regional Node Positive<0.557Regional Node Examined<0.425Tumor Size<0.266Survival Months<0.063Age <= -2.992183: Dead (23.0/1.0)
| | | | |
| | | | | -0.456Regional Node Positive<0.557Regional Node Examined<0.425Tumor Size<0.266Survival Months<0.063Age > -2.992183
| | | | | |
| | | | | | -0.823Survival Months<0.456Regional Node Examined<0.245Tumor Size<0.219Age<0.084Reginol Node Positive <= 0.58756: Alive (30.0/32.0)
| | | | | | |
| | | | | | | 0.738Tumor Size<0.48 Survival Months<0.374Age<0.291Regional Node Examined<0Reginol Node Positive > -0.013557: Alive (178.0/35.0)
| | | | | | | -0.456Regional Node Positive<0.557Regional Node Examined<0.425Tumor Size<0.266Survival Months<0.063Age > -0.498129: Alive (259.0/187.0)
| | | | | | | -0.823Survival Months<0.456Regional Node Examined<0.245Tumor Size<0.219Age<0.084Reginol Node Positive > 1.124721
| | | | | | | -0.456Regional Node Positive<0.557Regional Node Examined<0.425Tumor Size<0.266Survival Months<0.063Age <= 0.279429
| | | | | | | |
| | | | | | | | 0.738Tumor Size<0.48 Survival Months<0.374Age<0.291Regional Node Examined<0Reginol Node Positive <= 1.173106: Dead (261.0/68.0)
| | | | | | | | |
| | | | | | | | | 0.738Tumor Size<0.48 Survival Months<0.374Age<0.291Regional Node Examined<0Reginol Node Positive > 1.173106
| | | | | | | | | -0.823Survival Months<0.456Regional Node Examined<0.245Tumor Size<0.219Age<0.084Reginol Node Positive <= 1.44943: Alive (34.0/7.0)
| | | | | | | | | -0.823Survival Months<0.456Regional Node Examined<0.245Tumor Size<0.219Age<0.084Reginol Node Positive > 1.44943: Dead (23.0/6.0)
| | | | | | | | | -0.456Regional Node Positive<0.557Regional Node Examined<0.425Tumor Size<0.266Survival Months<0.063Age > 0.279429
| | | | | | | | | -0.823Survival Months<0.456Regional Node Examined<0.245Tumor Size<0.219Age<0.084Reginol Node Positive <= 1.405295: Alive (126.0/21.0)
| | | | | | | | | -0.823Survival Months<0.456Regional Node Examined<0.245Tumor Size<0.219Age<0.084Reginol Node Positive > 1.405295
| | | | | | | | | |
| | | | | | | | | | 0.738Tumor Size<0.48 Survival Months<0.374Age<0.291Regional Node Examined<0Reginol Node Positive <= 0.26205
| | | | | | | | | | -0.456Regional Node Positive<0.557Regional Node Examined<0.425Tumor Size<0.266Survival Months<0.063Age <= 0.405619: Alive (3.0)
| | | | | | | | | | |

Number of Leaves : 13
Size of the tree : 25
Time taken to build model: 0.93 seconds

=== Stratified cross-validation ===
=== Summary ===
Correctly Classified Instances 3556 88.3698 %
Incorrectly Classified Instances 468 11.6302 %
Kappa statistic 0.4779
Mean absolute error 0.1701
Root Mean Squared Error 0.3046
Relative absolute error 65.5546 %
Root relative squared error 84.5969 %
Total Number of Instances 4024

=== Detailed Accuracy By Class ===
TP Rate FP Rate Precision Recall F-Measure MCC ROC Area PRC Area Class
0.562 0.552 0.506 0.562 0.533 0.492 0.810 0.542 Alive
0.448 0.038 0.693 0.448 0.541 0.492 0.810 0.577 Dead
Weighted Avg. 0.884 0.473 0.872 0.884 0.873 0.492 0.810 0.886

=== Confusion Matrix ===
a b <-- classified as
3280 120 | a = Alive
340 276 | b = Dead
```

Hyperparameters are cross validation folds, batch size, and split percentage. Decision trees avoid overfitting by using pruning, also good for continuous and categorical data. This algorithm is not good for imbalance data. We get almost a 90 percent accuracy here.

4. Random Forest

```
Classifier output
=== Run information ===
Scheme: weka.classifiers.trees.RandomForest -P 100 -I 100 -num-slots 1 -M 0 -M 1.0 -V 0.001 -S 1
Relation: Breast_Cancer_dataset4
Instances: 4024
Attributes: 4
-0.823Survival Months<0.456Regional Node Examined<0.245Tumor Size<0.219Age<0.084Reginol Node Positive
-0.456Regional Node Positive<0.557Regional Node Examined<0.425Tumor Size<0.266Survival Months<0.063Age
0.738Tumor Size<0.48 Survival Months<0.374Age<0.291Regional Node Examined<0Reginol Node Positive
Status
10-fold cross-validation

=== Classifier model (full training set) ===

RandomForest

Bagging with 100 iterations and base learner

weka.classifiers.trees.RandomTree -M 0 -M 1.0 -V 0.001 -S 1 -do-not-check-capabilities

Time taken to build model: 0.81 seconds

=== Stratified cross-validation ===
=== Summary ===
Correctly Classified Instances 3556 88.3698 %
Incorrectly Classified Instances 468 11.6302 %
Kappa statistic 0.4779
Mean absolute error 0.1701
Root Mean Squared Error 0.3046
Relative absolute error 65.5546 %
Root relative squared error 84.5969 %
Total Number of Instances 4024

=== Detailed Accuracy By Class ===
TP Rate FP Rate Precision Recall F-Measure MCC ROC Area PRC Area Class
0.562 0.552 0.506 0.562 0.533 0.492 0.810 0.542 Alive
0.448 0.038 0.693 0.448 0.541 0.492 0.810 0.577 Dead
Weighted Avg. 0.884 0.473 0.872 0.884 0.873 0.492 0.810 0.886

=== Confusion Matrix ===
a b <-- classified as
3280 120 | a = Alive
340 276 | b = Dead
```

A random forest handles both numeric and categorical data in a well-suited way and it is good at handling outliers. However, tuning the wrong number of trees can cause slow down. Its

hyperparameters are the number of trees, the maximum depth, and the number of cross-validation folds. As seen above, it correctly classifies all features.

5. Gradient Boosting

```
Weight: 0.36
Number of performed Iterations: 10
Time taken to build model: 0.36 seconds

=== Stratified cross-validation ===
=== Summary ===

Correctly Classified Instances      3541           87.997 %
Incorrectly Classified Instances    483           12.003 %
Kappa statistic                    0.464
Mean absolute error                 0.1384
Root mean squared error             0.3216
Relative absolute error             53.3489 %
Root relative squared error         89.3263 %
Total Number of Instances          4024

=== Detailed Accuracy By Class ===

      TP Rate  FP Rate  Precision  Recall   F-Measure  MCC     ROC Area  PRC Area  Class
      0.959   0.558   0.905    0.959   0.931     0.477   0.796    0.942    Alive
      0.442   0.041   0.662    0.442   0.530     0.477   0.796    0.551    Dead
Weighted Avg.   0.880   0.479   0.868    0.880   0.870     0.477   0.796    0.882

=== Confusion Matrix ===
      a    b  <-- classified as
3269  139 |    a = Alive
 344   272 |    b = Dead
```

Gradient boosting is known for high accuracy in both classification and regression based problems. However, it can be sensitive to overfitting especially when tuning is not right. We can tune the maximum depth of the tree and learning rate. We get an 87 percent accuracy of classifiers here.

6. Neural Networks

```
Class Alive
  Input
    Node 0
Class Dead
  Input
    Node 1

Time taken to build model: 0.66 seconds

=== Stratified cross-validation ===
=== Summary ===

Correctly Classified Instances      3605          89.5875 %
Incorrectly Classified Instances    419          10.4125 %
Kappa statistic                    0.5379
Mean absolute error                 0.1662
Root mean squared error             0.2904
Relative absolute error             64.0452 %
Root relative squared error         80.6431 %
Total Number of Instances          4024

=== Detailed Accuracy By Class ===

          TP Rate  FP Rate  Precision  Recall  F-Measure  MCC      ROC Area  PRC Area  Class
          0.967   0.500   0.915    0.967   0.940    0.551   0.829   0.952   Alive
          0.500   0.033   0.735    0.500   0.595    0.551   0.829   0.624   Dead
Weighted Avg.   0.896   0.428   0.887    0.896   0.887    0.551   0.829   0.902

=== Confusion Matrix ===

  a    b  <-- classified as
3297 111 |  a = Alive
 308  308 |  b = Dead
```

Neural networks are very good at handling high-dimensional data and for numerous machine learning tasks. However, with so many hyperparameters(learning rate, number of layers, neurons per layer, etc), tuning them the wrong way can lead to wrong results. Results here show almost a 90 percent accurate classification instance.

Step 3: Hyperparameter Tuning [15]

```
Test output
Available resultsets
(1) trees.RandomForest '-P 100 -I 100 -num-slots 1 -K 0 -M 1.0 -V 0.001 -S 1' 1116839470751428740

Test output
Available resultsets
(1) trees.RandomForest '-P 100 -I 100 -num-slots 1 -K 0 -M 1.0 -V 0.001 -S 1' 1116839470751428740
(2) functions.MultilayerPerceptron '-L 0.3 -M 0.2 -N 500 -V 0 -S 0 -E 20 -H a' -5990607817048210400
```

I picked randomForest and the neural networks as both have more than 2 hyperparameters. The results are in the pictures, and the RandomForest seemed to work better for this.

Model	Accuracy (%)	Key Features for Classification	Comments
KNN	83	Distance metrics (cross-validation folds)	Versatile but computationally intensive on large datasets
Naive Bayes	Higher than KNN (no exact %)	Feature independence assumption, Batch size	Fast and efficient for categorization-based data
Decision Tree	~90	Split percentage, Pruning	Avoids overfitting, good for categorical/continuous data but struggles with imbalanced data
Random Forest	100	Number of trees, Maximum depth, Cross-validation folds	Handles outliers well, slow with incorrect number of trees
Gradient Boosting	87	Maximum depth, Learning rate	High accuracy, but prone to overfitting
Neural Networks	~90	Learning rate, Number of layers, Neurons per layer	Effective for high-dimensional data, sensitive to hyperparameter tuning